



**HAL**  
open science

# Data-driven AI techniques for fashion and apparel retailing

Chandadevi Giri

► **To cite this version:**

Chandadevi Giri. Data-driven AI techniques for fashion and apparel retailing. Machine Learning [cs.LG]. Université de Lille; Högskolan i Borås (Suède); Soochow University (Suzhou, China), 2021. English. NNT : 2021LILUB012 . tel-03533483

**HAL Id: tel-03533483**

**<https://theses.hal.science/tel-03533483>**

Submitted on 18 Jan 2022

**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.



UNIVERSITY OF BORÅS  
SCIENCE FOR THE PROFESSIONS



## THESE

En cotutelle avec

l' Université de Lille (France), l'Université de Borås (Suède) et l'Université de Soochow (Chine)

Présentée en vue d'obtenir le grade de

Docteur de l'Université de Lille

**Spécialité : Automatique, Génie informatique, Traitement du Signal et des Images**

Par

**Chandadevi GIRI**

Titre de la thèse :

**Techniques d'IA axées sur les données destinées au secteur de la mode et de l'habillement**

*Soutenue le 15 Octobre 2021 devant la commission d'examen*

**Jury:**

Prof. Christine BALAGUE	IMT-Business School, France	Rapporteur
Prof. Thomas MULLERN	Jönköping International Business School, Suède	Rapporteur
Prof. Margherita PAGANI	SKEMA Business School, Paris, France	Examineur
Prof. Zhenzhen ZHAO	SKEMA Business School, Paris, France	Examineur
Prof. Yan CHEN	l'Université de Soochow (Chine)	Examineur
Prof. Sebastien THOMASSEY	ENSAIT, France	Examineur
Prof. Xianyi ZENG	ENSAIT, France	Directeur de thèse
Prof. Ulf JOHANSSON	l'Université de Borås ,Suède	Co-directeur de thèse



UNIVERSITY OF BORÅS  
SCIENCE FOR THE PROFESSIONS



## THESIS

*Jointly organized by*

University of Lille (France), University of Borås (Sweden) and Soochow University (China)

Presented in view to obtain a doctoral degree at the University of Lille

**Specialty: Automation, Computer Engineering, Signal and Image Processing**

by

**Chandadevi GIRI**

Title of the thesis:

**Data-driven AI Techniques for Fashion and Apparel Retailing**

*Defended on 15 October 2021*

### PhD Committee:

Prof. Christine BALAGUE	IMT-Business School, France	Reviewer
Prof. Thomas MULLERN	Jönköping International Business School, Sweden,	Reviewer
Prof. Margherita PAGANI	SKEMA Business School, Paris, France	Examiner
Prof. Zhenzhen ZHAO	SKEMA Business School, Paris, France	Examiner
Prof. Yan CHEN	Soochow University, China	Examiner
Prof. Sebastien THOMASSEY	ENSAIT, France	Examiner
Prof. Xianyi ZENG	ENSAIT, France	Thesis Director
Prof. Ulf JOHANSSON	University of Borås, Sweden	Thesis Co-director



# Abstract

## Data-Driven AI Techniques for Fashion and Apparel Retailing

Digitalisation allows companies to develop many new ways of interacting with customers and other stakeholders. These digital interactions typically generate data that can be stored and later processed for different objectives. Currently, the fashion and apparel industry is undergoing a disruptive transformation due to digitalisation, including a rapid increase in the generation of data in various parts of the supply chain. While most data may not be stored with data mining or other analyses in mind, collected data frequently contain very valuable information that can be exploited. Analytics, in particular the use of data-driven AI techniques, is therefore becoming a pervasive tool that is used for a large variety of purposes and in many different processes. While the popularity of Artificial Intelligence (AI) as an advanced tool for improved decision support is increasing, applications of AI within the fashion and apparel industry have historically been rather limited.

With this in mind, the overall purpose of this thesis is to, after presenting an overview of research on applications of data-driven AI in the fashion and apparel industry, demonstrate how various data sets and AI techniques can be utilised for improved decision support in different scenarios.

Whilst the thesis first investigates the impact of AI on different parts of the supply chain, the empirical work focuses on fashion and apparel retailing. Here, different AI techniques are explored in a set of case studies covering several applications in fashion and apparel retailing, thus showing the potential of data-driven AI for decision support in that domain.

One important learning outcome, found in several of the studies, is the need to combine several data sources and techniques in the projects. Another takeaway is the benefit of interpretable models, which allow for inspection and analysis of the discovered relationships. From an applied perspective, approaches like RFM modelling can be utilised as a pre-step to predict customer churn, add sentiment analysis to short-term sales forecasting and build campaign and simulation engines from historical data, which could potentially be used by many retailers.

In conclusion, this thesis has, mainly through a set of case studies addressing real-world problems and utilising real-world data sets, demonstrated how data-driven AI techniques can support and improve fashion and apparel retailers' decision-making.

**Keywords:** Digitalisation, artificial intelligence, fashion and apparel industry, churn prediction, sales forecasting, campaign analysis, data-driven AI decision-making

# Sammanfattning

## Datadrivna AI-tekniker för mode- och klädhandeln

Digitaliseringen möjliggör nya sätt för företag att interagera med kunder och andra intressenter. Dessa digitala interaktioner genererar data som kan lagras och senare processeras för olika ändamål. Modeindustrin genomgår just nu en disruptiv transformation på grund av digitaliseringen, vilket har lett till en snabb ökning av den mängd data som genereras i olika delar av värdekedjan. Även om avsikten med merparten av den data som lagras inte är att den ska nyttjas för data mining eller annan datadriven analys, så innehåller den potentiellt värdefull information som kan utnyttjas av företagen. Datadrivna analysmetoder, framför allt AI, har därför fått ett stort genomslag, och används nu för en mängd olika syften och processer. Även om användandet av AI som verktyg för förbättrat beslutsstöd ökar, så har användningen inom modebranschen historiskt sett varit relativt begränsad.

Med detta som bakgrund, är denna avhandling övergripande syfte att presentera en genomgång av forskning kring applikationer av data-driven AI i modeindustrin samt att visa hur olika datamängder och AI tekniker kan användas för att skapa förbättrat beslutstöd i ett antal scenarion.

Även om avhandlingen börjar med en översikt över AI i olika delar av värdekedjan, så är det empiriska arbetet fokuserat på detaljhandeln inom modebranschen. Här utvecklas och prövas olika AI-tekniker i ett antal fallstudier som spänner över en mängd tillämpningar inom modehandeln.

En viktig lärdom från flera av avhandlingens studier är behovet att kombinera olika typer av datamängder och tekniker. En annan generell slutsats är fördelarna med tolkningsbara modeller, vilka möjliggör granskning och analys av identifierade samband. Utifrån ett mer tillämpat perspektiv bör vissa tillvägagångssätt, som till exempel att utnyttja RFM-analys vid churn-prediktion, att berika försäljningsdata med sentimentanalys vid korttidsprognoser samt att skapa ett simuleringsverktyg för kampanjer från historisk data, visa sig värdefulla för aktörer inom detaljhandeln.

**Nyckelord:** Digitalisering, Artificiell intelligens, Modeindustrin, Churnprediktion, Försäljningsprognoser, Kampanjanalys, Datadriven AI, Beslutsstöd

## Résumé

### Techniques d'IA axées sur les données destinées au secteur de la mode et de l'habillement

La digitalisation offre aux entreprises la capacité de développer de multiples et nouvelles modalités d'interaction avec les clients et les autres intervenants. Ces interactions numérisées engendrent typiquement des données qui peuvent être sauvegardées et traitées plus tard pour divers objectifs. L'industrie de l'habillement subit actuellement une transformation disruptive due à la digitalisation, notamment une augmentation rapide du volume de données générées par les différentes parties de la chaîne d'approvisionnement. Tandis que la plupart des données ne sont pas nécessairement stockées dans l'optique de l'exploration des données ou autres formes d'analyse, les données collectées contiennent fréquemment des informations très intéressantes qui peuvent être exploitées. Ainsi, les outils d'analyse, en particulier les techniques d'IA axées sur les données, deviennent omniprésents et sont utilisés à des fins très variées et dans de nombreux processus différents. Bien que la popularité de l'IA en tant qu'outil performant pour une meilleure aide à la décision ne cesse de croître, les applications de l'IA dans le secteur de la mode et de l'habillement restent relativement limitées. Dans cette optique, l'objectif général de cette thèse consiste, après avoir donné un aperçu général des applications de l'IA basée sur les données dans le secteur de la mode et de l'habillement, à démontrer la possibilité d'utiliser divers ensembles de données et techniques d'IA pour améliorer l'aide à la décision dans le cadre de divers scénarios.

Cette thèse examine initialement l'impact de l'IA sur différentes parties de la chaîne d'approvisionnement, le travail empirique se concentre sur le secteur de la mode et de l'habillement. Différentes techniques d'IA sont alors explorées dans une série d'études de cas couvrant plusieurs applications dans ce secteur, révélant ainsi le potentiel de mise en œuvre de l'IA basée sur les données pour la prise de décision.

L'un des acquis importants tiré de diverses études est la combinaison de multiples sources de données et de techniques dans les projets. Un autre constat général est l'avantage des modèles interprétables, permettant l'inspection et l'analyse des corrélations trouvées. Sur le plan pratique, certaines approches, comme la prédiction du comportement de désabonnement des clients, l'ajout de l'analyse des sentiments aux prévisions de ventes à court terme et la construction d'un moteur de simulation et de publicité à partir de données historiques, pourraient être utilisées par de nombreuses enseignes de prêt-à-porter.

