

# Des Textes Communautaires à la Recommandation

Damien Poirier

Soutenance de thèse pour obtenir le grade de Docteur en Informatique

Orléans, le 11 février 2011



## Contexte et problématique

## Description de la méthodologie

- Extraction des textes

- Inférence de notes sur les textes

- Recommandation par filtrage collaboratif

## Évaluation

## Conclusion

## Contexte et problématique

### Description de la méthodologie

Extraction des textes

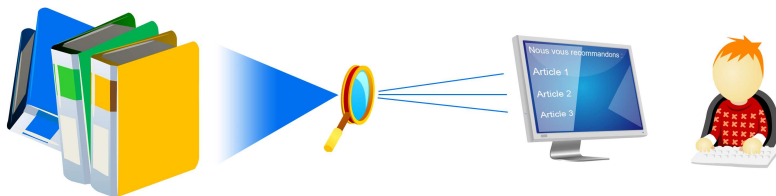
Inférence de notes sur les textes

Recommandation par filtrage collaboratif

### Évaluation

### Conclusion

## Qu'est-ce que la recommandation automatique ?



**Objectif** : sélectionner des articles (appelés items) dans un catalogue pour en faciliter la lecture.

# Qu'est-ce que la recommandation automatique ?

Deux grands types de recommandations possibles :

## Contextuelle :

Consiste à proposer des items en lien avec l'item observé



## Personnalisée :

Consiste à proposer des items adaptés à chaque utilisateur



# Qu'est-ce que la recommandation automatique ?

Deux grands types de recommandations possibles :

## Contextuelle :

Consiste à proposer des items en lien avec l'item observé



## Personnalisée :

Consiste à proposer des items adaptés à chaque utilisateur



## Exemple d'application : la VOD



## Recommandation personnalisée

L'intégration d'un service de recommandation personnalisée nécessite tout d'abord des **informations sur les utilisateurs...**

On peut alors collecter des données du type :

- ▶ L'utilisateur a regardé tel ou tel film
- ▶ L'utilisateur a aimé/pas aimé tel ou tel film









## Recommandation personnalisée

L'intégration d'un service de recommandation personnalisée nécessite tout d'abord des **informations sur les utilisateurs...**

On peut alors collecter des données du type :

- ▶ L'utilisateur a regardé tel ou tel film
- ▶ L'utilisateur a aimé/pas aimé tel ou tel film

⇒ **Matrice d'usages**

	Item 1	Item 2	Item 3	Item 4	...
Utilisateur 1					...
Utilisateur 2					...
Utilisateur 3					...
Utilisateur 4					...
...					...

## Recommandation personnalisée

Ainsi que des **informations sur les items**.

Deux grandes approches sont possibles (Nageswara Rao et al.<sup>1</sup>) :

- ▶ **Filtrage collaboratif** (basé sur les profils des utilisateurs)
- ▶ **Filtrage basé sur le contenu** (basé sur des descripteurs d'items)

---

1. Application domain and functional classification of recommender systems a survey, *In Desidoc journal of library and information technology, volume 28, pages 17-36, ACM Press, 2008*

# Principe du filtrage collaboratif

	Item 1	Item 2	Item 3	Item 4	...
Utilisateur 1	☹️		😊		...
Utilisateur 2	😊			😊	...
Utilisateur 3		☹️			...
Utilisateur 4				😊	...
...					...

- ▶ 1<sup>ère</sup> étape :
  - ▶ Calcul des similarités entre utilisateurs
  - ▶ ou similarités entre items
- ▶ 2<sup>ème</sup> étape :
  - ▶ Prédiction des notes manquantes à l'aide des similarités et des profils utilisateurs
- ▶ 3<sup>ème</sup> étape :
  - ▶ Recommandation des meilleures notes prédites

⇒ Très performant mais nécessite une grande quantité de données

# Principe du filtrage collaboratif

	Item 1	Item 2	Item 3	Item 4	...
Utilisateur 1	☹️		😊		...
Utilisateur 2	😊			😊	...
Utilisateur 3		☹️			...
Utilisateur 4				😊	...
...					...

