



HAL
open science

Le vandalisme de l'information géographique volontaire : analyse exploratoire et proposition d'une méthodologie de détection automatique

Thérèse Quy Thy Truong

► To cite this version:

Thérèse Quy Thy Truong. Le vandalisme de l'information géographique volontaire : analyse exploratoire et proposition d'une méthodologie de détection automatique. Sciences de l'information et de la communication. Université Paris-Est, 2020. Français. NNT : 2020PESC2009 . tel-02928979

HAL Id: tel-02928979

<https://theses.hal.science/tel-02928979>

Submitted on 3 Sep 2020

HAL is a multi-disciplinary open access archive for the deposit and dissemination of scientific research documents, whether they are published or not. The documents may come from teaching and research institutions in France or abroad, or from public or private research centers.

L'archive ouverte pluridisciplinaire **HAL**, est destinée au dépôt et à la diffusion de documents scientifiques de niveau recherche, publiés ou non, émanant des établissements d'enseignement et de recherche français ou étrangers, des laboratoires publics ou privés.

THÈSE

pour obtenir le grade de docteur de l'Université Paris-Est
Spécialité : Sciences et Technologies de l'Information Géographique

Thérèse Quy Thy Truong

Le vandalisme de l'Information Géographique Volontaire

Analyse exploratoire et proposition d'une méthodologie de
détection automatique

Soutenue publiquement le 8 janvier 2020

Composition du jury :

Mme Marie-Jeanne Lesot, Sorbonne Université	Rapportrice
M. Rodolphe Devillers, Institut de Recherche pour le Développement	Rapporteur
Mme Mireille Batton-Hubert, Mines Saint-Étienne	Examinatrice
M. Jean-François Girres, Université Paul Valéry	Examineur
M. Guillaume Touya, Université Gustave Eiffel	Directeur de thèse
M. Cyril de Runz, Université de Tours	Co-directeur de thèse

Résumé

La qualité de l'information géographique volontaire est actuellement un sujet qui questionne autant les utilisateurs de données géographiques que les producteurs officiels de données voulant exploiter les bienfaits de la démarche collaborative. En effet, si les données cartographiques collaboratives présentent l'intérêt d'être ouvertes, contrairement à certaines bases de données géographiques officielles, celles-ci sont néanmoins sujettes à des erreurs voire à des dégradations volontaires, provoquées par des contributeurs malintentionnés. Dans ce dernier cas, on parle de vandalisme cartographique ou de carto-vandalisme. Ce phénomène est un travers de la démarche collaborative, et bien qu'il ne concerne qu'une faible part des contributions, il peut constituer un obstacle à l'utilisation des données cartographiques participatives. Dans une démarche de qualification de l'information géographique volontaire, ce travail de thèse a plus précisément pour objectif de détecter le vandalisme dans les données collaboratives cartographiques. Dans un premier temps, il s'agit de formaliser une définition du concept de carto-vandalisme. Puis, en partant du principe que les contributions volontairement dégradées proviennent de contributeurs malveillants, nous cherchons à démontrer que la qualification des contributeurs aide à évaluer leurs contributions. Enfin, nos expériences explorent la capacité des méthodes d'apprentissage automatique (*machine learning*) à détecter le carto-vandalisme.

Mots-clés : Information Géographique Volontaire, carto-vandalisme, qualité, apprentissage automatique

Abstract

The quality of Volunteered Geographic Information (VGI) is currently a topic that questions spatial data users as well as authoritative data producers who are willing to exploit the benefits of crowdsourcing. Contrary to most authoritative databases, an advantage of VGI is to provide open access to spatial data. However, VGI is prone to errors, even to deliberate defacement perpetrated by ill-intended contributors. In the latter case, we may speak of cartographic vandalism or carto-vandalism. This phenomenon is one of the main downsides of crowdsourcing, and despite the small amount of incidents, it may be a barrier to the use of collaborative spatial data. This thesis follows an approach based on VGI quality. In particular, the objective of this work is to detect vandalism in spatial collaborative data. First, we formalize a definition of the concept of carto-vandalism. Then, assuming that corrupted spatial data come from malicious contributors, our position is that assessing VGI contributors enables to assess VGI contributions. Finally, the experiments explore the ability of learning methods to detect carto-vandalism.

Key words : Volunteered Geographic Information, carto-vandalism, quality, machine learning

Remerciements

Tout d’abord, mes remerciements vont à mes directeurs de thèse, Guillaume Touya et Cyril de Runz, sans qui cette thèse n’aurait pu ni voir le jour, ni se conclure dans de si bonnes conditions. La qualité d’une thèse de doctorat s’explique, certes, par la qualité du doctorant, mais également par la qualité de ses encadrants. C’était pour moi une réelle chance d’avoir été encadrée par les chercheurs passionnés, réactifs, attentionnés, patients et bienveillants qu’ils ont été vis-à-vis de moi et de mon travail. Je les remercie donc chaleureusement pour ces trois années de collaboration, ainsi que pour leurs conseils avisés et leurs encouragements qui m’ont appris à repousser mes limites et à rester motivée jusqu’au bout. J’adresse également mes remerciements aux rapporteurs et aux examinateurs de mon jury de thèse, en particulier pour leur enthousiasme et le réel intérêt qu’ils ont porté à mon travail, à travers leurs remarques précieuses et pertinentes, ainsi que leurs critiques constructives à la lecture de ce mémoire et durant la soutenance.

Je souhaite remercier tous mes collègues du LASTIG de m’avoir accueillie au sein du laboratoire, et de m’avoir manifesté, chacun à leur manière, leur soutien et leur sympathie pendant ces trois années. En particulier, je suis reconnaissante envers tous ceux sur qui j’ai pu compter lorsque j’ai eu besoin d’aide, que ce soit pour résoudre des problèmes informatiques – merci Imran, Marie-Dominique et Yann –, relire mon mémoire de thèse – un grand merci Ana-Maria –, me faire répéter mon oral – merci Cécile, Ana-Maria, Bénédicte et Yann –, m’aider dans l’organisation et le bon déroulement de la soutenance – merci Ibrahim, Cécile, Ana-Maria, Catherine et Julien –, ou simplement pour m’avoir encouragée durant les moments difficiles.

Au cours de cette thèse, j’ai eu l’occasion de faire de belles rencontres avec des personnes qui, par leur gentillesse et leur amitié, ont marqué mon parcours professionnel : merci à Jérôme pour nos échanges pendant et après les journées d’enseignement à l’ENSG ; à Madiha pour ces super moments passés ensemble en conférences à L’Aquila et à Melbourne ; et merci à Jennifer pour toutes ces pauses midi ensemble au *stretch* ou pendant nos déjeuners en tête à tête ! Enfin, je tiens à remercier chaleureusement tous mes proches, famille et amis, qui, sans forcément le savoir, ont grandement contribué à la réussite de ce doctorat par leur présence réconfortante et leur soutien indéfectible. À mes parents, Quy Tân, Quy Thao, Pierre, Théo et Alice, au groupe des « ING13 » – Adélaïde, Alban, Amina, Chloé, Gabin, Hanane, Philémon, Quentin, Valentin –, mais aussi à Camille, Élisabeth, Jaamini, Myriam, Pierre-Louis, Segil, Sœur Marie-Judith et Sophie : merci de m’avoir entourée, de près comme de loin, de votre amour et de votre amitié, et de m’avoir rendu ces années de dur labeur plus douces à vivre.. !

Table des matières

Introduction	9
I État de l'art	15
1 Le vandalisme : origines et caractéristiques	16
1.1 Étymologie et évolution du vandalisme	16
1.2 Composantes du vandalisme	16
1.3 Typologie du vandalisme	17
2 Le vandalisme numérique	18
2.1 Le vandalisme dans les bases de données ouvertes	18
2.2 Le vandalisme cartographique	19
2.3 Que dégrade le carto-vandalisme ?	20
2.4 Typologie du carto-vandalisme	21
2.5 Éléments de qualité de l'information géographique	24
3 État de l'art sur la qualité de l'information géographique volontaire . .	32
3.1 Approches d'assurance qualité et méthodes de qualification	32
3.2 Mesures, indicateurs et éléments de qualité de l'information géo- graphique volontaire	33
3.3 Typologies d'évaluation	35
II Qualité des contributeurs	39
1 De la qualité du contributeur à la qualité des contributions	40
1.1 Le contributeur comme source de données à qualifier	40
1.2 L'assurance qualité par l'approche sociale	41
1.3 Collaboration et interaction	43
1.4 De la confiance des contributeurs à la fiabilité de l'information . .	50

2	Modèle de graphe social multiplexe	54
2.1	Principe et définitions	54
2.2	Modélisation des collaborations	60
2.3	Intérêts, limites du modèle et précautions	67
3	Qualification des contributeurs du projet OpenStreetMap	68
3.1	Implémentation d'un réseau social multiplexe	71
3.2	Identification de profils de contributeurs	73
3.3	Discussion	78
3.4	Conclusion et perspectives de l'étude multiplexe	85
4	Évaluation de la confiance du contributeur	86
4.1	Positionnement du problème	86
4.2	Proposition de nouvelles méthodes d'évaluation	88
4.3	Confiance et participation des contributeurs OSM	94
4.4	Comparaison des méthodes d'évaluation de la fiabilité	96
4.5	Discussion et conclusions de l'évaluation de la fiabilité des contributeurs	101
III Apprentissage du carto-vandalisme		103
1	Modes d'apprentissage automatique	104
1.1	Apprentissage non-supervisé	104
1.2	Apprentissage supervisé	105
1.3	Indicateurs de performance	106
2	Apprentissage du vandalisme dans les bases de données ouvertes	106
2.1	Métriques et méthodes	107
2.2	Cas du vandalisme cartographique	111
3	Démarche expérimentale	111
4	Construction d'un corpus de carto-vandalisme	113
4.1	État de l'art des méthodes de construction d'un corpus de vandalisme	113
4.2	Modélisation d'un corpus de carto-vandalisme	117
4.3	Construction d'un premier corpus	118
5	Le carto-vandalisme comme anomalie à détecter	120

5.1	Mise en place de descripteurs	121
5.2	Sélection des descripteurs optimaux	125
5.3	Étude de la dépendance des descripteurs optimaux avec la zone d'étude	131
5.4	Analyse des faux positifs	134
5.5	Détection d'anomalies sur un jeu de données dépourvu de carto- vandalisme synthétique	139
5.6	Paramétrage du système de détection	142
5.7	Bilan de la détection non-supervisée du carto-vandalisme	143
6	Détection du carto-vandalisme par apprentissage supervisé	144
6.1	Forêts aléatoires (<i>random forest</i>)	144
6.2	Réseau de neurones à convolution (CNN)	153
7	Vers une amélioration du corpus de carto-vandalisme	164
	Conclusion	166
	Conclusion Générale	167
	Bibliographie	174
	Annexes	187
	A Exploration des utilisateurs bannis du projet OSM	189
	B Étude de la fiabilité de 30 contributeurs OSM	195
	C Ajout de profils synthétiques de carto-vandales	249

Introduction

Contexte général

Ces dernières années, le développement des technologies de l'information et de la communication a bouleversé la société quant à l'usage des données numériques. D'une part, alors que l'avènement du Web 2.0 a grandement facilité le partage et la diffusion de l'information sur Internet, une politique d'ouverture des données numériques (*open data*), tant au niveau national¹ qu'international², a permis de libérer peu à peu l'accès et l'utilisation des données publiques et privées. D'autre part, les progrès techniques ont permis de développer des appareils de plus en plus connectés et sophistiqués, si bien que leur démocratisation a conduit à une explosion de la production de données numériques, en quantité et en fréquence.

Dans ce contexte, de nombreuses initiatives de type sciences participatives ont émergé, en exploitant la capacité de la foule à produire des données, dans le but de répondre à des problématiques diverses. On peut mentionner des projets collaboratifs qui visent à résoudre des problèmes sociétaux³ ou environnementaux⁴, à alerter sur des incidents⁵, ou encore à faire avancer la recherche scientifique⁶. Par ailleurs, il arrive que dans certains projets collaboratifs, les producteurs de données soient également les utilisateurs de ces données. C'est, par exemple, le cas de l'encyclopédie libre Wikipédia⁷, qui permet à tout utilisateur de consulter des articles, ou d'enrichir la base de données ouverte par la rédaction et la correction d'articles. Les données collaboratives sont donc autant utilisées dans des cadres professionnels que privés.

Dans le cas de l'information géographique, les projets collaboratifs se sont particulièrement développés face à des situations où les données géographiques institutionnelles – *i.e.* produites par des institutions officielles au niveau local ou national – faisaient défaut, soit par leur absence (Rollason *et al.*, 2018), soit par leur obsolescence. Un des projets de cartographie collaboratifs les plus populaires est OpenStreetMap⁸ qui, depuis 2004, a pour vocation de constituer une base de données cartographique du monde. Ce projet a notamment permis de motiver des groupes

1. <https://www.etalab.gouv.fr/>

2. <https://www.opengovpartnership.org/policy-area/digital-governance/>

3. www.openideo.com

4. www.climatecolab.org/

5. www.safecity.in/

6. www.zooniverse.org

7. fr.wikipedia.org

8. www.openstreetmap.org

de contributeurs pour des actions humanitaires⁹, par exemple, sur des zones à cartographier d'urgence dans des situations de crise (Scholz *et al.*, 2018).

L'information géographique volontaire

On désignera par « information géographique volontaire » la donnée spatialisée saisie de manière active dans le cadre d'un projet collaboratif. Toutefois, il existe différents types d'information géographique volontaire, qui peut provenir d'une collecte plus ou moins active de la part du contributeur, et dont le caractère spatial est plus ou moins fort. La Figure 2 illustre notre délimitation de l'information géographique volontaire dans l'ensemble du contenu généré par les utilisateurs d'Internet.

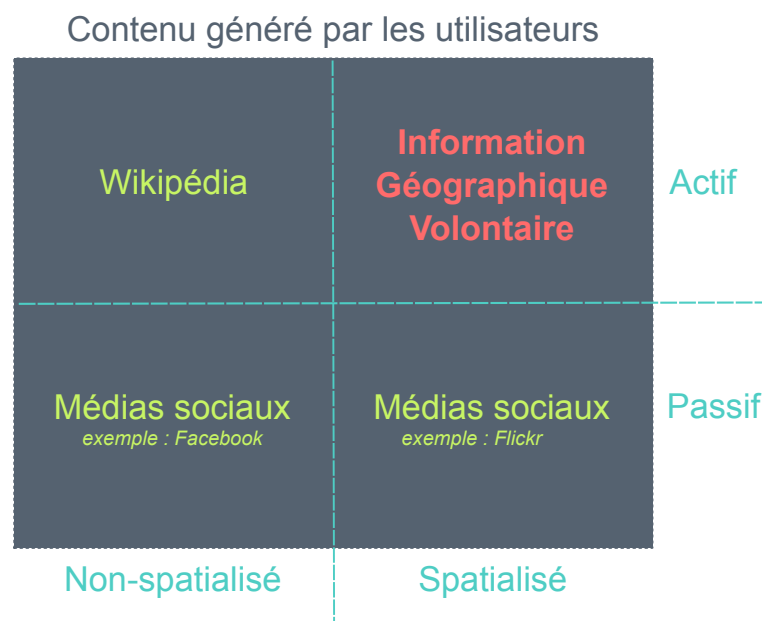


FIGURE 2. Délimitation de l'information géographique volontaire, tiré et adapté de (See *et al.*, 2016)

Des données spatialisées

Une donnée est dite « spatialisée » lorsqu'elle est référencée géographiquement, autrement dit, lorsqu'elle est géo-référencée ou géo-localisée. Dans notre travail, nous nous plaçons dans le cas où l'information géographique volontaire revêt un haut degré de spatialité, c'est-à-dire dans le cas des données spatiales. Parmi les plateformes collaboratives de données spatiales, nous pouvons citer OpenStreetMap, Wikimapia et Google Map Maker. En revanche, la spatialité de l'information géographique volontaire peut varier dans d'autres plateformes. Par exemple, l'application de partage d'images Flickr est une plateforme d'information géographique volontaire où les utilisateurs peuvent charger des photographies géolocalisées, qui constituent donc des données spatialisées. Cependant, ces images ne sont pas des données spatiales, au

9. www.hotosm.org/

sens que l'information spatiale des images est secondaire par rapport au contenu photographique de ces données.

La spatialité de l'information géographique volontaire peut être évaluée en s'interrogeant sur l'objectif premier de la plateforme collaborative dont elle provient. Dans le cas de Flickr, l'objectif premier de la plateforme est de partager des photographies avec ses proches : l'aspect géographique des données est donc secondaire. En revanche, OpenStreetMap, Wikimapia et Google Map Maker sont des plateformes où la dimension spatiale des données est au cœur du projet collaboratif. Dans ce travail, nous parlerons de « données géographiques collaboratives » pour désigner les contributions issues de ce type de projet d'information géographique volontaire.

Des données saisies volontairement

Dans cette thèse, l'information géographique volontaire correspond aux données qui ont été saisies de manière active par des contributeurs volontaires. Cette participation active est à distinguer d'un mode passif de collecte de données, où les contributeurs partagent simplement les données saisies par les capteurs mobiles de leurs *smartphones* par exemple. En particulier, il existe aujourd'hui de nombreuses applications dans lesquelles les utilisateurs sont appelés à partager leurs positions géographiques, telles que Foursquare pour recommander des lieux d'intérêts, ou Strava pour enregistrer les traces GPS de leurs activités sportives. Bien que les utilisateurs de ces applications contribuent effectivement de l'information géographique volontairement, certains n'ont pas conscience des utilisations qui peuvent être faites à leurs dépens, telles que la publicité personnalisée par exemple. Dans ce dernier cas, on peut parler d'information géographique involontaire (Fischer, 2012) ou d'information géographique obtenue par la force (McKenzie et Janowicz, 2014).

Problématique

L'information géographique volontaire a permis de pallier certains manques des données géographiques d'autorité en donnant la possibilité à toute personne volontaire d'alimenter des bases de données géographiques ouvertes, par l'ajout ou la mise à jour d'information. La facilité d'accès et d'édition de ces bases de données géographiques ouvertes offre la possibilité de réutiliser une information géographique riche et actuelle. Cette nouvelle source d'information est devenue une alternative qui, lorsqu'elle co-existe avec des données institutionnelles, peut même rivaliser avec ces dernières. Face au fort potentiel de l'information géographique volontaire, les producteurs institutionnels de données géographiques ont alors porté un réel intérêt à la démarche collaborative, dans le but d'en tirer parti afin d'améliorer le système traditionnel d'acquisition et de mise à jour de leurs bases de données (Olteanu-Raimond *et al.*, 2017, 2018).

Les données géographiques institutionnelles ou industrielles présentent des erreurs et des imprécisions d'origines diverses dont l'évaluation constitue un sujet de recherche à part entière (Girres, 2012; Batton-Hubert et Pinet, 2019). Par ailleurs, l'information géographique volontaire est également entachée d'une imperfection qui

nécessite d'être évaluée. Bien que la démarche collaborative présente un fort potentiel pour produire des données géographiques de qualité, la saisie de données par des volontaires comprend également des risques qui peuvent compromettre cette qualité.

En effet, la qualité de l'information géographique volontaire dépend fortement de la participation de ses contributeurs. Or, ces derniers présentent des profils assez variés, à commencer par leur niveau d'expérience, puisqu'ils peuvent être des non-professionnels voire des non-spécialistes (Neis et Zipf, 2012). De plus, les contributeurs peuvent se différencier par les connaissances qu'ils ont à partager, leurs modes de participation, ainsi que par leurs diverses motivations. L'information géographique volontaire qui en résulte est donc de qualité très variable. Pour faire bon usage d'une information, il est primordial d'avoir étudié sa qualité au préalable. Par conséquent il convient de chercher à qualifier l'information géographique volontaire.

En particulier, l'information géographique volontaire est potentiellement sujette à des erreurs d'édition, qui peuvent s'expliquer par la possible incompétence de certains contributeurs. Il existe cependant des contributions provenant d'actes malveillants, qui consistent à saisir volontairement des données erronées. Ce phénomène de vandalisme cartographique – aussi appelé carto-vandalisme – peut représenter une réelle menace pour la qualité de ces données. En parallèle, l'information géographique volontaire nécessite d'être évaluée selon une ou plusieurs typologies données, et l'identification du vandalisme dans ces données est une forme de qualification de l'information géographique volontaire selon une typologie spécifique du vandalisme. Par ailleurs, pour être utilisable, il est nécessaire que l'information géographique volontaire présente un niveau de qualité minimal. La détection du vandalisme dans l'information géographique volontaire pourrait donc constituer une première garantie *a minima* de la qualité, indépendamment de l'utilisation qui en sera faite.

Ainsi, la problématique de cette thèse porte sur la détection du vandalisme dans l'information géographique volontaire. Ce travail implique de définir le phénomène de vandalisme dans ce contexte particulier, et de chercher à détecter les contributions qui en relèvent. Ce travail de thèse s'inscrit donc dans une démarche de qualification de l'information géographique volontaire.

Objectifs

Le premier objectif de ce travail consiste à définir précisément ce qu'est le vandalisme cartographique. Il s'agit d'identifier les caractéristiques, les causes de ce phénomène et la manière dont il se manifeste dans l'information géographique volontaire. Pour observer empiriquement ce phénomène, nous utiliserons, dans nos expériences, les contributions issues du projet cartographique OpenStreetMap, bien que notre recherche porte sur l'information géographique volontaire en général.

Le phénomène de vandalisme cartographique est perpétré par des contributeurs malveillants. Par conséquent, le deuxième objectif consiste à qualifier les contributeurs d'information géographique, pour qu'à partir de l'identification des contributeurs malveillants, nous puissions détecter les contributions de carto-vandalisme.

Le troisième objectif vise à trouver des méthodes adéquates à la détection des

contributions cartographiques vandalisées. Nous étudierons différentes approches d'apprentissage automatique dans le but d'évaluer leur capacité à qualifier l'information géographique volontaire selon cette typologie. Cela suppose de s'interroger sur les métriques à utiliser pour modéliser l'information géographique volontaire dans le cadre de la détection du vandalisme cartographique.

Enfin, la détection du vandalisme cartographique est étudiée dans le but de qualifier l'information géographique volontaire. Notre travail vise ultimement à contribuer à une méthodologie globale de qualification de l'information géographique volontaire. Cette méthodologie, appliquée ici au cas particulier de la détection du vandalisme cartographique, doit être également valable pour d'autres typologies de qualification de l'information géographique volontaire.

Plan

Le chapitre 1 propose un état de l'art sur le vandalisme cartographique et la qualité de l'information géographique volontaire. Le phénomène de vandalisme cartographique étant encore peu connu et peu étudié, nous cherchons à le définir en revenant sur les différentes définitions du vandalisme depuis ses origines historiques jusqu'à son apparition dans le domaine numérique et cartographique. Par ailleurs, comme le vandalisme cartographique affecte la qualité des données, nous revenons sur les concepts de qualité de l'information géographique volontaire et les méthodes d'évaluation proposés dans la littérature scientifique.

Le chapitre 2 se focalise sur la qualité des contributeurs de données géographiques. Comme le carto-vandalisme provient de contributeurs malveillants, nous cherchons à qualifier les contributeurs de données géographiques collaboratives. Nous proposons de modéliser les interactions entre les contributeurs sur la plateforme cartographique par un graphe social multiplexe. L'étude de ce modèle permet d'identifier des profils de contributeurs qui pourraient témoigner d'un certain niveau de confiance. Puis, à partir des éléments de qualité des contributeurs issus de l'étude de notre modèle multiplexe, nous cherchons à évaluer la confiance des contributeurs. Il s'agit en particulier d'étudier les métriques et les méthodes qui permettent de qualifier le plus précisément possible la fiabilité et la compétence des contributeurs. De cette manière, nous pourrions détecter le vandalisme cartographique parmi les contributions qui proviennent de contributeurs qualifiés comme étant peu fiables.

Le chapitre 3 traite de la détection du vandalisme cartographique par apprentissage automatique. Nous passons d'abord en revue les méthodes existantes utilisées pour détecter le vandalisme dans d'autres projets collaboratifs comme Wikipédia. Puis, en menant des expériences d'apprentissage automatique sur des données du projet OpenStreetMap, nous cherchons à mieux comprendre comment le carto-vandalisme se manifeste dans ces données, et à étudier le potentiel de méthodes d'apprentissage supervisé et non-supervisé pour le détecter automatiquement.

Chapitre I

État de l'art sur le vandalisme cartographique et la qualité des données géographiques volontaires

Dans ce chapitre, nous proposons une définition formelle du carto-vandalisme. Nous abordons le concept général du vandalisme en revenant sur les origines historiques et les motivations psychologiques de ce phénomène. Puis, en dégagant les composantes du vandalisme exprimées dans le domaine législatif, nous les utilisons comme critères de qualification du vandalisme de contributions OpenStreetMap dont les utilisateurs ont été bannis de la plateforme. Cette mise en situation permet ainsi de valider la définition du carto-vandalisme, qui sera considérée dans la suite de notre exposé. Enfin, puisque la détection du vandalisme constitue une forme de qualification des données collaboratives cartographiques, nous dressons un état de l'art des méthodes d'évaluation de la qualité de l'information géographique volontaire.

1 Le vandalisme : origines et caractéristiques

1.1 Étymologie et évolution du vandalisme

Historiquement, le terme « vandalisme » tire son origine des actes de dégradation perpétrés par la tribu des Vandales, un peuple germanique réputé pour avoir pillé et saccagé les œuvres d'art et les monuments lors de leur invasion de l'Occident au V^e siècle. Plus tard, durant la période révolutionnaire, ce terme est repris par l'abbé Henri Grégoire pour dénoncer les dommages causés sur « les objets nationaux, qui, n'étant à personne, sont la propriété de tous » (Grégoire, 1794). Selon l'abbé Grégoire, ces actes « *ont pour cause l'ignorance ; il faut l'éclairer : la négligence ; il faut la stimuler : la malveillance [...]; il faut [la] comprimer.* » Aujourd'hui, le vandalisme est un crime sanctionné par la justice française, qui le définit ainsi :

« L'acte de vandalisme consiste à détruire, dégrader, ou détériorer volontairement le bien d'autrui. [...] L'acte de vandalisme doit être commis sans motif légitime. Il est par exemple permis de briser une vitre pour sauver une personne en danger »¹

1.2 Composantes du vandalisme

Les composantes du vandalisme varient en fonction de la définition qui lui est attribuée (Figure I.1). Dans le contexte législatif, un acte est qualifié de vandalisme s'il remplit les trois conditions suivantes :

- il correspond à une **détérioration du bien d'autrui** ;
- la dégradation a été **voulue** par son auteur ;
- **aucun élément de contexte ne justifie la légitimité** de l'acte.

Notons que, d'après la première condition, le vandalisme inclut les détériorations faites sur les biens privés dans le cas où l'auteur du vandalisme n'en est pas le propriétaire. La définition législative est donc plus large que celle initialement proposée par l'abbé Grégoire, qui se limitait uniquement aux biens publics. Toutefois, dans le cas du vandalisme de biens privés, il est nécessaire de s'interroger sur l'auteur de l'infraction afin de vérifier qu'il ne soit pas aussi le propriétaire du bien détérioré, à la différence du vandalisme sur le patrimoine culturel.

Par ailleurs, quel que soit le type de bien dégradé, la deuxième condition exige de s'intéresser à l'auteur de la dégradation afin de déterminer s'il a agi intentionnellement. En effet, à la différence de la définition historique de l'abbé Grégoire, la définition législative du vandalisme ne tient pas compte des cas où l'auteur de la dégradation a agi par ignorance ou négligence, mais considère uniquement le cas où l'intention du vandale est engagée. L'identité du potentiel vandale peut suggérer l'intentionnalité de l'acte. Par exemple, faire dérailler un train en déposant des obstacles sur les rails sera considéré comme une farce si cela a été fait par des enfants – en supposant que ces derniers n'ont pas encore conscience de l'ampleur de leurs actes – mais comme du vandalisme si cela a été commis par des adultes (Zimbardo, 1971).

1. <https://www.service-public.fr/particuliers/vosdroits/F1514>

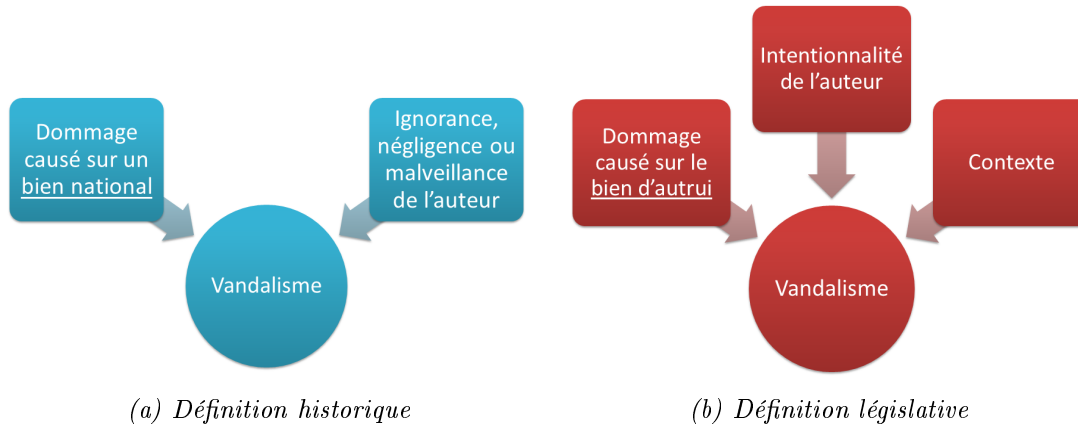


FIGURE I.1. Évolution des composantes du vandalisme (selon le contexte français)

Enfin, la troisième condition indique la prise en compte d'éléments de contexte pour qualifier un acte de vandalisme. Il s'agit de déterminer si l'acte volontaire de dégradation du bien d'autrui a été commis pour de bonnes raisons. En reprenant l'exemple donné dans la définition législative, briser une vitre pour sauver une personne en danger n'est pas un crime puisque cet acte est justifié par une bonne cause, alors que briser une vitre pour un cambriolage est considéré comme du vandalisme.

L'évaluation des trois composantes du vandalisme est ardue. En effet, il convient de s'interroger sur le sens de la dégradation d'un bien : est-ce rendre ce bien inutilisable ou simplement méconnaissable ? De plus, établir les raisons et l'intentionnalité d'un potentiel vandale à partir de l'observation d'un bien dégradé fait appel à de véritables compétences juridiques. Par conséquent, qualifier un acte de vandalisme est loin d'être trivial.

1.3 Typologie du vandalisme

D'après l'état de l'art du vandalisme dressé par [van Vliet \(1984\)](#) dans le domaine socio-psychologique, la typologie du vandalisme proposée par [Cohen \(1973\)](#) est la plus aboutie. Celle-ci se compose de 5 catégories de vandalisme, chacune correspondant à un but différent :

1. **Le vandalisme cupide** (*acquisitive vandalism*) : il s'agit de détériorer un objet pour s'emparer des biens qui y sont contenus. Briser une vitre pour cambrioler une maison est un exemple de vandalisme cupide.
2. **Le vandalisme tactique ou idéologique** (*tactical/ideological vandalism*) : la dégradation est produite dans le but d'attirer l'attention sur une injustice, une cause idéologique ou forcer une réaction. Par exemple, des graffitis ou des bris de vitrines produits par des manifestants pour exprimer leurs revendications sont des formes de vandalisme idéologique.
3. **Le vandalisme vindicatif** (*vindictive vandalism*) : l'objectif est de se venger sur le propriétaire du bien dégradé. C'est par exemple le cas lorsque, face à une situation vécue comme une injustice par des étudiants, ces derniers vandalisent leur établissement scolaire pour se venger.

4. **Le vandalisme ludique** (*play vandalism*) : il s'agit ici de dégrader volontairement le bien d'autrui dans un simple but d'amusement.
5. **Le vandalisme malicieux** (*malicious vandalism*) : dans ce cas, l'acte de dégradation est une fin en soi. Il est inspiré par un sentiment d'ennui, de frustration ou de colère.

À cette classification, Sutton (1987) y ajoute un autre type de vandalisme motivé par le statut auprès des pairs, correspondant à un acte causé pour obtenir ou maintenir le statut de son auteur vis-à-vis de ses pairs. Cette dernière catégorie concerne surtout les actes malveillants produits pas des gangs.

Toutefois, qualifier les actes de vandalisme selon cette typologie est délicat. En effet, il est difficile de deviner les intentions d'un vandale à partir de la seule observation des dégradations perpétrées. Par exemple, l'œuvre intitulée *Graffiti is a crime* (Figure I.2) de l'artiste britannique Banksy peut être sujet à plusieurs interprétations : d'une part, ce graffiti peut sembler provocateur en revendiquant, avec ironie, l'art présent dans cette forme de vandalisme, ce qui le rapproche d'un vandalisme idéologique ; d'autre part, le ton humoristique de ce graffiti peut aussi amener à le classer comme un vandalisme de type ludique. Par ailleurs, le vandalisme de Banksy est intéressant car la dimension artistique de ses œuvres soulèvent des questionnements sur le caractère dégradant de ses actes de vandalisme.



FIGURE I.2. « *Graffiti is a crime* » – Source : personnelle, septembre 2019. Reproduction du pochoir original de Banksy situé sur un mur à Manhattan.

2 Le vandalisme numérique

2.1 Le vandalisme dans les bases de données ouvertes

Si le vandalisme portait historiquement sur des dégradations du monde physique, le développement du domaine numérique a entraîné la propagation de ce phénomène sur les données numériques. En effet, l'évolution de l'espace numérique comme lieu de partage et d'échange a éveillé une forme de désinhibition du comportement chez

certaines internautes (Suler, 2004). En particulier, les plateformes collaboratives sont particulièrement sujettes à des actes de vandalisme, car elles ouvrent à des personnes volontaires la saisie et la modification de ces bases de données.

Les projets de *crowdsourcing*, tels que Wikipédia et Wikidata, se sont rapidement trouvés confrontés au problème de vandalisme de l'information. En effet, la liberté accordée aux contributeurs de données laisse la possibilité aux plus malintentionnés de participer en incorporant des informations – principalement textuelles – qui dégradent ces bases de données ouvertes. Wikipédia² et Wikidata³ sont tous deux des projets dont les bases de données ouvertes appartiennent à la Wikimedia Foundation⁴. Bien que ces deux plateformes soient liées, la première est une base de données ouverte non-structurée alors que la seconde est structurée (Heindorf *et al.*, 2016; Sarabadani *et al.*, 2017).

Le vandalisme dans Wikipédia est défini de la manière suivante⁵ :

Le terme de vandalisme désigne un comportement visant délibérément à porter atteinte à l'objectif encyclopédique du projet [...] Une modification erronée ou maladroite mais visant de bonne foi à améliorer le projet, où l'intention de dégradation n'est pas manifeste, ne doit cependant pas être considérée comme du vandalisme, quand bien même elle ne respecterait pas un consensus préexistant.

De nombreux exemples sont donnés sur la page Wikipédia dédiée au sujet du vandalisme, de manière à le distinguer des cas d'erreur ou de non-respect des règles de bonne conduite (participer à des controverses, avoir un comportement obstiné, *etc.*). Le vandalisme sur Wikipédia se produit par une dégradation des articles, à travers la suppression de donnée ou l'ajout d'informations inexactes, voire offensantes, typiquement par l'utilisation de propos grossiers. Les classifications du vandalisme sur Wikipédia ont principalement été proposées selon les types d'actions permises sur la plateforme de contribution (Chin *et al.*, 2010; Mola Velasco, 2011).

2.2 Le vandalisme cartographique

Le vandalisme cartographique désigne les actes de vandalisme provoqués sur les plateformes collaboratives cartographiques. Le terme de carto-vandalisme a été introduit pour désigner la dégradation intentionnelle d'éléments cartographiques numériques dans le contexte de l'information géographique volontaire (Ballatore, 2014). La plateforme cartographique collaborative Wikimapia définit le vandalisme comme toute action délibérée qui vise à corrompre l'information contribué⁶.

En revanche, le vandalisme dans le projet cartographique OpenStreetMap (OSM) comprend également l'ignorance des règles de consensus mises en place sur la plateforme⁷. Cette définition du carto-vandalisme dans le projet OSM se différencie donc

2. fr.wikipedia.org

3. www.wikidata.org

4. wikimediafoundation.org

5. <https://fr.wikipedia.org/wiki/Wikip%C3%A9dia:Vandalisme> (Consulté le 27 août 2019)

6. <http://wikimapia.org/user/tools/guidelines/>

7. <https://wiki.openstreetmap.org/wiki/FR:Vandalisme>

de la définition traditionnelle dans laquelle il était initialement question de l'intentionnalité de la dégradation des données⁸. Dans cet exposé, nous nous plaçons dans la continuité des travaux menés sur le vandalisme dans les bases de données ouvertes. Nous emploierons donc le terme de carto-vandalisme pour désigner toute action **délibérée** qui vise à porter atteinte à un projet collaboratif cartographique.

2.3 Que dégrade le carto-vandalisme ?

Les plateformes d'information géographique volontaire permettent la saisie, la modification et la suppression de données géographiques. Le vandalisme étant défini comme un acte intentionnel de dégradation du bien d'autrui et sans motif légitime, il convient de s'interroger sur la nature du bien dégradé par le carto-vandalisme. En effet, une dégradation comporte la notion de transformation d'un objet d'un état initial vers un état de qualité réduite. Si ce bien dégradé correspondait à la contribution – *i.e.* l'objet cartographique – alors le carto-vandalisme ne concernerait que les éléments cartographiques existants qui diminuent en qualité. Or, le carto-vandalisme fantaisiste consiste à créer des objets cartographiques totalement fictifs : dans ce cas, la notion de transformation d'un état initial vers un état endommagé ne serait pas observée sur ces objets. En revanche, la création de tout objet géographique transforme l'espace cartographique.

Nous considérons donc que le bien dégradé par tout acte de carto-vandalisme est l'espace cartographique de la plateforme collaborative. En effet, l'espace cartographique est une modélisation du monde réel, avec un certain degré d'abstraction, dans le but d'être utilisé dans différentes applications, comme par exemple la visualisation cartographique ou l'analyse spatiale. Cet espace étant ouvert, il peut être considéré comme un bien commun : il est, par excellence, le bien d'autrui, puisqu'il n'appartient à personne. En effet, l'espace cartographique, en tant que base de données ouvertes (spécifique au domaine géographique), fait partie des biens communs numériques que forment les connaissances, l'information et le réseau numérique internet (Le Crosnier, 2011). Ainsi, toute édition cartographique – que ce soit un ajout, une modification, ou une suppression de donnée – induit une transformation de l'espace cartographique. Une contribution faite intentionnellement de mauvaise qualité sans motif légitime dégrade l'espace cartographique, et constitue donc un acte de carto-vandalisme.

À l'image de la tragédie des biens communs théorisée par Hardin (1968) pour expliquer la destruction des ressources naturelles par la poursuite des intérêts individuels, des travaux plus récents sur la théorie des biens communs ont poussé la réflexion vers une tragédie des biens communs numériques (Greco et Floridi, 2004). L'infosphère⁹ – c'est-à-dire le monde des biens numériques – souffre aussi d'un problème de pollution, lié dans ce cas à une utilisation irraisonnée des ressources numériques conduisant à une surproduction des données, qui corrompt la communication d'information par du bruit, tel que le spam.

8. « In the case of OSM, vandalism can occur intentional and unintentional, contradicting the traditional definition of the term "vandalism". » (Neis *et al.*, 2012)

9. Ce néologisme est à comprendre à l'image de la biosphère qui désigne le monde des biens naturels

L'extension de la tragédie des biens communs numériques soulève des questions éthiques, cependant, les problématiques autour de la conservation des ressources naturelles ne sont pas toutes transposables aux biens numériques (Hess et Ostrom, 2007). L'équilibre de la biosphère nécessite une sauvegarde des biens communs, et ne pas agir en ce sens participe à la destruction de l'environnement. Or, ce n'est pas tout à fait le cas dans l'infosphère, où les « actes manqués » n'altèrent pas l'équilibre de cet écosystème numérique (Greco et Floridi, 2004). Ainsi, dans le cadre du carto-vandalisme, le fait de s'abstenir volontairement de contribuer à l'amélioration de l'espace collaboratif cartographique n'est pas un acte de dégradation. En effet, ne pas contribuer des données qui pourraient compléter ou améliorer la qualité de l'espace cartographique ne va pas y induire de transformation et donc de dégradation. Par conséquent, le carto-vandalisme « par omission » n'existe pas, mais il ne concernera que les contributions actives, c'est-à-dire l'ajout, la modification et la suppression de données géographiques.

Les conséquences du carto-vandalisme sur l'espace cartographique peuvent se situer à deux niveaux différents. D'une part, la dégradation peut affecter la capacité de l'espace cartographique à représenter le monde réel de manière fidèle. Cela concerne, par exemple, les contributions qui introduisent des données cartographiques imaginaires. D'autre part, la dégradation peut affecter la cohérence spatiale de l'espace cartographique, et le rendre inutilisable dans certaines applications (navigation, analyse de biodiversité, de morphologie urbaine, *etc.*).

2.4 Typologie du carto-vandalisme

a) Typologie des motivations

À partir de l'observation empirique d'incidents produits sur des plateformes collaboratives cartographiques Wikimapia et OpenStreetMap, Ballatore (2014) propose une classification du carto-vandalisme basée sur la typologie de Cohen (1973), présentée dans la Section 1.3. Cette classification liste des types de carto-vandalisme à partir des motivations du vandale :

1. **Le carto-vandalisme ludique** : il est équivalent au vandalisme ludique dans le contexte collaboratif cartographique. Ici, les mauvaises éditions cartographiques (suppression d'objets, distorsions géométriques) du carto-vandale sont faites dans un but d'amusement.
2. **Le carto-vandalisme idéologique** : il est équivalent au vandalisme idéologique. Les dégradations portent sur les objets cartographiques représentant des lieux à caractère politique, religieux ou ethnique.
3. **Le carto-vandalisme fantaisiste** : il s'agit de cartographier des éléments fictifs, de type ville imaginaire, qui n'existent pas dans la réalité.
4. **Le carto-vandalisme artistique** : il correspond à l'expression artistique de certains contributeurs qui utilisent une plateforme collaborative comme un éditeur graphique pour y dessiner des objets de forme artistique.
5. **Le carto-vandalisme industriel** : il concerne les dégradations perpétrées dans le but de réduire la réputation ou la crédibilité de projets cartographiques ouverts. Par exemple, des groupes industriels investis dans le marché

des données spatiales pourraient user de ce genre de pratique pour favoriser l'utilisation de leurs données.

6. **Le carto-spam** : il correspond à des éditions massives de messages promotionnels qui n'ont pas de lien avec la zone géographique éditée. Par exemple, il peut s'agir d'indiquer sur les données attributaires d'objets géographiques des liens commerciaux (vers des enseignes, des marques de produits, des lieux touristiques, *etc.*) pour gagner en visibilité.

Bien qu'elle repose sur des cas réels de carto-vandalisme, cette typologie ne permet pas de classer chaque cas de carto-vandalisme dans une catégorie exclusive. Autrement dit, les catégories peuvent se chevaucher : par exemple, le carto-vandalisme artistique peut également être vu comme étant du carto-vandalisme fictif, puisque les objets artistiques sont aussi imaginaires.

Par ailleurs, en comparant cette typologie à celle de [Cohen \(1973\)](#), nous pouvons observer qu'il n'y a pas de bijection entre les deux. La Figure I.3 schématise la correspondance des deux typologies aux motivations. Le carto-vandalisme dit « ludique » peut être motivé par un sentiment d'ennui ou de frustration. Or, ces deux types de motivations entraînent respectivement des actes de vandalisme ludique et malicieux de [Cohen \(1973\)](#). De même, le carto-vandalisme fantaisiste peut être vu comme du vandalisme ludique car tous deux sont justifiés par les mêmes motivations (l'humour ou l'expression de soi). L'enchevêtrement entre ces deux classifications s'explique par le fait que les motivations ne sont pas rangées dans les mêmes classes.

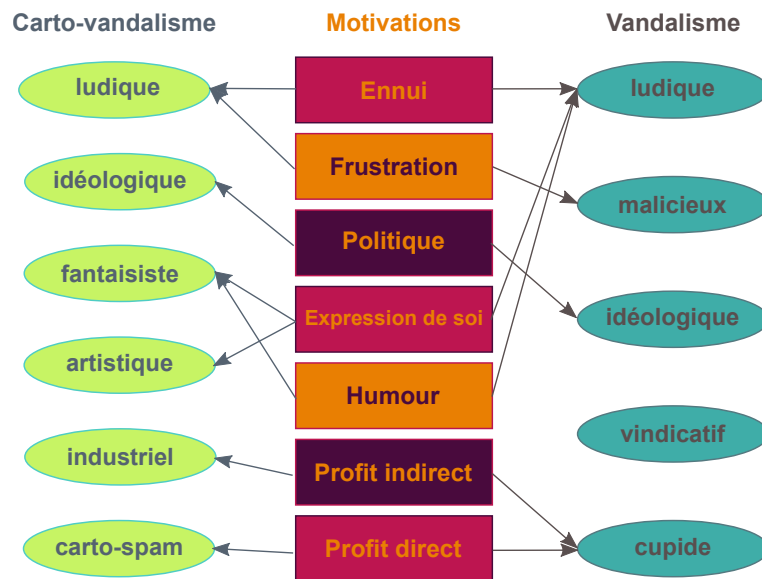


FIGURE I.3. Schéma de comparaison des typologies du carto-vandalisme ([Ballatore, 2014](#)) et du vandalisme ([Cohen, 1973](#))

Contrairement à la typologie [Cohen \(1973\)](#) où chaque motivation est associée à un seul type de vandalisme, on observe que la motivation « expression de soi » peut se manifester par du carto-vandalisme fantaisiste ou du carto-vandalisme artistique. Or, traditionnellement, tout phénomène peut se définir à partir de 4 « causes », d'après la philosophie d'Aristote ¹⁰ :

10. Aristote, Physique, II, 7

1. la **cause matérielle** définit ce dont une chose est faite (*de quoi est-elle faite ?*);
2. la **cause efficiente** définit ce qui donne l'existence à la chose (*par quoi/par qui est-elle faite ?*);
3. la **cause formelle** définit la forme de cette chose (*comment est-elle faite ?*);
4. la **cause finale** définit la destination de la chose (*pour quoi est-elle faite ?*)

Le carto-vandalisme est un phénomène qui se produit sur l'information géographique volontaire (cause matérielle) par des contributeurs malveillants (cause efficiente). La cause formelle du carto-vandalisme – c'est-à-dire ce qui définit sa forme – se trouve en partie décrite dans la typologie de [Ballatore \(2014\)](#). En effet, certaines catégories regroupent une forme particulière sous laquelle se manifeste le carto-vandalisme : le carto-vandalisme artistique désigne l'insertion de géométries complexes ; le carto-vandalisme fantaisiste concerne la cartographie d'objets fictifs. En revanche, les autres catégories de cette typologie décrivent la cause finale du phénomène, c'est-à-dire la raison de son existence. En effet, le carto-vandalisme ludique correspond aux actes expliqués par le sentiment d'ennui des contributeurs ; le carto-vandalisme industriel est motivé par la volonté de discréditer le projet collaboratif cartographique, notamment par certains groupes industriels ; le carto-spam est typiquement motivé par des objectifs de marketing poursuivis par des établissements commerciaux. Par conséquent, la typologie de [Ballatore \(2014\)](#) veut se fonder principalement sur la cause finale du carto-vandalisme, mais nous y trouvons également des types qui spécifient sa cause formelle.

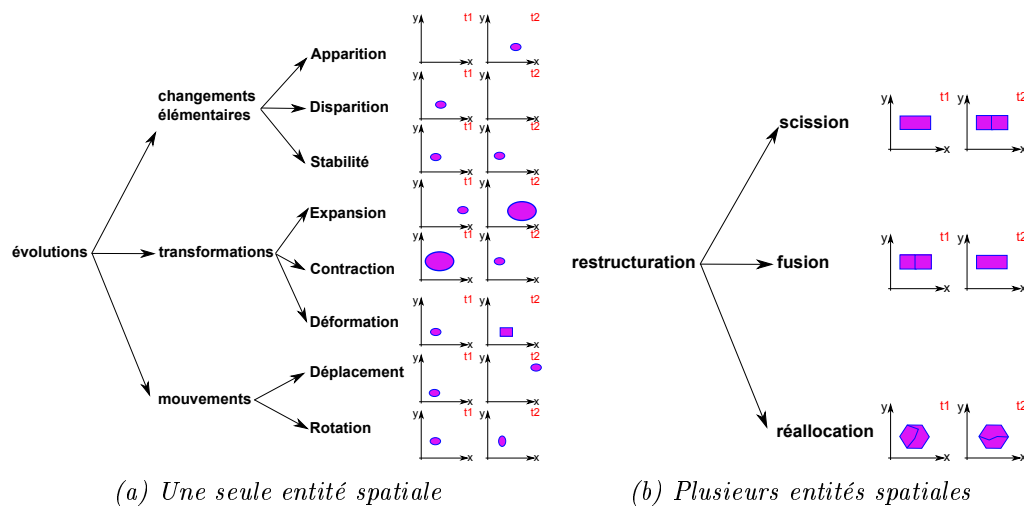


FIGURE I.4. Typologies des évolutions spatio-temporelles des données géographiques, tirées de [Claramunt et Thériault \(1996\)](#) et [Costes \(2016\)](#)

b) Typologie des formes de carto-vandalisme

Dresser une classification du carto-vandalisme à partir de sa cause formelle (*i.e.* les différentes formes prises par le carto-vandalisme) implique de considérer le phénomène comme une dégradation de l'espace cartographique collaboratif. Or, une dégradation est une transformation de l'espace cartographique à partir de la modification des données géographiques qui le compose. Sur ce sujet, des typologies

listant les formes de modification des données géographiques ont été proposées dans les travaux de [Duménieu \(2015\)](#); [Costes \(2016\)](#) et [Claramunt et Thériault \(1996\)](#) pour modéliser les changements spatio-temporels des entités spatiales. Ces typologies se subdivisent principalement en deux catégories : les évolutions désignant les modifications sur une seule entité spatiale et les restructurations désignant les modifications sur plusieurs entités spatiales. La Figure 1.4 schématise ces typologies. Les données géographiques étant constituées d'information attributaire, ces typologies sont donc valables pour décrire les modifications géométriques et sémantiques des entités spatiales.

2.5 Éléments de qualité de l'information géographique

Pour vérifier qu'une modification de données géographiques dégrade l'espace cartographique ou non, il est nécessaire d'étudier sa qualité. La qualité de l'information géographique peut être évaluée selon différents critères, aussi appelés des éléments de qualité. Différents éléments de qualité de l'information géographique ont été formulés dans la littérature, mais jusqu'ici aucun consensus n'a été trouvé pour évaluer les données spatiales selon des éléments de qualité précis ([Girres, 2012](#)). Il existe toutefois une norme ISO, la dernière version étant [ISO19157 \(2003\)](#), pour l'évaluation de la qualité de l'information géographique. Celle-ci spécifie six éléments de qualité :

- (i) l'exactitude de position ;
- (ii) l'exactitude thématique ;
- (iii) l'exactitude temporelle ;
- (iv) la complétude ;
- (v) la cohérence logique ;
- (vi) l'utilisabilité.

Les cinq premiers éléments évaluent l'information selon un point de vue producteur, c'est-à-dire selon la concordance des données acquises avec les spécifications du producteur. L'utilisabilité, quant à elle, permet d'évaluer l'information selon un point de vue utilisateur, c'est-à-dire l'adéquation des données produites avec les besoins des consommateurs de données géographiques.

Les éléments de qualité de l'information géographique sont des critères qui permettent de déterminer si une donnée dégrade l'espace cartographique. Nous détaillons ci-après chacun de ces éléments en l'illustrant par des incidents réels provenant du projet OpenStreetMap (OSM). Certains incidents sont relevés à partir de la liste des contributeurs bannis par le *Data Working Group* (DWG), qui est un groupe de contributeurs dont le rôle est de bannir temporairement les comptes utilisateurs OSM présentant des activités douteuses, notamment celles qui produisent du carto-vandalisme. Cette liste est consultable sur une page web dédiée¹¹.

En explorant la page des contributeurs bannis d'OSM¹², nous cherchons à comprendre les raisons qui ont poussé le DWG à bannir certains contributeurs de la

11. www.openstreetmap.org/user_blocks/

12. Voir Annexe A

plateforme et à qualifier les incidents qui relèvent du carto-vandalisme, selon notre définition. Cette exploration a été effectuée manuellement, en parcourant les messages émis par les modérateurs. Ces messages peuvent contenir l'identifiant de la contribution ou du groupe de modifications (*changeset*¹³) qui a mené au bannissement du contributeur, mais cette information n'est pas systématiquement renseignée. L'historique des données a été obtenu avec l'application OSM Deep History¹⁴. Cette phase exploratoire aurait pu se faire de manière automatique en recourant à des techniques de traitement automatique de texte, pour analyser les messages entre les modérateurs et les contributeurs. Cela aurait permis de récupérer une liste plus exhaustive d'actes de carto-vandalisme. Toutefois, cette exploration manuelle a permis d'identifier rapidement quelques cas intéressants et de donner un aperçu des différentes raisons qui peuvent pousser les modérateurs à bannir les contributeurs OSM.

Pour lire les Tables I.1 à I.6 : Ces tables donnent un extrait de l'historique des tags de quelques objets OSM. Un tag est un couple de clé/valeur qui permet d'apporter une description attributaire à un élément cartographique dans OSM. Dans une table, chaque colonne décrit les valeurs de tags pour une version donnée de l'objet. Les cellules colorées indiquent un changement par rapport à la version précédente : une **cellule verte** signifie l'ajout d'un nouveau tag, une **cellule jaune** signifie la modification de la valeur d'un tag, une **cellule rouge** correspond à la suppression d'un tag. Les versions dont le numéro est souligné correspondent aux contributions produites par les contributeurs bannis.



FIGURE I.5. Hôpital dont la géométrie chevauche plusieurs îles du Nord du Canada.

13. Dans OpenStreetMap, un *changeset* contient les identifiants de tout objet créé, modifié ou supprimé durant une session d'édition.

14. <https://osmlab.github.io/osm-deep-history/>

(i) **Exactitude de position :** L'exactitude de position de l'information géographique est un élément de qualité de la norme ISO qui a été particulièrement étudié pour qualifier les données collaboratives. Elle est définie comme la justesse du positionnement spatial des données par rapport à un système de référence spatial. La Figure I.5 illustre un exemple de dégradation de la carte OSM où un hôpital a été cartographié sur plusieurs centaines de kilomètres. Notons que l'objet a été saisi trois heures après la date de création de compte OSM de son auteur. De plus, le commentaire laissé par le contributeur sur cette session d'édition et ses contributions suivantes indiquent que celui-ci cherchait initialement à cartographier les éléments manquants aux alentours d'une université au Sri Lanka¹⁵. Ici, il s'agit probablement d'une erreur de débutant dans la manipulation du système de coordonnées. Quoiqu'il en soit, l'intentionnalité n'étant pas démontrée, cet incident ne correspond pas à un cas réel de carto-vandalisme.

(ii) **Exactitude thématique :** Elle est définie comme la justesse des attributs – quantitatifs et qualitatifs – et la justesse de classification des objets géographiques. La Figure I.6 et la Table I.1 illustrent le cas de l'île Hérald, transformée en parc (cf. le tag `leisure` à la version n°6) et renommée en « Jerryland » (cf. les tags `name` et `name:en` sur les versions n°6 et n°7). Ici, les modifications attributaires ont altéré l'exactitude thématique de l'objet géographique, qui ne représente plus la réalité.



FIGURE I.6. L'île Hérald est située entre la Russie et le Canada. Source : OpenStreetMap

TABLE I.1. Extrait de l'historique des tags de l'île Hérald de la Figure I.6.

Version Tag key	n° 5	n° 6	n° 7	n° 8	n° 9
<code>leisure</code>		<code>park</code>	<code>park</code>		
<code>name</code>	остров Геральд	Jerryland	Jerryland	Jerryland	остров Геральд
<code>name:en</code>	Herald Island	Herald Island	Jerryland	Herald Island	Herald Island
<code>place</code>	island		island	island	island
<code>wikipedia</code>	ru:Геральд (остров)	ru:Геральд (остров)			ru:Геральд (остров)

La Figure I.7 et la Table I.2 montrent un exemple de désaccord entre les contributeurs d'OSM sur la nature d'une forêt située à Riga, en Lettonie. En particulier,

15. <https://www.openstreetmap.org/changeset/53283079>

le débat a conduit au bannissement du contributeur qui insistait pour décrire le lieu avec le tag `leisure=park`. Or, en observant l'historique des éditions (Table I.2), nous pouvons remarquer la coexistence des tags `name` et `alt_name` à partir de la version n° 13, indiquant le fait que ce lieu pouvait être connu sous deux noms (Anniņmuižas mežs et Anniņmuižas parks). En menant des recherches sur ce lieu, il s'avère que ce lieu est un bois considéré comme un parc par les habitants de Riga¹⁶. Par conséquent, plutôt qu'une dégradation de l'espace cartographique, cet incident relève d'une ambiguïté sur la classification d'un objet géographique (Ali *et al.*, 2014; Glasze et Perkins, 2015; Mcnair et Arnold, 2016). Dans ce cas précis, le contributeur a été banni pour avoir participé à une guerre d'édition, et non pour avoir réellement dégradé l'espace cartographique. En effet, nous re-précisons que les guerres d'éditions ne sont pas autorisées sur OSM¹⁷, celles-ci sont même considérées comme étant du vandalisme, d'après le règlement du projet.

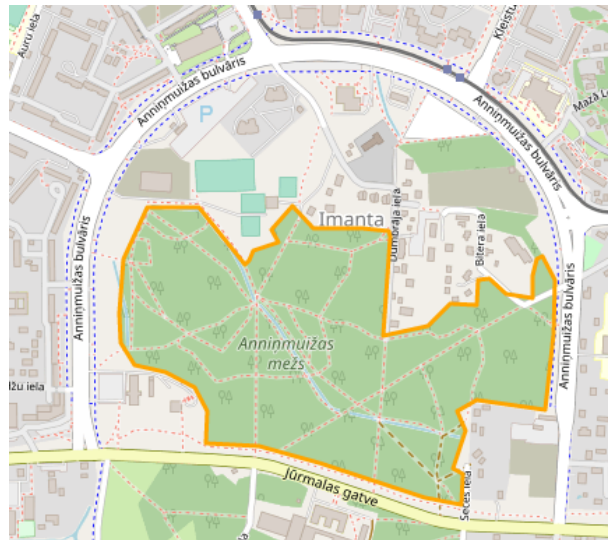


FIGURE I.7. Anniņmuižas mežs / Anniņmuižas parks. Source : OpenStreetMap.

TABLE I.2. Extrait de l'historique des tags de l'objet de la Figure I.7.

Version	n° 7	n° 8	n° 9	n° 10
Tag key				
alt_name				
landuse	forest		forest	
leisure		park	park	park
name	Anniņmuižas parks	Anniņmuižas parks	Anniņmuižas parks	Anniņmuižas parks
Version	n° 11	n° 13	n° 14	n° 15
Tag key				
alt_name		Anniņmuižas parks	Anniņmuižas parks	Anniņmuižas parks
landuse	forest	forest	forest	forest
leisure		park		park
name	Anniņmuižas mežs	Anniņmuižas mežs	Anniņmuižas mežs	Anniņmuižas mežs

N.B. : La version n° 12 n'apparaît pas car elle correspond à une édition purement géométrique, par conséquent, les tags à cette version sont identiques à ceux de la version n° 11.

La Figure I.8 est un incident détecté par le DWG comme étant un carto-vandalisme issu d'un détournement pour Pokémon Go, où le contributeur a tenté à plusieurs reprises de transformer une zone militaire en zone industrielle (Table I.3). L'espace

16. <https://www.spottedbylocals.com/riga/anninmuizas-mezs/>

17. <https://wiki.openstreetmap.org/wiki/Disputes>

cartographique est donc dégradé par ces éditions puisqu'il ne représente plus la réalité de manière fidèle. Cependant, la dégradation ici n'a probablement pas été produite par ennui – ou du moins pas uniquement – mais pour servir à des fins totalement personnelles¹⁸. En ce sens, ce type de carto-vandalisme peut se rapprocher du vandalisme cupide de Cohen (1973).



FIGURE I.8. Delaware National Guard Bethany Beach Training Site

TABLE I.3. Extrait de l'historique des tags de Delaware National Guard Bethany Beach Training Site de la Figure I.8.

Version Tag key	n° 1	n° 2	n° 3	n° 4	n° 5
landuse	military	industrial	military	industrial	military
name	Delaware National Guard Bethany Beach Training Site	Delaware National Guard Bethany Beach Training Site	Delaware National Guard Bethany Beach Training Site	Delaware National Guard Bethany Beach Training Site	Delaware National Guard Bethany Beach Training Site
type	multipolygon	multipolygon	multipolygon	multipolygon	multipolygon

Table I.4 dépeint un désaccord entre les contributeurs sur la toponymie de la ville la ville de Malishevë/Mališevo. Le contributeur de la version n° 32 de l'objet a été banni pour avoir retiré l'orthographe serbe du tag `[name]`. Après avoir été banni, la version n° 32 a été annulée (*i.e.* la version n° 33 correspond à la version n° 31). Vu les multiples versions de l'objet – plus de 30 – celui-ci semble être un sujet controversé pour les contributeurs OSM. D'après l'Organisation pour la Sécurité et la Coopération en Europe¹⁹, la ville de Malishevë/Mališevo est habitée majoritairement par des albanais. Cet élément de contexte expliquerait pourquoi le contributeur de la version n° 32 a conservé l'écriture albanaise de la ville sur le tag `[name]` principal. Notons que, même si sur la version n° 32, l'écriture serbe ne figure pas sur le tag `[name]`, celui-ci reste indiqué sur le tag `[name:sr-Latn]`. Dans cette situation, le contexte politique semble expliquer la dégradation des données. D'une certaine manière, cette édition n'a pas vraiment dégradé la qualité de l'espace cartographique,

18. Dans le cas de Pokémon Go, il s'agit de pouvoir attraper un maximum de créatures pour gagner des niveaux dans le jeu vidéo.

19. <https://www.osce.org/kosovo/13137?download=true>

si ce n'est qu'elle structure l'information attributaire de manière différente. Or, la structuration de l'information sémantique est variable parce que la spécification des données dans OSM n'a pas été définie une fois pour toute. Elle peut évoluer dans le temps (Antoniou et Skopeliti, 2017). Par conséquent, cet incident ne correspond pas à un acte de carto-vandalisme, car la dégradation de l'espace cartographique est ici équivoque. Le contributeur de la version n° 32 a donc été banni pour non-respect des règles de bonne conduite du projet.

TABLE I.4. Extrait de l'historique des tags (versions n° 28-n° 33) d'une ville du Kosovo (ID node : 1920392800) dont la toponymie est sujette à controverse, l'albanais et le serbe étant les deux langues officielles.

Version Tag key	n° 28	n° 29	n° 30	n° 31	n° 32	n° 33
name	Malishevë - Mališevo	Malishevë	Malishevë	Malishevë/ Mališevo	Malishevë	Malishevë/ Mališevo
name:sq	Malishevë	Malishevë	Malishevë	Malishevë	Malishevë	Malishevë
name:sr-Latn		Mališevo	Mališevo	Mališevo	Mališevo	Mališevo

(iii) **Qualité temporelle** : Elle est définie comme la qualité des attributs temporels et les relations temporelles entre les éléments. Elle inclut : l'exactitude temporelle entre la date de saisie, c'est-à-dire la différence entre la date de saisie d'une donnée par rapport à sa référence temporelle ; la cohérence temporelle, c'est-à-dire la justesse dans l'ordre chronologique des événements reportés ; la validité temporelle, c'est-à-dire l'actualité temporelle des données saisies.

(iv) **Complétude** : Elle correspond à la présence et l'absence de données, de leurs attributs et de leurs relations. Dans l'historique OSM de l'île Hérald, le lien vers la page Wikipédia de ce lieu, renseigné par le tag `wikipedia`, a été supprimé à la version n°6 (Table I.1). Cette suppression d'information nuit donc à la complétude de l'objet géographique. Par ailleurs, les dégradations perpétrées par le même contributeur sur cet objet géographique touchent plusieurs éléments de qualité de l'information géographique, et portent à croire que le contributeur a agi délibérément. Par conséquent, les versions n°6 à n°8 de l'île Hérald correspondent à des actes de carto-vandalisme.

(v) **Cohérence logique** : Elle est définie comme le degré d'adhérence à des règles logiques de structuration des données, d'attribution des données et de relations entre les données. La Figure I.9 et la Table I.5 représentent une zone commerciale transformée en lac sur sa version n°2. Or, la présence d'un lac sur une zone où se trouvent déjà un parking et des commerces paraît assez invraisemblable : cette contribution dégrade donc l'espace cartographique, d'une part car la carte ne décrit plus la réalité ; d'autre part car cette contribution affecte la cohérence logique entre les différents objets cartographiques. De plus, le fait de nommer le lac par un smiley (tag `name = :)`) suggère le caractère intentionnel du contributeur. En observant les autres contributions de ce compte utilisateur, ce dernier a ajouté, de la même manière, d'autres lacs en zone urbaine. À notre connaissance, aucun élément de

contexte n'a permis de justifier la saisie d'un tel élément. Par conséquent, cette contribution est un acte de carto-vandalisme, potentiellement de type ludique.

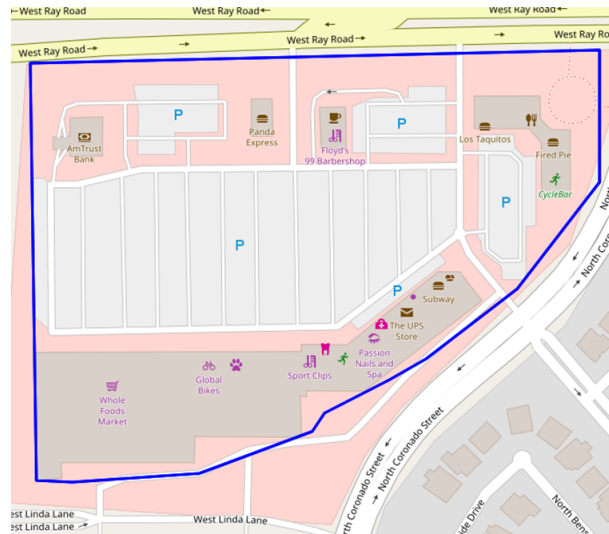


FIGURE I.9. Exemple de carto-vandalisme OSM sur une zone commerciale (transformée en lac par le contributeur banni)

TABLE I.5. Historique des tags de l'objet de la Figure I.9.

Version	n°1	n° 2
Tag key		
landuse	retail	
name		:)
natural		water
water		lake

(vi) **Utilisabilité** : L'utilisabilité (*fitness for use*) de l'information va dépendre des besoins de l'utilisateur. L'évaluation de l'utilisabilité dépendra donc du cas d'application considéré, qui déterminera des typologies différentes de qualification de l'information. Nous développons plus en détail cet élément de qualité dans la Section 3.3.

La Figure I.10 et la Table I.6 correspondent à une contribution dont l'auteur a été banni pour avoir importé massivement des objets décrits par des tags peu conventionnels. Par exemple, pour les tags de la première version dans la Table I.6, `description1` et `description2` sont des clés de tag qui ne sont pas communément utilisées par les autres contributeurs OSM, et donc par les algorithmes qui traitent les données OSM pour leur réutilisation dans d'autres applications. De même, il a été reproché à ce contributeur d'utiliser des majuscules sur les clés `Is_in` et `Ref:FR:SIREN` (au lieu de `is_in` et `ref:FR:SIRET` respectivement). Ce manque de rigueur dans le renseignement des tags rend l'information inutilisable, par exemple pour des applications de navigation qui se basent sur les cartes OSM mais qui ne gèrent pas la sensibilité à la casse sur les clés de tag : l'adresse renseignée avec cette clé de tag ne serait pas visible. Toutefois, les versions ultérieures de cet objet conservent les informations indiquées par le contributeur banni : la contribution de ce dernier a été utile, même si elle nécessitait d'être améliorée. Le bannissement

de ce contributeur n'avait donc pas pour but d'empêcher une réelle dégradation de l'espace cartographique, mais elle peut se comprendre comme un avertissement de la part des modérateurs envers ce contributeur, pour l'inciter à être plus rigoureux dans ses contributions.

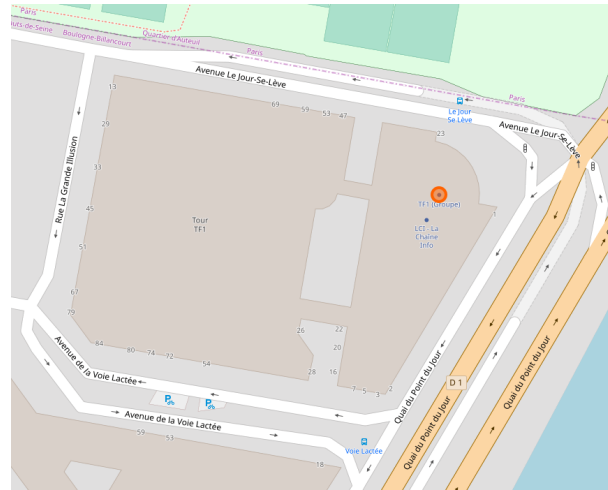


FIGURE I.10. Le point orange localise le siège du groupe industriel TF1.

TABLE I.6. Extrait de l'historique des tags de l'objet de la Figure I.10.

Version Tag key	n° 1	n° 2	n° 3	n° 4
is_in	1 QUAÏ DU POINT DU JOUR			
Ref:FR:SIREN	32630015900067			
description1	COMMUNICATION - PRESSE - EDITION	COMMUNICATION - PRESSE - EDITION	COMMUNICATION - PRESSE - EDITION	COMMUNICATION - PRESSE - EDITION
description2	Audiovisuel (diffusion - production)	Audiovisuel (diffusion - production)	Audiovisuel (diffusion - production)	Audiovisuel (diffusion - production)
is_in		1 QUAÏ DU POINT DU JOUR	1 QUAÏ DU POINT DU JOUR	1 QUAÏ DU POINT DU JOUR
name	TF1 (GROUPE)	TF1 (GROUPE)	TF1 (GROUPE)	TF1 (Groupe)
office	compagny	compagny	compagny	compagny
phone	01 41 41 12 34	01 41 41 12 34	+33 1 41 41 12 34	+33 1 41 41 12 34
ref:FR:SIRET		32630015900067	32630015900067	

Bilan de l'exploration des contributeurs bannis d'OSM

L'exploration des contributeurs bannis d'OSM montre que le bannissement ne concerne pas seulement les contributeurs ayant commis des actes de vandalisme sur les données OSM, mais également ceux qui ne respectent pas les règles de bonne conduite sur le projet : ne pas répondre aux commentaires adressés par d'autres contributeurs, importer massivement des données présentant une licence incompatible, ou participer à des guerres d'éditations sont des formes de vandalisme pour le projet OSM qui n'entrent pas dans le cadre de notre définition du carto-vandalisme. Par conséquent, le carto-vandalisme, tel que nous l'avons défini, couvre des cas très marginaux des incidents répertoriés sur la page des contributeurs bannis d'OSM.

3 État de l'art sur la qualité de l'information géographique volontaire

Les multiples incidents relevés sur le projet OSM dans la section précédente montrent que le carto-vandalisme se manifeste par une dégradation de la qualité de l'information géographique volontaire. La détection du carto-vandalisme implique donc d'évaluer la qualité de l'information géographique volontaire. Dans cette partie, nous passons en revue les contributions scientifiques sur la qualité des données géographiques collaboratives en nous appuyant sur les derniers articles d'état de l'art sur ce sujet (Fonte *et al.*, 2017; Senaratne *et al.*, 2017).

3.1 Approches d'assurance qualité et méthodes de qualification

Goodchild et Li (2012) ont formulé trois approches d'assurance qualité de l'information géographique volontaire. Ces approches sont à comprendre comme des mécanismes qui visent à tirer avantage des caractéristiques propres à la démarche collaborative cartographique afin d'assurer et améliorer la qualité des données acquises par ce moyen. Chaque approche repose sur une loi ou une hypothèse qui permet d'assurer la qualité de l'information géographique.

(i) **L'approche participative (*crowdsourcing*)** : Cette approche porte sur l'aspect participatif de l'information géographique volontaire. Elle reprend **la loi de Linus**²⁰ qui sous-entend que plus il y a de contributeurs, plus il y a de chances que les erreurs de saisie soient corrigées. Ainsi, les données peuvent potentiellement atteindre un bon niveau de qualité, à condition qu'il y ait suffisamment de contributeurs actifs.

(ii) **L'approche sociale** : Elle est basée sur l'hypothèse selon laquelle la qualité des données collaboratives « repose sur une hiérarchie d'individus de confiance agissant en tant que modérateurs ». Cette approche est complémentaire à l'approche *crowdsourcing*, puisque la participation des contributeurs de confiance va favoriser la correction d'erreurs et l'amélioration de la qualité des données.

(iii) **L'approche géographique** : Cette approche porte sur la dimension spatiale des données collaboratives saisies. Elle s'appuie sur la **première loi de la géographie** (Tobler, 1970), selon laquelle « tout interagit avec tout mais deux choses voisines ont plus de chances d'entrer en interaction que deux choses lointaines ». Cette loi signifie qu'il existe une relation de dépendance entre un objet géographique et son voisinage. Cette relation est généralement soumise à des règles de cohérence, qui peuvent se situer à différents niveaux :

20. « *given enough eyeballs, all bugs are shallow* »

- Cohérence géométrique : des objets géographiques qui sont à proximité les uns des autres ont des chances de présenter des similitudes de taille ou de forme.
- Cohérence topologique : des objets géographiques qui sont à proximité les uns des autres suivent généralement des règles topologiques précises. Par exemple, il est impossible que deux bâtiments se chevauchent, alors qu'un bâtiment peut se trouver à l'intérieur d'une forêt.
- Cohérence thématique : la nature des éléments géographiques suit certaines règles de vraisemblance. En reprenant la zone cartographiée sur la Figure I.9, il n'est pas étonnant de voir des commerces cartographiés à proximité d'un parking. En revanche, il sera moins probable de trouver une prairie dans la même zone.

Toutefois, la réalité du monde est complexe et imprévisible, par conséquent ces règles de cohérence ne sont pas absolues, même si elles peuvent guider la qualification des contributions, à partir de l'évaluation de la plausibilité de certaines situations (Ali *et al.*, 2014).

À ces trois approches d'assurance qualité, Senaratne *et al.* (2017) y ajoute l'approche fouille de données (*data mining*). Celle-ci est indépendante des trois mécanismes d'assurance qualité, c'est-à-dire qu'elle ne repose sur aucune loi ni aucune hypothèse permettant d'assurer la bonne qualité des données. En revanche, elle regroupe les méthodes d'apprentissage automatique (*machine learning*) permettant de qualifier l'information géographique volontaire. Or, il est possible de qualifier les données selon une approche d'assurance qualité (*crowdsourcing*, sociale ou géographique) en ayant recours à une technique d'apprentissage automatique. Par conséquent, l'approche *data mining* n'est pas une autre approche d'assurance qualité, mais plutôt une catégorie particulière de méthodes de qualification des données, qui peuvent très bien répondre à l'un des trois mécanismes présentés précédemment.

Par ailleurs, les méthodes de qualification de l'information géographique volontaire peuvent également se classer dans deux catégories. D'une part, les méthodes dites extrinsèques consistent à qualifier les données collaboratives en les comparant à des données de référence (Haklay, 2010; Hung *et al.*, 2016; Zhang et Malczewski, 2018). D'autre part, les méthodes intrinsèques évaluent la qualité des données collaboratives en se basant uniquement sur ces données (Barron *et al.*, 2014; Degrossi *et al.*, 2018). Les méthodes intrinsèques sont particulièrement intéressantes dans le cas où aucune donnée de référence n'est accessible. Il ne s'agit pas de chercher à choisir entre ces deux méthodes, car chacune d'elles présente des avantages et des inconvénients. En revanche, ces méthodes étant complémentaires, il est intéressant de combiner les méthodes extrinsèques et intrinsèques pour qualifier l'information géographique volontaire (Touya *et al.*, 2017b).

3.2 Mesures, indicateurs et éléments de qualité de l'information géographique volontaire

Généralement, les méthodes de qualification de l'information géographique volontaire consistent à calculer des mesures permettant d'estimer la qualité des données. Antoniou et Skopeliti (2015) font une distinction entre les mesures et les indicateurs

de qualité. Les mesures de qualité désignent les métriques issues de la comparaison des données collaboratives avec des données d'autorité. Les indicateurs de qualité correspondent aux métriques calculées à partir des seules données collaboratives. Par conséquent, nous parlerons de mesures de qualité pour les méthodes extrinsèques, et d'indicateurs de qualité pour les méthodes intrinsèques.

a) Mesures de qualité

D'après [Fonte *et al.* \(2017\)](#), les éléments de qualité de l'information géographique traditionnelle peuvent être repris pour évaluer l'information géographique volontaire. De nombreuses mesures ont été proposées pour évaluer les données géographiques collaboratives selon ces éléments de qualité. En effet, des méthodes extrinsèques d'appariement avec des données d'autorité ont été développées pour mesurer l'exactitude de position des données OSM, en particulier les bâtiments ([Fan *et al.*, 2014b](#); [Brovelli et Zamboni, 2018](#)) et les routes ([Brovelli *et al.*, 2017](#); [Ivanovic, 2018](#); [Castro *et al.*, 2019](#)). De plus, la comparaison thématique des données collaboratives avec des données d'autorité a été faite principalement pour évaluer la justesse de description des données dans le cadre l'utilisation et l'occupation des sols ([Dorn *et al.*, 2015](#); [Jokar Arsanjani *et al.*, 2015](#)).

En ce qui concerne la temporalité des données collaboratives, [Fonte *et al.* \(2017\)](#) soulève leur hétérogénéité temporelle entre leur date de saisie et leur validité. Par ailleurs, l'hypothèse selon laquelle les données collaboratives peuvent être en avance sur les données d'autorité a été admise, notamment dans les travaux de recherche qui traitent de la mise à jour des données géographiques institutionnelles à partir de l'utilisation de traces GPS participatives ([Ivanovic *et al.*, 2020](#)). Une étude statistique sur la capacité des données OSM à refléter les changements du monde réel a également été menée en considérant les ouvertures et fermetures de cafés à Manhattan ([Zhang et Pfoser, 2019](#)).

L'évaluation de la complétude des données géographiques collaboratives a souvent été faite de manière conjointe à celle de l'exactitude de position ([Fan *et al.*, 2014b](#); [Brovelli et Zamboni, 2018](#)). En effet, la complétude est généralement étudiée à partir de méthodes extrinsèques ([Haklay, 2010](#); [Hecht *et al.*, 2013](#); [Zielstra *et al.*, 2013](#); [Mahabir *et al.*, 2017](#)), grâce à des techniques d'appariement de données ([Chehreghan et Ali Abbaspour, 2018](#)). Il s'agit d'évaluer si l'information géographique volontaire est au moins aussi complète que les données d'autorité. Enfin, l'appariement des données OSM en Chine avec un jeu de référence a permis d'évaluer non seulement l'exactitude de position, la complétude mais aussi la cohérence logique de celles-ci ([Yang *et al.*, 2018](#)).

b) Indicateurs de qualité

L'évaluation des éléments de qualité de l'information géographique pour les données collaboratives ne signifie pas que celle-ci se réalise obligatoirement par des méthodes extrinsèques, impliquant une comparaison avec des données géographiques de référence. En ce qui concerne l'exactitude thématique des données collaboratives, des méthodes intrinsèques ont été développées pour évaluer leur sémantique à par-

tir de calculs de similarité sémantique (Mülligann *et al.*, 2011; Majic *et al.*, 2017). De plus, il est également possible d'évaluer la complétude par des méthodes intrinsèques, comme l'étude du développement du réseau routier dans OSM jusqu'à sa saturation (Barrington-Leigh et Millard-Ball, 2017).

Notons que l'évaluation de la cohérence logique des données géographiques collaboratives a été principalement étudiée à partir des méthodes intrinsèques. En particulier, Touya et Brando-Escobar (2013) ont abordé le problème des incohérences de niveaux de détail des données collaboratives qui entraînent des incohérences spatiales sur les relations topologiques entre objets géographiques. Pour étudier le niveau de détail des objets OSM, des indicateurs tels que la résolution et la granularité ont été proposés. Ali *et al.* (2014) se sont concentrés sur la cohérence thématique des données saisies, en traitant la question de la plausibilité et de vraisemblance thématique des objets géographiques.

Un objet cartographique collaboratif peut ne pas être de bonne qualité dès sa première version. Cependant, l'approche *crowdsourcing* assure que sa qualité peut être améliorée au moyen des corrections apportées par les contributeurs qui éditent cet objet. L'étude de l'historique des contributions – c'est-à-dire les différentes versions qui ont été produites sur chaque objet – permet alors de déterminer si un consensus a été atteint sur la saisie d'une donnée (Mooney et Corcoran, 2012a). L'accord des contributeurs sur la représentation cartographique d'un objet géographique peut également être évalué par un score de confirmation (Kessler et de Groot, 2013; Fogliaroni *et al.*, 2018). La stabilité des éditions s'apparente également à la notion de maturité des données (Gröchenig *et al.*, 2014; Maguire et Tomko, 2017).

Selon l'approche sociale, la présence de contributeurs agissant comme des modérateurs assure la qualité de l'information géographique volontaire. Cela implique d'étudier le comportement des contributeurs de manière à pouvoir identifier les profils de ceux qui pourraient être assimilés à celui de modérateur. Les études faites sur les contributeurs d'information géographique volontaire permettent d'évaluer, entre autres, la réputation du contributeur (Fogliaroni *et al.*, 2018), ses motivations (Budhathoki *et al.*, 2010) ou encore son expertise (Yang *et al.*, 2016).

3.3 Typologies d'évaluation

Les éléments de qualité (développés dans la Section 3.2 précédente) permettent d'évaluer la qualité interne de l'information géographique volontaire, c'est-à-dire la qualité de production des données. La question est maintenant de vérifier si ces éléments de qualité permettent à un utilisateur de données géographiques de connaître l'utilisabilité de l'information géographique volontaire par rapport à son cas d'utilisation : si oui, y a-t-il des éléments de qualité à préférer selon le type d'utilisation ? Sinon, quelle(s) typologie(s) faut-il adopter pour qualifier l'information ?

D'après Servigne *et al.* (2005), l'évaluation de l'adéquation des données aux besoins nécessite d'évaluer l'exhaustivité (ou la complétude) des données et l'exhaustivité du modèle. L'exhaustivité des données est mesurable (voir Section 3.2 sur la complétude) et indépendante de l'application, alors que l'exhaustivité du modèle évalue l'adéquation aux besoins de l'application. En effet, celle-ci compare le modèle

de données à qualifier en termes d'abstraction du monde réel avec celui correspondant à l'application considérée. [Ballatore et Zipf \(2015\)](#) parlent de "golfe sémantique" pour désigner les différentes abstractions du monde réel entre producteurs et consommateurs de données géographiques. Pour faciliter la communication entre ces deux catégories d'acteurs, ces auteurs développent la notion de qualité conceptuelle, qui comprend six dimensions conceptuelles – exactitude, granularité, complétude, cohérence, conformité et richesse – permettant de qualifier l'utilisabilité des données géographiques volontaires.

Plusieurs méthodes de qualification des données selon le point de vue de l'utilisabilité des données ont été développées. L'outil *iOSMAnalyzer* développé par [Barron et al. \(2014\)](#) permet de qualifier les données OSM de manière intrinsèque en calculant des indicateurs différents selon le cas d'utilisation considéré. Ces indicateurs de qualité peuvent reprendre les éléments de qualité interne, tels que la complétude du réseau routier dans le cas de la navigation, la complétude attributaire dans le cadre de recherche de points d'intérêt, la cohérence logique pour les applications de visualisation cartographique. [Jonietz et Zipf \(2016\)](#) proposent une méthode générique d'évaluation de l'utilisabilité des points d'intérêt (POI) sur OSM. Il s'agit de déterminer les fonctions basiques des POI : en l'occurrence, les POI ont pour fonctions principales le géo-référencement (localisation par coordonnées géographiques) et l'objet-référencement (localisation par noms). Les cas d'utilisation des POI vont privilégier l'une ou l'autre fonction. Ainsi, l'évaluation de l'utilisabilité revient à évaluer la fonction basique considérée, en choisissant des mesures d'exactitude de position et thématique adaptées.

Cependant, il existe des cas d'utilisation qui nécessitent de qualifier l'information géographique volontaire selon d'autres critères. Dans un contexte de gestion de crises, telles que les inondations, l'information géographique volontaire peut être évaluée par son niveau de crédibilité ([Hung et al., 2016](#)). En revanche, pour produire automatiquement des cartes topographiques à partir de données collaboratives, il s'agit plutôt de qualifier les données selon leur niveau de détail, notamment dans le but de gérer les problématiques de généralisation cartographique ([Touya et Brando-Escobar, 2013](#); [Touya et al., 2017a](#)). Enfin, pour étudier le vandalisme dans les données OSM, [Neis et al. \(2012\)](#) s'appuient entre autres sur la réputation du contributeur pour qualifier l'information. Ces différents éléments de qualité peuvent tenir compte de certaines mesures qui ont été développées dans le cadre des éléments de qualité interne de l'information géographique. L'analyse des besoins permet de choisir les métriques à utiliser pour qualifier l'information selon la typologie adaptée au cas d'application.

Les méthodes et métriques issues de la littérature pour qualifier l'information géographique volontaire sont généralement développées pour un cas d'application précis ou pour évaluer un élément de qualité précis. Nous constatons ici un manque d'approche globale de la qualité des données dans la littérature, c'est-à-dire qu'il n'existe pas – à notre connaissance – de contributions sur une méthodologie générique de qualification de l'information géographique volontaire. Par conséquent, le travail de qualification du carto-vandalisme implique de développer une approche spécifique à ce cas d'application.

Conclusion du chapitre

Ce premier chapitre a permis de définir le carto-vandalisme à partir de l'analyse des différentes définitions de la notion du vandalisme, depuis ses origines historiques à sa définition législative actuelle. Le carto-vandalisme correspond donc à une dégradation intentionnelle de l'espace cartographique collaboratif, sans motif légitime. À partir de cette définition, nous avons pu qualifier les incidents menant au bannissement de certains contributeurs de la plateforme OpenStreetMap. Cependant, nous avons constaté que, parmi les contributeurs bannis d'OSM, très peu ont réellement provoqué du carto-vandalisme. Cela s'explique par le fait que la plateforme ne cherche pas uniquement à détecter les carto-vandales, mais à alerter – au moyen du bannissement – tout contributeur qui ne respecte pas les règles de bonne conduite ou qui produit des erreurs à grande échelle.

Dans ce chapitre, nous avons dressé un état de l'art, présentant les nombreux travaux qui traitent de la qualité de l'information géographique volontaire. En particulier, les différentes approches d'assurance qualité ont été théorisées à partir de lois ou d'hypothèses qui exploitent les caractéristiques propres à la démarche collaborative cartographique, à savoir la participation des contributeurs (approche participative), l'existence de contributeurs de confiance agissant comme des modérateurs (approche sociale) et l'existence de règles spatiales (approche géographique). De nombreuses métriques et méthodes ont été développées dans le but d'évaluer l'information géographique volontaire selon différents éléments de qualité, du point de vue des producteurs de données et des utilisateurs. Sur ce dernier point, la question de l'utilisabilité des données a été abordée, car elle oriente le choix des métriques ainsi que la typologie selon laquelle qualifier l'information géographique volontaire.

Nous avons vu dans ce chapitre que le carto-vandalisme est un phénomène qui nécessite de qualifier non seulement les données saisies, pour détecter celles qui dégradent l'espace cartographique collaboratif, mais également les contributeurs, afin d'identifier ceux qui participent consciemment à dégrader la base de données géographiques ouverte. C'est pourquoi le chapitre suivant se concentre sur l'approche sociale : en particulier, il s'agit de qualifier le contributeur dans le but de qualifier les contributions qui relèveraient du carto-vandalisme.

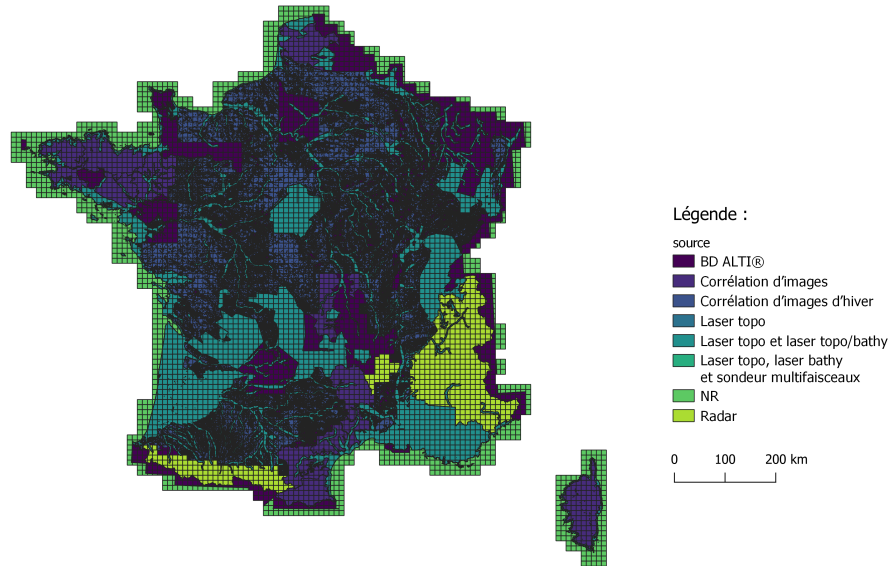
Chapitre II

Du vandale cartographique au carto-vandalisme

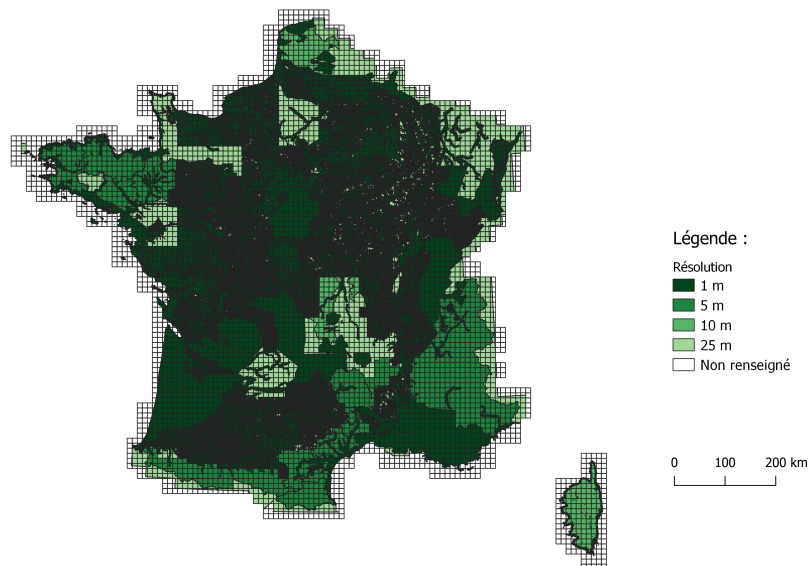
La qualification du carto-vandalisme ne passe pas seulement par l'évaluation de la qualité de l'information géographique volontaire, mais elle nécessite également d'estimer l'intentionnalité du contributeur ainsi que le contexte de saisie. Ce chapitre se concentre sur la qualification du contributeur, en particulier son comportement de collaboration. En revenant sur l'approche participative et l'approche sociale, nous montrons l'importance de qualifier le contributeur pour qualifier les données collaboratives. Puis, nous présentons une méthode de qualification du contributeur à partir d'un modèle de réseau social multiplexe. Ce modèle permet de représenter son comportement interactif sur la plateforme collaborative, et de mettre en valeur des profils comportementaux particuliers. En implémentant ce modèle sur des données OSM, les résultats de plusieurs cas d'étude permettent de valider notre méthodologie. Enfin, en exploitant les métriques de qualité du contributeur issues du modèle multiplexe, nous explorons de nouvelles méthodes d'évaluation du contributeur qui tiennent mieux compte des différents facteurs d'influence de sa fiabilité.