En guise de conclusion, la présente thèse a démontré, principalement par le biais d'une série d'études de cas abordant des problèmes du quotidien et utilisant des bases de données réelles, que les techniques d'IA axées sur les données peuvent soutenir et améliorer la prise de décision du secteur de la mode et de l'habillement

**Mots-clés :** Digitalisation, intelligence artificielle IA, industrie de la mode et de l'habillement, prédiction de désabonnement, prévision des ventes, analyse des promotions, Prise de décision par IA axée sur les données

# 时尚和服装产品零售业的数据驱动人工智能技术

## 摘要

数字化技术能够为服装公司提供了多种新方式与客户和其他利益相关者进行互动。这些数字交互通常会产生的数据，可以将这些数据进行存储和处理以实现不同的目标。目前，时装及成衣业正经历数字化技术带来的颠覆性转变，其中包括在供应链各环节中生成的数据量急剧增加。

尽管大部分数据可能不是为了数据挖掘或其他分析而进行存储的，但通常收集获得的数据中包含可以利用的、非常有价值的信息。各种算法，尤其是使用数据驱动的人工智能技术的算法，正在成为一种普遍应用的工具，正在用于各种各样场合和许多不同过程。虽然人工智能作为改善决策支持的先进工具越来越受欢迎，但人工智能在时装和服装业的应用历来相当有限。

本论文的总体目标是在有关数据驱动人工智能在时装和服装行业的应用研究成果基础上，分析利用各种数据集和人工智能技术针对不同情景下的优化提供决策支持。

本文首先研究了人工智能对供应链不同环节的影响作用，针对服装和服装零售业进行实证研究工作。对人工智能技术在一系列案例进行分析研究，这些案例研究涵盖了时装和服装零售领域的不同应用场景，从而证实了利用数据驱动的人工智能在该领域进行决策支持的应用潜力。

本研究的重要成果之一就是从小项相关案例分析中发现需要结合多个数据来源和多项技术来进行项目研究。另外，研究结果证实了解释模型的优点，能够通过检查和分析发现相互之间的关系。从应用的角度来看，本研究提出的一些方法，例如预测客户流失、在短期销售预测中加入情绪分析、以及根据历史数据建立活动和模拟引擎，可以满足于许多零售商的使用需求。

总之，本研究主要通过针对一系列现实问题的案例研究，提出解决方案，并利用真实的数据集合，论证了数据驱动的人工智能技术在支持和优化服装零售商决策方面的应用。

**关键词：**数字化，人工智能，服装产业，客户流失预测，销售预测，竞争分析，数据驱动的人工智能决策



## Acknowledgements

As much as I had imagined that this PhD journey would go smoothly, there were definitely a lot of downs that—I am sure all of us have experienced in these tragic pandemic times—made me stumble a little. At times, when I felt my willpower was drying up and my confidence was evaporating, there were amazing people who unconditionally supported me and made this difficult journey a smooth, joyful and memorable one. Before acknowledging them, I would like to extend my gratitude to the European Commission’s SMDTex project, which funded this PhD and made this research possible.

Firstly, I would like to thank Prof. Ulf Johansson, my main supervisor and a great mentor, whom I cannot thank enough for his guidance, patience and all the pains he took to ensure that I stay this difficult course and complete my thesis. I thank my co-supervisor, Jenny Balkow, who constantly taught me to conduct research from different lenses and to be more careful about every tiny detail. She supported me during the most difficult times and inspired me to always do better and improve. I truly enjoyed working with you both. I shall always remember our creative discussions and friendly collaborations, for which I owe you my utmost gratitude. I also thank Prof. Tuwe Löfström for providing me with the great opportunity to collaborate with him for the research and for sharing his valuable perspectives, from which I learned a lot and felt even more confident.

I would like to thank Prof. Xianyi Zeng, my main supervisor at ENSAIT, who put his confidence in me throughout this four years’ journey and has immensely helped me in finally accomplishing this daunting yet exciting project. I thank my co-supervisor, Prof. Sebastien Thomassey, for his valuable guidance, encouragement, innovative perspectives, and most importantly, for taking all the efforts to acquire the data for my PhD research, without which this thesis would not have been possible. I am grateful to Prof. Yan Chen and Prof. Lichuan Wang for the smooth collaboration, even in uncertain times, and for providing their continuous feedback and support during the remote mobility.

I greatly appreciate the support provided by the Evo Pricing and Ellos, and I am thankful to them for providing the data to carry out my PhD research. I would also like to express my sincere gratitude to Eva Gustafsson (the director of studies), Tina Carlson Ingdahl (the head of department), and Petri Granroth (research officer) at the University of Borås for extending consistent support in all academic and administrative matters and for being extremely helpful. I would additionally like to thank the amazing staff at ENSAIT for providing me with timely support whenever required. I would like to thank Marie Hombert for making my life easier during my stay in France by always helping me with patience and great compassion.

My amazing SMDTEX colleagues have been a great support, and I thank them from the bottom of my heart for lending their shoulders when I needed them most. Especially, I would like to thank Sheenam

for always being there to help me and for being a wonderful collaborator in my research work. I am grateful to Balkis, Melissa, Sweta, Zenglei, Mengru, Vijay, Tarun, Neaz, Ajinkya, Ashik and Shahood for all your support and for all the joyful conversations, a good excuse for keeping PhD work aside. I am thankful to my wonderful PhD colleagues, Ann, Emelie, Lars and Sara, for our exciting conversations during the brainstorming sessions and to Vaishnavi for cheering me up during our wonderful walks and coffee meets and for always being helpful. I am also grateful to the beautiful nature in my vicinity that embraced me in its lap and kept me engaged and motivated to always look forward to better times.

Finally, I would like to thank my father, Karunshankar Giri, and mother, Vidyavati Giri, for their sacrifices and unconditional love, and I dedicate this PhD thesis to them for believing in me and supporting me in my decisions. I thank my sister Manju, and my lovely brothers Shivam and Rishi for all their love, irritation, and moral support. You have been great fun to talk and fight with, and forgive me for staying away from you due to my studies and career. I thank my nephew Anirudh for lighting up my day with his cute smile and funny hellos.

Chandadevi Giri

Borås, Sweden

September, 2021

## List of Appended Papers

This thesis is built upon five research papers (below), which are enclosed at the end of the manuscript. In this thesis, these research papers are also called case studies or research articles.

### *Paper I*

**Giri, C., Jain, S., Zeng, X., & Bruniaux, P.** (2019). A detailed review of artificial intelligence applied in the fashion and apparel industry. *IEEE Access*, 7, 95376–95396.

### *Paper II*

**Giri, C., Thomassey, S., & Zeng, X.** (2019). Customer analytics in fashion retail industry. In *Functional textiles and clothing* (pp. 349–361). Springer, Singapore.

### *Paper III*

**Giri, C., & Johansson, U.** (2021). *Data-driven business understanding in the fashion and apparel industry*. In Proceedings of the 18th International Conference on Modeling Decisions for Artificial Intelligence (MDAI, 2021), Umeå, Sweden.

### *Paper IV*

**Giri, C., Thomassey, S., & Zeng, X.** (2019). Exploitation of social network data for forecasting garment sales. *International Journal of Computational Intelligence Systems*, 12(2), 1423–1435.

### *Paper V*

**Giri, C., Johansson, U., & Löfström, T.** (2019). Predictive modeling of campaigns to quantify performance in fashion retail industry. In: *2019 IEEE International Conference on Big Data (Big Data)* (pp. 2267–2273). IEEE.

## Table of Contents

<b>1</b>	<b><i>Introduction</i></b>	<b>1</b>
1.1	Problem statement	4
1.2	Research purpose and research questions	5
1.3	Thesis outline	6
<b>2</b>	<b><i>Conceptual framework</i></b>	<b>7</b>
2.1	Data-driven AI	7
2.2	Fashion, apparel and textile management	34
2.3	Data-driven decision-making in the F&A industry	42
<b>3</b>	<b><i>Research methodology</i></b>	<b>43</b>
3.1	Research strategy	43
3.2	Structure of PhD programme	45
3.3	Datasets	46
3.4	Methods	48
3.5	Model evaluation	50
<b>4</b>	<b><i>Summary of the appended papers</i></b>	<b>52</b>
4.1	Paper I: A detailed review of artificial intelligence applied in the fashion and apparel industry	53
4.2	Paper II: Customer analytics in fashion retail industry	58
4.3	Paper III: Data-driven business understanding in the fashion and apparel industry	59
4.4	Paper IV: Exploitation of social network data for forecasting garment sales	60
4.5	Paper V: Predictive modelling of campaigns to quantify performance in fashion retail industry	62
<b>5</b>	<b><i>Conclusions</i></b>	<b>64</b>
5.1	Concluding RQ1	64
5.2	Concluding RQ2	64
<b>6</b>	<b><i>Discussion and future work</i></b>	<b>68</b>
	<b><i>References</i></b>	<b>69</b>
	<b><i>Appended papers</i></b>	<b>75</b>

## List of Figures

Figure 1.1: Thesis structure outline	6
Figure 2.1: The KDD process	7
Figure 2.2: Virtuous cycle of data mining	8
Figure 2.3: The decision tree	18
Figure 2.4: The random forests	21
Figure 2.5: The gradient boost model	22
Figure 2.6: Fuzzy triangular membership	25
Figure 2.7: An example of a confusion matrix	28
Figure 2.8: A sample ROC curve	31
Figure 2.9: The traditional F&A supply chain	37
Figure 2.10: Product lifecycle in the F&A industry compared to other markets	41
Figure 3.1: Experimental design of the thesis	44
Figure 4.1: Overall distribution of articles by data types	55
Figure 4.2: Distribution of articles by data and F&A supply chain stages	56