- ▶ 1<sup>ère</sup> étape :
  - ▶ Calcul des similarités entre utilisateurs
  - ▶ ou similarités entre items
- ▶ 2<sup>ème</sup> étape :
  - ▶ Prédiction des notes manquantes à l'aide des similarités et des profils utilisateurs
- ▶ 3<sup>ème</sup> étape :
  - ▶ Recommandation des meilleures notes prédites

⇒ Très performant mais nécessite une grande quantité de données

# Principe du filtrage collaboratif

	Item 1	Item 2	Item 3	Item 4	...
Utilisateur 1	☹️		😊		...
Utilisateur 2	😊			😊	...
Utilisateur 3	↓	☹️			...
Utilisateur 4	😄			😊	...
...					

- ▶ 1<sup>ère</sup> étape :
  - ▶ Calcul des similarités entre utilisateurs
  - ▶ ou similarités entre items
- ▶ 2<sup>ème</sup> étape :
  - ▶ Prédiction des notes manquantes à l'aide des similarités et des profils utilisateurs
- ▶ 3<sup>ème</sup> étape :
  - ▶ Recommandation des meilleures notes prédites

⇒ Très performant mais nécessite une grande quantité de données

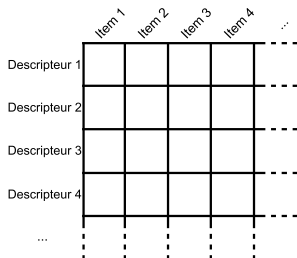
# Principe du filtrage collaboratif

	Item 1	Item 2	Item 3	Item 4	...
Utilisateur 1	☹️		😊		...
Utilisateur 2	😊			😊	...
Utilisateur 3		☹️			...
Utilisateur 4	😊	←		😊	...
...					...

- ▶ 1<sup>ère</sup> étape :
  - ▶ Calcul des similarités entre utilisateurs
  - ▶ ou **similarités entre items**
- ▶ 2<sup>ème</sup> étape :
  - ▶ Prédiction des notes manquantes à l'aide des similarités et des profils utilisateurs
- ▶ 3<sup>ème</sup> étape :
  - ▶ Recommandation des meilleures notes prédites

⇒ Très performant mais nécessite une grande quantité de données

# Principe du filtrage basé sur le contenu



- ▶ Exemples de descripteurs de films :
  - ▶ Genre, réalisateur, acteurs principaux, etc.
- ▶ 1<sup>ème</sup> étape :
  - ▶ Calcul des similarités entre items
- ▶ 2<sup>ème</sup> étape :
  - ▶ Prédiction des notes
- ▶ 3<sup>ème</sup> étape :
  - ▶ Recommandation des meilleures notes prédites

⇒ Moins performant que le filtrage collaboratif

# Problématique

Quelle que soit l'approche utilisée, des données sont nécessaires :

- ▶ Données d'usages pour le filtrage collaboratif
- ▶ Et données descriptives pour le filtrage basé sur le contenu

**Que faire lorsqu'elles ne sont pas disponibles ?**



# Problématique

Le Web regorge de ressources, principalement textuelles, qui contiennent un grand nombre d'avis et d'opinions d'utilisateurs sur des sujets de toutes sortes.

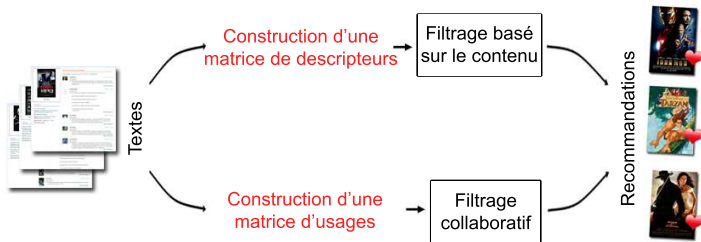
(blogs, forums, boîtes à réactions, etc.)



# Problématique

## Objectif principal de la thèse

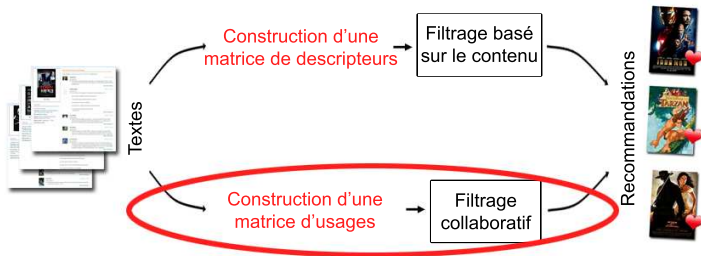
Montrer que l'on peut extraire, à partir de textes communautaires non structurés, des informations utiles aux moteurs de recommandation existants.