1 De la qualité du contributeur à la qualité des contributions

1.1 Le contributeur comme source de données à qualifier



(a) Sources utilisées pour les données du RGE ALTI®.



(b) Résolution spatiale des données altimétriques du RGE ALTI®

FIGURE II.1. Qualité des données altimétriques du RGE ALTI® en France métropolitaine

L'information géographique traditionnelle peut être produite par plusieurs sources de données différentes, qui influencent la qualité des données de manières différentes. C'est le cas, par exemple, du RGE ALTI® de l'IGN : ce jeu de données altimétriques

est issu de la fusion de différentes sources, ce qui explique l'hétérogénéité de la résolution spatiale des données (Figure II.1).

La saisie de l'information géographique volontaire se fait au moyen de différents outils de saisie (GPS, orthophoto, *etc.*), par conséquent, la qualité de l'information va également dépendre de la qualité de ces outils. Toutefois, le contributeur joue un rôle prépondérant dans la saisie des données géographiques. En effet, la qualité de données cartographiées à partir d'une orthophoto dépendra du soin mis par le contributeur à les tracer manuellement. Dans le cas des données issues de traces GPS, le contributeur influence la qualité de la collecte dans sa manière de porter l'appareil de mesure (Ivanovic, 2018), et éventuellement dans l'effort mis pour nettoyer ou corriger les données avant leur mise en ligne sur la plateforme collaborative. La compilation et l'interprétation des données étant des tâches propres au contributeur, elles font de lui un capteur de données des plus intéressants (Goodchild, 2007).

La qualité de l'information géographique volontaire va donc aussi dépendre des caractéristiques du contributeur qui la saisit : son expérience du terrain, son expérience de l'outil de saisie, mais aussi sa volonté de participer à enrichir la base avec des données de qualité. Cette notion de volonté est fortement liée à l'idée d'intentionnalité qui est inhérente à la définition du carto-vandalisme. Plusieurs questions se posent : comment qualifier le contributeur ? Quels sont les éléments de qualité du contributeur à considérer, dans le but de qualifier les données géographiques collaboratives par la suite ?

1.2 L'assurance qualité par l'approche sociale

Dans le chapitre précédent, nous avons présenté les trois approches formulées par Goodchild et Li (2012) pour assurer la qualité de l'information géographique volontaire. Parmi elles, l'approche sociale suppose que l'existence de contributeurs de confiance permet de structurer la communauté des contributeurs, et ainsi de garantir la qualité des données collaboratives. Cette approche concerne à la fois les contributeurs prédéfinis dans le projet qui ont un rôle privilégié de modération, et ceux dont le comportement dénote une confiance suffisante pour assimiler leur participation à un rôle de modération et/ou de contrôle qualité.

Le contrôle qualité de l'information géographique volontaire peut être effectué par différents moyens. La Table II.1 décrit les types de contrôle qualité effectués dans certains projets cartographiques collaboratifs ainsi que les personnes assignées à ces tâches. D'après Fonte *et al.* (2015), le contrôle qualité des données collaboratives peut être effectué par :

1. Des **méthodes automatiques** (pas d'intervention humaine) ;
2. La foule, *i.e.* la communauté des **volontaires** ;
3. Des **modérateurs**, *i.e.* des contributeurs de confiance identifiés dans la communauté ;
4. Des **experts**.

Dans la plupart des projets cartographiques – du moins ceux qui sont référencés dans la Table II.1 – la qualité de l'information géographique volontaire est assurée

par l'approche sociale dans une certaine mesure : certaines tâches de contrôle qualité sont réalisées par des modérateurs et des experts¹, tandis que d'autres sont laissées aux volontaires. Dans ce dernier cas, l'approche sociale peut être approfondie : en étudiant le rôle social de ces contributeurs *a priori* non qualifiés, l'évaluation de leur niveau de confiance permettra de mettre en évidence des rôles implicites qui permettent d'assurer la qualité des données.

TABLE II.1. *Le contrôle qualité dans des projets cartographiques collaboratifs, tiré et adapté de (Fonte et al., 2015)*

Nom du projet	Type de contrôle qualité	Tâche	Qui effectue la tâche ?
Geograph / Panoramio	Identification des erreurs et prévention des usages abusifs	Analyser le contenu des images	Modérateurs
		Vérifier à partir de cartes et d'images satellites	Modérateurs
		Suggérer des modifications	Volontaires
Flickr	Prévention des usages abusifs	Analyser le contenu des images	Modérateurs
Geo-Wiki	Créer des points de contrôle	Classer des images satellites	Experts
	Assigner les classes d'utilisation des sols dans la classe majoritaire	Choisir la classe la plus fréquente lorsqu'il n'y a pas de données de contrôle	Automatique
	Évaluation de la crédibilité des contributeurs	Comparer les contributions avec des données de contrôle	Automatique
		Comparer les contributions avec les valeurs majoritaires	Automatique
	Identification de la confiance de classification	Les contributeurs indiquent la confiance de leur saisie	Volontaires
	Identification du type d'imagerie utilisée	Indiquer la résolution de l'image utilisée	Volontaires
	Identification des contributions erronées	Indiquer que la classe assignée au pixel est incorrecte	Experts
Volontaires			
OpenStreetMap	Identification des erreurs	Lever des alertes sur des erreurs potentielles	Automatique
			Volontaires
	Prévention des incidents de dégradation massive	Bannissement des contributeurs	Administrateurs
	Correction des erreurs	Édition des données erronées	Volontaires

Il existe plusieurs manières de contribuer à un projet d'information géographique. Rehr et al. (2013) identifient quatre activités à travers lesquelles un contributeur peut s'investir :

1. **Contribuer de l'information géographique** : par exemple, ajouter, éditer des données ;
2. **Construire des structures de communautés** : par exemple, organiser des carto-parties (*mapping parties*), gérer les *mailing lists* ;
3. **Travailler sur les règles et les normes** : par exemple, contribuer à éditer les définitions sur le Wiki du projet, s'il y en a un ;
4. **Travailler sur les outils** : par exemple, développer les outils d'édition des données.

Comme nous cherchons à qualifier les données collaboratives, l'activité de contribution des données géographiques est celle qui nous intéresse ici. Il s'agit donc, dans

1. *N.B.* : Les experts ne sont pas forcément des professionnels. Par exemple, un contributeur expert peut être un contributeur non-professionnel expérimenté.

la suite, d'étudier le rôle social des contributeurs, à partir de leur comportement, à contribuer des données géographiques, afin d'en dégager des rôles particuliers qui garantissent la qualité de ces données.

1.3 Collaboration et interaction

a) Les opérations élémentaires de contribution cartographique

Vocabulaire : Un **objet géographique** désigne un élément géographique du monde réel qui a été saisi dans une base de données géographiques. Lorsqu'il est visible dans l'espace cartographique, nous le désignons par le terme d'**objet cartographique**. Dans OSM, un objet cartographique peut être supprimé de la carte, dans ce cas, il n'est plus visible dans l'espace cartographique. Toutefois, l'objet géographique correspondant reste enregistré sur la base de données (il porte alors l'attribut `visible=false` qui le rend invisible sur la carte).

En fonction du nombre de modifications qui lui sont faites sur la plateforme collaborative, il peut y avoir plusieurs **versions** d'un objet géographique. Nous utiliserons le terme de **contribution** pour désigner une version donnée d'un objet géographique. Ainsi, un contributeur est l'auteur d'une ou plusieurs versions d'un objet géographique.

Les travaux de [Rehrl *et al.* \(2013\)](#) proposent une conceptualisation des éléments de contributions de l'information géographique volontaire. Dans ce modèle, l'objet cartographique est composé d'une géométrie et d'un ensemble d'attributs. Dans OSM, un attribut est un tag, défini par une paire de clé-valeur. Un objet cartographique peut avoir un ensemble d'attributs vide (donc il n'est décrit par aucun attribut), en revanche, il est obligatoirement composé d'une géométrie, simple ou complexe. Pour plus de détails sur le modèle conceptuel de l'information géographique, nous renvoyons le lecteur à l'article de [Rehrl *et al.* \(2013\)](#). La contribution de données géographiques s'effectue à partir de trois opérations élémentaires, à savoir : l'ajout, la mise à jour et la suppression de données. Les opérations élémentaires atomiques portent sur l'objet dans sa globalité.

(i) **Ajout d'un nouvel objet cartographique :** Cette opération consiste à saisir la première version d'un objet géographique dans la base de données. *A priori*, cette opération ne manifeste pas de collaboration directe entre deux contributeurs. Toutefois, lorsqu'un contributeur ajoute un nouvel objet cartographique en se basant sur des contributions produites par un autre contributeur, on peut y voir une forme de collaboration entre contributeurs. Dans le cas du projet OSM, il est possible de cartographier une nouvelle route en s'appuyant sur un point d'intersection qui existait déjà sur la carte. Dans ce cas, le contributeur de la route manifeste une certaine confiance envers le contributeur qui avait saisi le point d'intersection.

(ii) **Mise à jour d'un objet cartographique :** Elle consiste à produire une nouvelle version d'un objet cartographique qui préexiste dans l'espace cartogra-

phique collaboratif. Mis à part le cas où le contributeur met à jour ses propres contributions, le comportement collaboratif est assez évident dans l'opération de mise à jour. Toutefois, selon le type de modification opéré, la collaboration portera une signification différente. [Lodigiani et Melchiori \(2016\)](#) divisent la mise à jour de données en deux sous-opérations, à savoir la correction et la complétion de données :

- la correction signifie que la version précédente de l'objet était incorrecte ;
- la complétion implique que les éléments pré-existants de la version précédente de l'objet géographique sont justes.

La mise à jour des données peut aussi être divisée en trois sous-opérations que sont l'ajout, la modification et l'élimination d'information ([Rehrl et al., 2013](#)). Ces sous-opérations peuvent porter sur la géométrie et les attributs de l'objet. Par exemple, la suppression des tags d'un objet OSM est une mise à jour de type élimination d'information. L'élimination d'une géométrie est possible dans le cas d'un objet de géométrie complexe, par exemple pour une ligne composée d'un ensemble nœuds reliant les tronçons de la ligne : dans ce cas, l'élimination du nœud met à jour l'objet ligne. Toutefois, au niveau du nœud éliminé, cette opération est une véritable suppression.

L'opération de correction de données est composée d'une élimination suivie d'une modification, puisqu'un élément de la version précédente disparaît pour être remplacé par un autre élément. Dans l'extrait de l'historique de la Table I.6 du Chapitre 1, la version n°3 est un exemple de correction de données. En revanche, l'opération de complétion de données correspond à un ajout d'information. Par exemple, le Grand Bassin Rond au Jardin des Tuileries (Figure II.2) a été complété sur trois versions successives (Table II.2) : à chaque nouvelle version, un nouveau tag est ajouté.



FIGURE II.2. Grand Bassin Rond du Jardin des Tuileries à Paris. Source : *OpenStreetMap*.

TABLE II.2. Extrait de l'historique des tags du Grand Bassin Rond de la Figure II.2.

Version Tag key	n° 10	n° 11	n° 12
barrier			wall
height			0.4
name		Grand Bassin Rond	Grand Bassin Rond
natural	water	water	water
source	bing; survey	bing; survey	bing; survey

(iii) **Suppression de données** : Elle consiste à retirer un objet de l'espace cartographique collaboratif. Précisons que la suppression des données concerne l'objet dans sa globalité, et cette opération n'est pas à confondre avec la mise à jour par élimination d'information (géométrique ou attributaire). La suppression d'un objet ne signifie pas nécessairement que celui-ci soit complètement éliminé de la base de données, mais l'objet n'est plus visible sur la carte. C'est le cas des données OSM : un objet qui est supprimé de la carte reste enregistré dans la base de données, cependant son attribut `visible`, fixé à la valeur `false`, le rend invisible sur la carte. Lorsqu'un contributeur supprime la contribution produite par une autre personne, il signifie que cette contribution n'est pas pertinente voire qu'elle est totalement fausse. Cela peut donner un renseignement sur la fiabilité du contributeur dont l'objet a été supprimé. La Figure II.3 donne un exemple de carto-vandalisme fantaisiste repéré par des contributeurs OSM : l'opération de suppression de l'objet fictif met en évidence le manque de fiabilité de l'auteur de cette contribution.



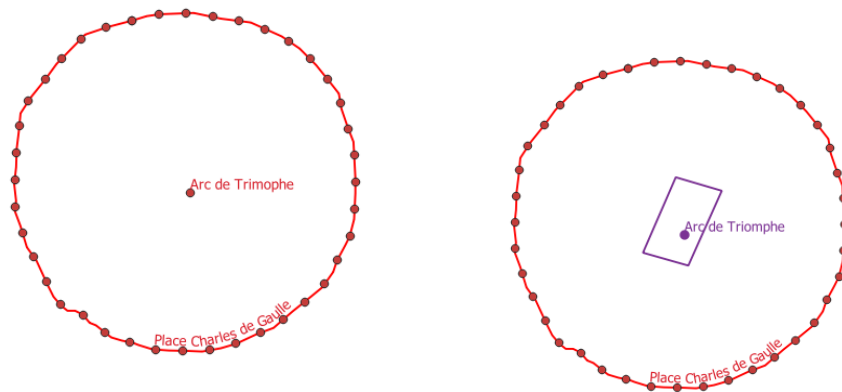
FIGURE II.3. Exemple de carto-vandalisme repéré par les contributeurs OSM

Pour résumer, les opérations élémentaires qui constituent l'activité de contribution cartographique sont :

- l'ajout d'un nouvel objet cartographique à partir de rien ;
- l'ajout d'un nouvel objet cartographique par réutilisation de données préexistantes dans la base de données ;
- la modification d'un objet cartographique (au niveau attributaire ou géométrique) ;
- la complétion d'information dans un objet cartographique ;
- la suppression d'information dans un objet cartographique.

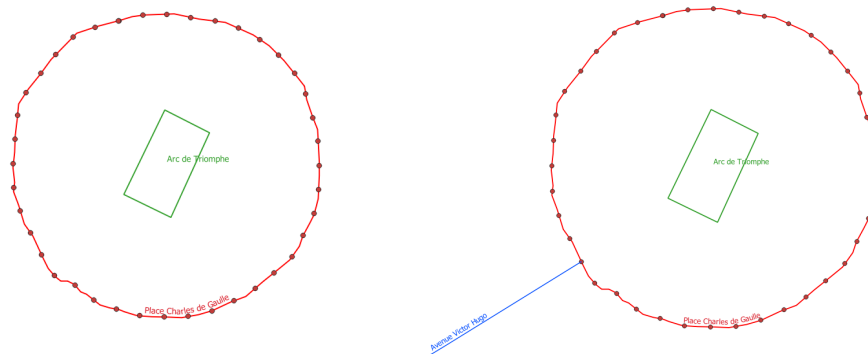
La Figure II.4 donne un exemple de cartographie collaborative par quatre contributeurs :

1. *Utilisateur1* ajoute deux objets : la Place Charles de Gaulle et l'Arc de Triomphe. Ce dernier est nommé « Arc de Trimophe ».
2. *Utilisateur2* corrige le nom de l'Arc de Triomphe et crée un bâtiment autour du point.
3. *Utilisateur3* modifie la géométrie du bâtiment, supprime le point nommé « Arc de Triomphe », et nomme le bâtiment par « Arc de Triomphe ».
4. *Utilisateur4* crée un nouvel objet nommé « Avenue Victor Hugo » qui s'appuie sur un nœud de la Place Charles de Gaulle.



(a) *Utilisateur1* : Ajout d'un point d'intérêt portant le tag « Arc de Trimophe » et d'un objet portant le tag « Place Charle de Gaulle »

(b) *Utilisateur2* : Ajout d'un bâtiment et correction du tag « Arc de Trimophe » en « Arc de Triomphe »



(c) *Utilisateur3* : Modification géométrique du bâtiment, complétion du bâtiment par le tag « Arc de Triomphe » et suppression du point d'intérêt

(d) *Utilisateur4* : Ajout de l'avenue Victor Hugo par réutilisation d'un nœud composant la place Charles De Gaulle

FIGURE II.4. Opérations élémentaires

À partir des différentes opérations élémentaires effectuées dans cet exemple, nous pouvons modéliser des interactions entre ces quatre contributeurs (Figure II.5). Ces interactions témoignent de différentes formes de collaboration, dont l'analyse permet de qualifier le comportement des contributeurs.

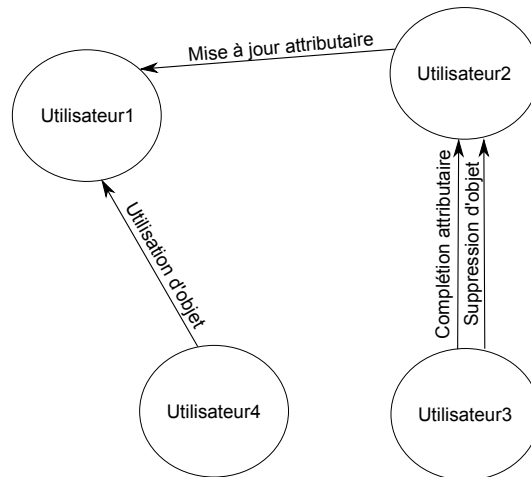


FIGURE II.5. Modélisation des interactions entre les quatre contributeurs de la Figure II.4.

b) Modéliser des interactions pour analyser les collaborations

La plupart des travaux portant sur la collaboration dans les projets d'information géographique volontaire ont été réalisés à partir des données OSM. En considérant l'historique des éditions de données, ceux-ci modélisent les interactions entre contributeurs sous la forme de réseaux sociaux. L'avantage de cette modélisation est l'exploitation des principes et des raisonnements issus de la théorie des graphes. La Figure II.6 donne un exemple d'édition de données entre trois contributeurs A, B et C sur trois objets cartographiques. Cet exemple de base est repris dans la suite pour illustrer les différents graphes proposés dans la littérature scientifique.

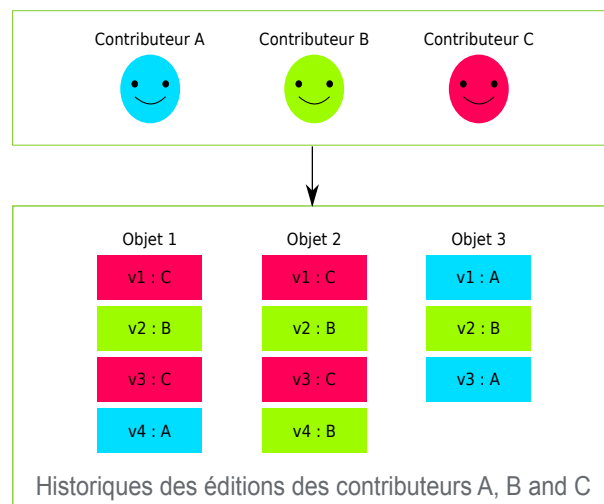


FIGURE II.6. L'historique des éditions des trois objets par trois contributeurs A, B et C.

Mooney et Corcoran (2012b, 2014) construisent un graphe de co-édition pour étudier l'existence d'une structure sociale – c'est-à-dire similaire à celle d'un réseau social – dans la communauté des contributeurs des données OSM à Londres. Dans un graphe de co-édition, deux contributeurs interagissent lorsqu'ils co-éditent un objet, c'est-à-dire que l'un édite une nouvelle version d'un objet cartographique dont l'autre avait édité la version précédente. Dans ce cas, la relation d'interaction est orientée du contributeur de la version v au contributeur de la version $v - 1$.

L'intensité de l'interaction – c'est-à-dire le poids de l'arc – est alors pondérée par le nombre de co-éditions qui se sont produites entre les deux contributeurs. En cela, le graphe de co-édition ne représente que les éditions directes entre les contributeurs. La Figure II.7 donne le graphe de co-édition qui résulte de l'exemple proposé à la Figure II.6.

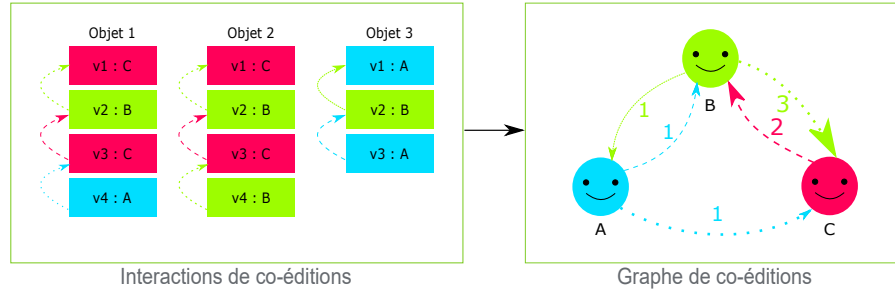


FIGURE II.7. Graphe de co-édition, adapté de *Mooney et Corcoran (2014)*

Stein et al. (2015) ont également analysé les opérations d'éditations entre contributeurs pour mettre en évidence des motifs de contribution. Dans ce cas, la modélisation des interactions tient compte des éditions indirectes entre contributeurs. Celle-ci se base sur la notion de communication imbriquée (*interlocking communication*), où il s'agit non seulement de relier le contributeur d'une version v d'un objet au contributeur de la version $v - 1$, mais également à tous les contributeurs des versions qui précèdent jusqu'à la version $v - k$ si la version $v - k - 1$ a été saisie par le contributeur de la version courante v . Si le contributeur de la version v édite pour la première fois cet objet, alors on considèrera qu'il interagit avec tous les contributeurs des versions précédentes jusqu'à la version initiale (*Stein et Blaschke, 2010*).

L'intensité des interactions peut être quantifiée selon les concepts de largeur de collaboration ou de profondeur de collaboration. La largeur de collaboration consiste à compter le nombre d'objets cartographiques différents sur lesquels deux contributeurs ont interagi. En revanche, la profondeur de la collaboration quantifie le nombre maximal d'interactions entre deux contributeurs sur un même objet cartographique. La Figure II.8 donne les graphes de largeur et de profondeur collaboration obtenus à partir de l'exemple de base de la Figure II.6.

En analysant conjointement ces graphes, *Stein et al. (2015)* ont mis en évidence différents profils de contributeurs selon l'intensité de leur comportement collaboratif (Table II.3). Notons que le profil des contributeurs ponctuels est caractérisé par une collaboration peu profonde. En revanche, ces derniers peuvent modifier un grand nombre de données.

TABLE II.3. Profils de contributeurs identifiés par les graphes de collaboration

Largeur \ Profondeur		Profondeur	
		Faible intensité	Forte intensité
Faible intensité		Experts	Contributeurs chevronnés
Forte intensité		Contributeurs ponctuels	Contributeurs ponctuels

Ma et al. (2015) construisent un graphe de co-contribution dans lequel les éditions

directes et indirectes sont également considérées : les contributeurs interagissent à partir du moment où ils ont édité le même objet (Figure II.9). Cependant, contrairement aux graphes précédents, le graphe de co-contribution est non orienté et non pondéré. Le graphe de co-contribution a été construit pour démontrer la distribution inégale entre les contributeurs actifs et inactifs sur OSM en termes de collaboration. Contrairement aux autres graphes d'interaction, il n'avait pas vocation à permettre d'identifier des profils de contributeurs.

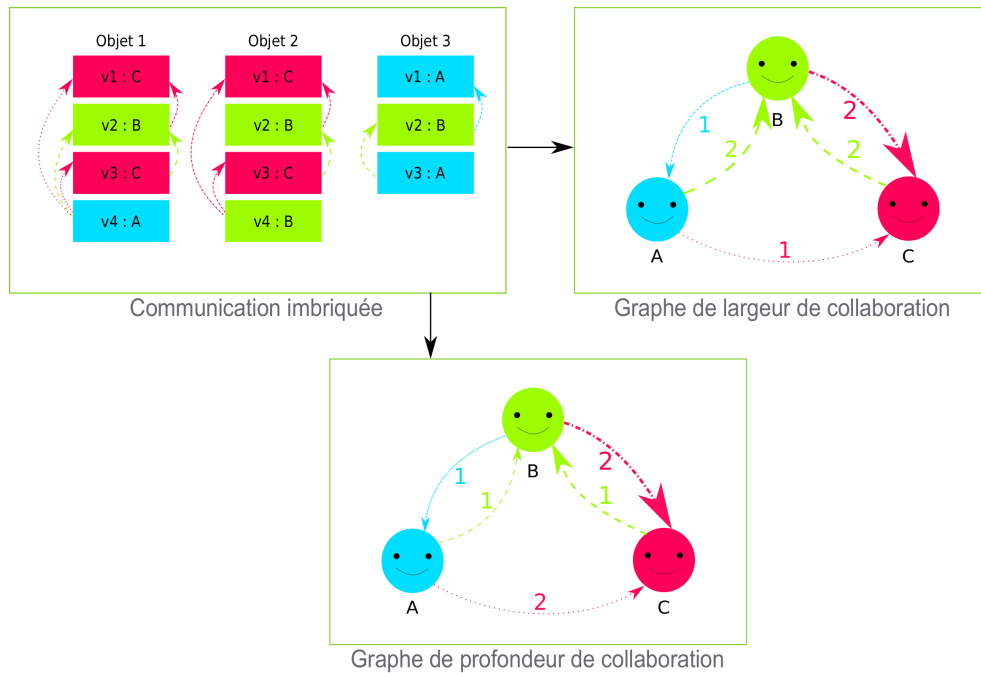


FIGURE II.8. Graphes de collaboration, adaptés de Stein et al. (2015)

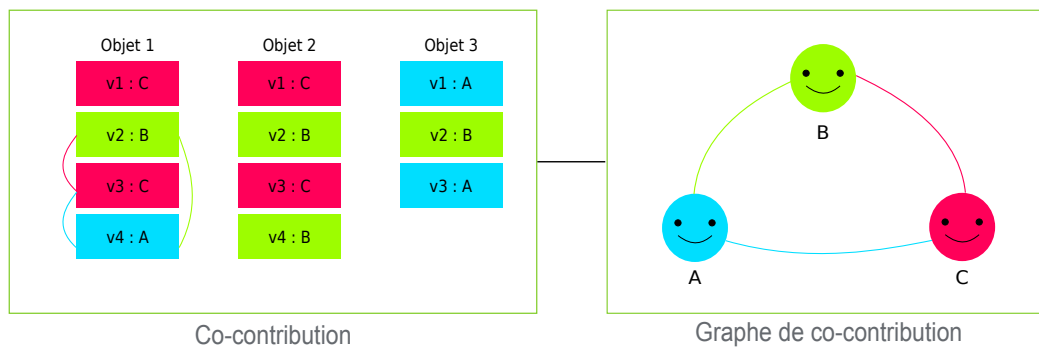


FIGURE II.9. Graphe de co-contribution, adapté de Ma et al. (2015)

c) Interprétation des collaborations

D'après Ikeda *et al.* (2016), il existe trois modes de collaboration. Le mode de collaboration séquentielle permet aux contributeurs de participer sur la base des contributions des uns et des autres afin d'améliorer leur qualité itérativement. Le mode de collaboration simultanée consiste à laisser les contributeurs participer de manière indépendante des uns et des autres. Enfin, le mode de collaboration hybride combine les deux premiers modes de collaboration présentés. Le fonctionnement du

projet OpenStreetMap repose sur un mode de collaboration hybride où les contributeurs peuvent saisir de manière simultanée tout en ayant la possibilité de corriger les contributions produites par d'autres.

1.4 De la confiance des contributeurs à la fiabilité de l'information

La représentation des interactions entre contributeurs sous forme de graphes sociaux permet d'étudier le comportement collaboratif de ces derniers. À présent, la question est de définir ce que l'on cherche à évaluer à travers l'analyse de ces interactions. Rappelons que le but de cette étude est de qualifier le contributeur à travers son comportement collaboratif.

Dans la littérature scientifique, l'étude des contributeurs a permis de les regrouper selon plusieurs typologies. Celles-ci se basent sur différents critères tels que leur niveau de participation, leur zone d'activité cartographique ou encore leur motivation (Neis et Zielstra, 2014). Neis et Zipf (2012) proposent une typologie des contributeurs basée sur la quantité de contribution et leur activité spatio-temporelle. Les contributeurs, selon leur niveau de participation cartographique, peuvent être qualifiés de : cartographe senior, cartographe junior, cartographe non-récurrent et cartographe sans édition. Cette typologie implique que les données de qualité proviennent des contributeurs qui ont une forte participation dans le projet. Dans la même idée, la typologie proposée par Bégin *et al.* (2016) qualifie les contributeurs selon leur ancienneté dans le projet collaboratif cartographique.

Les contributeurs peuvent également être catégorisés selon leur mode de contribution. Il s'agit de considérer, par exemple, la préférence des contributeurs à contribuer selon une ou plusieurs opérations d'édition : certains contributeurs peuvent participer uniquement par l'ajout de nouveaux objets, d'autres en ne faisant que des modifications géométriques, ou attributaires voire exclusivement les deux (Mooney et Corcoran, 2012b). Bégin *et al.* (2013) considèrent la préférence des contributeurs OSM à éditer des objets précis, tels que les routes, les aménités ou les adresses. En analysant les préférences cartographiques des contributeurs, il sera alors possible d'évaluer la complétude des données sur une zone cartographique par type d'objets selon les contributeurs qui l'ont éditée. Par exemple, un contributeur qui n'édite que des routes rendra l'espace cartographique plus complet en ce qui concerne le réseau routier que sur les autres types de données (les bâtiments, les points d'intérêts, *etc.*)

De plus, l'expertise du contributeur est un facteur d'influence de la qualité des données saisies. Coleman *et al.* (2010) dressent 5 catégories de contributeurs, selon leur compétence dans le domaine de l'information géographique : néophyte, amateur intéressé, amateur expert, professionnel expert et autorité experte. Néanmoins les frontières entre ces catégories ne sont pas clairement délimitées, et aucune méthodologie n'a été proposée pour classifier les contributeurs selon cette typologie. Dans la même perspective d'évaluation de la compétence des contributeurs, Comber *et al.* (2016) comparent la qualité des contributions du projet GeoWiki en se basant sur l'expertise des contributeurs et leurs origines géographiques². L'idée est de montrer

2. *N.B.* Ici, les contributeurs avaient eux-mêmes fourni leurs informations personnelles

que les experts et les contributeurs locaux produisent généralement des contributions de qualité. Cette idée est partagée par [Stein et al. \(2015\)](#) qui assimilent les contributeurs locaux d'une zone cartographique à des experts. [Yang et al. \(2016\)](#) proposent une méthodologie permettant de déterminer l'expertise des grands contributeurs OSM à partir des indicateurs de pratique, de compétence et de motivation. Enfin, le genre, l'âge et la profession des contributeurs sont d'autres facteurs de qualité de l'information géographique volontaire ([Schmidt et Klettner, 2013](#); [Duféal et al., 2016](#)).

La Figure II.10 détaille les différentes typologies de contributeurs qui viennent d'être présentées. Ces nombreuses typologies sont intéressantes, cependant, elles ne tiennent compte que d'un seul critère à la fois. Par conséquent, ces typologies sont incomplètes au sens où elles ne parviennent pas à refléter la complexité des comportements humains. Plutôt que de les remettre en cause, nous proposons de les combiner pour offrir une représentation plus complète des profils de contributeurs.

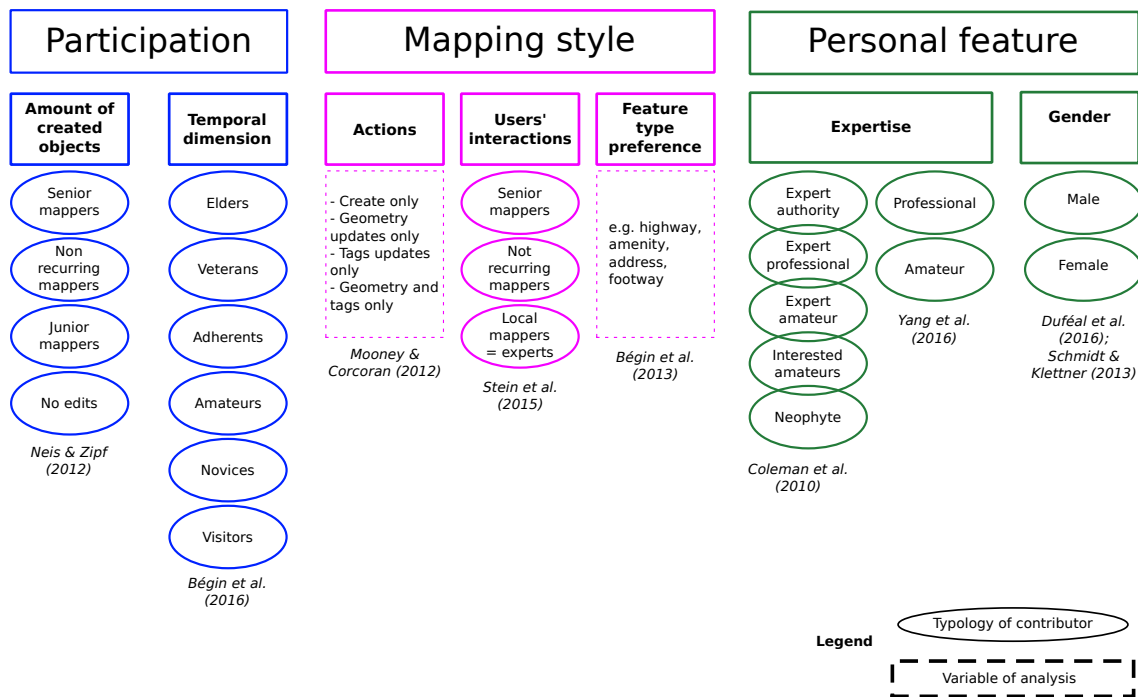


FIGURE II.10. Les typologies de contributeurs proposées dans l'état de l'art.

La détection du carto-vandalisme selon une approche sociale consiste à identifier les « carto-vandales », c'est-à-dire des contributeurs qui ne sont pas dignes de confiance. Il s'agit donc d'évaluer, à travers les graphes d'interaction, le niveau de confiance que l'on peut accorder aux contributeurs. Dans le but de la quantifier, [Bishr et Janowicz \(2010\)](#) ainsi que [Golbeck \(2005\)](#) reprennent une définition sociologique de la confiance comme étant « un pari sur les éventuelles actions futures d'autres personnes »³. [Golbeck \(2005\)](#) la vulgarise en formulant la définition suivante :

« Alice a confiance en Bob si elle s'engage dans une action basée sur la conviction que les actions futures de Bob auront de bonnes conséquences ».

3. 'a bet about the future contingent actions of others'

Cette confiance « interpersonnelle » (entre deux personnes) peut alors se transformer en une confiance dite « informationnelle ». [Bishr et Janowicz \(2010\)](#) parlent à ce propos de transitivité personne-objet pour décrire le passage d’une confiance interpersonnelle (entre deux personnes) à une confiance informationnelle ([Bishr et Janowicz, 2010](#)). En effet, si Alice a confiance en Bob alors Alice a confiance dans les informations fournies par Bob. Dans le cas de l’information géographique volontaire, la confiance informationnelle correspondrait donc à la relation de confiance entre un contributeur (Alice) et les données géographiques collaboratives produites par un autre contributeur (Bob).

Nous constatons donc que la confiance est un indicateur de qualité d’une personne – qui peut être vue comme une source d’information – et du contenu de l’information qu’elle apporte. Par ailleurs, la notion de confiance est rattachée à d’autres concepts tels que la crédibilité, la fiabilité, l’honnêteté, de sorte que ces termes sont souvent utilisés de manière indifférente ([Skarlatidou et al., 2011](#)). Un tel d’abus de langage introduit alors une certaine ambiguïté sur ces différents concepts, car ceux-ci peuvent porter sur la source ou le contenu de l’information. Cependant, dans d’autres domaines des sciences informatiques tels que celui des interactions homme-machine, les informations produites par des machines ont incité les chercheurs être plus rigoureux dans la définition des concepts liés à la confiance des données et des sources de données ([Fogg et Tseng, 1999](#)).

[Revault d’Allonnes \(2013\)](#) modélise la confiance d’une information comme un ensemble de quatre dimensions : la fiabilité et la compétence portent sur la source de l’information, tandis que la plausibilité et la crédibilité portent sur le contenu de l’information (Table II.4). Dans cette modélisation, la fiabilité et la plausibilité sont deux dimensions subjectives de la confiance : elles dépendent de la personne qui consulte l’information. En revanche, la compétence et la crédibilité sont des dimensions objectives de la confiance, mais qui dépendent d’un contexte donné, contrairement à la fiabilité et la plausibilité.

TABLE II.4. *Dimensions de la confiance, d’après Revault d’Allonnes (2013)*

	Général	Contexte
Source de l’information	Fiabilité	Compétence
Contenu de l’information	Plausibilité	Crédibilité

[Muttaqien et al. \(2018\)](#) proposent une méthode pour évaluer l’expertise – ou la compétence – du contributeur afin de qualifier la crédibilité de l’information géographique, car celle-ci est moins complexe à modéliser que la fiabilité. Dans le contexte du carto-vandalisme, il s’agit de détecter les contributions produites intentionnellement pour dégrader l’espace collaboratif. Par conséquent, un contributeur carto-vandale n’est pas incompetent mais non fiable : il s’agit donc d’évaluer en particulier la fiabilité des contributeurs dans notre cadre d’étude.

En plus des dimensions de fiabilité et de compétence, la confiance liée à la source d’information peut se modéliser de manière plus complète en considérant 3 composantes supplémentaires ([Lesot et Revault d’Allonnes, 2017](#)) : la première est l’intention de la source, c’est-à-dire son inclination à dire la vérité. Autrement dit, elle correspond à un degré de sincérité. La seconde est l’implication de la source,

c'est-à-dire son degré d'engagement. Pour le cas des contributeurs de données collaboratives, cela correspond au niveau de participation et d'investissement dans le projet collaboratif. La troisième composante concerne le degré de « propos rapportés » par cette source, c'est-à-dire le nombre de sources successives dont provient l'information communiquée. Dans le cadre de l'information géographique volontaire, les contributeurs peuvent importer des données provenant d'autres bases de données (par exemple le cadastre français ou des données européennes d'occupation du sol), et dans ce cas ils ne sont pas les sources « primaires » de ces données. En revanche, lorsque les contributeurs de données géographiques collaboratives ajoutent des informations provenant directement de leurs connaissances personnelles sur un élément géographique, ceux-ci sont les sources « primaires » de ces données.

La recherche en psychologie sociale affirme que la confiance du contributeur⁴ est le fruit d'une perception sociale : la cognition sociale humaine repose sur les notions d'agrément (*warmth*) et de compétence, qui sont deux dimensions universelles, la première étant subjective et la seconde objective (Fiske *et al.*, 2007). Par conséquent, la difficulté à évaluer la confiance d'un contributeur de données géographiques est due à la prise en compte de caractéristiques non seulement objectives mais aussi subjectives (Flanagin et Metzger, 2008).

Pour traiter ce problème de subjectivité, une solution possible est de considérer la réputation du contributeur. En effet, la réputation d'un contributeur résulte de la confiance que lui accorde chaque membre d'une communauté. La réputation d'un contributeur peut donc être modélisée comme une moyenne calculée à partir des relations de confiance entre celui-ci et chacun des autres contributeurs de la communauté (Kessler et de Groot, 2013; Fogliaroni *et al.*, 2018). Ainsi, la réputation du contributeur est un indicateur de sa confiance, qui servira à son tour d'indicateur de qualité à ses contributions. Le lien entre la réputation du contributeur et la confiance informationnelle est double. En effet, la confiance peut se calculer à partir des contributions pour évaluer la réputation du contributeur (D'Antonio *et al.*, 2014; Forati et Karimipour, 2016), mais la confiance informationnelle peut elle-même découler de la réputation de son contributeur (Kessler et de Groot, 2013). La Figure II.11 récapitule les différents indicateurs de confiance et de réputation considérés dans la littérature scientifique pour évaluer la qualité de l'information géographique volontaire.

Dans ce travail, nous cherchons à qualifier le comportement du contributeur pour évaluer la confiance des données et ainsi qualifier le carto-vandalisme de l'information géographique volontaire. Notre démarche se rapproche de celle de Neis *et al.* (2012), où il s'agit de tenir compte de la réputation du contributeur pour déterminer si une contribution relève d'un acte de carto-vandalisme. Par ailleurs, la mise en place d'un système de réputation pour détecter les contributeurs malicieux sur les plateformes collaboratives est une solution qui a été proposée dans des travaux antérieurs (Kamvar *et al.*, 2003; Huang *et al.*, 2010). La première étape de détection du carto-vandalisme consiste donc en une analyse des relations de confiance des contributeurs à travers leurs interactions cartographiques.

4. La confiance du contributeur est à entendre comme la confiance que l'on peut accorder à ce contributeur.

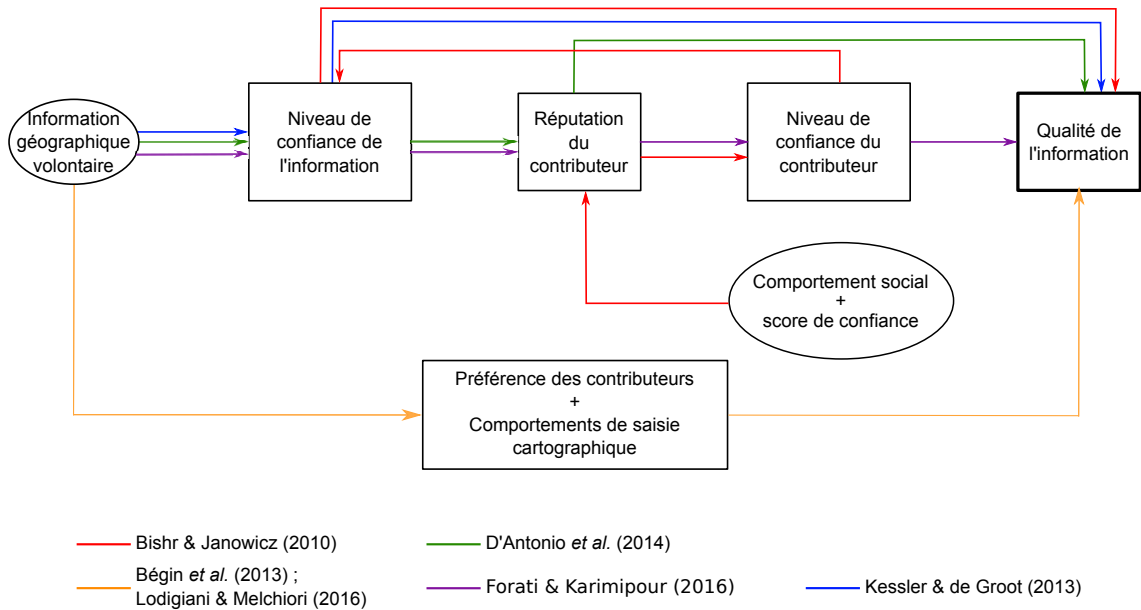


FIGURE II.11. Confiance et qualité des contributeurs et des contributions

2 Modèle de graphe social multiplexe

2.1 Principe et définitions

Notre état de l'art a montré que la modélisation par des graphes est un moyen adapté pour étudier le comportement collaboratif des contributeurs de données géographiques. Par ailleurs, les graphes décrits dans la Section 1.3 précédente présentent tous un certain intérêt, chacun d'entre eux mettant en évidence un aspect particulier du comportement de collaboration des contributeurs. Cependant, comme notre objectif est d'étudier le comportement de collaboration des contributeurs dans sa globalité, il nous faut considérer un modèle de graphe social qui tienne compte de toutes les interactions pouvant se produire entre les contributeurs, quels que soient leurs types.

Dans le domaine de l'analyse des réseaux sociaux, le concept de réseau multiplexe (Kivelä *et al.*, 2014), appelé aussi réseau multidimensionnel (Kazienko *et al.*, 2011), a été introduit pour modéliser les différentes natures de relations qui peuvent exister entre des individus, comme l'illustre la Figure II.12. Un réseau social multiplexe est, tout compte fait, un multigraphe dans lequel deux personnes seront reliées par autant d'arcs qui représentent les différentes relations qu'elles partagent. Dans le cadre de l'étude du comportement de collaboration des contributeurs d'information géographique, nous considérons un réseau social multiplexe, où deux contributeurs peuvent être reliés par plusieurs arcs correspondant aux différentes interactions qui se produisent à partir de leur activité cartographique.

a) Définition formelle du réseau multiplexe de collaboration

Nous reprenons la définition du réseau multiplexe telle qu'elle a été formulée par Kivelä *et al.* (2014), en l'adaptant à notre étude des contributeurs d'information

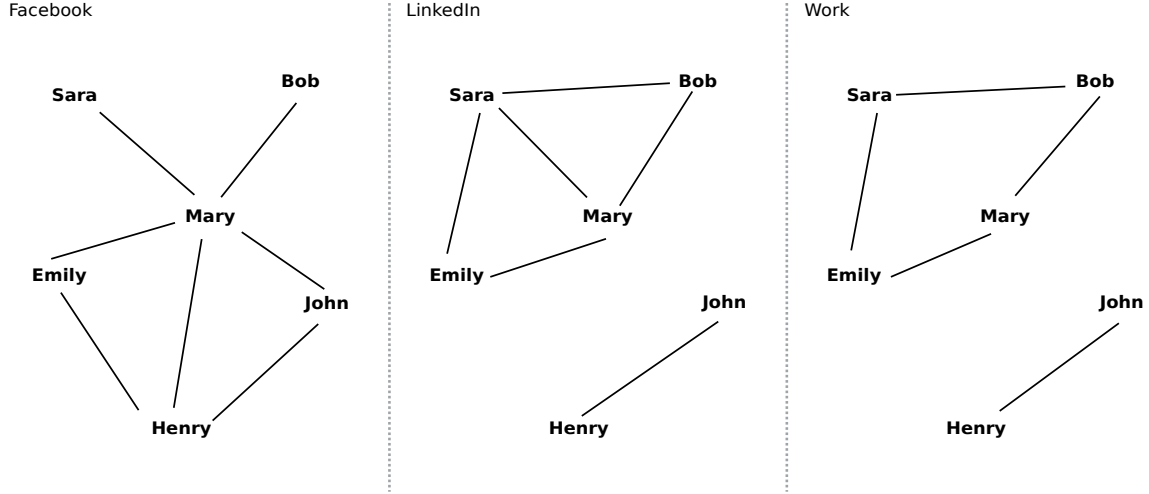


FIGURE II.12. Graphes sociaux des mêmes personnes dans différents environnements

géographique.

(i) **Graphe social des contributeurs** : Soit G_k un graphe social décrivant des interactions entre des contributeurs d'information géographique volontaire :

$$G_k = (V_k, E_k) \quad (\text{II.1})$$

L'ensemble des sommets V_k correspond à l'ensemble des contributeurs d'information géographique, et l'ensemble des arcs E_k correspond aux interactions qui se produisent entre eux. E_k est défini comme un ensemble de triplets (i, j, w) :

$$E_k = \{(i, j, w)\} \quad \forall i, j \in V_k \quad w \in \mathbb{R} \quad (\text{II.2})$$

où le poids w de l'arc (i, j) représente l'intensité de l'interaction entre deux contributeurs i et j . Si G_k est un graphe non pondéré, alors $w = 1$ sur chaque arc.

(ii) **Réseau multiplexe** : Le réseau multiplexe G se définit comme un ensemble de α graphes G_k (aussi appelés couches ou dimensions de G) :

$$G = \{G_k\}_{k=1}^{\alpha} = \{(V_k, E_k)\} \quad (\text{II.3})$$

où l'ensemble des sommets V_k est identique dans toutes les couches du réseau multiplexe G :

$$V_k = V_l = V \quad \forall k, l = 1..{\alpha} \quad (\text{II.4})$$

Si $\alpha = 1$, on dit que G est un **réseau monoplexe**.

b) Métriques multiplexes

Pour étudier les caractéristiques structurelles d'un réseau multiplexe, nous présentons ici certaines métriques définies par [Battiston et al. \(2014\)](#).

(i) **Matrice d'adjacence** : Chaque graphe peut se décrire par une matrice d'adjacence, qui indique, pour chaque sommet du graphe, la composition de son voisinage. Soit $A^{[k]}$ la matrice d'adjacence du graphe G_k :

$$A^{[k]} = \{a_{ij}^{[k]}\} \quad (\text{II.5})$$

Chaque élément $a_{ij}^{[k]}$ indique l'intensité avec laquelle le contributeur i interagit avec le contributeur j :

$$a_{ij}^{[k]} = \begin{cases} w & \text{si } (i, j, w) \in E_k \\ 0 & \text{sinon} \end{cases} \quad (\text{II.6})$$

La matrice d'adjacence A du réseau multiplexe G se définit alors comme l'ensemble des matrices d'adjacence de chacune de ses couches :

$$A = \{A^{[1]}, \dots, A^{[\alpha]}\} \quad (\text{II.7})$$

(ii) **Centralité de degré** : Soit m une couche du réseau multiplexe, et soit i un sommet du graphe m . Le degré $k_i^{[m]}$ est le nombre d'arcs entrants et sortants du sommet i dans le graphe m . Précisons que le degré $k_i^{[m]}$ est non-pondéré, c'est-à-dire que les arcs connectés au sommet i sont comptés indépendamment de leurs poids. Le degré quantifie donc le nombre d'interactions du contributeur i dans le graphe m .

(iii) **Degré de recouvrement** : Soit $k_i^{[m]}$ le degré de i dans le graphe m . Le degré de recouvrement (*overlapping degree*) o_i du sommet i est défini par :

$$o_i = \sum_m k_i^{[m]} \quad (\text{II.8})$$

Le degré de recouvrement o_i quantifie donc le nombre total d'interactions du contributeur i dans tout le réseau multiplexe.

(iv) **Coefficient de participation** : Soit un réseau multiplexe composé de α couches. Le coefficient de participation P_i du sommet i est défini par :

$$P_i = \frac{\alpha}{\alpha - 1} \left[1 - \sum_{m=1}^{\alpha} \left(\frac{k_i^{[m]}}{o_i} \right)^2 \right] \quad (\text{II.9})$$

où $k_i^{[m]}$ est le degré du sommet i et o_i le degré de chevauchement du sommet i . Le coefficient de participation P_i prend une valeur dans $[0, 1]$: $P_i = 0$ lorsque tous les arcs du sommet i sont concentrés dans une seule couche ; $P_i = 1$ lorsque le sommet i a le même nombre d'arcs dans chaque couche de réseau multiplexe.

Remarque : La définition du coefficient de participation impose que :

- $\alpha \neq 0$: le réseau multiplexe doit être constitué de deux couches au minimum ;
- $o_i \neq 0$: le sommet i a au moins un voisin dans l'une des couches du réseau.

(v) **Coefficient de clustering** : Le coefficient de *clustering* est un indicateur de la connectivité du voisinage d'un sommet (Watts et Strogatz, 1998). Il mesure l'intensité avec laquelle les voisins d'un sommet donné sont également reliés entre eux. Battiston *et al.* (2014) définissent le coefficient de *clustering* d'un sommet i par la formule suivante :

$$C_i = \frac{\sum_{j \neq i, m \neq i} a_{ij} a_{jm} a_{mi}}{k_i(k_i - 1)} \quad (\text{II.10})$$

où a_{ij} est un élément de la matrice d'adjacence du graphe et k_i est le degré du sommet i .

c) Détection de communautés dans les graphes

La détection de communautés dans les graphes est une méthode de classification non-supervisée qui consiste à regrouper les sommets d'un graphe en fonction de leur similarité. Chaque groupe de sommets forme une communauté. Techniquement, le partitionnement des sommets du graphe doit prendre en considération la structure du graphe, de manière à ce que les sommets d'une même communauté soient reliés par un grand nombre d'arcs, tandis que relativement peu d'arcs relient les communautés entre elles (Schaeffer, 2007). La Figure II.13 donne un exemple de détection de communautés dans un graphe social. La détection de communautés révèle des informations sociales telles que des liens de parenté, des relations amicales ou professionnelles. Toutes ces relations comportent une dimension de confiance qui semble également transposable à notre étude des contributeurs de données collaboratives et qui serait intéressante d'exploiter.

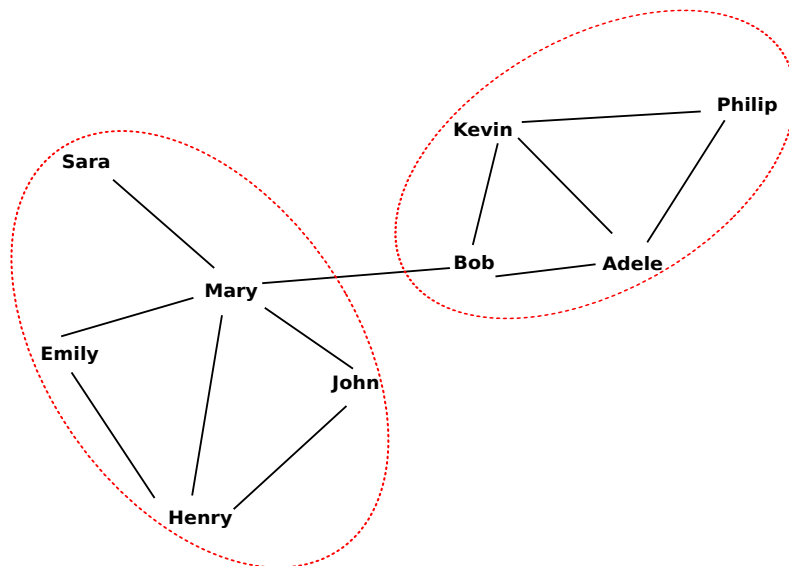


FIGURE II.13. Deux communautés détectées dans un graphe social de type Facebook.

Nous avons vu que la détection de communautés devait considérer la structure des arcs du graphe, en maximisant le nombre d'arcs dans les communautés et en minimisant les arcs entre ces communautés. Or, dans le cas où un graphe ne présente pas une structure de communauté, le regroupement des sommets en communautés n'aura pas beaucoup de sens. La modularité Q est un indicateur qui évalue la détection de communautés dans un graphe. Elle mesure si la répartition des sommets

dans les communautés respecte la structure du graphe (Clauset *et al.*, 2004). Soit un graphe G décrit par la matrice d'adjacence (a_{ij}) . Sa modularité Q est définie par :

$$Q = \frac{1}{2m} \sum_{i,j} \left(a_{ij} - \frac{k_i k_j}{2m} \right) \delta(c_i, c_j) \quad (\text{II.11})$$

avec a_{ij} l'élément de la matrice d'adjacence du graphe G ; c_i et c_j les communautés des sommets i et j respectivement; $\delta(c_i, c_j)$ une fonction valant 1 si $i = j$ et 0 sinon; m le nombre d'arc du graphe; k_i la centralité de degré du sommet i ; et $\frac{k_i k_j}{2m}$ la probabilité qu'il existe un arc entre les sommets i et j si les connexions sont aléatoires tout en respectant le degré des sommets du graphe.

Clauset *et al.* (2004) ont développé un algorithme basé sur la maximisation de la modularité Q , et affirment que la détection de communautés dans un graphe est significative dès lors que $Q > 0.3$. Plus récemment, Blondel *et al.* (2008) ont développé l'algorithme de Louvain, plus performant pour détecter les communautés dans un graphe en se basant également sur l'optimisation de la modularité. L'algorithme de Louvain est implémenté dans le logiciel *open source* Gephi, et a été utilisé pour détecter des communautés dans le graphe de normes de Wikipédia (Heaberlin et DeDeo, 2016).

Un réseau multiplexe peut être composée de plus d'une couche, c'est-à-dire de plusieurs graphes. Il existe trois approches possibles pour détecter des communautés dans un réseau multiplexe (Hmimida et Kanawati, 2015) : par agrégation des couches du réseau multiplexe, par agrégation successive des partitions et par exploration simultanée des couches.

(i) Approche d'agrégation de couches : Il s'agit d'agréger les différentes couches d'un réseau multiplexe, avant d'effectuer la détection de communautés sur le graphe agrégé obtenu (Figure II.14). Le graphe agrégé peut être orienté ou non, pondéré ou non, selon les critères fixés pour agréger les couches. Par exemple, l'agrégation de couches peut simplement consister à relier par un arc, dans un graphe agrégé non-pondéré, deux sommets qui sont reliés dans au moins n couches ($n \in \mathbb{N}^*$) du réseau multiplexe. Le seuil n permet de représenter dans le graphe agrégé les connexions les plus significatives dans le réseau multiplexe. La construction d'un graphe agrégé pondéré est également un moyen de quantifier l'intensité des connexions entre les sommets dans les différentes couches du réseau multiplexe. La pondération peut suivre différentes règles : il peut s'agir de pondérer les arcs du graphe agrégé par la moyenne ou la somme des poids observés dans chaque couche, ou bien par le nombre de couches dans lesquelles les arcs apparaissent.

La détection de communautés sur le graphe agrégé s'effectue alors en y appliquant une méthode de regroupement (*clustering*), telle que, pour en citer quelques-unes, la classification ascendante hiérarchique ou la méthode des k -moyennes. Puis, il s'agit de projeter les communautés détectées du graphe agrégé vers le réseau multiplexe. Sur ce sujet, Berlingerio *et al.* (2011) proposent une fonction de restauration des connexions entre les sommets d'un graphe agrégé vers un réseau multiplexe, afin de pouvoir évaluer la détection des communautés dans le réseau multiplexe.

L'approche par agrégation de couches permet de résoudre un problème mono-

plexe classique. De plus, la construction du graphe agrégé permet de fixer l'intensité des connexions à étudier. Cependant, l'inconvénient de l'agrégation de couches est la perte d'information quant à la nature hétérogène des arcs. De plus, comme il a été soulevé précédemment, il faut veiller à ce que l'agrégation de couches ait un sens, car les poids des arcs ne sont pas tous exprimés dans la même unité de mesure.

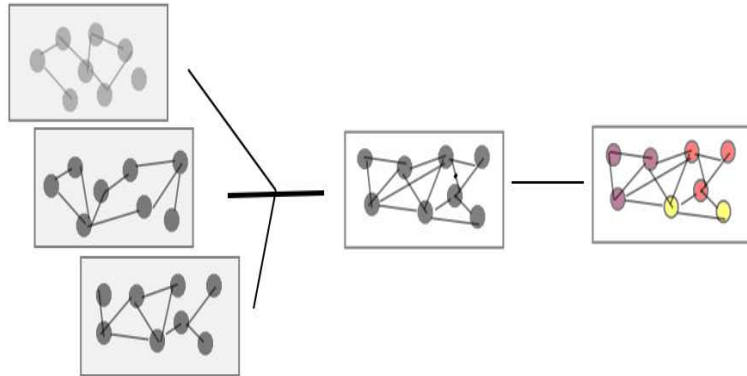


FIGURE II.14. Agrégation des couches d'un réseau multiplexe, tirée de *Hmimida et Kanawati (2015)*.

(ii) Approche par agrégation successive des partitions : Cette approche consiste à détecter des communautés dans chacune des couches du réseau multiplexe, puis à construire un graphe de consensus à partir des différentes communautés obtenues dans les différentes couches (Figure II.15). Le graphe consensus s'obtient en lançant à nouveau un algorithme de partitionnement sur les différentes communautés détectées, ou bien en reliant les sommets qui appartiennent à une même communauté.

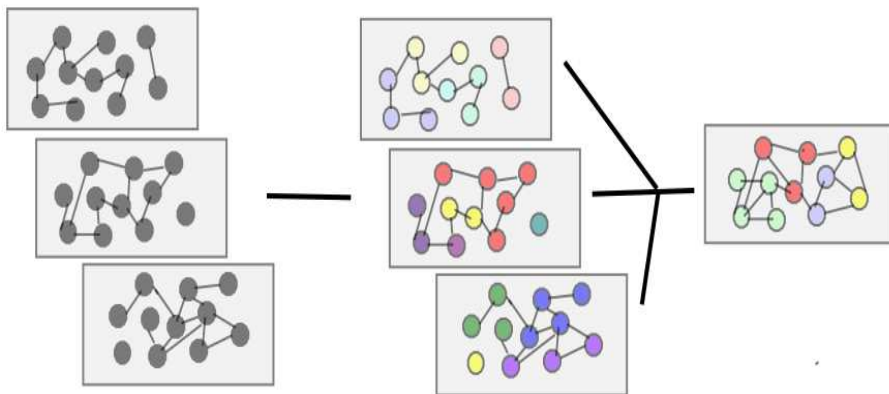


FIGURE II.15. Agrégation des communautés d'un réseau multiplexe, tirée de *Hmimida et Kanawati (2015)*.

Cette approche présente également l'avantage de revenir à un problème plus simple de détection de communautés dans le cas monoplexe. De plus, il permet de détecter des communautés intermédiaires par couche, qui ont chacune un sens. L'inconvénient ici est d'agréger des communautés qui peuvent être très différentes, au risque de regrouper des sommets qui ne partagent aucun lien.

(iii) **Approche « centrée graine »** : Cette approche ascendante consiste à sélectionner des nœuds (des « graines ») autour desquels des communautés locales peuvent être détectées. La sélection de ces graines peut se baser sur la centralité de degré multiplexe, qui est une mesure d'entropie semblable au coefficient de participation multiplexe (Battiston *et al.*, 2014). En effet, il s'agit de détecter des communautés en explorant simultanément les couches du réseau multiplexe. Issam *et al.* (2015) ainsi que Hmimida et Kanawati (2015) ont proposé une méthode d'agrégation des préférences qui consiste à calculer, pour chaque sommet, un vecteur de préférences d'appartenir à la communauté de chaque graine de départ.

2.2 Modélisation des collaborations

Dans cette section, nous proposons des modèles théoriques de graphes d'interactions qui peuvent être intégrés comme des couches dans un réseau multiplexe de contributeurs d'information géographique volontaire.

a) Couches d'interactions basées sur les opérations d'édition cartographique

(i) **Graphes d'édition** : Les graphes présentés dans la Section 1.3 de ce chapitre sont tous basés sur l'historique des éditions d'objets cartographiques. Leurs modélisations étant différentes, chacun de ces graphes peut constituer une couche à intégrer dans un réseau multiplexe de collaboration. Nous proposons ici une définition formelle de ces graphes.

Soient i et j deux contributeurs d'information géographique volontaire. Si le contributeur i a édité les versions $v + 1$ de w objets tels que les versions v de ces objets ont été saisies par le contributeur j alors ces contributeurs sont connectés dans le graphe de co-édition $G_{co-edition}$ (Mooney et Corcoran, 2012a) :

$$(i, j, w) \in E_{co-edition} \quad (\text{II.12})$$

Si le contributeur i a édité une version $v + 1$ de w objets dont une version $v - n$ ($n \in \mathbb{N}$) précédente a été saisie par le contributeur j alors ces contributeurs sont connectés dans le graphe de largeur de collaboration $G_{largeur}$ (Stein *et al.*, 2015) :

$$(i, j, w) \in E_{largeur} \quad (\text{II.13})$$

Par ailleurs, les contributeurs i et j sont également connectés dans le graphe de profondeur de collaboration $G_{profondeur}$. Parmi les objets cartographiques sur lesquels le contributeur i a édité à la suite du contributeur j , soit w' le nombre maximal d'éditions de i qui ont été produits à la suite (de manière directe ou indirecte) des éditions de j , alors :

$$(i, j, w') \in E_{profondeur} \quad (\text{II.14})$$

Nous pouvons remarquer que si $(i, j, w) \in E_{largeur}$ alors :

$$\exists w' | (i, j, w') \in E_{profondeur} \quad (\text{II.15})$$

La réciproque est vraie.

Si les contributeurs i et j ont contribué sur un même objet cartographique alors ils sont connectés dans le graphe de co-contribution (Ma *et al.*, 2015). Le graphe de co-contribution étant non-pondéré et non-orienté, on notera $(i, j, 1) = (i, j) = (j, i)$, tel que :

$$(i, j) \in E_{co-contribution} \quad (\text{II.16})$$

Des graphes supplémentaires peuvent être modélisés en considérant les interactions issues des opérations élémentaires présentées dans la Section 1.3 de ce chapitre.

(ii) Graphe d'utilisation : Le réseau de norme de Wikipédia (Heaberlin et DeDeo, 2016) est un graphe dont les articles Wikipédia constituent l'ensemble des sommets, et un arc orienté d'un article $A1$ vers un article $A2$ signifie qu'il existe un lien hypertexte dans l'article $A1$ faisant référence à l'article $A2$. Sur le même principe, si le contributeur i utilise w contributions produites par le contributeur j pour saisir ses propres contributions, alors ces contributeurs sont connectés dans le graphe d'utilisation $G_{utilisation}$:

$$(i, j, w) \in E_{utilisation} \quad (\text{II.17})$$

La construction de ce graphe n'est possible que si le projet collaboratif permet aux contributeurs de réutiliser les contributions faites par d'autres, et si cette information est accessible. Dans le cas des données OSM, la réutilisation des objets est tout à fait possible, en particulier pour produire des objets cartographiques complexes à partir de primitives plus élémentaires. Par exemple, si le contributeur i cartographie dans OSM une ligne de bus à partir de 10 arrêts de bus qui ont été préalablement saisis par le contributeur j , alors ce dernier semblera plus fiable qu'un contributeur k dont un seul nœud a été utilisé pour créer une intersection de route. Le graphe d'utilisation permet de considérer les contributeurs dont les données sont réutilisées comme des contributeurs de confiance. Le niveau de confiance peut, dans ce cas, être une fonction du nombre de contributions réutilisées.

(iii) Graphe de suppression : Si le contributeur i supprime w objets cartographiques dont le contributeur j était l'auteur de la dernière version, alors ces contributeurs sont connectés dans le graphe de suppression $G_{suppression}$:

$$(i, j, w) \in E_{suppression} \quad (\text{II.18})$$

Le graphe de suppression est intéressant pour plusieurs raisons. D'une part, il permet de mettre en évidence les contributeurs dont les contributions sont supprimées, car celles-ci dégradent potentiellement l'espace cartographique, ce qui correspondrait à l'une des trois composantes d'un acte de carto-vandalisme. D'autre part, il mettra également en évidence les contributeurs qui participent activement selon ce mode d'édition, soit dans le but de modérer les contributions d'autrui, soit dans le but de dégrader l'espace cartographique. Dans le premier cas, ce profil se rapprocherait de celui d'un modérateur, alors que dans le second, il se rapprocherait de celui d'un carto-vandale.

b) Couches d'interactions basées sur les éditions attributaires

Le réseau social multiplexe introduit par Kazienko *et al.* (2011) pour analyser les utilisateurs de Flickr contient une couche qui modélise les relations sociales à partir de l'usage des tags. Dans cette couche, deux contributeurs sont reliés dès qu'ils utilisent un tag en commun pour décrire leurs photographies. Dans le cas du projet OSM, il est également possible d'éditer des tags sur les objets cartographiques. Pour rappel, un tag sur OSM est un attribut de forme clé-valeur. On peut alors considérer deux types d'interactions basées sur l'édition attributaire :

1. Deux contributeurs interagissent dès qu'ils éditent la même clé de tag ;
2. Deux contributeurs interagissent dès qu'ils éditent le même tag, c'est-à-dire qu'ils contribuent la même clé et la même valeur de tag.

Dans les deux cas, il n'est pas nécessaire que l'édition attributaire porte sur le même objet cartographique. Dans le premier cas, les contributeurs qui éditent les mêmes clés de tag semblent suivre le même standard de description des objets cartographiques. Dans le second, les contributeurs qui éditent les mêmes tags forment d'une certaine manière un groupe de spécialistes dans une action donnée.

On pourra, par exemple, relier tous les contributeurs qui participent en éditant le tag `building=yes`. Ainsi, ces contributeurs formeront une communauté caractérisée par un comportement spécifique d'édition. En particulier, la cartographie sur OSM des bâtiments en France est généralement réalisée par des scripts automatiques qui labélisent les objets par le tag `building=yes`. Ce modèle de graphe permettrait alors d'identifier les contributeurs *bots* dont le mode de participation est de cartographier automatiquement des bâtiments sur OSM.

Soit $G_{key-based}$ un graphe de relation basé sur les clés de tag. Soit NK_{ij} le nombre de clés de tag que les contributeurs i et j ont édité en commun, et soit NK_i le nombre de clés de tag uniques éditées par le contributeur i . Si $NK_{ij} \neq 0$ alors $(i, j, w) \in E_{key-based}$ tel que :

$$w_{key-based} = \frac{NK_{ij}}{NK_i} \quad (\text{II.19})$$

De même, soit $G_{tag-based}$ un graphe de relation basé sur les tags. Soit NT_{ij} le nombre de tags que les contributeurs i et j ont édité en commun, et soit NT_i le nombre de tags uniques édités par le contributeur i . Si $NT_{ij} \neq 0$ alors $(i, j, w) \in E_{tag-based}$ tel que :

$$w_{tag-based} = \frac{NT_{ij}}{NT_i} \quad (\text{II.20})$$

Précisons que les graphes $G_{key-based}$ et $G_{tag-based}$ ainsi modélisés sont pondérés et orientés.

c) Couches d'interactions basées sur les éditions géométriques

De même que les éditions attributaires, les éditions géométriques permettent de révéler des informations sur le style des contributeurs dans leur mode de contribution. Nous considérons un graphe de granularité spatiale qui relie des contribu-

teurs dont les contributions font augmenter la granularité spatiale des objets cartographiques, par exemple en ajoutant des nœuds à un tronçon de route (Touya et Brando-Escobar, 2013). Soit $G_{granularite}$ un graphe de granularité; $nGeom_i$ le nombre d'éditions géométriques faites par le contributeur i qui sont des lignes ou des polygones⁵; $nodes_{ij}$ le nombre de nœuds ajoutés sur la géométrie des contributions de l'utilisateur j . Le graphe $G_{granularite}$ est orienté. Si $nGeom_i \neq 0$ alors $(i, j, w_{granularite}) \in E_{granularite}$ tel que :

$$w_{granularite} = \frac{nodes_{ij}}{nGeom_i} \quad (II.21)$$

Dans la même idée, le concept de maturation d'un objet cartographique consiste à créer un nouvel objet cartographique d'un objet géographique mais sous une forme plus complexe (Maguire et Tomko, 2017). Par exemple, une ville initialement cartographiée par un point aura une représentation plus « mature » sous la forme d'un polygone fermé correspondant aux frontières de cette ville. La Figure II.16 illustre la maturation de la cartographie de l'Arc de Triomphe d'un point vers un polygone. Soit $G_{maturation}$ un graphe de maturation. On note $nTagName_i$ le nombre d'éditions du contributeur i dans lesquelles celui-ci ajoute le tag `name` à un polygone, et $nMaturation_{ij}$ le nombre de nœuds du contributeur j portant le tag `name` qui ont servi au contributeur i à nommer les polygones par cet attribut. Si $nTagName_i \neq 0$ alors $(i, j, w_{maturation}) \in E_{maturation}$ tel que :

$$w_{maturation} = \frac{nMaturation_{ij}}{nTagName_i} \quad (II.22)$$



(a) Date t_1 : l'Arc de Triomphe est représenté par un point dans un polygone (b) Date t_2 : l'Arc de Triomphe est représenté par le polygone lui-même

FIGURE II.16. Maturation de la cartographie de l'Arc de Triomphe.

d) Couches d'interactions spatio-temporelles

Généralement, les projets d'information géographique volontaire ne demandent pas aux contributeurs d'indiquer les relations sociales qu'ils partagent réellement ou virtuellement. Toutefois, pour évaluer le niveau de confiance des contributeurs,

5. Dans OSM, les lignes et les polygones sont des chemins (*ways*) respectivement ouverts et fermés.

la connaissance de ces relations sociales pourrait s'avérer utile. Le modèle de co-occurrence spatio-temporelle a été proposé initialement pour inférer les liens sociaux qui peuvent exister entre les utilisateurs de la plateforme de partage d'images Flickr (Crandall *et al.*, 2010). La méthode repose sur l'hypothèse selon laquelle deux personnes qui ont pris une photographie au même endroit et au même moment se connaissent potentiellement. En effet, la modélisation des liens de co-occurrence se fait à partir d'un découpage régulier de l'espace géographique et temporel. Ainsi, deux utilisateurs qui ont pris une photographie qui se situe dans la même cellule spatiale et le même intervalle de temps partagent un lien de co-occurrence. Les utilisateurs de Flickr peuvent partager leurs images en y indiquant la date et le lieu de la prise de vue. Pour ceux dont les photographies sont prises par des *smartphones*, les images peuvent être géolocalisées et datées grâce aux différents récepteurs dont leurs appareils sont dotés (horloge numérique, GPS). Par conséquent, ces informations permettent de restituer précisément les liens de co-occurrence entre les utilisateurs.

Certains projets d'information géographique volontaire comportent des campagnes de saisie dans lesquels les contributeurs peuvent se rendre sur le terrain pour éditer des données ensemble : c'est le cas du projet OpenStreetMap qui organise souvent des carto-parties ; dans le cadre du projet LandSense⁶, les contributeurs étaient également invités à participer sur le terrain. Dans ces situations, la modélisation des co-occurrences pourrait refléter des interactions significatives quant aux relations réelles des contributeurs.

En revanche, pour les contributeurs d'information géographique, la saisie des données ne se fait pas nécessairement sur le terrain, et un contributeur peut tout à fait éditer des données géographiques sur une région du monde dans laquelle il ne se trouve pas. De plus, la date à laquelle les données cartographiques sont chargées dans la base de données ne correspond pas toujours au moment de la collecte des données : dans ce cas, modéliser une relation de co-occurrence ne serait pas conforme à la réalité. La modélisation des interactions à partir de l'étude spatio-temporelle des contributions est donc discutable. La Figure II.17 schématise les différentes situations de saisie des données dans OSM, qui peuvent entraîner une différence entre la date de saisie des données et de chargement dans la base, où *date_chargement* est la date renseignée dans la base. Les relations de co-occurrence sont donc modélisables dans les situations où la collecte a été réalisée sur le terrain et à l'aide d'outils de chargement automatique dans la base. L'utilisation des métadonnées sur ces contributions et/ou sur les sessions de modifications (*changesets*) est un moyen d'accéder au contexte de saisie des données et de filtrer les contributions qui correspondent à cette situation.

Mis à part dans les événements de carto-parties, il peut toutefois exister une forme de collaboration implicite entre les contributeurs – quelle que soit leur position géographique – qui contribuent sur la même zone cartographique. De même, des contributeurs qui saisissent leurs contributions dans la même période temporelle collaborent implicitement à l'actualisation de la base de données collaborative. Nous proposons donc de modéliser les interactions spatio-temporelles à partir de deux graphes, l'un portant sur les co-occurrences spatiales et l'autre sur les co-occurrences temporelles.

6. landsense.eu

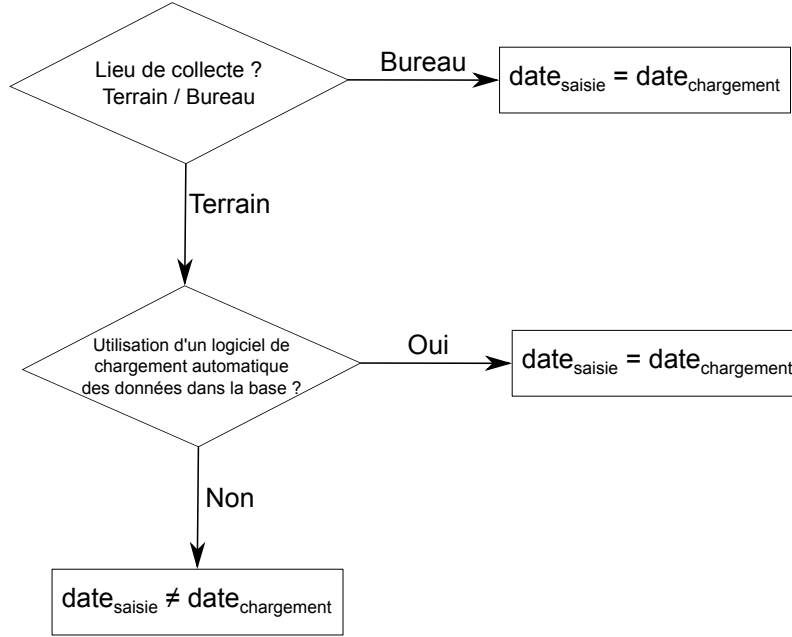


FIGURE II.17. Temporalité dans la saisie des données OSM

(i) **Graphe de co-temporalité** : Soit $E_{co-temporalite} = (i, j, w_{co-temporalite})$ l'ensemble des arcs du graphe de co-temporalité. $E_{co-temporalite}$ est l'ensemble des interactions entre les contributeurs dont les activités cartographiques se recouvrent. Soit $schedule_i$ les périodes temporelles durant lesquelles le contributeur i a édité des données. Dans OSM, les sessions d'éditations sont renseignées par une date de début et de fin de saisie. Les périodes temporelles peuvent donc être obtenues à partir de ces informations. Dans le cas où ces informations ne sont pas accessibles directement, il est possible d'obtenir les période temporelles à partir des dates des contribution, en regroupant celles qui sont temporellement proches, par exemple. En considérant les périodes temporelles de contributions $schedule_j$ du contributeur j , on définit $w_{co-temporalite}$ le recouvrement relatif des deux séries temporelles $schedule_i$ et $schedule_j$. Le recouvrement $w_{co-temporalite}$ s'obtient à partir d'une fonction $totalduration(.)$ qui présente les propriétés suivantes :

$$schedule_i \cap schedule_j = \emptyset \Rightarrow totalduration(schedule_i \cap schedule_j) = 0 \quad (II.23)$$

$$totalduration(schedule_i \cup schedule_j) = totalduration(schedule_i) + totalduration(schedule_j) \quad (II.24)$$

$$w_{co-temporal} = \frac{totalduration(schedule_i \cap schedule_j)}{totalduration(schedule_i \cup schedule_j)} \quad (II.25)$$

La Figure II.18 illustre l'obtention d'un graphe de co-temporalité à partir du recouvrement des périodes d'édition de trois contributeurs. Dans le cas où il existe des contributeurs qui feraient usage d'outils automatiques pour éditer massivement des données de manière continue, ceux-ci se retrouveront certainement connectés à un très grand nombre de contributeurs dans un graphe de co-temporalité. L'identification de tels contributeurs peut être intéressante dans le cas où des contributeurs provoquent du carto-vandalisme en masse en utilisant ces outils automatiques.

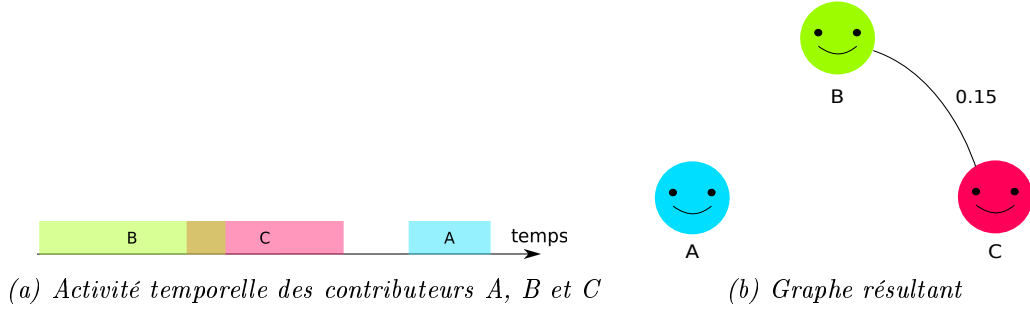


FIGURE II.18. Graphe de co-temporalité

(ii) **Graphe de co-location** : Soit $E_{co-location} = (i, j, w_{co-location})$ l'ensemble des arcs du graphe de co-location. $E_{co-location}$ est l'ensemble des interactions entre les contributeurs dont les zones cartographiques se recouvrent. Soit $ActivityArea_i$ l'ensemble des zones cartographiques sur lesquelles le contributeur i a produit ses éditions. Sur ce sujet, plusieurs méthodes ont été développées pour déterminer la zone d'activité des contributeurs (Neis et Zipf, 2012; Zielstra *et al.*, 2014).

Soit $S(\cdot)$ une fonction qui renvoie l'aire totale d'un ensemble de zones spatiales. $S(\cdot)$ possède les mêmes propriétés que la fonction *totalduration* définie précédemment, et permet d'obtenir le recouvrement total relatif $w_{co-location}$ entre deux ensembles d'activité spatiale :

$$w_{co-location} = \frac{S(ActivityArea_i \cap ActivityArea_j)}{S(ActivityArea_i \cup ActivityArea_j)} \quad (II.26)$$

La Figure II.19 illustre l'obtention d'un graphe de co-location à partir du recouvrement des zones d'activité de trois contributeurs A, B et C.

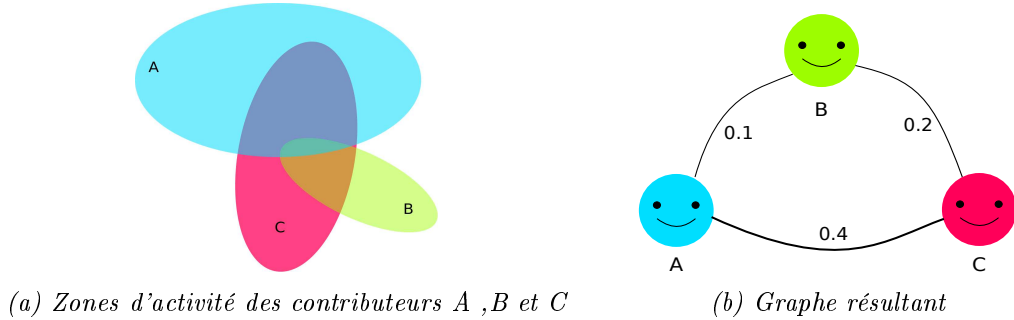


FIGURE II.19. Graphe de co-location

Cas des carto-parties : Pour faire apparaître les collaborations entre les contributeurs qui participent à des événements de type carto-partie, il faut redéfinir $w_{co-temporalite}$ et $w_{co-location}$. Soient deux fonctions sim_T et sim_S qui renvoient respectivement un score de similarité entre deux intervalles temporels et deux zones d'activités. La fonction sim_T compare les dates de début $schedule_i$ et les dates de fin $schedule_j$ des sessions d'éditions entre les contributeurs i et j . Plus les différences sont petites, plus leurs sessions d'édition sont similaires temporellement.

$$w_{co-temporalite} = \max_{\substack{d_i \in schedule_i \\ d_j \in schedule_j}} (sim_T(d_i, d_j)) \quad (II.27)$$

La fonction sim_S étant une extension en deux dimensions de la fonction sim_T , le poids $w_{co-occurrence}$ se définit de manière analogue à l'équation II.27.

2.3 Intérêts, limites du modèle et précautions

Les réflexions menées sur la modélisation des différents types d'interactions entre les contributeurs montrent que le comportement collaboratif est un objet d'étude très complexe qui peut se réaliser sur plusieurs plans. En synthétisant les différents graphes d'interactions dans un réseau multiplexe, nous espérons représenter dans un seul modèle et de manière exhaustive les relations de collaboration entre les contributeurs. Le premier intérêt est de pallier les insuffisances des graphes d'interactions lorsqu'ils sont pris en compte séparément, et d'exploiter les complémentarités qui peuvent exister entre les différentes couches du réseau multiplexe. Par exemple, si deux contributeurs ne sont pas reliés dans un graphe de co-contribution alors qu'ils le sont dans un graphe de co-location et de co-temporalité, cela peut révéler une forme d'accord entre les deux contributeurs sur les contributions de l'un et de l'autre. Au contraire, s'ils co-contribuent sur les mêmes objets cartographiques tout en étant reliés sur les graphes d'interactions spatio-temporelles, les graphes de largeur et de profondeur de collaboration peuvent apporter des informations supplémentaires quant à leur type de collaboration. Par exemple, ils permettront de comprendre si la collaboration consiste en des améliorations de la part de l'un sur les contributions de l'autre, ou bien si elle représente des désaccords entre ces deux contributeurs.

Ce modèle repose sur la modélisation des interactions issues des contributions cartographiques. Par conséquent, les interactions qui peuvent se produire en dehors de l'activité cartographique ne seront pas prises en compte dans le réseau multiplexe. Par exemple, certains contributeurs peuvent interagir à travers des outils de discussions (forums, mailing lists, commentaires, etc.). Ces médias sont intéressants car ils peuvent parfois indiquer explicitement les désaccords entre certains contributeurs. Par exemple, les contributeurs d'OSM peuvent faire des commentaires sur les sessions d'édition des uns et des autres, en particulier lorsque ces sessions contiennent des contributions qui dégradent l'espace cartographique. L'étude des interactions issues de ces commentaires peut être utile pour détecter les potentiels contributeurs de carto-vandalisme. Cela est réalisable par des moyens de traitement automatique du langage naturel, cependant nous n'avons pas envisagé cette piste dans le cadre de cette thèse.

Alors que les couches d'un réseau multiplexe présentent potentiellement des complémentarités, celles-ci peuvent également apporter des informations redondantes pouvant entraîner des biais par l'exagération de certaines collaborations. Il faut pouvoir identifier ces redondances pour les interpréter et équilibrer la composition des couches du réseau. Selon la définition de chaque graphe, l'ensemble des arcs pourra être pondéré par des valeurs qui ne seront pas forcément exprimées dans les mêmes unités. Par exemple, le graphe de co-temporalité pourra être pondéré en unités temporelles, alors que le graphe de co-location sera pondéré en unités spatiales. Toutefois, il est également possible de normaliser les poids des arcs pour pouvoir comparer les interactions provenant de différentes couches.

Certain modèles de graphes sont inspirés de plateformes qui ne sont pas tout à fait

des projets d'information géographique volontaire, comme Flickr par exemple. En conséquence, ces modèles mériteraient d'être réadaptés pour mieux servir à l'étude des collaborations entre contributeurs d'information géographique volontaire. Enfin, l'utilité de chacune des couches dans le réseau multiplexe nécessite d'être vérifiée expérimentalement, afin d'établir une configuration optimale du réseau dans le cadre de l'étude des collaborations entre contributeurs. Dans la suite de cet exposé, nous cherchons à démontrer expérimentalement l'intérêt d'utiliser un modèle multiplexe pour qualifier les contributeurs du projet OpenStreetMap.

3 Qualification des contributeurs du projet OpenStreetMap

Nous avons choisi de construire un réseau multiplexe à partir des données cartographiques collaboratives du projet OpenStreetMap, pour étudier les interactions qui se produisent entre les contributeurs de ce projet et ainsi déduire des profils particuliers de comportement. L'objectif est de montrer expérimentalement l'intérêt d'utiliser le modèle de réseau multiplexe pour qualifier les contributeurs, en particulier leur niveau de confiance. Notre expérience porte principalement sur les données OSM du quartier de l'Île de la Cité, à Paris, entre le 1^{er} janvier 2013 et le 31 décembre 2015. Des expériences additionnelles sont également menées à Stuhr (Allemagne) et à Kathmandou (Népal) pour démontrer la généralité de notre méthode.

Préparation des données OSM

Nous développons ici la phase de préparation des données OSM, qui a été nécessaire pour toutes les expériences présentées dans cette thèse : nous décrivons les étapes de stockage de l'historique des données OSM et le principe de chargement des données utilisées pour nos expériences. Toutefois, la lecture de ce cadre n'est pas nécessaire à la bonne compréhension de cette Section 3.

a) Récupération des fichiers historiques

Nous avons vu que la qualification du carto-vandalisme repose souvent sur l'étude de l'historique des objets cartographiques. De plus, l'étude des comportements des contributeurs peut se baser sur leurs interactions dans l'historique des objets cartographiques (voir Section 1.3). Il est possible d'obtenir l'historique de la base de données OSM de la planète entière ^a sous forme d'un fichier binaire (format `.pbf`). Toutefois, ce fichier étant assez lourd (sa taille actuelle est de 120 GB), nous choisissons de travailler sur des petites zones géographiques.

Le site Geofabrik ^b permet de télécharger des fichiers historiques de données OSM extraites par zone géographique (continent, pays ou région). Cependant, depuis 2018, les fichiers historiques contenant des données personnelles ne sont plus accessibles publiquement. En effet, les données historiques qui nous intéressent contiennent des informations personnelles sur les contributeurs, telles que leurs pseudonymes, l'identifiant correspondant, ainsi que les identifiants

des sessions d'éditions. L'usage de ces données étant soumis à la politique de protection des données de l'Union Européenne, le téléchargement de ces fichiers historiques est autorisé à condition d'avoir créé un compte OpenStreet-Map. Pour des raisons de confidentialité, tous les identifiants des contributeurs qui sont exposés dans ce travail sont anonymisés.

b) Traitement et stockage dans une base de données historiques

Les données OSM se composent de trois types de primitives : les nœuds (*nodes*), les chemins (*ways*) et les relations. Les fichiers historiques contiennent les métadonnées suivantes pour chaque contribution cartographique :

- son identifiant (*id*)
- son numéro de version (*version*)
- l'identifiant du contributeur (*uid*)
- le pseudonyme du contributeur (*username*)
- l'identifiant de la session d'édition (*changeset*)
- la date de modification (*datemodif*)
- les tags qui le décrivent (*tags*)
- sa visibilité sur la carte (*visible*)

Selon le type de primitive par lequel est représenté un objet OSM, des attributs supplémentaires compléteront ces métadonnées. Les nœuds sont les éléments de base du système OSM : ce sont des points dotés de coordonnées géographiques. Un objet OSM dont la primitive est un nœud aura deux attributs *lat* et *lon* indiquant respectivement sa latitude et sa longitude. Les chemins sont définis par une liste de nœuds : ils peuvent être ouverts ou fermés. Dans le cas d'un chemin fermé, le premier et le dernier nœud de la liste qui composent le chemin sont les mêmes. Un objet OSM dont la primitive est un chemin aura un attribut *iscomposedof* contenant la liste des identifiants des nœuds qui le composent. Une relation est une primitive complexe dont les membres peuvent être des nœuds et/ou des chemins, et/ou des relations. C'est pourquoi notre base de données contient une classe qui décrit les membres de chaque relation OSM (classe *relationmember*), dans laquelle sont précisés la nature de chacun des membres d'une relation : le type de primitive (*typemb*) et le rôle du membre (*rolemb*) de la relation (si, pour une surface qui possède un trou, il s'agit d'un anneau intérieur ou extérieur).