## List of Tables

Table 2.1: Sample customer transaction data _____	10
Table 3.1: PhD program by mobility _____	46
Table 3.2: A summary of data _____	48
Table 3.3: Summary of tasks and techniques by research paper _____	50
Table 3.4: The evaluation metrics used in the appended papers _____	51
Table 4.1: Categorisation of the appended papers according to the research questions _____	52
Table 4.2. Scope of the appended papers _____	52
Table 4.3: Contributions of the authors in the appended papers _____	53
Table 4.4 : Categorisation of the data used in the articles as per the F&A supply chain stages _____	56

# 1 Introduction

The information age and the advanced data-driven digital technologies of the 21st century, marked by fast-paced transformations across industries, have led to the invention of many digital artefacts that we use in our daily lives, such as smartphones, self-driving cars and smart translators. These innovations have been made possible by significant advances in Artificial Intelligence (AI) (Haenlein & Kaplan, 2019). In the past decade, AI has been widely applied in many industries, including health, transportation and manufacturing (Thrall et al., 2018; Xiong et al., 2015; Xu & Hua, 2017) as well as in service industries, such as telecommunication, banking, finance and insurance (Olson, 2007). The digital revolution has had many positive impacts on the consumer in terms of convenience, relevance experience, savings and empowerment (Reinartz et al., 2019). It has opened up new ways for companies to interact with the consumer. For example, enhancing consumer online interaction channels through mobile phones allows for new ways of brand value co-creation involving multiple value chain stakeholders (Ramaswamy & Ozcan, 2016). Overall, the digital transformation has facilitated interaction between different actors, such as customers and other stakeholders, across the value chain segments and the new digital technologies (Schallmo et al., 2020). The impact of this digital transformation not only actively links actors in the supply chain network but also makes available new combinations of resources and creates new bonds among actors, who were not previously linked, thus facilitating value creation by optimising and efficiently coordinating business processes and by enhancing customer experiences (Pagani & Pardo, 2017). The diffusion of digital technologies, while forcing some industries to challenge themselves and succeed in the race to gain competitive advantages by embracing digital transformation (Fitzgrald et al., 2013).

In addition to making the purchases more convenient for the customer and making interaction among multiple stakeholders within the value chain possible, digital interactions also generate data that can be stored and later processed. Customer interactions have led to a high amount of customer-generated data, both in terms of content and of activities, that goes beyond the pure transaction. Rapidly generated data can be effectively stored, organised, processed and analysed with the help of digital technologies; therefore, they are becoming critical to the management of supply chain processes (OECD, 2019). The size of the retail sector, along with its dynamic nature and above-average availability of data, makes the retail industry

especially interesting for researchers on big data analytics and AI, the combination of new technologies along with big data/predictive analytics is expected to cause a quantum leap in retailers' understanding of the shopping process (Dekimpke, 2020). For new digital technologies to have a greater impact on businesses, there is a demand for new skills. A study by (Schallmo & Williams, 2018) suggest that new skills can help with acquisition, analysis and conversion of data and with generating actionable insights from it, which will transform the existing and traditional processes, models and decision-making activities, and further, will improve the performance of different operational activities within the firms.

Just as in the rest of retailing, the fashion and apparel (F&A) industry is also undergoing a digital transformation. As argued by Griva et al. (2018), digitalisation has necessitated that F&A apparel retailers respond to the dynamic preferences and tastes of modern consumers. This is because consumer decisions are influenced by increased awareness gained through various offline and online channels (Griva et al., 2021). In addition, customer satisfaction is critical to building brand loyalty. Therefore, the F&A industry competes to retain its customers for increasing business revenue by providing them with high-quality products at the right time (Grabher et al., 2008). Traditional statistical decision models, such as segmentation and customer lifetime value (CLV), are incapable of handling the challenges that arise from the uncertain customer behaviour and market trends and from the complex data (Dahana et al., 2019). Therefore, they fail to support the accurate and improved decision-making that is integral to the success of the F&A retail industry. Consequently, novel decision models for real-time customer analytics and consumer behaviour prediction that are based on newly available customer transaction data are indispensable to the efficient management and growth of the F&A retail business.

Due to the growing impact of digital technologies, the F&A industry is becoming increasingly dynamic, as data is being generated when every new piece of apparel is designed, produced and sold (Jain et al., 2017). Powered by advanced digital technologies, such as cloud computing and database technologies coupled with AI, the F&A industry is able to integrate customers into the value chain segments through the adoption of e-commerce platforms. The customer actively adds value to the digitalisation of the F&A industry because the customer data can be used to improve products and services (Bhardwaj & Fairhurst, 2010).



AI is one of the technologies listed by the OECD (2019) as a main driver of digital transformation. The innovative transformations in the technology sector have been driven by significant advances in AI, a technology with the potential to solve complex modern problems (Haenlein & Kaplan, 2019). Today, AI has become an important topic in multiple domains, which means that definitions of the term might differ. As the pertinent research studies take a data analytics business perspective, the following definition will be considered through this thesis. AI, as per Kaplan and Haenlein (2019, p.17), is “a system's ability to interpret external data correctly, to learn from such data, and to use those learnings to achieve specific goals and tasks through flexible adaptation”.

The popularity of AI applications and the rapid growth of data from various sources, such as retail transaction data and social media data, offers plenty of opportunities to develop decision models based on AI that could enhance operation and management efficiency in the F&A industry (Zou et al., 2019). Hence, this thesis aims to exploit F&A retail data using AI techniques in order to improve decision-making in the F&A industry. As discussed previously, new data sources are valuable repositories of insights and knowledge that could aid F&A retailers in understanding customer behaviour in real-time and in increasing the quality and precision of their supply chain activities.

Thus, AI and big data are key technologies for making efficient, accurate and robust decisions in various business domains to maximise business profitability. Therefore, significant investments in such technologies are integral to the competitiveness of the F&A industry, as they enhance the efficiency of critical decisions, such as demand forecasting, product recommendations, customer retention, market trend prediction and customisation (Thomassey & Zeng, 2018).

Hence, AI is an obvious technology for exploiting the massive amount of data and deriving significant knowledge from it. As companies invest in AI to gain a competitive edge in the market, it is crucial for the F&A industry to tap into this knowledge and leverage it to make more informed and accurate decisions. However, the decision domains within the F&A industry that serve as crucial avenues for studying AI's potential remain limited. For instance, fashion retail data is a significant source for the study of customers' purchasing behaviour that could, in turn, enable F&A retailers to optimise their decisions regarding the sourcing of raw materials and manufacturing of F&A products.

However, the F&A industry has not yet extensively adopted the cutting-edge data technologies, such as AI and big data analytics, that are necessary for solving critical management-related decision problems, such as those in sales, retailing and marketing (Acharya et al., 2018). The understanding of the decision context and data is a main prerequisite of AI applications because they are both important elements in enhancing an AI's algorithmic intelligence (Stuart & Norvig, 2009). Therefore, it is essential for the F&A industry to collect data from all business transactions by customers, suppliers and other supply chain actors. In this respect, the F&A industry is facing a need to adapt to the rapid phenomenon of big data generation in their information repositories in order to derive significant insights from the data and make efficient decisions for improved business growth and performance (Bradlow et al., 2017; Ferraris et al., 2019)

## **1.1 Problem statement**

As mentioned above, the highly digitalised nature of the F&A industry generates massive amounts of data at various parts of its supply chains. Digitalisation has been considered as the most influential factor in enabling the F&A industry to better understand customer preferences and align them with the upstream supply chain stages, such as raw material sourcing, manufacturing, distribution, logistics, and sales and retailing, through centralised digital platforms. Data play an essential role in these processes, as the digital interactions among these supply chain stages, including customers, lead to massive data proliferation.

These rapidly growing data sets could, through advanced analytics, lead to improved support in a variety of decision processes. The interconnected and centralised information systems that store the datasets allow the F&A industry to fully understand and map out the range of decision problems related to supply chain operations and further allow them to make informed decisions using advanced data analytics. These decisions' efficiency and accuracy are of great importance for the sustained growth and success of the F&A industry.

In the above context, traditional decision models do not address the complexity of the newly generated datasets or extract insights from them to support decisions. These datasets are the significant factors in the decision-making processes of the F&A industry. Consequently, data-driven AI has replaced traditional statistical methods in various decision areas. However, a

systematic categorisation that simultaneously incorporates tasks addressed, data sets used and AI-techniques utilised in different phases of the F&A supply chain is missing. Therefore, it is worthwhile to investigate how decisions and the utilised analytical models in the F&A industry have evolved with respect to increased digitalisation, data proliferation and AI influence.

There also seems to be a lack of understanding of how the combination of data from different sources, using data-driven analytics, can support managers in decisions of strategic importance in the F&A industry. Since the supply chain processes in the F&A industry are data-driven, the complex problem of utilising data to make accurate and informed decisions has not yet been addressed in the literature.

## 1.2 Research purpose and research questions

In line with the problem statement outlined in the previous section, the overall purpose of this thesis is:

- to **study** existing research on applications of data-driven AI in the F&A industry in order to identify the data types and techniques that are frequently used.
- to **demonstrate** how different data sources and data mining techniques can be used and combined in novel ways to improve decision support in the F&A industry.

The problem statement discussed in Section 1.1 leads to the following research questions (RQs) that aim to achieve the purpose of this thesis:

1. What are the widely used data types and key data-driven AI applications in the F&A industry?
2. How can data-driven AI techniques support and improve F&A retailers' decision-making?
  - a. How can fashion customer transaction data be utilised to analyse and predict customer behaviour?
  - b. How can social media data be utilised to improve sales forecasts?
  - c. How can historical campaign data be utilised to design successful campaigns and promotions?