# Problématique

## Objectif principal de la thèse

Montrer que l'on peut extraire, à partir de textes communautaires non structurés, des informations utiles aux moteurs de recommandation existants.



Contexte et problématique

Description de la méthodologie

Extraction des textes

Inférence de notes sur les textes

Recommandation par filtrage collaboratif

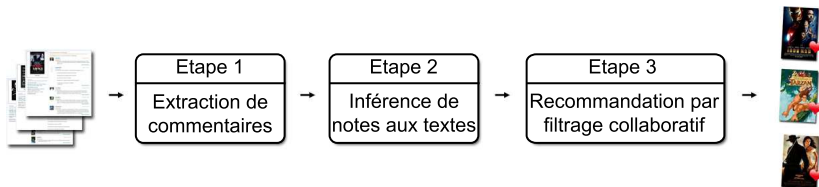
Évaluation

Conclusion

# Méthodologie

## Principale contribution :

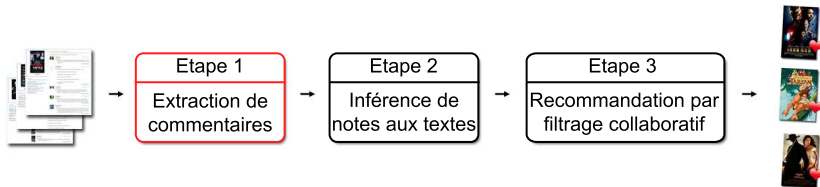
Mise en place d'une méthodologie complète, **automatisée** et en grande partie **reproductible** (étapes 2 et 3)



# 1<sup>ère</sup> étape : Extraction des textes

## Objectif de la tâche :

Acquérir des textes exprimant des opinions sur les produits de notre catalogue



# Provenance des commentaires

Domaine d'application : les films

- ▶ Construction d'un crawler pour le site communautaire *Flixster*<sup>2</sup>
  1. Recherche les films souhaités à l'aide du moteur de recherche du site
  2. Aspiration des pages contenant les commentaires sur les films
  3. Extraction des commentaires
- ▶ Nettoyage des textes (suppression des urls, des balises html, des lignes vides, des caractères inconnus, etc.)

---

2. [www.flixster.com](http://www.flixster.com)

# Provenance des commentaires

The screenshot shows the Flixster website interface. At the top, there is a navigation bar with the Flixster logo, a search bar containing 'Movies, Actors, Directors...', and links for 'Sign Up' and 'Login'. Below the navigation bar are tabs for 'Home', 'Movies', 'Profile', 'Friends', 'Meet People', 'Fun & Games', and 'Watch Now'. The main content area is titled 'Dude, Where's My Car? Reviews and Ratings'. On the left, there is a movie poster for 'Dude, Where's My Car?' featuring Ashton Kutcher and Seann William Scott. Below the poster is the movie title. To the right of the poster is a 'Summary' section with the following text: 'Dude, Where's My Car? Summary', 'Starring: Ashton Kutcher, Seann William Scott, Kristy Swanson, Jennifer Garner, Marla Sokoloff', 'Directed by: Danny Leiner', 'Genres: Comedy', 'Release Date: December 15, 2000', and 'DVD Release Date: June 26, 2001'. The right side of the page displays a list of user reviews. Each review includes a user profile picture, the username, the date of the review, a star rating, and the review text. The reviews are as follows: 1. User 'iluvspastley13' (January 26, 2011) gave a 4-star rating and wrote 'Really weird and random, but funny'. 2. User 'JuRn' (January 5, 2011) gave a 4-star rating and wrote 'I got quite a few laughs from dis good comedy movie'. 3. User 'TheFilmApache' (December 5, 2010) gave a 1-star rating and wrote 'The worst comedy I have ever seen...'. 4. User 'Quasim0' (November 28, 2010) gave a 1-star rating and wrote 'This movie is pointless and a waste of time'. 5. User 'incz' (November 25, 2010) gave a 4-star rating and wrote 'silly but still funny! :D'. 6. User 'Superman1984' (November 8, 2010) gave a 4-star rating and wrote 'I Liked the part at near the end. Daddy I wanna ride that ride. me too son.' Each review has a 'Share This Review' link.