Nous construisons donc une base de données PostgreSQL destinée à stocker ces informations. Cette base de données est dotée de l'extension *hstore* pour contenir les tags sous forme de clé-valeur, et de l'extension *PostGIS* pour contenir la géométrie des contributions cartographiques. La Figure II.20 illustre le modèle conceptuel de la base de données. Le chargement des données historiques OSM s'effectue avec un script Python. Un tutoriel du chargement des données historiques dans une base de données PostgreSQL est disponible en ligne^c.

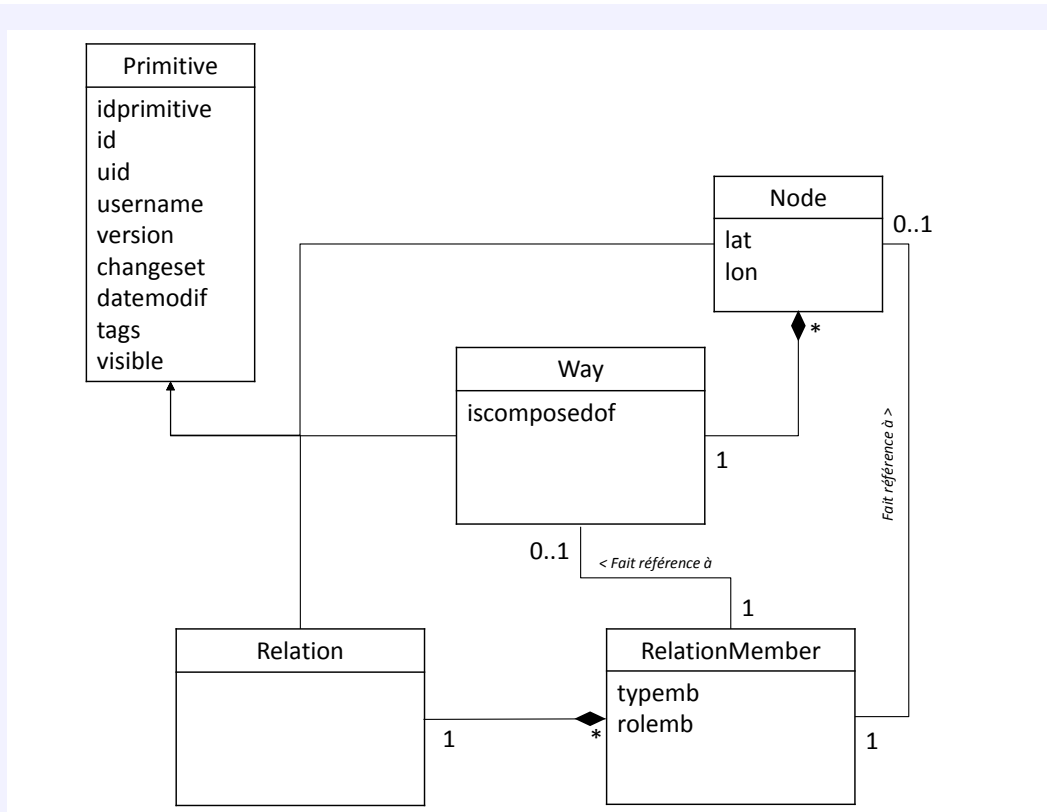


FIGURE II.20. Schéma conceptuel de la base de données historiques OSM

(i) **Reconstruction des géométries avec PostGIS :** Une colonne de géométrie est ajoutée à la classe `node`, dans laquelle les géométries des nœuds sont reconstruites à partir de leurs coordonnées géographiques. Pour reconstruire les géométries de la classe `way`, le traitement est différent. Rappelons qu'un chemin est référencé par la liste des identifiants des nœuds qui le composent. Le problème est donc de trouver, pour chaque nœud, la version valide à la date de la version du chemin. Par ailleurs, l'édition d'une nouvelle version d'un nœud composant un chemin, n'entraîne pas une nouvelle version du chemin. La géométrie d'un chemin, contrairement à un nœud, n'est donc pas fixée pour une version donnée, mais elle peut varier selon les versions de ses composants (les nœuds). Par conséquent, nous avons choisi de reconstruire approximativement l'emprise spatiale des chemins, en calculant la boîte englobante de chaque chemin, à partir des coordonnées extrêmes de ses composants à la date de modification du chemin.

c) Chargement des données selon une fenêtre spatio-temporelle

La construction des graphes d'interaction est réalisée sur la plateforme GeOxygene, qui permet de développer des applications SIG en langage Java (Bucher *et al.*, 2012). Nous avons, pour cela, développé une méthode de chargement des données historiques OSM d'une base de données PostgreSQL vers GeOxygene. La méthode de chargement des données historiques consiste à faire une requête spatio-temporelle sur la base de données, en indiquant une fenêtre spatiale et une fenêtre temporelle $[date_I, date_F]$. La sélection des données est

structurée en deux parties :

1. Un *snapshot* de la carte OSM sur la fenêtre spatiale au début de l'intervalle temporel : il s'agit de sélectionner la dernière version valide à la date $date_I$ des objets OSM visibles dont la géométrie est complètement contenue dans la fenêtre spatiale ;
2. L'historique des données du *snapshot* : il s'agit de charger les versions suivantes des contributions visibles dans la fenêtre spatiale à $date_I$ jusqu'à $date_F$.

Le schéma de la Figure II.21 résume la chaîne de traitement des données OSM pour parvenir à la construction de graphes sociaux. Nous avons également développé des méthodes de construction et d'analyse de réseau multiplexe sur GeOxygene.

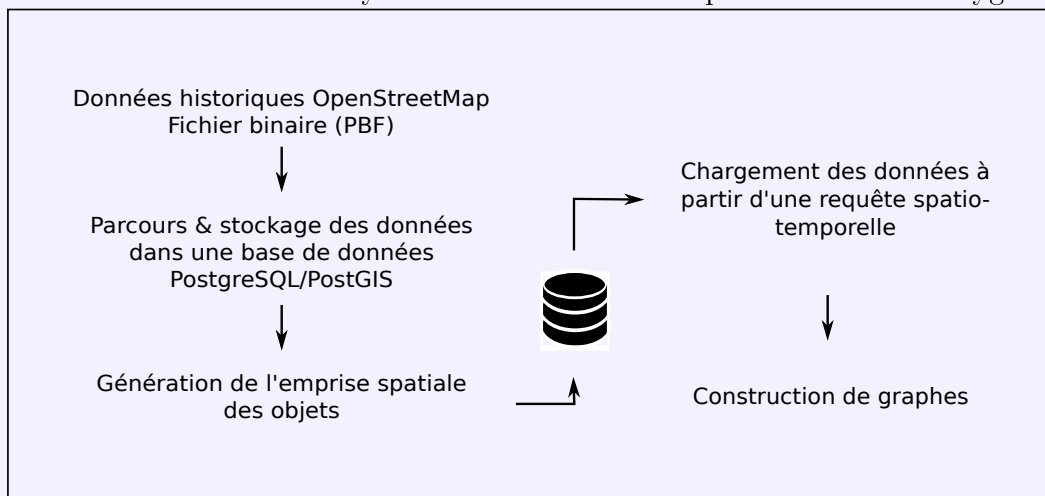


FIGURE II.21. Chaîne de traitement des données OSM

- a. <https://planet.osm.org/>
- b. <http://download.geofabrik.de/>
- c. <https://github.com/quythytruong/OSM-history-parser>

3.1 Implémentation d'un réseau social multiplexe : cas de l'Île de la Cité

a) Composition du réseau multiplexe

Dans notre expérimentation, nous considérons un réseau multiplexe des contributeurs OSM sur le quartier de l'Île de la Cité. Ce réseau, noté $R1$, est composé de cinq graphes d'interaction :

- un graphe de co-édition $G_{co-edition}$;
- un graphe de largeur de collaboration $G_{largeur}$;
- un graphe de profondeur de collaboration $G_{profondeur}$;
- un graphe de co-contribution $G_{co-contribution}$;
- un graphe d'utilisation $G_{utilisation}$.

b) Expériences

Pour détecter les communautés de contributeurs dans notre réseau multiplexe, nous avons adopté l’approche par agrégation de couches car elle permet de simplifier le problème en ne détectant les communautés qu’une seule fois, et d’appliquer des méthodes classiques de classification. Par ailleurs, [Kazienko *et al.* \(2011\)](#) ont également suivi une approche par agrégation de couches pour analyser les relations sociales entre les utilisateurs de Flickr dans un modèle de réseau social multidimensionnel.

Nous construisons un graphe agrégé $G_{agr} = (E_{agr}, V)$ à partir du réseau multiplexe $R1$. Pour conserver autant que possible l’intensité des interactions, le graphe agrégé est pondéré et non orienté tel que, s’il existe un arc (i, j, w) dans G_{agr} , alors l’arc (i, j) ou l’arc (j, i) existe dans w couches du réseau multiplexe $R1$. En d’autres termes, la pondération du graphe agrégé G_{agr} indique le nombre de graphes dans lesquels deux contributeurs interagissent.

La détection de communautés sur le graphe G_{agr} en utilisant l’algorithme de Louvain avec le logiciel Gephi a retourné 13 communautés avec une modularité $Q = 0.37$. La valeur de modularité ici trouvée est supérieure à la valeur seuil indiquant l’existence d’une structure communautaire dans un graphe. La Figure II.22 illustre les communautés détectées dans le graphe agrégé G_{agr} .

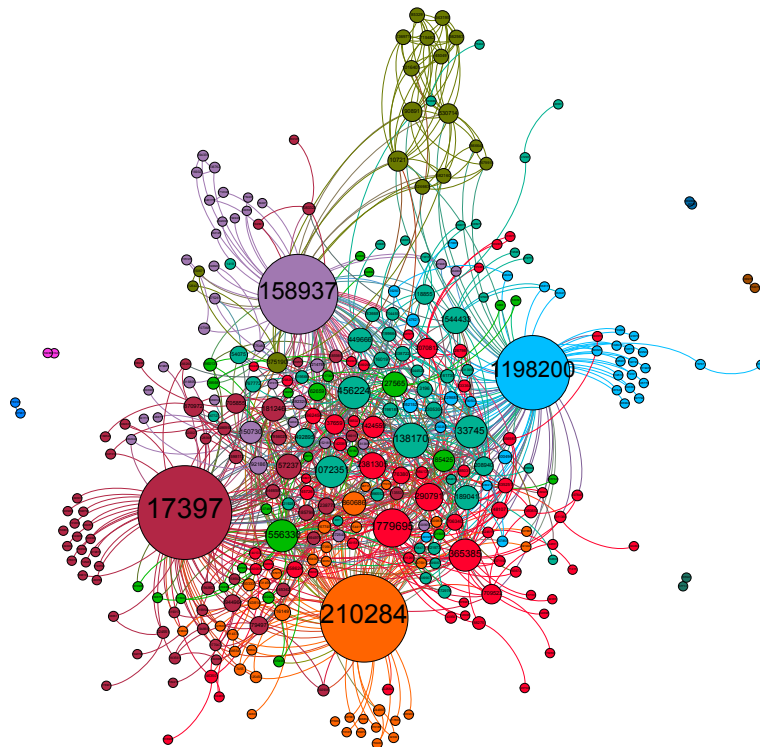


FIGURE II.22. *Communautés de contributeurs OSM détectées par l’algorithme de Louvain. Les sommets de même couleur appartiennent à la même communauté. La taille est proportionnelle au degré de degré. Les labels des sommets sont des identifiants anonymisés des contributeurs OSM.*

Parmi les 13 communautés ainsi détectées, nous observons cinq communautés de type « couple », qui sont constituées uniquement de deux sommets. L’observation

visuelle des huit communautés restantes nous permet de distinguer deux types de configurations : la première est vue comme communauté centrée sur un « noyau », dans laquelle des contributeurs sont généralement tous connectés à un seul contributeur sans être reliés entre eux (Figure II.23a). La deuxième configuration correspond à un groupe moins organisé (Figure II.23b).

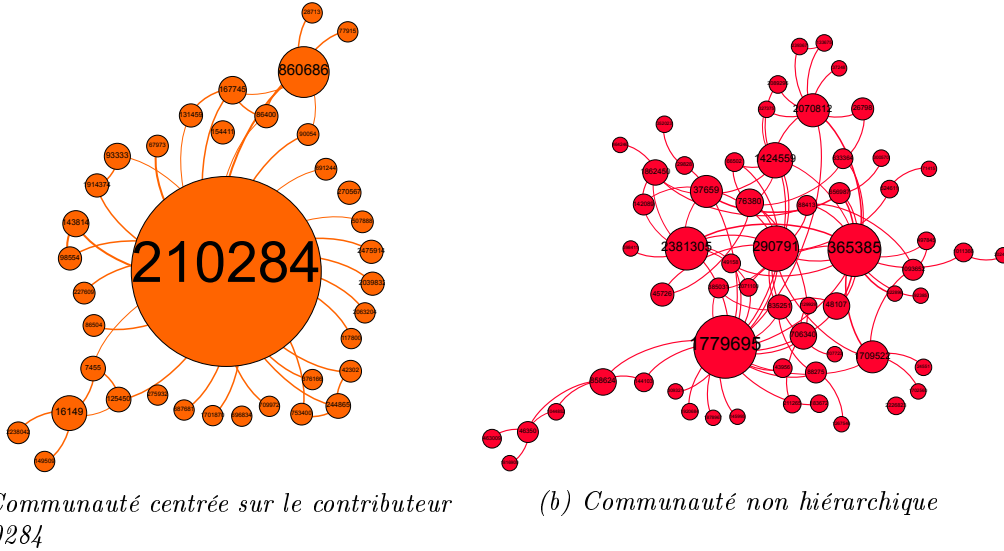


FIGURE II.23. Configurations caractéristiques dans les communautés détectées

Notre étude s’est principalement focalisée sur les contributeurs « noyaux » autour desquels s’organisent les communautés centrées. En effet, nous pouvons repérer quatre contributeurs « noyaux » dans la Figure II.22, dont les identifiants anonymisés sont #210284, #158937, #1198200 et #17397. L’analyse peut être menée sur différents niveaux, à savoir :

1. sur chaque couche individuellement ;
2. sur tout ou partie du réseau multiplexe ;
3. sur les communautés détectées à partir du graphe agrégé.

Selon les différents niveaux d’étude de ce réseau, nous pouvons identifier différents profils de contributeurs à l’aide des métriques de graphes définies dans la Section 2.1 de ce chapitre. L’analyse de ces profils consiste à étudier leurs caractéristiques et évaluer leur niveau de confiance qui peut se décliner en l’évaluation de la fiabilité et de la compétence des contributeurs, puisque ceux-ci peuvent être vus comme des sources d’information géographique (Revault d’Allonnes, 2013).

3.2 Identification de profils de contributeurs

a) Modérateur

Nous avons pu identifier des profils de type modérateur dans les trois plus gros « noyaux » du graphe agrégé G_{agr} . Pour ne pas alourdir la lecture, au lieu d’utiliser les identifiants anonymisés de ces contributeurs, nous désignons les trois modérateurs par $N1$, $N2$ et $N3$. Nous pouvons caractériser leurs profils de modérateur en étudiant les couches du réseau $R1$ et en analysant leurs contributions.

(i) **Analyse de la fiabilité :** La Figure II.24 présente la distribution des degrés entrants et sortants des contributeurs appartenant à une communauté centrée autour d'un modérateur dans les couches de co-édition et de collaboration⁷. Les graphes des Figures II.24a et II.24b ont été divisés en quatre cadrans : nous pouvons repérer le contributeur noyau de la communauté qui se trouve seul dans le deuxième cadran (en haut à droite), tandis que tous les autres sont concentrés dans le troisième cadran (en bas à gauche). Le contributeur modérateur se caractérise par un fort degré entrant et sortant dans les deux couches du réseau multiplexe, qui révèle une forte interaction avec les autres contributeurs.

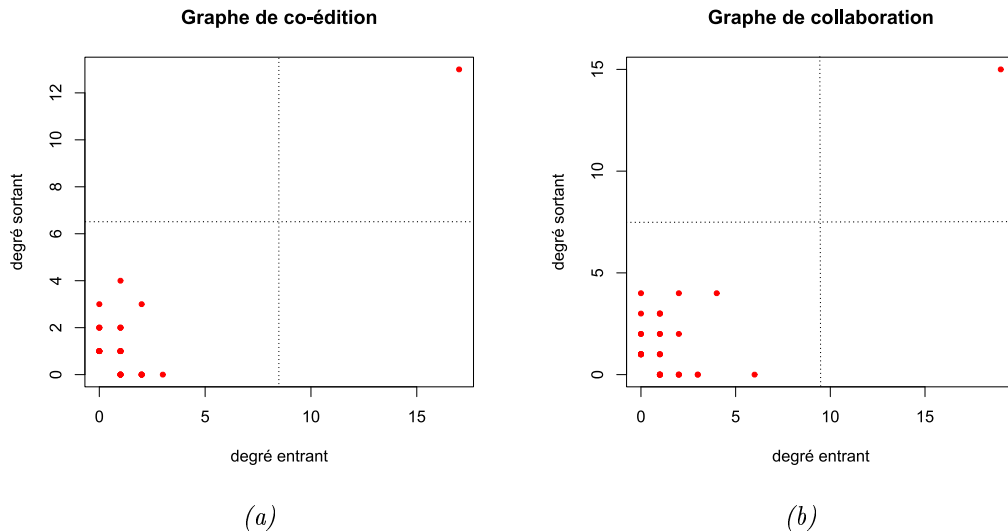


FIGURE II.24. Distribution des degrés entrants et sortants des contributeurs d'une communauté centrée sur un modérateur dans deux couches de $R1$

Au niveau de la couche d'utilisation $G_{utilisation}$, nous pouvons remarquer que les contributions de $N1$, $N2$ et $N3$ sont largement utilisées par rapport à la moyenne des contributeurs OSM de l'Île de la Cité (Table II.5). Cette forte utilisation de données montre que la communauté semble accorder une certaine confiance à leurs contributions. Les contributeurs $N1$, $N2$ et $N3$ sont donc vus comme fiables par la communauté des contributeurs OSM.

TABLE II.5. Degrés entrants des contributeurs noyaux de la couche d'utilisation $G_{utilisation}$

	Degré entrant	Degré sortant
Valeur min.	0	0
Médiane	1	2
Moyenne	3.343	3.343
Valeur max.	46	30
Contributeur $N1$	34	30
Contributeur $N2$	5	15
Contributeur $N3$	15	22
Contributeur $N4$	46	0

7. *N.B.* : En toute rigueur, la Figure II.24b est issue du graphe de largeur de collaboration mais d'après la propriété II.15, elle vaut aussi pour le graphe de profondeur de collaboration.

Par ailleurs, les contributeurs N_1 , N_2 et N_3 ont également un fort degré sortant sur la couche $G_{utilisation}$. Ces contributeurs ont donc une activité intense de contribution de données. De plus, la fiabilité de ces contributeurs entraîne que leurs contributions présentent aussi un certain niveau de confiance d'après le mécanisme de confiance informationnelle⁸ (Bishr et Janowicz, 2010).

(ii) Analyse de la compétence : Une manière d'évaluer la compétence des contributeurs consiste à déterminer si leur expertise vient du fait que leur profession est liée à la contribution d'information géographique ou s'ils sont plutôt des « amateurs experts » (Coleman *et al.*, 2010), c'est-à-dire que leur forte implication dans le projet collaboratif fait d'eux des contributeurs expérimentés et donc compétents. Ce type d'information personnelle n'est pas accessible dans le projet OSM, néanmoins, une analyse temporelle des contributions peut permettre de déduire cette information. Dans la Table II.6 renseigne la répartition temporelle de l'activité des contributeurs. Nous pouvons observer que 64% des contributions de N_1 sont faites en journées, et 90% de son activité s'effectue en semaine.

TABLE II.6. *Activité des contributeurs « noyaux »*

Contributeur	N_1	N_2	N_3	N_4
Journée - Soirée (%)	64 - 46	19 - 81	56 - 44	25 - 75
Semaine - Weekend (%)	90 - 10	72 - 28	25 - 75	63 - 37

Un professionnel aura tendance à contribuer durant la journée en semaine, tandis qu'un amateur contribuera plutôt pendant son temps libre. Selon cette hypothèse, nous pouvons remarquer que le contributeur N_1 se comporte comme un professionnel au sens où ses contributions se font principalement durant les journées en semaine. En revanche, l'activité du contributeur N_2 a principalement lieu en soir de semaine, et le contributeur N_3 est surtout actif le weekend. Par conséquent ces deux contributeurs ont un profil qui se rapproche plus de l'amateur expert.

(iii) Analyse des types de modération : Une exploration manuelle des sessions d'édition des contributeurs N_1 , N_2 et N_3 montre que la plupart de leurs modifications sont généralement des mises à jour de données automatiques, par exemple pour ajouter ou mettre à jour les valeurs de limitations de vitesse ou les adresses. Les éditions automatiques ne requièrent pas la présence active des contributeurs, mis à part au moment du lancement des scripts qui reste une opération manuelle. Par conséquent, notre interprétation de l'expertise des contributeurs peut devenir discutable. Une solution serait de ne considérer que les dates de début des sessions d'édition qui ont été lancées automatiquement. Pour cela, une étude plus approfondie sur les contributions serait donc nécessaire, afin de distinguer les éditions manuelles et automatiques. Cette étude sur les contributions permettrait alors d'évaluer la compétence des contributeurs de manière plus exacte.

8. Ce mécanisme est expliqué dans la partie 1.4 de ce chapitre.

b) Pionnier

Le noyau N_4 ne partage pas les mêmes caractéristiques que les autres noyaux qui ont été identifiés comme des modérateurs dans l'analyse de la partie a) précédente. En effet, l'analyse du comportement de ce contributeur dans le réseau multiplexe montre qu'il présente un profil de pionnier, c'est-à-dire un contributeur qui a initialement participé au projet en comblant le vide de l'espace cartographique. Ce profil peut être assimilé au « cartographe solitaire » que [Stein et al. \(2015\)](#) décrivent comme un contributeur qui produit un grand nombre de données sans interagir avec aucun membre de la communauté.

(i) **Analyse de la fiabilité :** En observant la distribution des degrés entrants et sortants des couches $G_{co\text{-}édition}$ et $G_{utilisation}$ (Figure II.25), nous localisons le noyau N_4 dans le dernier cadran (en bas à droite) de ces graphes. En particulier, ce contributeur est complètement inactif dans l'édition des données puisqu'il présente un degré sortant nul dans la couche $G_{co\text{-}édition}$ (Figure II.25a). En revanche, il est caractérisé par un degré entrant élevé dans la couche $G_{co\text{-}édition}$: cela signifie donc que de nombreux contributeurs éditent ses contributions.

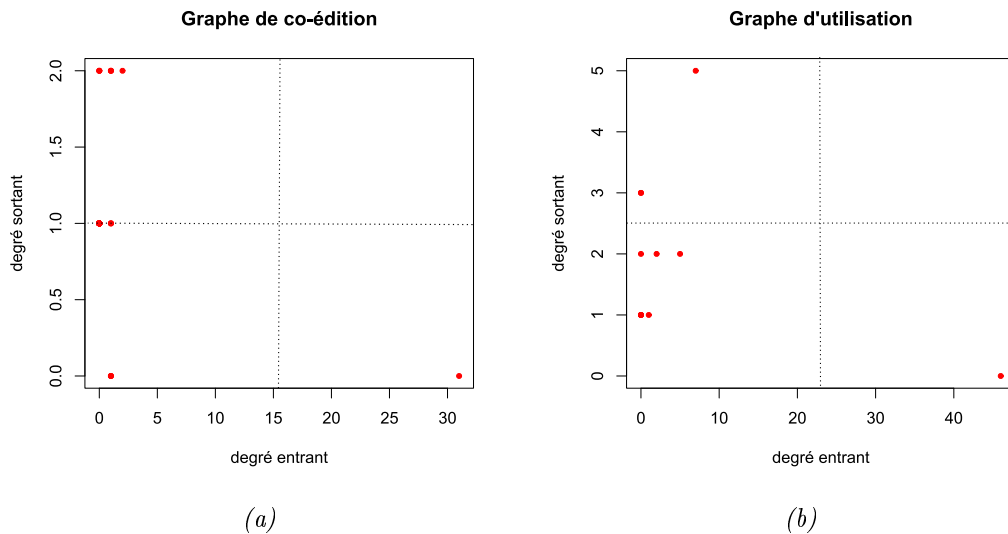


FIGURE II.25. Distribution des degrés entrants et sortants des contributeurs d'une communauté centrée sur un pionnier dans deux couches de R_1

Le faible degré sortant de N_4 dans la couche d'utilisation $G_{utilisation}$ (Figure II.25b) confirme l'inactivité de ce contributeur dans la zone spatio-temporelle étudiée. Cependant, son fort degré entrant – de valeur 46 – réalise la valeur maximale de la distribution du degré entrant dans la couche d'utilisation $G_{utilisation}$ (Table II.5), c'est-à-dire que 46 contributeurs différents ont utilisé ses contributions pour composer d'autres objets cartographiques. De plus, deux de ces 46 contributeurs ont utilisé plus de mille fois des nœuds provenant de N_4 pour créer des chemins. Par conséquent, la participation en amont du contributeur N_4 a dû favoriser la participation des contributeurs suivants par la modification et la réutilisation de ses données. En cela, le contributeur N_4 se comporte comme un pionnier. Les contributions de N_4 étant fortement réutilisées, ce contributeur semble être considéré comme fiable par

le reste de la communauté.

(ii) Analyse de la compétence : Le contributeur *N4* participe principalement en soirée et durant la semaine. Considérant cette seule information temporelle, le contributeur semble être un amateur expert. Toutefois, de la même manière que pour les modérateurs, il est nécessaire de mener une étude plus approfondie sur les contributions du pionnier *N4* pour comprendre à quel point sa participation est automatisée, et en quoi cela influence l'évaluation de son expertise. En particulier, le fait que les données de ce contributeur aient été largement éditées par la suite dénote potentiellement un manque de soin dans leur saisie.

c) Autres profils

L'étude des noyaux a permis de mettre en évidence deux profils particuliers chez les contributeurs d'OSM. Cependant, d'autres profils peuvent encore être identifiés : ils constituent autant de perspectives dans ce travail d'analyse du réseau multiplexe des contributeurs OSM. Dans cette partie, nous identifions d'autres groupes de contributeurs à étudier pour dégager potentiellement de nouveaux profils parmi les contributeurs qui n'ont pas été qualifiés jusqu'ici.

(i) Contributeurs isolés : Ce profil est détectable en filtrant les contributeurs qui ne prennent part à aucune interaction dans les couches du réseau multiplexe. Le pourcentage de sommets isolés dans le réseau multiplexe est de 39%. Autrement dit, 39% des contributeurs de l'Île de la Cité ont participé sans interagir avec aucun autre contributeur. Ces contributeurs ont donc ajouté de nouveaux objets qui n'ont pas été édités ni réutilisés, du moins dans l'intervalle (2013-2015), pour deux raisons possibles :

- ces contributions n'ont pas été vérifiées dans l'intervalle temporel étudié ;
- ces contributions sont de bonne qualité dès leur première version.

Ces contributeurs isolés ont un profil qui se rapproche de celui du pionnier et du cartographe solitaire. À nouveau, une étude approfondie sur leurs contributions permettrait d'obtenir des informations supplémentaires pour mieux qualifier leur fiabilité et leur compétence. Il se pourrait que, parmi ces contributeurs isolés, nous trouvions d'autres pionniers dont la fiabilité ne se caractérise pas par une forte réutilisation telle que détectée pour le contributeur *N4*, mais par d'autres indicateurs de fiabilité.

(ii) Contributeurs non-isolés : Les profils de modérateurs et de pionniers ont été détectés parmi ces contributeurs non-isolés, c'est-à-dire ceux qui interagissent dans au moins une couche du réseau multiplexe. Cependant, ces profils ont été identifiés parmi les noyaux des communautés centrées. Il reste donc à étudier les contributeurs qui forment les « signaux faibles » de la communauté d'OSM, à savoir :

- ceux qui s'organisent autour des noyaux dans les communautés centrées ;
- ceux qui appartiennent aux communautés non hiérarchiques.

3.3 Discussion

a) Influence de la composition des couches

La composition du réseau multiplexe semble déséquilibrée en faveur des interactions basées sur les éditions. En effet, sur les cinq couches du réseau, quatre d'entre elles sont des graphes d'édition. Par conséquent, nous considérons un sous-ensemble du réseau multiplexe composé uniquement des graphes de co-édition et d'utilisation. Ainsi, cette sous-sélection peut être considérée comme un second réseau multiplexe, dans lequel chaque couche correspond à une opération d'édition de données particulière. Nous noterons $R1$ le réseau multiplexe à cinq couches, et $R2$ le second réseau. La question est de savoir dans quelles mesures la composition des couches influence l'étude du réseau multiplexe.

La Figure II.26 montre la distribution des valeurs du coefficient de participation dans $R1$ et dans $R2$. On observe que la distribution dans $R1$ est désaxée sur la droite (Figure II.26a), à la différence de celle de $R2$ qui est désaxée sur la gauche (Figure II.26b). Comme les graphes de co-édition, de collaboration (en largeur et en profondeur) et de co-contribution modélisent la même interaction d'édition dans OSM, un contributeur qui interagit par édition de données va nécessairement apparaître dans chacun de ces graphes. Par conséquent, le coefficient de participation sera grand pour un tel contributeur. L'utilisation du réseau $R2$ permet alors de réduire la redondance des couches d'éditions. La Figure II.26b montre que la plupart des contributeurs interagissent à travers une seule opération d'édition, et généralement dans la couche $G_{co-edition}$.

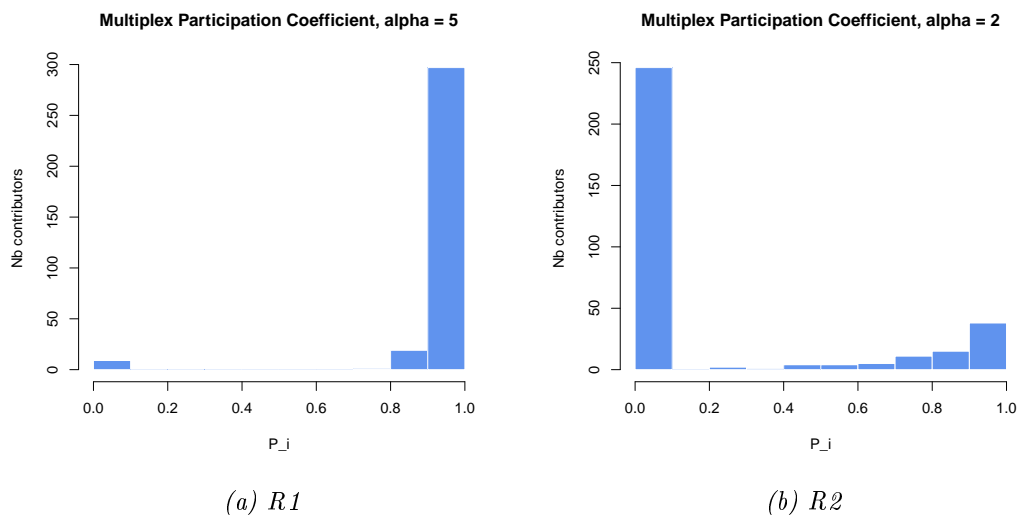


FIGURE II.26. Histogramme du coefficient de participation dans $R1$ et $R2$ (α désigne le nombre de couches dans le réseau considéré)

La composition des couches d'un réseau multiplexe va donc nécessairement impacter les métriques calculées sur ce réseau, ainsi que l'identification des profils de contributeurs. Stein *et al.* (2015) ont mis en évidence des profils de contributeurs (Table II.3) en analysant conjointement $G_{profondeur}$ et $G_{largeur}$. De même, il serait intéressant de chercher les combinaisons de couches du réseau à considérer pour révéler des profils particuliers.

b) Validité des configurations des communautés détectées

Parmi les communautés détectées par l'algorithme de Louvain, nous avons identifié trois types de configurations, à savoir :

- le couple ;
- la communauté centrée sur un noyau ;
- la communauté non hiérarchique.

L'identification de ces configurations était d'abord simplement visuelle. Bien que les communautés de type « couple » méritent d'être étudiées, nous ne nous sommes pas concentrés sur leur structure, car celle-ci est facilement repérable visuellement. En revanche, pour les deux autres configurations, il est opportun de valider ces configurations en étudiant la structure de ces sous-graphes. En particulier, il s'agit de s'assurer qu'il y a bien une différence entre ces deux types de configurations.

La Figure II.27 donne les distributions des coefficients de *clustering* C_i ⁹ des deux types de communautés étudiées. Les coefficients de *clustering* de la communauté centrée sur un noyau prennent des valeurs ou très faibles – ce sont des sommets qui sont uniquement connectés au noyau central – ou très élevées – ce sont les sommets dont les voisins sont également voisins entre eux.

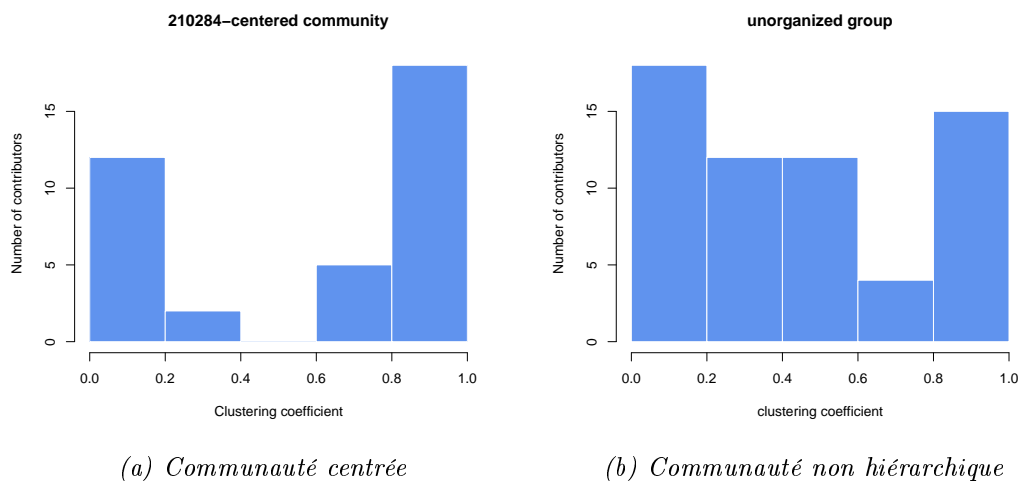


FIGURE II.27. Histogramme des coefficients de clustering

Par ailleurs, dans une communauté centrée, les sommets ayant un fort coefficient C_i (supérieur à 0.8) sont plus nombreux que les sommets dont la valeur de C_i est faible (inférieure à 0.2). La communauté centrée sur un noyau est plus soudée qu'elle n'en a l'air visuellement. La distribution du coefficient de *clustering* pour la communauté non hiérarchique est plus homogène que la communauté centrée, où il y a des sommets ayant un coefficient de *clustering* intermédiaire (entre 0.2 et 0.8). Finalement, la différence entre les deux configurations observées provient de l'existence de sommets ayant des centralités de degré et des valeurs de C_i intermédiaires.

La distribution des centralités de degré pour la communauté centrée¹⁰ est donnée dans la Figure II.28a. Cette distribution confirme l'observation visuelle de la

9. cf. définition II.10

10. celle de la Figure II.23a

configuration de la communauté : un contributeur de degré 30 forme un noyau auquel les autres contributeurs se rattachent sans être, pour la plupart, connectés entre eux, puisqu'ils sont de degré 1. Quant à la communauté non hiérarchique¹¹, la distribution des centralités de degré est donnée en Figure II.28b. Cette communauté contient également une majorité de sommets de degré 1. Cependant, nous observons à la fois des sommets présentant une forte centralité de degré (supérieure à 20) mais également des sommets de centralité de degré intermédiaire (entre 5 et 20).

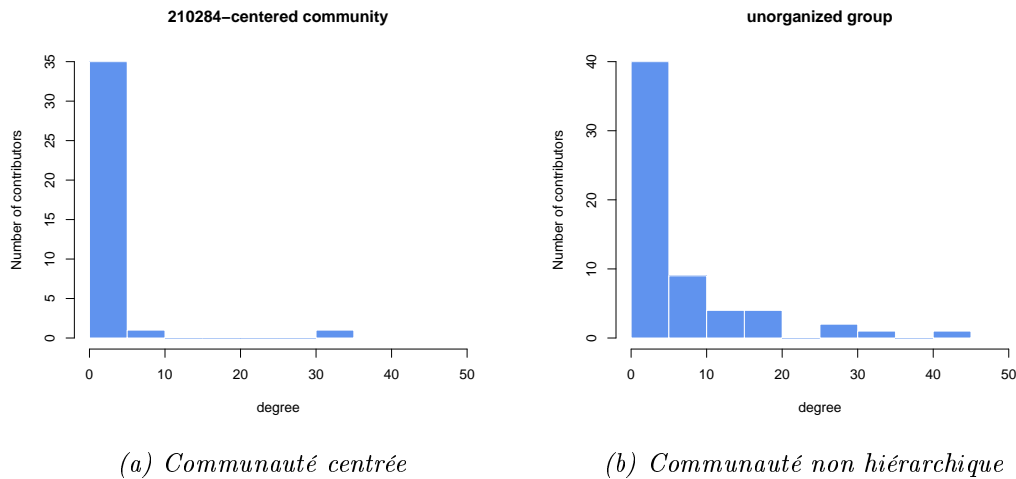


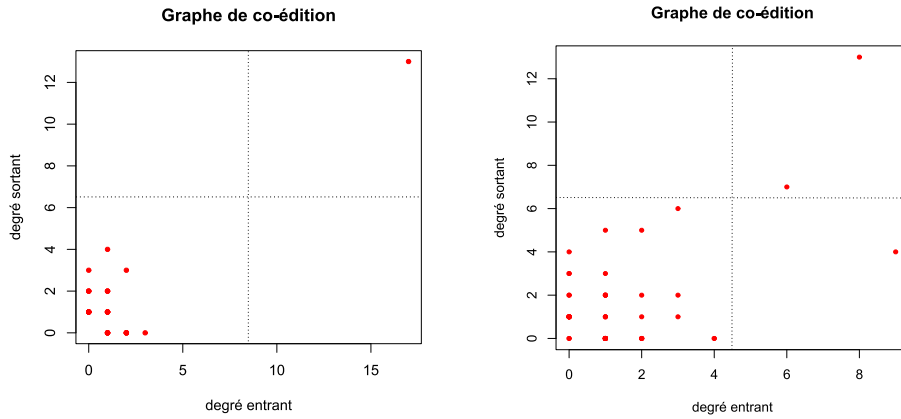
FIGURE II.28. Histogramme des centralités de degré des sommets dans chaque communauté.

Dans le réseau multiplexe $R1$, hormis le graphe de co-contribution, toutes les autres couches sont des graphes orientés. Afin d'observer plus précisément les interactions dans les communautés détectées, nous étudions les répartitions des degrés entrants et sortants des sommets de chaque graphe orienté du réseau multiplexe (Figures II.29, II.30 et II.31). En divisant les graphes de distribution des degrés entrants et sortants en quatre quadrants tel que sur la Figure II.29, nous pouvons observer que, dans les différentes couches du réseau multiplexe, la communauté centrée conserve une configuration assez homogène (Figures II.29a, II.30a et II.31a). En effet, dans ces trois distributions, la majorité des sommets se trouvent dans le quadrant inférieur gauche. Ces sommets présentent de faibles degrés entrants et sortants, tandis que le contributeur central présente de forts degrés entrants et sortants. Quant à la communauté non hiérarchique, ses sommets sont éparpillés dans les quatre quadrants du graphe de collaboration (Figure II.30b).

L'analyse de ces deux structures a permis de mettre en avant les caractéristiques propres à deux configurations de communautés détectées par l'algorithme de Louvain. Par conséquent, nous pouvons valider les configurations observées dans les communautés détectées par l'algorithme de Louvain, à savoir :

- la structure de communauté centrée autour d'un grand contributeur ;
- la structure de communauté non hiérarchique, qui présente une plus grande hétérogénéité dans les interactions et des contributeurs de centralité et de connectivité intermédiaires.

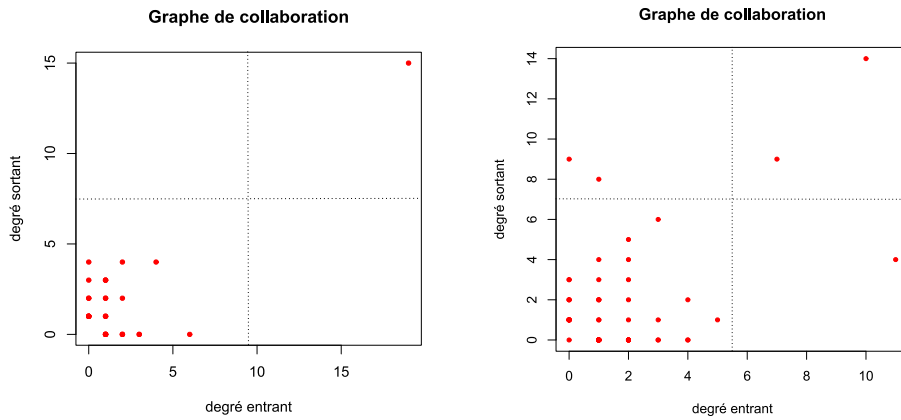
11. celle de la Figure II.23b



(a) Communauté centrée

(b) Communauté non hiérarchique

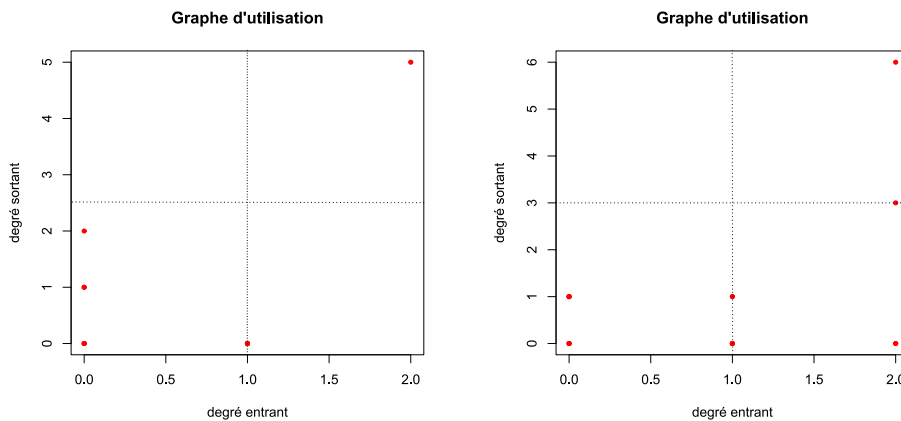
FIGURE II.29. Distribution des degrés entrants et sortants dans $G_{co\text{-}édition}$.



(a) Communauté centrée

(b) Communauté non hiérarchique

FIGURE II.30. Distribution des degrés entrants et sortants dans $G_{largeur}$.



(a) Communauté centrée

(b) Communauté non hiérarchique

FIGURE II.31. Distribution des degrés entrants et sortants dans $G_{utilisation}$.

c) Indépendance des configurations identifiées avec l'approche adoptée

La détection de communautés effectuée sur le graphe agrégé G_{agr} a permis d'identifier deux configurations de communautés caractéristiques : les communautés centrées sur un contributeur noyau et les communautés non hiérarchiques. Or, pour nous assurer que ces configurations ne sont pas conditionnées par la structure du graphe agrégé G_{agr} , nous procédons à la construction d'un second graphe G'_{agr} dans lequel sont agrégés seulement les arcs les plus significatifs de chaque couche. Il s'agit d'appliquer un seuil sur les poids des arcs, pour ne filtrer que les interactions les plus fortes. Ainsi, les interactions les plus faibles ne sont pas prises en compte dans le graphe agrégé G'_{agr} .

Nous fixons les seuils suivants :

- Couche de co-édition : filtrage des arcs (i, j, w) où $w > 3$
- Couche de largeur de collaboration : filtrage des arcs (i, j, w) où $w > 3$
- Couche de profondeur de collaboration : filtrage des arcs (i, j, w) où $w > 2$

Cette étape de filtrage permet de se débarrasser des redondances portées par les couches de co-édition, de largeur et de profondeur de collaboration. Le graphe agrégé G'_{agr} est alors construit à partir de ces graphes filtrés ainsi que les autres couches du réseau multiplexe $R1$.

En appliquant l'algorithme de Louvain sur le nouveau graphe agrégé G'_{agr} , nous observons, dans les nouvelles communautés détectées, des configurations plutôt proches de celles qui ont été trouvées pour G_{agr} . La communauté à laquelle appartient le contributeur #210284 conserve globalement une configuration de communauté centrée sur ce même noyau (Figure II.32a). Cependant, cette communauté présente une structure moins centrée que celle observée initialement (Figure II.23a) car nous y observons notamment la présence de sommets de degrés intermédiaires.

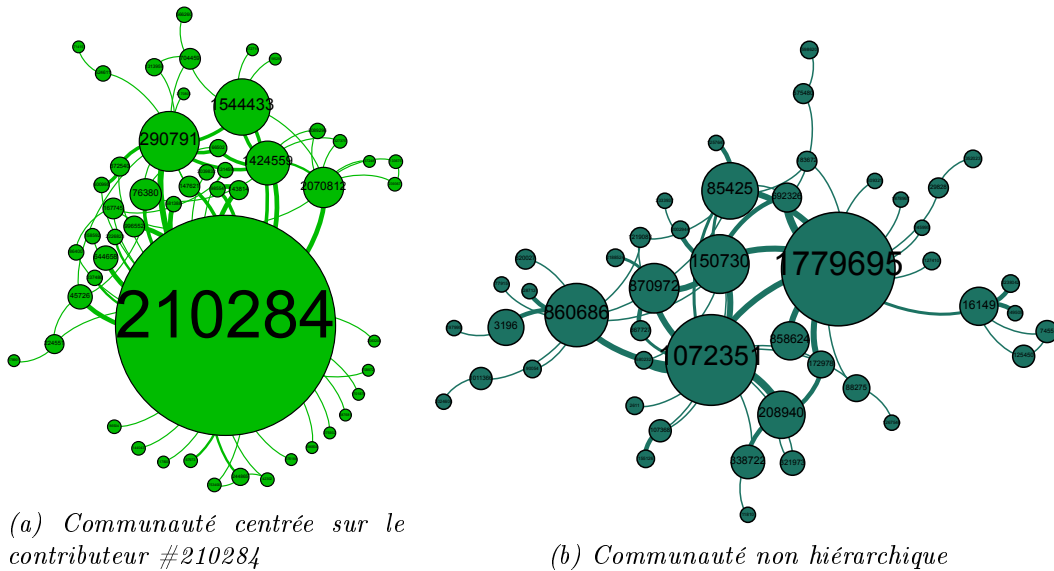


FIGURE II.32. Configurations caractéristiques dans les communautés détectées dans le nouveau graphe agrégé G'_{agr}

Par ailleurs, la structure de communauté non hiérarchique est également observée dans les communautés détectés sur G'_{agr} (Figure II.32b). Toutefois, les sommets

de très faible degré sont moins nombreux, car les communautés non hiérarchiques détectées sur G'_{agr} sont plus fortement connectées. En effet, le seuillage des arcs dans les couches de $R1$ qui a précédé la construction de G'_{agr} explique les fortes collaborations dans ces communautés.

Comme nous l'avons mentionné dans les autres profils à étudier, les communautés non hiérarchiques mériteraient d'être plus particulièrement étudiées, car elle contiennent la « masse » de contributeurs qui n'a pas encore été qualifiée jusqu'à présent. Les travaux menés par Bégin *et al.* (2018); Ma *et al.* (2015); Neis et Zielstra (2014) sur la communauté du projet OSM ont montré que cette dernière suit la règle des 90-9-1 de Nielsen (2006), où :

- 90% des membres enregistrés ne contribuent pas ;
- 9% de la communauté participent occasionnellement ;
- 1% de la communauté produit la plus grande partie du contenu de la base de données collaborative.

Selon cette classification, l'étude des contributeurs « noyaux » a permis d'identifier des profils des contributeurs les plus actifs, qui appartiennent à la catégorie minoritaire du projet OSM. Par ailleurs, le résultat de l'analyse quant à leur bonne fiabilité est en lien avec le fait qu'ils soient, d'une certaine manière, les fondateurs du projet en termes de quantité de données produites.

En étudiant les contributeurs qui appartiennent aux communautés non hiérarchiques, nous pourrions alors qualifier cette majorité de contributeurs moyennement actifs dont le niveau de confiance reste pour l'instant inconnu. De plus, cette structure non hiérarchique formée par de fortes collaborations rassemble potentiellement des contributeurs qui ont conscience de l'action des uns et des autres : typiquement, les désaccords entre certains contributeurs – tels que les guerres d'édition – pourraient être identifiés dans ces communautés. Par conséquent, l'analyse de ce type de structure permettrait de mettre en évidence des éléments pertinents à l'évaluation de la confiance des contributeurs.

d) Généricité de la méthode

Pour montrer la généricité de la méthode d'analyse du réseau multiplexe, nous construisons deux réseaux multiplexes $R3$ et $R4$ sur deux nouvelles fenêtres spatio-temporelles : $R3$ est construit à partir des données OSM collectées sur la ville de Stuhr en Allemagne entre 2009 et 2013, et $R4$ est construit sur les données de la ville de Katmandou au Népal entre 2014 et 2017.

Ces deux zones ont été choisies car elles ont une superficie relativement petite, ce qui permet de charger un nombre limité de données, et donc limiter le temps de traitement des données. Toutefois, nous cherchions des zones qui contenaient un nombre suffisant de contributions pour construire les graphes. Entre 2009 et 2013, une très grande partie des contributions – entre 26% et 50% des données de la base mondiale – a été produite sur l'Allemagne (Neis et Zipf, 2012). Par conséquent, nous avons choisi de travailler sur une zone allemande. En ce qui concerne Katmandou, la ville a subi un important tremblement de terre en 2015, ce qui a incité une vague de contributeurs OSM à rejoindre des projets collaboratifs pour cartographier

d'urgence la zone sinistrée¹². C'est pourquoi nous avons choisi d'étudier cette zone cartographique sur une fenêtre temporelle qui contient la date du tremblement de terre.

Les réseaux multiplexes $R3$ et $R4$ sont composés des graphes de co-édition, de largeur et de profondeur de collaboration, de réutilisation, de suppression, de co-location et de co-temporalité. Les graphes de co-location et de co-temporalité sont respectivement pondérés par le recouvrement spatial et temporel. Le principe de ces graphes est schématisé dans la Section 2.2 aux Figures II.18 et II.19).

En appliquant le même seuillage sur les graphes de co-édition et de collaboration (en largeur et en profondeur) que sur $R1$, la détection de communautés renvoie des configurations plutôt similaires à celles trouvées avec $R1$. Parmi les communautés détectées dans les réseaux multiplexe $R3$ et $R4$, nous observons des communautés centrées et des communautés non hiérarchiques (Figures II.33 et II.34).

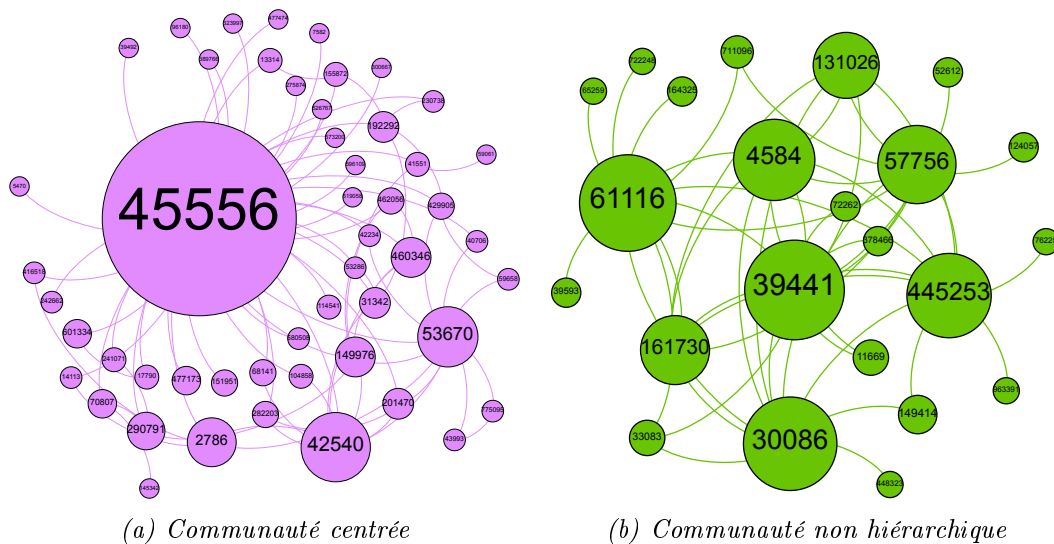


FIGURE II.33. Communautés détectées dans le réseau multiplexe $R3$ des contributeurs OSM à Stuhr.

De plus, le calcul du coefficient de *clustering* sur ces communautés confirme les observations qui ont été faites précédemment avec les communautés de $R1$. D'après la Table II.7, nous observons que, dans les communautés centrées autour d'un noyau, une grande partie des sommets ont un faible coefficient de *clustering* (environ 40% ont un coefficient compris entre 0 et 0.1). Au contraire, dans les communautés non hiérarchiques, une grande partie des sommets ont un fort coefficient de *clustering* (environ 40% ont un coefficient compris entre 0.8 et 1). Cette expérience initiale montre la généricité de la méthode sur d'autres zones spatio-temporelles.

Bien que les mêmes configurations aient été observées, cette identification ne suffit pas pour conclure sur l'existence de modérateurs et de pionniers dans une zone. Cette étape nécessite des analyses supplémentaires sur les interactions des contributeurs « noyaux » dans les différentes couches du réseau multiplexe, sur leur mode d'activité et sur les données éditées, de la même manière que pour les contributeurs

12. https://www.hotosm.org/updates/2015-04-30_helping_to_map_nepal_getting_started

de l'île de la Cité. Une étude approfondie sur ces réseaux multiplexes pourrait être intéressante pour mettre en évidence des collaborations propres aux zones de Stuhr et Katmandou, qui peuvent être différentes de celles de l'île de la Cité à Paris.

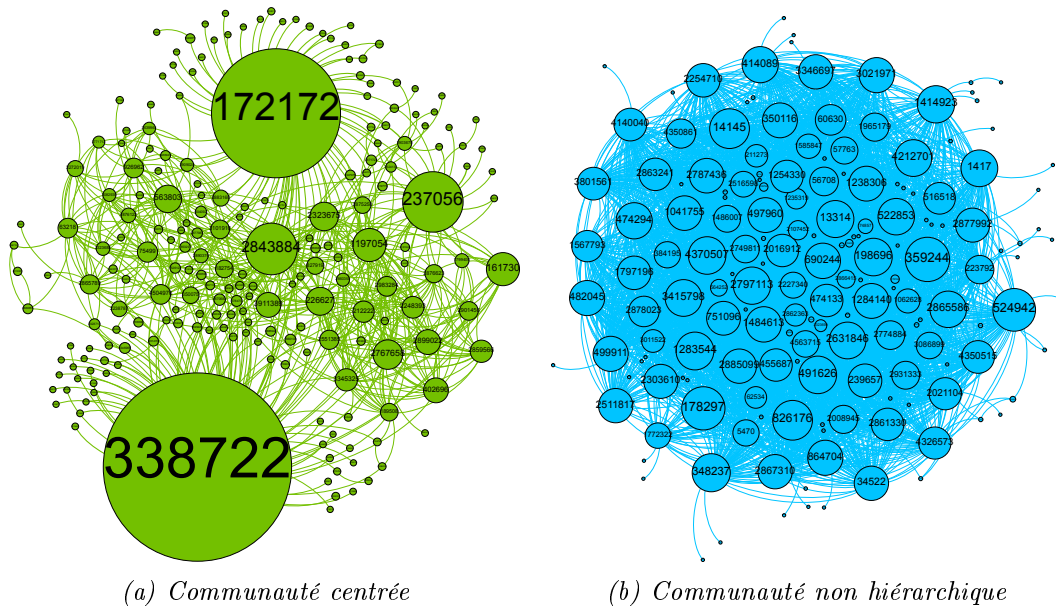


FIGURE II.34. Communautés détectées du réseau multiplexe RA des contributeurs OSM à Katmandou.

TABLE II.7. Distribution des valeurs du coefficient de clustering dans les différentes communautés de R3 et R4 représentées dans les Figures II.33 et II.33.

Configuration		Valeurs de C_i			
		[0, 0.1[[0.1, 0.5[[0.5, 0.8[[0.8, 1]
Centrée	Stuhr	41.2%	11.8%	9.8%	37.3%
	Katmandou	46.4%	14.2%	14.6%	24.9%
Non hiérarchique	Stuhr	21.6%	17.4%	17.4%	39.1%
	Katmandou	32.7%	0.6%	18.5%	48.2%

3.4 Conclusion et perspectives de l'étude multiplexe

Les expériences menées sur le réseau multiplexe des contributeurs d'OSM sur l'île de la Cité a permis de mettre en évidence les collaborations entre les contributeurs. En particulier, l'analyse de ce réseau a permis de détecter des profils comportementaux de modérateurs et de pionniers. Ces profils de contributeurs sont intéressants pour garantir la qualité des données selon l'approche sociale, car ils identifient des contributeurs présentant un certain niveau confiance, alors que ces derniers ne sont pas explicitement désignés comme étant des modérateurs.

Notre travail de recherche sur ce sujet a été publié dans une revue internationale en géomatique (Truong *et al.*, 2019), et a fait l'objet de présentations lors d'un atelier scientifique (Truong *et al.*, 2018a) et d'une conférence internationale en géomatique (Truong *et al.*, 2018b). Toutefois, l'étude du réseau multiplexe mérite encore d'être

approfondie sur plusieurs aspects. Bien que les perspectives de ce travail soient nombreuses, nous pouvons en relever quelques-unes :

1. Il serait utile de développer une méthode de composition ou de sélection des couches du réseau, de manière à ne pas introduire des redondances d'information, ou bien de les prévoir et déterminer les couches à sélectionner pour le calcul de certains indicateurs. Par exemple, nous avons vu que le calcul du coefficient de participation donnait des résultats différents entre les réseaux *R1* et *R2*.
2. Il serait également intéressant d'étudier les communautés qui seraient détectées par des approches différentes que par l'agrégation de couches. En effet, de nouvelles configurations pourraient apparaître, et d'autres collaborations pourraient être identifiées.
3. Mis à part les noyaux des communautés centrées dont les profils ont été analysés, une autre piste d'approfondissement serait d'analyser les autres contributeurs. Certains contributeurs, ayant une participation moindre, peuvent éventuellement avoir des comportements de modérateur ou manifester d'autres signes de fiabilité.
4. Il serait intéressant d'étudier la capacité de ces contributeurs qui ajoutent peu de nouvelles données mais dont les éditions, bien que moins nombreuses, sont tout aussi importante pour maintenir la base de données OSM ([Anderson et al., 2018](#)).
5. À partir des profils de modérateurs et de pionniers que nous avons pu identifier manuellement, il s'agira dans le futur de mettre en place des processus de qualification automatique de ces profils de contributeurs.

Au-delà de ces perspectives, notre étude sur les graphes d'interaction a permis de montrer qu'il était possible de qualifier les contributeurs à partir de leurs interactions sur les différentes opérations d'édition. Les résultats issus des graphes d'interaction peuvent déjà être ré-exploités pour évaluer la confiance des contributeurs.

4 Évaluation de la confiance du contributeur

4.1 Positionnement du problème

L'analyse des graphes d'interaction à travers le modèle de réseau multiplexe a proposé une première qualification des contributeurs OSM. En particulier, les informations extraites des opérations d'édition sont de bons indicateurs du niveau de confiance des contributeurs. Dans le but de qualifier le carto-vandalisme des données géographiques collaboratives, il faut tenir compte de la fiabilité du contributeur. Par conséquent, nous cherchons à mettre en place un indicateur de fiabilité du contributeur à partir des résultats issus des graphes d'interaction.

Pour détecter le vandalisme dans les données OSM, [Neis et al. \(2012\)](#) ont développé l'outil *OSMPatrol*, leur permettant de trouver que 76% du vandalisme OSM provient de nouveaux contributeurs. Rappelons que, selon [Neis et al. \(2012\)](#), le vandalisme est réduit à une dégradation de données uniquement, qu'elle soit intentionnelle ou non. Cette définition du vandalisme est donc différente de celle que nous

adoptons. Le système de réputation mis en place par Neis *et al.* (2012) évalue donc la fiabilité des contributeurs en distinguant les nouveaux contributeurs des contributeurs experts. Cependant, 50% des contributeurs experts – c’est-à-dire ceux qui ont bonne réputation d’après le système – provoquent près de la moitié des dégradations enregistrées. De plus, les auteurs admettent qu’*OSMPatrol* détecte un grand nombre de faux positifs, c’est-à-dire que le système détecte des contributions qui ne sont pas réellement dégradées. L’expertise du contributeur, évaluée à partir du nombre d’éditations, ne semble donc pas être un bon indicateur de qualification du carto-vandalisme.

L’évaluation de la qualité de participation des contributeurs est une problématique qui a été traitée dans le cadre du projet Wikipédia (Geiger et Halfaker, 2013). Une évaluation de la qualité de participation des contributeurs de Wikipédia montre l’influence du critère considéré sur le classement des contributeurs (Geiger et Halfaker, 2013). En assimilant la qualité de participation aux nombre d’éditations, les nouveaux contributeurs sont alors défavorisés puisqu’ils comptent peu de contributions à leur actif. En revanche, lorsque la participation est évaluée en fonction du temps passé à contribuer, l’investissement des contributeurs qui lancent des scripts automatiques est sous-estimé, car leurs contributions sont généralement de courte durée, aussi nombreuses soient-elles.

D’après la règle des 90-9-1 de Nielsen (2006), nous avons vu que 1% des contributeurs est à l’origine de la plupart des données de la base ouverte OSM. Or, l’étude présentée dans la Section 3 a permis de montrer que cette catégorie de contributeurs prolifiques contient en effet des profils de confiance, tels que des modérateurs et des pionniers. L’hypothèse selon laquelle un contributeur prolifique est fiable est donc valide pour les profils que nous avons étudiés jusqu’ici.

Toutefois, cette hypothèse ne permet pas d’en dire plus sur la fiabilité des profils moins prolifiques en termes de production de données. De plus, les grands contributeurs peuvent faire usage de scripts automatiques pour éditer un grand nombre de données, comme cela a pu être le cas pour des contributeurs noyaux étudiés dans la partie précédente. Or, l’édition massive et automatique de données n’entraîne pas forcément des contributions de bonne qualité. La Figure II.35 illustre un exemple d’erreur d’import automatique de données où des milliers de bâtiments africains se sont trouvés cartographiés par erreur en Amérique.

Le problème ici est qu’une contribution de bonne qualité provient d’un contributeur fiable, mais un contributeur prolifique ne produit pas forcément des données de qualité. Nous explorons donc des moyens d’évaluer la fiabilité des contributeurs qui surpasseraient les limites de l’indicateur du nombre de contributions. Dans la Section 4.2 suivante, nous proposons deux nouvelles méthodes d’évaluation de la fiabilité des contributeurs, à partir des métriques issues des modèles de graphe de collaboration présentés dans ce chapitre. La Section 4.3 est une étude de la fiabilité des contributeurs OSM, qui permet de vérifier expérimentalement si le nombre de contributions est révélateur de la fiabilité des contributeurs OSM. Enfin, la Section 4.4 compare les classements des contributeurs issus de nos méthodes d’évaluation de la fiabilité avec le classement par nombre de contributions.



FIGURE II.35. Erreur d'import automatique de données dans OSM.

4.2 Proposition de nouvelles méthodes d'évaluation

Dans un premier temps, nous listons les différentes métriques qui, dans une certaine mesure, indiquent le niveau de confiance d'un contributeur de données. Dans un second temps, nous proposons d'évaluer la confiance du contributeur selon deux méthodes différentes :

1. en calculant un score moyen des métriques proposées ;
2. en classant les contributeurs selon une méthode de décision multicritère, chaque métrique valant un critère.

a) Les indicateurs de fiabilité

La Table II.8 décrit l'ensemble des métriques à considérer théoriquement pour évaluer la fiabilité des contributeurs OSM.

TABLE II.8. Indicateurs de fiabilité du contributeur

Notation	Métrique
<i>nbContrib</i>	Nombre de contributions
<i>pModif</i>	Pourcentage de modifications faites par le contributeur
<i>pDelete</i>	Pourcentage de suppressions faites par le contributeur
<i>pUsed</i>	Pourcentage de contributions qui ont été utilisées par d'autres contributeurs
<i>pDeleted</i>	Pourcentage de contributions supprimées
<i>pModified</i>	Pourcentage de contributions modifiées
<i>nbWeeks</i>	Nombre de semaines
<i>focalisation</i>	Focalisation moyenne des <i>changesets</i> sur la zone d'étude
<i>profil</i>	Profil du contributeur (identifié d'après l'étude du réseau multiplexe)

Mis à part l'indicateur *profil*, tous les autres indicateurs ont été implémentés à partir de données OSM de la ville d'Aubervilliers, en France. Depuis quelques

années, Aubervilliers subit d'importants travaux d'aménagements¹³. Au regard des changements réels dans cette ville, nous nous attendions à ce qu'il y ait une certaine activité de contribution sur la zone cartographique correspondante dans OSM qui refléterait les évolutions du terrain.

(i) **Nombre de contributions** : La métrique $nbContrib$ est calculée pour chaque contributeur. Elle correspond au nombre de contributions produites sur la zone d'étude, c'est-à-dire que seules les contributions produites sur la fenêtre spatiale étudiée sont comptées. Bien que cette métrique présente des limites pour évaluer la confiance du contributeur, il nous semble important de la prendre en compte.

(ii) **Modes de contribution** : Les métriques $pModif$, $pDelete$, $pUsed$, $pDeleted$ et $pModified$ sont calculées à partir des graphes d'interaction. Elles donnent une indication sur le mode d'édition pour chaque contributeur : à quel point le contributeur édite la base de données par la modification et la suppression d'information, et comment ses éditions sont perçues par les autres membres du projet, à travers la réutilisation de leurs contributions ou au contraire, la suppression de celles-ci. La Figure II.36 schématise le processus d'obtention de ces différentes métriques.

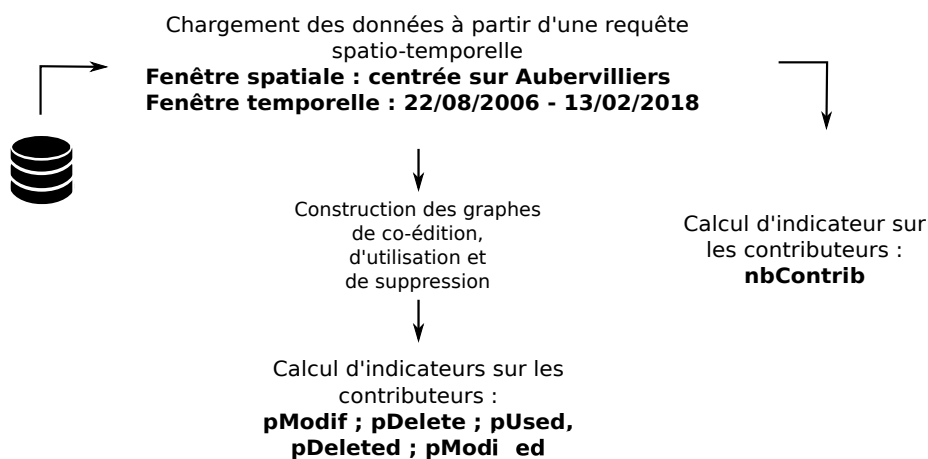


FIGURE II.36. Chaîne de traitement des données pour le calcul des indicateurs sur les contributeurs d'OSM.

(iii) **Investissement temporel** : Comme cela a été soulevé pour les données Wikipédia (Geiger et Halfaker, 2013), il est important de tenir compte de l'investissement en temps des contributeurs. L'étude de l'investissement temporel des contributeurs sur une échelle globale a permis d'identifier des profils qui témoignent de leur niveau de survie sur le projet de l'ordre du jour, de mois ou de l'année (Bégin et al., 2018). Or, quel que soit l'investissement global d'un contributeur, il est intéressant de mesurer à une échelle plus locale la régularité de ses contributions sur une zone donnée. Par exemple, un contributeur qui édite quelques objets sur une même zone de manière régulière n'est pas forcément moins fiable qu'un contributeur qui, en une seule fois, édite massivement une zone sans ne jamais y retourner. De plus,

13. <https://plainecommune.fr/projets/grands-projets-urbains/>

un contributeur qui édite régulièrement une zone a des chances d’être un observateur local ou du moins de posséder une certaine connaissance sur cette zone. Bien que ces hypothèses nécessitent d’être vérifiées, il nous semble que la connaissance locale constitue un point fort de la démarche collaborative, car elle peut permettre de capturer les changements du monde réel en temps réel (Goodchild, 2007).

Le bon niveau d’échelle à adopter pour mesurer la régularité des contributions n’est pas une question triviale. Si la mesure de l’investissement temporel se fait au niveau de l’année, elle ne permettra pas de différencier un contributeur qui a édité une seule fois dans l’année du contributeur qui participe chaque mois. Si cette mesure se fait au niveau de la journée, la régularité des contributeurs qui éditent à une fréquence hebdomadaire ou mensuelle ne sera pas aussi bien capturée que les contributeurs très actifs (et rares) dont l’activité continue peut s’expliquer par des techniques d’édition automatique. Par conséquent, nous considérons une métrique d’investissement temporel *nbWeeks* comptant le nombre de semaines durant lesquelles un contributeur a participé sur la plateforme OSM. Cette métrique se base uniquement sur les contributions extraites de la fenêtre spatio-temporelle étudiée.

(iv) **Focalisation moyenne :** De la même manière que l’investissement temporel, l’activité spatiale de contribution peut indiquer le niveau de connaissance local des contributeurs. L’étude des collaborations entre contributeurs a permis d’identifier des contributeurs locaux considérés comme des experts (Stein *et al.*, 2015). En effet, un contributeur local est censé posséder une certaine connaissance de son environnement local (bien que cela ne soit pas prouvé). Par conséquent, il sera plus fiable qu’un contributeur qui édite massivement sur un grand espace qui engloberait la zone d’étude, car dans ce dernier cas, on pourrait soupçonner le contributeur de ne pas être aussi soigneux que le premier contributeur. Pour favoriser les profils de contributeurs locaux dans l’évaluation de la fiabilité, nous mettons donc en place une métrique de focalisation moyenne des contributions. Cette métrique, notée *focalisation*, est calculée à partir de l’enveloppe spatiale des sessions d’éditations du contributeur, également appelées *changesets*. Le calcul de l’enveloppe spatiale à partir des différents *changesets* d’un contributeur OSM est expliqué sur le schéma de la Figure II.37.

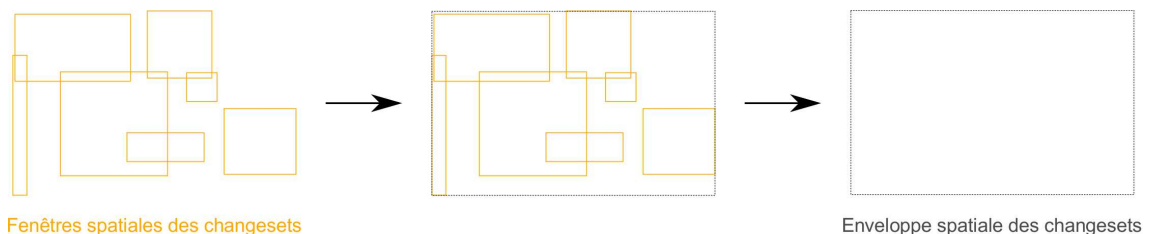


FIGURE II.37. Méthode de calcul de l’enveloppe spatiale des différents *changesets* d’un contributeur OSM.

Le principe de cette métrique est de mesurer à quel point l’activité d’un contributeur est ciblée sur une zone d’étude. La métrique *focalisation* est une quantité comprise entre 0 et 1 : *focalisation* = 0 lorsque la zone de contribution est très peu focalisée sur la zone d’étude, et *focalisation* = 1 lorsque les contributions sont concentrées dans la zone d’étude. La Figure II.38 présente les situations possibles à

considérer pour calculer la métrique de focalisation moyenne.

Soit *enveloppe* l'enveloppe spatiale sur laquelle a participé un contributeur, et soit *zoneEtude* la fenêtre spatiale de la zone d'étude. Dans le cas où l'activité spatiale est ciblée à l'intérieur de la zone d'étude (Figure II.38a), la focalisation est maximale, soit $focalisation = 1$. Dans le cas contraire, l'activité spatiale se trouve en partie à l'extérieur de la zone d'étude (Figure II.38b). La focalisation se calcule alors par la formule suivante :

$$focalisation = \frac{aire(enveloppe \cap zoneEtude)}{aire(enveloppe)} \quad (II.28)$$

De cette manière, plus l'enveloppe spatiale s'étend au-delà de la zone d'étude, moins les contributions sont ciblées sur cette zone, et moins le contributeur est susceptible d'être un observateur local.

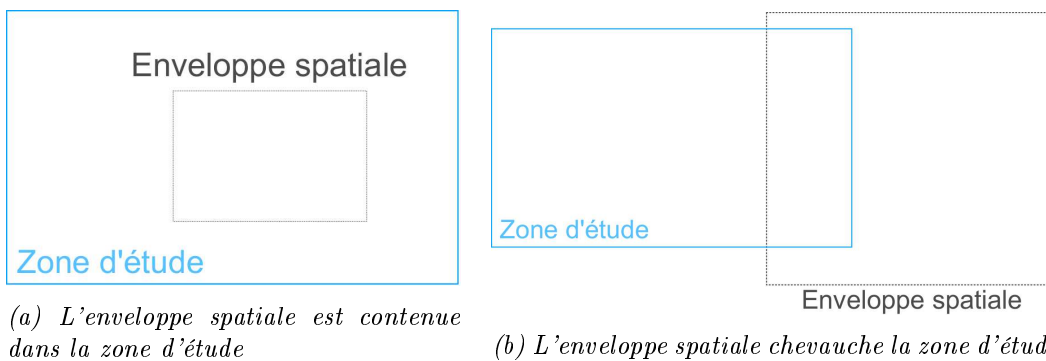


FIGURE II.38. Positionnements possibles de l'enveloppe spatiale par rapport à la zone d'étude.

b) Méthode n° 1 : calcul de la moyenne des métriques

TABLE II.9. Poids utilisés pour calculer la moyenne pondérée du score de fiabilité.

Métrique de fiabilité	Poids
$nbContrib$ (indicateur normalisé)	2
$pModif$	1
$pDelete$	1
$pUsed$	2
$pNonDeleted$	2
$pNonModified$	2
$nbWeeks$ (indicateur normalisé)	3
$focalisation$	3

Pour évaluer la fiabilité des contributeurs en tenant compte de ses différentes métriques, une première solution est de considérer la moyenne de ces métriques. Plus un contributeur aura une moyenne élevée, plus il sera vu comme fiable. Pour cela, nous considérons les métriques :

$$pNonModified = 1 - pModified \quad (II.29)$$

et

$$pNonDeleted = 1 - pDeleted \quad (\text{II.30})$$

de manière à ce que la variation des métriques ait le même sens, à savoir que la fiabilité du contributeur augmente avec le nombre de ses contributions qui ne sont pas modifiées (respectivement supprimées). Par ailleurs, ce score moyen peut se calculer en pondérant les métriques. Nous avons donc considéré deux scores : l'un est une moyenne simple, l'autre tient compte d'une pondération totalement arbitraire. Toutefois, la pondération des métriques mérite d'être étudiée pour que ce score moyen reflète au mieux la fiabilité des contributeurs. La Table II.9 indique les poids assignés *a priori* à chaque métrique pour le calcul de la moyenne pondérée.

c) Méthode n° 2 : classement selon la méthode PROMETHEE II

PROMETHEE II est une méthode de décision multicritère. Elle consiste à proposer un classement des décisions candidates en optimisant les critères donnés en entrée (Brans et Mareschal, 2005). Cette méthode a été notamment appliquée en géomatique pour traiter le problème de référencement des données thématiques sur des données topographiques, tout en respectant les règles de cohérence topologique (Jaara *et al.*, 2014).

Dans notre cas, les contributeurs constituent des décisions candidates à classer selon nos différentes métriques de fiabilité, qui modélisent les critères en entrée du système de classification. La méthode PROMETHEE II consiste à faire un classement des contributeurs en les comparant deux à deux sur chaque critère. Contrairement à la méthode du score moyen des métriques, la méthode PROMETHEE II ne tient pas compte directement de la valeur des métriques, mais de la différence des valeurs de chaque métrique entre deux contributeurs. Le système va classer les contributeurs en fonction de la distance calculée sur chaque métrique. Par ailleurs, cette méthode offre également la possibilité de pondérer les critères, c'est-à-dire qu'elle tient compte de l'importance accordée à chaque métrique – ou plutôt, à la différence de valeurs calculée sur chaque métrique – pour classer les contributeurs.

Soit $U = \{user_1, user_2, \dots, user_n\}$ l'ensemble des contributeurs d'OSM. Soit $C = \{c_1(user), c_2(user), \dots, c_k(user) | user \in U\}$ l'ensemble des métriques, c'est-à-dire les critères à optimiser. Soit d_j la différence des valeurs de métrique c_j de deux contributeurs :

$$d_j(user_a, user_b) = c_j(user_a) - c_j(user_b) \quad (\text{II.31})$$

La comparaison entre deux contributeurs s'effectue en définissant, pour chaque critère j , une fonction de préférence P_j quant à la différence d_j observée sur la métrique j :

$$P_j(user_a, user_b) = P_j(d_j(user_a, user_b)) \quad (\text{II.32})$$

Si $P_j(user_a, user_b) = 0$, alors il n'y a pas de préférence entre les contributeurs $user_a$ et $user_b$. Si $P_j(user_a, user_b) = 1$ alors le contributeur $user_a$ est préféré au contributeur $user_b$ sur le critère j . Dans notre cas, nous avons choisi arbitrairement de définir chaque fonction de préférence sous la forme :

$$P_j(user_a, user_b) = P_j(d_j(user_a, user_b)) \quad (\text{II.33})$$

$$P_j(\text{user}_a, \text{user}_b) = \begin{cases} 0 & \text{si } d_j \leq 0 \\ 1 - e^{-\frac{d_j^2}{2s^2}} & \text{si } d_j > 0 \end{cases} \quad (\text{II.34})$$

où s est un seuil à fixer. La Figure II.39 donne l'allure de la courbe gaussienne de la fonction de préférence.

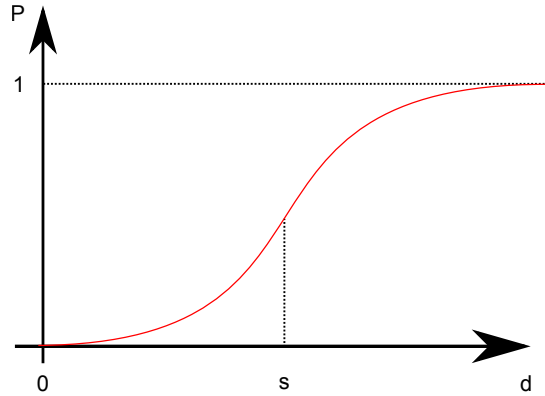


FIGURE II.39. Fonction de préférence de type gaussienne

En calculant la préférence du contributeur user_a sur le contributeur user_b sur chaque critère – *i.e.* sur chaque métrique de fiabilité – l'indice de préférence agrégée $\pi(\text{user}_a, \text{user}_b)$ va exprimer la préférence globale du contributeur user_a sur le contributeur user_b :

$$\pi(\text{user}_a, \text{user}_b) = \sum_{j=1}^k P_j(\text{user}_a, \text{user}_b) w_j \quad (\text{II.35})$$

où w_j est le poids donnant une importance relative à chaque critère, telle que $\sum_{j=1}^k w_j = 1$. En calculant l'indice de préférence agrégé entre chaque paire de contributeurs, un graphe de surclassement complet est généré, comme illustré sur la Figure II.40.

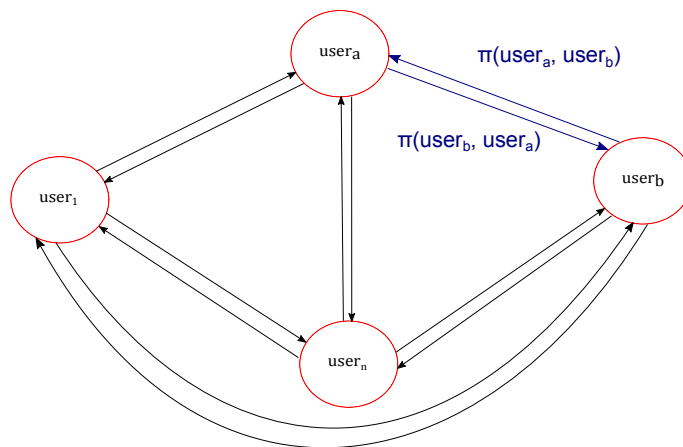


FIGURE II.40. Graphe de surclassement complet

Chaque contributeur est alors évalué à partir des deux flux de surclassement ϕ^+ et ϕ^- . Le flux de surclassement positif ϕ^+ exprime le pouvoir d'un contributeur user

à surclasser tous les autres :

$$\phi^+(user) = \frac{1}{n-1} \sum_{x \in U} \pi(user, x) \quad (\text{II.36})$$

Le flux de surclassement négatif ϕ^- exprime le caractère d'un contributeur *user* à être surclassé par les autres contributeurs :

$$\phi^-(user) = \frac{1}{n-1} \sum_{x \in U} \pi(x, user) \quad (\text{II.37})$$

Le flux net de surclassement du contributeur *user* est défini par :

$$\phi(user) = \phi^+(user) - \phi^-(user) \quad (\text{II.38})$$

Le système de décision multicritère va alors ordonner les contributeurs selon leur flux net de surclassement.

Le paramétrage du modèle PROMETHEE II consiste à fixer, pour chaque critère : le type de fonction de préférence, le seuil de la fonction de préférence et le poids relatif. Pour les mêmes raisons que l'étude de la pondération des métriques pour le calcul du score moyen de fiabilité, il est nécessaire d'étudier le paramétrage de la méthode PROMETHEE II. Cette étude n'a pas été effectuée dans le cadre de cette thèse, mais elle en constitue une perspective. Toutefois, nous avons appliqué la méthode PROMETHEE II en fixant arbitrairement les seuils et les poids de nos critères (Table II.10). Par ailleurs, toutes les fonctions de préférences P_j ont été choisies comme des fonctions gaussiennes de manière totalement arbitraire. Il existe d'autres types de fonction de préférences – par exemple des fonctions linéaires – qu'il faudrait considérer pour étudier le bon paramétrage de notre modèle.

TABLE II.10. Paramétrage du modèle de décision multicritère PROMETHEE II

Métrique de fiabilité	Poids w_j	Seuil s
<i>nbContrib</i>	0.01	1
<i>pModif</i>	0.13	0.2
<i>pDelete</i>	0.13	0.1
<i>pUsed</i>	0.13	0.1
<i>pDeleted</i>	0.1	0.01
<i>pModified</i>	0.1	0.01
<i>nbWeeks</i>	0.2	2
<i>focalisation</i>	0.2	0.3

4.3 Fiabilité et participation des contributeurs : une étude expérimentale sur OpenStreetMap

Cette étude a pour but de vérifier expérimentalement l'hypothèse selon laquelle le nombre de contributions est un indicateur de fiabilité des contributeurs. Pour cela, nous considérons l'ensemble des contributeurs qui ont édité des données sur

la ville d'Aubervilliers (93) entre le début du projet OSM (le 9 août 2004) et la date du 13 février 2018. Au total, 316 contributeurs ont participé dans cette fenêtre spatio-temporelle.

Nous considérons un échantillon de 30 contributeurs ayant participé sur la zone, et dont le nombre de contributions varie entre 1 et 93657. Nous avons choisi aléatoirement ces contributeurs selon le nombre de leurs contributions sur la zone, de manière à ce que l'échantillon soit constitué d'individus de profils différents en termes de nombre de contributions. La Table II.11 indique les caractéristiques de la distribution du nombre de contributions de cet échantillon. En explorant le contenu des contributions de chaque contributeur, nous les classons manuellement¹⁴ selon trois catégories :

- « fiable » ;
- « non fiable » ;
- « ne sait pas ».

La catégorie « ne sait pas » permet de ne pas nous prononcer sur la fiabilité d'un contributeur lorsque ses contributions n'apportent pas d'information à ce sujet.

TABLE II.11. *Résumé statistique du nombre de contributions des 30 contributeurs étudiés.*

Valeur min.	1 ^{er} quartile	Médiane	Moyenne	3 ^{ème} quartile	Valeur max.
1.0	7.5	92.0	3663.9	737.2	93657.0

En découpant le classement des contributeurs selon le nombre de leurs contributions, la qualification des différents contributeurs est indiquée dans la Table II.12. Nous observons que la totalité des contributeurs du dernier quartile (Q_4) a été qualifiée comme fiable. Par ailleurs, la proportion de contributeurs fiables augmente avec le nombre de contributions : elle passe de 50% de contributeurs fiables parmi les plus petits contributeurs (Q_1) – ceux qui ont effectué entre 1 et 7 éditions sur la zone – à 100% de contributeurs fiables parmi les plus grands contributeurs (Q_4) de l'échantillon étudié. La proportion des contributeurs non fiables est la plus élevée dans le quartile Q_1 des plus petits contributeurs. Par conséquent, la quantité de contributions semble être un bon indicateur de la fiabilité des contributeurs.

TABLE II.12. *Qualification manuelle de la fiabilité des contributeurs OSM.*

Quartile	Q1	Q2	Q3	Q4
Nombre de contributions	[[1, 8[[[8, 92[[[92, 737[[[737, 93657]
Nombre de contributeurs	8	7	7	8
Contributeurs fiables (%)	50.0	71.4	71.4	100.0
Contributeurs non fiables (%)	25.0	14.3	14.3	0
Contributeurs « ne sait pas » (%)	25.0	14.3	14.3	0

Cependant, notons que, dans le groupe des petits contributeurs, la moitié a été qualifiée comme étant de confiance. En réalité, les petits contributeurs ne sont pas caractérisés par un manque de fiabilité, mais plutôt par un manque de données qui

14. Voir Annexe B pour l'annotation manuelle des 30 contributeurs

ne permet pas de les qualifier facilement sur ce critère. En effet, puisque nous nous basons sur les contributions pour évaluer la fiabilité, cela explique pourquoi nous trouvons moins de contributeurs fiables dans le groupe des petits contributeurs.

Par ailleurs, il existe des contributeurs non fiables dans les quartiles Q_2 et Q_3 . En particulier, notre exploration a permis d'identifier un profil de vandale dans Q_3 , c'est-à-dire parmi des contributeurs plutôt prolifiques. Par conséquent, l'utilisation du nombre de contributions pour indiquer la fiabilité des contributeurs comporte donc le risque d'accorder une confiance dans des contributeurs prolifiques qui, en réalité, n'en sont pas dignes. Or, dans le cadre de la détection du carto-vandalisme, nous avons besoin de qualifier précisément les potentiels carto-vandales, c'est-à-dire les contributeurs non fiables. Une qualification trop hâtive des contributeurs prolifiques comme fiables entraîne un risque de ne pas détecter le carto-vandalisme qu'ils pourraient provoquer en grande quantité.

Par conséquent, bien que le nombre de contributions soit un bon indicateur de fiabilité des contributeurs, il ne semble pas suffisant pour qualifier la fiabilité dans le cadre du carto-vandalisme. En effet, il nous faut considérer un indicateur de fiabilité qui permette de mieux identifier les contributeurs non-fiables, afin de mieux qualifier les contributions qui relèvent du carto-vandalisme. Nous proposons donc, dans la partie suivante, deux méthodes d'évaluation de la fiabilité du contributeur, dans lesquelles d'autres métriques sont considérées en plus de l'indicateur du nombre de contributions. En effet, le nombre de contributions, bien qu'il présente des limites pour notre cas d'application, reste toutefois un indicateur d'influence de la fiabilité des contributeurs.

4.4 Comparaison des méthodes d'évaluation de la fiabilité

Pour désigner chaque classement de contributeurs, nous utilisons la notation suivante : C_1 désigne le classement obtenu avec le nombre de contributions, C_2 le classement obtenu par moyenne des métriques de fiabilité, C_3 le classement obtenu par moyenne pondérée des métriques de fiabilité et C_4 le classement obtenu avec la méthode de décision multicritère PROMETHEE II.

a) Étude du recouvrement entre chaque quartile

Dans la comparaison des classements des contributeurs obtenus par ces différentes méthodes, il ne s'agit pas de vérifier si chaque contributeur garde exactement le même rang dans ces classements. En effet, puisque les méthodes d'évaluation de la fiabilité sont différentes, celles-ci vont certainement entraîner des variations dans les classements des contributeurs. La Figure II.41 illustre cette idée, où trois contributeurs $user_1$, $user_k$ et $user_n$ sont positionnés à des rangs différents dans les quatre classements. Nous cherchons plutôt à observer si, d'un classement à l'autre, les petits contributeurs – en termes de nombre de contributions – sont classés comme peu fiables et si les grands contributeurs sont classés comme fiables.

Il s'agit de comparer la composition des quartiles selon les différents classements. Notons $Q_k(C_l)$ la composition du quartile Q_k selon le classement C_l . Nous étudions

respectivement le recouvrement de :

- $Q_1(C_1)$ dans tous les quartiles de C_2 , C_3 et C_4 ;
- $Q_4(C_1)$ dans tous les quartiles de C_2 , C_3 et C_4 .