### 1.3 Thesis outline

This thesis is composed of six chapters and five appended papers. The outline of the thesis is shown in Figure 1.1.

**Chapter 1** presents the introduction, research problem, purpose and research questions.

**Chapter 2** introduces the conceptual framework of the thesis by describing the key concepts used in the appended articles. Section 2.1 describes the key concepts of data-driven AI and different terminologies that support the thesis work. This section explains the data, data mining tasks, techniques and model evaluation. Section 2.2 provides the key concepts related to the F&A industry, supply chain and retailers, and finally, Section 2.3 discusses data-driven decision-making.

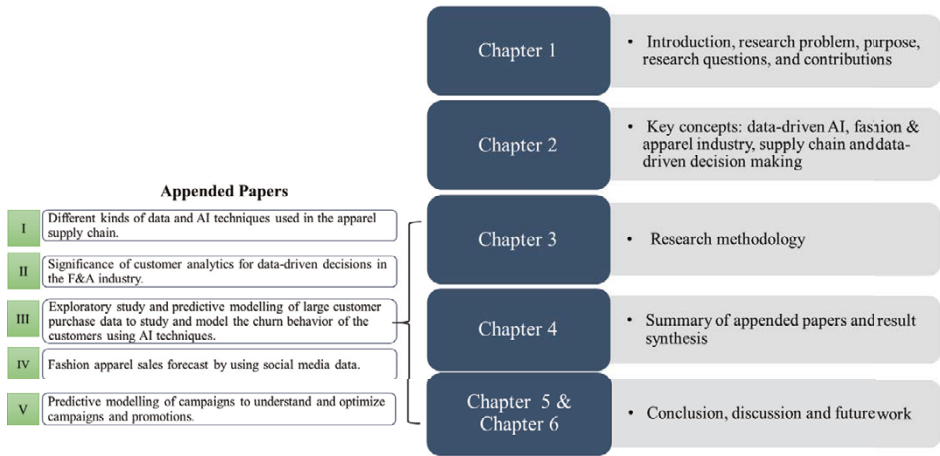


Figure 1.1: Thesis structure outline

**Chapter 3** explains the research methodology by outlining the research strategy, research process, datasets and scientific methods used to conduct the thesis research work.

**Chapter 4** summarises the five appended articles independently.

**Chapter 5** presents the conclusion by answering and discussing RQs and purpose of the thesis.

**Chapter 6** presents a discussion and suggests possible future research avenues.

## 2 Conceptual framework

The theoretical foundation of the papers and articles in this thesis are developed from two fields of research—data-driven AI and textile management. The aim of this chapter is to create a conceptual framework for discussion of the results. Thus, the chapter will first cover the main concepts of data-driven AI (Section 2.1) and then those of textile management in the context of the F&A industry (Section 2.2).

### 2.1 Data-Driven AI

The primary goal of data-driven AI is to utilise data to provide value in the form of knowledge and to support automated decision-making with or without human intervention (Russell & Norvig, 2009; Thamm et al., 2020). Data-driven AI systems contribute to the process of knowledge extraction from large quantities of data through activity referred to as data mining.

Data mining, as defined by Berry and Linoff (2004), is a business process for exploring large amounts of data to discover meaningful patterns and rules. The overall goal of data mining is to turn data into actionable insights that support decision-making across the organisations. In the varying contexts wherein data mining is applied, the goal is often to create value from the data. Therefore, data mining is a goal-driven business process. The data play a central role in the process of knowledge extraction.

Data mining constitutes the key step in the larger process known as *knowledge discovery in databases* (KDD). KDD is also a computer science field that mainly focuses on the extraction of hidden patterns and rules from data, enabling the discovery of new knowledge. According to the standard CRISP-DM methodology for data mining developed by Chapman et al. (1999), the KDD process is comprised of three broad phases, as illustrated in Figure 2.1.



Figure 2.1: The KDD process

Data from a multitude of sources and in varying formats constitute the main input to the KDD process. Data pre-processing precedes the data mining step because data need to be transformed into an appropriate format before being utilised for data mining. In the data pre-processing step, relevant data from the different sources are combined, transformed into a structured format and cleaned to remove missing and extreme values, also known as *outliers*. In this way, the data pre-processing step produces data in a standard format that can be used for descriptive or predictive data mining. The data post-processing step then ensures the selection of only the valid and useful results from the data mining step, and it often entails data visualisation and statistical analysis, such as hypothesis testing, data summarisation, etc. The specific context and purpose of the decision task often drive the KDD process; therefore, KDD is a context-driven process. From the business perspective, there are two crucial goals: firstly, to identify specific business decision areas and problems where data mining can produce value, and secondly, to extract actionable insights out of data, using data mining techniques, in order to support decision-making, i.e. to actually perform KDD. These two activities are part of a larger process known as the virtuous cycle of data mining, illustrated in Figure 2.2 (Berry & Linoff, 2004). Depending on the context, data mining can refer to either the complete virtuous cycle or a part of the KDD process. In this thesis, the application of data mining techniques to extract valuable and actionable knowledge from the data is considered within the business context.

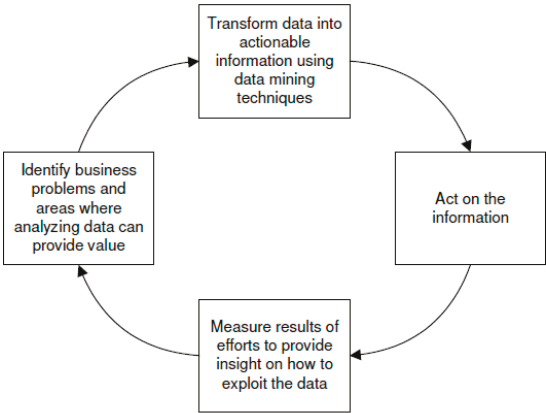


Figure 2.2: Virtuous cycle of data mining (Berry & Linoff, 2004)

Although data mining can be applied to a number of decision problems, these problems can almost always be mapped into a few generic data mining tasks. On a high level, these tasks

can be divided into descriptive and predictive tasks (Pang-Ning Tan et al., 2016). The overall purpose of the descriptive task is to perform an exploratory analysis that aims to produce meaningful information about the data in terms of high-level data summarisation, identification of regularities or patterns and relationships between the attributes. The data mining techniques used for descriptive tasks typically include association rules and clustering (Williams, 2011). The main purpose of the predictive task is to predict or explain a value of one variable, often its future value, based on the values of other variables (Hand, 2007). Predictive models are often built upon historical data, i.e. a set of examples consisting of input and target variable values (Peng et al., 2008).

Data mining tasks, both descriptive and predictive, require high-quality data that is relevant and has been pre-processed for the task (Berry & Linoff, 2004; Piatetsky-Shapiro & Parker, 2011). A key capability of organisations utilising data mining techniques is therefore the ability to map a business problem to one of the generic tasks while also ensuring to have access to the necessary data in a suitable format. Once the problem is identified and the data is harvested and subsequently pre-processed and transformed into a suitable format, the actual algorithm needs to be chosen based on the exact task (e.g. classification, forecasting, text mining, etc.), the characteristics of the problem and any other requirements (Kuhn & Johnson, 2013). Next, the evaluation of the chosen data mining techniques is performed to ensure that the objectives of the data mining tasks are fulfilled. Evaluation of data mining is actually two different activities; *internal evaluation* uses rigorous protocols and metrics to compare different techniques and parameter settings, while *external evaluation* somehow measures the effects of the entire project.

In the next section, we will present in detail a number of concepts that are related to data-driven AI: data, tasks, techniques, and model evaluation. This section will elaborate on the data mining process and give some insights into the different terminologies that are used in the data mining-based case studies.

### **2.1.1 Data**

Any data mining task requires data. Data come from different sources (such as relational databases, transactions, sensors, social media, demographic data, etc.) and in multiple formats, including text, numbers, videos, images, etc. In general, data can be *structured*, *unstructured* or *semi-structured*:

- Structured Data – Data in a standard format or structure, e.g. data from relational databases.
- Unstructured Data – Data without any defined structure that is stored in different formats, e.g. text documents, social media data, images and videos.
- Semi-Structured Data – Data with a discernible pattern or data that are only partially structured but contain tagged information, e.g. XML, JSON documents and emails.

In the business context, e.g. the F&A industry, the majority of all data analysis uses structured data, often in a standard tabular format. The table representing the data contains rows and columns. Data in a tabular format entails data objects, also known as *records*, *observations* or *instances*. These data objects, which are represented in the rows, are described by a number of *features*, also known as *attributes* or *variables*, that are represented in the columns. For example, Table 2.1 shows customer transaction data; and the attributes represented in the columns describe the customer.

Table 2.1: Sample customer transaction data

Customer	Age	Gender	Date	Product _Name	Purchased_Items	Cost (Euros)
Cust_1	25	M	26/02/2018	Shirt	1	25
Cust_2	38	F	26/02/2018	T-shirt	4	10
Cust_3	40	M	26/02/2018	Jeans	2	30

In the above table, the instance (*Cust\_1*) is described by attributes (*Age*, *Gender*, *Date*, *Product\_Name*, *Purchased\_Items* and *Cost*) that represent characteristics that vary from one customer to other. The attributes in the table are both numerical (*Age*, *Date*, *Purchased\_Items*, and *Cost*) and categorical (*Gender* and *Product\_Name*). Numerical attributes are represented by a number and exhibit the properties of numbers, while categorical attributes represent the qualitative aspect of the data objects (in Table 2.1, the customers). The attributes *Gender* and *Product\_Name* take up only a small number of possible values (*Male(M)* and *Female(F)* and *Shirt*, *T-shirt* and *Jeans*, respectively).