# Provenance des commentaires

**Flixster** Sign Up | Login  
Movies, Actors, Directors... Search

Home Movies Profile Friends Meet People Fun & Games Watch Now

### Dude, Where's My Car? Reviews and Ratings

First | Previous | 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 | 9 | 10 | 11 | Next | Last

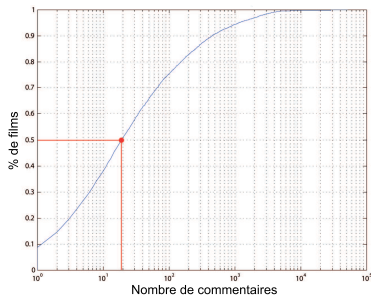
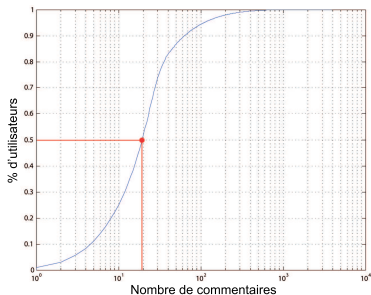
Avatar	Username	Date	Rating	Comment	Share
	iluvspastley13	January 26, 2011	☆☆☆☆☆	Really weird and random, but funny	Share This Review
	JuRn	January 5, 2011	☆☆☆☆☆	I got quite a few laughs from dis good comedy movie	Share This Review
	TheFilmApache	December 5, 2010	☆☆☆☆☆	The worst comedy I have ever seen...	Share This Review
	suasim0	November 28, 2010	☆☆☆☆☆	This movie is pointless and a waste of time	Share This Review
	incz		☆☆☆☆☆	silly but still funny! :D	Share This Review
	Superman1984	November 8, 2010	☆☆☆☆☆	I Liked the part at near the end. Daddy I wanna ride that ride. me too son.	Share This Review

**Données extraites :  
Id Utilisateur - Id Film - Commentaire  
- Note**

# Description du corpus

## Chiffres :

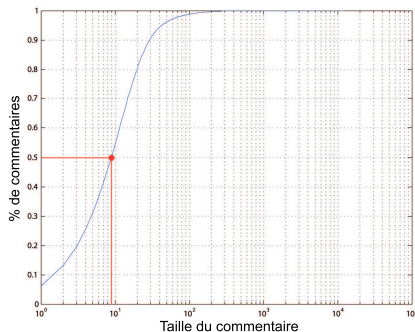
- ▶ Nombre : > 3 millions de commentaires  
(sous la forme de 4-uplets Id utilisateur-Id film-Texte-Note)
- ▶ Utilisateurs concernés :  $\simeq 100\ 000$ 
  - ⇒ Moyenne : **34** com./utilisateurs
  - ⇒ Médiane : **20** com./utilisateurs
- ▶ Films concernés :  $\simeq 10\ 000$ 
  - ⇒ Moyenne : **318** com./films
  - ⇒ Médiane : **20** com./films



# Nature des textes

## 1<sup>ère</sup> particularité : la taille

- ▶ Nombre moyen d'unités (mots + ponctuations + symboles + chiffres, etc.) par commentaire : 15
- ▶ 90% des commentaires font moins de 30 mots
- ▶ 50% des commentaires font moins de 9 mots



# Nature des textes

2<sup>ème</sup> particularité : le style d'écriture

⇒ Emploi du **Discours Électronique Médié** (DEM) (Panckhurst<sup>3</sup>)

- ▶ Didascalies électroniques
  - ▶ Smileys : " ;-) "
  - ▶ Abus de majuscules : " NON "
  - ▶ Étirement de mots : " coool "
  - ▶ Répétitions de ponctuation : " ??? "
- ▶ Fautes d'orthographe
- ▶ Néographie
  - ▶ Substitutions : " g " (j'ai), " mwa " (moi)
  - ▶ Réductions : " dsl ", " jaten "
  - ▶ Suppressions : " salut cava "
  - ▶ Ajouts : " oki "

---

3. Le discours électronique médié : bilan et perspectives, 2006

## Nature des textes

Particularités moins visibles :

- ▶ Moins de verbes
- ▶ Moins de pronoms
- ▶ Langage familier
- ▶ **Dialecte évolutif**

Exemples :

" weeeeeiiiiiiiiirrrrrddddd movie "

" i liked this cuz this shit can really happen. "