Chaque classement C_k des 316 contributeurs est divisé en quatre parties égales, autrement dit en quartiles, notées Q_1 à Q_4 . Ainsi, dans le classement C_1 , le quartile Q_1 (resp. Q_4) correspond au premier quartile contenant les 79 contributeurs dont le nombre de contributions est le plus faible (resp. le plus important). Dans les classements C_2 et C_3 , le quartile Q_1 (resp. Q_4) contient les 79 contributeurs de score moyen de fiabilité le plus bas (resp. de score moyen de fiabilité le plus haut). De même, dans le classement C_4 , le quartile Q_1 (resp. Q_4) contient les 79 contributeurs classés par PROMETHEE II comme étant les moins fiables (resp. les plus fiables). En d'autres termes, nous cherchons à comprendre comment sont classés les plus petits contributeurs (resp. les plus grands contributeurs) dans les autres classements.

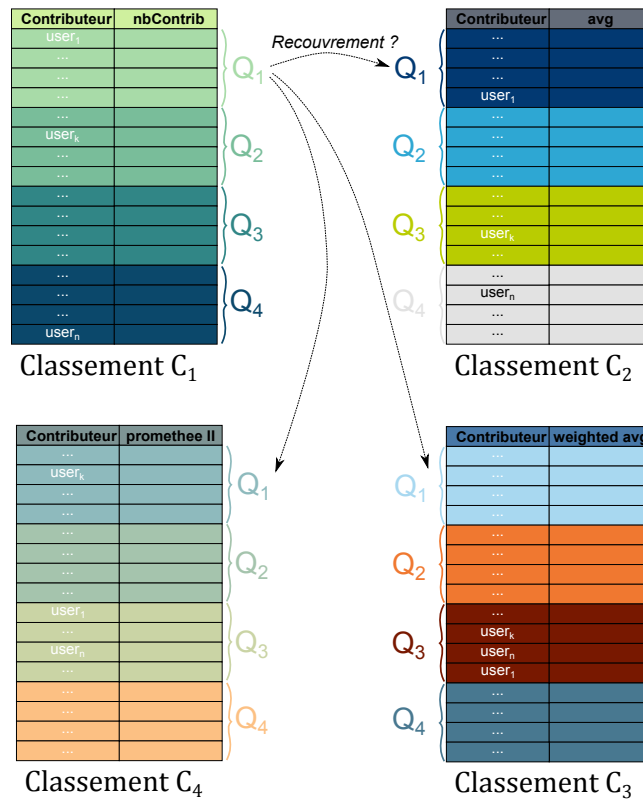


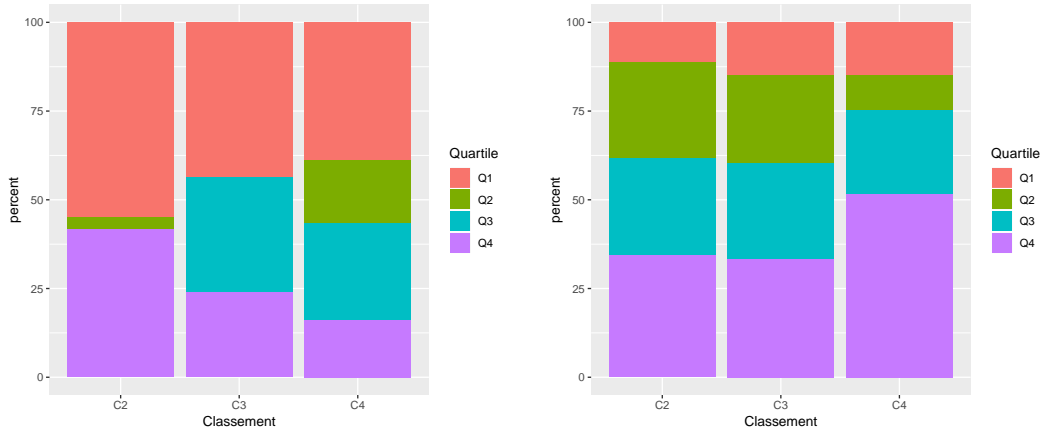
FIGURE II.41. Comparaison des différents classements par quartile

b) Résultats

La Figure II.42 donne la répartition des plus grands contributeurs et des plus petits contributeurs (selon le classement C_1) dans les trois autres classements. Les quartiles étant déterminés à partir des valeurs rangées dans l'ordre croissant, Q_1 correspond aux contributeurs classés non fiables (fin du classement) alors que les contributeurs fiables se trouvent dans Q_4 (tête du classement).

(i) Comparaison des classements C_1 et C_2 :

- Figure II.42a : plus de la moitié (55%) des petits contributeurs de $Q_1(C_1)$ se trouvent aussi dans $Q_1(C_2)$. Ce recouvrement indique le faible niveau de confiance de certains petits contributeurs, qui semble se confirmer à travers ces deux classements. Toutefois, on observe que 42% des contributeurs $Q_1(C_1)$ se trouvent dans $Q_4(C_2)$, c'est-à-dire en tête du classement C_2 . Cette portion est proche de celles des petits contributeurs fiables que nous avons qualifiés manuellement sur l'échantillon des 30 contributeurs dans la partie 4.3 précédente. Bien qu'il reste à vérifier que la méthode du classement C_2 évalue correctement les petits contributeurs comme fiables, nous pouvons dire, pour l'instant, que cette méthode permet de favoriser certains petits contributeurs dans l'évaluation de leur fiabilité.
- Figure II.42b : parmi les grands contributeurs de $Q_4(C_1)$, 11% sont classés dans $Q_1(C_2)$. En revanche, seulement 34% de ces grands contributeurs sont classés dans $Q_4(C_2)$. La méthode du classement C_2 permet de nuancer l'évaluation des contributeurs prolifiques, voire même d'en classer certains dans la catégorie des contributeurs les moins fiables de la communauté étudiée.



(a) Répartition des petits contributeurs de $Q_1(C_1)$ (b) Répartition des grands contributeurs de $Q_4(C_1)$

FIGURE II.42. Comparaison du classement C_1 avec les autres classements

(ii) **Comparaison des classements C_1 et C_3 :**

- Figure II.42a : 43.6% des contributeurs de $Q_1(C_1)$ se trouvent aussi dans $Q_1(C_3)$. Près de la moitié des petits contributeurs sont rangés dans la catégorie des contributeurs les moins fiables selon le classement C_3 . Cette observation est comparable au recouvrement observé entre $Q_1(C_1)$ et $Q_1(C_2)$, à la légère différence que le C_3 surclasse légèrement les petits contributeurs par rapport à C_2 . Par ailleurs, aucun contributeur de $Q_1(C_1)$ ne se trouve dans $Q_2(C_3)$, mais 24% d'entre eux sont classés en tête dans $Q_4(C_3)$. Cette proportion de petits contributeurs classés comme fiables est similaire à celle observée dans l'échantillon des 30 contributeurs qualifiés manuellement (Table II.11), ce qui peut donner du poids à la méthode d'évaluation de C_3 .
- Figure II.42b : les grands contributeurs de $Q_4(C_1)$ se trouvent répartis de manière globalement homogène dans $Q_4(C_2)$ (24.7%), $Q_4(C_3)$ (27.2%). Un tiers d'entre eux (33.3%) se trouve en tête dans $Q_4(C_4)$, et une faible proportion

se trouve dans $Q_4(C_1)$ (14.8%). De manière similaire à C_2 , la méthode d'évaluation du classement C_3 permet aussi de nuancer le niveau de confiance de certains grands contributeurs. Les pondérations utilisées dans le calcul de ce score moyen permettent de sous-classer plus de grands contributeurs dans la catégorie des contributeurs les moins fiables.

(iii) Comparaison des classements C_1 et C_4 :

- Figure II.42a : 38.7% des petits contributeurs de $Q_1(C_1)$ se trouvent aussi dans $Q_1(C_4)$. Ce recouvrement observé est le plus faible, en comparaison avec le recouvrement de $Q_1(C_1)$ respectivement avec $Q_1(C_2)$ et $Q_1(C_3)$. De plus, seulement 16.1% des petits contributeurs se trouvent dans $Q_4(C_1)$, le reste étant réparti dans le milieu du classement (17.7% dans $Q_4(C_2)$ et 24.7% dans $Q_4(C_3)$). La proportion de petits contributeurs classés comme très fiables est plus faible que dans les classements C_2 et C_3 . La méthode multicritère, telle que paramétrée, permet donc de surclasser les petits contributeurs à un niveau de confiance intermédiaire.
- Figure II.42b : 51.9% des contributeurs de $Q_4(C_1)$ se retrouvent en tête du classement $Q_4(C_4)$. Ce recouvrement est le plus important, en comparaison avec le recouvrement de $Q_4(C_1)$ avec respectivement $Q_4(C_2)$ et $Q_4(C_3)$. Toutefois, dans ce classement, 14.8% des grands contributeurs de $Q_4(C_4)$ sont dans $Q_4(C_1)$: la méthode PROMETHEE II permet donc de classer la même proportion de grands contributeurs comme non fiables que le classement par moyenne pondérée.

L'analyse comparative de ces classements montre comment la prise en compte d'autres métriques sur les contributeurs permet de modifier l'évaluation de leur niveau de fiabilité. En particulier, ces méthodes (agrégation des métriques, décision multicritère) permettent de classer :

- certains petits contributeurs comme fiables ;
- certains grands contributeurs comme peu fiables.

Toutefois, il convient de confronter ces classements avec une « vérité terrain », c'est-à-dire un ensemble de contributeurs dont la qualification est déjà connue.

c) Suivi des contributeurs identifiés

Pour vérifier la justesse de ces méthodes d'évaluation de fiabilité, nous observons les rangs pris par les 26 contributeurs qualifiés manuellement dans les classements C_2 , C_3 et C_4 , en ne tenant pas compte des 4 contributeurs qui ont été rangés dans la classe « ne sait pas ». Nous suivons donc le positionnement de 22 contributeurs fiables et 4 contributeurs non fiables dans les différents quartiles des classements. La Figure II.43 donne les résultats du suivi de ces contributeurs.

Dans la Figure II.43a, nous pouvons observer comment les 22 contributeurs qualifiés comme fiables se positionnent dans les différents classements. Le nombre de contributions permet de classer un maximum de contributeurs fiables (17 contributeurs, soit 77%) dans $Q_4(C_1)$, c'est-à-dire en tête du classement de C_1 , tandis que le modèle PROMETHEE II en classe un minimum (13 contributeurs, soit 59%) dans

$Q_4(C_4)$. Le nombre de contributions permet donc d'identifier avec plus de précision les contributeurs fiables, en comparaison avec les autres méthodes. Toutefois, nous remarquons également que le modèle PROMETHEE II permet de classer un maximum de contributeurs fiables dans la première moitié du classement C_4 (19 contributeurs, soit 86%)

Dans la Figure II.43b, nous observons où se positionnent les 4 contributeurs qualifiés comme non fiables dans les différents classements. Comme il a été soulevé précédemment, le nombre de contributions présente la limite de classer des contributeurs prolifiques non fiables en tête du classement C_1 . Les méthodes d'évaluation de la fiabilité par moyenne simple des métriques contributeurs permettent de classer 75% des contributeurs non fiables dans $Q_1(C_2)$, c'est-à-dire en fin de classement de C_2 . Le modèle PROMETHEE II classe un minimum de contributeurs non fiables en fin de classement, toutefois, nous remarquons que cette méthode ne range aucun de ces contributeurs dans la première moitié du classement, là où se trouvent les contributeurs fiables. En conséquence, les méthodes d'agrégation de métriques et PROMETHEE II évaluent de manière plus précise les contributeurs moins fiables, car elles parviennent à les ranger en fin de classement.

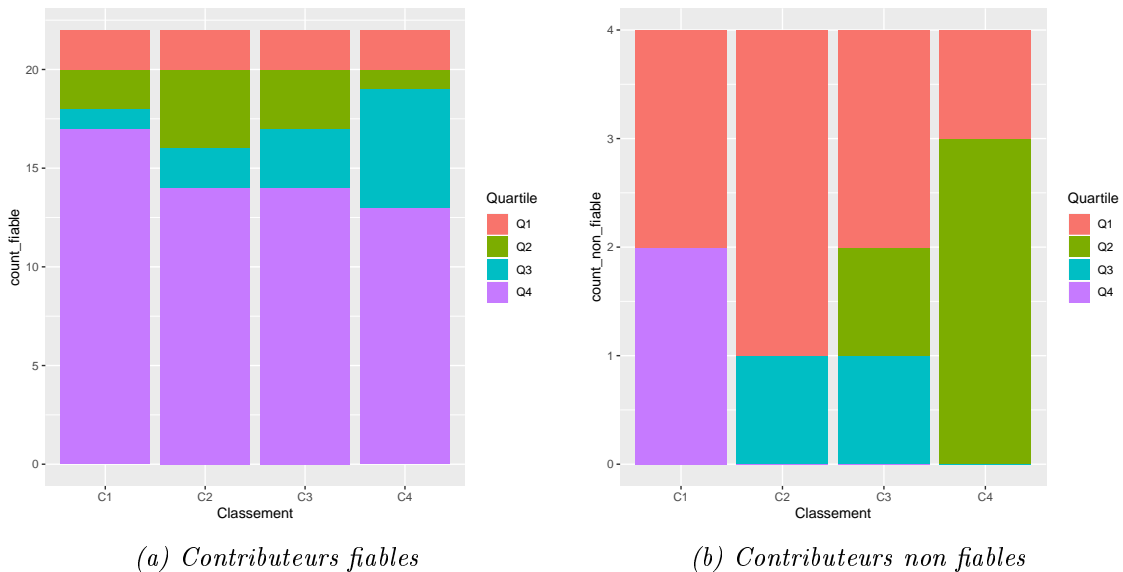


FIGURE II.43. Répartition des contributeurs qualifiés manuellement dans les classements

La Table II.13 décrit plus finement le positionnement des contributeurs qualifiés manuellement qui n'ont produit qu'une seule contribution sur la zone. Nous pouvons voir que les méthodes d'agrégation des métriques contributeurs et PROMETHEE II permettent de ranger un petit contributeur en tête de classement (dans Q_4). En revanche, contrairement à C_1 et C_2 , les classements C_3 et C_4 autorisent le surclassement d'un petit contributeur non fiable : celui-ci se trouve respectivement dans $Q_3(C_3)$ et $Q_2(C_4)$. Par conséquent, en ce qui concerne l'évaluation de la fiabilité des petits contributeurs, le classement C_2 semble être le plus juste par rapport aux autres classements.

TABLE II.13. *Positionnement des petits contributeurs dans les quartiles de C_1 à C_4 .*

Classement C_k	C_1	C_2	C_3	C_4
Fiables (total : 2)	Q_1	Q_1, Q_4	Q_1, Q_4	Q_1, Q_4
Non-fiables (total : 2)	Q_1	Q_1	Q_1, Q_3	Q_1, Q_2

4.5 Discussion et conclusions de l'évaluation de la fiabilité des contributeurs

Dans cette partie, notre étude sur la fiabilité des contributeurs s'est appuyée sur un échantillon de contributeurs qualifiés manuellement. Or, nous avons conscience que la qualification a été faite sur un petit échantillon (30 contributeurs OSM). Pour mieux rendre compte du potentiel des méthodes d'évaluation proposées (par moyenne des métrique ou par PROMETHEE II), il faudrait, dans l'idéal, qualifier plus de contributeurs.

Il faudrait également s'intéresser au paramétrage des modèles utilisés. En effet, le classement des contributeurs avec PROMETHEE II dépend de la définition de chaque fonction de préférence P_j . Dans notre cas, nous avons choisi de faire varier les préférences selon une fonction gaussienne, mais il existe d'autres types de fonctions de préférence. Par exemple, nous aurions pu choisir de faire varier linéairement les préférences. Par ailleurs, les seuils des fonctions choisies ici ont été fixés arbitrairement. Pour optimiser le classement, il faudrait générer des classements supplémentaires en faisant varier les seuils des fonctions de préférence. De même, la pondération des fonctions de préférence dans PROMETHEE II et celle des métriques de fiabilité dans le calcul de la moyenne pondérée pourraient permettre d'améliorer l'évaluation de la fiabilité des contributeurs.

Ce travail a permis de montrer qu'en réduisant la fiabilité des contributeurs à la quantité de données produites, nous risquons de ne pas détecter les carto-vandales parmi les grands contributeurs de données. Bien que le nombre de contributions soit un bon indicateur de fiabilité, il n'est pas suffisant pour qualifier précisément les contributeurs non fiables. Par conséquent, l'intégration de métriques supplémentaires a permis d'améliorer cette évaluation.

L'analyse des différentes méthodes d'évaluation de la fiabilité des contributeurs à permis de :

- montrer l'existence de grands contributeurs non fiables et de petits contributeurs fiables ;
- tenir compte de nouvelles métriques de fiabilité ;
- mettre en place des techniques qui tiennent compte de ces métriques, et qui permettent de corriger les limites de l'évaluation la fiabilité selon le nombre de contributions.

Conclusion du chapitre

Dans ce chapitre, nous avons vu qu'il existe une complémentarité entre l'approche participative, définie comme la capacité de la communauté de contributeurs à corriger les différentes contributions des uns et des autres, et l'approche sociale, qui assure la qualité des données par la présence de contributeurs de confiance ayant un rôle de modérateur. Alors que la première approche porte sur le nombre de contributeurs, la seconde s'intéresse à la qualité de chaque contributeur. Or, nous avons vérifié expérimentalement que le nombre de contributeurs ne suffit pas toujours pour assurer la qualité des données géographiques collaboratives, il faut encore pouvoir qualifier leur comportement de participation.

Nous avons étudié le comportement collaboratif des contributeurs à travers les différentes opérations d'édition qui leur sont permises pour saisir des données. Pour cela, la construction d'un réseau multiplexe de graphes d'interactions a permis de prendre en compte dans un seul modèle les différentes collaborations qui peuvent se produire sur différents plans. L'analyse de ce réseau multiplexe a mis en évidence des communautés particulières de contributeurs, parmi lesquelles nous avons identifié des profils de type modérateurs et pionniers, permettant de rendre compte de la qualité de leurs contributions.

Nous avons par la suite montré que la fiabilité d'un contributeur ne dépendait pas uniquement du nombre de ses contributions. En expérimentant des méthodes d'évaluation de la fiabilité qui prennent en compte des métriques issues des graphes d'interaction en plus du nombre de contributions, nous sommes parvenus à qualifier de manière plus précise les contributeurs non fiables, dont les agissements peuvent relever du carto-vandalisme. Ces résultats nous incitent donc à considérer plusieurs métriques de fiabilité du contributeur pour détecter le carto-vandalisme, comprenant non seulement le nombre de contributions, mais aussi des métriques de collaboration (modification, suppression, utilisation de la part d'autres contributeurs), l'investissement temporel, et la focalisation spatiale de la participation.

Chapitre III

Détection du carto-vandalisme par apprentissage automatique

Ce chapitre porte sur la détection du carto-vandalisme par des méthodes d'apprentissage automatique. Après avoir présenté les grandes familles d'apprentissage automatique, nous dressons un état de l'art sur les métriques et les méthodes d'apprentissage utilisées pour traiter des problématiques similaires. Puis, nous présentons les expériences que nous avons menées pour détecter le carto-vandalisme dans les données OSM à partir de différentes méthodes d'apprentissage. L'analyse de nos résultats expérimentaux permet d'enrichir notre définition théorique du carto-vandalisme (développée au Chapitre 1) et d'étudier le potentiel de ces méthodes d'apprentissage automatique à détecter le vandalisme cartographique.

1 Modes d'apprentissage automatique

L'apprentissage automatique (*machine learning*) est un sous-domaine de l'intelligence artificielle dans lequel un système – une machine – apprend à réaliser des tâches à partir d'observations sur des données. La plupart du temps, l'apprentissage automatique sert à effectuer des tâches de classification. Dans cette partie, nous expliquons le principe des deux grandes familles d'apprentissage automatique – à savoir l'apprentissage non-supervisé et supervisé – appliqué à la problématique de détection du vandalisme dans les données géographiques collaboratives.

1.1 Apprentissage non-supervisé

L'apprentissage non-supervisé est une approche de classification qui consiste à regrouper les données à partir de leurs variables descriptives (*clustering*). Le principe des algorithmes d'apprentissage non-supervisé est de regrouper dans une même classe les données qui présentent des caractéristiques similaires, selon les variables descriptives considérées. Le choix des variables descriptives sera donc déterminant pour la classification non-supervisée. L'approche non-supervisée est exploratoire, c'est-à-dire qu'elle ne nécessite pas de connaître *a priori* la classe des données à prédire, mais elle permet de mettre en évidence les structures présentes dans un jeu de données à classer.

Dans le cadre de la détection du vandalisme des données géographiques collaboratives, l'utilisation d'une telle méthode consiste à donner en entrée de l'algorithme un ensemble de métriques qui décrivent chaque donnée géographique collaborative à qualifier. L'intérêt ici est de pouvoir séparer les contributions qui relèvent du carto-vandalisme du reste du jeu de données. Cela suppose de connaître les caractéristiques de ce phénomène et de pouvoir les modéliser par des métriques pertinentes.

Le carto-vandalisme peut prendre plusieurs formes, comme nous l'avons vu au chapitre 1 (voir la partie 2.4 sur les typologies de carto-vandalisme). L'utilisation d'une méthode non-supervisée pour classer les données selon les différentes classes d'une typologie de carto-vandalisme – par exemple celle qui a été proposée par [Ballatore \(2014\)](#) – suppose de décrire les données géographiques par des métriques qui reflètent correctement chacune de ces catégories de carto-vandalisme. Par exemple, pour détecter le carto-vandalisme de type artistique, il serait intéressant de mettre en place des métriques géométriques (comme l'aire, le périmètre, *etc.*) afin d'identifier les contributions qui présentent des formes géométriques inhabituelles. Par ailleurs, le carto-vandalisme étant caractérisé par des contributeurs non fiables, il convient de tenir compte de métriques indiquant le niveau de confiance des contributeurs des données à qualifier. Toutefois, comme notre connaissance concernant les caractéristiques du carto-vandalisme et les différentes formes qu'elles peuvent prendre sont limitées, nous ne connaissons pas *a priori* toutes les métriques à mettre en place pour le détecter par une méthode non-supervisée.

1.2 Apprentissage supervisé

Contrairement à l'approche non-supervisée où le modèle regroupe des données sans connaître *a priori* les classes en sortie de la classification, l'approche supervisée consiste à construire un modèle capable de regrouper les données à prédire dans des classes pré-définies, à partir des variables descriptives sur ces données. Dans le cas de la détection du carto-vandalisme, un modèle de classification supervisé pourra classer les données selon deux classes possibles : la première contiendra toutes les données de carto-vandalisme, la seconde toutes les données qui n'en sont pas.

La construction d'un modèle d'apprentissage supervisé comporte deux étapes : l'entraînement et l'évaluation. Ces deux étapes nécessitent de disposer d'un jeu de données annotées, c'est-à-dire des données dont on connaît par avance la classification (carto-vandalisme ou non). L'étape d'entraînement consiste à donner en entrée du modèle une partie du jeu de données annotées, représentées par leurs variables descriptives ainsi que le label de la classe à laquelle elles appartiennent. Le modèle va apprendre à classer les données à partir de ces exemples d'entraînement. Dans notre cas, il s'agira de donner à un modèle de classification supervisé des exemples de données de carto-vandalisme et des données normales.

Puis, l'étape de test consiste à utiliser le modèle entraîné pour prédire la classe des données annotées qui n'ont pas servi pendant la phase d'entraînement. Autrement dit, il s'agit de vérifier si le modèle entraîné parvient à différencier les données annotées comme carto-vandalisme des données normales. Le résultat de la détection de ces données de test permettra alors d'évaluer la performance de classification du modèle. L'évaluation d'un modèle d'apprentissage peut être effectuée en calculant des indicateurs de performance. La Figure III.1 schématise les différentes étapes de l'apprentissage supervisée.

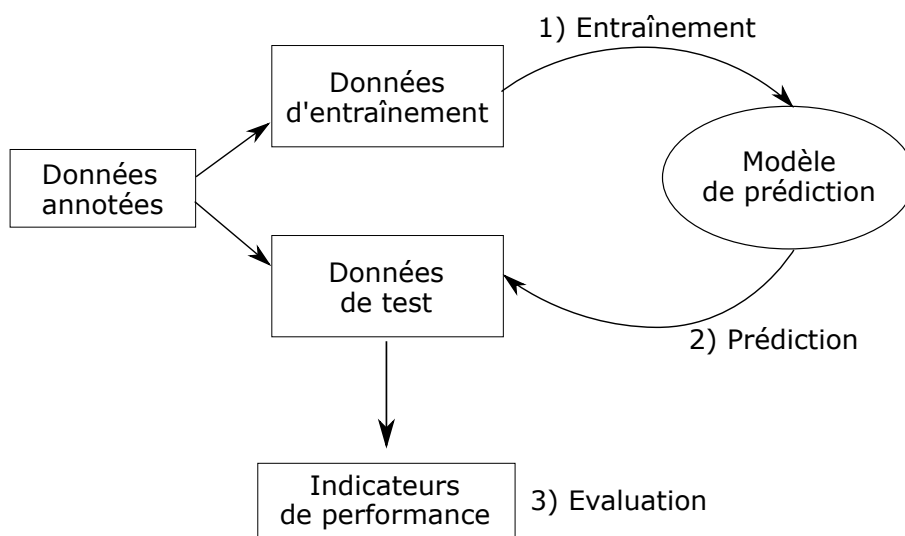


FIGURE III.1. Démarche de l'apprentissage supervisé

1.3 Indicateurs de performance

Dans notre situation, la tâche de classification consiste à distinguer les contributions qui relèvent du carto-vandalisme et celles qui n'en sont pas. Par conséquent, les indicateurs de performance vont permettre de qualifier la qualité de la détection du carto-vandalisme par un système d'apprentissage. Nous considérons les indicateurs de précision, de rappel et d'erreur pour évaluer nos méthodes d'apprentissage.

La précision est à entendre ici comme la précision de la détection du vandalisme. Elle se définit par la formule suivante :

$$precision = \frac{TP}{TP + FP} \quad (\text{III.1})$$

où TP est le nombre de vrais positifs et FP le nombre de faux positifs. Les vrais positifs sont les contributions de carto-vandalisme avéré qui ont été détectées alors que les faux positifs sont détectés par le modèle alors qu'ils ne relèvent pas du carto-vandalisme. Dans l'idéal, un modèle de classification doit être précis (la précision tend vers 1), c'est-à-dire qu'il détectera le carto-vandalisme en minimisant le nombre de faux positifs.

Le rappel se définit comme la part de vrais positifs (*true positive rate*), à savoir :

$$rappel = \frac{TP}{TP + FN} \quad (\text{III.2})$$

où FN est le nombre de faux négatifs, c'est-à-dire les cas avérés de carto-vandalisme qui n'ont pas été classés comme tels par le système. Dans l'idéal, un modèle de prédiction doit détecter tous les cas de carto-vandalisme sans se tromper, c'est-à-dire en minimisant le nombre de faux négatifs. Un rappel maximal vaut 1.

L'erreur correspond au pourcentage global de données mal classées :

$$erreur = \frac{FP + FN}{TP + TN + FP + FN}$$

où TN est le nombre de vrais négatifs, c'est-à-dire les contributions de non vandalisme qui ont été prédites comme tel par le modèle. Contrairement au rappel et à la précision, un modèle de prédiction doit, dans l'idéal, minimiser son erreur de classification vers 0.

2 Apprentissage du vandalisme dans les bases de données ouvertes

Cette partie dresse un état de l'art des métriques mises en place et des méthodes d'apprentissage utilisées pour détecter automatiquement le vandalisme dans Wikipédia. À partir des contributions faites dans la littérature sur ce sujet, nous cherchons à reprendre ou réadapter des métriques ou des méthodes qui seraient transposables dans le cas du carto-vandalisme.

Par ailleurs, de nombreux exemples sont donnés sur la page Wikipédia dédiée au sujet du vandalisme pour le distinguer des cas d'erreur ou de non-respect des

règles de bonne conduite (participer à des controverses, avoir un comportement obstiné, *etc.*). En particulier, les sujets controversés n'étant pas considérés comme du vandalisme, Kittur *et al.* (2007) ont développé une métrique à partir des articles annotés comme « controversés » pour évaluer le niveau de conflit sur Wikipédia.

2.1 Métriques et méthodes

La Table III.1 récapitule les méthodes et les métriques proposées dans l'état de l'art pour détecter le vandalisme dans les bases de données ouvertes. L'information contributive dans les bases de données ouvertes de type wiki est principalement sous forme textuelle. Par conséquent, les méthodes de détection du vandalisme sur ces bases de données ouvertes comprennent toujours des métriques textuelles. Adler *et al.* (2011) proposent un certain nombre de métriques linguistiques, ce que Heindorf *et al.* (2016) appellent des descripteurs de contenu (*content features*). Ceux-ci portent sur le contenu textuel à l'échelle du caractère, du mot ou de la phrase. Parmi ces descripteurs de contenu, Mola Velasco (2011) distingue ceux qui sont liés au langage, situés au niveau des mots, et ceux qui ne le sont pas, plutôt situés au niveau des caractères. Les descripteurs linguistiques quantifient la fréquence et l'impact de termes vulgaires. En revanche, les métriques textuelles non linguistiques se concentrent notamment sur la portion de caractères spéciaux, de majuscules ou de chiffres.

Dans les exemples réels de vandalisme que nous avons présentés dans le Chapitre 1, certains cas de carto-vandalisme consistent justement à éditer les tags des objets, qui sont des données textuelles. Nous avons pu identifier le cas où un tag `name=)` a été introduit pour nommer un lac par un *smiley*. Par ailleurs, comme les valeurs de tags sont des champs de texte libre, les carto-vandales peuvent potentiellement y insérer du texte injurieux, grossier ou encore des liens vers des sites pornographiques (Ballatore, 2014). Dans le cas du projet OSM, les contributeurs ont la possibilité de commenter leurs sessions d'édition ainsi que celles des autres contributeurs. Ces commentaires sont également des informations textuelles qui, potentiellement, peuvent indiquer des cas de vandalisme. Par conséquent, les métriques textuelles proposées dans la littérature peuvent être exploitées pour détecter le carto-vandalisme portant sur les tags des contributions cartographiques et sur les commentaires des sessions d'édition. Ces métriques peuvent aussi bien être syntaxiques (*i.e.* au niveau des caractères textuels) que sémantiques (*i.e.* au niveau du sens des mots).

Le vandalisme n'étant pas réduit à un acte de dégradation, il faut aussi prendre en compte l'aspect contextuel de la dégradation. En particulier, l'auteur de l'édition (Heindorf *et al.*, 2016; Sarabadani *et al.*, 2017) constitue un élément contextuel important puisque son évaluation permet de déterminer son niveau de confiance et par conséquent, sa tendance à provoquer du vandalisme. La réputation des contributeurs a notamment été utilisée à cet effet (Adler *et al.*, 2011). Notons ici que cette prise en compte des caractéristiques du contributeur rejoint l'approche sociale de la qualité de l'information géographique, qui a été développée dans le chapitre précédent. Le développement de métriques portant sur le niveau de confiance du contributeur de données est renforcée par le fait que celles-ci ont été considérées pour détecter le vandalisme dans les bases de données Wikipédia.

Par ailleurs, l'historique des données permet également de situer le contexte temporel dans lequel une contribution a été produite (Heindorf *et al.*, 2016; Sarabadani *et al.*, 2017). Il peut notamment mettre en évidence des situations de controverses, qui ne sont pas forcément du vandalisme, mais qui correspondent à des situations de forts désaccords entre des contributeurs pouvant déclencher par la suite des actes de vandalisme. Il est donc intéressant de pouvoir identifier ces cas de controverses. À ce sujet, Kittur *et al.* (2007) ont développé une métrique à partir des articles annotés comme « controversés » pour évaluer le niveau de conflit sur Wikipédia. Les situations de conflits et de controverses peuvent également se produire sur les projets d'information géographique volontaire, en particulier sur des territoires où le contexte géopolitique n'est pas clair. Par conséquent, nous pourrions également considérer des métriques d'historiques pour détecter les contributions de carto-vandalisme provoquées dans ce type de conflits.

Certaines métriques sont proposées pour détecter des formes spécifiques de vandalisme. En effet, pour détecter des actes de vandalisme sur Wikipédia tels que le spam, la désinformation et la suppression massive d'information, Chin *et al.* (2010) adoptent des modèles statistiques linguistiques. Pour cela, ils proposent des indicateurs statistiques comme métriques de détection du vandalisme. Sarabadani *et al.* (2017) développent des métriques pour détecter les motifs particuliers de vandalisme qui ont été identifiés par la communauté de contributeurs de Wikidata : par exemple la modification de la date de naissance, la nationalité ou le sexe d'une personnalité renseignée dans la base. De la même manière, des métriques de non-vandalisme ont été mis en place pour éviter de confondre le vandalisme avec des modifications drastiques, comme la restauration de données ou la création d'un nouvel item. La restauration d'une donnée ou l'annulation d'une version est une opération qui peut également être effectuée dans les projets comme OSM. En particulier, les actes de carto-vandalisme repérés par les membres du DWG d'OSM subissent des opérations de restauration/d'annulation, par conséquent, il serait intéressant de développer des métriques qui tiennent compte de ces types d'opération. Toutefois, le monde étant en constante évolution, il peut arriver que des données cartographiques soient mises à jour de manière plus ou moins importante. La suppression de données cartographiques obsolètes peut se produire fréquemment. Par conséquent, il faudra tenir compte de ces situations pour développer ce type de métriques.

Dans les travaux cités, la détection du vandalisme des données wiki est réalisée grâce à des méthodes d'apprentissage supervisé. En effet, l'existence de plusieurs corpus de vandalisme pour les données Wikipédia et Wikidata, c'est-à-dire un jeu de données annotées, permet d'entraîner un modèle d'apprentissage et d'évaluer ses performances de classification. Parmi les algorithmes d'apprentissage supervisé, l'utilisation des forêts aléatoires (*random forests*) a démontré une certaine efficacité pour détecter le vandalisme dans les données wiki (Adler *et al.*, 2011; Heindorf *et al.*, 2016; Mola Velasco, 2011). Chin *et al.* (2010) ont obtenu de bons résultats en utilisant d'autres méthodes comme les arbres de décision, la régression logistique et les Séparateurs à Vaste Marge (SVM). Quel que soit le choix de l'algorithme, la détection supervisée du vandalisme dans les bases de données ouvertes nécessite le développement de métriques pertinentes – *i.e.* des caractéristiques du vandalisme des données – et de disposer d'un corpus de vandalisme. De plus, pour étudier le potentiel des méthodes d'apprentissage non-supervisé pour détecter le carto-vandalisme, nous

avons également besoin de données annotées de carto-vandalisme. Par conséquent, un corpus de données annotées pour le carto-vandalisme est utile quelle que soit la méthode d'apprentissage considérée. Or, comme un tel corpus de données annotées n'existe pas encore, nous discutons ultérieurement des méthodes de construction des corpus de vandalisme dans la Section 4 de ce chapitre, dans le but d'en construire un.

Référence	Métriques utilisées	Méthodes d'apprentissage
Chin <i>et al.</i> (2010)	<ul style="list-style-type: none"> — Métriques statistiques linguistiques — Métriques portant sur la révision 	<ul style="list-style-type: none"> — Arbre de décision (J48) — Régression logistique — Séparateur à vastes marges (SVM)
Adler <i>et al.</i> (2011)	4 catégories de métriques : métadonnées ; texte ; réputation ; langage.	Forêts aléatoires (<i>random forests</i>)
Mola Velasco (2011)	Métriques textuelles : linguistiques ; indépendantes du langage	Forêts aléatoires (<i>random forests</i>)
Heindorf <i>et al.</i> (2016)	<ul style="list-style-type: none"> — Métriques de contenu — Métriques de contexte 	Forêts aléatoires (<i>random forests</i>)
Sarabadani <i>et al.</i> (2017)	<ul style="list-style-type: none"> — Métriques générales — Motifs typiques de vandalisme — Motifs types de non-vandalisme — Caractéristiques contributeurs 	Web service ORES, développé par <i>Wikimedia scoring</i> , qui fournit des classifieurs supervisés pour notamment détecter le vandalisme.
Martinez-Rico <i>et al.</i> (2019)	Métriques linguistiques ; métriques issues de l'état de l'art	Forêts aléatoires (<i>random forests</i>)

TABLE III.1. Méthodes et métriques utilisées dans la littérature pour détecter le vandalisme dans les bases de données ouvertes

2.2 Cas du vandalisme cartographique

Contrairement aux données wiki, la détection automatique du vandalisme dans les données cartographiques a été très peu étudiée dans la littérature scientifique. Le prototype le plus abouti pour détecter le vandalisme dans les données OSM a été proposé par Neis *et al.* (2012). Ce système, appelé *OSMPatrol* est construit à partir d'un ensemble de règles de décision visant à détecter les actes de dégradation dans OSM. En effet, Neis *et al.* (2012) adoptent une définition du vandalisme cartographique plus large que celle que nous avons fixée pour le carto-vandalisme : selon leur définition, le vandalisme englobe tous les actes de dégradation, aussi bien intentionnels qu'accidentels.

Les règles de décisions prennent en considération la réputation du contributeur et le type d'action réalisé pour produire sa contribution. En effet, *OSMPatrol* s'appuie sur une exploration manuelle d'incidents OSM provenant de 204 contributeurs bannis par le DWG, permettant de relever 51 cas de vandalisme cartographique (selon la définition de Neis *et al.* (2012)), parmi lesquels se trouvent :

- 33.3% cas de carto-vandalisme fantaisiste,
- 33.3% cas de modifications géométriques anormales,
- 43.1% cas de suppression de données.

La définition de Neis *et al.* (2012) permet donc de détecter certains cas de carto-vandalisme. Par ailleurs, leur analyse montre que 76.4% des cas relevés sont issus de nouveaux contributeurs. *OSMPatrol* évalue donc la réputation des contributeurs, de manière à défavoriser les débutants qui, à cause de leur inexpérience, sont plus susceptibles de contribuer des données de faible qualité. Bien que cette méthode soit discutable, nous notons qu'elle prend en compte le profil du contributeur pour qualifier les contributions.

Toutefois, les résultats présentés par Neis *et al.* (2012) révèlent que le système détecte plus de cas de vandalisme qu'il n'y en a réellement. En effet, 44% des contributions OSM testées par *OSMPatrol* sont identifiées comme vandalisme potentiel. Les auteurs déclarent que, parmi ces résultats de détection, il est possible de trouver des cas avérés de carto-vandalisme. Le système n'a donc pas été conçu pour détecter avec précision le carto-vandalisme. À ce sujet, les performances de ce prototype n'ont pas été évaluées car il n'existe pas de corpus de données de carto-vandalisme OSM. Par conséquent, nous ne pouvons pas savoir si *OSMPatrol* est capable de détecter entièrement le carto-vandalisme dans OSM ou s'il existe des cas indétectables par ce modèle.

3 Démarche expérimentale

Les expériences que nous menons sont effectuées selon une démarche purement exploratoire. En effet, une définition du carto-vandalisme a été présentée au Chapitre 1 et constitue une grande partie de notre contribution à ce sujet. Toutefois, celle-ci reste principalement théorique, c'est-à-dire que nous ne sommes pas encore capables de mesurer la différence entre une donnée de mauvaise qualité et une donnée de carto-vandalisme. Bien que nous ayons relevé des exemples réels de carto-vandalisme

dans OSM (voir au Chapitre 1), nous ne sommes pas certains que ces exemples soient suffisamment nombreux et représentatifs de toutes les formes de vandalisme cartographique qui peuvent se trouver dans la base OSM.

La définition du carto-vandalisme mérite donc d'être enrichie par une analyse sur un plus grand nombre de données géographiques réelles, pour mieux comprendre ce qu'est le vandalisme cartographique et comment il se présente concrètement dans une base de données géographiques ouverte. Les expériences menées n'ont donc pas pour but de trouver la meilleure technique de détection du carto-vandalisme, puisque nous cherchons encore à définir empiriquement ce phénomène. En revanche, nous espérons que l'exploration de ces différentes méthodes mette en lumière le potentiel qu'elles renferment pour traiter notre problème. Ces résultats d'analyse permettront d'envisager, dans des travaux futurs, de tirer le meilleur parti de ces méthodes pour développer des processus de détection du carto-vandalisme.

L'expérimentation d'une méthode d'apprentissage non-supervisé permet de révéler les structures présentes dans un jeu de données à classer (voir Section 1.1 précédente). Par conséquent, en classifiant un jeu de données géographiques qui contient du carto-vandalisme par une méthode non-supervisée, nous cherchons à identifier les éléments caractéristiques du carto-vandalisme. Plus précisément, il s'agit d'identifier les métriques qui révèlent le mieux le carto-vandalisme des données, c'est-à-dire les descripteurs qui permettent de les discriminer du reste du jeu de données.

Comme les contributions qui relèvent du carto-vandalisme sont peu nombreuses, nous commençons par adopter l'hypothèse selon laquelle ces données sont considérées comme des anomalies à détecter. Par exemple, le carto-vandalisme artistique correspond à des contributions qui présentent des formes géométriques inhabituelles. Par conséquent, ces données seront vues comme des anomalies en considérant des indicateurs géométriques de forme. Le carto-vandalisme qui consiste à insérer des propos grossiers peut être vu comme une anomalie selon un point de vue sémantique. Par exemple, si l'on étudie la sémantique des tags des données OSM, et en particulier la quantité de termes violents, grossiers ou offensifs qui s'y trouvent, ce type de vandalisme apparaîtra comme une anomalie car les données OSM se veulent neutres et objectives. À partir de cette première hypothèse de modélisation, il s'agira d'expérimenter une méthode non-supervisée de détection d'anomalies, où l'anomalie à détecter est la donnée géographique vandalisée.

Cependant, modéliser le carto-vandalisme comme une anomalie trouve sa limite dans la potentielle existence de carto-vandalisme « discret » à tout point de vue (*i. e.* quel que soit le descripteur), auquel cas les données correspondantes passent inaperçues par une méthode de détection d'anomalies. L'expérimentation d'une méthode d'apprentissage supervisé nous permet donc d'étudier dans quelles mesures cette méthode parvient à mieux détecter le carto-vandalisme, en particulier des contributions pernicieuses et discrètes qui ne sont pas détectées comme des anomalies. De plus, il s'agit, dans ces deux expérimentations, d'évaluer les forces et les faiblesses de ces méthodes pour la détection automatique du carto-vandalisme.

4 Construction d'un corpus de carto-vandalisme

4.1 État de l'art des méthodes de construction d'un corpus de vandalisme

À notre connaissance, il n'existe pas de corpus de données de carto-vandalisme permettant d'évaluer des modèles de détection tels que OSMPatrol, à la différence des données wiki. En effet, les travaux de recherche sur la détection du vandalisme dans les bases de données ouvertes contient de nombreuses réflexions sur les méthodes de récupération d'exemples de vandalisme pour constituer des corpus de données annotées. Dans cette partie, nous étudions les méthodes proposées dans les travaux existants en vue de constituer un corpus spécifique à notre problématique de détection du carto-vandalisme. La Table III.2 récapitule les caractéristiques de différents corpus de vandalisme issus de l'état de l'art.

Le premier corpus de vandalisme des données Wikipédia, appelé Webis-WVC-07, a été construit par annotation manuelle de 940 contributions sur Wikipédia, dont 32% sont des exemples de vandalisme. Ce premier corpus a ensuite été remplacé par le corpus PAN-WVC-10, obtenu grâce au service *Amazon's Mechanical Turk*, où 753 personnes ont été payées pour annoter un jeu de données beaucoup plus important que le premier : celui-ci contient près de 33000 exemples de contributions, dont 2391 sont des cas de vandalisme (Potthast, 2010). L'annotation de ces données par cette méthode participative repose sur un système de vote. Une contribution est donc annotée comme vandalisme parce que la majorité des annotateurs l'ont annotée comme telle. Toutefois, cette méthode n'est pas infaillible – une majorité peut avoir tort – et elle nécessite, par ailleurs, des ressources humaines et économiques, comme c'est le cas du service *Amazon's Mechanical Turk*. Par ailleurs, Potthast (2010) précise que l'utilisation du service *Amazon's Mechanical Turk* a exigé un effort de conception du formulaire d'annotation, pour que celui-ci soit assez claire et simple à remplir. En effet, une tâche d'annotation conçue de manière trop rébarbative ou trop complexe comporte le risque que les annotateurs ne la réalisent pas sérieusement, et donc pas correctement.

La proportion de cas de vandalisme dans le corpus PAN-WVC-10 est moindre que dans le premier corpus de vandalisme Wikipédia : celle-ci diminue de 32% à 7.4%. D'un côté, un corpus de données annotées qui contient une faible proportion d'exemples de vandalisme est réaliste, car le vandalisme est un phénomène rare qui affecte peu de contributions. Un tel corpus permettrait d'apprendre le caractère exceptionnel du vandalisme à un modèle de détection automatique, et donc de lui apprendre à détecter le vandalisme en faible proportion dans un jeu de données test. D'un autre côté, plus un modèle de détection est entraîné sur un grand nombre d'exemples de vandalisme, plus il sera capable de reconnaître ce qu'est une donnée de vandalisme. Dans l'idéal, il faudrait donc que le corpus contienne un grand nombre de données, de manière à ce que, malgré leur faible proportion, les exemples de vandalisme soient en quantité suffisante pour entraîner correctement les modèles de détection du vandalisme.

Or, les exemples de vandalisme étant rares, leur obtention est une étape coûteuse. Chin *et al.* (2010) ont eu recours à une technique d'apprentissage actif pour annoter

automatiquement des exemples supplémentaires de vandalisme, et permettre ainsi d'enrichir le jeu de données d'entraînement. La méthode consiste à entraîner initialement un modèle d'apprentissage supervisé sur le corpus PAN-WVC-10. À partir d'un ensemble test de données non annotées, le modèle produit un classement des contributions qui sont potentiellement du vandalisme. Un humain va alors annoter manuellement les 50 premières contributions du classement, et enrichir le corpus de données. Toutefois, les données de test utilisées dans la situation de *Chin et al.* (2010) sont des articles Wikipédia connus pour avoir été fortement vandalisés. Par conséquent, l'usage d'une technique d'apprentissage actif suppose d'avoir localisé les objets qui sont les plus susceptibles d'être vandalisés. Dans le cas de Wikipédia, il existe une page¹ listant les articles les plus sujets au vandalisme. Dans notre situation, cette méthode pourra être envisagée seulement après avoir construit un premier jeu de données annotées de carto-vandalisme, dans le but de l'enrichir.

Le corpus de vandalisme Wikipédia le plus récent est WP_Vandal. Celui-ci est de taille équivalente à PAN-WVC-10, mais il contient autant de cas de vandalisme que de cas normaux (*Martinez-Rico et al.*, 2019). WP_Vandal a été construit automatiquement en sélectionnant les révisions de Wikipédia marquées comme révoquées (*reverted*) par des contributeurs humains ou des outils automatiques. En effet, les contributeurs de Wikipédia ont la possibilité de remplir un champ « commentaires » lorsqu'ils produisent une révision d'un article, notamment pour indiquer lorsqu'ils révisent un article qui a précédemment été vandalisé. Par conséquent, les exemples de non-vandalisme correspondent à la version qui précède immédiatement la révision annotée comme vandalisme (sauf si celle-ci a déjà été annotée comme vandalisme). L'annotation des données est ici effectuée à partir de l'historique des articles vandalisés de Wikipédia : les exemples et contre-exemples du corpus WP_Vandal portent donc sur les mêmes articles. Cette méthode de récupération d'exemples est intéressante car, en entraînant un modèle sur deux versions – l'une vandalisée, l'autre non – d'un même objet, celui-ci pourra apprendre à reconnaître le vandalisme à partir de la différence entre ces deux versions.

Les corpus de vandalisme pour les données Wikidata ont été produits notamment pour la détection de vandalisme issu de contributions faites par des humains (en opposition à des contributions produites par des scripts automatiques). Le corpus WDVC-2015 contient 24 millions de données annotées (*Heindorf et al.*, 2015). Le processus d'annotation a consisté à marquer automatiquement comme vandalisme les révisions qui ont subi une opération de *rollback*. Cette opération n'est réservée qu'à des contributeurs privilégiés de Wikidata, et consiste à annuler toutes les révisions consécutives d'un même contributeur sur un item donné. Ce corpus a permis de mener une étude approfondie sur les formes de vandalisme des données Wikidata, et sur les profils des vandales.

Toutefois, l'équivalence entre une contribution de vandalisme et une contribution issue d'une opération de *rollback* n'est pas exacte. Une vérification manuelle du corpus WDVC-2015 montre que 86 % des données annotées *vandalisme* sont réellement du vandalisme : cela signifie que 14% des données annotées comme tel n'en sont pas réellement (ce sont des faux positifs). De plus, il a été observé que 62% de contributions Wikidata issues des opérations d'annulation ou de restauration sont

1. https://fr.wikipedia.org/wiki/Wikip%C3%A9dia:Pages_souvent_vandalis%C3%A9es

réellement du vandalisme, alors que celles-ci ne sont pas marquées comme tel dans le corpus (ce sont des faux négatifs). Par conséquent, ce processus d’annotation peut entraîner des erreurs et des imprécisions (*i.e.* annoter du non-vandalisme comme vandalisme et *vice versa*) lors de l’annotation de données de vandalisme.

Le manque de précision dans l’annotation des données de vandalisme n’est donc pas sans importance : l’entraînement des modèles de détection sur un jeu de données qui contient des erreurs ne va pas leur apprendre à détecter précisément du vandalisme, mais également des dégradations involontaires. Finalement, cela reviendrait à développer des outils de détection du vandalisme selon la définition de Neis *et al.* (2012). Or, l’objet de ce travail est justement de tendre à détecter avec précision les cas de vandalisme cartographique, et non les erreurs de saisie cartographique.

Une autre méthode de construction de corpus de données Wikidata a été proposée pour surmonter les limites du corpus WDVC-2015 (Sarabadani *et al.*, 2017). En considérant un ensemble de 500 000 éditions manuelles, l’annotation a consisté à labéliser comme vandalisme les contributions qui ont subi une opération de révocation (*revert*) par des contributeurs non fiables. À partir d’un échantillon du corpus ainsi construit, une étude manuelle indique que 32% des données annotées comme vandalisme n’en sont pas réellement (ce sont des faux positifs), et 1% des données annotées comme non vandalisme sont des cas réels de vandalisme (ce sont des faux négatifs). La méthode d’annotation automatique est donc encore imprécise, puisque le corpus contient des faux exemples de vandalisme. La construction de ce corpus avait pour but d’aider les modérateurs de Wikidata à détecter le vandalisme. Autrement dit, le corpus a été créé dans le cadre d’une application d’aide à la décision : un modèle appris sur ce corpus va permettre de filtrer un ensemble de contributions à destination des modérateurs de Wikidata, qui auront la charge de les qualifier par la suite. Cependant, dans le cadre d’une détection totalement automatisée du vandalisme, il faudrait concevoir un corpus de vandalisme dans lequel l’annotation des données est réalisée de manière plus exacte. Par conséquent, il nous semble plus judicieux de ne pas envisager une méthode d’annotation complètement automatique.

Nom du corpus / Référence	Méthode d'annotation	Nombre d'exemples	Nombre de cas de vandalisme	Application
Webis-WVC-07 (Pothast <i>et al.</i> , 2007)	Humaine	940 éditions	301 (32%)	Détection automatique
PAN-WVC-10 (Pothast, 2010)	Humaine	32 452 éditions	2391 (7.4%)	Détection automatique
WDVC-2015 (Heindorf <i>et al.</i> , 2015)	Automatique	24 millions d'éditions produites manuellement	103205 (0.4%)	Analyse du vandalisme
(Sarabadani <i>et al.</i> , 2017)	Automatique	500 000 éditions produites manuellement	622 (0.12%)	Aide à la décision
WP_Vandal (Martinez-Rico <i>et al.</i> , 2019)	Automatique	36315 éditions	18506 (51%)	Détection automatique

TABLE III.2. Corpus de vandalisme issus de l'état de l'art

4.2 Modélisation d'un corpus de carto-vandalisme

Un corpus de carto-vandalisme de données OSM devrait contenir des exemples et des contre-exemples de carto-vandalisme. Ces exemples correspondraient donc à des contributions pouvant avoir été produites pour dégrader (ou non) intentionnellement (ou non) l'espace cartographique. Nous proposons une modélisation d'un corpus de carto-vandalisme de données OSM dans la Figure III.2. Notre corpus est composé de fenêtres spatiales de l'espace cartographique à une date donnée. On appellera cette fenêtre spatiale un *snapshot*. Celui-ci contient tous les objets cartographiques visibles sur cette zone, sous une version valide par rapport à la date du *snapshot* : autrement dit, ce sont les contributions valides à la date du *snapshot*. Ces contributions seront labélisées comme étant du carto-vandalisme ou non. Le corpus de carto-vandalisme ainsi modélisé pourra contenir des *snapshots* de différentes zones géographiques, à différentes dates.

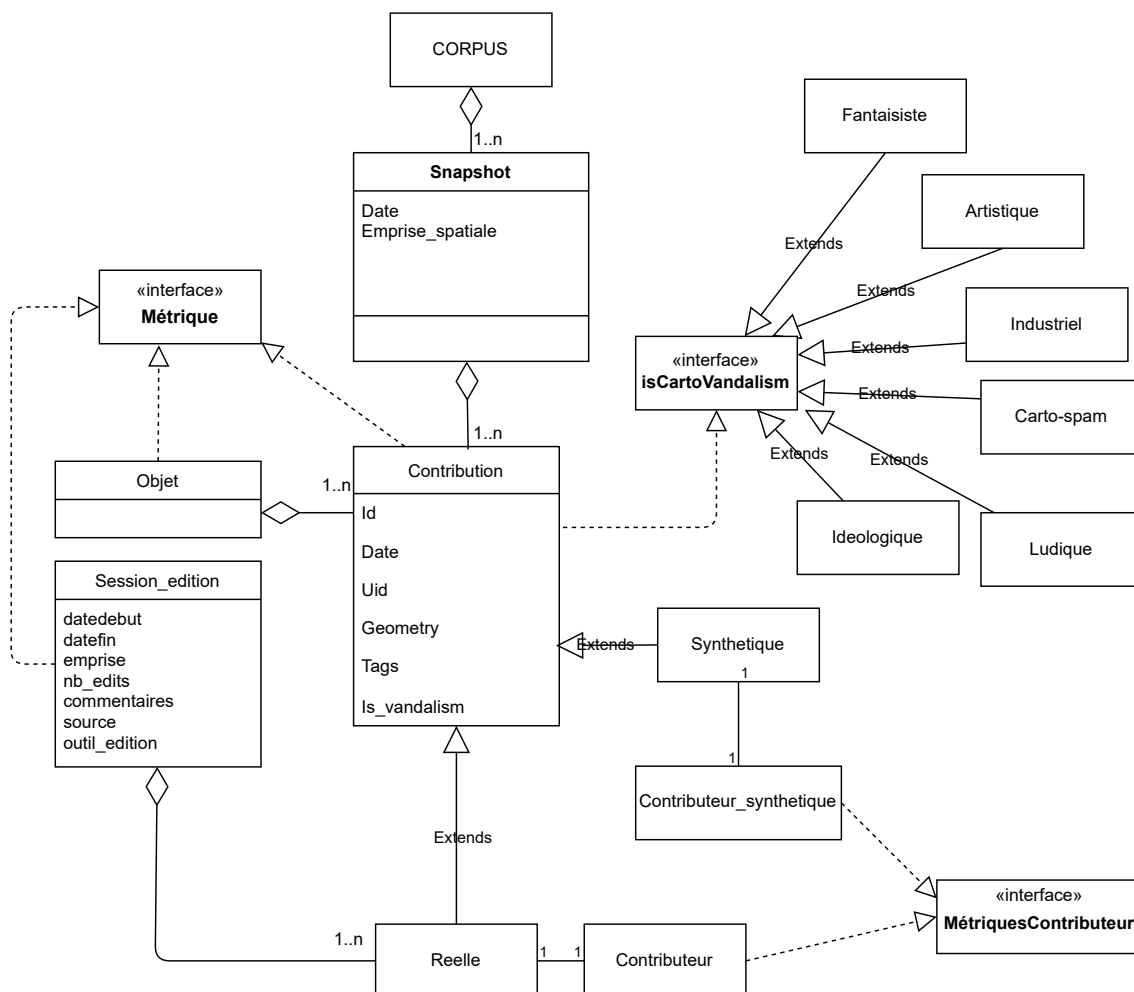


FIGURE III.2. Modélisation d'un corpus de carto-vandalisme pour les données OSM

Comme l'exploration de la page des contributeurs bannis du projet OSM ne nous a pas permis de trouver suffisamment de cas de carto-vandalisme, nous avons choisi d'insérer des cas synthétiques de carto-vandalisme. Comme nous ne connaissons *a priori* ni les lieux ni les dates auxquels les contributions de carto-vandalisme sont visibles sur l'espace cartographique, l'insertion de contributions synthétiques permet d'avoir des exemples connus de carto-vandalisme sur n'importe quel *snapshot*.

Par ailleurs, pour entraîner des modèles sur ces différentes contributions, il faut pouvoir calculer un certain nombre de métriques sur ces données. Plus concrètement, ces métriques constituent les descripteurs utilisés comme données d'entrée des algorithmes de détection du carto-vandalisme. L'interface « Métriques » réalisée par les différentes classes de la Figure III.2 montre que les descripteurs peuvent être calculés à plusieurs niveaux. En effet, les métriques calculées peuvent porter sur la contribution en elle-même, *i.e.* sous sa version courante dans le *snapshot*, ou sur l'objet cartographique qu'elle représente, ou sur la session d'édition dans laquelle elle a été produite. De plus, la contribution peut être décrite par des métriques portant sur le contributeur qui l'a produite.

Comme il s'agit de représenter au mieux le carto-vandalisme, il convient de s'interroger sur les métriques à mettre en place pour les données du corpus. Le carto-vandalisme artistique est typiquement identifiable sur des critères géométriques : il serait donc pertinent d'implémenter des indicateurs géométriques quantifiant la taille et la forme des contributions cartographiques. Le carto-vandalisme fantaisiste peut se produire en ajoutant des données spatiales par-dessus des éléments cartographiques existants : des indicateurs topologiques pourraient être envisagés pour mettre en évidence des incohérences spatiales, comme un bâtiment fictif qui serait cartographié au milieu de l'océan ou dans une zone naturelle (Touya et Brando-Escobar, 2013).

Le carto-vandalisme industriel, idéologique, et le carto-spam se produisent souvent à travers la modification de tags. Des métriques syntaxiques et textuelles pourraient permettre de mettre en évidence le contenu offensif ou insolite présent dans les tags. Par exemple, l'utilisation abusive de caractères spéciaux dans les tags d'une contribution OSM est une métrique possible. Elle permettrait de mettre en évidence le cas de vandalisme de la zone commerciale transformée en lac nommé par le tag `name=:)`². De plus, les indicateurs sur les contributeurs peuvent permettre de détecter les cas de carto-vandalisme discret, c'est-à-dire lorsque ceux qui ne sont pas identifiable à partir d'aucune métrique considérée jusqu'ici. Les indicateurs sur les contributeurs ont justement permis de détecter le vandalisme dans les bases de données ouvertes (Adler *et al.*, 2011).

4.3 Construction d'un premier corpus

Une analyse qualitative d'incidents sur OSM menant au bannissement de contributeurs nous a permis de récupérer moins de 30 cas de carto-vandalisme. Or, il y avait 100 fois plus d'exemples de vandalisme dans Webis-WVC-07 (voir Table III.2), le premier corpus de vandalisme pour Wikipédia. Il nous fallait donc récolter plus d'exemples de carto-vandalisme pour constituer un premier corpus. Par ailleurs, plus le corpus contient d'exemples, plus les analyses à partir de celui-ci seront significatives et plus les modèles entraînés sur ce corpus seront performants dans la détection du carto-vandalisme dans OSM. Comme nous n'étions pas assurés de trouver suffisamment d'exemples réels de carto-vandalisme sur OSM, nous avons choisi de constituer un corpus à partir d'exemples artificiels.

2. Cas réel de vandalisme dans OSM provoqué par un contributeur banni, que nous avons identifié au Chapitre 1 (Figure I.9)

En nous inspirant des types de vandalisme relevés par Ballatore (2014) et de ceux que nous avons identifiés sur les contributeurs bannis, nous avons inséré des contributions synthétiques de vandalisme au sein d'une base de données OSM. En limitant le carto-vandalisme uniquement sur les objets cartographiques de type bâti, nous avons uniquement créé et modifié des bâtiments dans la base de données. Pour ne pas perturber le bon fonctionnement du projet, les contributions synthétiques de vandalisme sont effectuées sur un extrait de la base OSM importé localement dans une base de données PostGIS. Nous avons choisi d'extraire les données OSM visibles au 13 février 2018 sur 4 *snapshots* situés en France (Aubervilliers et Lannilis) et en Allemagne (Stuhr et Heilsbronn). Nous avons ensuite inséré dans chacun des 4 *snapshots* des contributions synthétiques représentant du vandalisme cartographique. Les données, artificielles ou non, sont stockées localement dans une base de données différente selon chaque zone d'étude.

Comme la plupart des contributeurs OSM éditent les données à partir d'une imagerie aérienne de type Bing Maps, nous avons tenté de reproduire ce mode de contribution pour produire le vandalisme. Ainsi, nous avons utilisé le logiciel QGIS pour afficher une image aérienne Bing et la couche vectorielle des bâtiments réellement cartographiés dans OSM sur ces 4 *snapshots*. Puis, nous avons inséré du vandalisme ludique en ajoutant de nouveaux bâtiments de taille immense sur des bâtiments existants. Le carto-vandalisme fantaisiste a consisté à cartographier des bâtiments imaginaires sur des zones dépourvues de bâtiments OSM. L'insertion de carto-vandalisme artistique a consisté à cartographier des objets de forme étrange. La Figure III.3 illustre des bâtiments artificiels insérés pour « carto-vandaliser » OSM sur la ville d'Aubervilliers. Enfin, nous avons également modifié certains bâtiments commerciaux, industriels et religieux en remplaçant la valeur du tag `name` par des absurdités telles que « ... » ou « :) ». La Table III.3 récapitule le carto-vandalisme synthétique qui a été produit dans les quatre villes en quatre catégories auxquelles nous nous référerons dans la suite. Ces contributions synthétiques de carto-vandalisme ont été réalisées avec l'aide d'Étienne le Bihan et Adriano Marzec, deux étudiants ingénieurs en géomatique de l'École Nationale des Sciences Géographiques que nous avons encadrés durant un projet d'initiation à la recherche.

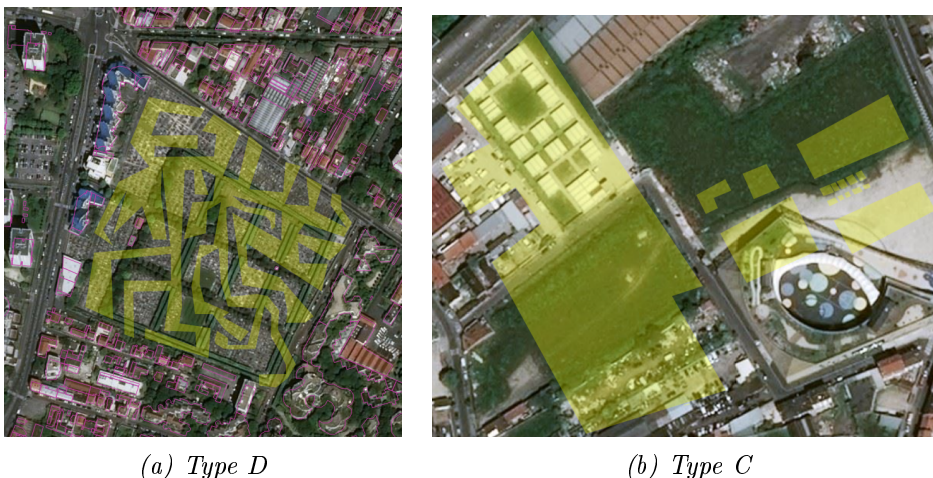


FIGURE III.3. Carto-vandalisme synthétique à Aubervilliers

Dans chaque zone, le carto-vandalisme a été inséré de manière à ce que le carto-

vandalisme représente moins de 1% du nombre total de bâtiments OSM (Table III.4). Ce faible pourcentage de données de vandalisme a été choisi intentionnellement pour reproduire le caractère exceptionnel du carto-vandalisme sur l'ensemble des données contribuées. Toutefois, ce choix a été fait de manière totalement arbitraire, puisque nous ne savons pas estimer le pourcentage de données vandalisées à l'échelle d'une ville. Ce pourcentage pourrait être éventuellement déterminé en relevant des cas de vandalisme connus dans certaines villes. En particulier, une étude récente sur les cas de vandalisme provoqués par les joueurs de Pokémon Go pourrait permettre d'obtenir une estimation de ce type (Juhász *et al.*, 2019).

Type	Description
A	Bâtiments existants, souvent des magasins ou des lieux de culte, dont les tags ont été modifiés. La modification attributaire peut contenir des smileys
B	Bâtiments fictifs ajoutés dans des zones inhabituelles (étang, cours d'eau)
C	Bâtiments fictifs ajoutés dans des espaces vides
D	Bâtiments de forme et/ou de taille inhabituelle

TABLE III.3. Décompte du carto-vandalisme synthétique dans les zones d'étude

Zone	Nombre de bâtiments	Carto-vandalisme synthétique
Aubervilliers	10250	71
Stuhr	6274	44
Lannilis	8066	32
Heilsbronn	6254	31

TABLE III.4. Décompte du carto-vandalisme synthétique dans les zones d'étude

5 Le carto-vandalisme comme anomalie à détecter

Dans le domaine de la fouille de données, la détection d'anomalies consiste à identifier des motifs non-conformes à un comportement attendu au sein d'un jeu de données (Chandola *et al.*, 2009). Le carto-vandalisme est un phénomène qui rend certaines contributions non-conformes à ce qui est attendu au sein d'une base de données géographiques collaborative. Par ailleurs, le caractère exceptionnel de ce phénomène entraîne que ces contributions sont en faible nombre par rapport au reste des données contribuées. Notre hypothèse consiste à modéliser les données de carto-vandalisme comme des anomalies à détecter. Ces anomalies peuvent être alors mise en évidence sous réserve de se placer dans le bon espace de description, c'est-à-dire dans un espace où le carto-vandalisme apparaîtra comme aberrant³.

3. Ici, le terme « aberrant » est employé au sens statistique qui désigne une anomalie (*outlier* en anglais) (Crettaz de Roten et Helbling, 1996).

5.1 Mise en place de descripteurs

Les descripteurs relevés dans l'état de l'art (voir Section 2.1) pour détecter le vandalisme dans les bases de données ouvertes incitent à considérer des métriques qui rendent compte du contenu et du contexte des contributions cartographiques collaboratives. Dans l'espoir de capturer le carto-vandalisme sous diverses formes, nous proposons un certain nombre de métriques qui permet de décrire différents aspects des données géographiques, et en particulier ceux sous lesquels les données de vandalisme prendront des valeurs aberrantes qui les rendront identifiables.

a) Descripteurs de contenu

(i) **Indicateurs géométriques** : Les indicateurs portant sur la géométrie des contributions OSM ont pour but de repérer les actes de carto-vandalisme qui se manifestent au niveau de la forme et la taille des géométries de ces objets. Le carto-vandalisme désigné par le type D dans la Table III.3 pourra typiquement être révélé à travers des indicateurs de taille et de forme. La taille d'un bâtiment cartographié peut se mesurer par son aire et son périmètre. Ainsi, les contributions vandalisées de la Figure III.4 pourront être mises en évidence à travers ces indicateurs de taille.

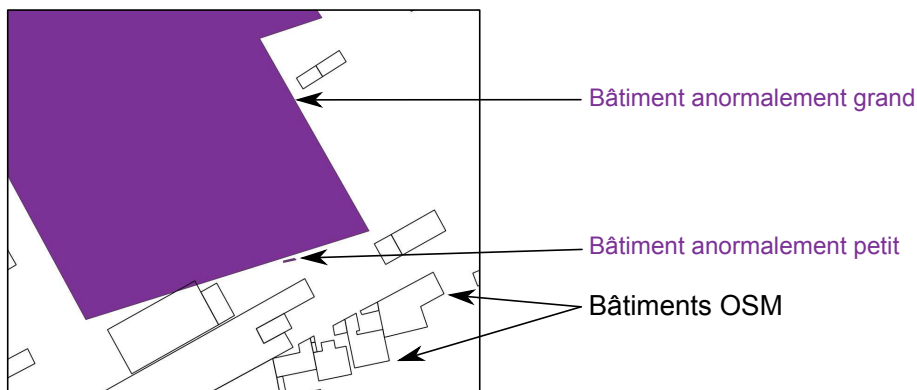


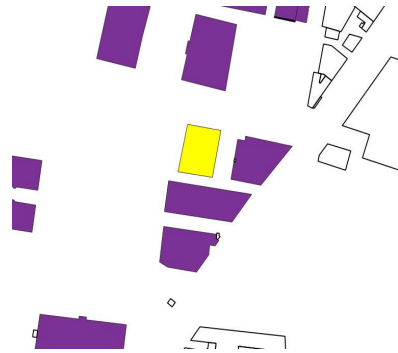
FIGURE III.4. Carto-vandalisme (bâti violet) détectable par des indicateurs de taille

La forme du bâtiment peut être évaluée à travers sa granularité, donnée par la longueur du plus petit côté du bâtiment (Girres et Touya, 2010) ainsi que la valeur médiane des longueurs des côtés du bâtiment. Cette dernière métrique apporte en plus une idée de la régularité de la forme du polygone. Le bâtiment illustré dans la Figure III.5 présente une forme géométrique régulière : la médiane des longueurs de ses côtés (*longueur_médiane*) est relativement proche de la longueur de son plus petit côté (*longueur_min*), alors que le bâtiment de la Figure III.6 possède une forme irrégulière qui se manifeste par un déséquilibre entre la valeur médiane des longueurs des côtés et celle du plus petit côté.

La forme d'un bâtiment peut également être évaluée par sa compacité, calculée par l'indice de Miller :

$$compacite = \frac{4 * \pi * aire(bati)}{perimetre(bati)^2} \quad (III.3)$$

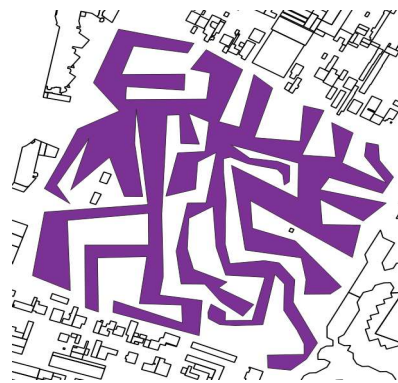
où : *bati* désigne la géométrie de l'objet bâtiment ; *compacite* = 1 lorsque la géométrie est un cercle. Le bâtiment de la Figure III.6 a une géométrie peu compacte, contrairement au bâtiment de la Figure III.5. L'indicateur de compacité peut donc permettre de mettre en évidence des objets de forme inhabituelle.



longueur_min = 30 m
longueur_mediane = 39 m

compacite = 0.77
elongation = 0.76
convexite = 1

FIGURE III.5. Carto-vandalisme de type A (les indicateurs portent sur le bâti jaune)



longueur_min = 5 m
longueur_mediane = 24 m

compacite = 0.02
elongation = 0.99
convexite = 0.51

FIGURE III.6. Carto-vandalisme (bâti violet) détectable par des indicateurs de forme

L'élongation est un indicateur de forme géométrique dont le calcul suit la formule suivante :

$$elongation = \frac{largeur(PPRE_{bati})}{longueur(PPRE_{bati})} \quad (III.4)$$

où $PPRE_{bati}$ correspond au plus petit rectangle englobant la géométrie du bâti. $elongation = 1$ dans le cas où la géométrie est un carré, et tend vers 0 lorsque la géométrie est allongée. Pour le bâti anormalement petit illustré à la Figure III.4, $elongation = 0.14$: cela indique que ce petit bâti présente une forme allongée. En revanche, pour le bâti de la Figure III.5, l'indicateur $elongation$ a une valeur proche de 1.

La convexité d'un polygone est calculée par la formule suivante :

$$convexite = \frac{aire(bati)}{aire(PPRE_{bati})} \quad (III.5)$$

L'indicateur de convexité permet de mettre en évidence les bâtiments de forme étrange dont l'aire est très différente de celle de leur plus petit rectangle englobant. Par exemple, le bâtiment de la Figure III.6 est très peu convexe, alors qu'un bâti de géométrie plus ordinaire a une convexité plus proche de 1 (voir Figure III.5).

(ii) **Indicateurs syntaxiques** : Pour caractériser les dégradations attributaires des données OSM, on considère un indicateur qui mesure le taux maximal de caractères spéciaux dans les tags d’une contribution OSM (Heindorf *et al.*, 2016). Bien que cette métrique ne suffise pas pour identifier toutes les formes de carto-vandalisme qui peuvent être perpétrées sur les tags OSM, elle peut au moins permettre de repérer certaines situations spécifiques de carto-vandalisme, telles que le carto-vandalisme de type A (*cf.* Table III.3). Nous définissons la métrique de caractères spéciaux par :

$$max_special_char_ratio = \max_{v \in V} \left(\frac{n_special_char}{length(v)} \right) \quad (\text{III.6})$$

où V est l’ensemble des valeurs de tags d’une contribution OSM, $n_special_char$ est le nombre de caractères spéciaux dans une valeur de tag v , et la fonction $length$ retourne la taille de la valeur du tag v .

b) Descripteurs de contexte

(i) **Métadonnées** : Les métadonnées peuvent donner des indications sur le contexte d’une contribution OSM. Par exemple, le nombre de tags n_tags d’une contribution donne une indication sur la richesse de détails renseignés à propos de l’objet cartographique. Nous pouvons supposer qu’un objet cartographique décrit par de nombreux tags a plus de chances d’être de bonne qualité. Le nombre d’éditions de la session à laquelle est rattachée une contribution donne une information supplémentaire sur le contexte de saisie : si la session d’édition contient un grand nombre de données, cela peut signifier que le contributeur est prolifique (potentiellement compétent ou fiable) ou bien que les données sont chargées à partir d’un script automatique.

Les métadonnées peuvent également indiquer le type de logiciel utilisé par le contributeur pendant sa saisie. Par exemple, l’utilisation de l’outil d’édition JOSM nécessite une certaine maîtrise technique. Une donnée saisie par cet éditeur peut provenir d’un contributeur compétent (Yang *et al.*, 2016). Par ailleurs, la source de laquelle provient une contribution est une métadonnée intéressante pour qualifier l’information, car elle permet d’indiquer si les informations saisies proviennent d’une autre base de données ouverte ou si celles-ci proviennent de la connaissance personnelle d’un contributeur local sur son environnement. Or, cette connaissance locale est précieuse pour enrichir la base de données ouverte par des informations actuelles et de qualité (Goodchild, 2007).

(ii) **Contexte spatial** : Pour caractériser les relations spatiales d’un bâtiment OSM avec les éléments cartographiques environnants, nous considérons des indicateurs topologiques. Soit $bati$ la géométrie d’un bâtiment OSM, et $OSM_{natural}$ l’ensemble des éléments OSM qui correspondent à des objets géographiques naturels sur lesquels il n’est pas censé y avoir de bâtiment (eau, prairie, forêt, *etc.*). On définit deux métriques topologiques :

$$n_is_within_lulc = |\{nat \in OSM_{natural} / bati \subseteq nat\}| \quad (\text{III.7})$$

$$n_inter_lulc = |\{nat \in OSM_{natural}/bati \cap nat \neq \emptyset\}| \quad (III.8)$$

$n_is_within_lulc$ compte le nombre de zones naturelles dans lesquelles est contenu le bâtiment, alors que n_inter_lulc compte le nombre de zones naturelles dont la géométrie intersecte celle du bâtiment. Ces deux métriques topologiques sont spécifiques aux relations spatiales entre un bâtiment et des zones naturelles. Pour décrire les relations spatiales avec d'autres types d'objets (tels que les routes, par exemple), il faut considérer d'autres métriques.

Lorsque des données géographiques de référence sont disponibles, il peut être intéressant de considérer un indicateur d'appariement. L'appariement consiste à trouver une correspondance entre un objet géographique d'une base de données – en l'occurrence ici de la base OSM – avec le même objet dans une base de données de référence (Walter et Fritsch, 1999). Dans le cas de la ville d'Aubervilliers, comme nous disposons des données du thème bâti provenant de la BD TOPO® de l'IGN, nous avons mis en place un indicateur d'appariement qui quantifie l'éventualité qu'un bâtiment de la BD TOPO® soit homologue au bâti OSM, garantissant alors l'existence de l'objet géographique dans le monde réel. Nous définissons le score d'appariement $min_dist_surf_bati_bdtopo$ par :

$$min_dist_surf_bati_bdtopo = min_{bati_{BDTOPO}}(d_S(bati_{OSM}, bati_{BDTOPO})) \quad (III.9)$$

où d_S est la distance surfacique :

$$d_S(A, B) = 1 - \frac{aire(A \cap B)}{aire(A \cup B)} \quad (III.10)$$

Ce score renvoie la distance surfacique minimale entre un bâtiment OSM donné et tous les bâtiments de la BD TOPO® qui l'intersectent géométriquement. Ce score ne permet pas de trouver un homologue au bâtiment OSM, mais d'indiquer la probabilité que cet homologue existe.

Si $min_dist_surf_bati_bdtopo = 1$, alors le bâtiment OSM n'a pas d'homologue dans la base de référence. Si $min_dist_surf_bati_bdtopo = 0$ alors le bâtiment OSM a un homologue dans la base de référence, avec lequel il est parfaitement apparié spatialement. Le score d'appariement permet de mettre en évidence les bâtiments cartographiés sur OSM qui ne se trouvent pas dans la BD TOPO® : soit parce que, sur cette zone, la base OSM est plus à jour que la BD TOPO®, ou parce qu'une erreur de saisie voire un acte de carto-vandalisme s'est produit dans OSM.

(iii) Contexte temporel ou historique : Pour tenir compte de l'environnement temporel de la contribution, on observe la durée entre la version actuelle et la version précédente de l'objet cartographique, que l'on notera *timespan_to_previous*. Le nombre de contributeurs différents n_users peut également permettre de mettre en évidence des guerres d'éditions ou un contributeur qui s'auto-corrige plusieurs fois pour améliorer ses propres contributions.

(iv) Indicateurs sur les contributeurs : Nous considérons les scores de fiabilité *avg* et *weighted_avg* qui ont été précédemment développés dans le Chapitre

2, car ces deux métriques ont permis de qualifier plus précisément des contributeurs non-fiables. Comme ces scores ne sont pas les seuls indicateurs de contributeurs possibles, nous évaluons également le contributeur sur sa participation spatiale et temporelle. C'est pourquoi nous tenons compte du nombre de ses contributions $n_contributions$ ainsi que le nombre de semestres $n_semesters$ durant lesquels il a participé. Nous ne connaissons pas *a priori* les indicateurs à considérer pour qualifier les données collaboratives, par conséquent, nous choisissons de considérer toutes les métriques possibles, afin de filtrer, par la suite, celles qui sont les plus pertinentes pour notre problématique de détection du carto-vandalisme.

5.2 Sélection des descripteurs optimaux

Nous considérons deux jeux de données du corpus de vandalisme, portant sur les villes d'Aubervilliers et Stuhr. Rappelons que ces jeux de données contiennent des contributions OSM et des contributions synthétiques de carto-vandalisme. Il s'agit dans un premier temps de déterminer les descripteurs qui permettent de détecter le carto-vandalisme dans ces deux jeux de données. En effet, parmi tous les descripteurs proposés dans la partie précédente, nous en avons calculé 37 et 38 respectivement pour les données de Stuhr et Aubervilliers. Les données d'Aubervilliers sont décrites par un descripteur d'appariement (voir équation III.9) avec les données de bâti de la BD TOPO® de l'IGN. Comme nous ne disposons pas de données de référence en Allemagne, ce descripteur n'a pas été implémenté pour les données de Stuhr.

a) Détermination des descripteurs discriminants du carto-vandalisme

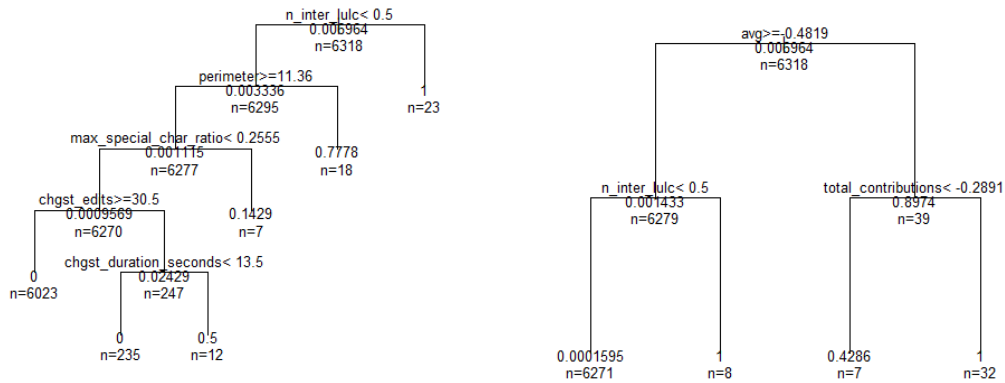
L'implémentation des descripteurs sur les deux jeux de données a été réalisée sans savoir *a priori* la pertinence de chaque métrique pour identifier le carto-vandalisme dans les données OSM. Pour comprendre ce qui permet de discriminer les contributions qui relèvent du carto-vandalisme de celles qui n'en sont pas, nous commençons par construire des arbres de décision sur ces deux jeux de données (Safavian et Landgrebe, 1990). En effet, les arbres de décision permettent de déterminer, parmi les 37 (respectivement 38) descripteurs, un sous-ensemble optimal de descripteurs qui permettent d'expliquer la structuration d'un jeu de données dans les classes **VRAI** (donnée carto-vandalisée) ou **FAUX** (donnée non carto-vandalisée). La construction des arbres de décision a été réalisée sous R, un utilisant le package `rpart`, qui contient une implémentation de l'approche CART (Breiman *et al.*, 1984).

Lecture des arbres de décision : Les Figures III.7 et III.8 illustrent des arbres de décision. Dans un arbre, chaque nœud correspond à une condition portant sur une variable descriptive du jeu de données en entrée de l'arbre. Il sépare le jeu de données en deux branches : la branche de gauche correspond aux données qui respectent la condition, et la branche de droite correspond aux données qui ne la respectent pas. Les feuilles sont décrites par deux valeurs : la valeur n indique le nombre d'individus contenus dans la feuille, et la seconde valeur indique la proportion d'individus qui appartiennent bien à la classe carto-vandalisme ($carto_vandalisme = \mathbf{VRAI}$).

Dans le jeu de données de Stuhr, nous observons sur la Figure III.7 que l'arbre de décision se construit différemment en fonction de la prise en compte ou non des descripteurs de contributeurs *avg*, *weighted_avg* et *total_contributions* (ce dernier donne le nombre total d'éditions produites par le contributeur sur la zone d'étude). En ne tenant pas compte de ces métriques contributeurs, l'arbre de décision de la Figure III.7a montre que les descripteurs discriminants du carto-vandalisme sont :

- *n_inter_lulc*, qui compte le nombre d'intersections de l'objet bâti avec des zones naturelles ;
- *perimeter*, qui indique le périmètre du bâti ;
- *max_special_char_ratio*, qui quantifie le taux de caractères spéciaux dans les tags de l'objet bâti ;
- *chgst_edits*, qui indique le nombre d'éditions produites dans la session durant laquelle l'objet bâti a été édité ;
- *chgst_duration_seconds*, qui indique la durée de la session d'édition.

Vraisemblablement, le descripteur topologique *n_inter_lulc* permet de discriminer le carto-vandalisme de type B ; le descripteur géométrique *perimeter* explique le carto-vandalisme de type D et le descripteur syntaxique *max_special_char_ratio* explique le carto-vandalisme de type A. Par conséquent, le carto-vandalisme de type C est expliqué par les descripteurs de métadonnées *chgst_edits* et *chgst_duration_seconds*.



(a) Sans aucune métrique contributeurs (b) Avec toutes les métriques contributeurs

FIGURE III.7. Arbres de décision construits sur le jeu de données de Stuhr

Dans la Figure III.7b, nous pouvons observer que la prise en compte de tous les descripteurs de contributeurs retourne un arbre de décision dans lequel le carto-vandalisme s'explique principalement par trois descripteurs :

- le score moyen de fiabilité *avg* du contributeur ;
- le nombre total de contributions *total_contributions* du contributeur ;
- *n_inter_lulc*.

Cet arbre de décision montre que les descripteurs de contributeurs permettent d'expliquer en grande partie le carto-vandalisme dans ce jeu de données. Ce résultat

renforce donc l'idée de tenir compte de métriques contributeurs pour détecter le carto-vandalisme, car celui-ci dépend du niveau de confiance de ces derniers.

Cependant, pour les données d'Aubervilliers, la prise en compte ou non des métriques contributeurs n'a pas influencé la construction de l'arbre de décision résultant. Sur la Figure III.8, nous pouvons observer que le carto-vandalisme produit sur cette zone s'explique par les descripteurs :

- *max_special_char_ratio*, qui quantifie le taux de caractères spéciaux dans les tags de l'objet bâti ;
- *n_semesters_auberv*, qui compte le nombre de semestres de participation du contributeur sur la zone ;
- *n_users*, qui compte le nombre de contributeurs uniques ayant édité l'objet bâti (dans son historique) ;
- *lifespan_seconds*, qui indique la durée de vie de la version courante de l'objet bâti ;
- *v_contrib*, qui indique le numéro de version courant de l'objet bâti.

D'après cet arbre de décision, le carto-vandalisme sur la zone d'Aubervilliers s'explique donc par des descripteurs topologiques (*max_special_char_ratio*), contributeurs (*n_semesters_auberv*) et historiques. Selon la zone étudiée, nous remarquons que le carto-vandalisme ne s'explique pas par les mêmes descripteurs discriminants. Cela est certainement dû au fait que les propriétés géographiques et les caractéristiques de contribution sur OSM sont différentes à Stuhr et à Aubervilliers.

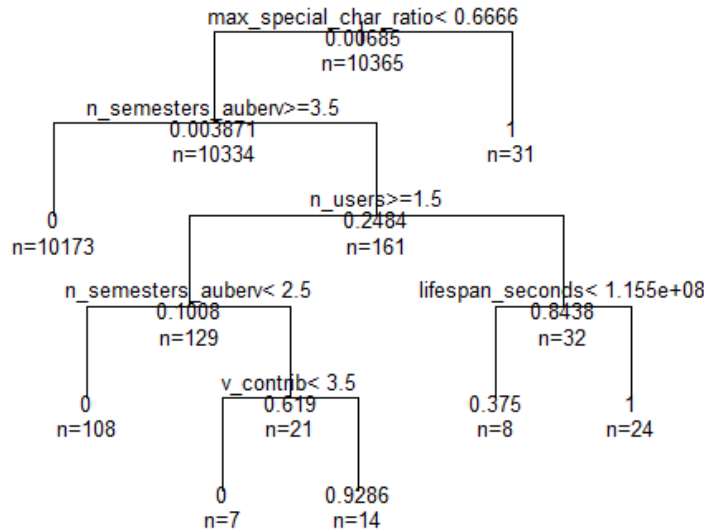


FIGURE III.8. Arbre de décision construit sur le jeu de données d'Aubervilliers

b) détection d'anomalies avec différentes combinaisons de descripteurs

L'analyse de ces arbres de décision a mis en évidence un sous-ensemble de descripteurs qui permettent d'expliquer le carto-vandalisme dans les données de Stuhr et Aubervilliers. Pour expérimenter la détection du vandalisme par une méthode non-supervisée de détection d'anomalies, nous commençons par considérer le sous-ensemble de descripteurs mis en évidence dans les arbres de décision présentés dans

la partie précédente. En effet, l'utilisation d'une méthode de détection d'anomalies impose de considérer un nombre limité de descripteurs pertinents, car l'introduction d'un trop grand nombre de descripteurs peut affecter la performance des algorithmes de classification. Ce phénomène, appelé fléau de la dimension (*curse of dimensionality*), constitue une limite des méthodes d'apprentissage à classifier des données lorsque celles-ci sont décrites dans un espace à trop haute dimension, une dimension étant modélisée par un descripteur (Verleysen et François, 2005).

Nous cherchons à détecter le carto-vandalisme avec une méthode de détection d'anomalies en commençant par un ensemble initial de descripteurs formé par ceux qui sont issus des arbres de décisions. Puis, en faisant varier la combinaison des descripteurs, il s'agit d'étudier la performance du système de détection d'anomalies à détecter le carto-vandalisme. Pour cela, nous utilisons DBSCAN, un algorithme non-supervisé de regroupement de données qui se base sur la notion de densité (Ester *et al.*, 1996). Cet algorithme va regrouper en classes les données qui seront proches en termes de distance dans l'espace des descripteurs. L'algorithme prend en entrée :

- un ensemble des descripteurs portant sur les bâtiments ;
- *eps*, la distance maximale entre deux points d'une même classe ;
- *minPts*, le nombre minimal d'objets requis pour pouvoir former une classe.

En sortie de DBSCAN, un individu sera classé dans une classe s'il existe d'autres données qui lui sont proches dans l'espace des descripteurs. Dans le cas contraire, c'est-à-dire s'il n'existe pas suffisamment de données proches (*i.e.* à une distance inférieure de *eps*) de celle-ci, elle sera considérée comme une anomalie par DBSCAN. Dans notre situation, nous espérons que les contributions de carto-vandalisme soient classifiées comme des anomalies par DBSCAN. Cela suppose donc de se placer dans un espace de description des données adéquat, et de paramétrer l'algorithme de manière à ce que les cas de carto-vandalisme ne soient regroupés dans aucune classe.

Afin de trouver les valeurs optimales de paramètres de DBSCAN – c'est-à-dire les valeurs de *eps* et *minPts* qui permettent de maximiser le rappel et la précision dans la détection d'anomalies – nous lançons DBSCAN de manière itérative sur les données d'Aubervilliers (sous R, avec le `package dbSCAN`) en faisant varier les paramètres. Cette expérience préliminaire a permis de fixer les paramètres (*eps*, *minPts*) = (0.75, 3) pour la détection des anomalies avec DBSCAN dans notre cas d'étude.

En faisant varier la combinaison des descripteurs, nous lançons la détection d'anomalies sur les données d'Aubervilliers. Les meilleurs résultats de performance de DBSCAN sont données dans la Table III.5, et ont été obtenus en considérant les descripteurs suivants :

- l'indicateur d'appariement *min_dist_surf_bati_bdtopo*,
- l'indicateur de taille *perimeter*,
- l'indicateur syntaxique *max_special_char_ratio*,
- l'indicateur topologique *n_is_within_lulc*,
- l'indicateur de contributeur *avg*.