Detailed descriptions of the most common data attribute types are given below:

- Numeric attributes: The *numeric* or *quantitative* attributes are represented by numbers and possess most of the properties of numbers. Numeric attribute can be *continuous* or



*discrete*. Continuous attributes take up values of real numbers, e.g. temperature, weight or height. Discrete attributes take up a finite or an infinite set of integer values, e.g. ID numbers, zip codes or number of days. A binary attribute is a special discrete attribute that takes up only two values, e.g. Male/Female, True/False, Yes/No, 0/1, etc. The properties of numbers, such as *addition* (or *subtraction*) and *multiplication* (or *division*), define the two common numeric attributes, viz. *Interval* and *Ratio*. The unit of measurement, for example the sum or difference between the values defined by addition and subtraction, are meaningful for *interval* attributes (dates in the calendar, temperature measured in Celsius or other units, etc.). For *ratio* attributes, both the differences and ratios defined by the multiplication or division operations are meaningful. For example, a variable such as *length* can be used to compare or measure the objects by their lengths. Other examples of *ratio* attributes include *age*, *quantities*, *mass*, and *electric current*.

- **Categorical attributes:** The categorical attributes lack most of the number properties and have a well-defined set of values without any natural order, e.g. country, gender, etc. In general, the values of categorical attributes are represented by strings that, if required, can be transformed into numbers. Defined by the specific number properties, viz. *order* and *distinctness*, the categorical attributes are categorised as *nominal* and *ordinal*. The values of *nominal* attributes are just different names that help to distinguish one value from the other, for example colour, gender, city names, etc. On the other hand, *ordinal* attribute values provide information to order objects, for example the performance of students (poor, average, better, excellent), grades, rankings, etc.

Apart from the structured data, there exist other special data types, such as temporal data (i.e. time series data) and unstructured data (e.g. text data), described below, that do not necessarily fit into the standard tabular format. Hence, they need to be transformed into a structured format for data mining tasks.

- **Time series data:** Sequential data represented by a series of measurements recorded over time, e.g. stock data or housing prices that vary over time.
- **Text data:** Unstructured data containing words or text, for example social media data and documents. These data need to be transformed into numerical values to train the data mining algorithms.

Using poorly processed data for the data mining tasks can lead to inappropriate results, and it is normally not possible if the data is not in a suitable format (Pyle, 1999). As previously mentioned, data is far from being perfect; therefore, data pre-processing is a crucial step that takes a considerable amount of time. Data must be inspected for quality by taking important steps, such as identifying missing data, detecting noise and outliers, and transforming categorical attributes into a numerical form in order to apply data mining techniques (Pang-Ning Tan et al., 2016). From the perspective of data quality, some problems that can affect the data mining results are described as follows:

- a. **Missing data:** a field without a value, which could occur due to many reasons and can be handled using various methods. Following are some of the methods for handling missing values:
  - Attributes with missing values can be removed, for instance removing attributes with more than 10 % missing values.
  - Filling the missing values, for example replacing them with the mean value or the most frequent value, which is known as mode.
- b. **Noisy data:** irrelevant data that cannot be easily interpreted by machines (Salgado et al., 2016), for example unstructured data, mislabelled attributes or measurement errors in the values of the attributes. The binning method, clustering and regression are some of the methods of handling noisy data (Pyle, 1999).
- c. **Outliers:** a data point that stands out from the rest of the data (Barnett, 1978). Anomalies, discordant and deviants are terms that describe outliers (Aggarwal, 2017). Most commonly, the outliers are handled by replacing them with appropriate values because outliers can influence the performance of predictive models.

For data mining techniques, data often have to be in a numerical format, meaning that if the data mining is applied to text data or data that contain categorical attributes, conversion of text data and categorical values into numbers is a necessary and crucial step.

## 2.1.2 Tasks

The specific tasks that data-driven AI aims to achieve are generally of two types: descriptive and predictive (Pang-Ning Tan et al., 2016). These are explained in more detail in the following sections.

### 2.1.2.1 Descriptive tasks

The data used for analysis or research typically contain hidden underlying relationships among its attributes. It may also contain complex patterns that could help identify the relationships in the datasets. The overall goal of descriptive data mining is therefore to identify hidden relationships between attributes, trends, correlations, clusters, etc. in the data and to explain its general properties. One basic technique used for descriptive tasks is simply to summarise the data by using either statistical measures or graphs. Statistical measures, such as *mean*, *mode*, *median*, *standard deviation*, and *variance*, are single number measures that describe the characteristics of the data attributes (Kuhn & Johnson, 2013). The measures such as *mean*, *mode* and *median* give the central tendencies of a large number of attribute values. The dispersion or spread of values of the numerical attributes is measured by *standard deviation*, *variance* and *range* (Pang-Ning Tan et al., 2016). The degree of symmetrical distribution of attribute values around the *mean* is measured by the statistical measure known as *skewness*. When there are two attributes to be analysed, the measure that captures the linear relationship between them is *correlation*.

Apart from these measures that summarise the numerical characteristics of data, there exist more subtle and complicated characteristics of the data set that also need to be described. Graphical techniques are used to describe these characteristics; they summarise the data in the form of graphs, such as histograms, bar plots, box plots, pie plots, scatter plots and stem and leaf plots (Pang-Ning Tan et al., 2016). Data summarisation becomes more complex and difficult as the dimensionality of the data set increases. One notable set of techniques that provides interactive analysis and enhanced visualisation capabilities to generate summary statistics of a multi-dimensional data set is *Online Analytical Processing* (OLAP; see Berry & Linoff, 2004). OLAP techniques aggregate multi-dimensional data across either attribute values or different dimensions. For example, if the transaction data set contains sales

information (as per attributes such as date of purchase, product, location and customer IDs), the OLAP technique can be applied to generate a data summary that describes the purchase by a particular customer ID by location, product category and date.

Another important task in the descriptive analysis is to divide the instances into small groups (segments) based on similarity measures or a variety of similarity factors. This is known as *segmentation*. The automatic data-driven segmentation of data, typically performed on large volumes of data, is known as clustering (Berry & Linoff, 2004). In clustering, the overall goal is to group the instances into clusters; the instances belonging to the same cluster should be as similar as possible, while each cluster should be as different as possible from the others. For example, clusters of customers can be formed based on similar purchased products. Another example of segmentation is RFM analysis, which is widely used to analyse customer behaviour for marketing purposes (Berry & Linoff, 2004). The RFM technique quantifies customer actions in terms of *Recency (R)*, *Frequency (F)* and *Monetary (M)*. *Recency* is the measure of time elapsed since the last transaction by the customer. *Frequency* is the number of transactions made by the customer in the analysed period, and *Monetary* is the corresponding total spending. One example of segmenting customers using RFM analysis is identifying the most profitable customers by looking for low *Recency*, high *Frequency* and high *Monetary* values (Ballings & Van den Poel, 2012).

For a special type of data without a standard format, for example the point-of-sale data that is often collected at retail counters, a possible descriptive task would be to identify the patterns that represent the strongly associated attributes. This task is known as association analysis. The identified pattern is represented in the form of feature subsets. The goal of association analysis is to efficiently discover the most interesting patterns in the data set. The widely used application of association analysis is called *Market Basket*, wherein point-of-sale data is analysed to find products that are frequently bought together by the customers (Berry & Linoff, 2004; Pang-Ning Tan et al., 2016). For example, association analysis could discover that the customers who buy shampoo also tend to buy conditioner.

### **2.1.2.2 Predictive tasks**

A predictive task generally has the main objective of predicting the value of the *dependent* or *target* variable based upon the values of other variables, known as *independent* or *explanatory*

variables (Pang-Ning Tan et al., 2016). Predictive tasks are most often handled by *predictive modelling* (PM), which is the process of building a model that defines a *target* variable in terms of a function of *explanatory* variables. PM is defined as the problem of approximation of a function  $f$  to map input variables  $X$  (*independent* variables) and the target variable  $Y$ , and the function approximation is on a very high level represented by Eq. (2.1).

$$Y = f(X) \quad (2.1)$$

PM builds a model using labelled data, i.e. examples containing values for both the target variable  $Y$  and the explanatory variables  $X$ , and then applies this model to a novel data set, i.e. the new unseen instances where values of the target variable are unknown, to predict the unknown values of  $Y$ . When PM uses labelled data to train the model using data mining algorithms, this is called *supervised learning*, and the model is said to be *trained*. The data used to build a predictive model is therefore referred to as the *training* set.

Depending on the type of target variable, the predictive tasks are categorised as predictive regression or predictive classification. When the target variable is categorical, then the task is predictive classification, and when the target variable is continuous, the task is predictive regression. Furthermore, if the prediction task is performed on temporal data, i.e. time series data, then the task is known as time series prediction or forecasting. Time series prediction is thus a special case in which the underlying assumption is that the values of a target variable are dependent on the previous values of either the target itself or of the explanatory variables. Typically, there are two types of time series prediction: *univariate* and *multivariate*. In univariate time series prediction, the value of the target variable is predicted based upon only its past values, while in multivariate time series prediction, the value is predicted based upon the historical values of several variables.