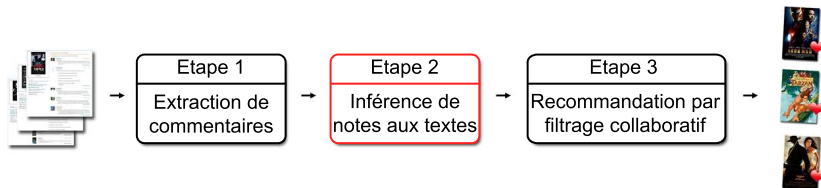
" I see it @ work everyday so this movie sucks! "

" OMG I LOVED THIS MOVIE!!!!!!!!!!!!!!!!!!!!<3 "

## 2<sup>ème</sup> étape : Inférence des notes

### Objectif de la tâche :

Attribuer des notes aux textes afin d'obtenir une matrice d'usages



# Classification d'opinion

## Objectif de la classification d'opinion

Classer les textes selon l'opinion exprimée par l'auteur (positif ou négatif)

**Deux grands types de méthodes** (Pang et al.<sup>4</sup>) :

- ▶ Les méthodes basées sur les lexiques
- ▶ Les méthodes basées sur l'apprentissage automatique

---

4. Opinion mining and sentiment analysis, *Found. Trends Inf. Retr.*, 2008

# Classification d'opinion : 1<sup>ère</sup> approche

## Méthodes basées sur les lexiques

**Objectif** : Construire des lexiques contenant des mots liés à l'expression de l'opinion et classer les textes en fonction de la présence de ces mots.

**Exemples** :

Mots positifs	awesome, great, best, love, ...
Mots négatifs	dislike, awful, hate, bad, ...

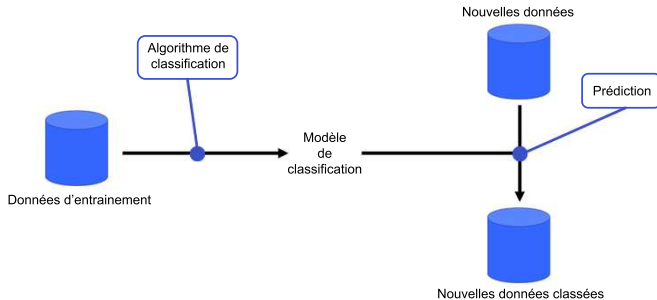
⇒ Trop de contraintes manuelles pour l'automatisation des traitements



## Classification d'opinion : 2<sup>ème</sup> approche

### Méthodes basées sur l'apprentissage automatique

**Objectif de la classification supervisée** : apprendre des modèles à partir d'exemples et classer les nouveaux textes à l'aide de ces modèles.



# Méthodes basées sur l'apprentissage automatique

Plusieurs choix à plusieurs niveaux :

Segmentation	Représentation
Mots	Binaire
N-grammes de mots	Fréquentielle
Unités	Fréquentielle normalisée
N-grammes de caractères	TF-Idf
...	...
Classifieur	Nombre de classes
SVM	2
(S) Naïve Bayes <sup>5</sup>	3
Réseaux de neurones	5
Règles de décision	Échelle continue
...	...

5. SNB : Selective Naïve Bayes

# Méthodes basées sur l'apprentissage automatique

## Pré-sélection basée sur la littérature :

- ▶ Pang et al. : Thumbs up? : sentiment classification using machine learning techniques. *In Proceedings of EMNLP '02, 2002*
- ▶ Wilson et al. : Just how mad are you? finding strong and weak opinion clauses. *In In Proceedings of AAAI, 2004*
- ▶ Nigam et al. : Towards a robust metric of polarity. *In Computing Attitude and Affect in Text : Theory and Applications, 2006*
- ▶ Kaiser et al. : Swarm intelligence for analyzing opinions in online communities. *Hawaii International Conference on System Sciences, 2010*
- ▶ Plantié et al. : Is a voting approach accurate for opinion mining? *In Proceedings of DaWaK '08, 2008*
- ▶ Maurel et al. : Classification d'opinions par méthodes symbolique, statistique et hybride. *In Actes du 3<sup>ème</sup> DEfi Fouille de Textes, 2007*
- ▶ etc.