D'après la Table III.5, la précision de 20% obtenue est très faible : elle s'explique par le fait que l'algorithme a classé 67 vrais positifs et 266 faux positifs. En revanche,

TABLE III.5. Détection du carto-vandalisme avec DBSCAN sur Aubervilliers

Précision	Rappel	Erreur	Type A	Type B	Type C	Type D
0.201	0.944	0.026	1	0.67	1	0.67

le fort rappel traduit le fait que le système est capable de détecter en grande partie le vandalisme synthétique (environ 94%). L'erreur (presque 3%) signifie que globalement, l'algorithme se trompe très peu dans la distinction entre le vandalisme (*i.e.* l'anomalie) et le non-vandalisme. Pour confirmer l'intérêt des descripteurs d'appariement et de fiabilité, nous étudions par la suite l'influence des résultats avec et sans les indicateurs respectifs.

c) Pertinence de l'indicateur de fiabilité du contributeur

La Table III.6 indique les résultats de la détection d'anomalies en fonction de la prise en compte de différents indicateurs de contributeur. De manière générale, la prise en compte d'un indicateur de contributeur entraîne une diminution de la précision et une augmentation de l'erreur : cela signifie que le nombre de faux positifs augmente. En revanche, on observe également une augmentation du rappel, ce qui signifie que le système parvient à classer correctement un plus grand nombre de contributions vandalisées (*i.e.* celles-ci sont classées comme des anomalies). Quel que soit le descripteur de fiabilité, notons que le vandalisme de type C est mieux détecté qu'en l'absence totale d'indicateur de fiabilité, puisque sans cet indicateur, le carto-vandalisme de type C n'est pas du tout détecté. Par conséquent, ces résultats montrent l'intérêt de prendre en compte la fiabilité du contributeur dans la détection de vandalisme, en particulier pour détecter les cas les moins évidents, tels que le carto-vandalisme de type C.

TABLE III.6. Performance de DBSCAN avec différents indicateurs de contributeur en entrée (les valeurs optimales de chaque colonne sont en gras).

Indicateur contributeur	Précision	Rappel	Erreur	Type C détecté
\emptyset	0.231	0.563	0.015	0
<i>total_contributions</i>	0.146	0.592	0.026	0.67
<i>n_semesters</i>	0.086	0.958	0.070	0.73
<i>weighted_avg</i>	0.201	0.944	0.026	1
<i>avg</i>	0.201	0.944	0.026	1

Parmi les trois indicateurs de contributeur introduits, on constate qu'avec les descripteurs *avg* et *weighted_avg*, le système détecte le carto-vandalisme de type C en intégralité. En comparaison, avec *total_contributions* et *n_semesters*, la prise en compte de *avg* ou *weighted_avg* donne une précision maximale et une erreur minimale. Quant au rappel de 94%, celui-ci est moins élevé qu'avec la prise en compte de *n_semesters* (96%), mais il reste meilleur en comparaison avec le résultat d'une détection qui ne tient compte d'aucun indicateur de contributeur, à 56% de rappel.

Les indicateurs *avg* et *weighted_avg* semblent donc les plus adaptés pour dé-

tecter le carto-vandalisme fantaisiste créé dans des espaces vides. Pour déterminer lequel des indicateurs *avg* et *weighted_avg* est le plus pertinent, il conviendrait d'étudier l'influence de la pondération du score de fiabilité sur la détection d'anomalies, ce que nous n'avons pas fait dans le cadre de cette thèse, mais qui peut constituer une perspective de ce travail.

d) Pertinence de l'indicateur d'appariement

La Table III.7 donne les résultats de la détection d'anomalies avec et sans prise en compte de l'indicateur d'appariement en entrée de DBSCAN. Avec l'indicateur d'appariement, la précision diminue de 3%, mais le rappel augmente de 3%. En d'autres termes, le système parvient à mieux détecter le vandalisme mais devient plus tolérant face aux faux positifs, d'où une plus grande erreur avec l'appariement. En observant les différents types de vandalisme détectés, on constate que l'indicateur d'appariement n'a pas d'effet sur la détection du vandalisme de type A, B et D. En revanche, il permet de détecter intégralement le vandalisme de type C, alors que celui-ci n'était détecté qu'à 87% sans appariement. Ces résultats permettent donc de valider la pertinence de ce descripteur pour détecter le vandalisme.

TABLE III.7. Performance de DBSCAN avec et sans prise en compte de l'indicateur d'appariement (les valeurs optimales de chaque ligne sont en gras).

Indicateur d'appariement	Sans	Avec
Précision	0.234	0.201
Rappel	0.916	0.944
Erreur	0.021	0.026
Type A détecté	1	1
Type B détecté	0.67	0.67
Type C détecté	0.87	1
Type D détecté	0.67	0.67

L'indicateur d'appariement développé ici repose sur l'équation III.9. Cette implémentation est très basique car elle ne repose que sur des critères géométriques et topologiques. En envisageant un indicateur d'appariement plus sophistiqué, les résultats seraient probablement meilleurs que ceux obtenus ici. À ce sujet, des travaux récents proposent des méthodes plus poussées pour apparier des données de bâti OSM avec des données de référence (Fan *et al.*, 2014a; Xu *et al.*, 2017).

5.3 Étude de la dépendance des descripteurs optimaux avec la zone d'étude

a) Détection du carto-vandalisme sur Stuhr avec les descripteurs optimaux d'Aubervilliers

Les résultats expérimentaux obtenus ont permis d'établir un ensemble de descripteurs de carto-vandalisme pertinents pour détecter les actes de vandalisme sur les bâtiments OSM d'Aubervilliers. Toutefois, nous pouvons questionner la validité de ces descripteurs optimaux pour détecter le carto-vandalisme sur d'autres zones géographiques. Pour cela, nous avons lancé des expériences de détection d'anomalies avec DBSCAN sur les données de la ville de Stuhr.

Les résultats de la détection sur le bâti OSM de Stuhr avec tous les autres descripteurs optimaux sont donnés dans la Table III.8. Les cas de carto-vandalisme sont détectés avec un plus faible rappel, mais la précision est légèrement meilleure que sur les données d'Aubervilliers. L'absence d'indicateur d'appariement pourrait expliquer pourquoi le carto-vandalisme de type C n'est détecté qu'à 55% sur Stuhr, alors qu'il est entièrement détecté à Aubervilliers.

TABLE III.8. Résultats de détection avec DBSCAN avec les descripteurs optimaux sur Aubervilliers et Stuhr.

Zone	Précision	Rappel	Erreur	Type A	Type B	Type C	Type D
Aubervilliers	0.201	0.944	0.026	1	0.67	1	0.67
Stuhr	0.224	0.75	0.020	0.78	1	0.55	1

b) Détermination des descripteurs optimaux sur Stuhr

Nous lançons plusieurs détections d'anomalie avec DBSCAN en faisant varier la combinaison des descripteurs en entrée de l'algorithme. Les meilleurs résultats sont obtenus pour un ensemble différent de descripteurs de celui qui avait été trouvé pour détecter les anomalies dans les données d'Aubervilliers. La Table III.9 indique les descripteurs optimaux pour détecter les anomalies dans les jeux de données d'Aubervilliers et Stuhr respectivement.

L'indicateur *min_dist_surf_bati_bd_topo* n'a pas été calculé pour Stuhr car cet indicateur quantifie l'appariement avec des données de référence française (Thème « Bâti » de la BD TOPO®). Comme nous ne disposions pas de données de référence sur Stuhr, nous n'avons pas développé un tel indicateur sur ce jeu de données. Cela ne signifie donc pas que l'indicateur d'appariement ne soit pas utile pour ce jeu de données. Au contraire, si cet indicateur avait été calculé pour Stuhr, l'ensemble optimal des descripteurs aurait été probablement différent de celui décrit dans la Table III.9. Il serait intéressant d'étudier dans quelle mesure l'absence d'un descripteur (tel que l'appariement) peut être remplacé par d'autres indicateurs dans la détection du carto-vandalisme.

L'indicateur de fiabilité du contributeur optimal est *weighted_avg* pour Stuhr,

alors que sur les données d'Aubervilliers, la détection est plus performante en prenant en compte l'indicateur *avg*. De plus, pour détecter le carto-vandalisme de Stuhr, il a fallu considérer l'indicateur topologique *n_inter_lulc*, qui compte le nombre d'intersections avec des zones naturelles, en plus de l'indicateur *n_is_within_lulc*, qui compte le nombre de zones naturelles dans lesquelles le bâti est entièrement contenu.

TABLE III.9. Descripteurs optimaux pour la détection du carto-vandalisme dans les données d'Aubervilliers et Stuhr

Descripteur	Aubervilliers	Stuhr
<i>min_dist_surf_bati_bd_topo</i>	X	
<i>perimeter</i>	X	X
<i>max_special_char_ratio</i>	X	X
<i>n_is_within_lulc</i>	X	X
<i>n_inter_lulc</i>		X
<i>avg</i>	X	
<i>weighted_avg</i>		X

c) Intérêt d'un indicateur de fiabilité du contributeur pour la détection du carto-vandalisme

Pour chercher à confirmer l'intérêt d'un indicateur de contributeur pour détecter le carto-vandalisme indépendamment de la zone, on relance DBSCAN sur Stuhr avec les descripteurs suivants :

- l'indicateur de taille *perimeter* ;
- l'indicateur syntaxique *max_special_char_ratio* ;
- l'indicateur topologique *n_inter_lulc* ;
- un indicateur de contributeur.

Les résultats de la détection sont renseignés dans la Table III.10.

Indicateur contributeur	Précision	Rappel
\emptyset	0.24	0.55
<i>total_contributions</i>	0.13	0.61
<i>n_semesters</i>	0.07	0.64
<i>weighted_avg</i>	0.21	0.80
<i>avg</i>	0.22	0.80

TABLE III.10. Influence d'un indicateur de fiabilité sur Stuhr

On observe que la prise en compte d'un indicateur de contributeur, quel qu'il soit, augmente le rappel. Dans cette situation, il serait envisageable que le tri des vrais positifs des faux positifs soit réalisé manuellement par un humain. De plus, l'obtention de faux positifs est moins problématique qu'un grand nombre de faux

négatifs, c'est-à-dire des cas réels de carto-vandalisme qui n'ont pas été détectés comme anomalies par le système.

En revanche, la prise en compte d'un indicateur de contributeur réduit la précision de la détection. Par ailleurs, les indicateurs *avg* et *weighted_avg* parviennent à limiter la détection de faux positifs (la baisse de précision est de 1% à 2% par rapport à la non prise en compte d'aucun indicateur de fiabilité), en comparaison des autres indicateurs de fiabilité (où la baisse de précision est de 11% à 17%). Ces observations corroborent celles qui ont été faites pour Aubervilliers. Par conséquent, nous pouvons affirmer que l'utilisation d'un indicateur de fiabilité du contributeur de type *avg* ou *weighted_avg* permet d'optimiser la détection du carto-vandalisme.

d) Détection du carto-vandalisme sur Aubervilliers avec les descripteurs optimaux de Stuhr

La Table III.11 contient les résultats issus de cette expérience sur Stuhr et Aubervilliers, en utilisant l'ensemble optimal de descripteurs pour Stuhr en entrée de DBSCAN.

TABLE III.11. Résultats de détection avec DBSCAN avec les descripteurs optimaux de Stuhr

Zone	Précision	Rappel	Erreur	Type A	Type B	Type C	Type D
Aubervilliers	0.23	0.915	0.022	1	0.67	0.87	0.67
Stuhr	0.215	0.795	0.022	1	1	0.55	1

La détection est meilleure sur Stuhr que sur Aubervilliers au sens que le carto-vandalisme de type A est intégralement détecté. Cependant, sur les données d'Aubervilliers, le système ne parvient plus à détecter tout le carto-vandalisme de type C. L'ensemble optimal de descripteurs n'est donc pas universel pour détecter le carto-vandalisme, mais il varie selon la zone géographique étudiée. Néanmoins, d'après nos analyses, pour détecter le carto-vandalisme comme anomalie dans une autre zone, il est recommandé de prendre en compte :

- un score de fiabilité du contributeur (*avg* ou *weighted_avg*) ;
- un indicateur d'appariement (lorsque c'est possible) pour détecter le carto-vandalisme fantaisiste ;
- au moins un indicateur topologique pour évaluer les relations spatiales invraisemblables entre l'objet étudié et d'autres objets géographiques (tel que le chevauchement d'un bâtiment avec un lac dans notre étude) ;
- un indicateur syntaxique (type *max_special_char_ratio*) ;
- un indicateur géométrique (type *perimeter*).

Ces recommandations peuvent être considérées comme un ensemble initial de descripteurs qu'il convient de réajuster pour chaque nouvelle zone d'étude. Ce réajustement peut être fait en ajoutant de nouveaux descripteurs ou en revoyant l'implémentation de certaines métriques : par exemple, envisager une nouvelle méthode d'appariement ou un indicateur syntaxique différent.

5.4 Analyse des faux positifs

La détection des anomalies sur les données d'Aubervilliers avec les descripteurs optimaux (voir Table III.5) a permis de détecter le carto-vandalisme avec une faible précision de 20%. Cela signifie donc que les anomalies détectées contiennent 80% de faux positifs, c'est-à-dire de bâtis OSM réels détectés comme des anomalies par l'algorithme DBSCAN. Nous menons une analyse unidimensionnelle sur ces faux positifs pour vérifier si ces données sont aberrantes par rapport au reste du jeu de données d'Aubervilliers selon chaque variable descriptive. Une donnée est aberrante lorsqu'elle se trouve hors de l'intervalle $[Q1 - 1.5 \cdot IQR, Q3 + 1.5 \cdot IQR]$, où l'espace interquartile est défini par $IQR = Q3 - Q1$, $Q1$ et $Q3$ étant respectivement le premier et le troisième quartile. La Figure III.9 schématise la chaîne de traitement utilisée pour filtrer les faux positifs vus comme des anomalies selon chaque variable descriptive.

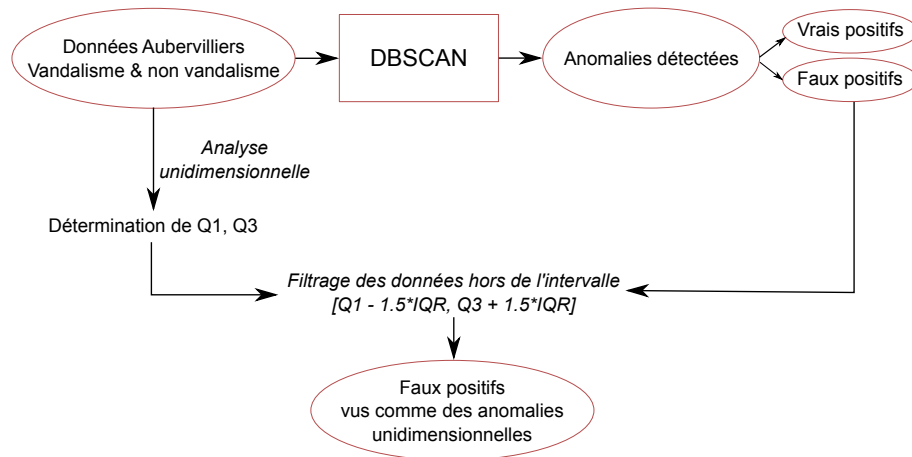


FIGURE III.9. Processus de filtrage des faux positifs sur chaque variable descriptive

a) Anomalies selon le descripteur géométrique de taille

Les faux positifs issus de DBSCAN qui sont également vus comme des données aberrantes selon le descripteur géométrique de taille (*perimeter*) sont visualisés sur une carte OSM (voir Figure III.10). Nous pouvons remarquer que ces faux positifs vus comme des anomalies selon la variable *perimeter* sont des bâtiments de très grande taille, parmi lesquels nous pouvons identifier un centre commercial, une école et une clinique (Figure III.11).

D'après la Figure III.12, la distribution du périmètre des bâtiments montre que les bâtiments les plus grands – dont le périmètre est compris entre 119 et 1250 mètres – sont vus comme des anomalies selon la variable *perimeter*. En effet, le jeu de données est constitué majoritairement de petits bâtiments résidentiels : ces bâtiments de plus petite taille sont observables sur la Figure III.11, alors qu'une école et une clinique sont deux bâtiments plus grands qui ont été détectés comme des anomalies. Par ailleurs, le bâti cartographié sur OSM provient en grande partie de l'import de données cadastrales : le jeu de données contient donc de nombreux morceaux de bâti présentant un faible périmètre. Les données de faible périmètre

étant plus nombreuses que les grands bâtiments, ces derniers apparaissent donc des données aberrantes selon la variable *perimeter*.

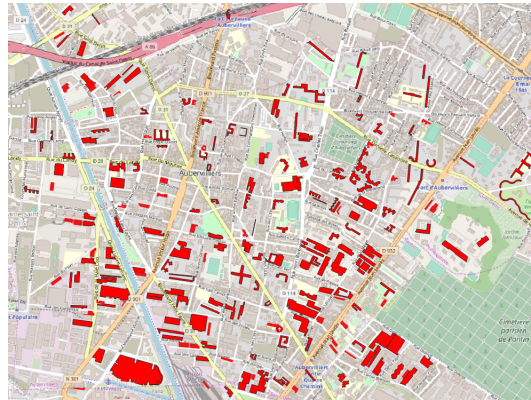


FIGURE III.10. Aperçu des faux positifs vus comme des points selon la variable *perimeter*



FIGURE III.11. Anomalies de périmètre détectées par DBSCAN

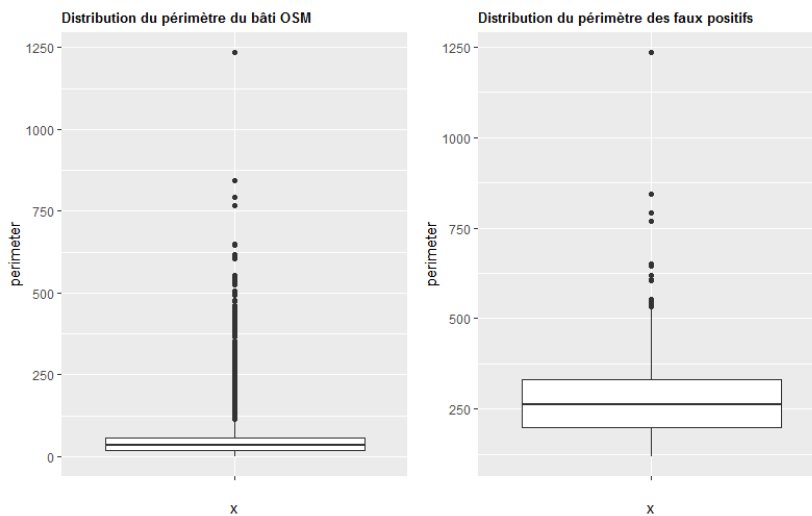


FIGURE III.12. Distribution du périmètre du bâti et des faux positifs

b) Anomalies selon le descripteur d'appariement

Le score d'appariement varie entre 0 et 1 : un score égal à 1 signifie que le bâti n'a pas d'homologue dans la base de référence ; au contraire, un score égal à 0 signifie que le bâtiment a été parfaitement apparié à un bâtiment de la BD TOPO®. D'après la Figure III.13, nous remarquons que les anomalies selon la variable d'appariement dans le jeu de donnée correspondent aux bâtiments qui n'ont pas été appariés (leur score d'appariement est compris entre 0.95 et 1).

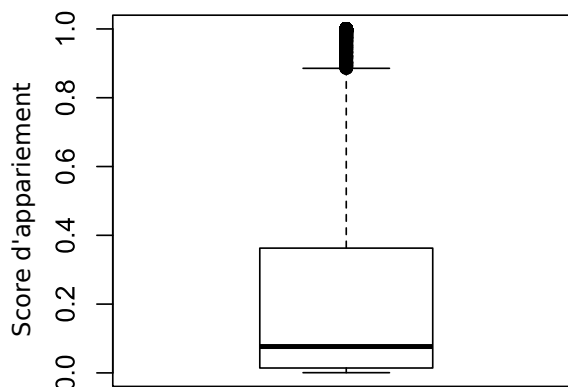


FIGURE III.13. Distribution de l'indicateur d'appariement du bâti OSM et de la BD TOPO® à Aubervilliers

Parmi les faux positifs détectés par DBSCAN qui sont également des données aberrantes selon le descripteur d'appariement, nous avons identifié des erreurs dans les données de bâti OSM. En particulier, la Figure III.14 illustre le cas d'un objet cartographié sur OSM comme étant un garage dont le tag `name` contient la description `dalle libre - toit du parking`. Or, ce bâti n'a aucun homologue dans la BD TOPO® (Figure III.14a), et la visualisation de ce bâti sur une image aérienne montre que la zone sur laquelle il se trouve semble être une cour d'immeuble. Ce bâti OSM n'existe donc pas sur le terrain (Figure III.14b). L'historique de cet objet indique que ce bâti a été importé dans OSM à partir des données cadastrales en 2010. Nous ne savons pas si, à cette époque, cet espace était effectivement construit. Quoi qu'il en soit, la contribution étudiée est une version modifiée du bâtiment dont les attributs `name` et `building` ont été modifiés, transformant ce bâti en toit de parking.

En étudiant le profil du contributeur dont provient cette contribution, nous observons que celui-ci n'a produit que deux éditions sur le projet OSM. Sa deuxième contribution est illustrée à la Figure III.15 et se situe sur la même zone que le bâti transformé en garage. Cet objet correspond en réalité à un projet de potager organisé par une association de la ville. Le contributeur est donc un habitant de cette ville, probablement un membre de cette association, car il contribue une information locale.

Quoi qu'il en soit, le garage détecté comme faux positif par DBSCAN est un objet qui n'existe pas dans la réalité. Toutefois, bien que cet objet soit fictif – et en cela il dégrade l'espace cartographique – nous ne pouvons pas qualifier cette contribution de carto-vandalisme. En effet, nous n'avons pas suffisamment d'éléments qui laisseraient à croire que ce contributeur a dégradé volontairement la carte. Il se peut

que le contributeur soit une personne qui, dans le cadre de son activité associative, aurait chargé des données concernant les aménagements faits pour des événements éphémères dans ce quartier, sans savoir que cela dégraderait la carte. Le cas de ce faux positif est intéressant car, même s'il ne correspond pas à du carto-vandalisme – du moins, nous ne pouvons pas le qualifier comme tel avec certitude – il correspond à une donnée qui mériterait d'être contrôlée par un opérateur humain.

Nous avons également identifié un autre cas de bâti OSM détecté comme anomalie par DBSCAN qui n'a pas été apparié avec aucun bâti de la BD TOPO® (Figure III.16). Ce bâti a été importé du cadastre en 2010 et n'a jamais été édité depuis. En visualisant une image aérienne de cette zone, il s'avère que cette dernière est actuellement une zone de construction. Par conséquent, ce bâti est un faux positif car, n'étant plus valide, il correspond à un objet fictif sur la carte.



(a) Le bâti ne trouve aucun homologue dans la BD TOPO®



(b) Affichage du bâti sur un fond de carte Bing

FIGURE III.14. Cas d'un bâti cartographié sur OSM qui n'existe pas dans la réalité



(a)

Attributs

landuse	farmland
name	104 Barbus
operator	projet de potager urbain d'access libre

(b) Tags associés à l'objet

FIGURE III.15. Deuxième édition du contributeur ayant édité un bâti fictif.



FIGURE III.16. Autre cas de bâti détecté comme faux positif qui n'existe plus sur le terrain (image aérienne Bing).

c) Anomalies selon le descripteur syntaxique

Le descripteur syntaxique a été implémenté pour capturer le carto-vandalisme affectant les tags des données OSM. Cependant, parmi les faux positifs détectés par DBSCAN, nous avons pu identifier des bâtiments apparaissant comme des données aberrantes selon ce descripteur syntaxiques car ceux-ci comportent un tag `website` indiquant un lien vers un site internet. La Figure III.17 illustre deux cas de faux positifs détectés comme anomalies selon le descripteur syntaxique. Ces bâtiments correspondent à un centre commercial et à une salle d'escalade dont le site internet est renseigné dans les tags de ces objets. Ces liens contiennent des caractères spéciaux tels que les points et les barres obliques (« / »), donnant une valeur anormalement élevée à la métrique $max_special_char_ratio$. Le descripteur syntaxique mérite donc d'être amélioré pour éviter de détecter des bâtiments publics pour lesquels un site internet ou des horaires d'ouverture sont susceptibles d'être renseignés dans les tags OSM, car ces informations contiennent souvent un certain nombre de caractères spéciaux.

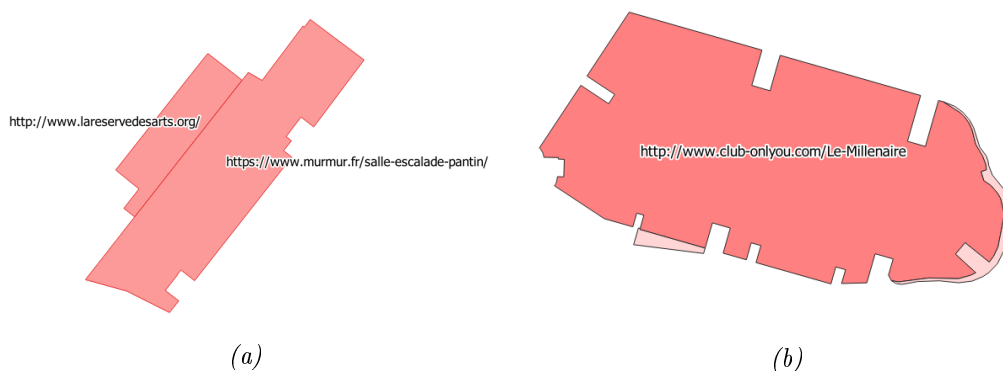


FIGURE III.17. Faux positifs détectés comme anomalies par le descripteur syntaxique

d) Anomalies selon le descripteur topologique

L'ensemble optimal de descripteurs pour le jeu de données d'Aubervilliers contient le descripteur $n_is_within_lulc$ qui compte le nombre de zones naturelles qui

contiennent la géométrie des bâtiments. En étudiant la distribution des valeurs prises par les données d'Aubervilliers sur cette métrique, nous observons que des bâtiments de la Figure III.18 sont détectés comme des anomalies par DBSCAN et apparaissent comme des données aberrantes selon le descripteur $n_is_within_lulc$. En effet, ces bâtiments se situent à l'intérieur d'un bois. Or, Aubervilliers étant une ville fortement urbanisée, il est très inhabituel que des bâtiments se trouvent à l'intérieur d'un bois ou d'une forêt. Par ailleurs, nous remarquons qu'un seul de ces bâtiments est apparié à un bâti de la BD TOPO®, et le feuillage du bois sur l'imagerie aérienne ne permet pas de vérifier que les autres bâtiments existent sur le terrain. Par conséquent, la détection de ces bâtiments comme faux positifs est intéressante car elle peut permettre de vérifier certains objets dans OSM dont la validité n'est pas évidente.

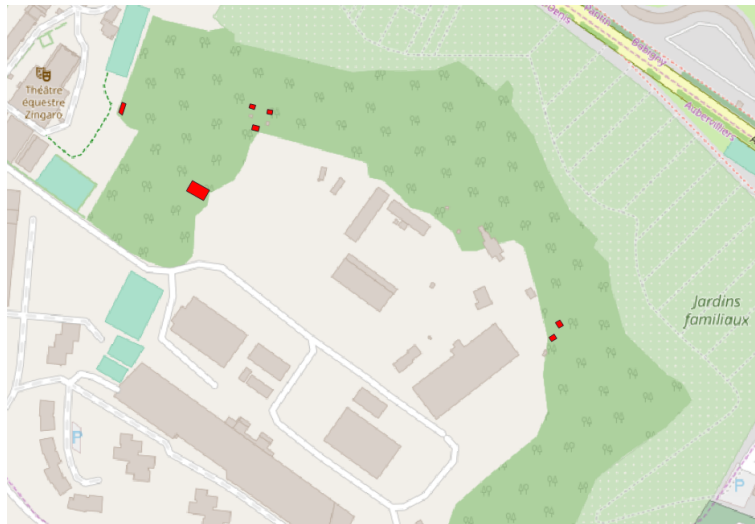


FIGURE III.18. Faux positifs apparaissant comme des anomalies selon le descripteur topologique $n_is_within_lulc$

5.5 Détection d'anomalies sur un jeu de données dépourvu de carto-vandalisme synthétique

Il s'agit ici d'étudier le potentiel de la méthode à détecter du carto-vandalisme réel. Nous lançons une détection d'anomalies avec DBSCAN sur des données bâti OSM sur la ville de Bondy (93) sans y insérer aucune contribution synthétique de carto-vandalisme, dans le but d'analyser les anomalies qui sont détectées. Le paramétrage de l'algorithme DBSCAN est le même que celui utilisé pour détecter les anomalies dans le jeu de données d'Aubervilliers. L'ensemble des descripteurs considéré est quasiment le même que celui d'Aubervilliers, la seule différence étant que nous considérons le descripteur de fiabilité du contributeur pondéré $weighted_avg$ au lieu de avg .

La détection d'anomalies sur les bâtiments OSM de Bondy classe 223 bâtiments comme anomalies, soit environ 2% de l'ensemble des données (Table III.12). Un aperçu des bâtiments détectés comme anomalies sur la Figure III.19 montre que la plupart d'entre eux sont de grands bâtiments. Cela s'explique notamment par la prise en compte du descripteur géométrique de taille $perimeter$.

TABLE III.12. Détection d'anomalies sur les bâtis OSM de Bondy avec DBSCAN

Zone	Nombre d'anomalies détectées	Nombre de bâtiments normaux
Bondy	223	10392

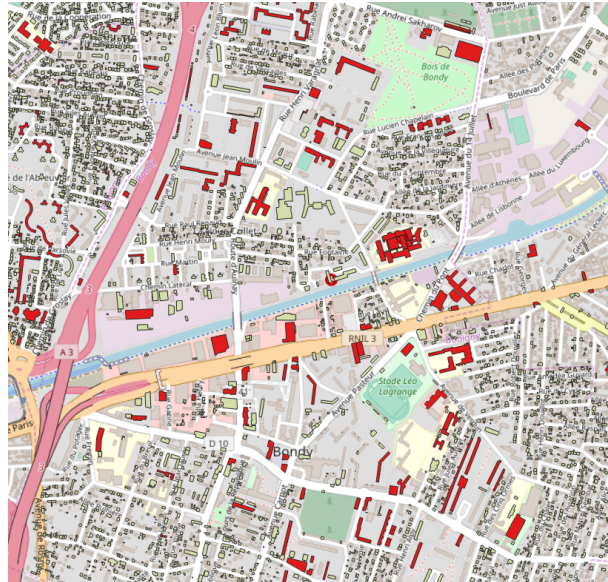
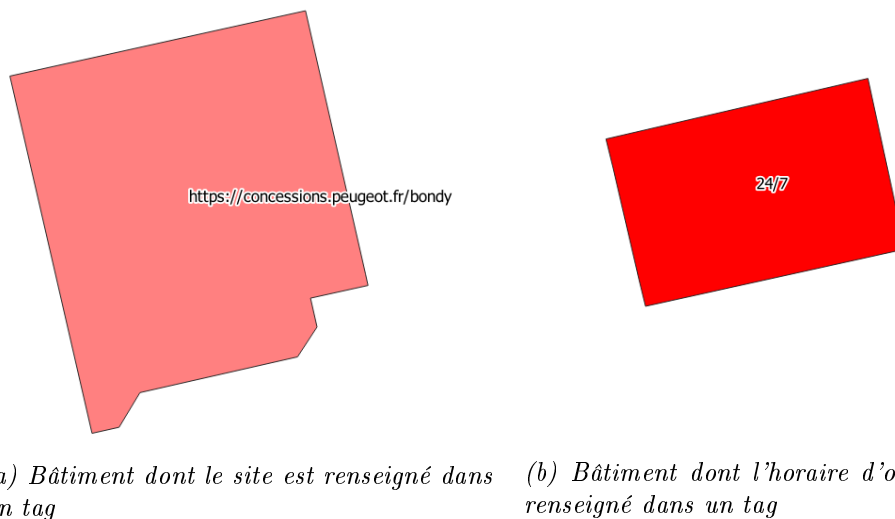


FIGURE III.19. Bâtiments détectés comme anomalies par DBSCAN sur les données de Bondy

De la même manière que pour les faux positifs détectés sur le jeu de données d'Aubervilliers, le descripteur syntaxique *max_special_char_ratio* entraîne que les bâtiments qui possèdent un tag indiquant un site internet ou des horaires d'ouverture ont tendance à être détectés comme anomalies car ces tags comportent des caractères spéciaux (Figure III.20).



(a) Bâtiment dont le site est renseigné dans un tag

(b) Bâtiment dont l'horaire d'ouverture est renseigné dans un tag

FIGURE III.20. Anomalies détectées à Bondy dont les tags contiennent des caractères spéciaux

Certaines anomalies détectées dans les bâtiments de Bondy peuvent s'expliquer par un mauvais appariement avec les données de bâti de la BD TOPO®. Nous avons

identifié le cas d'un bâtiment OSM détecté comme anomalie car il ne trouve aucun homologue dans la BD TOPO®), alors que le bâtiment existe bien sur le terrain (Figure III.21). Cette anomalie révèle donc un manque d'actualité de la base de référence.



FIGURE III.21. Bâtiment réel (en rouge) absent de la BD TOPO® (bâti jaune)

Cependant, nous avons relevé un autre cas où l'anomalie détectée est un bâtiment faiblement apparié avec les données de la BD TOPO®, alors que celles-ci sont correctes par rapport à la vérité terrain. En réalité, l'anomalie est un bâti importé des données du cadastre et qui n'a pas été mis à jour depuis 2011. Par conséquent, l'anomalie détectée permet de mettre en évidence une donnée OSM qui n'est plus actuelle et donc de mauvaise qualité, mais qui n'est pas un cas réel de carto-vandalisme.



(a) Trois bâtiments représentés correctement dans la BD TOPO®



(b) Le bâti OSM (en rouge) n'est pas conforme à la réalité terrain

FIGURE III.22. Cas d'une anomalie correspondant à une donnée obsolète

Notre analyse des anomalies détectées sur les bâtiments OSM de Bondy n'a pas permis d'identifier des cas réels de carto-vandalisme. En considérant ces 223 anomalies comme des faux positifs, l'erreur de classification du système est donc estimée à 2.1%. Cette erreur est proche de celle d'environ 2% qui a été estimée pour détecter le

carto-vandalisme sur les jeux de données d’Aubervilliers et Stuhr (voir Table III.11). Par conséquent, le système de détection semble posséder une erreur de base de 2% pour détecter le carto-vandalisme, ce qui n’est pas négligeable dans un grand jeu de données. Dans un jeu de données comportant 10000 objets, le système risque de détecter 200 contributions qui ne sont pas des cas réels de vandalisme cartographique. Par conséquent, cette méthode ne permet pas de détecter avec précision le carto-vandalisme, mais elle offre la possibilité de filtrer un sous-ensemble de données qui nécessitent d’être vérifiées, par un opérateur humain par exemple.

5.6 Paramétrage du système de détection

Nous avons choisi de détecter les anomalies avec l’algorithme non-supervisé DBSCAN, car celui-ci permet de regrouper des données qui partagent des caractéristiques proches en étant robuste au bruit. Comme nous l’avons présenté dans la Section 5.2, les paramètres d’entrée de l’algorithme sont : l’ensemble des descripteurs sur les données à classifier, un paramètre de distance *eps* et un paramètre *minPts* indiquant le nombre minimal de données pour former une classe. Les expériences présentées à la Section 5.2 avaient pour but d’optimiser l’ensemble des descripteurs pour permettre au système de détecter un maximum de cas synthétiques de carto-vandalisme. Toutefois, les performances du système de détection d’anomalies peuvent être améliorées en optimisant les paramètres *eps* et *minPts*.

Nous relançons plusieurs fois l’algorithme DBSCAN en faisant varier itérativement les valeurs de *eps* et *minPts* pour tenter de maximiser le rappel et la précision dans la détection d’anomalies sur les données de Stuhr et Aubervilliers. En paramétrant DBSCAN pour Stuhr et pour Aubervilliers, nous obtenons des résultats de performance différents. La Table III.13 donne les performance de l’algorithme pour différentes valeurs de *eps* et *minPts*. Le choix des paramètres doit maximiser les valeurs de précision et de rappel. Le paramétrage optimal pour Aubervilliers est donc (0.75,3) alors que celui de Stuhr est (0.9, 3). Ces résultats montrent que le paramétrage de l’algorithme ne sera pas le même en fonction de la zone géographique étudiée. Par conséquent, pour détecter le carto-vandalisme sur une nouvelle zone, nous pouvons reprendre par défaut les valeurs de paramètres déterminées sur Aubervilliers – comme nous l’avons fait pour détecter les anomalies sur les données de Bondy – ou sur Stuhr si le jeu de données porte sur une zone géographique plus proche de cette ville.

TABLE III.13. Paramétrage de DBSCAN sur Stuhr et Aubervilliers

Zone	eps	minPts	Précision	Rappel
Stuhr	0.75	3	0.193	0.795
	0.9	3	0.202	0.795
Aubervilliers	0.75	3	0.221	0.831
	0.9	3	0.225	0.704

5.7 Bilan de la détection non-supervisée du carto-vandalisme

La détection du carto-vandalisme par une méthode non-supervisée de détection d'anomalies a permis de détecter en moyenne 85% du carto-vandalisme synthétique. Par ailleurs, nos expériences ont permis de déterminer des descripteurs avec lesquels certaines contributions de carto-vandalisme apparaissent comme des anomalies, tel qu'un descripteur de fiabilité du contributeur, d'appariement, de syntaxe des tags et de relation topologique avec d'autres objets cartographiques.

Cependant, les résultats de la détection d'anomalies ont renvoyé un grand nombre de faux positifs. La méthode de détection du carto-vandalisme comme anomalie ne permet donc pas, pour l'instant, de détecter de manière automatique le carto-vandalisme. L'analyse de ces faux positifs nous incite à améliorer l'implémentation de certains descripteurs, de manière à éviter de faire apparaître certaines contributions normales comme des anomalies (par exemple, les grands bâtiments publics dont le site internet est renseigné dans un tag). Les faux positifs relevés sont intéressants car, même s'ils ne sont pas *a fortiori* des cas réels de carto-vandalisme, certains n'étaient pas non plus de bonne qualité et mériteraient un contrôle supplémentaire. En plus du grand nombre de faux positifs détectés, certains cas synthétiques de carto-vandalisme n'ont pas été détectés comme des anomalies (voir Figure III.23). Par conséquent, la méthode de détection d'anomalies ne permet pas, pour l'instant, de détecter automatiquement tous les cas de carto-vandalisme.

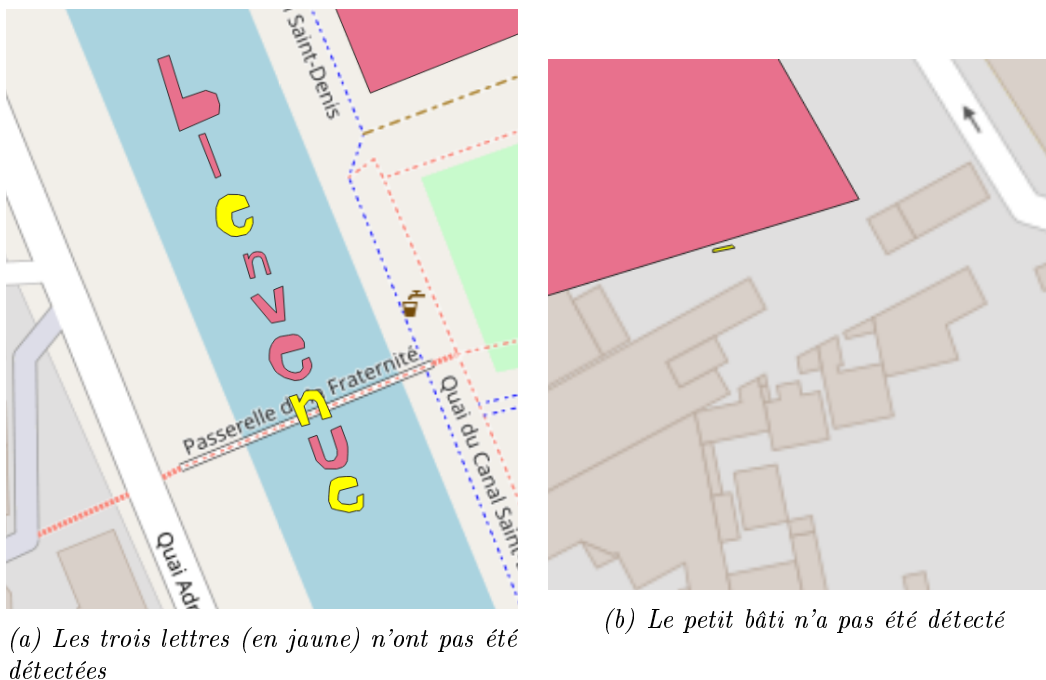


FIGURE III.23. Carto-vandalisme synthétique sur Aubervilliers non détecté avec DBSCAN

Une partie de notre travail de recherche sur la détection non-supervisée du carto-vandalisme dans OSM a été présenté à l'occasion de conférences francophones (Truong *et al.*, 2018c) et internationales (Truong *et al.*, 2018d) en géomatique. Par ailleurs, une version étendue de l'article (Truong *et al.*, 2018c) est actuellement en cours de publication dans la Revue Internationale en Géomatique.

6 Détection du carto-vandalisme par apprentissage supervisé

Pour détecter le carto-vandalisme par une méthode d'apprentissage non-supervisé, nous avons considéré l'hypothèse que les contributions de carto-vandalisme peuvent être vues comme des anomalies, à condition de se placer dans un espace de description adéquat. Cependant, les résultats de nos expériences de détection d'anomalies n'ont pas permis de détecter entièrement tous les cas synthétiques de vandalisme. En particulier, certaines contributions n'apparaissent comme des anomalies dans aucun espace de description : l'hypothèse de modélisation du carto-vandalisme comme anomalie n'est donc pas toujours valable.

Dans cette partie, nous explorons le potentiel des méthodes d'apprentissage supervisé à détecter le carto-vandalisme après une phase d'entraînement sur des exemples de carto-vandalisme synthétique. Il ne s'agit donc plus de détecter les contributions qui apparaissent comme des anomalies, mais d'entraîner un système à reconnaître les différentes caractéristiques du carto-vandalisme. Ainsi, nous espérons que ces méthodes supervisées parviennent à détecter tout les cas de carto-vandalisme, y compris les cas les plus discrets qui n'ont pas été relevés par la détection d'anomalies.

Nous considérons deux méthodes d'apprentissage supervisé : les forêts aléatoires (*random forests*) qui utilisent les métriques calculées sur les objets cartographiques du corpus, et les réseaux de neurones (*convolutional neural network* ou CNN) qui s'appuient sur des données images. Nos expériences visent à étudier la capacité des méthodes d'apprentissage à détecter le carto-vandalisme avec une meilleure précision, ainsi que leur capacité de transfert de l'apprentissage (*transfer learning*). En effet, ces systèmes détectent le carto-vandalisme après avoir été entraînés sur des exemples : il s'agit donc d'explorer leur capacité à généraliser la détection du carto-vandalisme sur différentes zones géographiques. Pour plus de détails sur la notion de transfert de l'apprentissage spatial, nous renvoyons le lecteur aux articles de Iddianozie et McArdle (2019) et Maggiori *et al.* (2017) qui traitent de ce sujet.

6.1 Forêts aléatoires (*random forest*)

L'algorithme des forêts aléatoires (*random forest* ou RF) est une méthode d'apprentissage supervisé. Elle consiste à construire de manière aléatoire une multitude d'arbres de décision qui sont entraînés sur des sous-ensembles des descripteurs en entrée de l'algorithme. La classification s'obtient par le décompte des prédictions faites par chaque arbre de décision. Cette méthode s'est montrée performante dans la détection du vandalisme dans les bases de données ouvertes (*cf.* Table III.1), par conséquent, nous cherchons à évaluer ses performances dans le cadre de la détection du carto-vandalisme.

Par ailleurs, comme il ne s'agit plus de déterminer un espace de description dans lequel le carto-vandalisme apparaît comme une anomalie, pour construire les modèles de forêts aléatoires, nous avons considéré un ensemble de descripteurs différent de celui utilisé pour la détection d'anomalies. L'idée ici est de décrire de manière exhaus-

tive les données à classifier. Les descripteurs considérés pour construire les modèles de forêts aléatoires (que l'on notera modèles RF) sont décrits dans la Table III.14.

TABLE III.14. *Descripteurs utilisés pour construire les modèles RF.*

Variable	Définition
<i>perimeter</i>	périmètre du polygone
<i>shortest_length</i>	longueur du plus petit côté du polygone
<i>median_length</i>	médiane des longueurs des côtés du polygone
<i>elongation</i>	élongation (voir formule III.4)
<i>convexite</i>	convexité du polygone (voir formule III.5)
<i>compacite</i>	compacité du polygone (voir formule III.3)
<i>n_is_within_lulc</i>	nombre de zones naturelles contenant le bâti
<i>n_inter_lulc</i>	nombre de zones naturelles qui intersectent le bâti
<i>max_special_ratio</i>	indicateur syntaxique de caractères spéciaux
<i>n_tags</i>	nombre de tags décrivant l'objet
<i>n_users</i>	nombre de contributeurs uniques de l'objet
<i>timespan_to_previous</i>	durée écoulée entre la version courante et la version précédente
<i>avg</i>	score moyen de fiabilité du contributeur
<i>weighted_avg</i>	score moyen pondéré de fiabilité du contributeur

Nous avons construit trois modèles de forêts aléatoires sur R, en utilisant le package `caret`, entraînés sur trois jeux de données, à savoir :

- Les données sur Aubervilliers
- Les données sur Stuhr
- La fusion des données sur Stuhr et Aubervilliers

Chaque modèle est entraîné avec une validation croisée sur 10 échantillons, répétée trois fois. Des exemples de chaque type de vandalisme sont utilisés pour entraîner le modèle, où 20% de ces exemples sont réservés pour la phase de prédiction uniquement. Pour étudier le transfert d'apprentissage des modèles entraînés, nous utilisons également des jeux de données sur les villes de Heilsbronn en Allemagne, et Lannilis en France. Ces jeux de données appartiennent au corpus que nous avons construit au préalable. Ils sont composés des contributions OSM réelles de bâtiments et de contributions synthétiques de carto-vandalisme.

a) Prédiction sur la zone d'entraînement

La Table III.15 contient les résultats de détection du carto-vandalisme sur Aubervilliers et Stuhr par des modèles entraînés sur chacune de ces zones. Nous constatons que ces deux systèmes détectent entièrement le vandalisme (le rappel est de 100%) sans erreur (la précision est de 100%). Par conséquent, le modèle RF parvient à détecter totalement et correctement le carto-vandalisme dans une zone après avoir été entraîné sur cette même zone.

La Table III.16 contient les résultats de détection sur Aubervilliers et Stuhr par un modèle entraîné sur des données provenant des deux zones. Rappelons que les

TABLE III.15. Résultats de détection avec des modèles RF entraînés sur Aubervilliers et Stuhr.

Zone entraînement / Zone test	Précision	Rappel	Erreur
Aubervilliers / Aubervilliers	1	1	0
Stuhr / Stuhr	1	1	0

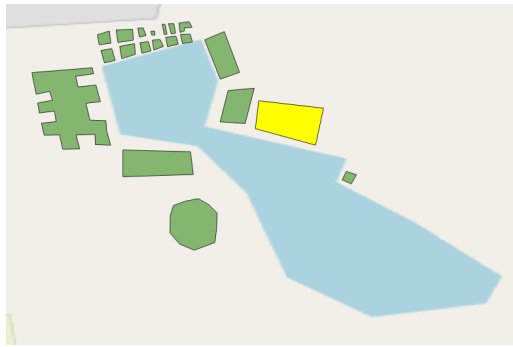
résultats de ces prédictions sont issus des 20% de données qui n’ont pas été utilisées pour entraîner ces modèles. Nous observons que le modèle ainsi construit détecte le carto-vandalisme sur Stuhr avec un rappel 67%, signifiant que certaines contributions synthétiques de carto-vandalisme n’ont pas été détectées. Au contraire, il parvient à détecter entièrement et correctement le vandalisme sur le jeu test d’Aubervilliers. Le modèle semble donc avoir mieux appris à détecter le carto-vandalisme sur Aubervilliers : cela peut s’expliquer par le fait que le jeu de données d’entraînement contenait plus de données sur Aubervilliers que sur Stuhr, et par conséquent, le modèle a appris à reconnaître le vandalisme tout en tenant compte des caractéristiques géographiques des données de la zone française.

TABLE III.16. Résultats de détection avec un modèle RF entraîné sur la fusion des données d’Aubervilliers et Stuhr.

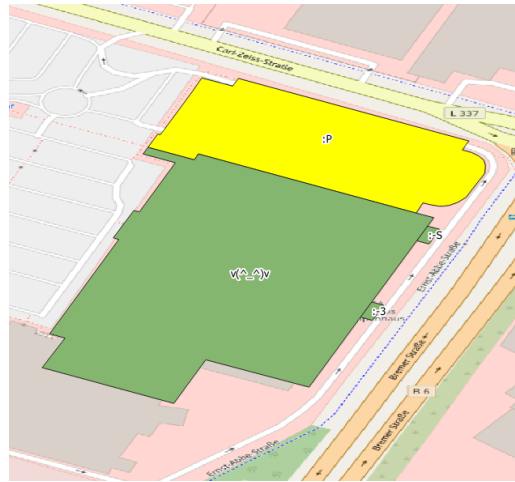
Zone entraînement / Zone test	Précision	Rappel	Erreur
Fusion / Aubervilliers	1	1	0
Fusion / Stuhr	1	0.667	0.002

Le jeu de données de test sur Stuhr contient six cas synthétiques de carto-vandalisme. Nous observons les deux cas classés comme faux négatifs à la Figure III.24. Ces faux négatifs correspondent à un cas de carto-vandalisme fantaisiste (Figure III.24a) et un cas de carto-vandalisme ludique (Figure III.24b). Dans le premier cas, nous avons cartographié ce bâtiment autour d’un lac avec d’autres bâtiments imaginaires, également visibles sur la Figure III.24a. Dans le second cas, nous avons carto-vandalisé un bâtiment commercial existant dans la base OSM en remplaçant le nom de l’enseigne par le tag `name=:P`. Par ailleurs, ce bâtiment jouxte d’autres bâtiments commerciaux que nous avons modifiés de manière similaire pour constituer des cas de carto-vandalisme ludique. Ces bâtiments sont également visibles sur la Figure III.24b. Les deux cas de faux négatifs observés sont donc localisés à proximité de vrais positifs.

En observant ces vrais positifs, nous ne savons pas expliquer pourquoi ces contributions ont pu être détectées par le système alors que d’autres bâtiments vandalisés de manière similaire à proximité ne l’ont pas été. Toutefois, la proximité spatiale des faux négatifs et des vrais positifs peut être exploitée pour détecter les faux négatifs lors d’un post-traitement. En effet, puisque la détection des vrais positifs permet de localiser des zones où l’espace cartographique a été vandalisé, nous pouvons étudier les objets qui se trouvent sur ces zones, et ainsi détecter d’éventuels cas de carto-vandalisme qui n’ont pas été détectés par le modèle RF.



(a) Carto-vandalisme de type fantaisiste



(b) Carto-vandalisme de type ludique

FIGURE III.24. Faux négatifs (en jaune) sur Stuhr obtenus avec le modèle RF entraîné sur une fusion des données de Stuhr et Aubervilliers

b) Prédiction sur des jeux de données du corpus de carto-vandalisme

Pour étudier le transfert de l'apprentissage du carto-vandalisme des modèles RF d'une zone à une autre, nous lançons :

- Une détection du carto-vandalisme sur Stuhr avec le modèle entraîné sur Aubervilliers ;
- Une détection du carto-vandalisme sur Aubervilliers avec le modèle entraîné sur Stuhr.

Les résultats indiqués sur la Table III.17 sont moins bons qu'une détection réalisée sur la zone d'entraînement (voir Table III.15) : cela montre que les modèles RF, en apprenant à détecter le carto-vandalisme, apprennent en même temps les caractéristiques cartographiques de la zone sur laquelle ils sont entraînés. Par conséquent, ils ne parviennent pas à généraliser leur connaissance du carto-vandalisme de manière à le reconnaître sur une autre zone.

Le modèle entraîné sur Aubervilliers détecte le carto-vandalisme de Stuhr avec une précision d'environ 60% et un rappel inférieur à 20%. En revanche, nous remarquons que les résultats de la détection sur Aubervilliers par le modèle entraîné sur Stuhr détecte le carto-vandalisme avec une précision de près de 80% et un rappel de 37%. Le modèle entraîné sur Stuhr semble donc mieux transférer l'apprentissage du carto-vandalisme sur Aubervilliers que le modèle entraîné sur Aubervilliers pour détecter le vandalisme sur Stuhr. Pour l'instant, nos analyses ne permettent pas d'expliquer la différence de sur-apprentissage entre ces deux modèles, mais cette question constitue une perspective de ce travail.

TABLE III.17. Détection sur des zones différentes de la zone d'entraînement

Zone entraînement / Zone test	Précision	Rappel	Erreur
Aubervilliers / Stuhr	0.615	0.182	0.0065
Stuhr / Aubervilliers	0.788	0.366	0.005

◦ les données de la ville de Lannilis avec le modèle RF entraîné sur Aubervilliers. Les résultats de la classification sont indiqués dans la Table III.18.

TABLE III.18. *Détection sur une zone du même pays que la zone d'entraînement*

Zone entraînement / Zone test	Précision	Rappel	Erreur
Stuhr / Heilsbronn	0	0	0.0051
Aubervilliers / Lannilis	0	0	0.0046

Le modèle RF entraîné sur Stuhr ne détecte aucun cas synthétique de vandalisme sur la ville de Heilsbronn. Toutefois, un bâtiment de la ville de Heilsbronn a été détecté comme faux positif, probablement parce qu'il chevauche une zone de broussailles (Figure III.26). Bien que cette contribution ne soit pas un cas réel de carto-vandalisme, sa détection indique que le modèle entraîné sur Stuhr permet de lever des alertes sur des contributions de Heilsbronn dont la qualité est à vérifier, en l'occurrence en ce qui concerne les relations spatiales avec des zones naturelles.

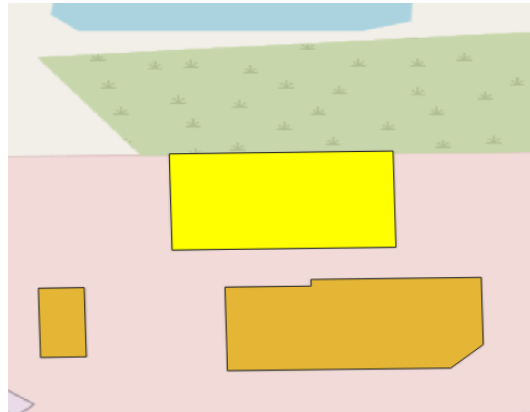


FIGURE III.26. *Faux positif de Heilsbronn (bâti jaune) détecté avec le modèle RF de Stuhr*

Nous remarquons que le transfert de l'apprentissage entre les deux villes françaises et les deux villes allemandes est encore moins bon que le transfert d'apprentissage entre Aubervilliers et Stuhr. De fortes différences géographiques entre la ville d'Aubervilliers et de Lannilis peuvent expliquer un faible transfert d'apprentissage entre ces deux villes (Figure III.27). En effet, Aubervilliers est une ville située en banlieue parisienne, et dont la forte urbanisation se caractérise par de nombreux bâtiments de grande taille correspondant à des immeubles, des bâtiments industriels et commerciaux. En revanche, Lannilis est une ville située en Bretagne, dont le paysage moins urbanisé se caractérise par des bâtiments plus épars, de plus petite taille, et présentant des formes plus régulières, la plupart étant des habitations individuelles. Par conséquent, bien que ces deux villes soient situées sur le même territoire, elles présentent probablement plus de différences urbaines qu'il n'y en a entre Aubervilliers et Stuhr, d'où un meilleur transfert d'apprentissage entre ces deux villes. Cette affirmation mérite toutefois d'être vérifiée en étudiant plus précisément les caractéristiques géographiques de ces villes.

Il s'agit à présent d'explorer si un modèle entraîné sur deux zones différentes parvient à s'abstraire des caractéristiques géographiques propres à chaque zone, et

donc à détecter le carto-vandalisme sur d'autres zones. Nous lançons donc une détection du carto-vandalisme sur les données de Heilsbronn puis sur celles de Lannilis avec le modèle entraîné sur la fusion des données d'Aubervilliers et Stuhr.



FIGURE III.27. Aperçu des couches de bâti (en violet) dans deux villes françaises.
 N.B. : Les deux images sont prises à la même échelle (1 :15000)

D'après les résultats de détection de la Table III.19, on observe que le modèle RF entraîné sur Stuhr et Aubervilliers ne détecte aucun cas synthétique de vandalisme sur Lannilis. En observant les faux positifs qui ont été détectés par ce modèle, nous relevons une caserne de pompier (Figure III.28a) et deux bâtiments de type logement (Figures III.28b et III.28c). Ces faux positifs ne sont pas des cas réels de carto-vandalisme, et par ailleurs, ils ne semblent pas présenter de problèmes de qualité. Ce modèle ne parvient donc pas à transférer l'apprentissage du carto-vandalisme des deux zones d'entraînement sur la ville de Lannilis. Cela est probablement dû au fait que le paysage urbain de cette ville est très différent de celui d'Aubervilliers et de Stuhr.

TABLE III.19. Détection sur Lannilis et Heilsbronn avec le modèle RF entraîné Stuhr et Aubervilliers

Zone entraînement / Zone test	Précision	Rappel	Erreur
Fusion / Lannilis	0	0	0.004
Fusion / Heilsbronn	0.33	0.03	0.0051

En revanche, le modèle RF entraîné sur Stuhr et Aubervilliers détecte le carto-vandalisme sur la ville de Heilsbronn avec une précision de 33% et un rappel de 3%, où nous identifions un vrai positif et deux faux positifs. Le cas synthétique de vandalisme détecté par ce modèle est de type C car c'est un bâti purement fictif (Figure III.29a). Les faux positifs (Figure III.26 et III.29b) sont des bâtiments d'habitation qui, à notre connaissance, ne présentent pas de problème de qualité. Ces résultats de détection, en comparaison avec la détection sur Lannilis, sont légèrement meilleurs. Toutefois, la détection d'un seul vrai positif n'est pas suffisante pour conclure que le modèle RF entraîné sur Stuhr et Aubervilliers parvient à transférer l'apprentissage du carto-vandalisme sur Heilsbronn.

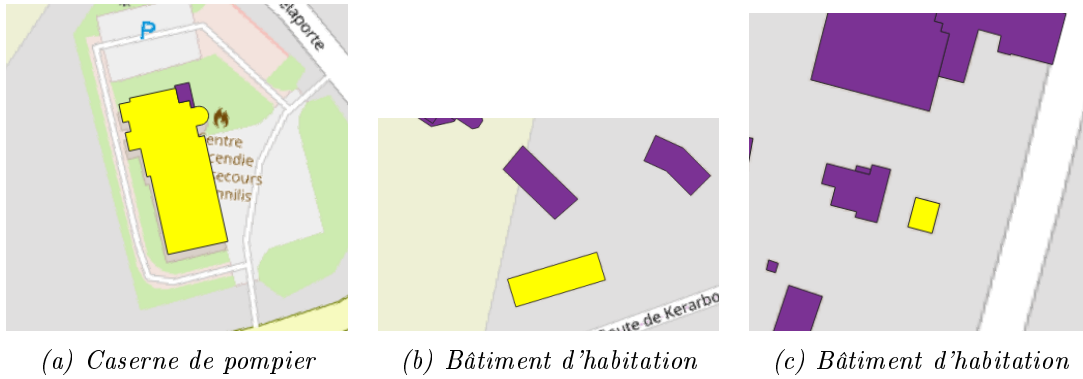


FIGURE III.28. Faux positifs de Lannilis avec le modèle *RF* entraîné sur la fusion Stuhr et Aubervilliers

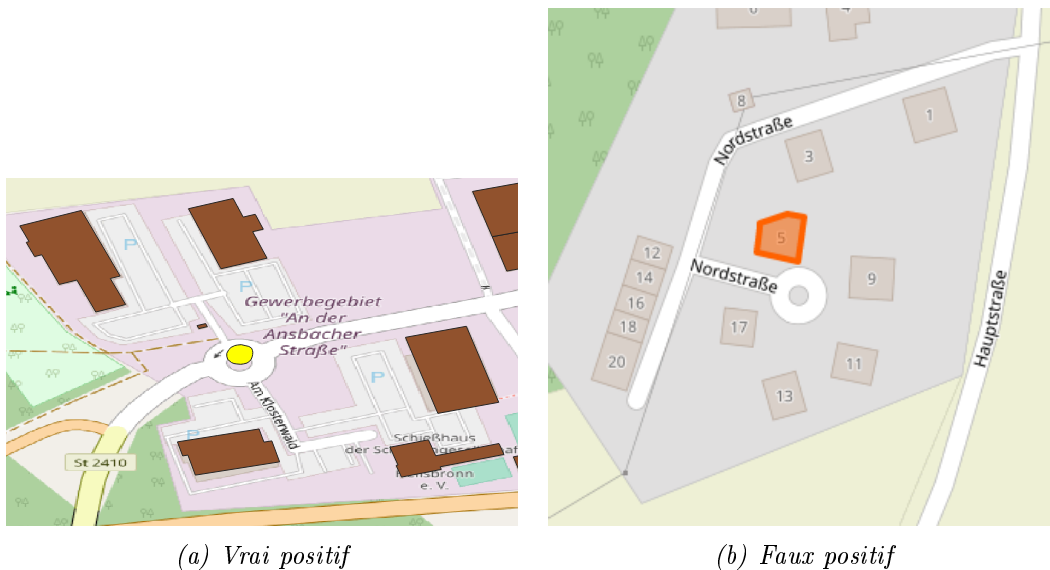


FIGURE III.29. Faux positifs de Heilsbronn avec le modèle *RF* entraîné sur la fusion Stuhr et Aubervilliers

c) Prédiction sur un jeu de données dépourvu de carto-vandalisme synthétique

Jusqu'ici, les résultats avec les modèles *RF* construits montrent que ceux-ci ne permettent pas de détecter le carto-vandalisme sur une zone différente de celle utilisée pour entraîner ces modèles. Or, pour qu'un système de détection du carto-vandalisme soit performant, il ne suffit pas qu'il détecte un maximum de cas réels, mais il faut également que celui-ci détecte un minimum de faux positifs. Pour étudier la capacité des modèles de *RF* à bien classer les données non vandalisées, nous lançons des expériences de détection sur des jeux de données *a priori* dépourvus de carto-vandalisme.

Nous choisissons de lancer une détection de carto-vandalisme sur les données de bâti de la ville de Bondy, qui se trouve dans la même région qu'Aubervilliers, et donc présente globalement les mêmes caractéristiques géographiques. Nous lançons également une détection avec les données de la ville de Fougères, située en Bretagne, comme Lannilis. Les résultats de la détection avec les différents modèles *RF* se

trouvent dans la Table III.20.

TABLE III.20. Détection du carto-vandalisme sur Bondy et Fougères

Zone entraînement / Zone test	Nombre de bâtis détectés
Aubervilliers / Bondy	0
Stuhr / Bondy	0
Fusion / Bondy	0
Aubervilliers / Fougères	0
Stuhr / Fougères	7
Fusion / Fougères	1

Nous observons que les trois modèles ne détectent aucun cas de vandalisme sur la ville de Bondy. En revanche, sur la ville de Fougères, le modèle RF entraîné sur Aubervilliers ne détecte aucun cas de vandalisme, alors que le modèle entraîné sur Stuhr détecte sept cas de vandalisme (voir Figure III.30). Ces cas correspondent à des bâtiments scolaires, un garage et un château. Par ailleurs, le garage (Figure III.30a) a également été détecté comme carto-vandalisme par le modèle entraîné sur Stuhr et Aubervilliers. Cependant, après vérification, aucun de ces objets n'est un cas réel de carto-vandalisme. Ces résultats montrent que les modèles de forêts aléatoires font très peu d'erreur sur la détection du vandalisme dans une zone qui, *a priori*, n'en contient pas. Le maximum d'erreur dans ces résultats est obtenu sur la détection des données de Fougères avec le modèle RF entraîné sur Stuhr, qui ne représente que 0.05% du jeu de données total, ce qui est faible.

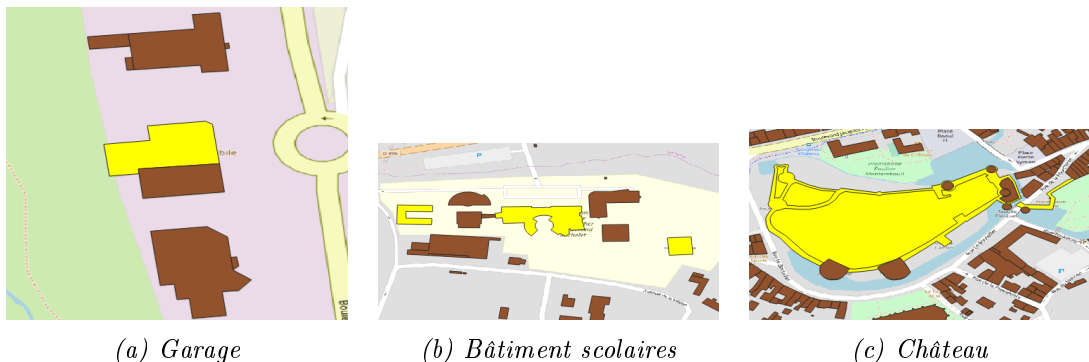


FIGURE III.30. Bâtiments détectés comme vandalisme sur la zone de Fougères avec le modèle RF entraîné sur la zone de Stuhr

d) Bilan de la détection du carto-vandalisme par des forêts aléatoires

La construction de modèles de forêts aléatoires a permis d'entraîner ces systèmes à détecter le carto-vandalisme à partir d'exemples synthétiques. À la différence de la méthode de détection d'anomalies, le choix des descripteurs n'est plus motivé par l'objectif de représenter le carto-vandalisme sous forme d'anomalie, mais par celui de décrire de manière exhaustive les données vandalisées et non vandalisées. En entraînant un modèle RF avec des exemples de vandalisme qui n'apparaissent

pas comme des anomalies sous aucun descripteur, nous espérons parvenir à détecter les cas de vandalisme les plus discrets.

D'après les résultats de nos expériences, les modèles de RF détectent correctement le carto-vandalisme lorsqu'ils sont préalablement entraînés sur la zone à prédire. Toutefois, la détection s'avère moins probante lorsque la zone à prédire s'éloigne de la zone d'entraînement de ces modèles : cela a donc mis en évidence les limites de ces systèmes à apprendre les caractéristiques du carto-vandalisme indépendamment de la zone sur laquelle il se trouve. Or, dans la définition de carto-vandalisme, nous avons vu que celui-ci est fortement lié au contexte (spatial et temporel) dans lequel il se produit. De ce point de vue, il est donc compréhensible que ces modèles peinent à transférer l'apprentissage du carto-vandalisme à d'autres zones. La détection automatique du carto-vandalisme suppose de tenir compte des éléments de contexte géographique des données, alors que le transfert d'apprentissage exige qu'un modèle soit capable de s'abstraire de ce contexte géographique. L'enjeu pour ces modèles d'apprentissage automatique est donc à la fois d'être capable de tenir compte d'éléments de contexte géographique, qui sont essentiels pour détecter le carto-vandalisme, et de transférer la connaissance du vandalisme appris sur une autre région géographique.

6.2 Réseau de neurones à convolution (CNN)

Le développement matériel des machines de calcul a permis à l'apprentissage profond de se révéler comme une méthode puissante pour les problématiques de traitement d'image (LeCun *et al.*, 2015). L'architecture des réseaux de neurones à convolution (CNN) a notamment prouvé son efficacité pour la reconnaissance d'objets dans les images (LeCun *et al.*, 1998). Le principe consiste à extraire les éléments saillants des images à classer, en leur appliquant des filtres de convolution. Les éléments caractéristiques de ces images sont matérialisés sous forme de descripteurs qui servent alors à la classification des images.

Dans le domaine de l'information géographique, l'apprentissage profond a notamment été utilisé pour la classification d'images de télédétection (Ma *et al.*, 2019; Zhu *et al.*, 2017). Récemment, les données géographiques collaboratives ont été exploitées au sein de ces méthodes pour améliorer la classification d'images d'occupation du sol (Audebert *et al.*, 2017; Srivastava *et al.*, 2018). Toutefois, l'utilisation de l'information géographique volontaire en apprentissage profond suppose de tenir compte de ses problèmes de qualité (Chen et Zipf, 2017). Les avantages présentés par des méthodes d'apprentissage profond et notre problématique de détection du carto-vandalisme dans l'information géographique volontaire nous ont donc conduits à explorer le potentiel de ces méthodes à classer les données géographiques collaboratives selon le contexte du carto-vandalisme.

Le principe de fonctionnement des réseaux de neurones à convolution est intéressant pour notre problématique de détection du carto-vandalisme : certains cas synthétiques de vandalisme peuvent être identifiés à l'œil nu en visualisant la couche vectorielle du bâti par-dessus une image aérienne de la zone géographique. Par conséquent, nous pouvons espérer qu'un modèle CNN parvienne à détecter le carto-vandalisme en extrayant les caractéristiques graphiques typiques du vandalisme. En

particulier, le carto-vandalisme de type fantaisiste, qui consiste à créer des éléments réalistes mais fictifs (tels que les bâtiments de la Figure III.24a), pourrait être détecté par un tel modèle. En effet, il pourrait être entraîné à détecter les incohérences graphiques produites par un objet vectoriel qui n'existe pas sur l'image aérienne.

a) Construction du modèle

Dans notre expérience, nous construisons un modèle CNN dont l'architecture est similaire à celle du réseau LeNet (LeCun *et al.*, 1998). La Figure III.31 schématise l'architecture de notre modèle. La phase d'extraction des descripteurs de l'image (*feature extraction*) est réalisée par l'application de filtres de convolution, d'un filtre non linéaire ReLU, puis d'une opération de sous-échantillonnage (Figure III.32). Ces trois opérations sont répétées quatre fois de manière à extraire les éléments les plus saillants de l'image. La phase de classification est réalisée grâce à une couche entièrement connectée (*fully connected layer*).

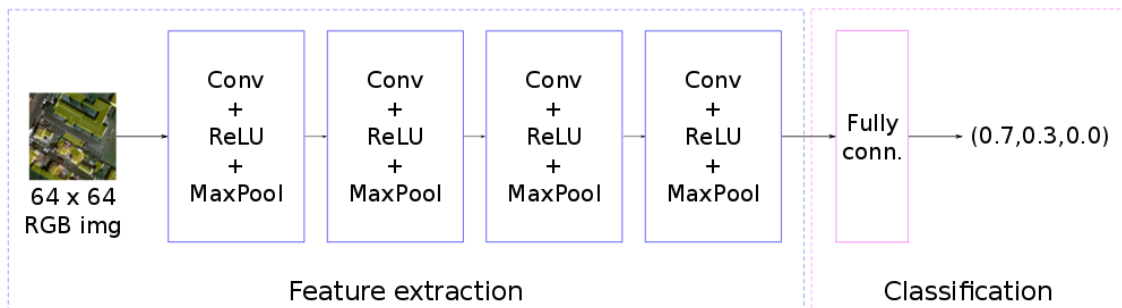


FIGURE III.31. Schéma de l'architecture du modèle CNN implémenté

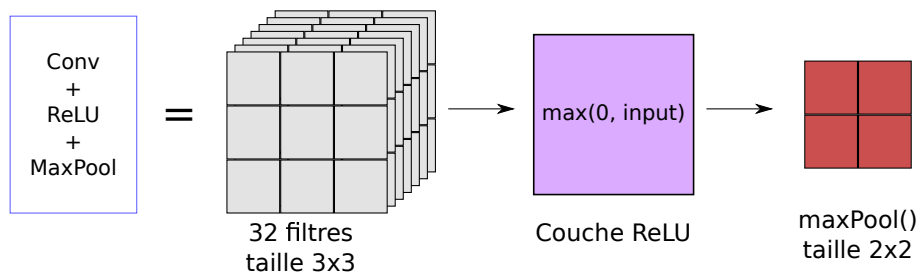


FIGURE III.32. Opérations pour l'extraction de descripteurs

Notre modèle CNN prend en entrée des images RGB de taille 64 pixels x 64 pixels. Nous avons choisi de générer des images de petite taille pour que le temps de traitement de ces images ne soit pas trop long. En effet, l'entraînement du modèle étant effectué sur un grand nombre d'images, plus les images sont grandes, plus le temps de calcul sera long.

Ces images contiennent un extrait d'une image aérienne Bing, sur laquelle se trouve la couche vectorielle du bâti OSM. Nous avons choisi d'afficher la couche vectorielle en jaune pour que ces objets vectoriels ressortent par rapport à l'image aérienne. Toutefois, pour que les pixels de l'image aérienne restent visibles sur les zones où se trouvent aussi les bâtiments vectoriels, nous avons affiché cette couche

avec une opacité de 30%. La génération des images a été effectuée en découpant la zone d'étude selon une grille régulière de 100 mètres \times 100 mètres. Cette échelle a été choisie arbitrairement, en considérant qu'un carreau de 100 mètres \times 100 mètres permettrait de capturer les bâtiments à une taille raisonnable. En effet, pour détecter le carto-vandalisme, il faut faire apparaître les bâtiments avec une taille optimale : de taille suffisante pour que le système puisse traiter les caractéristiques de l'objet bâti, mais la taille de celui-ci ne doit pas être trop grande pour que l'image puisse contenir le contexte géographique du bâtiment. Le fait de capturer des éléments environnant le bâtiment dans l'image permet aussi au modèle de réaliser sa classification : un bâtiment qui se trouve entouré d'autres bâtiments aura plus de chances d'exister réellement qu'un bâtiment cartographié dans un cours d'eau. La construction de ces images est effectuée sur QGIS (avec le plugin Atlas pour exporter chaque carreau sous forme d'image), et la Figure III.33 schématise les différentes couches utilisées pour générer les images.

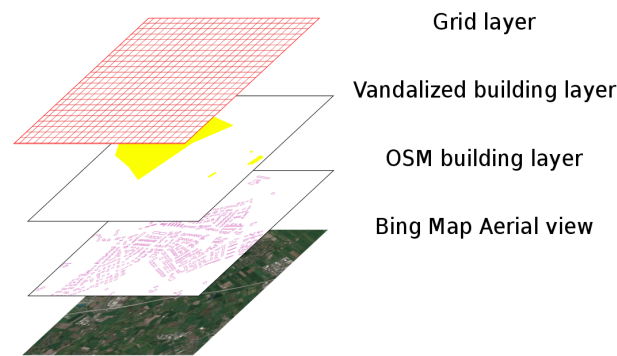


FIGURE III.33. Couches image et vecteurs utilisées pour générer les images à classifier

Ainsi, le modèle CNN est entraîné à classer correctement les images dans lesquelles le bâti OSM est correctement cartographié par rapport à l'image aérienne, et à détecter les images sur lesquelles cette situation n'est pas respectée pour cause de vandalisme. Le modèle CNN classe donc les images selon 3 catégories :

- **Classe 1** : L'image contient du carto-vandalisme
- **Classe 2** : L'image contient uniquement du bâti OSM ordinaire (*i.e.* non vandalisée)
- **Classe 3** : L'image ne contient pas de bâti OSM

De la même manière que les expériences menées avec les forêts aléatoires, nous avons construit trois modèles CNN entraînés sur :

- les images d'Aubervilliers ;
- les images de Stuhr ;
- les images de Stuhr et d'Aubervilliers.

Ces modèles ont été construits en Python, avec la bibliothèque Keras, sur l'environnement Colaboratory de Google⁴.

4. Google Colab met à disposition des GPU et TPU gratuits pour lancer des calculs d'apprentissage automatique.

Pour améliorer la classification, il est possible d’augmenter artificiellement le nombre de données d’entraînement par des opérations de translation et de rotation sur les images. Comme les données synthétiques de carto-vandalisme sont moins nombreuses que les objets OSM réels (voir Table III.21), l’augmentation artificielle d’images permet notamment de rééquilibrer le nombre d’images contenant du carto-vandalisme. Cependant, la classification se situe au niveau des images et non de l’objet cartographique, par conséquent, le paramétrage permet d’améliorer la classification des images sans toutefois entraîner nécessairement une amélioration de la détection des cas de vandalisme au niveau des objets cartographiques.

TABLE III.21. Nombre d’images utilisées pour l’entraînement des modèles CNN

Modèle	Classe 1	Classe 2	Classe 3	Total
Aubervilliers	104	1335	334	1773
Stuhr	538	7629	25281	33448
Fusion	664	8965	25617	35246

Le modèle prend en entrée des images et renvoie une classification de ces images. Or l’objectif est de classer les objets cartographiques OSM qui sont contenus dans ces images. Il faut donc effectuer une étape de transposition de cette classification au niveau de l’image vers une classification au niveau des objets cartographiques. Le carroyage de l’image par une grille de 100 mètres \times 100 mètres peut couper des objets qui ne se trouveront pas entièrement dans une image. Autrement dit, un bâtiment peut se retrouver partiellement dans plusieurs images. Or, ces images ne sont pas forcément classées dans la même catégorie par le modèle CNN. Compte tenu de la classe des images qui contiennent le même bâtiment, ce dernier est alors classifié selon les règles suivantes :

- dès qu’une des images contenant le bâtiment a été classée comme carto-vandalisme, alors le bâtiment est classé comme carto-vandalisme ;
- au contraire, si aucune de ces images n’est classée comme vandalisme, alors le bâtiment n’est pas classé comme vandalisme.

Même si les modèles CNN classent les images selon trois classes, les objets de bâtiments restent classés selon deux catégories (carto-vandalisme ou non carto-vandalisme). Les analyses de nos expériences portent sur la capacité des modèles CNN à classer correctement les objets de bâtiments.

b) Prédiction sur la zone d’entraînement

Comme pour les modèles de forêts aléatoires, nous commençons par étudier des modèles entraînés sur une zone, et leur capacité à détecter le carto-vandalisme synthétique situé sur la même zone. Les résultats de la prédiction du carto-vandalisme sur les zones de Stuhr et Aubervilliers avec les modèles CNN entraînés sur chacune de ces zones sont renseignés dans la Table III.22.

Le modèle CNN entraîné sur Stuhr parvient à détecter quatre cas de carto-vandalisme (sur 29 cas dans le jeu de données de test). En les visualisant aux Figures III.34 et III.36, nous remarquons qu’ils correspondent à des objets de grande

taille. Sur ces figures, l’affichage du carroyage (en jaune) permet d’observer que ces grands bâtiments s’étalent sur plusieurs carreaux. En particulier, les bâtiments de vandalisme des Figures III.34a et III.34b sont de si grande taille qu’ils apparaissent sur plus de 43 carreaux. De plus, sur certains de ces carreaux, la portion du bâtiment peut prendre une grande partie voire toute la surface du carreau. Dans ce dernier cas, l’image résultante peut apparaître complètement jaunée, telle que sur la Figure III.35.

TABLE III.22. Résultats des modèles CNN sur leur zone d’entraînement

Zone entraînement / Zone test	Précision	Rappel	Erreur
Aubervilliers / Aubervilliers	0.015	0.62	0.42
Stuhr / Stuhr	0.105	0.138	0.007

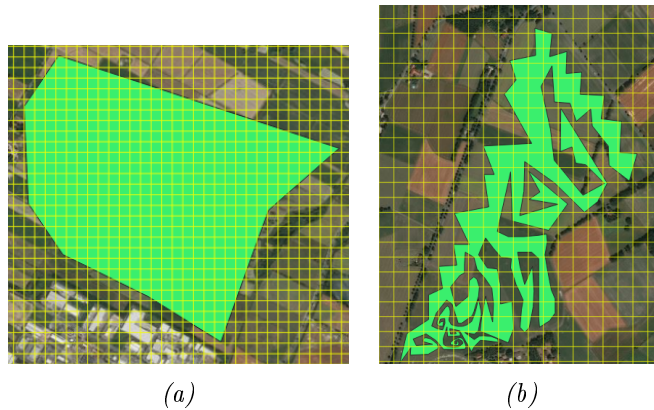


FIGURE III.34. Grands bâtiments détectés par le modèle CNN entraîné sur la zone de Stuhr.

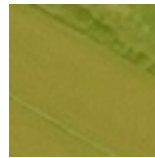


FIGURE III.35. Cas d’un bâtiment trop grand par rapport au carreau de 100*100m



FIGURE III.36. Carto-vandalisme non géométrique détecté par le modèle CNN entraîné sur la zone de Stuhr

Le modèle CNN apprend à classer des images contenant de grands objets de couleur jaune ou des images jaunies (comme celle de la Figure III.35) dans la classe 1. Or, ces cas correspondent au carto-vandalisme qui porte sur les grands bâtiments. Donc le modèle CNN apprend à détecter de grands bâtiments et non à détecter de carto-vandalisme. Par ailleurs, les cas de vrais positifs des Figures III.36a et III.36b sont des bâtiments existants sur lesquels la dégradation a été effectuée sur les tags. Graphiquement, ils ne présentent pas de problème de qualité, mais ce sont de grands bâtiments. Cela renforce l'idée que le modèle CNN a dû classer ces images comme vandalisme car ces grands bâtiments prenaient une grande place sur les images.

En observant le grand nombre de faux positifs détectés sur Aubervilliers avec le modèle entraîné sur cette zone (Figure III.37a), nous constatons qu'ils correspondent à des bâtiments de taille trop grande par rapport aux carreaux de 100 mètres de côté. Au contraire, les vrais négatifs classés par le modèle CNN entraîné sur Aubervilliers sont, pour la plupart, des bâtiments de taille plus petite par rapport à la taille des carreaux (Figure III.37b). L'observation des éléments détectés par ces modèles CNN (faux positifs et vrais positifs) montre à nouveau que ces modèles ont été entraînés à classer les images qui contiennent de grands bâtiments.

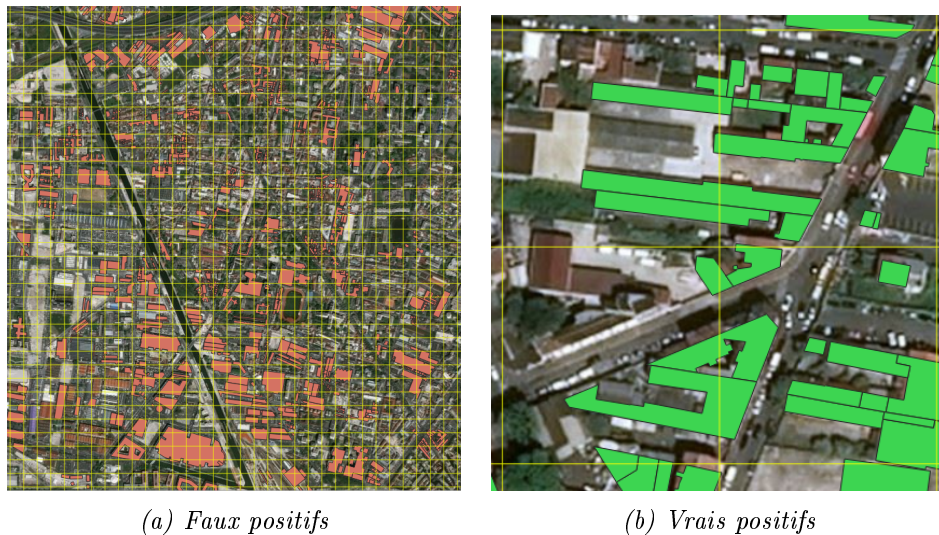


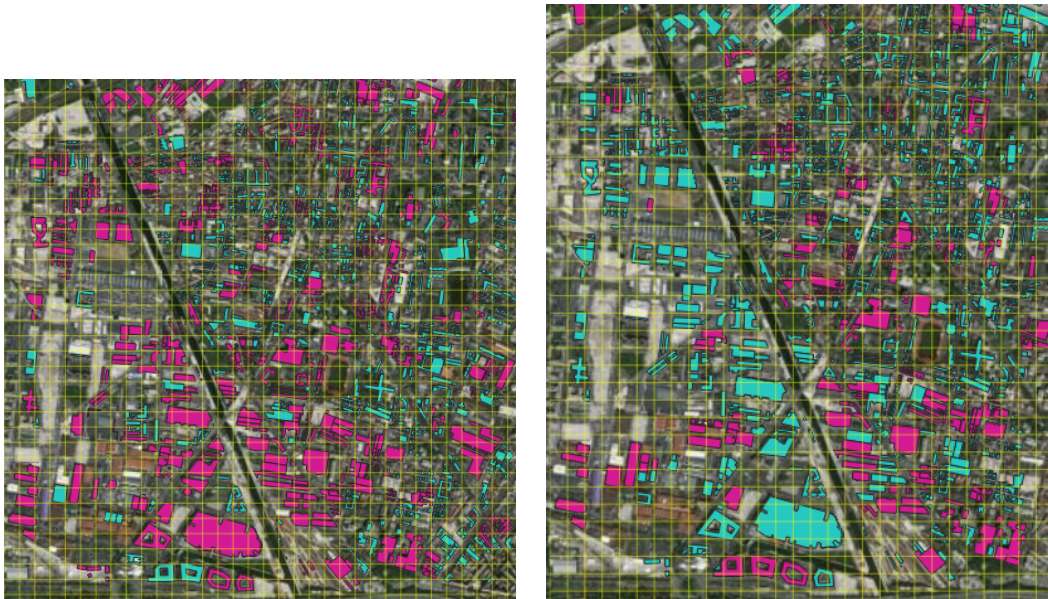
FIGURE III.37. Bâtiments à Aubervilliers classés avec le modèle CNN sur Aubervilliers

c) Prédiction avec un modèle entraîné sur deux zones

Pour étudier la capacité de transfert d'apprentissage des modèles CNN, nous entraînons un modèle sur la fusion des images d'Aubervilliers et de Stuhr. Les résultats de la prédiction sur Aubervilliers et Stuhr avec ce modèle sont indiqués sur la Table III.23. Le modèle entraîné sur les deux zones permet d'obtenir un meilleur rappel pour détecter le carto-vandalisme sur Aubervilliers qu'avec le modèle entraîné uniquement sur Aubervilliers. En effet, le modèle entraîné sur les deux zones détecte moins de faux positifs : sur la Figure III.38, nous observons que certains grands bâtiments normaux ne sont plus détectés comme carto-vandalisme. En cela, l'entraînement sur deux zones a permis au modèle CNN de faire moins d'erreurs sur la détection des bâtiments sur Aubervilliers.

TABLE III.23. Détection du carto-vandalisme avec un CNN entraîné sur Aubervilliers et Stuhr

Zone entraînement / Zone test	Précision	Rappel	Erreur
Fusion / Aubervilliers	0.021	0.38	0.19
Aubervilliers / Aubervilliers	0.015	0.62	0.42
Fusion / Stuhr	0.012	0.231	0.06
Stuhr / Stuhr	0.105	0.138	0.007



(a) Prédiction avec le modèle entraîné sur Aubervilliers

(b) Prédiction avec le modèle entraîné sur Aubervilliers et Stuhr

FIGURE III.38. Test sur Aubervilliers : faux positifs (rose) et vrais négatifs (bleu)

En revanche, le modèle CNN entraîné sur les deux zones détecte le vandalisme sur Stuhr avec moins de précision que le modèle CNN entraîné uniquement sur Stuhr. En effet, le modèle entraîné sur les deux zones détecte moins de faux positifs. L'observation des faux positifs détectés par ces deux modèles CNN à la Figure III.39 ne permet pas de comprendre pourquoi certains bâtiments normaux ont été détectés comme vandalisme par le modèle entraîné sur les deux zones. Quoiqu'il en soit, il semblerait que le modèle CNN n'ait pas été aussi bien entraîné pour détecter le vandalisme sur Stuhr que sur Aubervilliers. Pourtant, le jeu de données d'entraînement utilisé pour construire ce modèle contenait plus de données sur Stuhr que sur Aubervilliers (83% des images de la classe 1 utilisées pour l'entraînement étaient à Stuhr). Cela signifie que la qualité du jeu de données d'entraînement est à revoir : il s'agit non seulement d'entraîner un modèle CNN sur un grand nombre d'images, mais aussi sur des données de qualité.

Rappelons que les analyses de nos expériences précédentes avaient permis de démontrer que les images, telles que construites, ont eu pour effet d'apprendre aux modèles CNN à détecter des gros bâtiments sur les images plutôt qu'à détecter le carto-vandalisme des bâtiments apparaissant sur celles-ci. Ces résultats renforcent l'idée selon laquelle notre méthode de construction des images doit être revue. Plutôt

que d’extraire des données à partir d’un carroyage régulier sur la zone étudiée, nous pourrions adapter le niveau d’échelle de chaque carreau de manière à faire figurer chaque bâtiment en entier dans une image. Les questionnements sur la bonne méthode de construction des images sont essentielles pour exploiter au mieux l’apprentissage profond à des problématiques qui portent sur des données cartographiques vectorielles. Des réflexions actuelles sur ce sujet sont par ailleurs abordées dans d’autres contextes, tels que la généralisation cartographique (Touya *et al.*, 2019), et qui peuvent s’appliquer à notre problématique de détection du carto-vandalisme dans les données OSM.



(a) Prédiction avec le modèle entraîné sur Stuhr

(b) Prédiction avec le modèle entraîné sur Aubervilliers et Stuhr

FIGURE III.39. Test sur Stuhr : faux positifs (rose) et vrais négatifs (bleu)

Nous avons étudié la capacité du modèle CNN entraîné sur Aubervilliers et Stuhr à détecter le carto-vandalisme sur les zones de Heilsbronn et Lannilis. Ces deux zones sont inconnues au modèle CNN puisqu’il n’a pas été entraîné sur les images de ces villes. Les résultats de la détection sont indiqués dans la Table III.24. Globalement, sur des régions inconnues, le modèle parvient à détecter le vandalisme avec une précision inférieure à 10%, un rappel inférieur à 30% et une erreur de près de 8%. Ces faibles performances de classification ne sont pas surprenantes : d’une part, nous savions que l’extraction de descripteurs graphiques à partir d’opérations de convolution sur des images ne suffirait pas pour détecter le carto-vandalisme sur des images. D’autre part, les modèles CNN construits s’étaient déjà montrés peu performants pour détecter le carto-vandalisme sur des zones utilisées pour l’entraînement. Les résultats de la Table III.24 permettent de confirmer que les modèles CNN ne parviennent pas à détecter correctement le carto-vandalisme sur des zones inconnues.