The main objective of any of the predictive tasks is primarily to reduce the *error*, i.e. the difference between the predicted value of the target variable and its actual value. This difference is measured using various error metrics that evaluate how well a predictive model fits the training set. The model training on the training set is considered finalised when the error on the training set is sufficiently small, and then the model can be applied to unseen novel data to make predictions. The internal function that measures the model fitness on the training instances is known as the *score function* (Berry & Linoff, 2004; Pang-Ning Tan et al.,

2016). Although there are many score functions, no score function is optimal for all the predictive techniques. The overall goal is to generate a model that generalises well to instances *not* used for building the model. In order to do so, the obvious requirement is that the novel examples come from the same distribution as the training instances. However, these training examples are still just samples from this distribution and may very well contain noise. A powerful algorithm learning the training set “too well”, i.e. more or less memorising the examples, leads to larger generalisation errors. This is called overfitting. Underfitting, on the other hand, occurs when the model is not powerful enough to learn the relationship.

Since minimising the generalisation error is the goal of predictive modelling, being able to produce unbiased estimations of future errors is critical. With this in mind, the available labelled data set is usually split into several subsets before undergoing predictive modelling. Each subset serves a distinct purpose and is used in different ways. Unfortunately, the terminology used for the subsets is not completely consistent. The various subsets and their usage are described below.

*Training set:* This is the subset used for actually building, i.e. training, the predictive model.

*Validation set:* This is the subset used to select the parameters of the model. A validation set enables model selection by scoring it against the subset of data not used during the model building.

*Test set:* This is the subset used for model evaluation. The test set is not used at all during the model building phase. Test data, consequently, serve to estimate the expected results from new data.

*Production set:* This is the data set on which the finalised model is applied by using instances where the true target values are not known.

Training, validation and test data all have known target values for every example. For production data, as mentioned above, the target values are unknown.

In the research literature, the term *production set* is not commonly used, likely because predictive models are rarely applied to examples in which the targets are unknown. Instead, academic research on data mining focuses on results from *test sets* and sometimes uses *validation sets* for model ranking and selection. This terminology is also used throughout this

thesis, i.e. the *training set* is used for the model building, the validation set (if employed at all) for model and parameter selection and the *test set* for producing the final results.

Descriptive and predictive tasks are normally combined in a data-mining project. The descriptive tasks, being exploratory in nature, are the means to gaining insights into the data set, and they also help to identify any noise in the data. Therefore, predictive tasks often utilise different descriptive techniques as a preliminary (exploratory) step toward a better understanding of the data and the problem.

### **2.1.3 Techniques**

The predictive tasks are handled by numerous data mining techniques, and they all come with certain advantages and disadvantages. This section describes the data mining techniques used in this thesis.

#### ***2.1.3.1 Decision Trees***

Decision trees (DT) are predictive models that are used for both classification and regression tasks (Breiman, 2001). Therefore, DTs are categorised as *classification trees* and *regression trees*. DTs have become widely popular for data mining because they are reasonably accurate, fast to train, and most importantly, produce transparent models. DTs are a set of *if-then* rules that are arranged using a tree-like structure.

To illustrate how DTs work, the problem of classifying mammals from non-mammals is illustrated in Figure 2.3. To determine whether the species is mammal or non-mammal, a series of questions can be posed about the characteristics of the species. The first question that could be asked is whether the species is warm- or cold-blooded. If the species is warm-blooded, then it is either a mammal or a bird, otherwise it is definitely a non-mammal. Another follow-up question that could be asked is whether the females of the species give birth to their offspring or lay eggs. Species that lay eggs are most certainly non-mammals, while those who give birth to their young are definitely mammals. This series of questions and answers could be organised into a tree-like hierarchical structure that forms a DT, the components of which are a root node, branches and leaf nodes (see Figure 2.3). A root node has no incoming links but can have zero or more outgoing links. An internal node has exactly

one incoming and two or more outgoing links. A leaf node, also known as a terminal node, has exactly one incoming link and no outgoing link. Each leaf node is represented by a class label. The root and internal nodes, also known as non-terminal nodes, include attribute test conditions that can typically be defined by a single attribute. The possible outcomes of the attribute test conditions produce child nodes. For example, the attribute *body temperature* is used by the root node to define the test attribute conditions determining whether the species is warm- or cold-blooded, thereby resulting in two child nodes. DTs classify the data points from root node to some leaf nodes via internal nodes. Each node contains either a prediction or a Boolean condition (Yes/No) based upon which leaf nodes are created that are linked to target class labels.

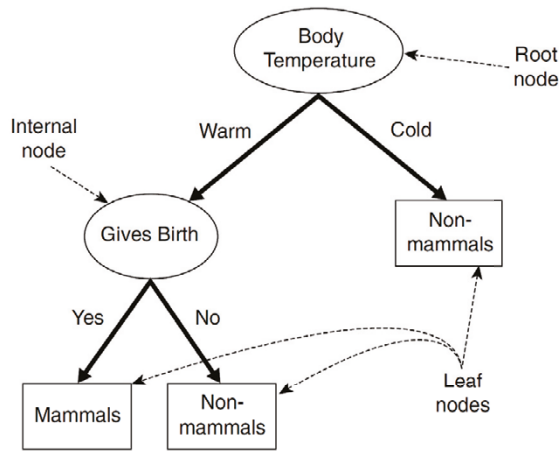


Figure 2.3: The decision tree (Pang-Ning Tan et al., 2016)

### i. Creation of DTs

The creation of DTs is based on the goal of minimising the prediction error on a training set. The tree is created in a recursive manner by partitioning the data set on the explanatory variables. A metric known as *purity gain* is used as the score function to evaluate the split of the data set  $D$  into its subsets  $N(D_1, D_2, \dots, D_i)$ . The *purity gain*, as shown in Eq. (2.2), is defined as the difference in purity between  $D$  and the subsets  $D_i$ .

$$gain(D, N) = purity(D) - \sum_{i=1}^N P(D_i) * purity(D_i) \quad (2.2)$$



$P(D_i)$  is the proportion of  $D$  in each subset  $D_i$ .

Based on Eq. 2.2, the purity gain of each split is calculated, and the one with the highest purity gain is selected. This process is continued recursively for each subset in the selected split.

The metrics that are used to measure the information uncertainty in a data set and class impurity in the nodes are *entropy* and *gini index*, respectively. The overall aim of using *entropy* and the *gini index* is to minimise the depth of the final tree by always choosing a number of splits that impacts the classification of an instance. Considering all the explanatory variables, the information gain is calculated for all the conditions of the split based on the probability of random selection of an instance of the target class, which is denoted as  $P_i$  in the *entropy* formula shown in Eq. (2.3). The *entropy* value ranges from 0 to 1, and its maximum value is reached when the probabilities of all classes are equal.

$$E(N) = \sum_{i=1}^n P_i * \log\left(\frac{1}{P_i}\right) \quad (2.3)$$

*Gini index*, as calculated using the formula shown in Eq. (2.4), calculates the class impurity in each node  $t$  based on the class probabilities ( $p^2$ ) of  $j$  classes.

$$Gini\ index = 1 - \sum_{i=1}^n p^2(j|t) \quad (2.4)$$

## ii. Pruning

Overfitting is a common problem in DTs, which means that when a fully grown tree is optimised, it memorises the training set, including noise, in such a way that it later yields a high generalisation error. The standard way of handling the DT overfitting problem is to prune the trees. Pruning, simply put, limits the growth of the large trees, typically by iteratively merging leaf nodes.

Pruning can be handled by two approaches. Pre-pruning prevents the growth of the tree before the weaker branches grow, while post-pruning occurs after the tree is fully grown. Generally, subtrees created by the post-pruning are evaluated on the training set or a test set that contains unseen instances. The ways in which these subtrees are created are different for every post-pruning algorithm; however, all the algorithms perform subtree replacement and raising. Subtree replacement, which begins in the leaves of the tree, replaces the selected subtree with single leaves. Subtree raising removes internal nodes and thus raises the subtree to a higher

level in its branches. Overall, pruning creates a large number of candidate subtrees, out of which the best performing tree on the test set becomes the final choice.

### 2.1.3.2 Ensembles

Ensembles of DTs provide more model accuracy than do single models. An ensemble is a technique that combines the prediction results of several different classification or regression models to minimise the generalisation error. Intuitively, classifier ensembles that use averaging to combine several models will reduce the error by eliminating any uncorrelated errors in the ensemble models. However, for the ensemble approach to work, the base classification models must produce their errors on different data instances. Nothing is gained by combining models if the errors are made on the same instances. Therefore, the tendency of a set of base models to commit their errors on different instances is a vital property of an ensemble, often referred to as *diversity*. There are many different metrics for diversity, but the overall idea is that a higher diversity will increase the ensemble accuracy as long as the base classifier accuracies are not affected. In practice, base classifier accuracy and diversity are highly correlated, creating a trade-off rather than a simple optimisation problem. Provided there is some diversity, the ensemble will have a higher accuracy than the average accuracy of the individual classifiers.

The techniques that create diversity for a set of base models are categorised as implicit and explicit techniques. Explicit techniques directly optimise the diversity measure, while implicit techniques produce diversity without targeting it. The most popular implicit techniques divide the training set by instances or attributes to produce slightly different training sets for each model. Bagging (Breiman, 1996) and Boosting (Schapire, 1990) are two well-known implicit techniques that are widely used for both regression and CTs.

Bagging is a bootstrap ensemble technique that creates  $N$  different models, each of which is trained on a different sample of data. Each sample is chosen at random with replacement, which means that the same sample could be used to train multiple  $M$  models. The final model makes a prediction based on the average of  $M$  models. Boosting is a method of combining multiple models in a sequential rather than parallel manner. The main idea is that the new model is created based on the previous model's performance, with a higher weight given to the instances that were misclassified by previous models. This improves performance by allowing more generalisation across the data.

Two popular ensembles that use a decision tree as the base model are *gradient boosting* (GB) and *random forest* (RF). *GB* uses the Boosting technique, whereas RF uses Bagging.