# Méthodes basées sur l'apprentissage automatique

Segmentation	Représentation
Mots	Binaire
N-grammes de mots	<b>Fréquentielle</b>
<b>Unités</b>	<b>Fréquentielle normalisée</b>
<b>N-grammes de caractères</b>	<b>TF-Idf</b>
...	...
Classifieur	Nombre de classes
<b>SVM</b> <sup>6</sup>	<b>2</b>
<b>(S) Naïve Bayes</b> <sup>7</sup>	<b>3</b>
Réseaux de neurones	<b>5</b>
Règles de décision	Échelle continue
...	...

6. SVM<sup>light</sup>

7. KHIOPS

# Méthodes basées sur l'apprentissage automatique

Nombre initial d'unités (mots + ponctuations...) = **150 000**

Choix concernant les pré-traitements :

► Pré-traitements de base

Pré-traitements	Taille du vocabulaire
Minusculation	≈ 60 000
Suppression des mots très peu fréquents	≈ 20 000

► Pré-traitements additionnels

Pré-traitements	Taille du vocabulaire
Corrections orthographiques et lemmatisation <sup>8</sup>	≈ 15 000
Suppression des mots vides <sup>9</sup>	≈ 19 700

---

8. Utilisation du logiciel TiLT

9. Évaluation de plusieurs listes prises sur Internet

# Méthodes basées sur l'apprentissage automatique

Conditions expérimentales :

- ▶ Construction d'un corpus d'apprentissage :
  - ▶ 175 000 nouveaux commentaires déjà classés
  - ▶ Classes équilibrées
  - ▶ Aucune intersection au niveau des films et des utilisateurs
  - ▶ Pas plus de 10 commentaires rédigés par le même auteur
- ▶ Mesure utilisée :

$$F_{score} = 2 * \frac{Précision * Rappel}{Précision + Rappel}$$

$$Précision = \sum_{i=1}^n \frac{\text{Documents correctement attribués à la classe } i}{\text{Nombre de documents attribués à la classe } i}$$

$$Rappel = \sum_{i=1}^n \frac{\text{Documents correctement attribués à la classe } i}{\text{Nombre de documents appartenant à la classe } i}$$

## Méthodes basées sur l'apprentissage automatique

Exemple d'évaluation :

$F_{score}$  obtenus pour chaque représentation avec les trois classifieurs pour une classification binaire

	Fréquentielle	Fréq. normalisée	TF-IDF
SVM	<b>0,741</b>	0,724	0,726
SNB	0,726	0,729	0,728
NB	0,695	0,708	0,708

⇒ Meilleurs résultats obtenus avec SVM et la représentation fréquentielle

## Méthodes basées sur l'apprentissage automatique

Exemple d'évaluation :

$F_{score}$  obtenus avec les pré-traitements minimaux et linguistiques et toutes les représentations possibles avec le classifieur SVM

	Fréquentielle	Fréq. normalisée	TF-IDF
Basiques	<b>0,741</b>	0,724	0,726
Anti-dictionnaire	0,720	0,700	0,703
Linguistiques	0,738	0,722	0,724

⇒ Corrections et lemmatisation n'entraînent pas d'amélioration des résultats



# Observations sur l'impact des pré-traitements courants

## Suppression de la ponctuation

ATTENTION : La **ponctuation** joue un rôle important

Cas du " ! "		
Fréquence	Négatif	Positif
$n \geq 2$	37%	<b>63%</b>

*"so funny!!!!u have 2 c it!!!!!"*

Cas du " . "		
Fréquence	Négatif	Positif
$n \geq 1$	<b>65%</b>	35%

*"so disappointed ....."*

Cas du " ? "		
Fréquence	Négatif	Positif
$n \geq 1$	<b>74%</b>	26%

# Observations sur l'impact des pré-traitements courants

## Suppression des mots vides

ATTENTION : Les **mots vides** jouent un rôle important

Cas de " but "		
Fréquence	Négatif	Positif
$n \geq 1$	<b>71%</b>	29%

*"i like science fiction but not this one"*

Cas de " it "		
Fréquence	Négatif	Positif
$n \geq 1$	23%	<b>77%</b>

*"it's tooooo funny"*

Cas de " or "		
Fréquence	Négatif	Positif
$n \geq 1$	<b>72%</b>	28%

# Observations sur l'impact des pré-traitements courants

## Lemmatisation ou stemmatisation

ATTENTION : La **conjugaison** des verbes joue également un rôle

Cas de " is "		
Fréquence	Négatif	Positif
$n \geq 1$	41%	<b>59%</b>

"I Love It. It is A Great Movie"