TABLE III.24. Détection du carto-vandalisme avec un modèle CNN entraîné sur Aubervilliers et Stuhr

Zone entraînement / Zone test	Précision	Rappel	Erreur
Fusion / Lannilis	0.097	0.281	0.076
Fusion / Heilsbronn	0.040	0.258	0.077

Les vrais positifs détectés sur Heilsbronn et Lannilis correspondent aux cas où le carto-vandalisme prend la forme de grands bâtiments imaginaires (Figure III.40). Comme les jeux de données sur Stuhr et Aubervilliers contenaient du carto-vandalisme sur des grands bâtiments, le modèle CNN entraîné sur ces zones parvient bien à détecter ce type de carto-vandalisme sur des zones inconnues.



(a) Heilsbronn

(b) Lannilis

FIGURE III.40. Vrais positifs (bâti bleu) détectés par le modèle CNN entraîné sur la fusion des zones de Stuhr et Aubervilliers



(a) Lannilis

(b) Lannilis

(c) Heilsbronn

(d) Heilsbronn

FIGURE III.41. Faux négatifs (bâti marron) classés par le modèle CNN entraîné sur Stuhr et Aubervilliers

En revanche, les cas de carto-vandalisme qui n'ont pas été détectés sur ces zones

(les faux négatifs) correspondent à un type de carto-vandalisme fantaisiste beaucoup moins évident que celui des vrais positifs détectés (Figure III.41). En effet, ces bâtiments imaginaires sont cartographiés à proximité de bâtiments existants, et leur forme géométrique est très similaire à celui du bâti OSM réel. Or, pour générer les images de classification, nous avons choisi d’afficher la couche des bâtiments en jaune avec une opacité de 30%, de manière à permettre aux modèles CNN de détecter les cas où le bâti vectoriel imaginaire est superposé à une image aérienne qui ne contient pas de bâtiment. Au regard de ces faux négatifs, il semblerait que ce choix de représentation ne soit pas adapté : il faudrait donc revoir l’affichage de la couche vecteur de bâti en modifiant sa couleur et/ou sa transparence. À nouveau, ces résultats montrent l’importance de la méthode de construction des images de classification pour apprendre à un modèle CNN à détecter du carto-vandalisme fantaisiste.

d) Prédiction sur des contributions non synthétiques

Nous avons lancé une détection du carto-vandalisme sur Bondy et Fougères, deux zones dépourvues de carto-vandalisme synthétique, par le modèle CNN entraîné sur Aubervilliers. À partir des éléments détectés, nous cherchons à étudier si le modèle détecte des cas réels de carto-vandalisme sur ces zones, et le cas échéant, d’estimer le niveau d’erreur du modèle. D’après les résultats de la Table III.25, le modèle CNN détecte 6% et 12% des bâtiments respectivement à Bondy et Fougères.

TABLE III.25. *Détection du carto-vandalisme sur Bondy et Fougères avec un modèle CNN entraîné sur Aubervilliers*

Zone entraînement / Zone test	Proportion de bâtiments détectés
Aubervilliers / Bondy	0.067
Aubervilliers / Fougères	0.124

D’après une analyse visuelle, ces bâtiments détectés ne semblent pas être des cas réels de carto-vandalisme. Par conséquent, nous pouvons estimer que le modèle CNN entraîné sur Aubervilliers produit une erreur de détection de 6 à 12% selon la zone de prédiction. Nous remarquons que ce taux d’erreur est plutôt faible comparé à l’erreur de prédiction de ce modèle sur Aubervilliers : celui-ci était de 42% et s’expliquait par un grand nombre de faux positifs à Aubervilliers. Or, d’après la Figure III.42, nous observons que les bâtiments qui n’ont pas été détectés comme vandalisme (les vrais négatifs) sont pour la plupart des petits bâtiments. En particulier, il semble que la ville de Bondy contienne plus de petits bâtis correspondant à des habitations individuelles, en comparaison avec Aubervilliers et Fougères. Or, nos analyses précédentes ont montré que le modèle CNN entraîné sur Aubervilliers avait surtout appris à détecter de grands bâtiments. Par conséquent, cela explique pourquoi le modèle CNN ait détecté moins d’objets dans la ville de Bondy, relativement à Aubervilliers et Fougères.

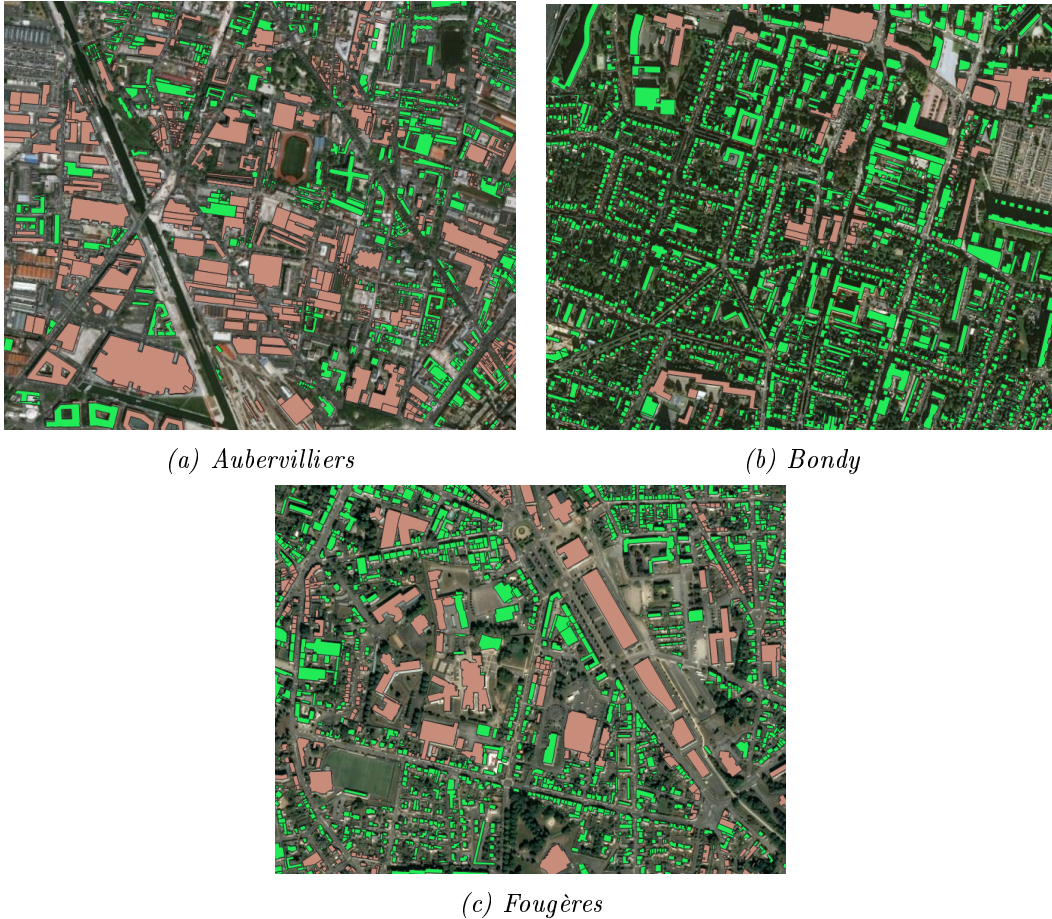


FIGURE III.42. Aperçu des vrais positifs (bâti vert) avec le modèle CNN entraîné sur Aubervilliers

e) Bilan de la détection du carto-vandalisme par des réseaux de neurones à convolution

Les expériences menées avec les modèles CNN entraînés sur des images montrent que ces modèles ne permettent pas de détecter les contributions synthétiques de carto-vandalisme sans erreur. La faible performance de ces modèles peut s'expliquer par la méthode de construction des images de classification. En effet, la génération des images à partir d'un carroyage de taille fixe n'a pas permis aux modèles CNN d'apprendre à détecter des bâtiments de carto-vandalisme mais plutôt des bâtiments de grande taille. Certains cas synthétiques de carto-vandalisme n'étant pas des bâtiments de grande taille, les modèles CNN ne les ont donc pas détectés. À l'inverse, de grands bâtiments ordinaires ont été détectés comme faux positifs.

Une autre méthode de construction des images de classification doit être envisagée pour améliorer la détection du carto-vandalisme par des modèles CNN. Une méthode possible consisterait à découper des images autour de chaque objet bâtiment, avec une échelle adaptée à l'objet. De cette manière, chaque image serait centrée sur un bâtiment qui serait alors complètement contenu dans une seule image. Puisque chaque image correspondrait à un objet cartographique (un bâtiment dans notre cas ici), il faudrait alors augmenter le nombre d'exemples de carto-vandalisme pour équilibrer les différentes classes représentées.

Par ailleurs, le mode d’affichage de la couche vectorielle des objets de bâtiments constitue un autre levier d’amélioration de la classification des images. Il faut que la visualisation de ces objets sur l’image aérienne soit ajustée pour permettre aux modèles de détecter les bâtiments imaginaires correspondant à du carto-vandalisme fictif. Pour cela, la modification de la couleur et la transparence de la couche vectorielle peut être une piste à considérer pour améliorer les images de classification.

Le carto-vandalisme ne se détecte pas uniquement sur des critères visuels : nous avons défini ce phénomène comme un acte de dégradation volontaire de l’espace cartographique. Par conséquent, il est normal que les modèles CNN, entraînés sur des descripteurs purement graphiques, ne parviennent pas à détecter ce phénomène. Toutefois, même si nos premières expériences ne se sont pas révélées efficaces pour détecter le carto-vandalisme par apprentissage profond, cela ne signifie pas que cette approche soit complètement inadaptée pour résoudre notre problématique. Au contraire, certaines pistes méritent encore d’être explorées.

Une première piste serait de construire des modèles d’apprentissage profond plus sophistiqués qui, en plus des descripteurs graphiques, prendraient en compte d’autres descripteurs (tels que ceux proposés dans la Section 5.1) pour réaliser la classification des objets. Une seconde piste serait d’exploiter l’apprentissage profond pour des tâches spécifiques sur lesquelles cette approche a prouvé son efficacité, telle que la segmentation sur des images (Bischke *et al.*, 2017; Marmanis *et al.*, 2016). Dans ce cas, un modèle d’apprentissage profond pourrait permettre d’extraire des objets de bâtiments à partir des images aériennes. Nous pourrions alors considérer une étape d’appariement des objets cartographiques OSM avec les objets issus de cette segmentation, afin de mettre en évidence des cas potentiels de carto-vandalisme fantaisiste ou artistique.

7 Vers une amélioration du corpus de carto-vandalisme

Les résultats expérimentaux présentés précédemment reposent sur les données du corpus de carto-vandalisme que nous avons construit initialement. Dans cette partie, nous discutons des choix méthodologiques de construction du corpus et de leur impact sur les résultats obtenus expérimentalement. La discussion pourra s’ouvrir sur une proposition de pistes d’amélioration du corpus de carto-vandalisme.

Le corpus de carto-vandalisme contient des contributions OSM auxquelles nous avons ajouté des contributions synthétiques de carto-vandalisme. Dans un premier temps, nous nous sommes limités à insérer du carto-vandalisme synthétique sur les objets de bâtiments. Bien que certains cas synthétiques de carto-vandalisme semblent facilement repérables (par exemple les objets de grande taille ou de forme étrange), ces contributions de carto-vandalisme auraient pu exister dans la base de données du projet OSM. En effet, nous aurions pu les charger réellement sur le serveur OSM, afin qu’elles ne soient plus des contributions synthétiques de carto-vandalisme. Pour respecter le bon fonctionnement du projet, nous avons choisi de ne pas les exporter dans le projet OSM, mais d’importer les données du projet dans une base de données PostgreSQL locale dans laquelle nous y avons inséré nos contributions synthétiques. Le réalisme de notre carto-vandalisme synthétique n’est

donc pas à remettre en cause. En revanche, nous ne pouvons pas affirmer que les contributions synthétiques de carto-vandalisme du corpus soit représentatives de tous les types de carto-vandalisme possibles.

L'implémentation des descripteurs de carto-vandalisme sur ces contributions synthétiques a nécessité de synthétiser des profils de contributeurs. Pour cela, nous avons créé de faux profils de contributeur en observant la distribution des indicateurs de contributeurs sur les profils réels de la zone d'étude. Ainsi, nous avons assigné aux contributions vandalisées de type C des profils de contributeurs non fiables pour que le système de détection d'anomalies puisse les retrouver à partir d'un score de fiabilité anormalement bas. En effet, cette méthode ne permet pas de détecter un cas de carto-vandalisme inaperçu aux niveaux géométriques et sémantiques et qui aurait été provoqué par un contributeur fiable.

Nous pouvons remarquer les erreurs faibles obtenues pour la détection avec les modèles RF alors que le nombre de faux positifs et/ou de faux négatifs était important. Cela vient du fait que nous avons créé 1% de carto-vandalisme sur l'ensemble des données de bâtiments OSM sur chaque zone. Il y a donc un fort déséquilibre entre le nombre d'exemples et de contre-exemples de carto-vandalisme dans ce corpus. Or, plus un modèle d'apprentissage est entraîné sur un grand nombre d'exemples de carto-vandalisme, plus il est capable de le détecter. Il est donc souhaitable de disposer d'un grand nombre d'exemples de carto-vandalisme pour améliorer significativement la performance de nos modèles d'apprentissage. Cependant, le nombre de contre-exemples ne doit pas diminuer, pour que les modèles apprennent à distinguer le carto-vandalisme d'une contribution ordinaire.

Les limites de ce premier corpus de carto-vandalisme résident dans le faible nombre d'exemples de carto-vandalisme et dans le manque de variété dans les formes de carto-vandalisme. Pour améliorer le corpus, nous pourrions commencer par intégrer les quelques exemples de carto-vandalisme rencontrés dans notre exploration des contributeurs bannis au Chapitre 1. Cela supposerait de reconstituer un *snapshot* de la zone sur laquelle l'objet de carto-vandalisme était encore visible. Ainsi, nous obtiendrions des fenêtres spatio-temporelles de données réelles contenant du carto-vandalisme réel connu *a priori*.

Plusieurs pistes peuvent être envisagées à partir des méthodologies proposées dans l'état de l'art pour récolter des exemples réels de carto-vandalisme en plus grand nombre. La première piste serait d'avoir recours à l'annotation humaine. Les services de *crowdsourcing*, de type *Amazon's Mechanical Turk*, peuvent être envisagées pour obtenir des données annotées en grande quantité et de qualité très correcte. Toutefois, cette piste est coûteuse en ressources humaines, et est également sujette aux risques et aux limites du *crowdsourcing*. Par ailleurs, il est important d'avoir conscience des questions éthiques liées aux conditions de travail des personnes embauchées pour réaliser de telles tâches (Tubaro et Casilli, 2019), avant de choisir d'utiliser un service de *crowdsourcing* pour annoter ses données.

La deuxième piste serait d'envisager des méthodes d'annotation automatique. La conception de ces méthodes nécessite de filtrer avec précision un maximum de cas réels de carto-vandalisme. Par exemple, dans le cas des données OSM, assimiler les opérations de restauration à du carto-vandalisme peut être une approximation qui

ne permet pas de les récupérer avec précision. La qualification des contributeurs en vue de récupérer les contributions produites par ceux qui ont été détectés comme carto-vandales est une autre possibilité. Toutefois, cette solution exige de mettre en place une méthode de qualification des contributeurs assez précise pour différencier les contributeurs non fiables des carto-vandales.

Conclusion du chapitre

Dans ce chapitre, nous avons dressé un état de l'art de la détection du vandalisme dans les bases de données ouvertes. Les méthodes et les métriques proposées dans les travaux existants ont permis d'identifier le besoin d'un corpus de données annotées spécifiquement pour le carto-vandalisme. Un tel corpus de carto-vandalisme est essentiel pour permettre d'évaluer les performances des méthodes d'apprentissage à détecter ce phénomène dans les contributions. Par conséquent, nous avons été amenés à produire un premier corpus de carto-vandalisme composé de données OSM réelles et synthétiques.

Notre démarche expérimentale s'est articulée autour de trois objectifs, à savoir : la détermination des bonnes variables descriptives du carto-vandalisme, l'évaluation des méthodes d'apprentissage pour la détection du carto-vandalisme et l'amélioration du corpus de données annotées. Nous avons montré expérimentalement que la détection du carto-vandalisme requiert des descripteurs sur les contributeurs : cela confirme l'intérêt de qualifier les contributeurs pour qualifier l'information géographique volontaire. De plus, nous avons également montré l'intérêt de tenir compte de données de référence lorsque cela est possible. En effet, l'appariement entre des données collaboratives et des données de référence permet de mettre en évidence d'éventuelles incohérences qui peuvent être difficilement détectées grâce à des descripteurs purement intrinsèques sur les données collaboratives.

Les résultats sur les expériences issues d'une méthode non-supervisée de détection d'anomalies ont permis de montrer que cette approche peut permettre d'aider les contributeurs à corriger le carto-vandalisme. La méthode supervisée des forêts aléatoires s'est révélée très intéressante pour détecter le carto-vandalisme, à condition que le modèle ait été préalablement entraîné sur la zone à prédire. Les résultats moins satisfaisants de détection par les réseaux de neurones convolutifs nous incitent à explorer davantage cette approche pour déterminer dans quelle mesure elle peut être exploitée pour détecter le vandalisme cartographique.

Enfin, les résultats expérimentaux ont permis d'identifier certaines limites provenant des contributions synthétiques de carto-vandalisme qui composent initialement le corpus de données annotées. Pour pouvoir améliorer nos résultats et valider les analyses qui en découlent, nous suggérons de construire un corpus qui contienne de nombreux cas réels de carto-vandalisme présents dans un projet d'information géographique volontaire. Cela implique de mettre en place une méthodologie d'annotation, manuelle ou automatique, de données réelles.

Conclusion Générale

Face à l'utilisation grandissante des données géographiques dans diverses applications et à l'intérêt porté aux données géographiques collaboratives grâce à leur accessibilité, le besoin de qualifier ces données est urgent. Dans une démarche de qualification de l'information géographique volontaire, ce travail de thèse s'est focalisé plus particulièrement sur la problématique du vandalisme cartographique.

Ce phénomène, qui corrompt la qualité des données collaboratives géographiques, est encore très peu étudié. Notre travail avait donc pour objectif premier de dégager une définition de ce phénomène dont les caractéristiques sont encore floues, dans le but de détecter les contributions cartographiques qui relèvent du carto-vandalisme. Le second objectif de cette thèse consistait en particulier à qualifier le niveau de confiance des contributeurs de données géographiques, car notre thèse repose sur le fait que la qualité de l'information géographique est fortement liée à la qualité de son contributeur. Enfin, dans le but ultime de détecter les données qui relèvent de carto-vandalisme, nous cherchions à étudier le potentiel des méthodes d'apprentissage automatique à réaliser cette tâche de classification.

En réponse aux objectifs qui viennent d'être rappelés, nous résumons ici les principales contributions de ce travail de recherche. Ce travail de thèse a permis de faire avancer la connaissance du vandalisme cartographique et sa détection dans l'information géographique volontaire. Cependant, notre recherche ne constitue qu'un point de départ à d'autres pistes de recherche sur ce sujet. Aussi, nous présentons des perspectives de notre travail qui pourront être considérées dans le futur.

Bilan des contributions

Définition du carto-vandalisme

Nous avons défini le carto-vandalisme à partir d'une recherche théorique et empirique. Notre recherche théorique a permis de mettre en lumière les évolutions de la définition du vandalisme dans différents contextes, à commencer par le contexte historique, législatif puis numérique. Nous avons dégagé trois composantes du vandalisme autour desquelles tourne la qualification de ce phénomène dans le cadre de l'information géographique volontaire. En effet, la difficulté à détecter le carto-vandalisme provient de la difficulté à qualifier les cas de dégradation de l'espace cartographique pour lesquels le contributeur a agi librement et volontairement. C'est précisément sur ce point que la détection du carto-vandalisme se distingue de la

qualification de la mauvaise qualité de l'information géographique volontaire. Une exploration de la liste des contributeurs bannis du projet OSM a permis de trouver empiriquement quelques cas réels de vandalisme cartographique qui correspondaient à notre définition théorique. De plus, notre analyse de la seule typologie de carto-vandalisme actuellement proposée dans la littérature montre que ce phénomène peut prendre une multitude de formes cartographiques et peut se justifier par diverses motivations, rendant sa classification ardue.

Modélisation du comportement collaboratif

Nous avons montré que la qualification du contributeur était primordiale pour qualifier les contributions, et en particulier dans le cadre de la détection du carto-vandalisme. Nous avons cherché à étudier le comportement collaboratif des contributeurs de données, à travers leurs interactions sur la plateforme de contribution OpenStreetMap, dans le but d'évaluer leur niveau de confiance. Pour étudier ces interactions, nous avons proposé un modèle de réseau multiplexe des contributeurs dans lequel il était possible de représenter les diverses façons de contribuer des données géographiques de manière collaborative. Le modèle proposé a permis d'identifier des profils de contributeurs qui témoignent d'un certain niveau de confiance. Ces contributeurs assurent la qualité de l'information géographique volontaire, par conséquent, leur identification est nécessaire à la détection du carto-vandalisme qui provient justement de contributeurs non fiables.

Évaluation de la confiance des contributeurs

La qualification des contributeurs de données géographiques collaboratives vise à évaluer leur niveau de confiance. Notre travail s'est donc attaché à étudier la qualité des différents indicateurs de confiance à évaluer les contributeurs sur cet aspect. En particulier, nous avons proposé de considérer des métriques qui révèlent leurs interactions sur la plateforme collaborative avec le reste de la communauté ainsi que leur investissement dans la durée et dans l'espace des contributions à qualifier. Ces nouvelles métriques ont été intégrées selon deux méthodes d'évaluation de la confiance :

1. sous forme d'un score moyen (pondéré ou non) de confiance du contributeur ;
2. sous forme d'un classement généré par un algorithme de décision multicritère.

En proposant de tenir compte de métriques supplémentaires, nous avons cherché à confronter notre évaluation des contributeurs à l'indicateur de confiance basé sur le nombre de leurs contributions. Nous avons montré que la prise en compte de ces métriques additionnelles permet d'évaluer les petits contributeurs dont la faible participation ne traduit pas nécessairement un faible niveau de confiance.

Création d'un corpus de données annotées de carto-vandalisme

Comme il n'existe actuellement pas de corpus de données contenant des exemples de carto-vandalisme, une partie de notre travail a permis de constituer un tel cor-

pus de données annotées pour le carto-vandalisme. Nous avons construit ce corpus de données en y récupérant des contributions réelles du projet OSM *a priori* non-vandalisées, et nous y avons ajouté des contributions de carto-vandalisme synthétique. Par ailleurs, notre qualification manuelle de la fiabilité des contributeurs a, d'une certaine manière, également permis de constituer un corpus de profils de contributeurs.

Mise en place d'indicateurs de détection du carto-vandalisme

Il s'agissait de choisir des descripteurs qui mettent en évidence les contributions produites intentionnellement pour dégrader l'espace cartographique. D'après cette définition du carto-vandalisme, nous avons donc proposé des descripteurs qui portent sur :

- la confiance des contributeurs de données géographiques collaboratives ;
- la qualité de l'information géographique volontaire.

Nous avons étudié l'intérêt de ces descripteurs à l'aide d'une méthode de détection d'anomalies, nous permettant de valider la pertinence d'un ensemble d'indicateurs de détection de carto-vandalisme sur les cas synthétiques contenus dans notre corpus de données annotées.

Évaluation des méthodes de détection automatique

Nous avons utilisé des méthodes d'apprentissage non-supervisé et supervisé sur notre corpus afin d'évaluer leur performances de détection du carto-vandalisme. Nos expériences sur des méthodes de détection d'anomalies, de forêts aléatoires et de réseaux convolutifs ont permis d'étudier les capacités de ces algorithmes à :

- détecter sans erreur un grand nombre de contributions de carto-vandalisme ;
- apprendre à détecter le carto-vandalisme à partir d'un nombre limité d'exemples d'entraînement ;
- apprendre à détecter le carto-vandalisme sur des zones qui présentent des propriétés géographiques différentes.

Application des mécanismes de qualité de l'information géographique volontaire

Notre état de l'art a permis de présenter les nombreux travaux scientifiques qui développent chacun des trois mécanismes d'assurance qualité de l'information géographique volontaire formulés par [Goodchild et Li \(2012\)](#), à savoir l'approche participative (*crowdsourcing*), l'approche sociale et l'approche géographique. Cependant, nous avons constaté un manque d'approche globale de ces méthodes : en général, elles se focalisent uniquement sur l'une de ces trois approches. Par conséquent, nous avons adopté une approche globale dans le cadre de notre problématique. Pour cela, la détection du carto-vandalisme dans les données géographiques collaboratives s'est

effectuée en combinant des critères qui portent sur chacun des mécanismes d'assurance qualité de l'information géographique volontaire.

Approche participative

Nous avons proposé des métriques rendant compte de l'aspect collaboratif des contributeurs sur la qualité des données géographiques. En particulier, nous avons modélisé des graphes de collaboration pour mesurer la participation des contributeurs à partir de l'historique des contributions faites sur les objets cartographiques. Par ailleurs, les métriques décrivant l'historique d'un objet cartographique font également partie de l'approche participative. En effet, ces métriques calculées sur l'historique des objets permettent, par exemple, de mesurer le nombre de contributeurs uniques sur un objet ou la durée écoulée entre deux versions d'une donnée. D'une certaine manière, ces métriques permettent de mesurer à quel point une donnée suit la loi de Linus, qui est à l'origine de l'approche participative.

Approche sociale

L'adoption d'une approche sociale pour la détection du carto-vandalisme est au cœur de notre thèse, puisqu'il s'agissait de qualifier le contributeur pour qualifier la contribution cartographique. Notre étude du réseau social multiplexe a permis notamment d'identifier des profils de modérateurs, qui permettent de garantir la qualité de l'information géographique volontaire, d'après l'approche sociale. Nos expériences de détection du carto-vandalisme par les méthodes d'apprentissage ont par ailleurs confirmé la nécessité de prendre en compte des métriques « sociales », en particulier pour modéliser l'aspect intentionnel de l'acte de carto-vandalisme.

Approche géographique

Notre méthode de détection du carto-vandalisme s'est également appuyée sur des critères géographiques. En particulier, nous avons proposé des métriques portant sur la géométrie des objets cartographiques et sur les relations spatiales avec d'autres objets. Ces métriques permettent de mettre en évidence des incohérences géométriques, topologiques, ou thématiques. Par ailleurs, l'appariement des données collaboratives avec des données de références fait partie de l'approche géographique. Cette méthode permet également de qualifier la cohérence géographique de l'information géographique volontaire par rapport à une base de référence.

Perspectives

Sur la qualification des contributeurs

Exploiter le réseau social multiplexe

La construction et l'étude d'un réseau social multiplexe sur les contributeurs OSM a permis d'identifier des profils de pionniers et de modérateurs. Une perspective de ce travail serait d'automatiser la détection de ce type de profil. Par ailleurs, pour mieux exploiter un réseau multiplexe, nous pourrions chercher à améliorer sa composition. En effet, dans ce travail, nous avons proposé un certain nombre de graphes de collaboration qui pourraient être considérés dans un modèle de réseau multiplexe. Toutefois, nous avons vu que certaines couches pouvaient être plus ou moins similaires. Par conséquent, l'étude de la similarité des couches d'un réseau pourrait permettre de valider les modèles construits (Bródka *et al.*, 2018). Le fait que deux couches soient similaires n'est pas un problème en soi, mais il s'agit d'avoir identifié cette similitude dans l'analyse du réseau, et par conséquent, dans la détection des communautés de contributeurs.

Nous avons identifié des profils de modérateurs et de pionniers dans un cas d'étude bien précis (sur l'Île de la Cité). Ce résultat pourrait permettre de chercher à vérifier, grâce à la construction et l'exploitation d'un réseau social multiplexe sur d'autres zones, s'il existe des modérateurs et/ou des pionniers quelle que soit la fenêtre spatiale étudiée. Dans le cas où une zone en serait dépourvue, il faudrait chercher des éléments provenant d'un autre mécanisme que celui de l'approche sociale pour garantir la qualité de l'information géographique volontaire.

Par ailleurs, l'analyse de ce type de réseau multiplexe pourrait être approfondie, pour permettre éventuellement d'identifier d'autres profils. Dans le cadre de la détection du carto-vandalisme, nous pourrions chercher s'il est possible de développer une méthode de détection des carto-vandales sur les réseaux multiplexes. Plus généralement, l'étude de ce réseau peut aussi aider à comprendre le comportement collaboratif des contributeurs « moyens » et des petits contributeurs dont la fiabilité n'est pas facile à évaluer, vu leur faible quantité de contributions.

Améliorer l'évaluation de la confiance du contributeur

Nous avons étudié deux méthodes d'évaluation de la confiance du contributeur. La première méthode est un score de confiance agrégé à partir de métriques sur les contributeurs, et la seconde est issue d'un système de décision multicritère. Ces deux méthodes se sont avérées utiles pour nuancer l'hypothèse trop forte selon laquelle un grand contributeur est fiable. Une des perspectives sur ce travail serait d'améliorer ces méthodes d'évaluation de la confiance du contributeur. Il s'agira de chercher la meilleure pondération des métriques de contributeurs pour la méthode de confiance agrégée, et de paramétrer l'algorithme PROMETHEE II pour la méthode multicritère.

Nous avons montré dans ce travail que la prise en compte de la confiance de

contributeur, en tant que source de données, permettait de mieux qualifier les contributions de carto-vandalisme. Une autre perspective d'amélioration de l'évaluation de la confiance du contributeur peut consister à améliorer la modélisation de cette confiance. Il s'agirait d'équilibrer les métriques de contributeurs selon les cinq dimensions de confiance d'une source de données que sont : l'intention, la sincérité, l'implication et la rumeur (Lesot et Revault d'Allonnes, 2017).

Le niveau de confiance d'un contributeur est une qualité qui n'est pas figée. Au contraire, elle évolue dans le temps et dans l'espace : un contributeur non fiable peut devenir fiable avec le temps, et inversement. De même, un contributeur peut être enclin à enrichir une zone, et être motivé à dégrader une autre zone. Par conséquent, une perspective d'amélioration de l'évaluation de la confiance du contributeur serait de chercher à modéliser l'aspect dynamique de la confiance. Il pourrait s'agir d'évaluer la confiance d'un contributeur sur différentes zones spatio-temporelles, dans le but de mieux connaître l'évolution de ses intentions et de ses compétences par exemple.

Sur la qualité du corpus de carto-vandalisme

Augmenter le nombre d'exemples annotés

L'entraînement et l'évaluation des méthodes de qualification des contributions et des contributeurs ont nécessité de disposer de données annotées. Or, les limites de nos méthodes et de nos analyses proviennent notamment du manque de données annotées (sur les contributeurs et sur les contributions). Il nous semble donc indispensable d'augmenter la taille de ce corpus, en particulier en y ajoutant un plus grand nombre d'exemples de carto-vandalisme. L'idéal serait d'intégrer un grand nombre d'exemples réels, avec des profils de contributeurs réels.

L'alimentation de ce corpus par des exemples réels peut s'effectuer de diverses manières. Nous pourrions commencer par reconstituer les *snapshots* correspondant aux cas réels de carto-vandalisme identifiés durant notre exploration des données produites par les contributeurs bannis d'OSM. Cela permettra également d'intégrer dans ce corpus des profils réels de carto-vandales dont proviennent ces cas de vandalisme.

Pour annoter un grand nombre de données, des méthodes d'annotation automatique ont été proposées dans l'état de l'art. Nous pourrions reprendre ces méthodes, mais nous avons vu que leur utilisation comportait des risques d'erreur d'annotation. Par conséquent, une solution serait de s'inspirer de ces méthodes automatiques pour filtrer un certain nombre de données « douteuses » qui seraient annotées manuellement par la suite. La récupération de ces données peut s'effectuer à partir des contributeurs qui ont été bannis au moins une fois, où à partir des sessions d'édition qui a suscité des commentaires de la part d'autres contributeurs, dans lesquels le terme « vandalisme » apparaît dans les discussions. Cela suppose de développer des outils de traitement textuels de ces échanges de commentaires.

De plus, nous pourrions également faire appel à la communauté des contributeurs du projet OSM pour aider à répertorier des cas réels de carto-vandalisme actuels

ou passés. Dans ce dernier cas, ces données ont pu être corrigés sans que le DWG n'ait banni le carto-vandale⁵. En demandant aux contributeurs d'OSM de fournir les lieux et les dates auxquels les cas de carto-vandalisme sont visibles sur la carte, nous pourrions reconstituer des *snapshots* de données qui contiennent réellement du vandalisme OSM. Comme Wikipédia possède une page qui répertorie tous les articles qui sont les plus vandalisés, une autre idée serait de faire appel à la communauté des contributeurs d'OSM pour leur demander de localiser des zones ou des objets qui sont les plus sujets au vandalisme cartographique.

Développer des métriques pertinentes

L'annotation d'un grand nombre de données ne suffit pas pour constituer un corpus de vandalisme de qualité. Ces données doivent être décrites par un ensemble de métriques pertinentes qui permettent de distinguer une contribution de carto-vandalisme d'une contribution ordinaire. Nos expériences de détection d'anomalie ont permis de déterminer les descripteurs optimaux dans lesquels les données de vandalisme apparaîtraient comme des anomalies.

Nous avons vu que cette hypothèse de modélisation du carto-vandalisme comme anomalie présente des limites. Toutefois, la nécessité de déterminer les métriques qui discriminent les deux types de données persiste. Une perspective possible serait d'étudier la notion de typicalité des données décrites par ces descripteurs pour mesurer leur niveau de représentativité dans la classe à laquelle ils appartiennent (carto-vandalisme ou non vandalisme), dans le but de sélectionner les descripteurs qui discriminent au mieux chacune des classes (Lesot *et al.*, 2008).

Sur la qualité de l'apprentissage automatique

La classification des données par des méthodes de forêts aléatoires trouve des limites dans le transfert d'apprentissage d'une zone à une autre. Mais ce constat a été fait intuitivement, en partant de l'hypothèse que deux villes appartenant au même pays présentent des caractéristiques géographiques communes. Or, il faudrait justement vérifier cette hypothèse, et dans quelles mesures elle reste valable. Concrètement, il faudrait étudier la similarité cartographique des différents *snapshots* utilisés, pour pouvoir évaluer plus précisément la capacité de transfert d'apprentissage des forêts aléatoires, et plus généralement des méthodes d'apprentissage supervisé.

Quant aux expériences menées par apprentissage profond, nos premiers résultats nous incitent à revoir l'étape de préparation des données. En particulier, il faudrait expérimenter d'autres méthodes de génération d'images sur les données, en envisageant d'autres modes de représentation des données vectorielles sur ces images, ou en essayant d'adapter l'échelle de zoom sur les objets capturés sur ces images. L'objectif est de générer des images qui contiennent assez de détails pour permettre aux modèles de CNN d'extraire des descripteurs pertinents à la détection de contributions de carto-vandalisme. Il faudra également chercher à mettre en place

5. Le DWG n'intervient pour bannir les contributeurs que dans des cas très sérieux de dégradation du projet.

une architecture d'apprentissage profond qui ne prend pas uniquement en compte les descripteurs graphiques extraits des images, mais également les autres que nous avons proposés.

Enfin, sur le long terme, il serait envisageable de combiner différentes méthodes d'apprentissage (qu'il soit supervisé ou non) au sein d'un processus hybride de détection du vandalisme. Dans ce cas, chaque méthode serait utilisée pour une tâche précise, dans laquelle son efficacité est prouvée. Par exemple, nous pourrions utiliser l'apprentissage profond pour effectuer uniquement l'appariement d'objets vectoriels avec les objets présents sur des images aériennes. Le résultat d'appariement pourrait alors être utilisé comme un descripteur supplémentaire à intégrer dans les données d'entrée de la classification par un modèle de forêt aléatoire. Le résultat de la classification pourrait alors être confirmé ou contredit par le résultat d'une détection d'anomalies sur ces données.

Du carto-vandalisme à la qualité de l'information géographique volontaire

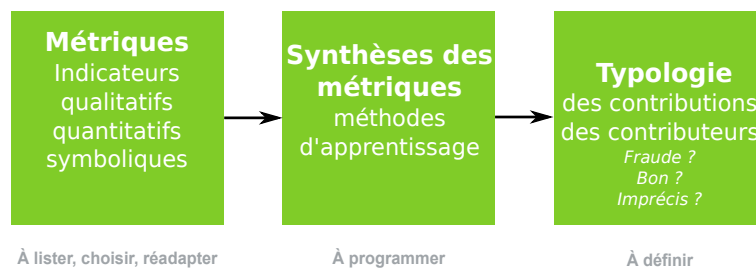


FIGURE 1. *Méthodologie globale de qualification de l'information géographique volontaire*

Ce travail de thèse s'est focalisé sur la détection du vandalisme cartographique. Toutefois, ce cadre d'étude est fortement lié à la question de la qualité des données géographiques collaboratives. Notre démarche de qualification des données de carto-vandalisme s'est attachée à la mise en place de métriques adaptées à notre problématique, puis à qualifier les données à partir de méthodes d'apprentissage selon une typologie propre au carto-vandalisme (Figure 1).

De la même manière, cette démarche a également été adoptée pour qualifier les contributeurs de données géographiques : à partir d'un ensemble de métriques sur les contributeurs, nous avons utilisé des méthodes d'apprentissage (non-supervisé en l'occurrence) pour identifier des profils de contributeurs. Cette méthodologie ne se limite donc pas seulement à la détection du carto-vandalisme et pourrait être reprise pour qualifier l'information géographique volontaire – et ses contributeurs – dans d'autres cas d'application.

Bibliographie

- ADLER, ALFARO, L., MOLA-VELASCO, S. M., ROSSO, P. et WEST, A. G. (2011). Wikipedia Vandalism Detection : Combining Natural Language, Metadata, and Reputation Features. In GELBUKH, A., éditeur : *Computational Linguistics and Intelligent Text Processing*, volume 6609 de *Lecture Notes in Computer Science*, pages 277–288. Springer Berlin Heidelberg, Berlin, Heidelberg.
- ALI, A. L., SCHMID, F., AL-SALMAN, R. et KAUPPINEN, T. (2014). Ambiguity and plausibility : Managing classification quality in volunteered geographic information. In *Proceedings of the 22nd ACM SIGSPATIAL International Conference on Advances in Geographic Information Systems - SIGSPATIAL '14*, pages 143–152, Dallas, Texas. ACM Press.
- ANDERSON, J., SODEN, R., KEEGAN, B., PALEN, L. et ANDERSON, K. M. (2018). The crowd is the territory : Assessing quality in peer-produced spatial data during disasters. *International Journal of Human–Computer Interaction*, 34(4):295 – 310.
- ANTONIOU, V. et SKOPELITI, A. (2015). Measures and Indicators of VGI Quality : An Overview. *ISPRS Annals of Photogrammetry, Remote Sensing and Spatial Information Sciences*, II-3/W5:345–351.
- ANTONIOU, V. et SKOPELITI, A. (2017). The Impact of the Contribution Micro-environment on Data Quality : The Case of OSM. In FOODY, G., SEE, L., FRITZ, S., MOONEY, P., OLTEANU-RAIMOND, A.-M., COSTA FONTE, C. et ANTONIOU, V., éditeurs : *Mapping and the Citizen Sensor*.
- AUDEBERT, N., SAUX, B. L. et LEFEVRE, S. (2017). Joint Learning from Earth Observation and OpenStreetMap Data to Get Faster Better Semantic Maps. In *2017 IEEE Conference on Computer Vision and Pattern Recognition Workshops (CVPRW)*, pages 1552–1560, Honolulu, HI, USA. IEEE.
- BALLATORE, A. (2014). Defacing the Map : Cartographic Vandalism in the Digital Commons. *The Cartographic Journal*, 51(3):214–224.
- BALLATORE, A. et ZIPF, A. (2015). A Conceptual Quality Framework for Volunteered Geographic Information. In FABRIKANT, S. I., RAUBAL, M., BERTOLOTTO, M., DAVIES, C., FREUNDSCHUH, S. et BELL, S., éditeurs : *Spatial Information Theory*, volume 9368 de *Lecture Notes in Computer Science*, pages 89–107. Springer International Publishing.
- BARRINGTON-LEIGH, C. et MILLARD-BALL, A. (2017). The world’s user-generated road map is more than 80% complete. *PLOS ONE*, 12(8):e0180698.

- BARRON, C., NEIS, P. et ZIPF, A. (2014). A Comprehensive Framework for Intrinsic OpenStreetMap Quality Analysis. *Transactions in GIS*, 18(6):877–895.
- BATTISTON, F., NICOSIA, V. et LATORA, V. (2014). Structural measures for multiplex networks. *Phys. Rev. E*, 89:032804.
- BATTON-HUBERT, M. et PINET, F. (2019). Formalisms and Representations of Imperfect Geographic Objects. In BATTON-HUBERT, M., DESJARDIN, E. et PINET, F., éditeurs : *Geographic Data Imperfection 1*, pages 73–103. Wiley, 1 édition.
- BÉGIN, D., DEVILLERS, R. et ROCHE, S. (2013). Assessing Volunteered Geographic Information (VGI) Quality Based on Contributors' Mapping Behaviours. In *8th International Symposium on Spatial Data Quality*, volume XL-2/W1, pages 149–154.
- BÉGIN, D., DEVILLERS, R. et ROCHE, S. (2016). The Life Cycle of Volunteered Geographic Information (VGI) Contributors : The OpenStreetMap Example. In MILLER, J., O'SULLIVAN, D. et WIEGAND, N., éditeurs : *Short Paper Proceedings of GIScience 2016*, pages 9–12, Montreal, Canada.
- BÉGIN, D., DEVILLERS, R. et ROCHE, S. (2018). The life cycle of contributors in collaborative online communities -the case of OpenStreetMap. *International Journal of Geographical Information Science*, 32(8):1611–1630.
- BERLINGERIO, M., COSCIA, M. et GIANNOTTI, F. (2011). Finding and Characterizing Communities in Multidimensional Networks. In *Advances in Social Networks Analysis and Mining (ASONAM), 2011 International Conference On*, pages 490–494. IEEE.
- BISCHKE, B., HELBER, P., FOLZ, J., BORTH, D. et DENGEL, A. (2017). Multi-Task Learning for Segmentation of Building Footprints with Deep Neural Networks. *arXiv :1709.05932 [cs]*.
- BISHR, M. et JANOWICZ, K. (2010). Can we Trust Information? - The Case of Volunteered Geographic Information. In *Towards Digital Earth : Search, Discover and Share Geospatial Data, Workshop at Future Internet Symposium*, Berlin, Germany.
- BLONDEL, V. D., GUILLAUME, J.-L., LAMBIOTTE, R. et LEFEBVRE, E. (2008). Fast unfolding of communities in large networks. *Journal of Statistical Mechanics : Theory and Experiment*, 2008(10):P10008+.
- BRANS, J.-P. et MARESCHAL, B. (2005). Promethee Methods. In *Multiple Criteria Decision Analysis : State of the Art Surveys*, volume 78, pages 163–186. Springer-Verlag, New York.
- BREIMAN, L., FRIEDMAN, J. H., OLSHEN, R. A. et STONE, C. J. (1984). Classification and Regression Trees. *Biometrics*, 40(3):874.
- BRÓDKA, P., CHMIEL, A., MAGNANI, M. et RAGOZINI, G. (2018). Quantifying layer similarity in multiplex networks : A systematic study. *Royal Society Open Science*, 5(8):171747.

- BROVELLI, M. A., MINGHINI, M., MOLINARI, M. et MOONEY, P. (2017). Towards an Automated Comparison of OpenStreetMap with Authoritative Road Datasets. *Transactions in GIS*, 21(2):191–206.
- BROVELLI, M. A. et ZAMBONI, G. (2018). A New Method for the Assessment of Spatial Accuracy and Completeness of OpenStreetMap Building Footprints. *ISPRS International Journal of Geo-Information*, 7(8):289.
- BUCHER, B., BRASEBIN, M., BUARD, E., GROSSO, E., MUSTIÈRE, S. et PERRET, J. (2012). GeOxygene : Built on Top of the Expertise of the French NMA to Host and Share Advanced GI Science Research Results. In BOCHER, E. et NETELER, M., éditeurs : *Geospatial Free and Open Source Software in the 21st Century*, Lecture Notes in Geoinformation and Cartography, pages 21–33. Springer Berlin Heidelberg, Berlin, Heidelberg.
- BUDHATHOKI, N. R., NEDOVIC-BUDIC, Z. et BRUCE, B. (2010). An Interdisciplinary Frame for Understanding Volunteered Geographic Information. *Geomatica*, 64(1):11–26.
- CASTRO, R., TIERRA, A. et LUNA, M. (2019). Assessing the Horizontal Positional Accuracy in OpenStreetMap : A Big Data Approach. In ROCHA, Á., ADELI, H., REIS, L. P. et COSTANZO, S., éditeurs : *New Knowledge in Information Systems and Technologies*, pages 513–523, Cham. Springer International Publishing.
- CHANDOLA, V., BANERJEE, A. et KUMAR, V. (2009). Anomaly detection : A survey. *ACM Computing Surveys*, 41(3):1–58.
- CHEHREGHAN, A. et ALI ABBASPOUR, R. (2018). An evaluation of data completeness of VGI through geometric similarity assessment. *International Journal of Image and Data Fusion*, 9(4):319–337.
- CHEN, J. et ZIPF, A. (2017). Deep Learning with Satellite Images and Volunteered Geographic Information. In *Geospatial Data Science Techniques and Applications*, volume 1, page 274.
- CHIN, S. C., STREET, W. N., SRINIVASAN, P. et EICHMANN, D. (2010). Detecting Wikipedia vandalism with active learning and statistical language models. In *Proceedings of the 4th Workshop on Information Credibility, WICOW '10*, pages 3–10, New York, NY, USA. ACM.
- CLARAMUNT, C. et THÉRIAULT, M. (1996). Toward semantics for modelling spatio-temporal processes within GIS. *Advances in GIS Research I*, pages 27–43.
- CLAUSET, A., NEWMAN, M. et MOORE, C. (2004). Finding community structure in very large networks. *Physical Review E*, 70(6):066111+.
- COHEN, S. (1973). Property Destruction : Motives and Meanings. *Architectural Press*, pages 23–53.
- COLEMAN, D. J., GEOGIADOU, Y. et LABONTE, J. (2010). Volunteered Geographic Information : The Nature and Motivation of Producers. *International Journal of Spatial Data Infrastructures Research*, 4:332 – 358.

- COMBER, A., MOONEY, P., PURVES, R. S., ROCCHINI, D. et WALZ, A. (2016). Crowdsourcing : It Matters Who the Crowd Are. The Impacts of between Group Variations in Recording Land Cover. *PLOS ONE*, 11(7):e0158329+.
- COSTES, B. (2016). *Vers la construction d'un référentiel géographique ancien : un modèle de graphe agrégé pour intégrer, qualifier et analyser des réseaux géohistoriques*. Thèse de doctorat, Paris-Est.
- CRANDALL, D. J., BACKSTROM, L., COSLEY, D., SURI, S., HUTTENLOCHER, D. et KLEINBERG, J. (2010). Inferring social ties from geographic coincidences. *Proceedings of the National Academy of Sciences*, 107(52):22436–22441.
- CRETIAZ DE ROTEN, F. et HELBLING, J.-M. (1996). Données manquantes et aberrantes : le quotidien du statisticien analyste de données. *Revue de Statistique Appliquée*, 44(2):12.
- D'ANTONIO, F., FOGLIARONI, P. et KAUPPINEN, T. (2014). VGI Edit History Reveals Data Trustworthiness and User Reputation. In *Proceedings of the 17th AGILE Conference on Geographic Information Science*, Castellon, Spain. AGILE.
- DEGROSSI, L. C., de ALBUQUERQUE, J. P., ROCHA, R. d. S. et ZIPF, A. (2018). A taxonomy of quality assessment methods for volunteered and crowdsourced geographic information. *Transactions in GIS*, 22(2):542–560.
- DORN, H., TÖRNROS, T. et ZIPF, A. (2015). Quality Evaluation of VGI Using Authoritative Data—A Comparison with Land Use Data in Southern Germany. *ISPRS International Journal of Geo-Information*, 4(3):1657–1671.
- DUFÉAL, M., JONCHÈRES, C. et NOUCHER, M. (2016). ECCE CARTO - DES ESPACES DE LA CONTRIBUTION A LA CONTRIBUTION SUR L'ESPACE - Profils, pratiques et valeurs d'engagement des contributeurs d'OpenStreetMap (OSM). Rapport technique, UMR 5319.
- DUMÉNIU, B. (2015). *Un système d'information géographique pour le suivi d'objets historiques urbains à travers l'espace et le temps*. Thèse de doctorat, Paris-Est.
- ESTER, M., KRIEGEL, H.-p., JORG, S. et XU, X. (1996). A density-based algorithm for discovering clusters in large spatial databases with noise. In *Proceedings of 2nd International Conference on KDD*, pages 226–231.
- FAN, H., ZIPF, A. et FU, Q. (2014a). Estimation of Building Types on OpenStreetMap Based on Urban Morphology Analysis. pages 19–35. Springer International Publishing, Cham.
- FAN, H., ZIPF, A., FU, Q. et NEIS, P. (2014b). Quality assessment for building footprints data on OpenStreetMap. *International Journal of Geographical Information Science*, 28(4):700–719.
- FISCHER, F. (2012). VGI as Big Data : A New but Delicate Geographic Data-Source. *GeoInformatics*, pages 46–47.
- FISKE, S. T., CUDDY, A. J. C. et GLICK, P. (2007). Universal dimensions of social cognition : Warmth and competence. *Trends in Cognitive Sciences*, 11(2):77 – 83.

- FLANAGIN, A. et METZGER, M. (2008). The credibility of volunteered geographic information. *GeoJournal*, 72(3-4):137–148.
- FOGG, B. J. et TSENG, H. (1999). The elements of computer credibility. In *Proceedings of the SIGCHI Conference on Human Factors in Computing Systems the CHI Is the Limit - CHI '99*, pages 80–87, Pittsburgh, Pennsylvania, United States. ACM Press.
- FOGLIARONI, P., D'ANTONIO, F. et CLEMENTINI, E. (2018). Data trustworthiness and user reputation as indicators of VGI quality. *Geo-spatial Information Science*, 21(3):213–233.
- FONTE, C. C., ANTONIOU, V., BASTIN, L., ESTIMA, J., ARSANJANI, J. J., LASO BAYAS, J.-C., SEE, L. et VATSEVA, R. (2017). Assessing VGI Data Quality. In FOODY, G., SEE, L., FRITZ, S., MOONEY, P., OLTEANU-RAIMOND, A.-M., COSTA FONTE, C. et ANTONIOU, V., éditeurs : *Mapping and the Citizen Sensor*, pages 137–163. Ubiquity Press Ltd.
- FONTE, C. C., BASTIN, L., SEE, L., FOODY, G. et LUPIA, F. (2015). Usability of VGI for validation of land cover maps. *International Journal of Geographical Information Science*, 29(7):1269–1291.
- FORATI, A. M. et KARIMIPOUR, F. (2016). A VGI Quality Assessment Method for VGI based on Trustworthiness. *GI_Forum*, 1:3–11.
- GEIGER, R. S. et HALFAKER, A. (2013). Using edit sessions to measure participation in Wikipedia. In *Proceedings of the 2013 Conference on Computer Supported Cooperative Work - CSCW '13*, CSCW '13, pages 861–870, New York, NY, USA. ACM Press.
- GIRRES, J.-F. (2012). *Modèle d'estimation de l'imprécision des mesures géométriques de données géographiques : Application aux mesures de longueur et de surface*. Thèse de doctorat, Paris-Est.
- GIRRES, J.-F. et TOUYA, G. (2010). Quality Assessment of the French OpenStreet-Map Dataset. *Transactions in GIS*, 14(4):435–459.
- GLASZE, G. et PERKINS, C. (2015). Social and Political Dimensions of the OpenStreetMap Project : Towards a Critical Geographical Research Agenda. In JO-KAR ARSANJANI, J., ZIPF, A., MOONEY, P. et HELBICH, M., éditeurs : *OpenStreetMap in GIScience*, pages 143–166. Springer International Publishing, Cham.
- GOLBECK, J. A. (2005). *Computing and Applying Trust in Web-Based Social Networks*. PhD Thesis, University of Maryland.
- GOODCHILD, M. F. (2007). Citizens as Sensors : The World of Volunteered Geography. *GeoJournal*, 69(4):211–221.
- GOODCHILD, M. F. et LI, L. (2012). Assuring the quality of volunteered geographic information. *Spatial Statistics*.
- GRECO, G. M. et FLORIDI, L. (2004). The tragedy of the digital commons. *Ethics and Information Technology*, 6(2):73–81.

- GRÉGOIRE, H. (1794). Rapport sur les destructions opérées par le Vandalisme, et sur les moyens de le réprimer. Rapport technique.
- GRÖCHENIG, S., BRUNAUER, R. et REHRL, K. (2014). Digging into the history of VGI data-sets : Results from a worldwide study on OpenStreetMap mapping activity. *Journal of Location Based Services*, 8(3):198–210.
- HAKLAY, M. (2010). How Good is Volunteered Geographical Information ? A Comparative Study of OpenStreetMap and Ordnance Survey Datasets. *Environment and Planning B : Planning and Design*, 37(4):682–703.
- HARDIN, G. (1968). The Tragedy of the Commons. *Science*, 162(3859):1243–1248.
- HEABERLIN, B. et DEDEO, S. (2016). The Evolution of Wikipedia’s Norm Network. *Future Internet*, 8(2):14+.
- HECHT, R., KUNZE, C. et HAHMANN, S. (2013). Measuring Completeness of Building Footprints in OpenStreetMap over Space and Time. *ISPRS International Journal of Geo-Information*, 2(4):1066–1091.
- HEINDORF, S., POTTHAST, M., STEIN, B. et ENGELS, G. (2015). Towards Vandalism Detection in Knowledge Bases. In *Proceedings of the 38th International ACM SIGIR Conference on Research and Development in Information Retrieval - SIGIR ’15*, pages 831–834. ACM Press.
- HEINDORF, S., POTTHAST, M., STEIN, B. et ENGELS, G. (2016). Vandalism Detection in Wikidata. In *Proceedings of the 25th ACM International on Conference on Information and Knowledge Management - CIKM ’16*, pages 327–336. ACM Press.
- HESS, C. et OSTROM, E. (2007). *Understanding Knowledge As a Commons : From Theory to Practice*, volume 59. The MIT Press.
- HMIMIDA, M. et KANAWATI, R. (2015). Community detection in multiplex networks : A seed-centric approach. *Networks and Heterogeneous Media*, 10(1):71–85.
- HUANG, K. L., KANHERE, S. S. et HU, W. (2010). Are You Contributing Trustworthy Data ? The Case for a Reputation System in Participatory Sensing. pages 14 – 22.
- HUNG, K.-C., KALANTARI, M. et RAJABIFARD, A. (2016). Methods for assessing the credibility of volunteered geographic information in flood response : A case study in Brisbane, Australia. *Applied Geography*, 68:37–47.
- IDDIANOZIE, C. et MCARDLE, G. (2019). A Transfer Learning Paradigm for Spatial Networks. In *Proceedings of the 34th ACM/SIGAPP Symposium on Applied Computing, SAC ’19*, pages 659–666, New York, NY, USA. ACM.
- IKEDA, K., MORISHIMA, A., RAHMAN, H., ROY, S. B., THIRUMURUGANATHAN, S., YAHIA, S. A. et DAS, G. (2016). Collaborative Crowdsourcing with Crowd4U. *Proc. VLDB Endow.*, 9(13):1497–1500.

- ISO19157 (2003). Information géographique — Qualité des données. Rapport technique, International Organization for Standardization.
- ISSAM, F., MANEL, H. et RUSHED, K. (2015). Une approche centrée graine pour la détection de communautés dans les réseaux multiplexes. *In EGC*, pages 377–382.
- IVANOVIC, S. (2018). *Une Approche Basée Sur La Qualité Pour Mettre à Jour Les Bases de Données Géographiques de Référence à Partir de Traces GPS Issues de La Foule*. PhD Thesis, Université Paris Est.
- IVANOVIC, S. S., OLTEANU-RAIMOND, A.-M., MUSTIÈRE, S. et DEVOGELE, T. (2020). Potential of Crowdsourced Traces for Detecting Updates in Authoritative Geographic Data. *In KYRIAKIDIS, P., HADJIMITSIS, D., SKARLATOS, D. et MANSOURIAN, A., éditeurs : Geospatial Technologies for Local and Regional Development*, Lecture Notes in Geoinformation and Cartography, pages 205–221. Springer International Publishing.
- JAARA, K., DUCHÊNE, C. et RUAS, A. (2014). Preservation and modification of relations between thematic and topographic data throughout thematic data migration process. *In BUCHROITHNER, M., PRECHTEL, N. et BURGHARDT, D., éditeurs : Cartography from Pole to Pole : Selected Contributions to the XXVIth International Conference of the ICA, Dresden 2013*, pages 103–117. Springer Berlin Heidelberg, Berlin, Heidelberg.
- JOKAR ARSANJANI, J., MOONEY, P., ZIPF, A. et SCHAUSS, A. (2015). Quality Assessment of the Contributed Land Use Information from OpenStreetMap Versus Authoritative Datasets. *In JOKAR ARSANJANI, J., ZIPF, A., MOONEY, P. et HELBICH, M., éditeurs : OpenStreetMap in GIScience*, pages 37–58. Springer International Publishing, Cham.
- JONIETZ, D. et ZIPF, A. (2016). Defining Fitness-for-Use for Crowdsourced Points of Interest (POI). *ISPRS International Journal of Geo-Information*, 5(9):149+.
- JUHÁSZ, L., HOCHMAIR, H., QIAO, S. et NOVACK, T. (2019). Exploring the Effects of Pokémon Go Vandalism on OpenStreetMap. *In Proceedings of the Academic Track*, page 4, Heildelberg, Germany.
- KAMVAR, S. D., SCHLOSSER, M. T. et MOLINA, H. G. (2003). The Eigentrust Algorithm for Reputation Management in P2P Networks. *In Proceedings of the 12th International Conference on World Wide Web, WWW '03*, pages 640–651, New York, NY, USA. ACM.
- KAZIENKO, P., MUSIAL, K. et KAJDANOWICZ, T. (2011). Multidimensional Social Network in the Social Recommender System. *IEEE Transactions on Systems, Man, and Cybernetics - Part A : Systems and Humans*, 41(4):746–759.
- KESSLER, C. et DE GROOT, R. T. A. (2013). Trust as a Proxy Measure for the Quality of Volunteered Geographic Information in the Case of OpenStreetMap. *In VANDENBROUCKE, D., BUCHER, B. et CROMPVOETS, J., éditeurs : Geographic Information Science at the Heart of Europe*, Lecture Notes in Geoinformation and Cartography, pages 21–37. Springer International Publishing.

- KITTUR, A., SUH, B., PENDLETON, B. A. et CHI, E. H. (2007). He says, she says : Conflict and coordination in Wikipedia. *In Proceedings of the SIGCHI Conference on Human Factors in Computing Systems*, CHI '07, pages 453–462, New York, NY, USA. ACM.
- KIVELÄ, M., ARENAS, A., BARTHELEMY, M., GLEESON, J. P., MORENO, Y. et PORTER, M. A. (2014). Multilayer networks. *Journal of Complex Networks*, 2(3):203–271.
- LE CROSNIER, H. (2011). Une bonne nouvelle pour la théorie des biens communs. *Vacarme*, 56(3):92.
- LECUN, Y., BENGIO, Y. et HINTON, G. (2015). Deep learning. *Nature*, 521:436.
- LECUN, Y., BOTTOU, L., BENGIO, Y. et HA, P. (1998). Gradient-Based Learning Applied to Document Recognition. *In Proceedings of the IEEE*.
- LESOT, M.-J. et REVAULT D'ALLONNES, A. (2017). Information quality and uncertainty. page 14.
- LESOT, M.-J., RIFQI, M. et BOUCHON-MEUNIER, B. (2008). Fuzzy Prototypes : From a Cognitive View to a Machine Learning Principle. *In BUSTINCE, H., HERRERA, F. et MONTERO, J., éditeurs : Fuzzy Sets and Their Extensions : Representation, Aggregation and Models*, volume 220, pages 431–452. Springer Berlin Heidelberg, Berlin, Heidelberg.
- LODIGIANI, C. et MELCHIORI, M. (2016). A PageRank-based Reputation Model for VGI Data. *Procedia Computer Science*, 98:566–571.
- MA, D., SANDBERG, M. et JIANG, B. (2015). Characterizing the Heterogeneity of the OpenStreetMap Data and Community. *ISPRS International Journal of Geo-Information*, 4(2):535–550.
- MA, L., LIU, Y., ZHANG, X., YE, Y., YIN, G. et JOHNSON, B. A. (2019). Deep learning in remote sensing applications : A meta-analysis and review. *ISPRS Journal of Photogrammetry and Remote Sensing*, 152:166 – 177.
- MAGGIORI, E., TARABALKA, Y., CHARPIAT, G. et ALLIEZ, P. (2017). Can semantic labeling methods generalize to any city ? the inria aerial image labeling benchmark. *In 2017 IEEE International Geoscience and Remote Sensing Symposium (IGARSS)*, pages 3226–3229, Fort Worth, TX. IEEE.
- MAGUIRE, S. et TOMKO, M. (2017). Ripe for the picking ? Dataset maturity assessment based on temporal dynamics of feature definitions. *International Journal of Geographical Information Science*, 31(7):1334–1358.
- MAHABIR, R., STEFANIDIS, A., CROITORU, A., CROOKS, A. T. et AGOURIS, P. (2017). Authoritative and Volunteered Geographical Information in a Developing Country : A Comparative Case Study of Road Datasets in Nairobi, Kenya. *ISPRS International Journal of Geo-Information*, 6(1):24.

- MAJIC, I., WINTER, S. et TOMKO, M. (2017). Finding equivalent keys in OpenStreetMap : Semantic similarity computation based on extensional definitions. *In 1st Workshop on Artificial Intelligence and Deep Learning for Geographic Knowledge Discovery*, pages 24–32. ACM.
- MARMANIS, D., SCHINDLER, K., WEGNER, J. D., GALLIANI, S., DATCU, M. et STILLA, U. (2016). Classification With an Edge : Improving Semantic Image Segmentation with Boundary Detection. *arXiv :1612.01337 [cs]*.
- MARTINEZ-RICO, J. R., MARTINEZ-ROMO, J. et ARAUJO, L. (2019). Can deep learning techniques improve classification performance of vandalism detection in Wikipedia? *Engineering Applications of Artificial Intelligence*, 78:248–259.
- MCKENZIE, G. et JANOWICZ, K. (2014). Coerced Geographic Information : The Not-so-voluntary Side of User-generated Geo-content. *In GIScience'2014*, page 3, Vienne, Autriche.
- MCNAIR, H. et ARNOLD, L. (2016). Crowd-sorting : Reducing bias in decision making through consensus generated crowdsourced spatial information. *International Conference on GIScience Short Paper Proceedings*, 1(1).
- MOLA VELASCO, S. M. (2011). Wikipedia vandalism detection. *In Proceedings of the 20th International Conference Companion on World Wide Web, WWW '11*, pages 391–396, New York, NY, USA. ACM.
- MOONEY, P. et CORCORAN, P. (2012a). The Annotation Process in OpenStreetMap. *Transactions in GIS*, 16(4):561–579.
- MOONEY, P. et CORCORAN, P. (2012b). How social is OpenStreetMap? *In GENSEL, J., JOSSELIN, D. et VANDENBROUCKE, D., éditeurs : Proceedings of the AGILE'2012 International Conference on Geographic Information Science*.
- MOONEY, P. et CORCORAN, P. (2014). Analysis of Interaction and Co-editing Patterns amongst OpenStreetMap Contributors. *Transactions in GIS*, 18(5):633–659.
- MÜLLIGANN, C., JANOWICZ, K., YE, M. et LEE, W.-C. (2011). Analyzing the Spatial-Semantic Interaction of Points of Interest in Volunteered Geographic Information. volume 6899, pages 350–370, Berlin, Heidelberg. Springer Berlin Heidelberg.
- MUTTAQIEN, B. I., OSTERMANN, F. O. et LEMMENS, R. L. G. (2018). Modeling aggregated expertise of user contributions to assess the credibility of OpenStreetMap features. *Transactions in GIS*, 22(3):823–841.
- NEIS, P., GOETZ, M. et ZIPF, A. (2012). Towards Automatic Vandalism Detection in OpenStreetMap. *ISPRS International Journal of Geo-Information*, 1(3):315–332.
- NEIS, P. et ZIELSTRA, D. (2014). Recent Developments and Future Trends in Volunteered Geographic Information Research : The Case of OpenStreetMap. *Future Internet*, 6(1):76–106.

- NEIS, P. et ZIPF, A. (2012). Analyzing the Contributor Activity of a Volunteered Geographic Information Project — The Case of OpenStreetMap. *ISPRS International Journal of Geo-Information*, 1(2):146–165.
- NIELSEN, J. (2006). The 90-9-1 Rule for Participation Inequality in Social Media and Online Communities.
- OLTEANU-RAIMOND, A.-M., JOLIVET, L., VAN DAMME, M.-D., ROYER, T., FRAVAL, L., SEE, L., STURN, T., KARNER, M., MOORTHY, I. et FRITZ, S. (2018). An Experimental Framework for Integrating Citizen and Community Science into Land Cover, Land Use, and Land Change Detection Processes in a National Mapping Agency. *Land*, 7(3):103.
- OLTEANU-RAIMOND, A.-M., LAAKSO, M., ANTONIOU, V., FONTE, C. C., FONSECA, A., GRUS, M., HARDING, J., KELLENBERGER, T., MINGHINI, M. et SKOPELITI, A. (2017). VGI in National Mapping Agencies : Experiences and Recommendations. In *Mapping and the Citizen Sensor*, pages 299–326. Ubiquity Press.
- POTTHAST, M. (2010). Crowdsourcing a wikipedia vandalism corpus. In *Proceedings of the 33rd International ACM SIGIR Conference on Research and Development in Information Retrieval*, SIGIR '10, pages 789–790, New York, NY, USA. ACM.
- POTTHAST, M., GERLING, R. et STEIN, B. (2007). Webis Wikipedia Vandalism Corpus (Webis-WVC-07).
- REHRL, K., GRÖECHENIG, S., HOCHMAIR, H., LEITINGER, S., STEINMANN, R. et WAGNER, A. (2013). A Conceptual Model for Analyzing Contribution Patterns in the Context of VGI. In KRISP, J. M., éditeur : *Progress in Location-Based Services*, Lecture Notes in Geoinformation and Cartography, pages 373–388. Springer Berlin Heidelberg.
- REVAULT D'ALLONNES, A. (2013). Architecture de l'évolution de la confiance : Définition et influence des dimensions nécessaires à la formation d'une opinion. In *L'évaluation de l'information*.
- ROLLASON, E., BRACKEN, L. J., HARDY, R. J. et LARGE, A. R. G. (2018). The importance of volunteered geographic information for the validation of flood inundation models. *Journal of Hydrology*, 562:267–280.
- SAFAVIAN, S. R. et LANDGREBE, D. (1990). A Survey of Decision Tree Classifier Methodology. page 50.
- SARABADANI, A., HALFAKER, A. et TARABORELLI, D. (2017). Building Automated Vandalism Detection Tools for Wikidata. In *Proceedings of the 26th International Conference on World Wide Web Companion - WWW '17 Companion*, pages 1647–1654. ACM Press.
- SCHAEFFER, S. E. (2007). Graph clustering. *Computer Science Review*, 1(1):27–64.

- SCHMIDT, M. et KLETTNER, S. (2013). Gender and Experience-Related Motivators for Contributing to OpenStreetMap. *In Online Proceedings of the International Workshop on Action and Interaction in Volunteered Geographic Information (ACTIVITY) at the 16th AGILE Conference on Geographic Information Science*. AGILE.
- SCHOLZ, S., KNIGHT, P., ECKLE, M., MARX, S. et ZIPF, A. (2018). Volunteered Geographic Information for Disaster Risk Reduction—The Missing Maps Approach and Its Potential within the Red Cross and Red Crescent Movement. *Remote Sensing*, 10(8):1239.
- SEE, L., MOONEY, P., FOODY, G., BASTIN, L., COMBER, A., ESTIMA, J., FRITZ, S., KERLE, N., JIANG, B., LAAKSO, M., LIU, H.-Y., MILČINSKI, G., NIKŠIČ, M., PAINHO, M., PÓDÖR, A., OLTEANU-RAIMOND, A.-M. et RUTZINGER, M. (2016). Crowdsourcing, Citizen Science or Volunteered Geographic Information? The Current State of Crowdsourced Geographic Information. *ISPRS International Journal of Geo-Information*, 5(5):55.
- SENARATNE, H., MOBASHERI, A., ALI, A. L., CAPINERI, C. et HAKLAY, M. M. (2017). A review of volunteered geographic information quality assessment methods. *International Journal of Geographical Information Science*, 31(1):139–167.
- SERVIGNE, S., LESAGE, N. et LIBOUREL, T. (2005). Composantes qualité et métadonnées. *In Qualité de l'information Géographique*, pages 213–245. Lavoisier édition.
- SKARLATIDOU, A., HAKLAY, M. et CHENG, T. (2011). Trust in Web GIS : The role of the trustee attributes in the design of trustworthy Web GIS applications. *International Journal of Geographical Information Science*, 25(12):1913–1930.
- SRIVASTAVA, S., VARGAS MUÑOZ, J. E., LOBRY, S. et TUIA, D. (2018). Fine-grained landuse characterization using ground-based pictures : A deep learning solution based on globally available data. *International Journal of Geographical Information Science*, pages 1–20.
- STEIN, K. et BLASCHKE, S. (2010). Interlocking Communication. *In MEMON, N. et ALHAJJ, R., éditeurs : From Sociology to Computing in Social Networks*, pages 231–252. Springer Vienna.
- STEIN, K., KREMER, D. et SCHLIEDER, C. (2015). Spatial Collaboration Networks of OpenStreetMap. *In JOKAR ARSANJANI, J., ZIPF, A., MOONEY, P. et HELBICH, M., éditeurs : OpenStreetMap in GIScience*, Lecture Notes in Geoinformation and Cartography, pages 167–186. Springer International Publishing.
- SULER, J. (2004). The Online Disinhibition Effect. *Cyberpsychology & Behavior*, 7(3):7.
- SUTTON, M. (1987). *Differential Rates of Vandalism in a New Town : Towards a Theory of Relative Place*. Thèse de doctorat, Lancashire Polytechnic.
- TOBLER, W. R. (1970). A Computer Movie Simulating Urban Growth in the Detroit Region. *Economic Geography*, 46:234.

- TOUYA, G., ANTONIOU, V., CHRISTOPHE, S. et SKOPELITI, A. (2017a). Production of Topographic Maps with VGI : Quality Management and Automation. *In Mapping and the Citizen Sensor*, pages 61–91. Ubiquity Press.
- TOUYA, G., ANTONIOU, V., OLTEANU-RAIMOND, A.-M. et VAN DAMME, M.-D. (2017b). Assessing Crowdsourced POI Quality : Combining Methods Based on Reference Data, History, and Spatial Relations. *ISPRS International Journal of Geo-Information*, 6(3):80+.
- TOUYA, G. et BRANDO-ESCOBAR, C. (2013). Detecting Level-of-Detail Inconsistencies in Volunteered Geographic Information Data Sets. *Cartographica : The International Journal for Geographic Information and Geovisualization*, 48(2):134–143.
- TOUYA, G., ZHANG, X. et LOKHAT, I. (2019). Is deep learning the new agent for map generalization ? *International Journal of Cartography*, pages 1–16.
- TRUONG, Q. T., DE RUNZ, C. et TOUYA, G. (2018a). Analyse du comportement des contributeurs dans l’Information Géographique Volontaire via la construction de réseaux sociaux. *In DE RUNZ, C., KERGOSIEN, É., GUYET, T. et SALLABERRY, C., éditeurs : 18ème Conférence Internationale Sur l’Extraction et La Gestion Des Connaissances (EGC 2018)*, pages 44–54, Paris, France.
- TRUONG, Q.-T., DE RUNZ, C. et TOUYA, G. (2019). Analysis of collaboration networks in OpenStreetMap through weighted social multigraph mining. *International Journal of Geographical Information Science*, 33(8):1651–1682.
- TRUONG, Q.-T., TOUYA, G. et DE RUNZ, C. (2018b). Building Social Networks in Volunteered Geographic Information Communities : What Contributor Behaviours Reveal About Crowdsourced Data Quality. *In FOGLIARONI, P., BALLATORE, A. et CLEMENTINI, E., éditeurs : Proceedings of Workshops and Posters at the 13th International Conference on Spatial Information Theory (COSIT 2017)*, Lecture Notes in Geoinformation and Cartography, pages 125–131. Springer International Publishing.
- TRUONG, Q. T., TOUYA, G. et DE RUNZ, C. (2018c). Le vandalisme dans l’information géographique volontaire : apprendre pour mieux détecter ? *In Actes de la conférence SAGEO 2018*, pages 61–76, Montpellier, France.
- TRUONG, Q.-T., TOUYA, G. et DE RUNZ, C. (2018d). Towards Vandalism Detection in OpenStreetMap Through a Data Driven Approach (Short Paper). *In WINTER, S., GRIFFIN, A. et SESTER, M., éditeurs : 10th International Conference on Geographic Information Science (GIScience 2018)*, volume 114 de *Leibniz International Proceedings in Informatics (LIPIcs)*, Dagstuhl, Germany. Schloss Dagstuhl–Leibniz-Zentrum fuer Informatik.
- TUBARO, P. et CASILLI, A. A. (2019). Micro-work, artificial intelligence and the automotive industry. *Journal of Industrial and Business Economics*, 46(3):333–345.
- VAN VLIET, W. (1984). Vandalism : An assessment and agenda. *In Vandalism, Behaviour and Motivations*. Claude Lévy-Leboyer.

- VERLEYSEN, M. et FRANÇOIS, D. (2005). The curse of dimensionality in data mining and time series prediction. In CABESTANY, J., PRIETO, A. et SANDOVAL, F., éditeurs : *Computational Intelligence and Bioinspired Systems*, pages 758–770, Berlin, Heidelberg. Springer Berlin Heidelberg.
- WALTER, V. et FRITSCH, D. (1999). Matching spatial data sets : A statistical approach. *International Journal of Geographical Information Science*, 13(5):445–473.
- WATTS, D. J. et STROGATZ, S. H. (1998). Collective dynamics of 'small-world' networks. *Nature*, 393(6684):440–442.
- XU, Y., CHEN, Z., XIE, Z. et WU, L. (2017). Quality assessment of building footprint data using a deep autoencoder network. *International Journal of Geographical Information Science*, 31(10):1929–1951.
- YANG, A., FAN, H. et JING, N. (2016). Amateur or Professional : Assessing the Expertise of Major Contributors in OpenStreetMap Based on Contributing Behaviors. *ISPRS International Journal of Geo-Information*, 5(2):21+.
- YANG, S., SHEN, J., KONEČNÝ, M., WANG, Y. et ŠTAMPACH, R. (2018). Study on the Spatial Heterogeneity of the POI Quality in OpenStreetMap. In *Proceedings*, Sozopol, Bulgaria. Bandrova T., Konečný M.
- ZHANG, H. et MALCZEWSKI, J. (2018). Accuracy Evaluation of the Canadian OpenStreetMap Road Networks. *International Journal of Geospatial and Environmental Research*, 5(2):1 – 14.
- ZHANG, L. et PFOSER, D. (2019). Using OpenStreetMap point-of-interest data to model urban change—A feasibility study. *PLOS ONE*, 14(2):e0212606.
- ZHU, X. X., TUIA, D., MOU, L., XIA, G., ZHANG, L., XU, F. et FRAUNDORFER, F. (2017). Deep Learning in Remote Sensing : A Comprehensive Review and List of Resources. *IEEE Geoscience and Remote Sensing Magazine*, 5(4):8–36.
- ZIELSTRA, D., HOCHMAIR, H., NEIS, P. et TONINI, F. (2014). Areal Delineation of Home Regions from Contribution and Editing Patterns in OpenStreetMap. *ISPRS International Journal of Geo-Information*, 3(4):1211–1233.
- ZIELSTRA, D., HOCHMAIR, H. H. et NEIS, P. (2013). Assessing the Effect of Data Imports on the Completeness of OpenStreetMap - A United States Case Study. *Transactions in GIS*, 17(3):315–334.
- ZIMBARDO, P. G. (1971). A Social-Psychological Analysis of Vandalism : Making Sense of Senseless Violence. Rapport technique, Stanford University, Department of Psychology.

Annexe A

Exploration des utilisateurs bannis du projet OSM

ID changeset	Zone géographique	Opérations effectuées	Nombre total de blocks	En quoi cela dégrade la base de données	Vandalisme ?
54637913	Saragosse, Espagne	Suppression d'une zone résidentielle "Área residencial de Zaragoza" (relation 1841510)	2 dont 1 actif	Il manque des informations attributaires essentielles qui permettent de décrire l'objet	Oui
54858210	Saragosse, Espagne	Modification de la zone résidentielle "Área residencial de Zaragoza" (relation 1841510) et suppression de nodes	3 dont 1 actif	Il manque des informations attributaires essentielles qui permettent de décrire l'objet	Oui
54813961	Dar es Salaam, Tanzanie	Création d'un bâtiment (node 5297502521) ayant des valeurs de tags aléatoires	1	Les informations attributaires sur l'objet n'ont aucun sens	Oui
54366258	Milville, Delaware, USA	Création d'un park nommé "Bay Forest" (way 544774374) qui est en réalité une zone résidentielle. Le contributeur est un joueur de Pokémon Go qui cherche à favoriser le développement de Pokémon dans cette zone. Il y a eu un edit war sur cet objet (donc il existe plusieurs changesets pour cet objet)	3 dont 1 actif	L'objet géographique ne décrit plus la même réalité	Oui
54511454	Milville, Delaware, USA	Création d'une piscine fictive (way 545841767). Le contributeur est un joueur de Pokémon Go qui cherche à favoriser le développement de Pokémon dans cette zone		L'objet géographique n'existe pas dans la réalité	Oui
54379511	Bethany Beach, Delaware, USA	Changement d'une zone militaire "Delaware National Guard Bethany Beach Training Site" en zone industrielle (relation 6385988). Pokemon Go Edit war		L'objet géographique ne décrit plus la même réalité	Oui
46358006	Valais, Suisse	N'a pas répondu aux commentaires que les autres contributeurs lui ont adressés sur ses changesets: http://resultmaps.neis-one.org/osm-discussion-comments?uid=5088068 . La plupart de ses actions ont été annulées	8	Il manque des informations attributaires essentielles qui permettent de décrire l'objet	Oui
51984229	Einsiedeln, Suisse	Dans ce changeset, création d'un objet de type polygone (way 524208725) qui n'est pas décrit par des tags, et d'une piste de ski qui a été supprimée après coup. On note que le changeset ne possède pas que des objets vandalisés par exemple l'utilisateur a modifié la 13ème version du way 30000144. Le contributeur a surtout été bloqué car il ne répondait à aucun commentaire adressé sur ses changesets	10	Il manque des informations attributaires essentielles qui permettent de décrire l'objet	Oui
53255494	Victoria Island, Canada	Création d'énormes bâtiments (comme way 535485873 ou way 535485874) qui n'existent pas puisqu'ils chevauchent plusieurs îles du Nord du Canada. Remarque : Est-ce que c'est réellement un cas de vandalisme (i.e. intentionnel) ou d'erreur de débutant ? On voit que le changeset date du jour de l'enregistrement du contributeur	1	L'objet géographique n'existe pas dans la réalité	Non

ID changeset	Zone géographique	Opérations effectuées	Nombre total de blocks	En quoi cela dégrade la base de données	Vandalisme ?
53137491	Herald Island, Russie	Ajout de données imaginaires. Renomme Herald Island (way 138623366) en JerryLand et lui donne le tag leisure=park. D'ailleurs il édite cet objet plusieurs fois successivement avant d'être revert	4	L'objet géographique ne décrit plus la même réalité	Oui
52420815	Nottingham, UK	L'utilisateur a ajouté des tags pour nommer des objets en russe en faisant de la translittération, ce qui est interdit dans OSM. Par exemple, le node 18712963 porte un tag name=Castle Bytham et le contributeur a ajouté un tag name:ru=Касл-Байтем qui est une translittération du nom anglais. L'utilisateur n'a pas répondu quand on lui a demandé la source qu'il a utilisé pour nommer ces objets en russe, d'où le block. Ces tags ont été annulés	1	Ajoute une information qui n'est pas très utile. Mais on ne peut pas vraiment dire que la base de données ait été vandalisée	Non
54868474	Haifa, Israël	Création d'un way 548731776 qui chevauche un chemin piéton déjà existant. Globalement, cet utilisateur a été bloqué car il ne répondait pas aux commentaires	1	Incompatibilité avec les objets géographiques pré-existants	Oui
55705923	Niamey, Niger	Création d'objets dont les tags n'ont pas de signification ("Les tags sur https://www.openstreetmap.org/node/5358438540#map=19/13.49525/2.09384 sont inutiles à OpenStreetMap"). On voit par exemple que le node 5358438213 possède des tag keys tout en majuscules et des tag values principalement chiffrés. <i>Remarque : Est-ce que c'est réellement un cas de vandalisme (i.e. intentionnel) ou d'erreur de débutant ?</i>	1 dont 1 actif	Il manque des informations attributaires essentielles qui permettent de décrire l'objet	Non
55652281	Niamey, Niger	Idem que pour le contributeur précédent (7517131). Les tags utilisés sont du même type. Ils semblent travailler sur un projet Hot OSM. <i>Remarque : Est-ce que c'est réellement un cas de vandalisme (i.e. intentionnel) ou d'erreur de débutant ?</i>	1	Il manque des informations attributaires essentielles qui permettent de décrire l'objet	Non
55838792	Phoenix, Arizona, USA	Modification des tags d'une zone commerciale (way 544402424) en natural=water, water=lake	1	L'objet géographique ne décrit plus la même réalité + incompatibilité avec les objets géographiques existants	Oui
55838543	Phoenix, Arizona, USA	Création de lacs dans une zone urbaine (way 556540255 et way 556540256). Le changeset contient d'autres opérations qui ne sont pas forcément du vandalisme	1	L'objet géographique ne décrit plus la même réalité + incompatibilité avec les objets géographiques existants	Oui
55593777	Puerto Princesa, Philippines	Suppression du tag name="Plaza Cuartel" désignant un parc	1	Il manque des informations attributaires essentielles qui permettent de décrire l'objet	Oui

ID changeset	Zone géographique	Opérations effectuées	Nombre total de blocks	En quoi cela dégrade la base de données	Vandalisme ?
56360014	Palerme, Italie	Ajout de nombreux noms de routes (avec le tag key name="") sans en préciser la source, mais il indique utiliser Bing Imagery. Il ne répond pas aux commentaires d'un contributeur qui lui dit qu'avec Bing Imagery, on ne peut pas avoir le nom des routes.	1 dont 1 actif	Il manque la source des informations entrées (incompatibilité potentielle avec la licence ODBL)	Non
56151423	Leposavić, Kosovo	Le node 448739405 correspond au centre de la ville Leposavić, qui est principalement habitée par des serbes. Donc son orthographe principale correspond à l'écriture serbe, mais il y a aussi un tag name_alt="Leposaviq / Albanik" donnant l'écriture albanaise de la ville. Ici, l'utilisateur a modifié la valeur du tag name principal pour y mettre l'écriture alternative. Il y a eu plusieurs tentatives de sa part pour imposer l'écriture alternative	1	La base de données perd en objectivité ? Ici, le nom serbe reste sur un autre tag, même s'il n'est plus sur le nom principal. Donc la base de données n'est pas vraiment dégradée ici	Non
56151459	Malishevë / Mališevo, Kosovo	Le node 1920392800 possédait sur son tag "name" principal les deux écritures (albanaise / serbe) puisque le Kosovo est habité et par des serbes et par des albanais. Ici, l'utilisateur a retiré l'écriture serbe du tag name pour ne garder que l'écriture albanaise. Cependant, il y a tout de même un tag "name:sr-Latn" où l'écriture serbe est indiquée		La base de données perd en objectivité ? Ici, le nom serbe reste sur un autre tag, même s'il n'est plus sur le nom principal. Donc la base de données n'est pas vraiment dégradée ici	Non
54630223	Columbia, South Carolina, USA	Suppression de nombreux polygones (way 546861742, way 546860888, way 546861739, way 546862239) qui délimitent des zones scolaires. Le contributeur semble être un joueur de Pokémon Go. Par la suite, dans le changeset 54661812 il crée sur la même zone des polygones de type natural=sand, landuse=grass et landuse=wood (cf changeset 54370252)	2	L'objet géographique ne décrit plus la même réalité	Oui
56598546	Denver, USA	Contributeur ayant ajouté massivement des données de type bâtiment sur l'Etat de Denver aux Etats-Unis alors que son changeset indique qu'il travaille pour un projet Hot OSM au Togo	2 dont 1 actif	L'objet géographique n'existe pas dans la réalité + Il manque des informations attributaires essentielles qui permettent de décrire l'objet	Non
56250619	Catania, Italie	Suppression d'un grand nombre de trottoirs (way 553856732, way 559895688) et de passages piétons (way 118656919)	1 dont 1 actif	L'objet géographique n'apparaît plus alors qu'il existe toujours dans la réalité	Oui

ID changeset	Zone géographique	Opérations effectuées	Nombre total de blocks	En quoi cela dégrade la base de données	Vandalisme ?
56253421	Bondeno, Italie	Ajout de tag name=pas un nom propre	2 dont 1 actif	L'objet géographique porte un nom erroné	Oui
50160595	Boulogne-Billancourt, France	Valeur des tags en majuscule (adresse, noms, etc.) et utilisation de key tags non conventionnels (exemple: description1, description2). Ce contributeur semble importer massivement des données issues de la BANO	3	Pour OSM, ce contributeur a été bloqué quelques temps et ses contributions ont été annulées car il ne répondait pas aux commentaires. Cependant, celui-ci continue de contribuer largement et ses données n'étaient pas complètement fausses.	Non
49337814	Lansing, Michigan, USA	Ajout d'un POI (node 4897853364) qui existait déjà dans un bâtiment existant avec un nom de rue incorrect	2 dont 1 actif	Créer des objets redondants (et avec des mauvaises descriptions)	Oui
49281498		Ajout d'un POI (node 4897972537) dans un bâtiment qui portait une fonction différente		Incompatibilité avec les objets géographiques existants	Oui
48186273	Chesterfiel, UK (entre autre)	Des noeuds de route (node 339830304 et node 324706111) deviennent des stations d'essence avec l'ajout de tags type "fuel:lgp=yes"	2	Incompatibilité avec les objets géographiques existants	Oui
48048355	Long Beach, Californie, USA	Suppression du Marina Green Park (node 4792395256, node 4792395284...) de Long Beach à cause de Pokemon Go.	4 dont 1 actif	L'objet géographique n'apparaît plus alors qu'il existe toujours dans la réalité	Oui
47768701	Autour de LA, Californie, USA	Création d'un chemin qui n'existe pas (way 486716027)	1	L'objet géographique n'existe pas dans la réalité	Oui
		Suppression du tag tiger:source sur un chemin (way 31929606) importé de la base TIGER		Il manque des informations attributaires essentielles qui permettent de décrire l'objet	Oui
		Suppression du Marina Green Park (way 251520438)		L'objet géographique n'apparaît plus alors qu'il existe toujours dans la réalité	Oui
47757296	Riga, Lettonie	Supprime le tag landuse=forest et ajoute leisure=park (way 87406131)	1	L'objet géographique ne décrit plus la même réalité	Oui
47850376		Sur le même objet (way 87406131) qui est sujet à un edit war : Laisse le tag landuse=forest et ajoute leisure=park. Après vérification, il y a une ambiguïté sur la nature de ce lieu. D'après un contributeur (du DWG?) c'est officiellement une forêt mais d'après http://www.spottedbylocals.com/riga/anninmuižas-mezs/ c'est un parc.		Je considère ici que ce n'est pas du vandalisme mais que cela met en avant les lieux qui possèdent une ambiguïté sur leur nature	Non
50182925		Sur le même objet (way 87406131) : le lieu était nommé selon le tag name=Anninmuižas mezs et name_alt=Anninmuižas parks. Dans ce changeset, l'utilisateur remplace la valeur de name par celle de name_alt		Il manque des informations attributaires essentielles qui permettent de décrire l'objet	Oui

ID changeset	Zone géographique	Opérations effectuées	Nombre total de blocks	En quoi cela dégrade la base de données	Vandalisme ?
47571821	Moscou, Russie	Suppression du tag name sur un restaurant (node 4218558692)	2	Il manque des informations attributaires essentielles qui permettent de décrire l'objet	Oui
47571994		Suppression de tous les tags name sur un restaurant (node 3325594248)			
47420276	Giardini-Naxos, Messine, Italie	Remplace le nom du bureau de poste par son opérateur (node 2838107107) alors que le tag name avait une valeur qui le distingue des autres bureaux de poste	1	Il manque des informations attributaires essentielles qui permettent de décrire l'objet	Oui
481094786	Près de Sirajganj, Bangladesh	Crée de nombreuses maisons (way 481094817, way 481094825) qui portent le tag name=house	1 dont 1 actif	Ici, ce n'est pas du vandalisme car ça semble être plutôt une erreur de saisie	Non
46946121	Dolgeville, New York, USA	Crée un chemin (way 481139845) au milieu d'un parc qui n'existe pas vu la géométrie improbable	1	L'objet géographique n'existe pas dans la réalité	Oui
46291205	Granite Falls, Whashington, USA	Modifie le nom du barrage (node 356549067) Sweet Dam en Crappy Dam, transforme les rues résidentielles en rues piétonnes (exemple : way 475568045) et supprime une école (exemple: node 3157628478)	2	Il manque des informations attributaires essentielles qui permettent de décrire l'objet + L'objet ne décrit plus la même réalité + l'objet n'apparaît plus alors qu'il existe toujours dans la réalité	Oui

Annexe B

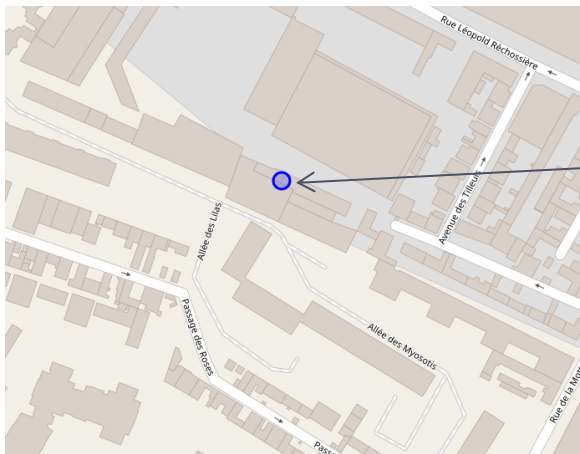
Étude de la fiabilité de 30 contributeurs OSM

Evaluation de la fiabilité des contributeurs OSM

Sur les contributeurs (anonymisés) ayant contribué sur la ville d'Aubervilliers jusqu'au 13/02/2018

Contributeur U1

nombre de contributions	% création	% modification	% suppression	% objets réutilisés par d'autres	% d'objets qui n'ont pas été édités	% d'objets qui n'ont pas été supprimés
1	1	0	0	0	1	0



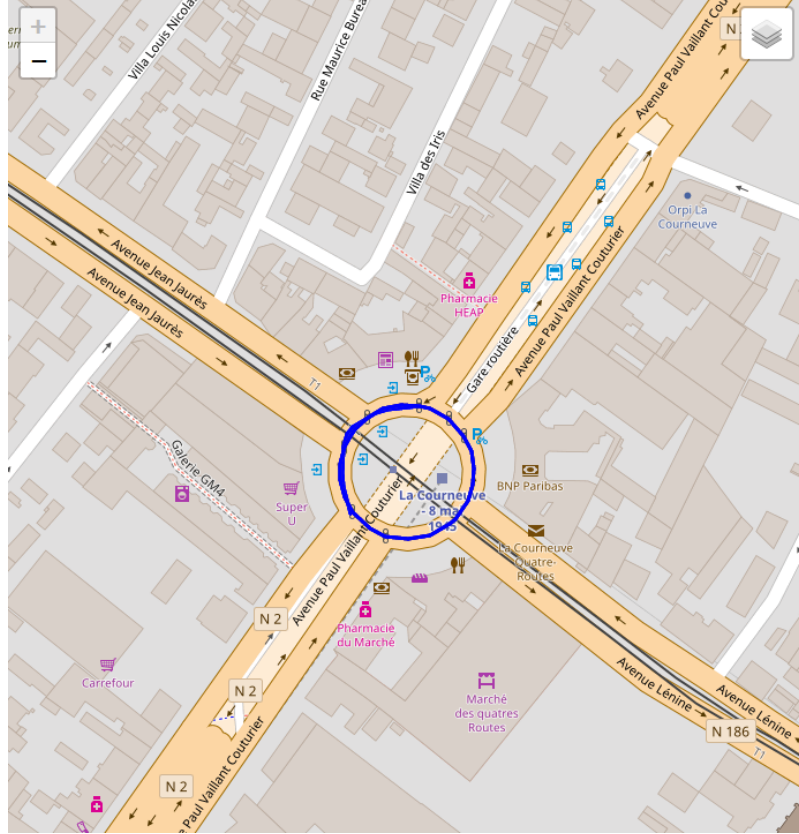
Time	February 15, 2011 10:12 PM
Changeset	
User	
Lat	48.911918
Lon	2.39554
Tags	
addr:housenumber	11
addr:street	Allée des myosotis
amenity	toilets
name	My Home

Ce contributeur a ajouté comme aménité les toilettes de chez lui... Avant d'être supprimé par un modérateur (du Data Working Group d'OSM).