**i. Random forest**

RF builds a set of multiple decision trees using bootstrapping on a training set and then combines their predictions. When creating individual trees, an arbitrary subset of attributes (hence “Random”) is selected, from which the best attribute for the split is chosen. The final model is based on the aggregation of the results from the individual trees. The structure of the RF is presented in Figure 2.4. The aggregation of the predictions from strong and de-correlated DTs enables *RF* to reduce the variance of the DTs without affecting their low bias, and this makes RF models more robust to the overfitting problem (Breiman, 2001). Moreover, as it considers only a small subset of attributes at each internal node, RF is computationally fast and robust when applied to a high dimensional data set.

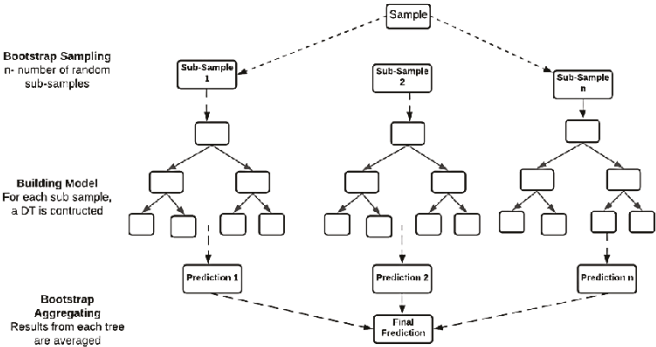


Figure 2.4: The random forests

Model performance can be optimised by using hyper-parameter tuning, which is used to control the learning process. The common hyper-parameters that need to be tuned are number of trees, depth of trees and the maximum number of feature attributes at each split.

The main advantage of using RF is that it avoids overfitting, i.e. it improves generalisation performance. However, the major drawback of the RF model is the lack of interpretability; therefore, it is considered a black-box model.

## ii. Gradient Boosting

*GB* constructs the model in a sequential manner and optimises the error in each iteration (Friedman, 2001; see Figure 2.5). When multiple trees are added, the fitted model minimises the error. Fitting the model too closely to the training set, on the other hand, can result in poor generalisation due to overfitting. Therefore, it is necessary to have the optimum number of gradient boost iterations to prevent this. Important hyper-parameters that need to be tuned to avoid overfitting are *number of trees*, *learning rate*, which determine the boosting learning rate and its value range (from 0 to 1) and the *depth of trees*, which determines the individual tree's maximum depth.

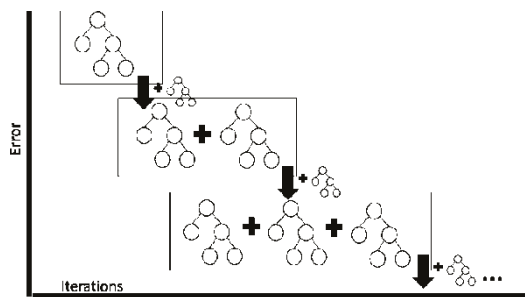


Figure 2.5: The gradient boost model

The advantage of *GB* is that it often provides high accuracy and is flexible due to multiple options for hyper-parameter tuning that optimise the model for improved performance.

Overfitting may also arise as *GB* continues to minimise errors after each iteration because the outliers are overemphasised. This can be handled by using a sampling technique, like cross validation. The disadvantage of *GB* is that it is time consuming and less interpretable.

In fact, interpretability is a common problem with all ensemble techniques. To handle this, *rule extraction* techniques are used that transform the opaque model into a comprehensible model without affecting its accuracy (Craven & Shavlik, 1999). The idea behind rule extraction is that an opaque model can reduce noise from the original data set, making prediction easier. The extracted model can either be used to explain the opaque model's behaviour or to make actual predictions. Any rule representation could be utilised for extraction as long as the decision maker understands it.

### 2.1.3.3 Naïve Bayes

Naïve Bayes (NB) is a classifier that estimates the conditional probability of the class, given the class label  $C$ , under the assumption that conditional independence exists among the attributes (Bishop, 2006). The *NB* classifier is based on the principle of the Bayes theorem, the aim of which is to identify the posterior probability, i.e. the probability of class  $C$  given the feature attributes  $P(L | \text{features})$  based on its prior probability, class-conditional probability and the evidence of features (Bishop, 2006). Bayes' formula for calculating posterior probability of class  $C$  based on the given feature  $f$  is shown in Eq. (2.5).

$$P(C|f) = \frac{P(C)P(f|C)}{P(f)} \quad (2.5)$$

$P(C|f)$  is the posterior probability of class  $C$ ,  $P(C)$  is the prior probability of class  $C$ ,  $P(f|C)$  is the likelihood (the probability) of the feature given class and  $P(f)$  is the prior probability of the feature.

The estimation of  $P(C)$  is based on the fraction of training instances belonging to each of the classes present in the training set. A good estimate of the conditional probability of each feature  $f_i$  given the class  $C$  does not require a very large training set. While performing the classification on a test set, *NB* classifies test instances by calculating the posterior probability for each class  $C$ , see Eq. (2.6).

$$P(C|f) = \frac{P(C) \prod_{i=1}^n P(f_i|C)}{P(f)} \quad (2.6)$$

As the  $P(f)$  is constant for every  $C$ , the numerator in Eq. (2.8), i.e.  $P(C) \prod_{i=1}^n P(f_i|C)$ , needs to be maximised for a given selection of class  $C$ .

*NB* classifiers are robust in the sense that they can handle noise in the data set. The issue of missing values in the data set can also be easily solved by *NB* classifiers because they ignore them during the model training. *NB* classifiers are also robust to irrelevant attributes since they have a uniform distribution of their likelihood from the distribution of instances belonging to class. The drawback of *NB* classifiers is that a violation of the underlying

assumption of conditional independence among the attributes can occur if they are correlated, which can degrade the performance of the classification (Pedregosa et al., 2011). *NB* is easy to use, as it can be simply implemented without hyper-parametric tuning, and it can effectively handle the large data sets.

#### 2.1.3.4 Fuzzy sets

Fuzzy set theory (FST) lends itself to the design of intelligent systems and their applications in AI (Hüllermeier, 2005). The methods developed in FST have the potential to perform all the steps in the knowledge discovery or model induction. FST is particularly appropriate for application in the data selection and preparation stages, for example, for vague data modelling using fuzzy sets (Viertl, 2011) and for creating fuzzy summaries of the data.

Fuzzy sets are used to map the uncertainty of real-world problems, which can be subjective or vaguely defined. The theory was introduced by Zadeh (1965) as an extended version of classical set theory, and it has been further extended and intensively applied since 1970 (Gottwald, 2010; Zadeh, 1975).

The fuzzy set can be defined as  $(U, m)$ , where  $U$  is defined as a universal set. Each member in this set is assigned a membership function  $m$  that lies between 0 and 1 and can be represented as  $U \rightarrow [0,1]$  (Dubois et al., 1980). For a fuzzy set  $X = (U, m)$ , the member function will be  $m = \mu(E)$ . Membership functions for each element  $e$  for a finite set  $U = \{e_1, \dots, e_n\}$  can be written as  $\{m(e_1)/e_1, \dots, m(e)/e_n\}$ . If  $e \in U$  and  $m(e) = 0$ , then the element  $e$  is not considered in  $X = (U, m)$ , and if  $m(e) = 1$ , it is considered fully only if  $0 < m(e) < 1$  (Beg & Ashraf, 2009).

Fuzzy sets allow scientific operations of *intersection*, *complement* and *union*. For the fuzzy set  $X$  and  $Y$ ,  $X, Y \subseteq U$ ,  $u \in U$ , the *complement* operation can be represented as  $\mu_{\neg X}(u) = 1 - \mu_X(u)$ , the *intersection* operation can be represented as  $\mu_{X \cap Y}(u) = \min\{\mu_X(u), \mu_Y(u)\}$  and the *union* of two fuzzy sets can be shown as  $\mu_{X \cup Y}(u) = \max\{\mu_X(u), \mu_Y(u)\}$  (Gottwald, 2010). Fuzzy sets have commutative, associative, distributive properties as well as an identity property and a transitive property (Gottwald, 2010). Triangular membership is a popular method that has been widely applied to solve engineering problems due to its flexibility (Pedrycz, 1994). For a fuzzy set  $Y$ , the degree of membership function, as illustrated in Figure 2.6 and defined in Eq. (2.7), represents the lower and upper limits,  $a$  and  $b$  as  $a < b < c$ .

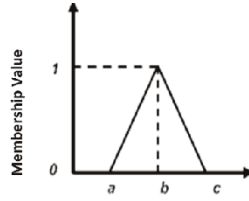


Figure 2.6: Fuzzy triangular membership

$$\mu_Y(x) = \begin{cases} 0 & x \leq a \\ \frac{x-a}{b-a} & a < x \leq b \\ \frac{c-x}{c-b} & b < x < c \\ 0 & x \geq c \end{cases} \quad (2.7)$$

There are many applications of *FST* that are widely used in various fields. From the data mining perspective, *FST* has been extended to handle clustering tasks where the degree to which a particular data instance belongs to different clusters within the data set at the same time is measured using the membership functions. Another interesting application of *FST* in data mining is the induction of models on the basis of rule-based inferencing. Both classification and regression functions can be represented by fuzzy rule-based inferencing, during which the crisp input attribute values are transformed into fuzzy ones. This process is known as *fuzzification*. It then maps the fuzzy output back to a crisp value, which is known as *defuzzification*. However, some models, for example the Takagi-Sugeno model, do not need to defuzzify the output since they directly generate output in terms of crisp values. The dependencies among the attributes in the dataset are often expressed in terms of *IF-THEN* rules, which can be extracted using fuzzy association analysis that produces rule-based models. For example, fuzzy association analysis could be used in the predictive classification to extract the *IF-THEN* rules to predict the label of the class. The class assignment to the instances in the data set could be a form of single rule defining the relationship between attributes  $X_i$  and the target variable  $Y$ , as given below:

$$\text{IF } (X_1 \in a) \text{ AND } (X_2 \in b) \text{ THEN (class label} = Y)$$

where  $a$  and  $b$  are fuzzy sets.