Cas de " was "		
Fréquence	Négatif	Positif
$n \geq 1$	<b>67%</b>	33%

"the movie was a waste of money"

Cas de " were "		
Fréquence	Négatif	Positif
$n \geq 1$	<b>69%</b>	30%

## Conclusion des expérimentations

Méthode sélectionnée :

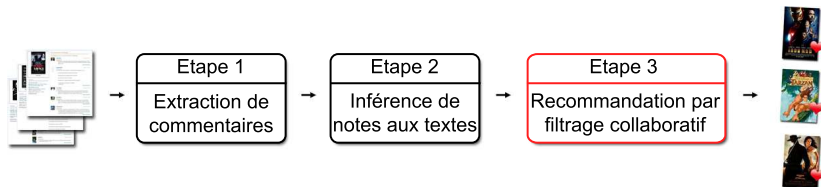
Critère	Choix
Pré-traitements	Basiques
Segmentation	Unités (Mots, ponctuations, nombres...)
Représentation	Fréquentielle
Classifieur	SVM
Nombre de classes	2

⇒ **Construction de la matrice d'usages**

## 3<sup>ème</sup> étape : Recommandations

### Objectif de la tâche :

Établir des recommandations à partir de la matrice d'usages



## Description du moteur

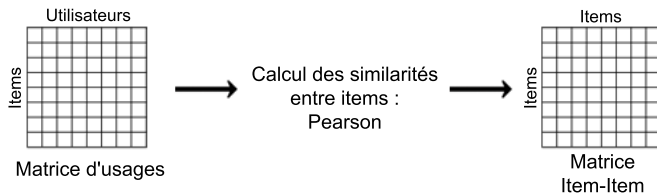
- ▶ Moteur développé chez Orange (Meyer<sup>10</sup>)
- ▶ Choix issus du contexte industriel
- ▶ Permet le filtrage collaboratif et basé sur le contenu
- ▶ Permet la recommandation contextuelle et personnalisée
- ▶ Performances au niveau de l'état de l'art

⇒ Fonctionnement en 2 étapes

---

10. Apport des données thématiques dans les systèmes de recommandation : hybridation et démarrage à froid, *In Extraction et gestion des connaissances (EGC'2011), Brest, France, 2011*

## Construction de la matrice item-item



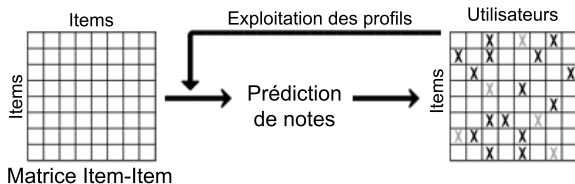
$$Pearson(i, j) = \frac{\sum_{\{u \in S_i \cap S_j\}} (r_{iu} - \bar{r}_i) \times (r_{ju} - \bar{r}_j)}{\sqrt{\sum_{\{u \in S_i \cup S_j\}} (r_{iu} - \bar{r}_i)^2 \sum_{\{u \in S_i \cup S_j\}} (r_{ju} - \bar{r}_j)^2}}$$

$S_i$  (resp.  $S_j$ ) : l'ensemble des utilisateurs ayant noté l'item  $i$  (resp.  $j$ )

$r_{iu}$  (resp.  $r_{ju}$ ) : la note donnée par l'utilisateur  $u$  sur l'item  $i$  (resp.  $j$ )

$\bar{r}_i$  (resp.  $\bar{r}_j$ ) la moyenne des notes obtenues par  $i$  (resp.  $j$ )

# Prédiction des notes



$$p_{ui} = \bar{r}_i + \frac{\sum_{\{j \in S_u\}} \text{Pearson}(i, j) \times (r_{uj} - \bar{r}_j)}{\sum_{\{j \in S_u\}} \text{Pearson}(i, j)}$$

$S_u$  : l'ensemble des notes connues données par l'utilisateur  $u$



Contexte et problématique

Description de la méthodologie

Extraction des textes

Inférence de notes sur les textes

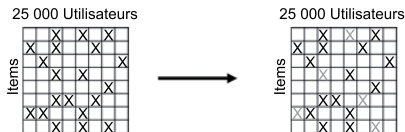
Recommandation par filtrage collaboratif