Par conséquent, ce contributeur n'est pas fiable.

Contributeur U2

nombre de contributions	% création	% modification	% suppression	% objets réutilisés par d'autres	% d'objets qui n'ont pas été édités	% d'objets qui n'ont pas été supprimés
1	0	1	0	0	0	1



Version	4 export	5 export	6 export
Time	August 1, 2017 1:07 AM	August 6, 2017 10:18 PM	September 5, 2017 10:58 PM
Changeset			
User			
Tags			
bicycle	yes	yes	yes
foot	yes	yes	yes
highway	primary	primary	primary
junction	roundabout	roundabout	roundabout
lanes	2	2	2
maxspeed	30	30	30
name	Place du 8 Mai 1945	Place du 8 Mai 1945	Place du 8 Mai 1945
oneway	yes		yes
sidewalk	both	both	both
surface	asphalt	asphalt	asphalt

Ce contributeur a supprimé le tag « oneway = yes » du rond-point. Cette contribution est en fait une erreur de la part du contributeur (cf. la discussion sur le changeset qui contient la contribution en question à la page suivante). On note que ce changeset est un bot automatique lancé sur une très grande zone.

Ce contributeur n'est donc pas fiable

Discussion

S'abonner

Commentaire de [redacted] il y a 11 mois

Dear [redacted]

Noticed you recently deleted many 'oneway'-tags from roundabout. While this might be OK for normal highways, I'm not so sure about cycleways.

For example, you deleted the oneway-tag from the following road: <https://www.openstreetmap.org/way/197537357>

In the Netherlands, many roundabout+cycleways are oneway streets, but there are many exceptions. Therefore, I think it's useful to include the oneway-tag on roundabouts. Furthermore, I'm not sure whether 'roundabout' implies 'oneway' for cycleways.

What do you think?

Kind regards,

[redacted]

Commentaire de [redacted] il y a 11 mois

Hello [redacted]

I found the information in the wiki.

<http://wiki.openstreetmap.org/wiki/Tag:junction%3Droundabout>

In the German translation roundabout implies the oneway. In the English Version, the definition is not so strong like in German.

In the case I am really wrong, any idea to fix my fault?

Best regards

[redacted]

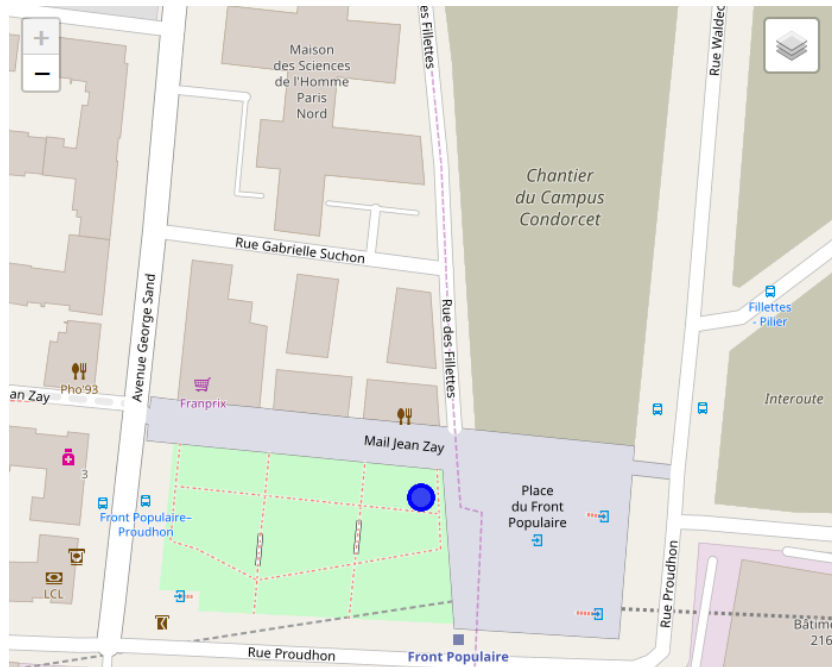
Commentaire de [redacted] il y a 11 mois

This changeset has been reverted fully or in part by changesets 51764077, 51763240 where the changeset comment is: Revert undiscussed mass removal of 'oneway=yes' on roundabouts performed by user Grauer after seeing this flagged on Osmose; after complaints from the Dutch community the user asked DWG to help with revert.



Contributeur U3

nombre de contributions	% création	% modification	% suppression	% objets réutilisés par d'autres	% d'objets qui n'ont pas été édités	% d'objets qui n'ont pas été supprimés
1	1	0	0	0	1	0



Version	4 export	5 export	6 export	7
Time	November 24, 2013 1:41 AM	January 7, 2014 10:53 PM	March 1, 2014 2:04 PM	June 19, 2016 9:44
Changeset				
User				
Lat	48.9068701	48.9068701	48.9068701	0
Lon	2.3657737	2.3657737	2.3657737	0
Tags				
STIF:zone	1	1	1	
name	Front Populaire	Front Populaire	Front Populaire	
railway	station	station	station	
start_date	2012-12-18	2012-12-18	2012-12-18	
station	subway	subway	subway	
type:RATP	metro	metro	metro	
wikipedia		fr:Front Populaire (m%C3%A9tro de Paris)	fr:Front Populaire (métro de Paris)	

Ce contributeur a simplement corrigé la valeur du tag « wikipedia » qui était mal écrite dans la version précédente. La suppression de cet objet n'est pas due au fait que la contribution était mauvaise. La contribution provient d'un changeset qui recouvre l'Île de France et qui avait pour but, via l'outil Osmose, de rétablir les valeurs de tags selon le formalisme d'OSM. **On ne peut donc rien dire de ce contributeur...**

Contributeur U4

nombre de contributions	% création	% modification	% suppression	% objets réutilisés par d'autres	% d'objets qui n'ont pas été édités	% d'objets qui n'ont pas été supprimés
1	0	1	0	0.005747126	0	1

Groupe de modifications

Debug - ligne de bus

Fermé il y a environ 2 ans par

Attributs

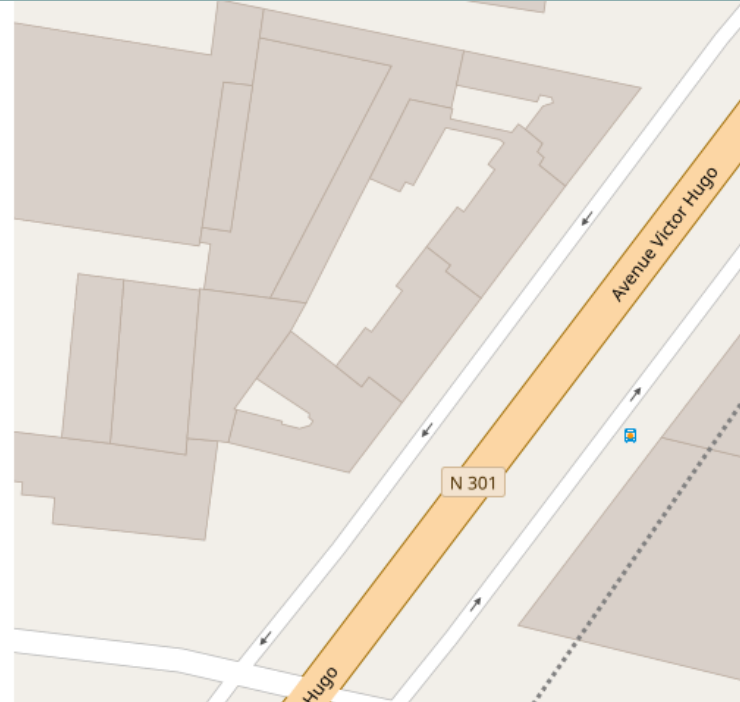
created_by

OpenBeerMap
javascript editor

Discussion

S'abonner

Commentaire



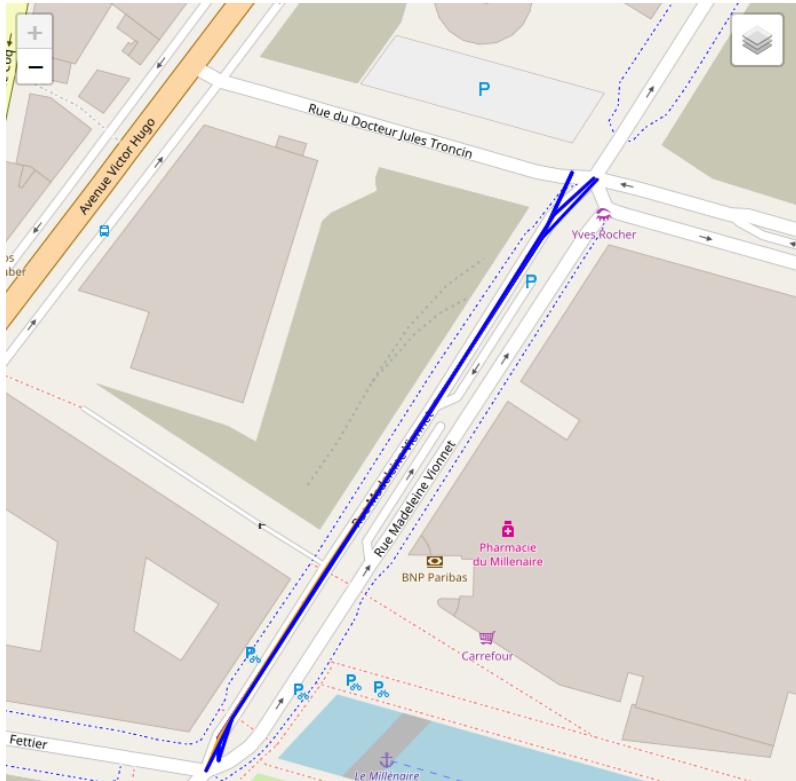
La seule contribution de ce contributeur a été faite sur un arrêt de bus (node). Il a ajouté une note permettant à l'utilisateur suivant d'ajouter ce node à la relation indiquée par ce contributeur, qui correspond aux arrêts du noctilien n°43. **On peut donc dire que ce contributeur est fiable.**

Ve	export	5	export	6	export
Time		May 26, 2016 9:19 PM		May 27, 2016 8:20 AM	
Changeset		[Redacted]			
User		Contributeur U4			
Lat		48.9065538		48.9065538	
Lon		2.3759265		2.3759265	
Tags					
created_by					
fixme:relation		add to 1057576		[Redacted]	
highway		bus_stop		bus_stop	
note		line 65		line 65	
public_transport				platform	
route_ref					
shelter					
tactile_paving					

Le contributeur suivant supprime son tag puisqu'il a bien ajouté le node à la relation. La note n'est donc plus utile à ce moment-là.

Contributeur U5

nombre de contributions	% création	% modification	% suppression	% objets réutilisés par d'autres	% d'objets qui n'ont pas été édités	% d'objets qui n'ont pas été supprimés
1	0	1	0	0.	0	1



Version	1	5	6	7
Time	1, 2014 6:06 PM	May 3, 2017 5:16 PM	December 15, 2017 3:18 PM	April 26, 2018 6:43 PM
Changeset		Contributeur U5		
User		Contributeur U5		
Tags				
highway	lential	residential	residential	residential
name	Madeleine Vionnet	Rue Madeleine Vionnet	Rue Madeleine Vionnet	Rue Madeleine Vionnet
old_name	de la Gare	Rue de la Gare	Rue de la Gare	Rue de la Gare
oneway		yes	yes	yes
surface		asphalt	asphalt	asphalt

Ce contributeur a ajouté le tag « surface = asphalt » à la route. Cette contribution fait partie d'un changeset dans lequel il a ajouté ce tag à plusieurs routes. Au vu de l'outil d'édition utilisé (StreetComplete¹) et la zone du changeset, on peut déduire que cette contribution a été faite sur le terrain. La contribution faite sur

¹ <https://wiki.openstreetmap.org/wiki/FR:StreetComplete>

Aubervilliers a été par la suite éditée par d'autres mais ces éditions n'ont pas remis en cause cet ajout de tag, ce qui semble même confirmer cette contribution. Par conséquent, **ce contributeur peut être vu comme fiable**.

Groupe de modifications :

Add road surfaces

Fermé il y a environ un an par [redacted]

Attributs

StreetComplete:quest	AddRoadSurface
_type	
created_by	StreetComplete 0.8
source	survey

Discussion [S'abonner](#)

Commentaire

Chemins (20)

- Boulevard MacDonald [redacted]
- Rue Lounès Matoub [redacted]
- Rue Chana Orlof [redacted]
- Boulevard MacDonald [redacted]
- Rue Gaston Tessier [redacted]
- Rue d'Aubervilliers [redacted]
- Rue d'Aubervilliers [redacted]
- Rue d'Aubervilliers [redacted]
- Rue d'Aubervilliers [redacted]
- Avenue de la Porte d'Aubervilliers [redacted]

v4)

Emprise du changeset contenant la contribution

Contributeurs ayant fait entre 3 et 10 contributions

Contributeur U6

nombre de contributions	% création	% modification	% suppression	% objets réutilisés par d'autres	% d'objets qui n'ont pas été édités	% d'objets qui n'ont pas été supprimés
3	0.7	0	0.3	0	0.33	1

Dans un même changeset, ce contributeur a ajouté 2 tronçons de rails (pour cartographier la ligne du RER B), et supprimé un tronçon. Les éditions suivantes montrent un enrichissement et une mise à jour de l'information contribué par ce contributeur. Par conséquent, **ce contributeur est fiable**.

Groupe de modifications :

Modifications mineures sur la relation du RER B

Fermé il y a environ 9 ans par

Attributs

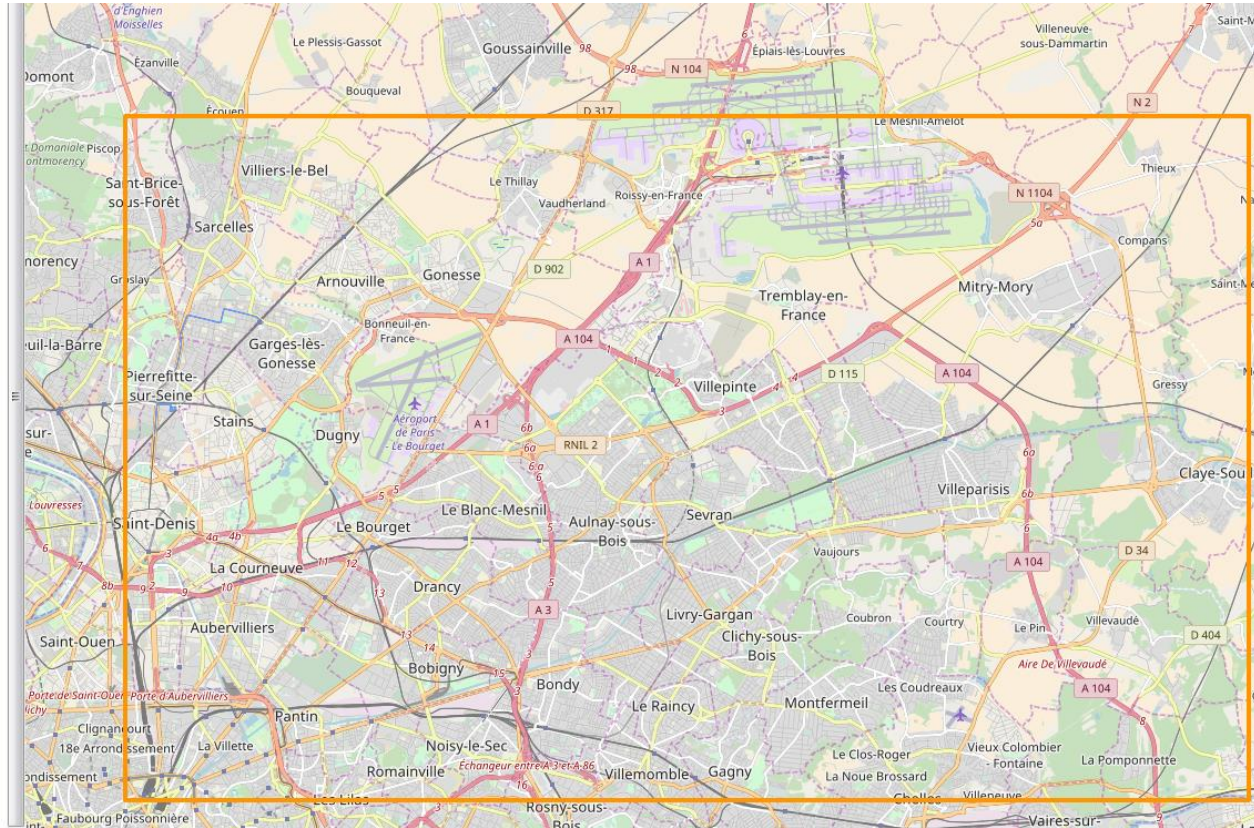
created_by JOSM/1.5 (1632 fr)

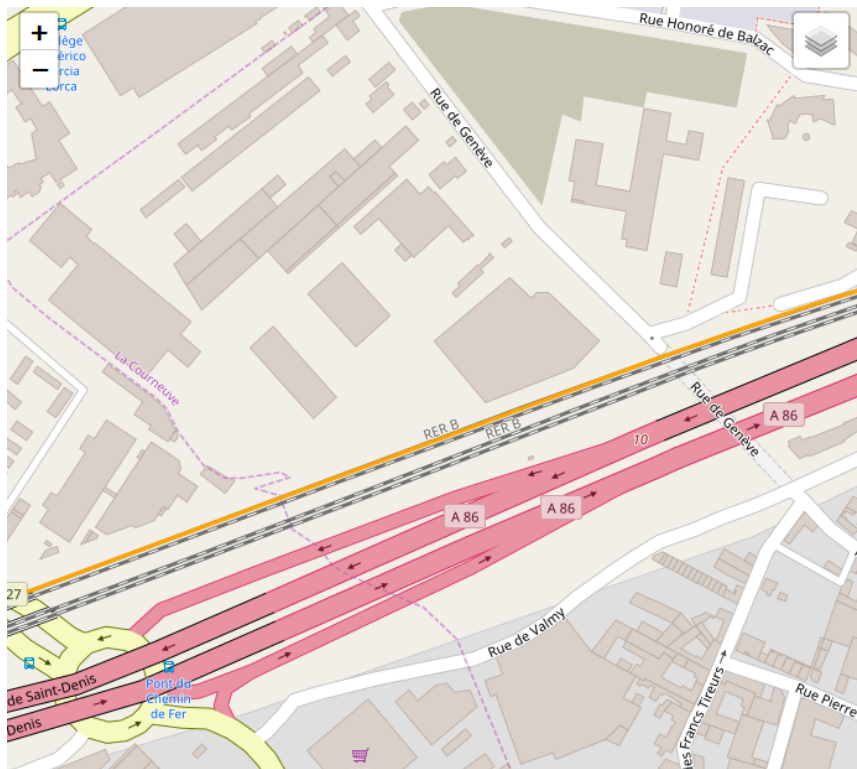
Discussion [S'abonner](#)

Commentaire

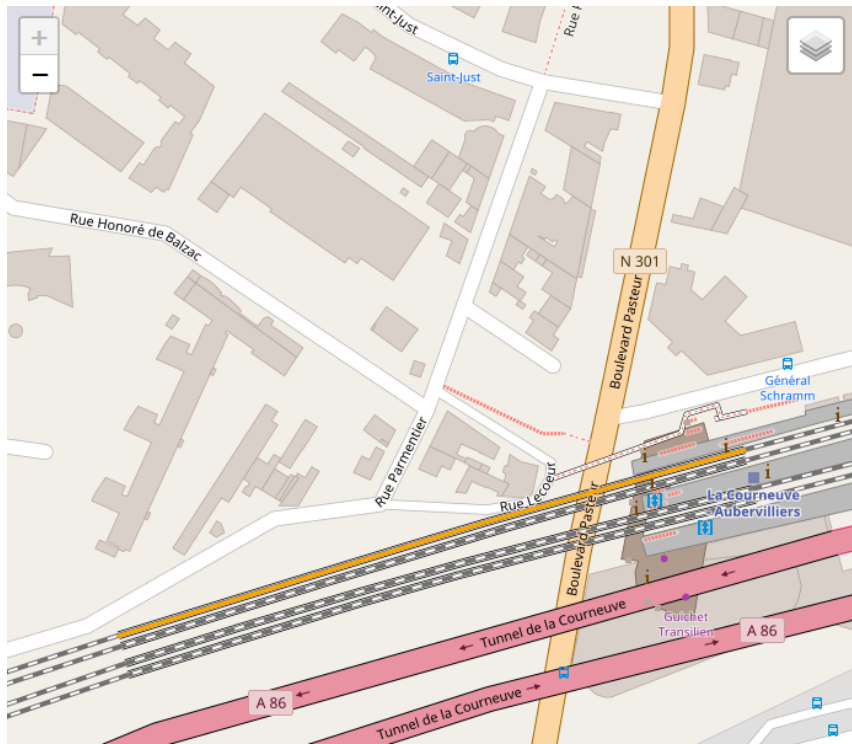
Chemins (10)

- RER B Ligne de la Plaine à Hirson et Anor
- RER B
- RER B
- RER B
- RER B Ligne de la Plaine à Hirson et Anor
- RER B Ligne de la Plaine à Hirson et Anor
- RER B Ligne de la Plaine à Hirson et Anor
- RER B Ligne d'Aulnay à Roissy





Version	1	2	3	4	5
Time	June 4, 2009 11:05 PM	May 22, 2011 9:10 AM	June 10, 2011 4:15 PM	December 5, 2011 2:25 PM	January 12, 2014
Changeset	Contributeur U6				
User					
Tags					
alt_name	Ligne de la Plaine à Hirson et Anor	Ligne de la Plaine à Hirson et Anor	Ligne de la Plaine à Hirson et Anor	Ligne de la Plaine à Hirson et Anor	Ligne de la Plaine à Hirson et Anor
electrified				contact_line	contact_line
frequency				50	50
gauge				1435	1435
importance					
maxspeed					90
name	RER B	RER B	RER B	RER B	RER B
oneway			-1	-1	-1
operator	SNCF	SNCF	SNCF	SNCF	SNCF
railway	rail	rail	rail	rail	rail
usage					
voltage				25000	25000



Version	1	2	3	4	5	6
Time	June 4, 2009 11:05 PM	June 10, 2011 4:15 PM	June 10, 2011 4:58 PM	December 5, 2011 2:12 PM	January 12, 2014 12:42 PM	January 12, 2014 12:42 PM
Changeset	Contributeur U6					
User	Contributeur U6					
Tags						
alt_name	Ligne de la Plaine à Hirson et Anor	Ligne de la Plaine à Hirson et Anor				
bridge	yes	yes	yes	yes	yes	yes
electrified				contact_line	contact_line	contact_line
frequency				50	50	50
gauge				1435	1435	1435
importance						
layer	1	1	1	1	1	1
maxspeed					90	90
name	RER B	RER B				
oneway		-1	-1	-1	-1	-1
operator	SNCF	SNCF	SNCF	SNCF	SNCF	SNCF
railway	rail	rail	rail	rail	rail	rail
usage						main
voltage				25000	25000	25000

Contributeur U7

nombre de contributions	% création	% modification	% suppression	% objets réutilisés par d'autres	% d'objets qui n'ont pas été édités	% d'objets qui n'ont pas été supprimés
5	0.7	1	0.3	0	1	1

Ce contributeur a enrichi certains points d'intérêt en y ajoutant l'ID Wikidata de l'objet correspondant. **Ce contributeur est plutôt fiable.**

Groupe de modifications :

fixing and cleaning up wikipedia/wikidata using
https://wiki.openstreetmap.org/wiki/Wikidata%2BOSM_SPARQL_query_service

Fermé il y a 10 mois par [redacted]

Attributs

created_by JOSM/1.5 (12712 en)

Discussion

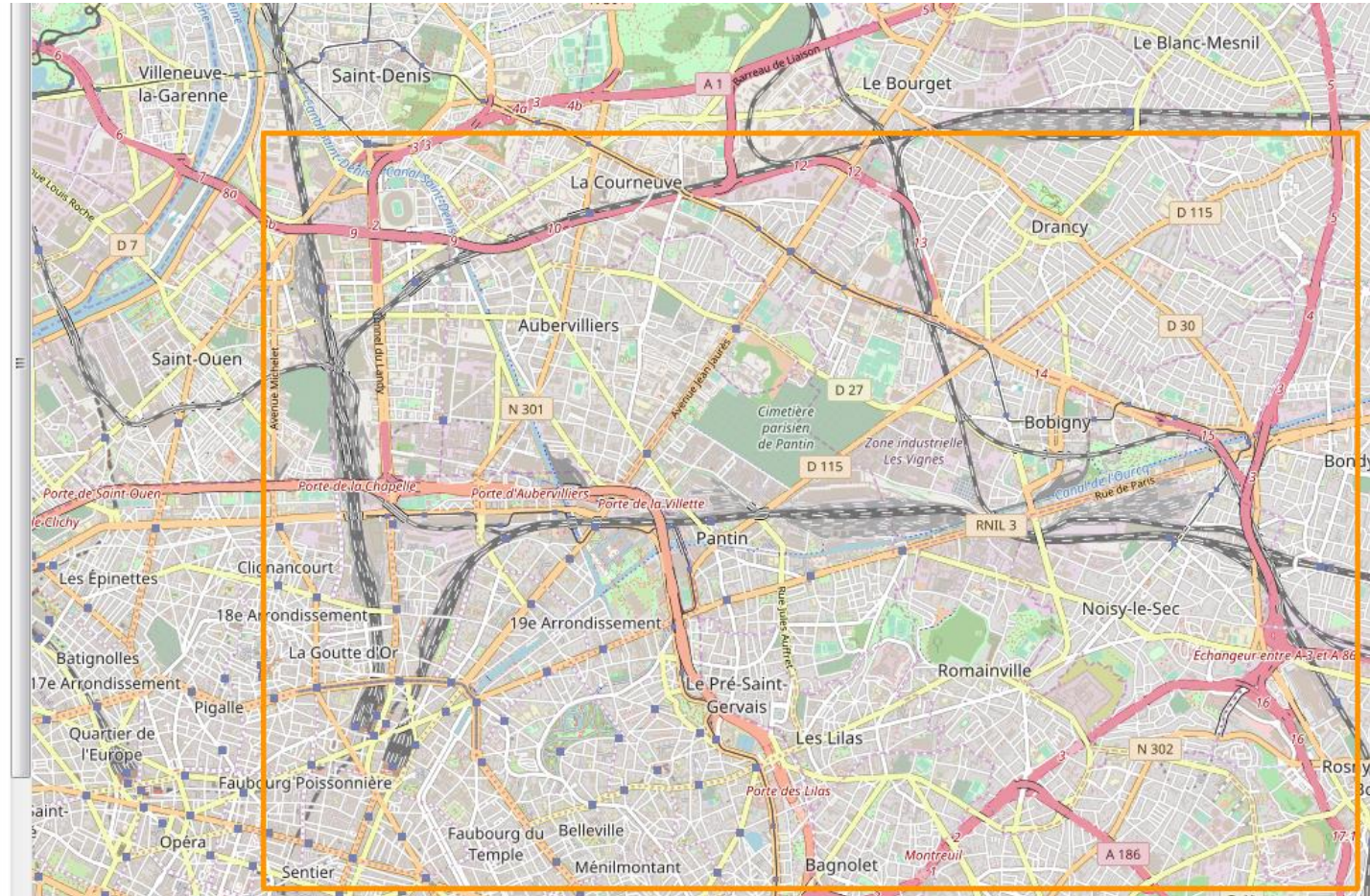
S'abonner

Commentaire

Nœuds (1 à 20 sur 51)

1 2 3

- Bondy [redacted]
- [redacted]
- QVC France [redacted]
- Bobigny [redacted]
- Autolib Pantin/Cartier Bresson/84 [redacted]
- v3) [redacted]
- [redacted]



Contributeur U8

nombre de contributions	% création	% modification	% suppression	% objets réutilisés par d'autres	% d'objets qui n'ont pas été édités	% d'objets qui n'ont pas été supprimés
7	0	0.4	0.6	0.006	1	1

Il semble que ses contributions soient issues d'un bot, vu l'étendue du changeset.

Groupe de modifications : [redacted] x

(aucun commentaire)

Fermé il y a plus de 7 ans par [redacted]

Attributs

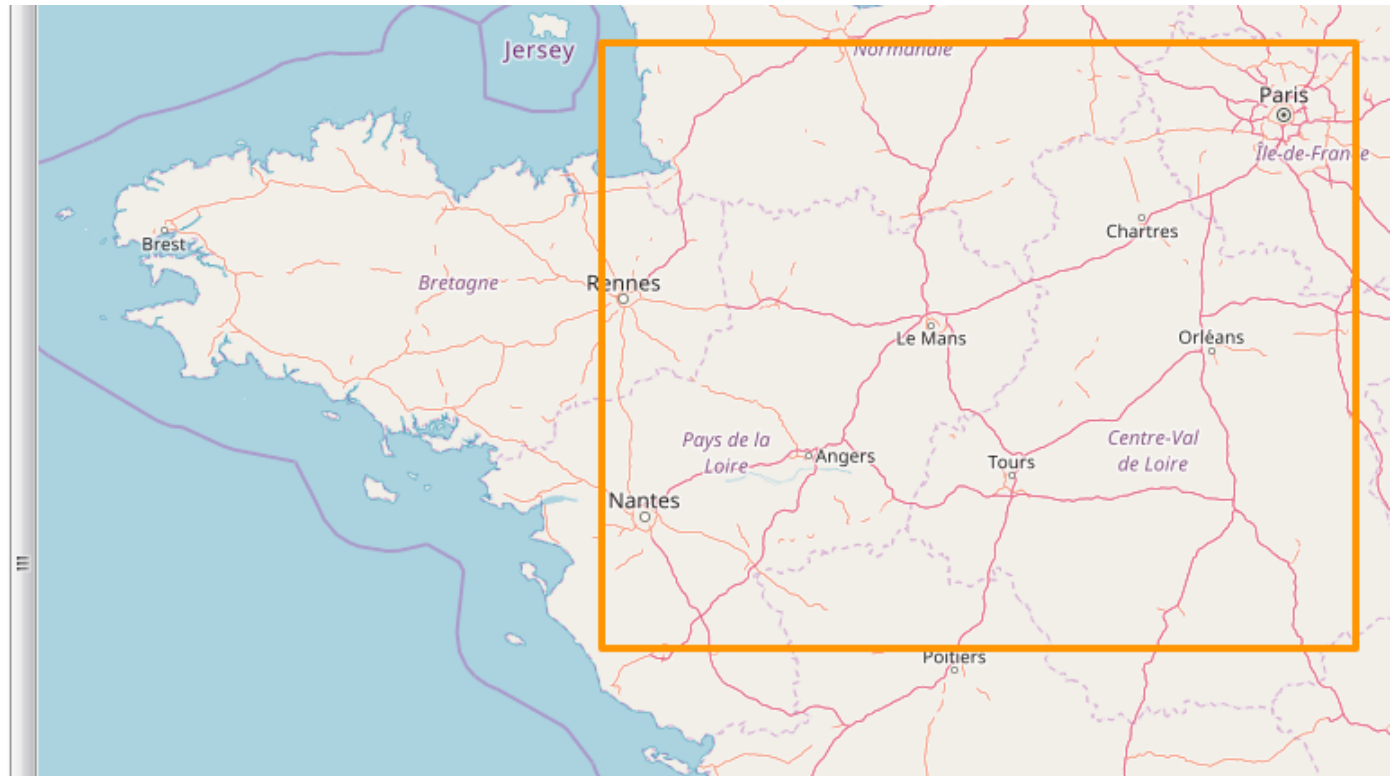
created_by

JOSM/1.5 (3790 fr)

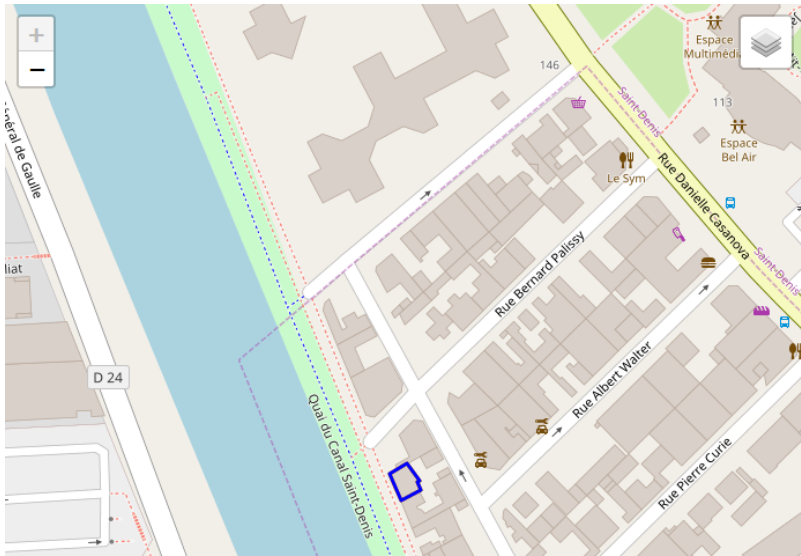
Discussion

S'abonner

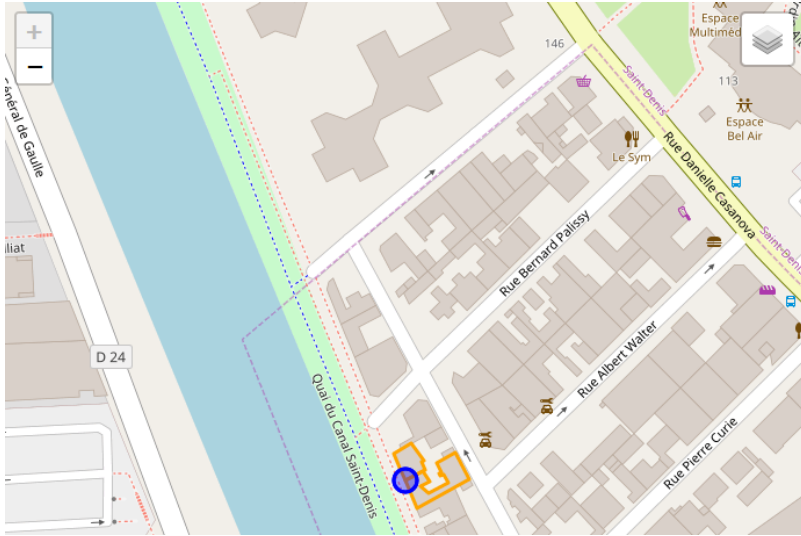
Commentaire



Editions : elles ne changent que très peu les données. Peut-être qu'il a voulu confirmer la validité de la donnée par une édition qui ne change rien ?

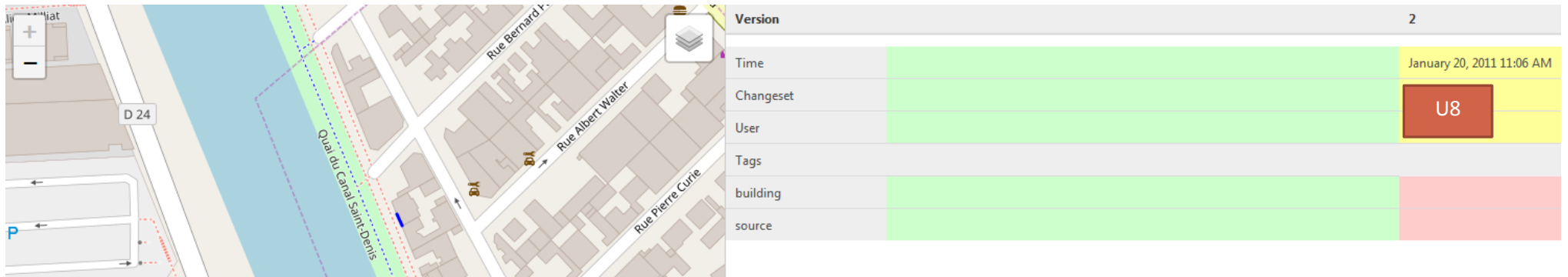


Version	1	export	2	export
Time	August 18, 2010 11:18 AM		January 20, 2011 11:06 AM	
Changeset			U8	
User				
Tags				
building	yes		yes	
source	extraction vectorielle v1 cadastre-dgi-fr source : Direction Générale des Impôts - Cadas. Mise à jour : 2010		extraction vect	



Version	2	export
Time	January 20, 2011 11:06 AM	
Changeset	U8	
User		
Lat	48.9231297	
Lon	2.3674462	
Tags		

Suppression :



The image shows a map interface with a street view of a canal area. The map includes labels for 'D 24', 'Quai du Canal Saint-Denis', 'Rue Bernard P.', 'Rue Albert Walter', and 'Rue Pierre Curie'. A metadata table is overlaid on the right side of the map, showing details for a specific version of the data.

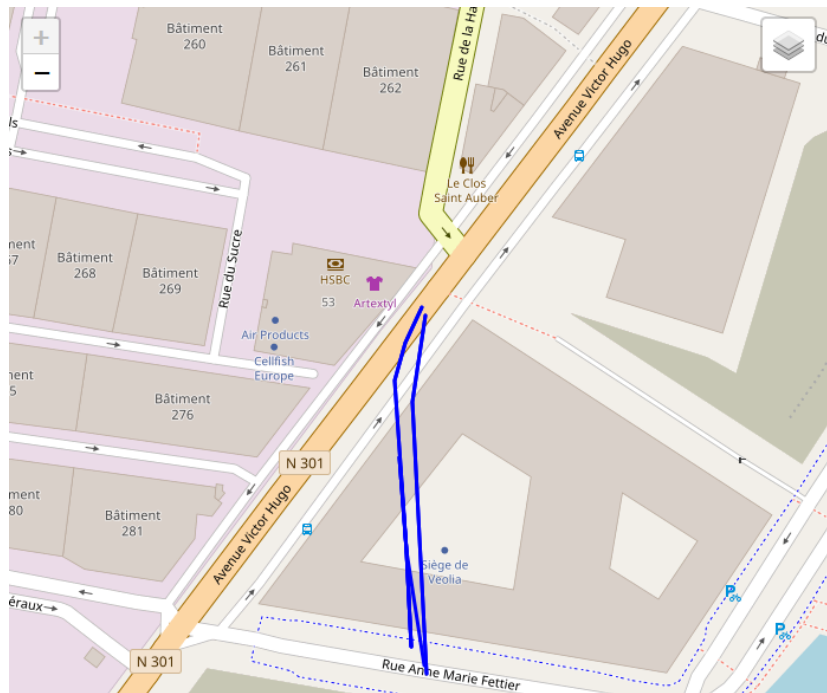
Version		2
Time		January 20, 2011 11:06 AM
Changeset		U8
User		
Tags		
building		
source		

Les contributions de ce compte ne sont pas assez explicites pour révéler la fiabilité du contributeur. Cependant, comme il n'y a pas eu de remise en cause de ses contributions, elles ne semblent pas être fausses. **On ne peut donc rien dire de ce contributeur.**

Contributeur U9

nombre de contributions	% création	% modification	% suppression	% objets réutilisés par d'autres	% d'objets qui n'ont pas été édités	% d'objets qui n'ont pas été supprimés
9	0.6	0.3	0.1	0.	0.7	0.9

Ce contributeur a créé une route, qui a été éditée par la suite. Vu l'évolution des valeurs de tag, on comprend que la suppression de l'objet est due à la présence de travaux sur la zone en 2014. La contribution de ce contributeur était donc pertinente quoi qu'imprécise d'après les modifications de tags sur la version suivante.



Version	1	2	3	4	5
Time	April 29, 2011 12:04 AM	April 29, 2011 3:48 PM	June 15, 2012 12:46 AM	July 19, 2012 10:48 AM	June 1, 2013 1
Changeset	U9				
User					
Tags					
access		no		no	
highway	primary	road	primary	primary	unclassified
name	Rue de la Haie Coq	Rue de la Haie Coq	Rue de la Haie Coq	Rue de la Haie Coq	Rue de la Hai
oneway	yes		yes	yes	
ref	N 301	N 301			

Ce contributeur a également ajouté deux POI : le premier indiquant la présence de bateau électrique (POI resté intact), le deuxième pour indiquer la présence d'un centre commercial (supprimé par la suite car l'information a été transférée sur les attributs du bâtiment matérialisé par un way). Globalement, **on peut donc dire que ce contributeur est fiable.**

Contributeur U10

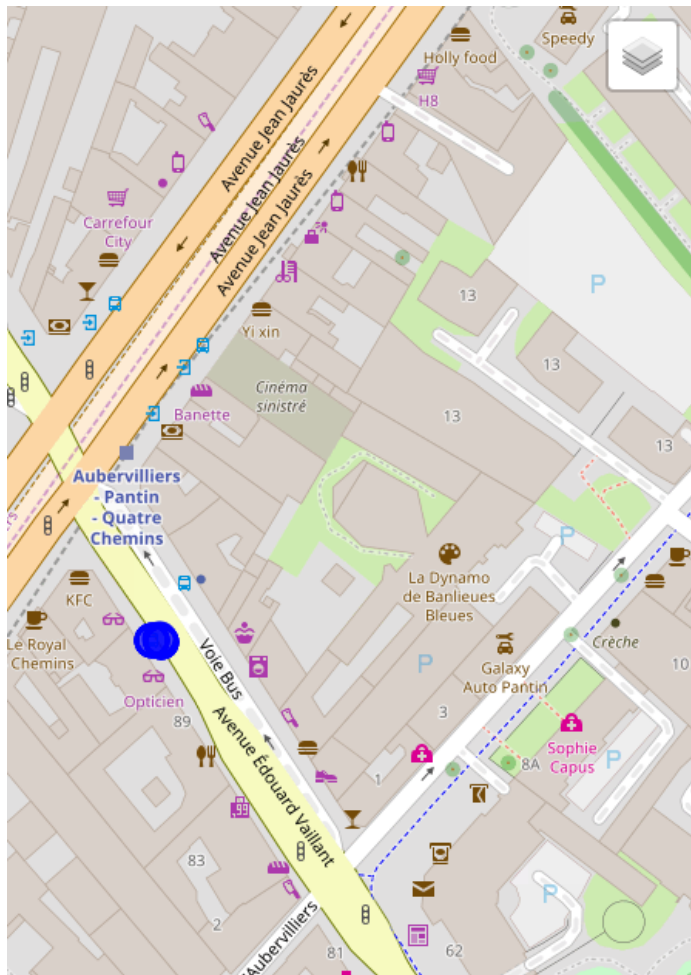
nombre de contributions	% création	% modification	% suppression	% objets réutilisés par d'autres	% d'objets qui n'ont pas été édités	% d'objets qui n'ont pas été supprimés
10	0	1	0	0	0.5	1

Ce contributeur n'a modifié que des objets de type way. Ses contributions sont toutes dans un changeset. Elles consistent à modifier la valeur de certains attributs, comme l'ajout d'un accent aigü sur un 'E' majuscule, réécrire 'A86' en 'A 86' en y insérant un espace. Au vu du nom du contributeur, on peut penser que ses contributions sont automatiques. Les contributions de ce contributeur ne sont pas incorrectes, **mais elles ne sont pas assez significatives pour affirmer la fiabilité de ce contributeur.**

Contributeur U11

nombre de contributions	% création	% modification	% suppression	% objets réutilisés par d'autres	% d'objets qui n'ont pas été édités	% d'objets qui n'ont pas été supprimés
31	0.9	0.1	0	0	0.9	1

Ce contributeur a notamment ajouté des informations sur les arrêts de bus, soit en ajoutant des attributs sur des objets existants, soit en créant de nouveaux arrêts de bus. Bien que ces objets aient été modifiés par la suite, il n'y a pas eu de remise en cause des informations que ce contributeur a ajoutées. Ce contributeur a également entré des informations locales telles que les adresses de bâtiment et des points de recyclage de verre. **Ce contributeur est donc fiable.**

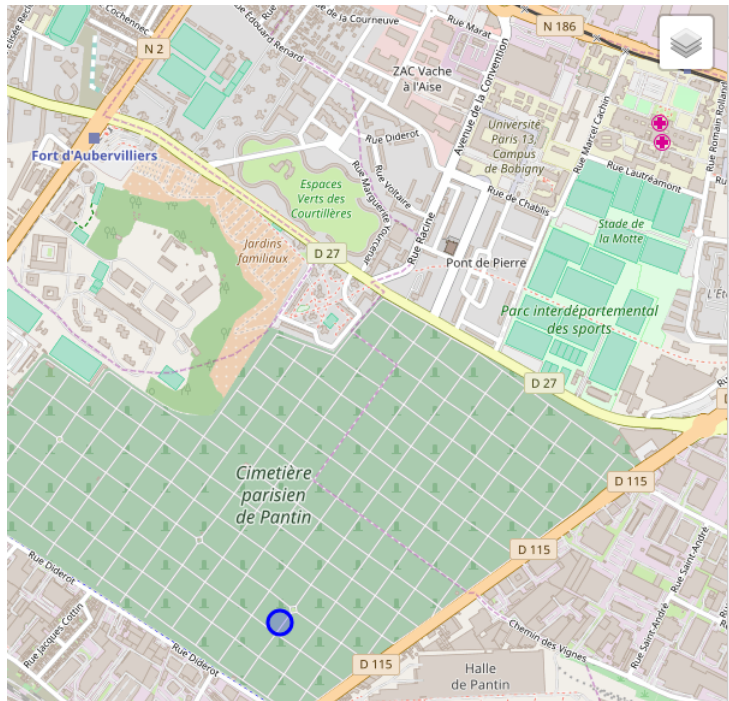


Version	1 <input type="button" value="export"/>	2 <input type="button" value="export"/>	3 <input type="button" value="export"/>	4 <input type="button" value="export"/>
Time	April 11, 2014 4:20 PM	April 25, 2014 2:46 AM	May 8, 2014 3:22 AM	September 26, 2015 3:31 PM
Changeset				U11
User				
Lat	48.9032053	48.9031977	48.9031977	48.9031977
Lon	2.3923645	2.392371	2.392371	2.392371
Tags				
STIF:zone				
bench				
bus				
gtfs_id				
highway	bus_stop	bus_stop	bus_stop	bus_stop
name				Quatre Chemins — Édouard Vaillant
public_transport				platform
ref:FR:STIF				
ref:FR:STIF:stop_id				
shelter			yes	yes
source				
source:position				
tactile_paving				

Contributeur U12

nombre de contributions	% création	% modification	% suppression	% objets réutilisés par d'autres	% d'objets qui n'ont pas été édités	% d'objets qui n'ont pas été supprimés
41	0.02	0	0.98	0	0.98	0.98

Le seul élément créé par ce contributeur est un hôpital en plein milieu d'un cimetière, qui a été supprimé un mois plus tard par un autre contributeur.



Version	1 export
Time	February 4, 2014 1:05 PM
Changeset	
User	
Lat	48.9028853
Lon	2.410445
Tags	
amenity	hospital

Le reste de ses contributions correspondent à la suppression du cimetière de Pantin. Toutes ses contributions ont été annulées par la suite. On ne le remarque pas à partir des indicateurs calculés car le nombre d'objets édités est compté à partir des objets visibles. **Ce contributeur n'est pas fiable, c'est même un vandale.**

Contributeur U13

uid	nombre de contributions	% création	% modification	% suppression	% objets réutilisés par d'autres	% d'objets qui n'ont pas été édités	% d'objets qui n'ont pas été supprimés
375079	52	0.1	0.4	0.5	0.1	0.9	1

Création :

Un bâtiment en 2011, qui n'a jamais été édité jusqu'à présent

Modification :

Correction de la géométrie de bâtiments initialement importés automatiquement à partir du cadastre

Suppression :

Suppression de doublons tels que deux bâtiments identiques (donc superposés) ayant été importés automatiquement

Ce contributeur est fiable.

Contributeur U14

nombre de contributions	% création	% modification	% suppression	% objets réutilisés par d'autres	% d'objets qui n'ont pas été édités	% d'objets qui n'ont pas été supprimés
64	0	0	1	0.005	1	1

Ce contributeur a exclusivement participé en supprimant des données. Il a supprimé des bâtiments près de l'impasse Michel qui, d'après la note de son changeset en 2016, ont été détruits. **Ce contributeur semble être fiable, puisqu'il n'a jamais été remis en cause.**

Contributeur U15


nombre de contributions	% création	% modification	% suppression	% objets réutilisés par d'autres	% d'objets qui n'ont pas été édités	% d'objets qui n'ont pas été supprimés
76	0.4	0.6	0	0.1	0.8	0.9

Création :

Des objets très locaux tels que des passerelles, des marches, des rues piétonnes, qui sont assez bien renseignés dès la 1^{ère} version. Ces objets ont été par la suite édités mais de manière à compléter ce que le contributeur avait ajouté.

Modification :

Mise à jour de noms de rue (comme la rue Rosa Parks), de la géométrie des tronçons de route. Enrichissement attributaire sur certains bâtiments tels que la crèche ci-dessous.



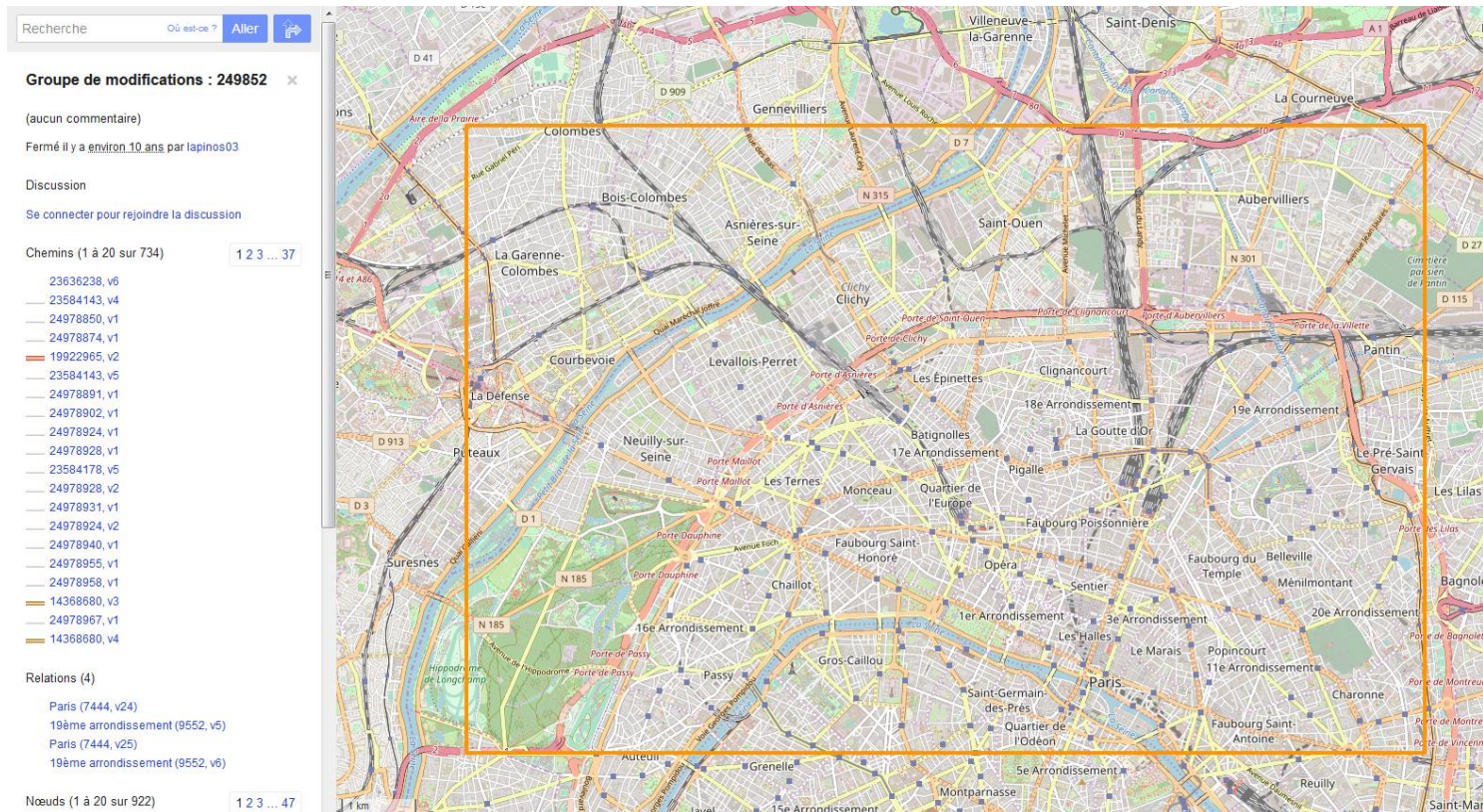
Version	1	2
Time	July 11, 2010 1:57 PM	February 17, 2013 6:22 PM
Changeset		
User		U15
Tags		
addr:housename		Crèche Départementale Jean Jaurès
addr:housenumber		110
addr:postcode		93120
addr:street		Avenue Jean-Jaurès
amenity		school
building	yes	yes
name		Crèche Départementale Jean Jaurès
source	extraction vectorielle v1 cadastre-dgi-fr source : Direction Générale des Impôts - Cadas. Mise à jour : 2010	http://www.seine-saint-denis.fr/Les-cresches-du-De

Ce contributeur est fiable.

Contributeur U16

nombre de contributions	% création	% modification	% suppression	% objets réutilisés par d'autres	% d'objets qui n'ont pas été édités	% d'objets qui n'ont pas été supprimés
108	0.7	0.3	0.01	0	0.93	0.97

Ce contributeur a surtout créé/modifié des ways, en particulier les routes. L'emprise du seul changeset qu'il a produit sur la zone est assez étendue, donc on peut supposer qu'il a contribué de manière automatique.



Création :

Création de routes qui ont été éditées/complétées par la suite, car elles étaient tracées sans être renseignées par de nombreux attributs (ce qui nous pousse à croire que ses contributions n'étaient pas manuelles)

Modification :

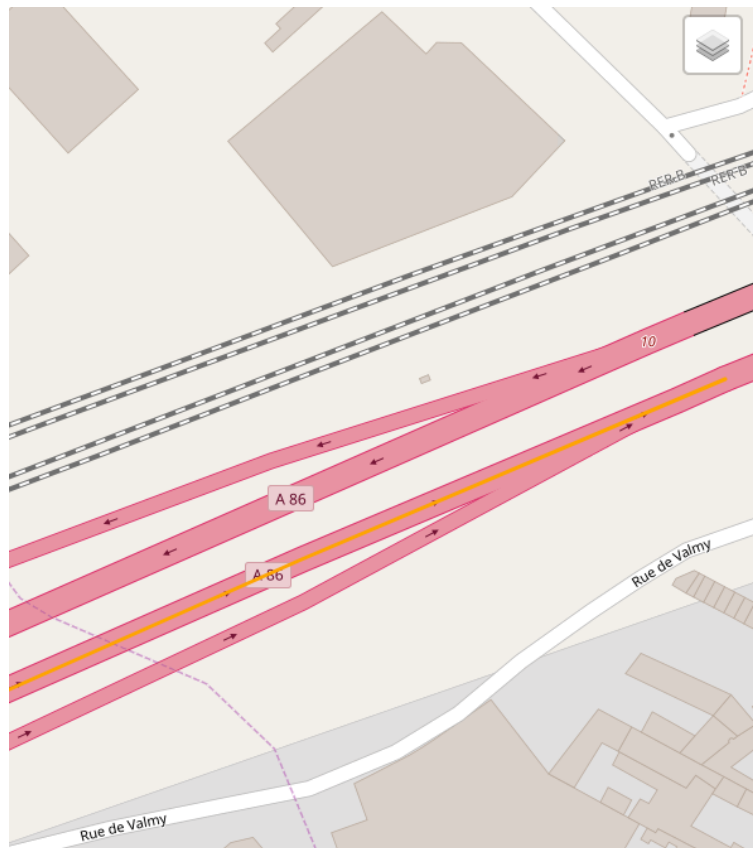
Les modifications faites portent sur la géométrie des routes. Ce contributeur a fortement participé à la structuration du réseau routier.

On ne peut rien dire sur ce contributeur. En effet, même s'il a beaucoup contribué, ses contributions ne sont pas d'une qualité notable.

Contributeur U17

nombre de contributions	% création	% modification	% suppression	% objets réutilisés par d'autres	% d'objets qui n'ont pas été édités	% d'objets qui n'ont pas été supprimés
201	0.01	0.99	0	0.07	0.97	0.99

Ce contributeur a modifié plusieurs fois des routes qu'ils avaient créées en premier lieu. Il a contribué en 2008, et ses contributions ont été conservées jusqu'en 2017. Les objets créés par ce contributeur sont plutôt riches sémantiquement (voir ci-dessous). Par conséquent, on peut dire que **ce contributeur est fiable**.

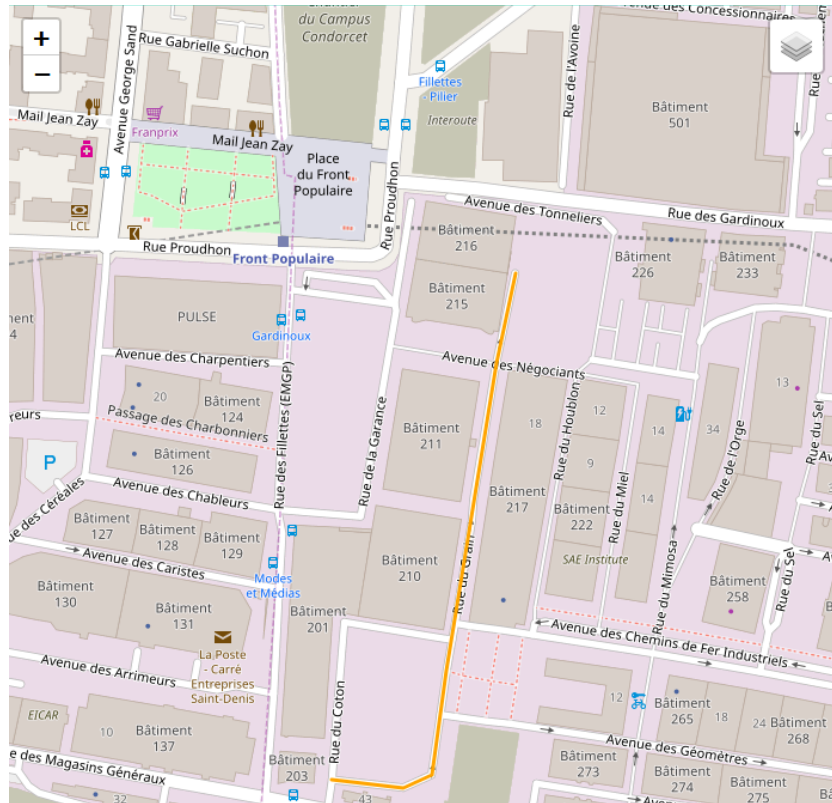


Version	1	2	3	4	5
Time	October 10, 2008 11:37 PM	March 11, 2009 10:20 PM	March 12, 2009 11:48 PM	June 16, 2009 3:37 PM	January 17, 2011 9:55 PM
Changeset	U17				
User					
Tags					
created_by	Potlatch 0.10d	Potlatch 0.10d	Potlatch 0.10d		
highway	motorway	motorway	motorway	motorway	motorway
lanes	3	3	3	3	3
lit					
maxspeed	90	90	90	90	90
name	Périphérique d'Île de France	Périphérique d'Île de France	Périphérique d'Île de France	Périphérique d'Île de France	Périphérique d'Île de France
oneway	true	true	true	yes	yes
operator					
ref	A86	A 86	A86	A 86	A 86
toll					

Contributeur U18

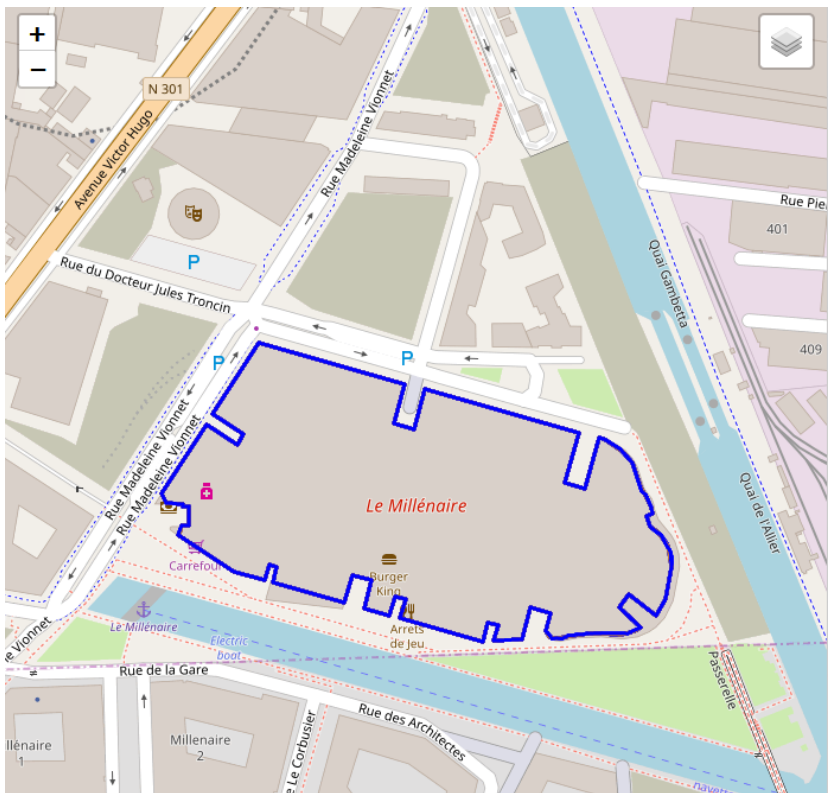
nombre de contributions	% création	% modification	% suppression	% objets réutilisés par d'autres	% d'objets qui n'ont pas été édités	% d'objets qui n'ont pas été supprimés
241	0.4	0.5	0.1	0.14	0.96	0.97

Ce contributeur a créé des routes, en a modifié d'autres souvent sur leurs géométries. Cependant les attributs renseignés pour les objets créés sont plutôt vagues :



Version	1	2	3	4	5	6
Time	July 4, 2012 12:45 AM	August 9, 2015 8:18 PM	October 8, 2015 6:00 AM	June 28, 2016 9:00 PM	September 5, 2016 12:15 PM	September 5, 2016 12:16 PM
Changeset	U18					
User						
Tags						
FIXME				check the oneway direction	check the oneway direction	check the oneway direction
highway	unclassified	unclassified	service	service	service	service
name	Rue N° 22	Avenue des Forgerons	Avenue des Forgerons	Avenue des Forgerons	Avenue des Forgerons	Rue du Grain
oneway		yes;-1	yes;-1	yes;-1	yes;-1	yes;-1

Il complète cependant la sémantique de certains objets en leur ajoutant des attributs pertinents :



Version	1	export	2	export
Time	March 13, 2012 2:07 PM		July 24, 2012 11:07 AM	
Changeset			U18	
User				
Tags				
addr:houseName			Le Millénaire	
addr:houseNumber			19	
addr:postcode			93300	
addr:street			Rue Madeleine Vionnet (ex rue de la Gare)	
architect				
building	yes		supermarket	
designation				
name	Le Millénaire		Le Millénaire	
old_addr:street				
shop	mall		mall	
source	cadastre-dgi-fr source : Direction Générale des Impôts - Cadastre. Mise à jour : 2012		cadastre-dgi-fr source : Direction Générale des Impôts - Cadastre. Mise à	
source:architect				
website				
wikidata				
wikipedia				

On peut dire que ce contributeur est fiable.

Contributeur U19

nombre de contributions	% création	% modification	% suppression	% objets réutilisés par d'autres	% d'objets qui n'ont pas été édités	% d'objets qui n'ont pas été supprimés
317	0.08	0.9	0.01	0.16	0.95	0.99

Ce contributeur a créé des routes qui perdurent encore aujourd'hui, et dont les attributs sont plus ou moins corrects : en effet, sur les rues qu'il crée, il se trompe sur la valeur du tag highway.



Version	1	2	3	4	5	6
Time	January 11, 2009 3:17 PM	January 11, 2009 3:59 PM	February 17, 2009 3:39 PM	February 19, 2009 9:33 AM	February 19, 2009 9:37 AM	February 20, 2009 8:32 AM
Changeset	U19	U19				
User						
Tags						
created_by	Merkaator 0.12	Merkaator 0.12	Potlatch 0.10f	Potlatch 0.10f	Potlatch 0.10f	Potlatch 0.10f
highway	tertiary	secondary	secondary	secondary	secondary	secondary
name	Avenue de la République	Avenue de la République	Avenue de la République	Avenue de la République	Avenue de la République	Avenue de la République
ref	D 20	D 20	D 20	D 20	D 20	D 20

Il s'autocorrige et corrige les autres aussi, notamment sur la géométrie des ways.

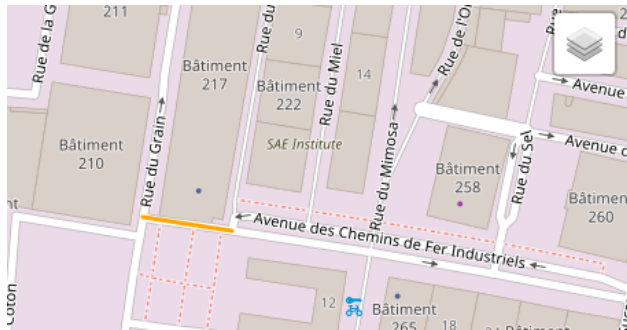
On peut dire que ce contributeur est globalement fiable

Contributeur U20

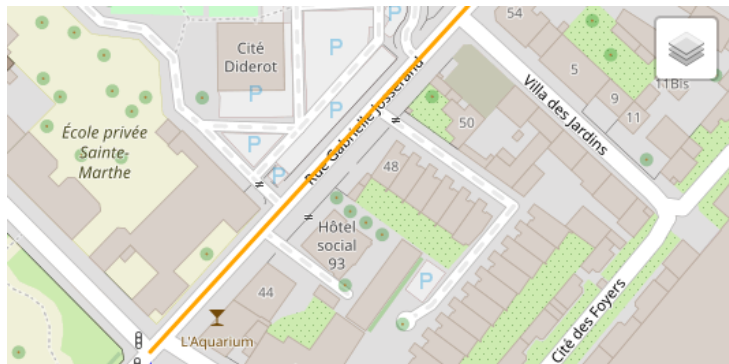
uid	nombre de contributions	% création	% modification	% suppression	% objets réutilisés par d'autres	% d'objets qui n'ont pas été édités	% d'objets qui n'ont pas été supprimés
705842	356	0.26	0.72	0.02	0.16	0.95	0.99

Création

Comme le contributeur précédemment étudié, celui-ci crée des objets qui perdurent encore aujourd'hui, mais dont les attributs sont initialement parfois pauvres.

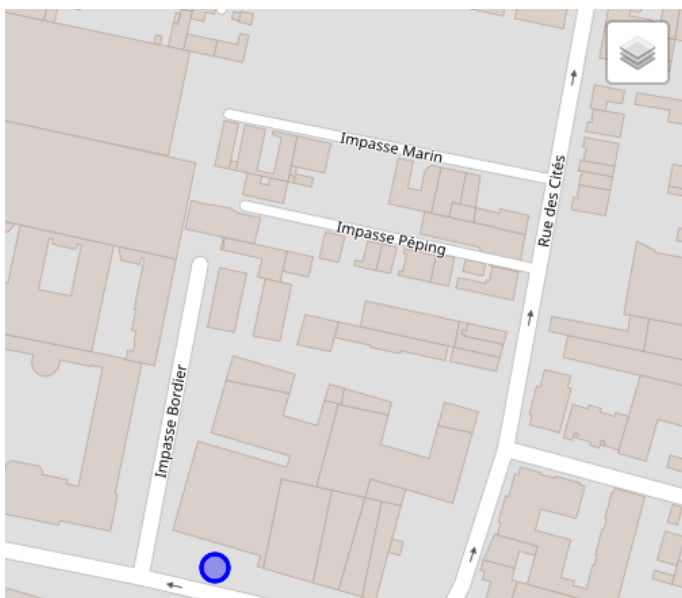


Version	1	2	3	4	5
Time	December 3, 2010 4:10 PM	October 8, 2015 6:00 AM	July 24, 2016 1:36 PM	June 4, 2017 10:42 AM	July 17, 2017 7:31
Changeset	U20				
User					
Tags					
highway	unclassified	service	service	service	service
name				Avenue des Chemins de Fer Industriels	Avenue des Chemins de Fer Industriels



Version	1	2	3	4	5
Time	February 3, 2011 6:01 PM	September 16, 2012 8:38 AM	December 23, 2012 3:42 PM	May 2, 2014 12:02 AM	September 4, 2015 12:57 PM
Changeset	U20				
User					
Tags					
highway	residential	residential	residential	residential	residential
name	Rue Gabrielle Josserand	Rue Gabrielle Josserand	Rue Gabrielle Josserand	Rue Gabrielle Josserand	Rue Gabrielle Josserand

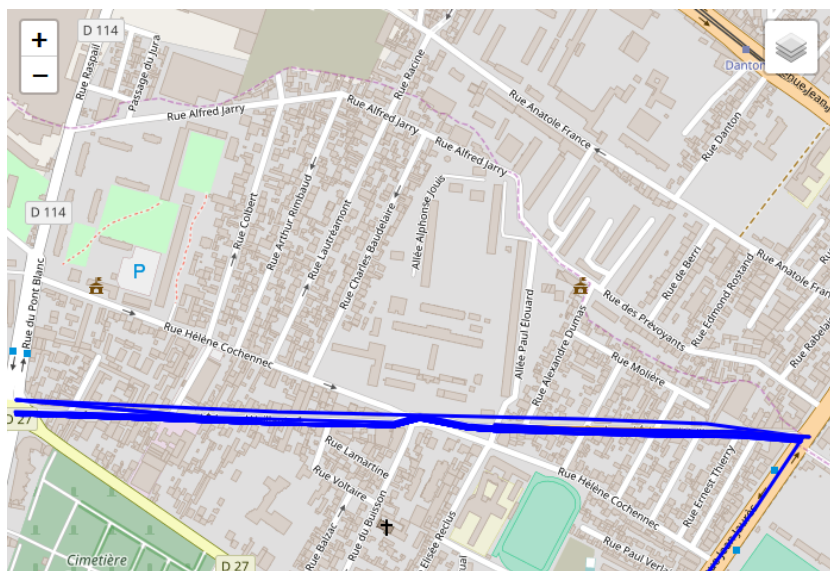
Le contributeur a également recensé les stations de Vélib à Aubervilliers, qui ont été mis à jour par la suite par d'autres contributeurs.



Version	1 <input type="button" value="export"/>	2 <input type="button" value="export"/>
Time	November 8, 2010 7:05 PM	April 26, 2013 10:19 PM
Changeset	U20	
User		
Lat	48.9037164	48.9037164
Lon	2.3856482	2.3856482
Tags		
amenity	bicycle_rental	bicycle_rental
capacity	45	45
network	velib'	Vélib'
operator		JCDecaux
ref	33003	33003

Modification

Il participe à la structuration du réseau routier en modifiant la géométrie des routes.



Vers	16	17	18	19	20	
Time	4, 2010 11:06 AM	March 24, 2011 7:59 PM	April 6, 2011 7:22 PM	February 17, 2014 12:47 PM	May 19, 2014 11:54 AM	August 3, 2014 1:05 PM
Changeset			U20			
User						
Tags						
created_by						
highway		residential	residential	residential	residential	residential
name	douart Vaillant	Boulevard Edouart Vaillant	Boulevard Edouart Vaillant	Boulevard Edouart Vaillant	Boulevard Edouard Vaillant	Boulevard Édouard Vaillant
ref						

Suppression

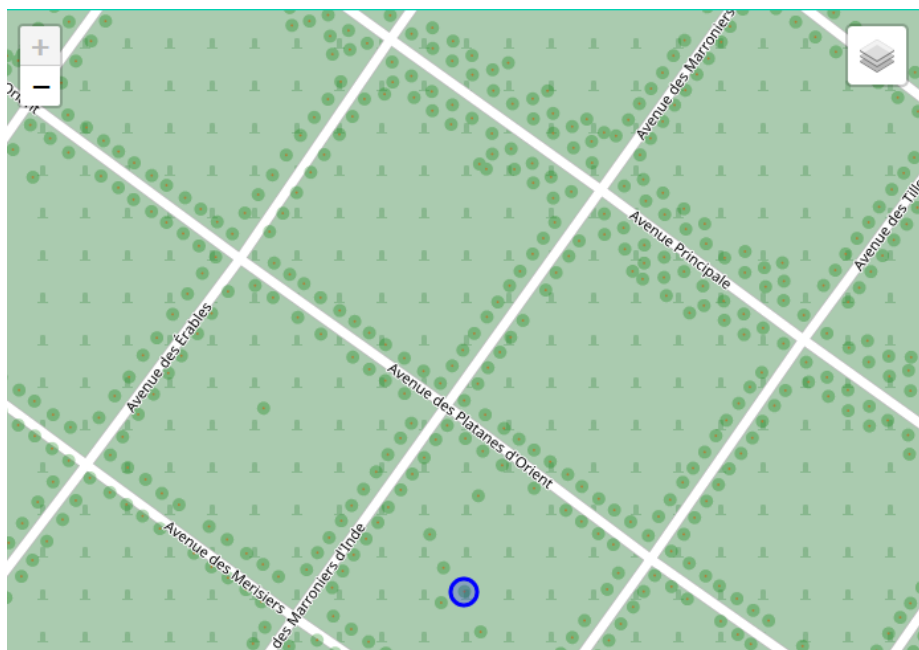
Les objets supprimés semblent être des nœuds de route

Ce contributeur est globalement fiable.

Contributeur U21

nombre de contributions	% création	% modification	% suppression	% objets réutilisés par d'autres	% d'objets qui n'ont pas été édités	% d'objets qui n'ont pas été supprimés
558	0.85	0.03	0.12	0.005	0.99	0.99

Ce contributeur a ajouté de la végétation dans la zone d'Aubervilliers, notamment des arbres sur le cimetière de Pantin. Il a surtout produit des objets de type node. Par la suite, le cimetière de Pantin a été mis à jour et certaines de ses contributions ont été modifiées ou supprimées.



Version	1 <input type="button" value="export"/>
Time	November 5, 2013 4:09 PM
Changeset	U21
User	
Lat	48.9045454
Lon	2.4053539
Tags	
natural	tree

D'après le changeset d'un contributeur ayant supprimé les objets de ce contributeur, on comprend que cette suppression vient du fait que ce dernier a ajouté des doublons : il a créé des éléments qui existaient déjà dans la base de données. On note également que ces 558 contributions sont issues d'un seul changeset : cela nous laisse à penser que cet utilisateur a utilisé un bot sans se préoccuper de la qualité de son import.

Groupe de modifications : 18738842 ✕

Duplicate nodes in way corrected

Fermé il y a plus de 4 ans par FVGordon

Attributs

created_by	JOSM/1.5 (6238 de)
------------	--------------------

Discussion

[Se connecter pour rejoindre la discussion](#)

Chemins (5)

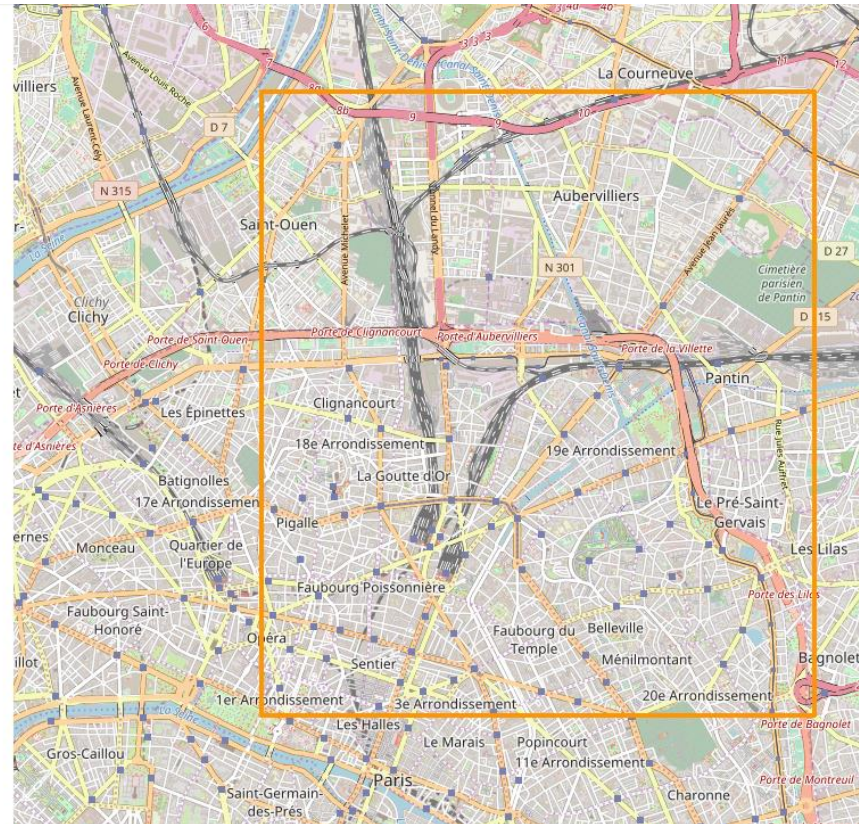
- Avenue Paul Vaillant Couturier (50406141, v8)
- Métro 7 (49993490, v15)
- Place du 8 Mai 1945 (50381224, v12)
67175462, v2
67178554, v3

Relations (2)

- 1061677, v2
- 1061684, v2

Nœuds (7)

- 2517989484, v2
- 640167143, v5
- 640167140, v6
- 640165713, v6
- 599915661, v6
- 2283276716, v2
- 2283276457, v3

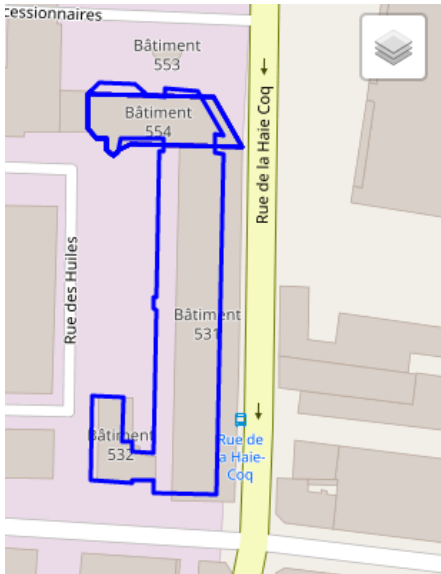


Ce contributeur n'est pas fiable.

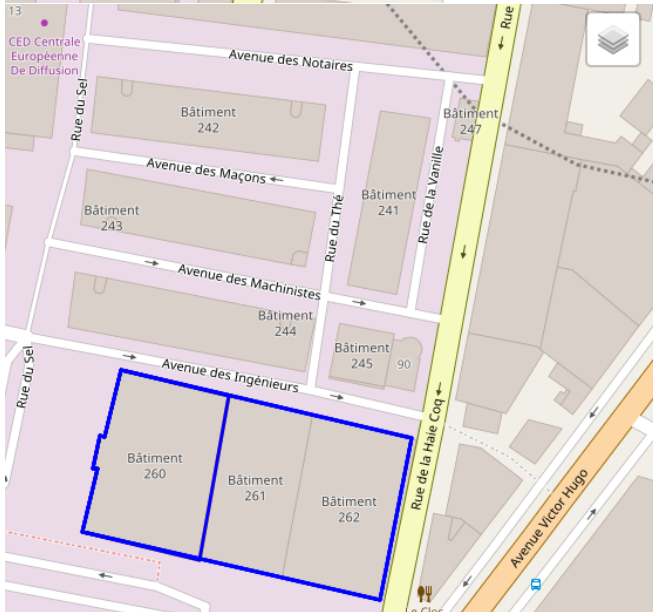
Contributeur U22

nombre de contributions	% création	% modification	% suppression	% objets réutilisés par d'autres	% d'objets qui n'ont pas été édités	% d'objets qui n'ont pas été supprimés
618	0.17	0.41	0.42	0.09	0.98	0.99

Ce contributeur a notablement séparé des groupes de bâtiments saisis à partir du cadastre : il a ajouté des noms et des adresses à ces bâtiments



Version	1	export	2	export
Time	August 18, 2010 11:13 AM		September 5, 2016 11:18 AM	
Changeset			U22	
User				
Tags				
addr:city			Aubervilliers	
addr:housenumber			52	
addr:postcode			93300	
addr:street			Rue de la Haie Coq	
building	yes		industrial	
name			Bâtiment 254	
source	extraction vectorielle v1 cadastre-dgi-fr source : Direction Générale des Impôts - Cadas. Mise à jour : 2010		extraction vectorielle v1 cadastre-dgi-fr source : Dir	



Version	3	export	4	export
Time	January 4, 2016 2:45 PM		September 4, 2016 8:23 PM	
Changeset			U22	
User				
Tags				
addr:city			Aubervilliers	
addr:housenumber				
addr:postcode			93300	
addr:street			Avenue des Chemins de Fer Industriels	
building	yes		industrial	
name	Bâtiment n°260-261-262		Bâtiment 260	
source	cadastre. Mise à jour : 2010		cadastre-dgi-fr source : Direction Générale des Impôts - Cadastre. Mise à jour : 2010	

Ce contributeur a également mis à jour le réseau routier ainsi que des noms de rue.

The map shows a street network with a blue line highlighting 'Rue du Grain'. The table below tracks changes to this street across six versions.

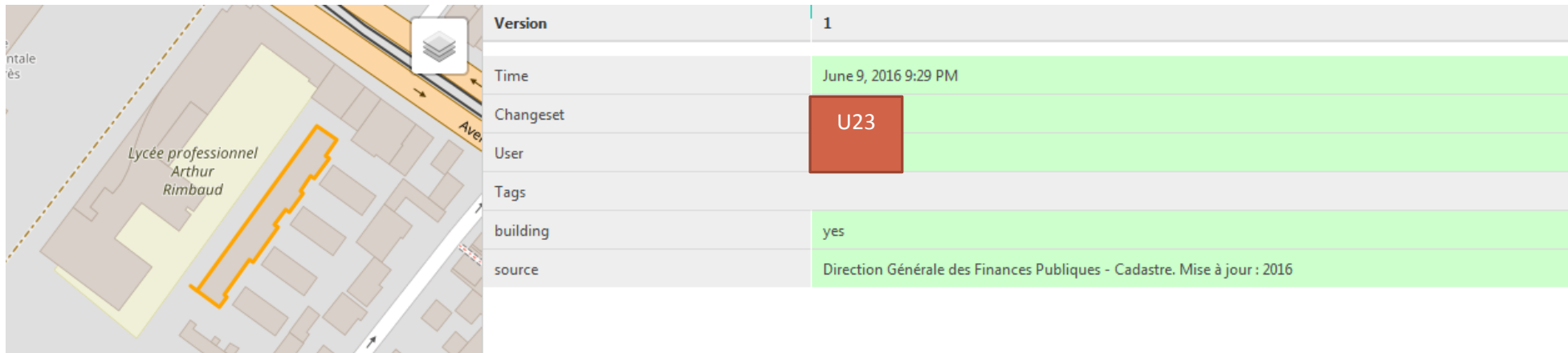
Version	4	5	6	7	8	9
Time	June 28, 2016 9:00 PM	September 5, 2016 12:15 PM	September 5, 2016 12:16 PM	September 5, 2016 12:17 PM	June 4, 2017 10:42 AM	October 6, 2017 10:42 AM
Changeset		U22	U22	U22	U22	
User						
Tags						
FIXME	check the oneway direction	check the oneway direction	check the oneway direction	check the oneway direction	check the oneway direction	check the oneway direction
highway	service	service	service	service	service	service
name	Avenue des Forgerons	Avenue des Forgerons	Rue du Grain	Rue du Grain	Rue du Grain	Rue du Grain
oneway	yes;-1	yes;-1	yes;-1	yes;-1	yes;-1	yes;-1

On peut dire que ce contributeur est fiable

Contributeur U23

nombre de contributions	% création	% modification	% suppression	% objets réutilisés par d'autres	% d'objets qui n'ont pas été édités	% d'objets qui n'ont pas été supprimés
777	0.6	0.1	0.3	0.04	0.99	1

Ce contributeur a lancé des mini bots sur des zones localisées à Aubervilliers et peu de ses contributions ont été éditées, aucune n'a été supprimée. Il utilise le cadastre (versions 2014 à 2016) pour mettre à jour les données. Il supprime le bâti qui n'existe plus dans les versions ultérieures du cadastre.



Version	1
Time	June 9, 2016 9:29 PM
Changeset	U23
User	
Tags	
building	yes
source	Direction Générale des Finances Publiques - Cadastre. Mise à jour : 2016

Il répare également la géométrie de certains bâtis et ajoute des noms aux rues créées par Esperanza86.

Version	1	export	2	export
Time	August 18, 2010 12:23 PM		January 19, 2015 5:42 PM	
Changeset	[Red Box]		U23	
User	[Red Box]		[Red Box]	
Tags				
building	yes		yes	
source	extraction vectorielle v1 cadastre-dgi-fr source : Direction Générale des Impôts - Cadas. Mise à jour : 2010		extraction vectorielle v1 cadastre-dgi-fr source : Direction G	

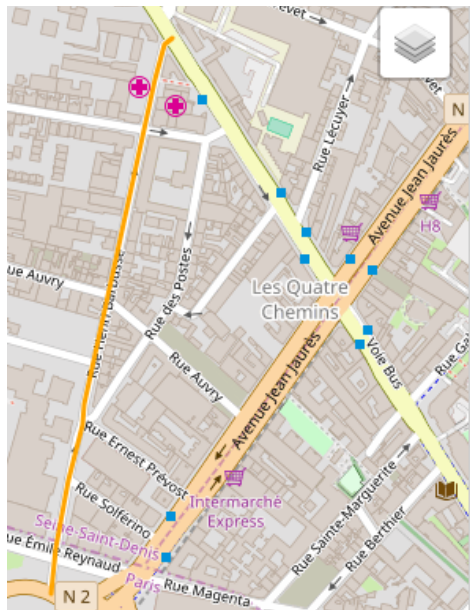
Version	2
Time	August 13, 2014 11:25 AM
Changeset	U23
User	[Red Box]
Tags	
highway	residential
name	Rue Bengali

Ce contributeur peut être considéré comme fiable.