Évaluation

Conclusion

## Conditions expérimentales

Données de test : extraites du challenge Netflix (Bennett et al.)<sup>11</sup>



X : 90% des notes dédiées aux profils utilisateurs

X : 10% des notes dédiées à la validation des notes prédites

Données d'entrée :

- ▶ Données d'usages obtenue à partir des commentaires **Flixster**
  - ▶ Notes issues de la classification d'opinion
  - ▶ Vraies notes
- ▶ Descripteurs **IMDb** pour le filtrage basé sur le contenu
- ▶ Données d'usages **Netflix** (profils) pour le filtrage collaboratif

11. The netflix prize, *San Jose, California, ACM, 2007*

## Méthode d'évaluation

Mesure utilisée : Erreur quadratique moyenne

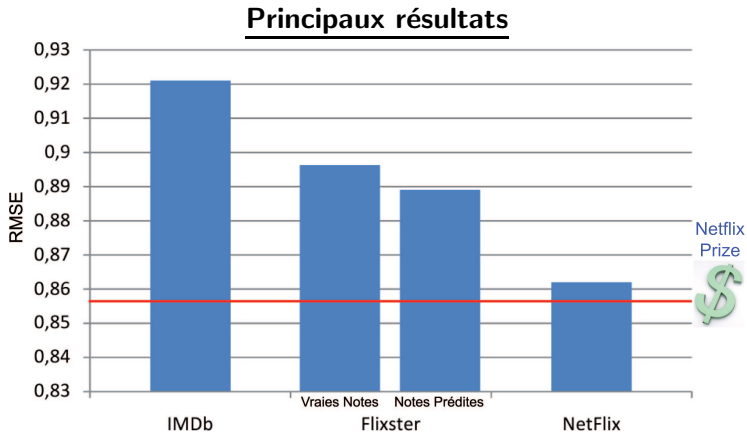
$$RMSE = \sqrt{\frac{\sum_{u,i} (p_{ui} - n_{ui})^2}{n}}$$

$n_{ui}$  = note réelle donnée par l'utilisateur  $u$  sur l'item  $i$

$p_{ui}$  = note prédite par le moteur de recommandation

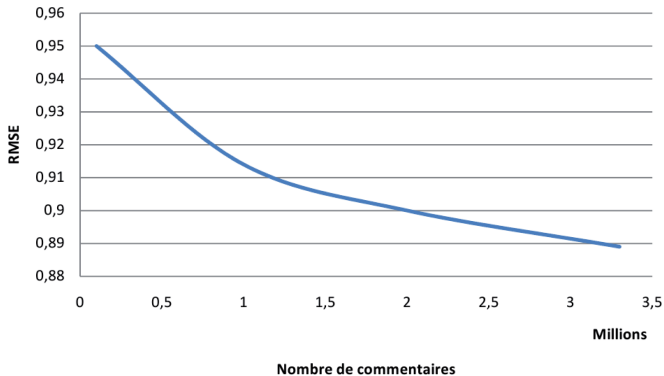
$n$  = nombre total de notes prédites

# Résultats



# Résultats

Une amélioration possible et simple :  
Augmenter le nombre de commentaires présents dans la base



## Contexte et problématique

## Description de la méthodologie

- Extraction des textes

- Inférence de notes sur les textes

- Recommandation par filtrage collaboratif

## Évaluation

## Conclusion

# Bilan

## Objectif atteint !

Démonstration faite que les textes non-structurés issus de blogs peuvent pallier le manque de données en recommandation

- ▶ Le choix de transformer les textes en données d'usages s'est avéré pertinent  
⇒ Meilleures performances qu'avec les descripteurs IMDb
- ▶ Validation et évaluation de chaque élément de la chaîne
- ▶ Étapes 2 (inférence des notes) et 3 (recommandation) automatisables et applicables à d'autres domaines que les films

# Perspectives

- ▶ À court terme :
  - ▶ Évaluation de la chaîne sur d'autres domaines que les films et d'autres langues que l'anglais
- ▶ À plus long terme :
  - ▶ Automatisation complète de la chaîne
    - ▶ Détection automatique des commentaires sur les sites Web
    - ▶ Extraction automatique du sujet traité dans le commentaire
  - ▶ Extraction d'informations sociales (liens entre utilisateurs)



**Merki !!!**  
Des questions ? :-)