Contributeur U24

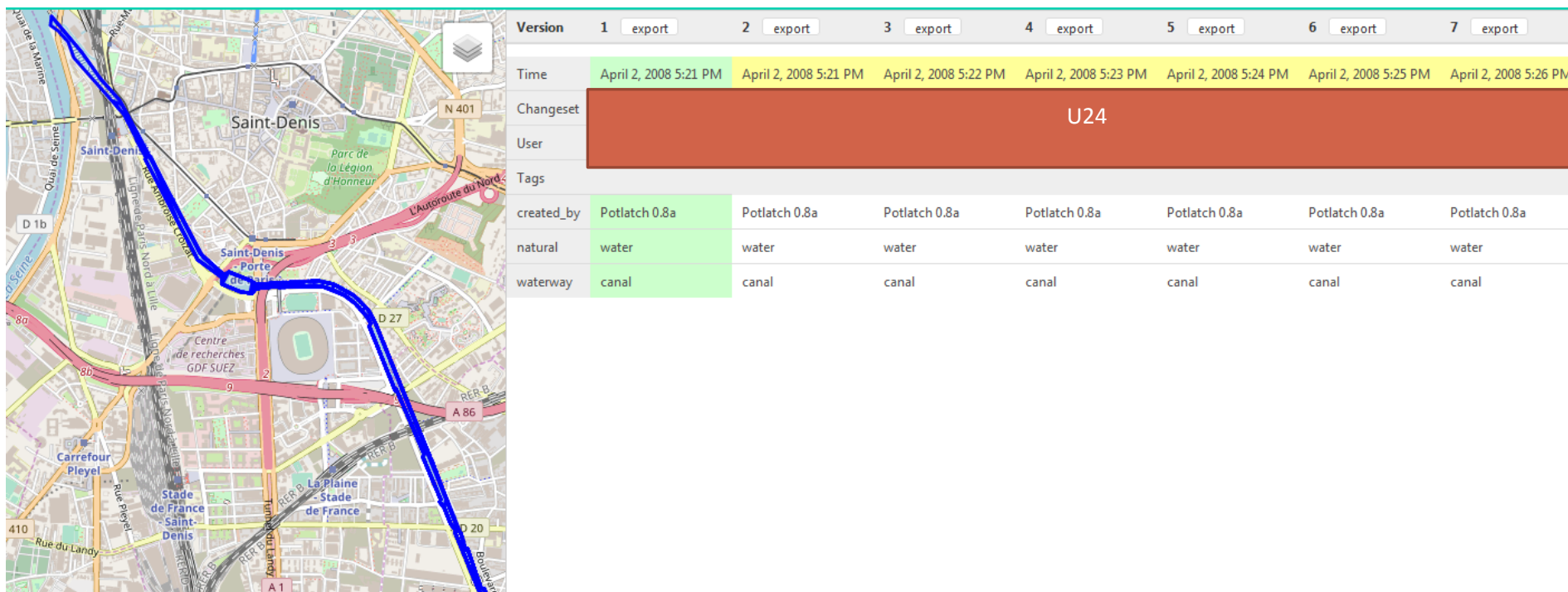
nombre de contributions	% création	% modification	% suppression	% objets réutilisés par d'autres	% d'objets qui n'ont pas été édités	% d'objets qui n'ont pas été supprimés
870	0.82	0.15	0.03	0.18	0.98	0.99

Ce contributeur a produit un grand nombre de routes en 2008, qui ont été par la suite améliorées par d'autres contributeurs, mais sans être supprimées. En effet, les routes qu'il crée sont pauvres en attributs.



Version	1	2	3	4	5	6
Time	April 7, 2008 11:58 PM	April 7, 2008 11:58 PM	June 20, 2008 11:27 PM	December 20, 2008 4:51 PM	February 17, 2009 3:39 PM	February 21, 2009 12:17 PM
Changeset	U24	U24				
User						
Tags						
created_by	Potlatch 0.8a	Potlatch 0.8a	Potlatch 0.9c	Potlatch 0.9c	Potlatch 0.10f	Potlatch 0.10f
highway	unclassified	unclassified	unclassified	unclassified	unclassified	unclassified
maxspeed						
name				Rue Henri Barbusse	Rue Henri Barbusse	Rue Henri Barbusse
oneway					yes	yes

Les modifications effectuées portent souvent sur ses propres contributions : correction géométrique des tronçons de routes ajoutées. De même que pour les suppressions : il supprime ses propres contributions



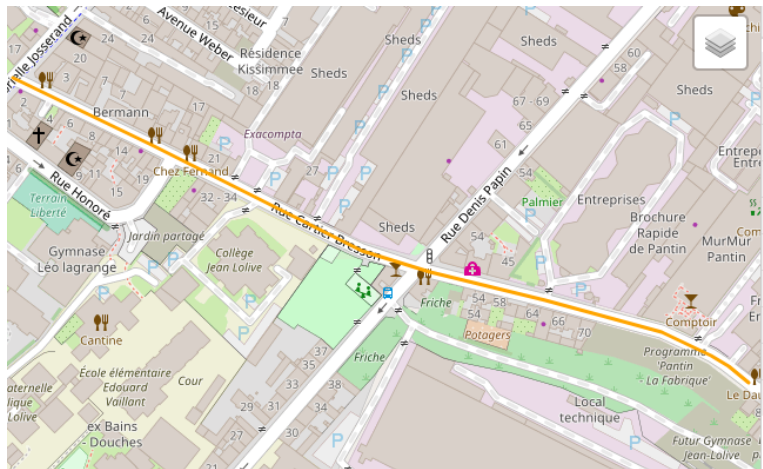
On observe ci-dessus que le contributeur a modifié à plusieurs reprises sa contribution. Finalement, il supprime cet élément. Ce contributeur semble donc peu sûr de lui. **On ne peut rien dire de la fiabilité de ce contributeur.**

Contributeur U25

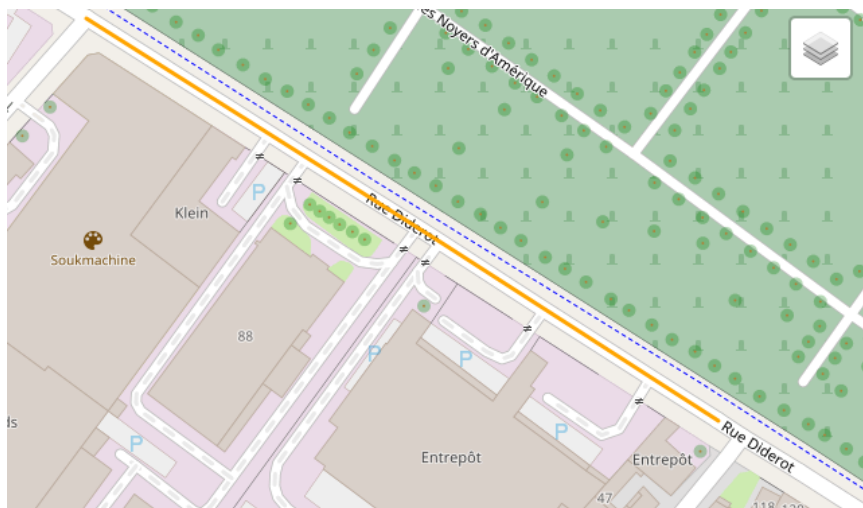
nombre de contributions	% création	% modification	% suppression	% objets réutilisés par d'autres	% d'objets qui n'ont pas été édités	% d'objets qui n'ont pas été supprimés
1014	0.39	0.46	0.15	0.23	0.98	0.99

Ce contributeur a participé entre 2012 et 2017.

Il a amélioré les routes en modifiant leur géométrie et leurs attributs, notamment en ajoutant des informations locales telles que des interdictions de passage de bus, des limitations de vitesse, etc.



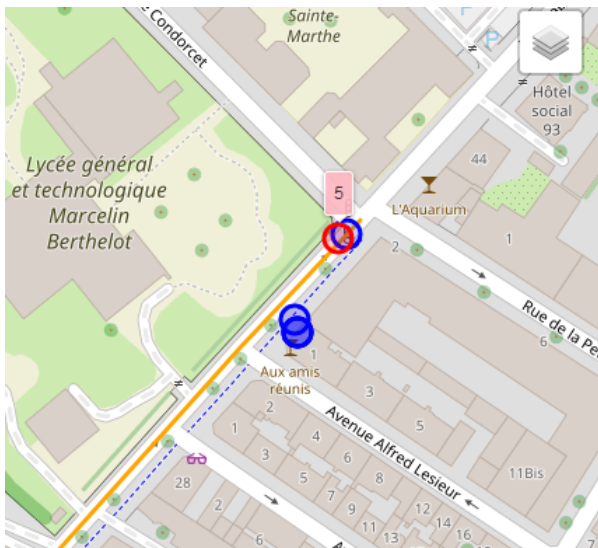
Version	25	26	27	
Time	3 PM	October 26, 2015 6:36 PM	March 15, 2016 4:24 PM	July 11, 2016 5:18 PM
Changeset	U25			
User	U25			
Tags	U25			
created_by	U25			
highway	residential	residential	residential	
maxspeed	50	50	50	
name	Rue Cartier Bresson	Rue Cartier Bresson	Rue Cartier Bresson	
prohibited	trucks (except local traffic), trailers, and tourist buses prohibited	trucks (except local traffic), trailers, and tourist buses prohibited	trucks (except local traffic), trailers, and tourist buses prohibited	
surface	asphalt	asphalt	asphalt	



Version	1	2
Time	December 8, 2012 2:31 AM	May 8, 2014 4:06 AM
Changeset	U25	U25
User	U25	U25
Tags	U25	U25
highway	unclassified	residential
maxspeed		50
name	Rue Diderot	Rue Diderot
oneway	no	no
surface		asphalt


On observe donc un comportement de maintenance chez ce contributeur.

Il ramène des nœuds de route près des intersections, pour indiquer la présence de feux tricolores.



Version	4 export	5 export	6 export
Time	February 3, 2011 4:21 PM	December 23, 2012 3:42 PM	October 24, 2015 7:11 PM
Changeset		U25	U25
User			
Lat	48.9046335	48.904835	48.9048458
Lon	2.3958031	2.3959713	2.3960013
Tags			
crossing			zebra
highway			traffic_signals
highway_1			crossing

Il cartographie également les dos d'ânes/zebra-crossings, les écoles, et il enrichit les bouches de métro ajoutées par le contributeur fiable U30 en indiquant la présence d'escalators.



Version	export	2 export
Time		November 21, 2017 2:53 PM
Changeset		U25
User		
Lat		48.9039269
Lon		2.3921173
Tags		
railway		subway_entrance
subway_entrance		escalator

Ce contributeur semble mettre à contribution sa connaissance locale sur cet espace : **ce contributeur est donc fiable.**

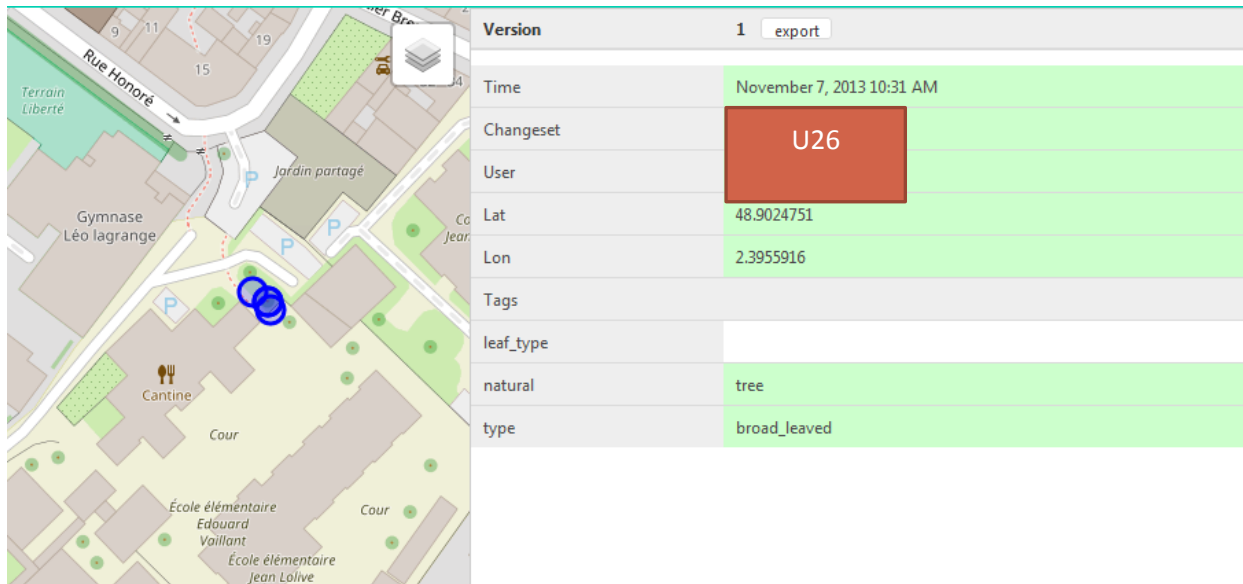
Contributeur U26

nombre de contributions	% création	% modification	% suppression	% objets réutilisés par d'autres	% d'objets qui n'ont pas été édités	% d'objets qui n'ont pas été supprimés
1360	0.7	0.1	0.2	0.06	0.99	0.99

Ce contributeur a notablement participé en 2013 puis en 2015.

Création

Il a amélioré la cartographie des espaces verts dans la zone du gymnase Léo Lagrange et le Collège Jean Lolive, en ajoutant beaucoup d'arbres.



Version	1	export
Time	November 7, 2013 10:31 AM	
Changeset	U26	
User	[redacted]	
Lat	48.9024751	
Lon	2.3955916	
Tags		
leaf_type		
natural	tree	
type	broad_leaved	

Il ajoute de nouveaux bâtiments selon la mise à jour du cadastre en 2015. Ces bâtiments n'ont pas été édités par d'autres contributeurs. Il indique aussi la présence de zones de construction.

Groupe de modifications :

Ajouts nouveaux bat selon māj cadastre

Fermé il y a plus de 2 ans par **U26**

Attributs

created_by JOSM/1.5 (9060 fr)

Discussion [S'abonner](#)

Commentaire

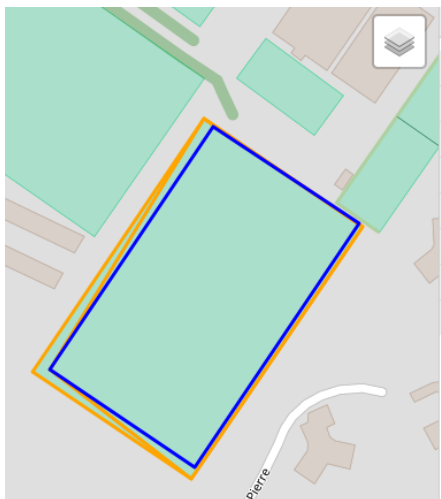
Chemins (1 à 20 sur 27)

12



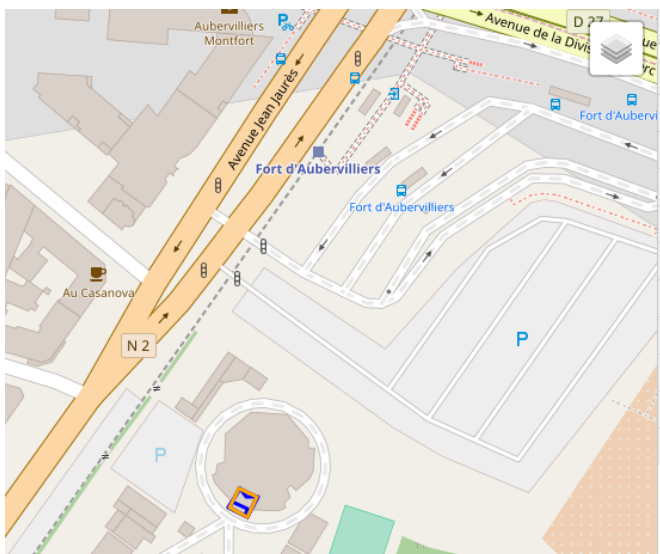
Modification

Il ajoute des tags aux objets existants.



Version		2	export
Time		November 5, 2013 6:26 PM	
Changeset		U26	
User			
Tags			
leisure		pitch	
source		Bing	
sport		soccer	
surface		compacted	

Il ajoute de la précision aux bâtiments importés massivement par le contributeur fiable U30 en ajoutant le tag « building :part » ainsi que des informations sur la taille du bâti, la forme du toit. **Ce contributeur est fiable.**

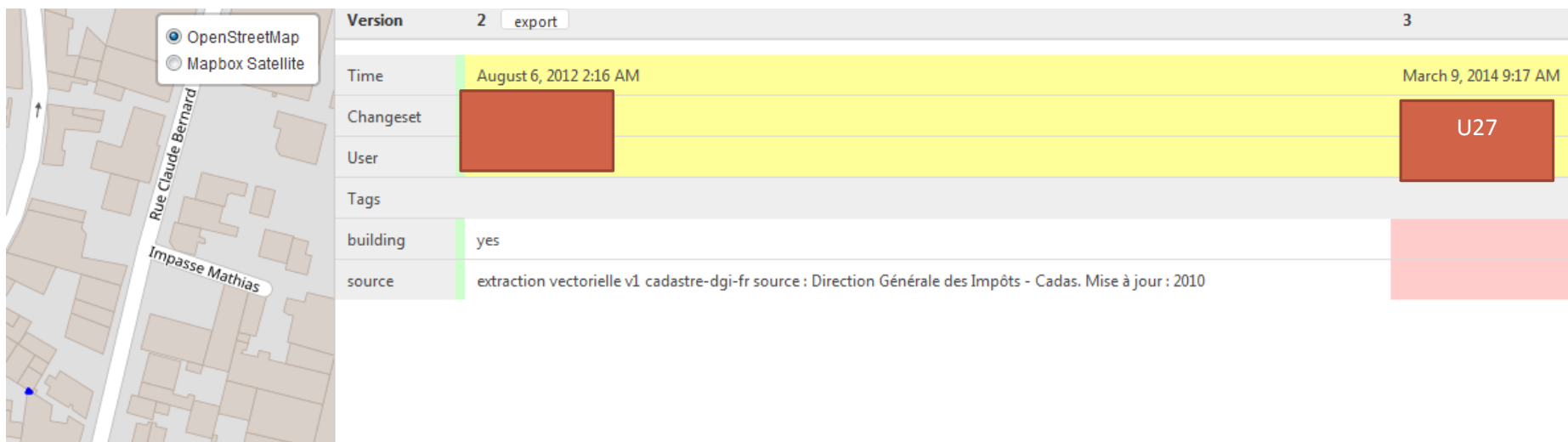


Version		1	export	2	export
Time		August 18, 2010 11:54 AM		November 26, 2013 4:57 PM	
Changeset		U30		U26	
User					
Tags					
building		yes			
building:part				yes	
height				6	
roof:height					
roof:shape				gabled	
source		extraction vectorielle v1 cadastre-dgi-fr source : Direction Générale des Impôts - Cadas. Mise à jour : 2010		extraction vectorielle v1 cadastre-dgi-fr source : Direction	
wall		no		no	

Contributeur U27

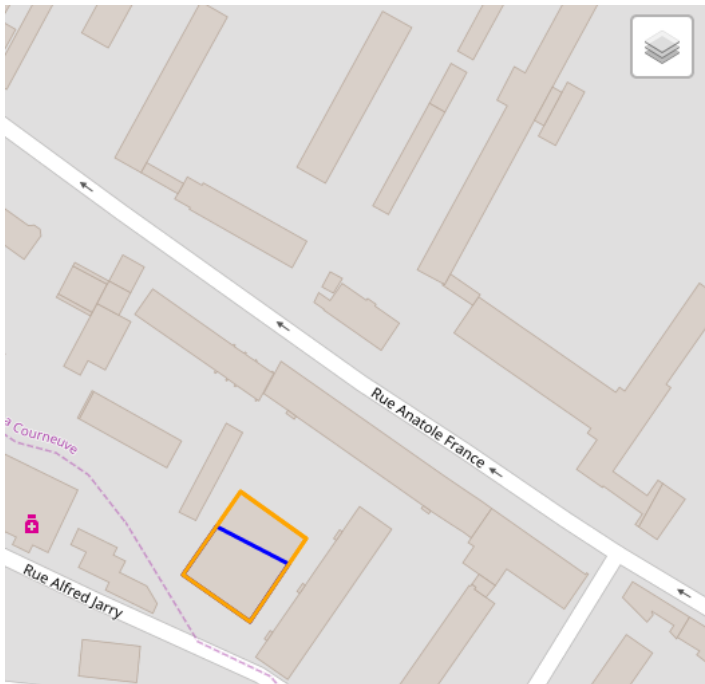
nombre de contributions	% création	% modification	% suppression	% objets réutilisés par d'autres	% d'objets qui n'ont pas été édités	% d'objets qui n'ont pas été supprimés
1559	0.04	0.38	0.58	0.1	0.99	0.99

Ce contributeur a notamment supprimé des objets. En effet, sa contribution a principalement pour but de fusionner les petits morceaux de bâtiments importés depuis le cadastre par le contributeur fiable U30.



The screenshot shows a map on the left with a blue dot indicating a location on Rue Claude Bernard, near Impasse Mathias. On the right, a comparison of two versions of a building object is displayed:

Version	2	export	3
Time	August 6, 2012 2:16 AM		March 9, 2014 9:17 AM
Changeset	[Redacted]		U27
User	[Redacted]		[Redacted]
Tags			
building	yes		[Redacted]
source	extraction vectorielle v1 cadastre-dgi-fr source : Direction Générale des Impôts - Cadas. Mise à jour : 2010		[Redacted]



Version 2 [export](#)

Time	March 9, 2014 9:16 AM
Changeset	U27
User	
Tags	
building	yes
source	extraction vectorielle v1 cadastre-dgi-fr source : Direction Générale des Impôts - Cadas. Mise à jour : 2010



Version 1 [export](#)

Time	March 9, 2014 9:32 AM
Changeset	U27
User	
Tags	
building	yes
source	bing

Il a aussi saisi des bâtiments depuis Bing, et ces bâtiments n'ont pas été édités par la suite.

Globalement, ce contributeur est fiable.

Contributeur U28

nombre de contributions	% création	% modification	% suppression	% objets réutilisés par d'autres	% d'objets qui n'ont pas été édités	% d'objets qui n'ont pas été supprimés
2597	0.3	0.2	0.5	0.5	0.98	0.99

Ce contributeur participe via des cartoparties, durant lesquelles il ajoute beaucoup de données locales, telles que des parkings, et des routes nommées. Il cartographie aussi les aménités : boulangeries, garages, passages piétons. Il modifie les arrêts de tram.

Groupe de modifications :

St Denis - cartopartie du 16 juin

Fermé il y a environ 6 ans par **U28**

Attributs

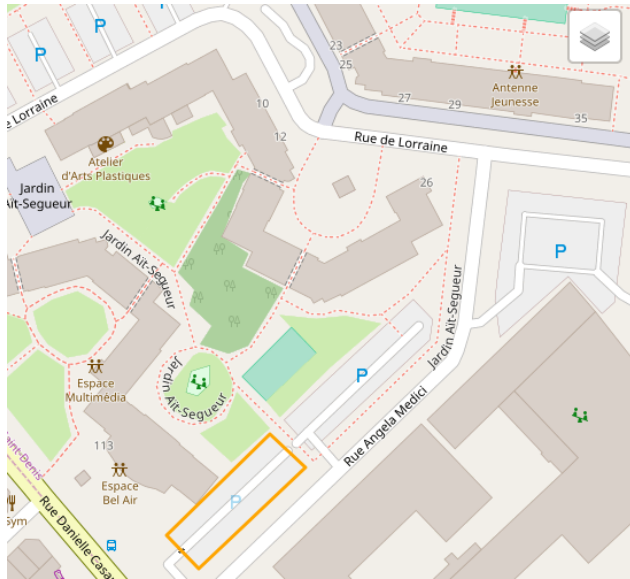
created_by JOSM/1.5 (5267 fr)

Discussion [S'abonner](#)

Commentaire

Chemins (17)

The screenshot shows a map editing interface. On the left, there is a sidebar for a group of modifications. The group is titled 'Groupe de modifications :'. Below the title, it says 'St Denis - cartopartie du 16 juin' and 'Fermé il y a environ 6 ans par U28'. There are sections for 'Attributs' (showing 'created_by JOSM/1.5 (5267 fr)'), 'Discussion' (with a 'S'abonner' button), 'Commentaire', and 'Chemins (17)'. The main part of the image is a map of a residential area in St Denis, France. The map shows streets such as Rue de Lorraine, Rue Danielle Casanova, Rue Albert Walter, Rue Bernard Palissy, Rue Pierre Curie, Rue Francis de Pressensé, and Rue Angela Medici. There are also labels for 'Jardin Alt-Segueur', 'Atelier d'Arts Plastiques', 'Espace Multimédia', 'Espace Bel Air', and 'Pharmacie du Grand Canal'. A red dashed line outlines a specific area on the map, and a blue box highlights a specific location. The map also shows a canal, 'Canal du Canal Saint-Denis', and a tram line with stops marked 'P'.

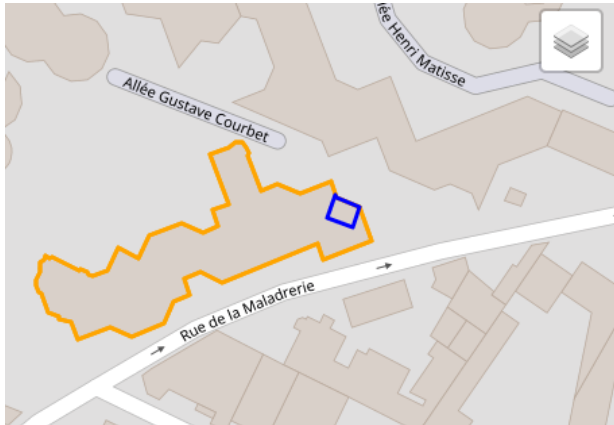


Version	1
Time	June 16, 2012 5:41 PM
Changeset	U28
User	
Tags	
access	private
amenity	parking
parking	surface




Version	1
Time	June 14, 2012 2:00 PM
Changeset	U28
User	
Tags	
highway	tertiary
name	Rue Francis de Pressensé
oneway	yes
ref	D 30
source	cadastre-dgi-fr source : Direction Générale des Impôts - Cadastre. Mise à jour : 2012

Il effectue également des fusions de bâti importé du cadastre :



Version	1	export	2	export
Time	August 18, 2010 11:43 AM		May 1, 2014 11:43 AM	
Changeset			U28	
User				
Tags				
building	yes		yes	
source	extraction vectorielle v1 cadastre-dgi-fr source : Direction Générale des Impôts - Cadas. Mise à jour : 2010		extraction vect	

Et il fait des mises à jour sur les bâtis qui disparaissent sur les zones de construction.



Version	2	export	3
Time	August 13, 2012 8:39 AM		May 1, 2014 6:06 PM
Changeset			U28
User			
Tags			
building	yes		
source	cadastre-dgi-fr source : Direction Générale des Impôts - Cadastre. Mise à jour : 2010		

Ce contributeur est fiable

Contributeur U29

nombre de contributions	% création	% modification	% suppression	% objets réutilisés par d'autres	% d'objets qui n'ont pas été édités	% d'objets qui n'ont pas été supprimés
5380	0.7	0.2	0.1	0.17	0.99	0.99

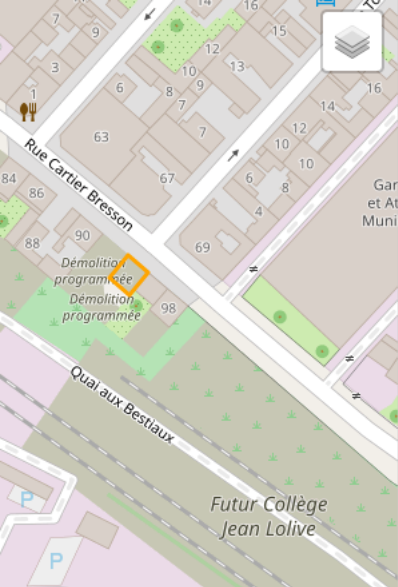
Ce contributeur ajoute des données locales : jardins, rues à double sens...



The map shows a street grid with buildings and green spaces. A specific area is highlighted in green, indicating a garden. The map includes street names like 'Marie Louise' and 'Rue Toffier Decaux'. A legend icon is visible in the top right corner of the map area.

Version	1
Time	February 27, 2016 1:13 PM
Changeset	U29
User	
Tags	
leisure	garden

Il fait un réel suivi de l'évolution géographique de l'espace.




Version	5
Time	July 13, 2016 6:03 PM
Changeset	U29
User	
Tags	
addr:city	Pantin
addr:house:number	94
addr:postcode	93500
addr:street	Rue Cartier Bresson
building	
building:levels	0
landuse	construction
name	Démolition programmée
source	tion Générale des Impôts - Cadas. Mise à jour : 2010 extraction vectorielle v1 cadastre-dgi-fr source : Direction Générale des Impôts - Cadas. Mise à jour : 2010

Ce contributeur est donc fiable.

Contributeur U30

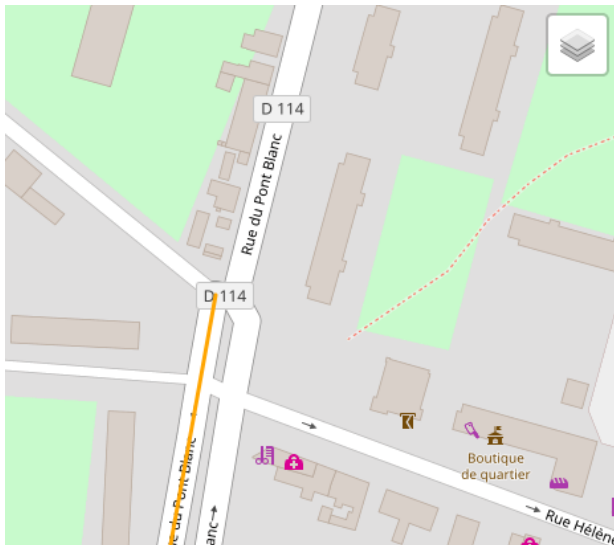
nombre de contributions	% création	% modification	% suppression	% objets réutilisés par d'autres	% d'objets qui n'ont pas été édités	% d'objets qui n'ont pas été supprimés
93657	0.96	0.03	0.01	1	0.99	0.99

Ce contributeur est principalement à l'origine de l'import massif de données du cadastre : en effet, en 2010 un centaine de bâtiments sont ajoutés d'un coup. Mais on note tout de même qu'il édite ses propres données.




Version	1	2	3	4	5
Time	February 20, 2009 8:34 PM	February 20, 2009 8:34 PM	April 6, 2011 6:57 PM	March 11, 2013 1:58 PM	March 11, 2013 1:59 PM
Changeset	U30	U30			
User					
Tags					
created_by	Potlatch 0.10f	Potlatch 0.10f			
highway	unclassified	unclassified	unclassified	unclassified	unclassified
maxspeed				30	30
name		Rue de Presles	Rue de Presles	Rue de Presles	Rue de Presles
surface					

Il édite également les routes produites par d'autres contributeurs, en les corrigeant, en y ajoutant des attributs...



Version	3	4	5	6	7	8
Time	January 3, 2008 3:20 PM	December 30, 2008 7:34 PM	February 17, 2009 3:46 PM	February 17, 2009 3:46 PM	August 16, 2009 6:48 PM	September 6, 2009 10:00 AM
Changeset			U30	U30	U30	U30
User						
Tags						
created_by	JOSM	JOSM	Potlatch 0.10f	Potlatch 0.10f	Potlatch 0.10f	
cycleway						
highway	secondary	secondary	tertiary	secondary	tertiary	tertiary
name			Rue du Pont Blanc	Rue du Pont Blanc	Rue du Pont Blanc	Rue du Pont Blanc
oneway						
ref	D114	D114	D114	D114	D 114	D 114

Enfin il effectue aussi une sorte de maintenance sur ses propres contributions. Par exemple, ci-dessus il ajoute une note pour informer de l'état d'une route qu'il avait créée :



Version	14	15	16	17	18
Time	May 22, 2011 11:23 PM	March 9, 2012 12:40 PM	June 1, 2013 10:39 AM	September 13, 2014 7:59 AM	February 26, 2015 11:33 AM
Changeset	U30	U30	U30	U30	U30
User					
Tags					
access		no	no		
bicycle		no	no		permissive
created_by					
cycleway					
fixme					
foot		no	no		
highway	track	track	track	track	track
name	Quai des Vertus	Quai des Vertus	Quai des Vertus	Quai des Vertus	Quai des Vertus
note		fermé la semaine jusqu'en 2013	fermé la semaine jusqu'en 2013	fermé la semaine jusqu'en 2013	fermé la semaine jusqu'en 2013
surface	pebblestone			dirt	dirt
tracktype					

Donc ce contributeur est fiable

Annexe C

Ajout de profils synthétiques de carto-vandales

Profils de vandales synthétiques - Aubervilliers

Types de vandalisme

- A. Edition ou suppression des tags 'name' des bâtiments commerciaux, industriels et lieux de culte (44 bâtiments)
- B. Bâtiments tracés dans un cours d'eau: 'BIENVENUE' (9 bâtiments)
- C. Bâtiments imaginaires (15 contributions)
- D. Bâtiments de taille et/ou de forme invraisemblable (3 contributions)

Construction de profils synthétiques pour le vandalisme produit sur Aubervilliers

Contributeurs

	Contributeur A	Contributeur B	Contributeur C	Contributeur D
total_contributions (1)	500	9	60	3
p_creation	0.01	1	0.5	1
p_modif	0.79	0	0.5	0
p_delete	0.20	0	0	0
n_semesters ¹ (2)	3	1	2	1
p_is_used (3)	0	0	0	0
p_is_edited* (4)	0.25	1	0.5	1
p_is_deleted* (5)	0.3	1	0.1	1
nbWeeks ² (6)	8	1	3	1
focalisation ^{3*} (7)	0.3	1	0.5	1

* Ces variables ont été inversées de sorte qu'une valeur proche de 0 corresponde au pourcentage de données qui n'ont pas été supprimées ou éditées, par rapport au nombre de total de contributions par contributeur.

```
(1) > summary(users_aubervilliers$total_contributions)
```

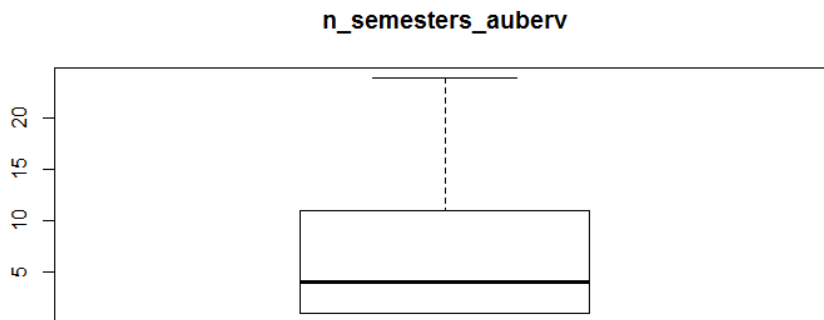
```
Min. 1st Qu. Median Mean 3rd Qu. Max.
  1      2      8    423    31  93657
```

¹ Calcul basé sur les dates des changesets qui contiennent des bâtiments. Pour avoir une distribution représentative du nombre de semestres contribués, il faut extraire un échantillon qui contient autant de lignes que de contributeurs (62 contributeurs de bâtiments dans Stuhr) et non pas autant de lignes que de bâtiments (6274)

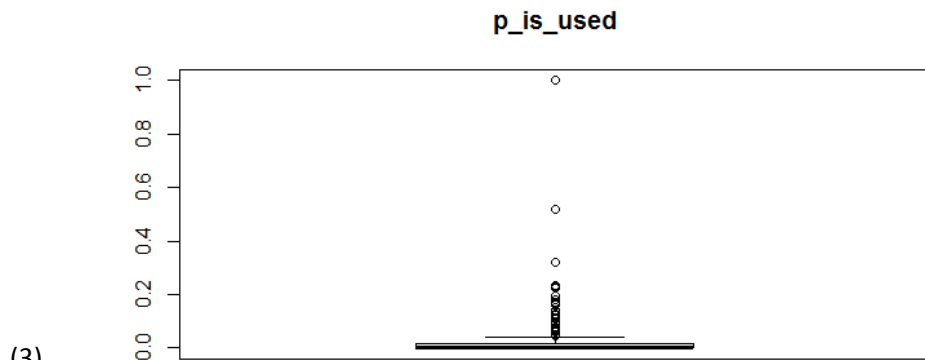
² Calcul basé sur les ways et les nodes contribués depuis l'an 2000 des contributeurs de la zone (Total : 378).

³ Calcul basé sur la comparaison entre la surface de la bbox spatio-temporelle utilisée pour charger les données et la surface moyenne des changesets des contributeurs. 0 : changeset très peu focalisé sur la zone d'étude, 1 et plus : changeset très focalisé (Total : 378)

(2) Distribution de n_semesters :



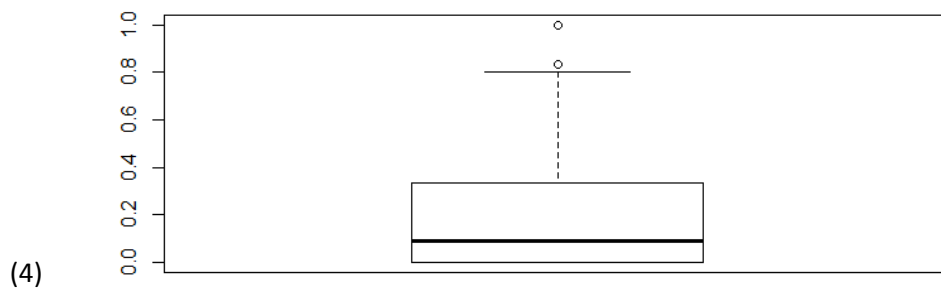
```
> users_n_semesters <- dbGetQuery(con, "SELECT DISTINCT ON (uid) uid, n_semesters_auberv from indicators.aubervilliers")
> summary(users_n_semesters$n_semesters_auberv)
  Min. 1st Qu.  Median    Mean 3rd Qu.    Max.
 1.000  1.000   4.000   6.887 11.000  24.000
```



(3)

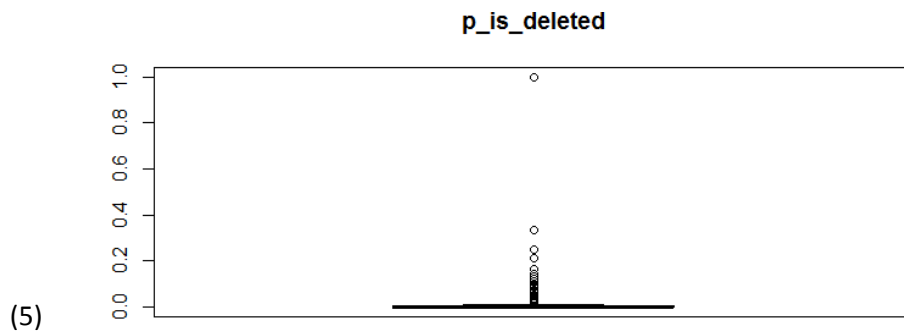
```
> summary(users_aubervilliers$p_is_used)
  Min. 1st Qu.  Median    Mean 3rd Qu.    Max.
0.00000 0.00000 0.00000 0.02479 0.01724 1.00000
```

p_is_edited

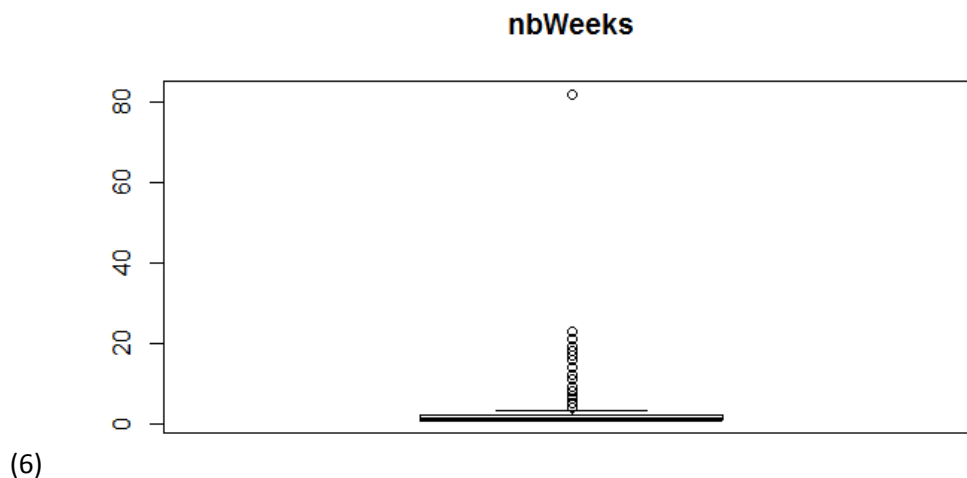


(4)

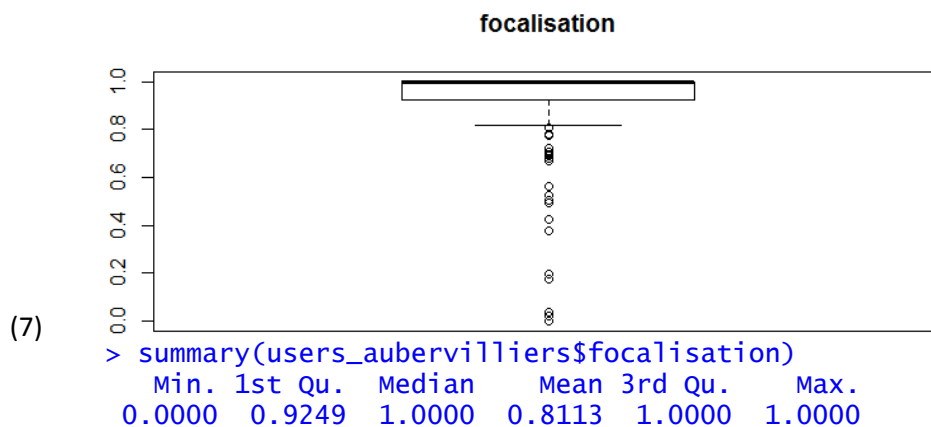
```
> summary(users$p_is_edited)
  Min. 1st Qu.  Median    Mean 3rd Qu.    Max.
0.00000 0.00000 0.09091 0.24317 0.33333 1.00000
```



```
> summary(users$p_is_deleted)
  Min. 1st Qu.  Median    Mean 3rd Qu.    Max.
0.000000 0.000000 0.000000 0.037870 0.006239 1.000000
```



```
> summary(users_aubervilliers$nbweeks)
  Min. 1st Qu.  Median    Mean 3rd Qu.    Max.
 1.000  1.000  1.000  2.608  2.000  82.000
```



```
> summary(users_aubervilliers$focalisation)
  Min. 1st Qu.  Median    Mean 3rd Qu.    Max.
0.0000 0.9249 1.0000 0.8113 1.0000 1.0000
```

Changesets

Pour simplifier, on assigne un changeset par groupe de vandalisme : changeset A pour le groupe de vandalisme A, changeset B pour le groupe de vandalisme B, etc.

On assimile aussi l'âge de la contribution à l'âge du changeset.

Ce changeset peut éventuellement contenir d'autres contributions.

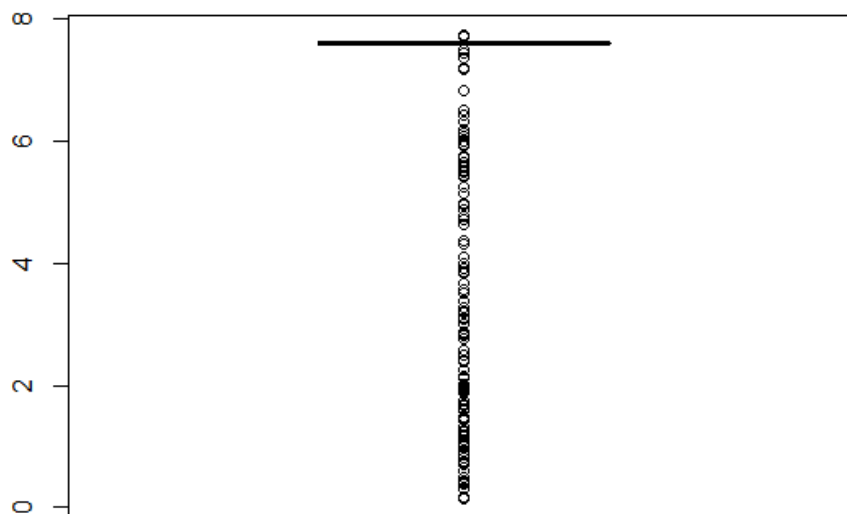
	Changeset A	Changeset B	Changeset C	Changeset D
nb éditions dans le changeset (1)	60	9	15	3
durée (en secondes) du changeset	1000	300	600	120
Âge de la contribution	31536000s (1 an)	126144000s (4 ans)	157680000s (5 ans)	15552000s (6 mois)

```
(1) > summary(aubervilliers$chgst_edits)
Min. 1st Qu. Median Mean 3rd Qu. Max.
  3  29229  29229 30826  29229  86405
```

```
(2) > summary(aubervilliers$chgst_duration_seconds)
Min. 1st Qu. Median Mean 3rd Qu. Max.
  0  5062  5062  4323  5062 12992
```

```
(3) > summary(aubervilliers$lifespan_seconds) #En secondes
Min. 1st Qu. Median Mean 3rd Qu. Max.
4324307 236387246 236433296 213988342 236435078 241094808
> summary(aubervilliers $lifespan_seconds/3600/24) # En jours
Min. 1st Qu. Median Mean 3rd Qu. Max.
 50.05 2735.96 2736.50 2476.72 2736.52 2790.45
```

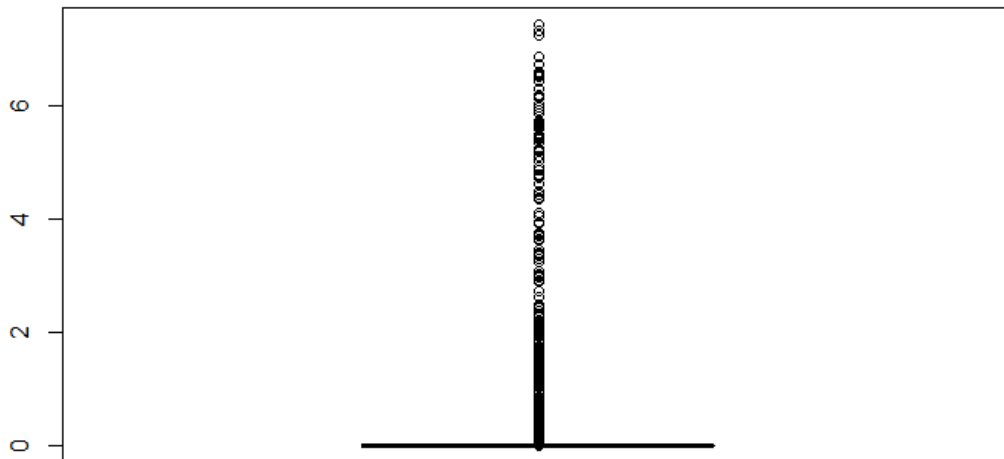
Age (Années)



Date

```
summary(aubervilliers$timespan_to_previous/3600/24/30/12)
Min. 1st Qu. Median Mean 3rd Qu. Max.
0.0000 0.0000 0.0000 0.4958 0.0000 7.4134
```

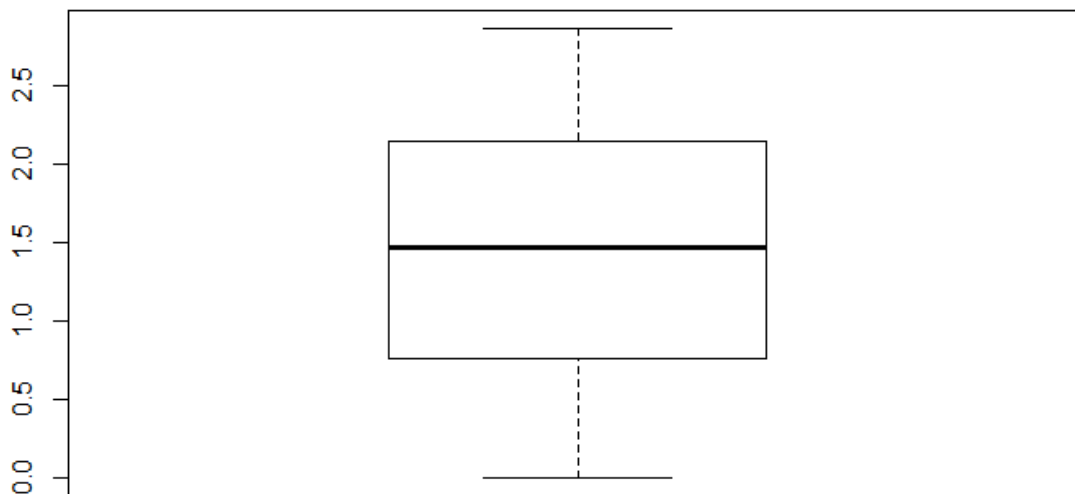
timespan to previous (années)



En ayant filtré les objets dont $v_contrib > 1$:

```
> summary(which(aubervilliers$timespan_to_previous>0)/3600)
  Min. 1st Qu.  Median    Mean 3rd Qu.    Max.
0.001667 0.761111 1.463611 1.446634 2.147500 2.859444
```

timespan to previous (heures)



Pour calculer `timespan_to_previous` : `changeset A` a été produit le 7 décembre 2017

Profils de vandales synthétiques - Stuhr

Types de vandalisme

- E. Edition des tags 'name' des bâtiments commerciaux, et modification de leurs valeurs par des smileys (9 bâtiments)
- F. Bâtiments tracés dans un étang : 'HAPPY NEW YEAR' (12 bâtiments)
- G. Bâtiments imaginaires appartenant à la ville 'Salt Lake Kingdom' (20 contributions)
- H. Bâtiments de taille et/ou de forme invraisemblable (3 contributions)

Construction de profils synthétiques pour le vandalisme produit sur Stuhr

Contributeurs

On assigne un contributeur par type de vandalisme : contributeur A pour le groupe de vandalisme A, contributeur B pour le groupe de vandalisme B, etc. (cf. fichier `stuhr_users_metrics.R` pour retrouver les plots).

	Contributeur A	Contributeur B	Contributeur C	Contributeur D
total_contributions (1)	10000 (très grand contributeur)	100 (grand contributeur)	30 (contributeur moyen)	3 (petit contributeur)
p_creation	0.33	0.5	1	1
p_modif	0.33	0.5	0	0
p_delete	0.33	0	0	0
n_semesters ⁴ (2)	5	10	1	1
p_is_used (3)	0.2	0.001	0	0
p_is_edited (4)	0.7	0.2	0.66	0
p_is_deleted (5)	0.99	0.9	0.99	1
nbWeeks ⁵ (6)	10	10	1	1
focalisation ⁶ (7)	0.5	0	1	1

Les valeurs des variables sont attribuées au regard de la distribution des valeurs (**non-centrées réduites**) des contributeurs

```
(8) > summary(stuhr$total_contributions)
```

```
Min. 1st Qu.  Median    Mean 3rd Qu.    Max.
1.00   2.00   6.00  292.33  28.75 40240.00
```

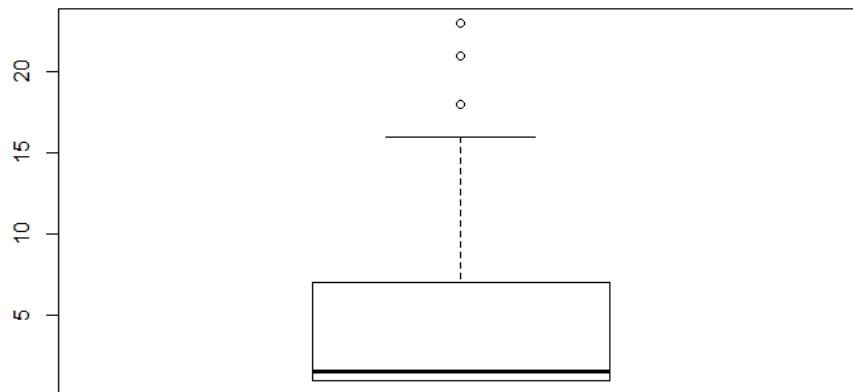
⁴ Calcul basé sur les dates des changesets qui contiennent des bâtiments. Pour avoir une distribution représentative du nombre de semestres contribués, il faut extraire un échantillon qui contient autant de lignes que de contributeurs (62 contributeurs de bâtiments dans Stuhr) et non pas autant de lignes que de bâtiments (6274)

⁵ Calcul basé sur les ways et les nodes contribués depuis l'an 2000 des contributeurs de la zone (Total : 378).

⁶ Calcul basé sur la comparaison entre la surface de la bbox spatio-temporelle utilisée pour charger les données et la surface moyenne des changesets des contributeurs. 0 : changeset très peu focalisé sur la zone d'étude, 1 et plus : changeset très focalisé (Total : 378)

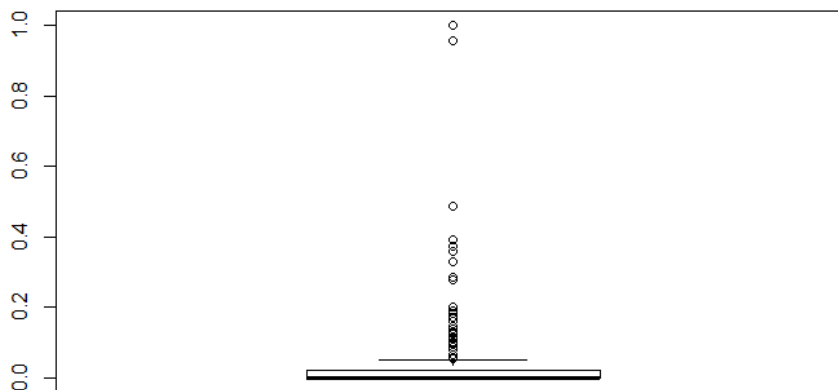
(9) Distribution de n_semesters :

Nombre de semestres



```
> users_n_semesters <- dbGetQuery(con, "SELECT DISTINCT ON (uid)
uid, n_semesters_bremen from indicators.stuhr")
> summary(users_n_semesters$n_semesters_bremen)
  Min. 1st Qu.  Median    Mean 3rd Qu.    Max.
 1.000  1.000   1.500   4.726  6.750  23.000
```

p_is_used



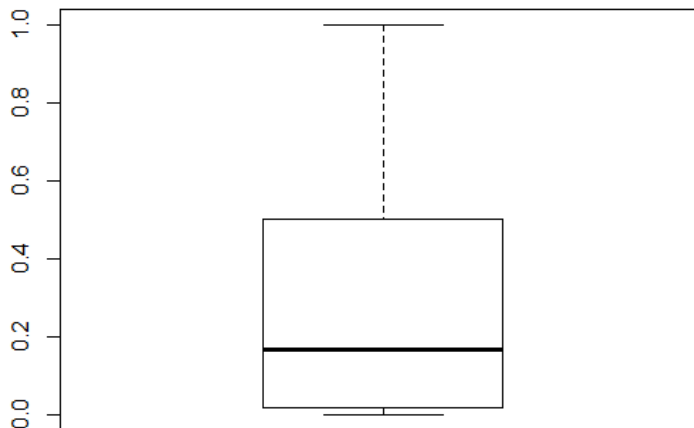
(10)

```
> summary(users$p_is_used)
  Min. 1st Qu.  Median    Mean 3rd Qu.    Max.
 0.00000 0.00000 0.02931 0.01838 1.00000
```

(11) Attention : il faut inverser l'ordre de cette variable ($p_is_edited = 1 - p_is_edited$) pour qu'elle varie dans le même sens que les autres variables. Ainsi, p_is_edited sera proche

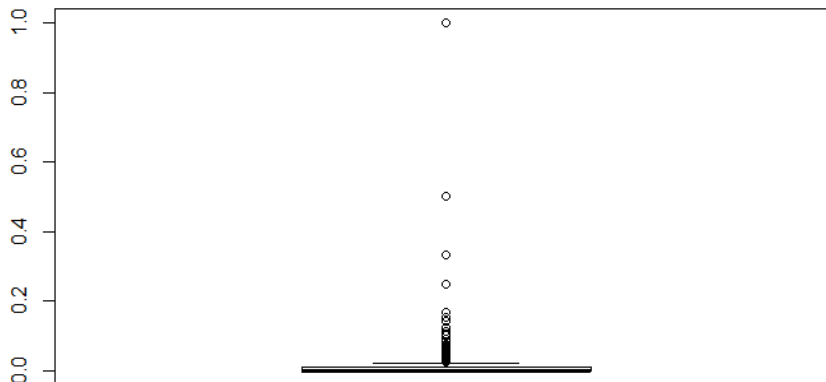
de 0 lorsqu'un contributeur est très édité, et proche de 1 lorsqu'il est peu édité.

p_is_edited



(12)

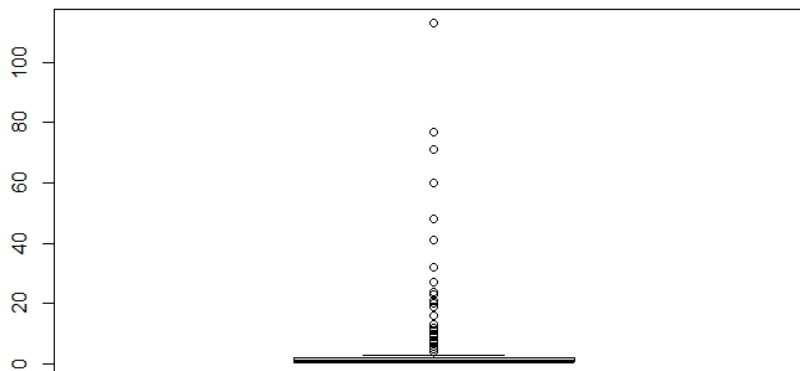
```
> users$p_is_deleted <- 1 - users$p_is_deleted
> summary(users$p_is_deleted)
  Min. 1st Qu.  Median    Mean 3rd Qu.    Max.
0.00000 0.01827 0.16667 0.31539 0.50000 1.00000
  p_is_deleted
```



(13)

```
> users$p_is_deleted <- 1 - users$p_is_deleted
> summary(users$p_is_deleted)
  Min. 1st Qu.  Median    Mean 3rd Qu.    Max.
0.00000 0.00000 0.00000 0.04393 0.00843 1.00000
```

Nb of weeks



(14)

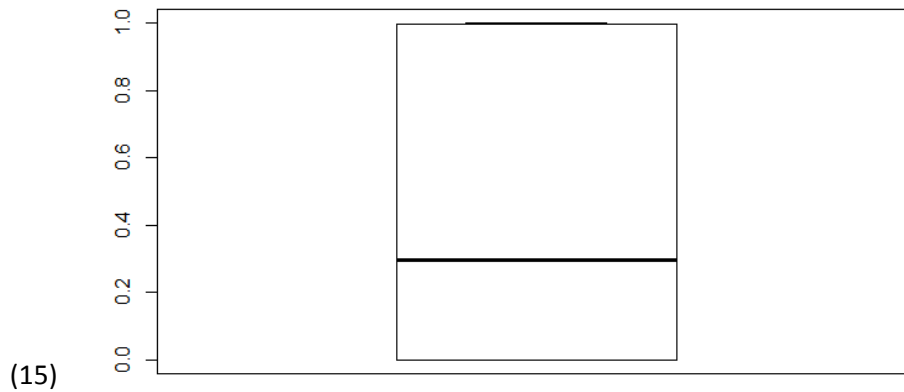
```
> summary(users$nbweeks)
```

```

Min. 1st Qu.  Median    Mean 3rd Qu.    Max.
1.000  1.000    1.000    3.688  2.000 113.000

```

focalisation



```
summary(stuhr$focalisation)
```

```

Min. 1st Qu.  Median    Mean 3rd Qu.    Max.
0.0000  0.0000  0.7265  0.6760  0.8415  1.5559

```

Changesets

Pour simplifier, on assigne un changeset par groupe de vandalisme : changeset A pour le groupe de vandalisme A, changeset B pour le groupe de vandalisme B, etc.

On assimile aussi l'âge de la contribution à l'âge du changeset.

Ce changeset peut éventuellement contenir d'autres contributions.

	Changeset A	Changeset B	Changeset C	Changeset D
nb éditions dans le changeset (1)	300	12	30	3
durée (en secondes) du changeset	5	100	20	60
Âge de la contribution	31536000s (1 an)	126144000s (4 ans)	157680000s (5 ans)	15552000s (6 mois)

```
(2) >summary(stuhr$chgst_edits)
```

```

Min. 1st Qu.  Median    Mean 3rd Qu.    Max.
  3     113     355    1967   3603    8997

```

```
(2) > summary(stuhr$chgst_duration_seconds)
```

```

Min. 1st Qu.  Median    Mean 3rd Qu.    Max.
 0.0    6.0    22.0   438.6   74.0 32028.0

```

```
(3) > summary(stuhr$lifespan_seconds) #En secondes
```

```

Min. 1st Qu.  Median    Mean 3rd Qu.    Max.
7634083 96615366 108770129 100034857 109674064 302007693

```

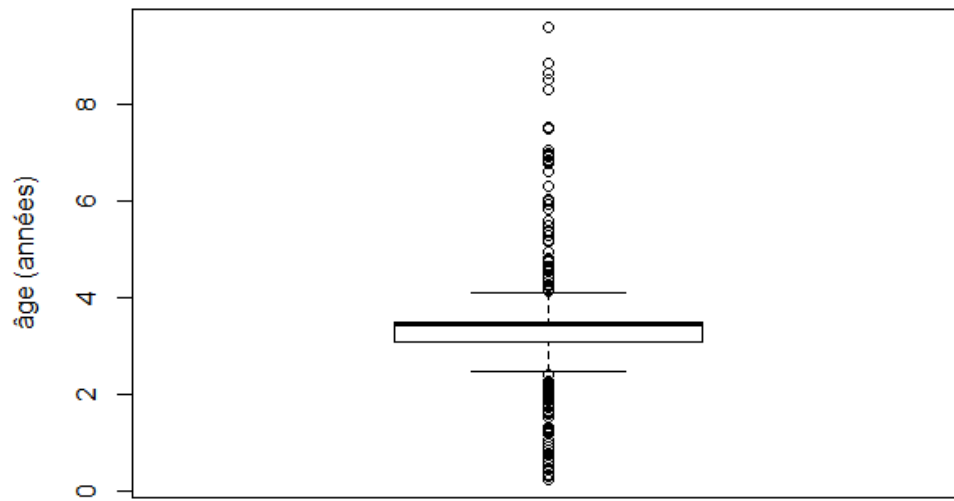
```
> summary(stuhr$lifespan_seconds/3600/24) # En jours
```

```

Min. 1st Qu.  Median    Mean 3rd Qu.    Max.
88.36 1118.23 1258.91 1157.81 1269.38 3495.46

```

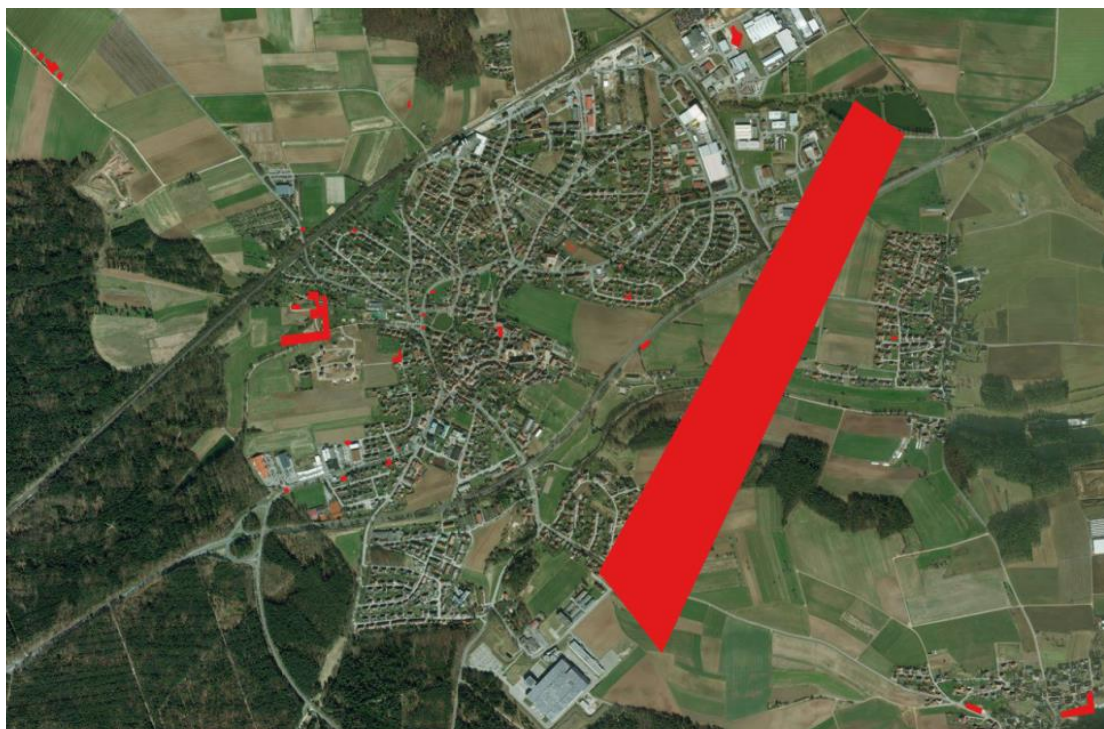
Distribution de l'âge de la contribution



Date

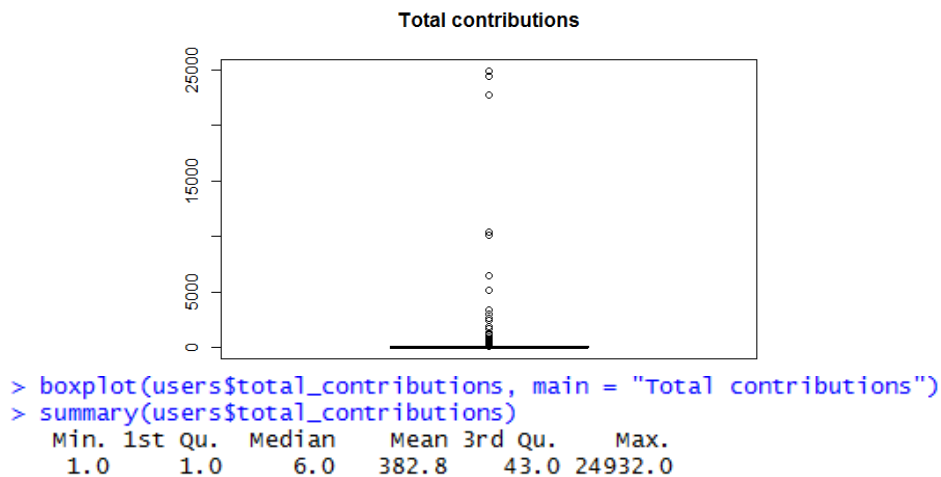
Changeset A a été produit le 7 décembre 2017

Profils de vandales synthétiques – Heilsbronn



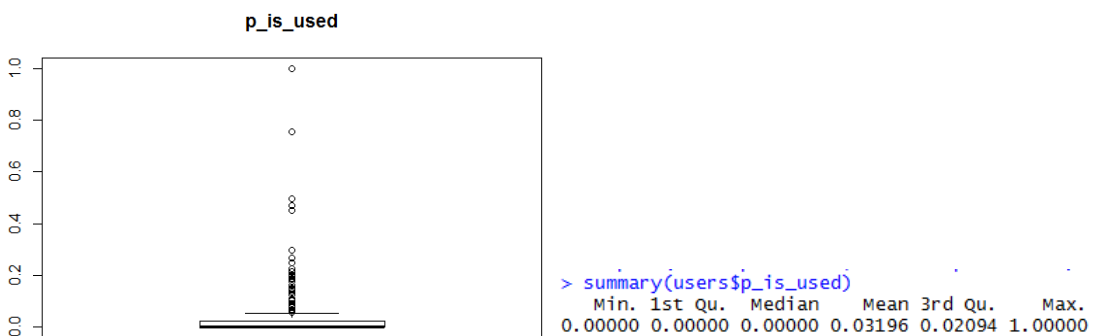
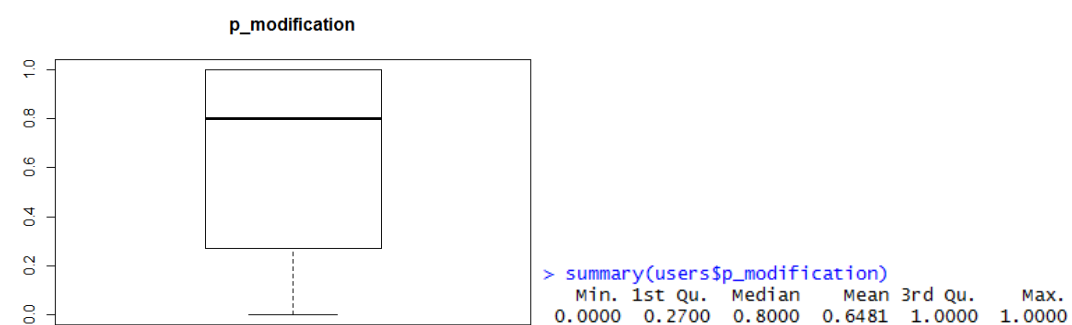
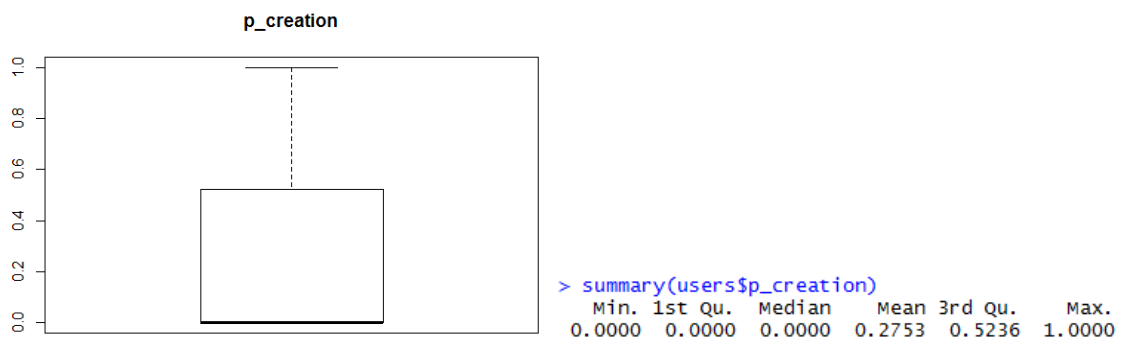
- 31 bâtiments vandalisés, ce sont de nouveaux éléments OSM (donc de version 1)
- On assigne un seul profile de vandale pour ces contributions de vandalisme
 - Il a réalisé un changeset contenant ces 31 contributions
 - En 120 secondes
 - Ces données existent depuis 500000 s soit ~6 jours
 - Uid du vandale = $\max(\text{uid des vrais contributeurs}) + 1 = 5614429$

Total contributions



- Nombre de contributions du vandale: 50
- On imagine que le vandale n'a pas uniquement vandalisé la base, et qu'il a contribué d'autres données non vandalisées

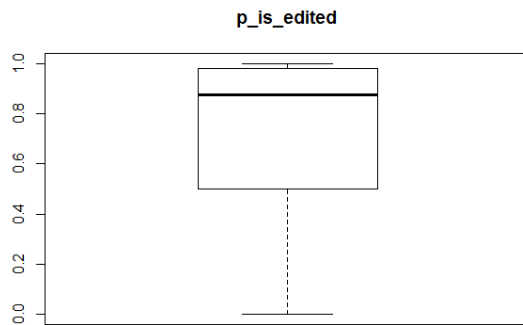
Ajout/Modification/suppression



On imagine que le vandale n'a fait que des ajouts, et qu'il n'a jamais été édité, ni utilisé

- Vandale p_creation = 1
- Vandale p_modification = 0
- Vandale p_is_used = 0
- Vandale p_is_edited = 0

Est édité / est supprimé



```
> summary(users$p_is_edited)
  Min. 1st Qu.  Median    Mean 3rd Qu.    Max.
0.0000 0.5000  0.8750  0.6847 0.9825  1.0000
```



```
> summary(users$p_is_deleted)
  Min. 1st Qu.  Median    Mean 3rd Qu.    Max.
0.0000 0.9869  1.0000  0.9541 1.0000  1.0000
```

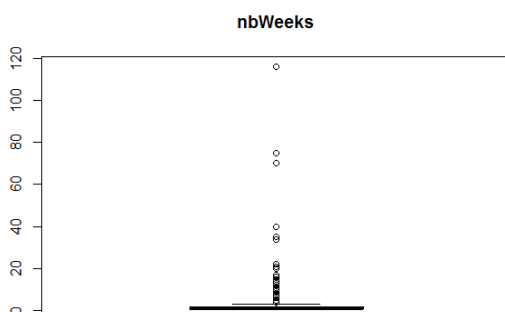
Si p_is_edited = 1, alors le contributeur n'a jamais été édité

Si p_is_deleted = 1 alors le contributeur n'a jamais été supprimé

On imagine que le vandale n'a jamais été édité ni supprimé :

- Vandale p_is_edited = 1
- Vandale p_is_deleted = 1

Semaines de contribution

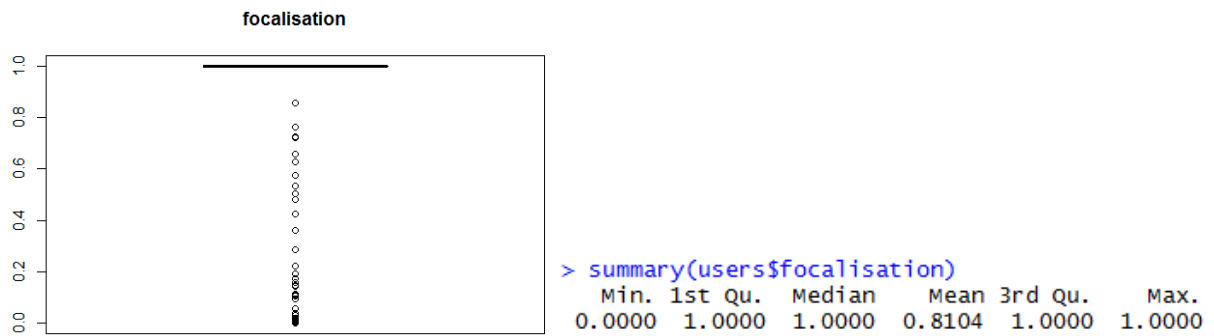


```
> summary(users$nbweeks)
  Min. 1st Qu.  Median    Mean 3rd Qu.    Max.
 1.000  1.000  1.000  3.419  2.000 116.000
```

On imagine que le vandale contribue depuis une semaine :

- Vandale nbWeeks = 1
- Vandale n_semesters = 1

Focalisation



Si focalisation = 1, alors le contributeur n'a participé que dans la zone

Si focalisation = 0, alors le contributeur a participé sur une région qui englobe la zone d'étude

On imagine que le vandale n'a contribué que dans Heilsbronn

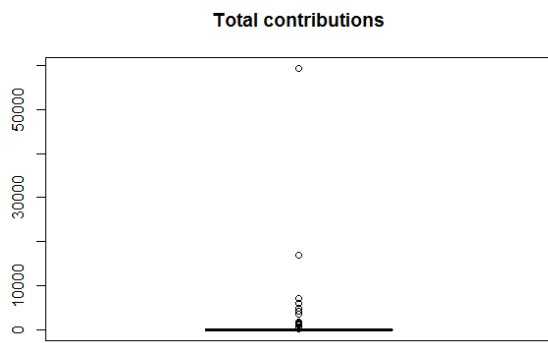
- Vandale focalisation = 1

Profils de vandales synthétiques - Lannilis



- 32 bâtiments vandalisés, ce sont de nouveaux éléments OSM (version 1)
- On assigne un seul profile de vandale pour ces contributions de vandalisme :
 - Il a réalisé un changeset contenant ces 32 contributions
 - En 120 secondes
 - Ces données existent depuis 500000 s soit ~6 jours
 - Uid du vandale = $\max(\text{uid des vrais contributeurs}) + 1 = 7585117$

Total contributions

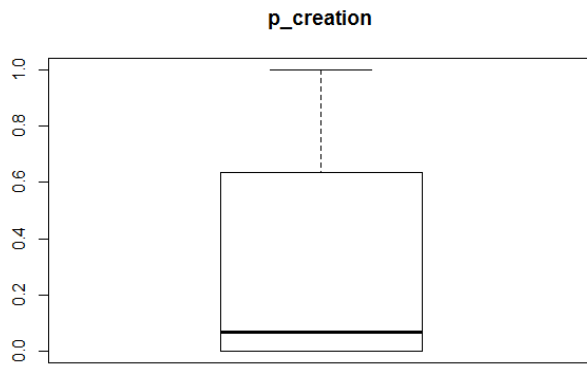


```
> summary(users$total_contributions)
  Min. 1st Qu.  Median    Mean 3rd Qu.    Max.
  1.0     2.0     9.5   620.4   75.0 59360.0
```

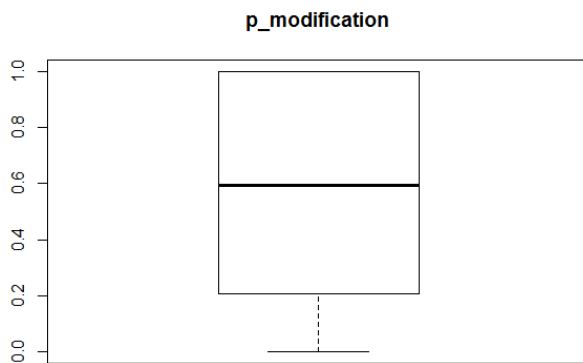
On imagine que le vandal n'a pas uniquement vandalisé la base

- Nombre de contributions du vandal : 100

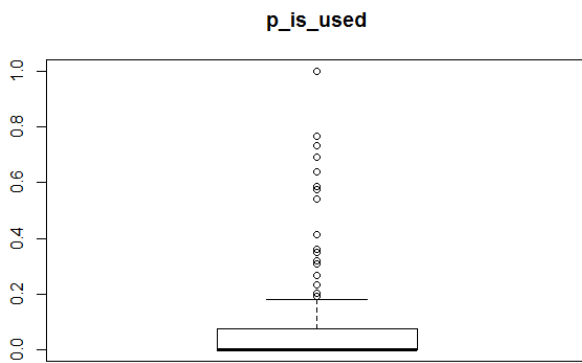
Ajout/Modification/Réutilisation



```
> summary(users$p_creation)
  Min. 1st Qu.  Median    Mean 3rd Qu.    Max.
0.00000 0.00000 0.06905 0.31070 0.62929 1.00000
```



```
> summary(users$p_modification)
  Min. 1st Qu.  Median    Mean 3rd Qu.    Max.
0.0000 0.2123 0.5950 0.5819 1.0000 1.0000
```

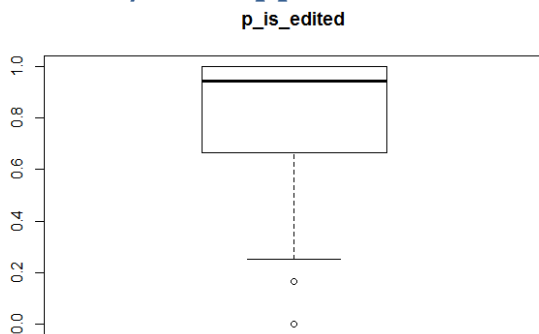


```
> summary(users$p_is_used)
  Min. 1st Qu.  Median    Mean 3rd Qu.    Max.
0.00000 0.00000 0.00000 0.07387 0.07447 1.00000
```

On imagine que le vandale a fait 32 ajouts, 68 modifications, il n'a jamais été utilisé

- Vandal p_creation = 0,32
- Vandal p_modification = 0,68
- Vandal p_is_used = 0

Est édité / est supprimé



```
> summary(users$p_is_edited)
  Min. 1st Qu.  Median    Mean 3rd Qu.    Max.
0.0000 0.6667  0.9438  0.7776 1.0000  1.0000
```



```
> summary(users$p_is_deleted)
  Min. 1st Qu.  Median    Mean 3rd Qu.    Max.
0.0000 0.9926  1.0000  0.9577 1.0000  1.0000
```

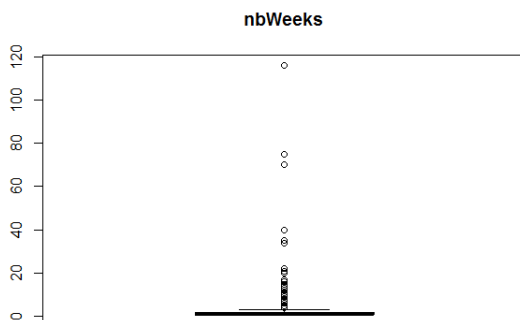
Si p_is_edited = 1, alors le contributeur n'a jamais été édité

Si p_is_deleted = 1 alors le contributeur n'a jamais été supprimé

On imagine que le vandale n'a jamais été édité ni supprimé :

- Vandale p_is_edited = 1
- Vandale p_is_deleted = 1

Semaines de contribution

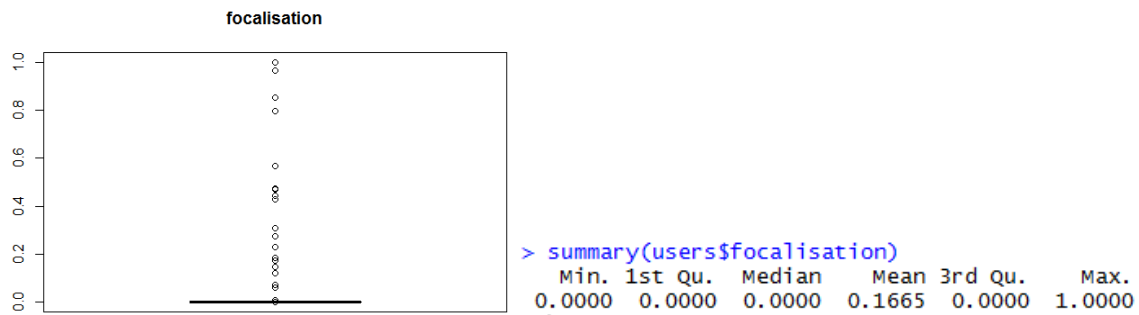


```
> summary(users$nbweeks)
  Min. 1st Qu.  Median    Mean 3rd Qu.    Max.
 1.000  1.000  1.000  3.419  2.000 116.000
```

On imagine que le vandale contribue depuis un mois :

- Vandale nbWeeks = 4
- Vandale n_semesters = 1

Focalisation



Si focalisation = 1, alors le contributeur n'a participé que dans la zone

Si focalisation = 0, alors le contributeur a participé sur une région qui englobe la zone d'étude

On imagine que le vandale n'a contribué a contribué sur une zone plus grande que Lannilis (3 fois plus grande que la surface de la ville)

- Vandale focalisation = 0,3

Table des figures

2	Délimitation de l'information géographique volontaire	10
I.1	Évolution des composantes du vandalisme (selon le contexte français)	17
I.2	« Graffiti is a crime »	18
I.3	Comparaison des typologies de (carto-)vandalisme	22
I.4	Typologies des processus spatio-temporels	23
I.5	Exemple d'erreur de débutant OSM	25
I.6	Exemple de carto-vandalisme OSM sur l'île Herald	26
I.7	Controverse sur une forêt en Lettonie	27
I.8	Exemple de carto-vandalisme Pokémon Go	28
I.9	Exemple de carto-vandalisme OSM sur une zone commerciale	30
I.10	Exemple de controverse OSM	31
II.1	Qualité du RGE ALTI®	40
II.2	Grand Bassin Rond du Jardin des Tuileries	44
II.3	Exemple de carto-vandalisme repéré par les contributeurs OSM	45
II.4	Opérations élémentaires	46
II.5	Modélisation des interactions	47
II.6	Historique des éditions de trois contributeurs A, B et C	47
II.7	Graphe de co-édition	48
II.8	Graphes de collaboration	49
II.9	Graphe de co-contribution	49
II.10	Typologies de contributeurs de la littérature	51
II.11	Confiance et qualité des contributeurs et des contributions	54
II.12	Exemple d'un réseau social multiplexe	55
II.13	Exemple de détection de communautés dans un réseau social	57

II.14 Agrégation de couches	59
II.15 Agrégation de couches	59
II.16 Exemple de maturation cartographique	63
II.17 Temporalité dans la saisie des données OSM	65
II.18 Graphe de co-temporalité	66
II.19 Graphe de co-location	66
II.20 Schéma conceptuel de la base de données historiques OSM	70
II.21 Chaîne de traitement des données OSM	71
II.22 Communautés OSM détectées par l’algorithme de Louvain	72
II.23 Configurations caractéristiques dans les communautés détectées	73
II.24 Degrés E/S d’une communauté centrée sur un modérateur	74
II.25 Degrés E/S d’une communauté centrée sur un pionnier	76
II.26 Histogrammes du coefficient de participation	78
II.27 Histogramme des coefficients de <i>clustering</i>	79
II.28 Histogramme des centralités de degré	80
II.29 Centralité de degré de la couche de co-édition	81
II.30 Centralité de degré de la couche de largeur de collaboration	81
II.31 Centralité de degré de la couche d’utilisation	81
II.32 Communautés du graphe agrégé	82
II.33 Communautés détectées du réseau multiplexe de Stuhr	84
II.34 Communautés détectées du réseau multiplexe de Katmandou	85
II.35 Erreur d’import automatique dans OSM	88
II.36 Chaîne de traitement des données	89
II.37 Méthode de calcul de l’enveloppe spatiale	90
II.38 Positionnement de l’enveloppe spatiale	91
II.39 Fonction de préférence de type gaussienne	93
II.40 Graphe de surclassement complet	93
II.41 Comparaison des différents classements par quartile	97
II.42 Comparaison du classement C_1 avec les autres classements	98
II.43 Contributeurs qualifiés manuellement dans les différents classements	100
III.1 Démarche de l’apprentissage supervisé	105

III.2	Modèle conceptuel d'un corpus de carto-vandalisme OSM	117
III.3	Carto-vandalisme synthétique à Aubervilliers	119
III.4	Carto-vandalisme de taille	121
III.5	Carto-vandalisme de type A	122
III.6	Carto-vandalisme de forme	122
III.7	Arbres de décision sur Stuhr	126
III.8	Arbre de décision sur Aubervilliers	127
III.9	Analyse unidimensionnelle des faux positifs	134
III.10	Anomalies selon le périmètre	135
III.11	Anomalies selon le périmètre	135
III.12	Distribution du périmètre du bâti et des faux positifs	135
III.13	Distribution de l'indicateur d'appariement	136
III.14	Arbres de décision sur Stuhr	137
III.16	Faux positif non apparié	138
III.18	Faux positifs anormaux selon le descripteur topologique	139
III.19	Anomalies détectées à Bondy	140
III.20	Anomalies syntaxiques dans les données de Bondy	140
III.21	Bâtiment réel absent la base de référence	141
III.24	Faux négatifs sur Stuhr avec le modèle Stuhr + Aubervilliers	147
III.25	Vrais et faux positifs d'Aubervilliers détectés avec le modèle RF de Stuhr	148
III.26	Faux positif de Heilsbronn détecté avec le modèle RF de Stuhr	149
III.27	Aperçu du bâti à Aubervilliers et Lannilis	150
III.28	Faux positifs à Lannilis avec le modèle RF Stuhr + Aubervilliers	151
III.29	Faux positifs à Heilsbronn avec le modèle RF Stuhr + Aubervilliers	151
III.30	Cas détectés sur Fougères par le modèle RF entraîné sur Stuhr	152
III.34	Grands bâtiments détectés par le modèle CNN de Stuhr	157
III.36	Faux vrais positifs détectés par le modèle CNN de Stuhr	157
III.37	Bâtiments à Aubervilliers classés avec le modèle CNN sur Aubervilliers	158
III.38	Faux positifs et vrais négatifs d'Aubervilliers	159
III.39	Faux positifs et vrais négatifs d'Aubervilliers	160
III.40	Vrais positifs sur Heilsbronn et Lannilis	161

III.41	Faux négatifs sur Lannilis et Heilsbronn	161
III.42	Aperçu des vrais positifs à Aubervilliers, Bondy et Fougères	163
1	Méthodologie globale de qualification de l'information géographique volontaire	174

Liste des tableaux

I.1	Historique des tags de l'île Hérald	26
I.2	Historique des tags de Anniņmuižas mežs / Anniņmuižas parks	27
I.3	Historique des tags d'une zone militaire	28
I.4	Historique des tags de Malishevë/Mališevo	29
I.5	Historique des tags de l'objet de la Figure I.9.	30
I.6	Extrait de l'historique des tags de l'objet de la Figure I.10.	31
II.1	Contrôle qualité dans les projets cartographiques collaboratifs	42
II.2	Historique des tags du Grand Bassin Rond	45
II.3	Profils de contributeurs identifiés par les graphes de collaboration	48
II.4	Dimensions de la confiance	52
II.5	Degrés entrants de $G_{utilisation}$	74
II.6	Activité des contributeurs « noyaux »	75
II.7	Coefficients de clustering de Suthr et Katmandou	85
II.8	Indicateurs de fiabilité du contributeur	88
II.9	Poids utilisés pour calculer la moyenne pondérée du score de fiabilité.	91
II.10	Paramétrage du modèle de décision multicritère PROMETHEE II	94
II.11	Résumé statistique du nombre de contributions des 30 contributeurs étudiés.	95
II.12	Qualification manuelle de la fiabilité des contributeurs OSM.	95
II.13	Positionnement des petits contributeurs dans les quartiles de C_1 à C_4	101
III.1	Récapitulatif des méthodes et métriques de détection du vandalisme	110
III.2	Corpus de vandalisme issus de l'état de l'art	116
III.3	Décompte du carto-vandalisme synthétique dans les zones d'étude	120
III.4	Décompte du carto-vandalisme synthétique dans les zones d'étude	120

III.5	Détection du carto-vandalisme avec DBSCAN sur Aubervilliers	129
III.6	Performance de DBSCAN avec différents indicateurs de contributeur	129
III.7	Influence de l'indicateur d'appariement sur DBSCAN	130
III.8	Résultats de détection avec DBSCAN avec les descripteurs optimaux sur Aubervilliers et Stuhr.	131
III.9	Descripteurs optimaux pour la détection du carto-vandalisme dans les données d'Aubervilliers et Stuhr	132
III.10	Influence d'un indicateur de fiabilité sur Stuhr	132
III.11	Résultats DBSCAN avec les descripteurs optimaux de Stuhr	133
III.12	Détection d'anomalies sur les bâtis OSM de Bondy avec DBSCAN . .	140
III.13	Paramétrage de DBSCAN sur Stuhr et Aubervilliers	142
III.14	Descripteurs utilisés pour construire les forêts aléatoires	145
III.15	Prédiction des modèles de RF sur les zones d'entraînement	146
III.16	Résultats avec un modèle RF entraîné sur deux zones	146
III.17	Détection sur des zones différentes de la zone d'entraînement	147
III.18	Détection sur une zone du même pays que la zone d'entraînement . .	149
III.19	Détection sur Lannilis et Heilsbronn par le modèle RF Stuhr + Au- bervilliers	150
III.20	Détection du carto-vandalisme sur Bondy et Fougères	152
III.21	Nombre d'images utilisées pour l'entraînement des modèles CNN . . .	156
III.22	Résultats des modèles CNN sur leur zone d'entraînement	157
III.23	Détection du carto-vandalisme avec un CNN entraîné sur Aubervil- liers et Stuhr	159
III.24	Détection par le modèle CNN Stuhr + Aubervilliers	160
III.25	Détection sur Bondy et Fougères avec un modèle CNN entraîné sur Aubervilliers	162