#### 2.1.4 Evaluation of predictive models

As previously discussed, there are two types of evaluation, external and internal, that are used to measure the performance of predictive models. In an external evaluation, the performance of predictive models is estimated based on their application in actual practice in the business context. This allows for comparison of models' performance in the testing and production phases. Take, for example, a scenario in which a set of customers to be targeted with a special offer is identified using an algorithm. The response behaviour of the customers could be studied against the reference group of customers that was not given any such offer. This approach can be compared with the widely used statistical technique known as A/B testing, in which two versions of the same variable are compared to study performance. One version is kept in a randomised control environment, another version is not. External evaluation of model performance helps to study how the models perform when applied in a production environment, i.e. when applied to a production set.

An internal evaluation is an evaluation of model performance during the training phase, i.e. before application to a production set. The main goal of this approach is to estimate the future performance of the models based on optimising their score functions using available data, i.e. so that they generalise well to unseen novel data.

When predictive models are trained on a training set, a minimisation of error is expected. However, this error may not be a good estimate of the performance of the predictive model when applied to a production set. The prediction of a target variable is based on the probability function created using the training set, which often does not represent all the possible instances of the novel data. In this sense, the model cannot estimate the true error. In order for the training error to be representative of true error of the predictive model when applied in practice, sampling approaches are used, such as holdout sample (test set creation), random sampling and cross validation (Pang-Ning Tan et al., 2016).

- i. **Holdout sample:** A sampling in which a subset of the training instances is removed randomly and inserted into test data prior to training. Typically, two-thirds of the data are used for training (training set) and one third is used for testing (test set). The drawback of this approach is that if the number of instances in the training data set is small, it will negatively impact the model performance. Besides, if the test set is too small, i.e. only a small number of instances are present in the test set, it will result in high variance, which in turn reduces the reliability of the estimated error. Moreover, since the distribution of



classes in both training and test sets can vary, the estimation of error and the model performance can be negatively affected.

- ii. **Random subsampling:** An extension of holdout sampling. Multiple holdouts randomly split the data into subsets several times. The overall error is the average of the errors of all the models. A major drawback of random subsampling is that it is difficult to control the number of times the instance is used in the *test* and *training* sets.
- iii. **Cross validation (CV):** This is one of the most common subsampling techniques, and it is considered to be a more systematic approach than is random subsampling. In CV, before the model training, the data is randomly split into equal sized  $k$  subsets, also known as *folds*, to which each instance is randomly assigned. The  $k$  number of models is then trained on the data. One of the  $k$  folds serves as a *test set*, while the remaining  $k-1$  folds are used as a training set. In this way, since each fold serves as a *test set* once, CV allows each instance to be used for training  $k-1$  times and to be used once for testing. Since the folds are randomly selected, they lack representation of all the instances. To ensure the complete representation of instances, the approach known as *leave one out*, the special case of CV, is utilised to ensure that a single instance is present in each fold. Leave one out utilises maximum instances for the training, while the *test set* remains unseen. In the predictive classification task, it is important to ensure that the created folds have representative class distribution. The technique that assigns an equal number of instances of each class to each fold is known as stratification. When an unequal number of instances are present in the folds, duplicates of existing instances are used to balance their proportion. Typically, 10-fold CV with stratification is used for the evaluation of predictive models.

For evaluating the predictive performance of the models, several accuracy metrics are used. These metrics are the measure of error, i.e. the difference between the predictive and the actual value of the target variable.

### 2.1.4.1 Evaluation metrics for predictive classification

Arguably the most important metric for evaluating the performance of classification models is the error, i.e. the proportion of incorrect predictions:

$$Error\ rate = \sum_{i=1}^n I(Y_i - \hat{y}_i) = \frac{Number\ of\ incorrect\ predictions}{Total\ number\ of\ predictions} \quad (2.8)$$

Most often, however, the classification accuracy, i.e. the proportion of correctly classified instances, is used instead:

$$Classification\ Accuracy = 1 - Error\ rate = \frac{Number\ of\ correct\ predictions}{Total\ number\ of\ predictions} \quad (2.9)$$

The standard way of representing the results of a classification model is in a visual format called a *confusion matrix*, in which the number of instances with a certain predicted label is shown along the rows and the number of instances with a certain correct target along the columns. The values in the diagonal cells of the matrix consequently represent the correctly predicted instances (TP and TN), whereas the values in off-diagonal cells represent the incorrectly predicted instances (FP and FN). An example of a confusion matrix that shows the results of a binary classification model predicting the labels of two classes, *Positive (1)* and *Negative (0)*, is shown in Figure 2.7.

		Actual Values	
		Positive (1)	Negative (0)
Predicted Values	Positive (1)	TP	FP
	Negative (0)	FN	TN

Figure 2.7: An example of a confusion matrix

The notations represented in the cells of a confusion matrix, as shown in the figure above, are defined as follows:

- $TP = True\ Positives$  : the number of positive instances predicted correctly by the model
- $TN = True\ negatives$  : the number of negative instances predicted correctly by the model
- $FP = False\ positives$  : the number of positive instances predicted wrongly by the model, i.e. predicted as negative
- $FN = False\ negatives$  : the number of negative instances predicted wrongly by the model, i.e. predicted as positive

These notations are used to calculate a number of performance metrics in terms of percentages. These alternative metrics are particularly useful when the class instances are imbalanced, thus making accuracy very blunt. Some metrics that are calculated based on the numbers represented in the confusion matrix are defined as follows:

- True positive rate ( $TPR$ ):  $TPR$  is the proportion of positive instances that are correctly predicted by the classifier.  $TPR$  is also known as *sensitivity*.

$$Sensitivity = TPR = \frac{TP}{TP+FN} \quad (2.10)$$

- True negative rate ( $TNR$ ):  $TNR$  is the proportion of negative instances that are correctly predicted by the classifier.  $TNR$  is known as *specificity*.

$$Specificity = TNR = \frac{TN}{TN+FP} \quad (2.11)$$

- False positive rate ( $FPR$ ):

$$FPR = \frac{FP}{TN + FP} \quad (2.12)$$

- False negative rate ( $FNR$ ):

$$FNR = \frac{FN}{TP + FN} \quad (2.13)$$

Two other very important and commonly used metrics are *Precision* and *Recall*. These are used when predicting one class correctly is more important than predicting the other. In other words, these metrics are typically used when the minority class (here the positive class) is given more importance. *Precision* measures the proportion of instances classified as positive that actually are positive, whereas *Recall* is the proportion of positive instances that are actually predicted as positive.

$$Precision = \frac{TP}{TP + FP} \quad (2.14)$$

$$Recall = \frac{TP}{TP + FN} \quad (2.15)$$

Another metric that combines *Precision* and *Recall* is the *F1* score, which is calculated as the harmonic mean of *Precision* and *Recall*. The *F1* score, consequently, provides a balance between them.

$$F1 = 2 \times \frac{Precision \times Recall}{Precision + Recall} \quad (2.16)$$

To graphically represent the trade-off between *TPR* and *FPR*, a *receiver operating characteristic* (ROC) curve is commonly used. ROC curves evaluate the accuracy of the classifier by plotting the *TPR* along the *Y*-axis and the *FPR* along the *X*-axis. The curve (see Figure 2.8) represents the model generated by the classifier. The classifier is considered to be a strong classifier when the curve is located as close as possible to the upper left corner. The area under the ROC curve, called *AUC*, measures the aggregated area under the ROC curve, and it can be used to comparatively evaluate the performances of different models. The classification threshold values range from (0, 0) to (1, 1), as shown in Figure 2.8. The value of *AUC* ranges between 0 and 1, where 1 is perfect and 0.5 corresponds to random guessing. The higher the *AUC* value, the better the prediction performance, or rather the ranking ability, of the model. ROC curves are commonly used when the classifier produces prediction

probabilities of the instances being one of the classes. In this case, the output of the classifier is continuous, and its value ranges between 0 and 1. The threshold value needs to be set in order to interpret the output of the classifier. For example, if the threshold value is set to be 0.5, then the prediction probability values above 0.5 can be considered *positive* and those below 0.5 *negative*.

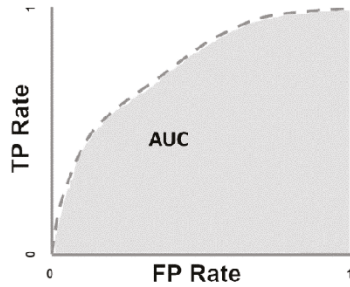


Figure 2.8: A sample ROC curve

#### **2.1.4.2 Evaluation metrics for predictive regression**

The standard metrics used for evaluating the performance of the predictive regression models are *mean absolute error* (MAE), *mean squared error* (MSE), *root mean absolute error* (RMSE), *correlation* ( $r$ ), *coefficient of determination* ( $R^2$ ) and *mean absolute percentage error* (MAPE) (Bishop, 2006; Pang-Ning Tan et al., 2016).

However, not all metrics are relevant for specific conditions, i.e. these metrics serve various purposes and some are more useful than others in a given situation. In order to compare different regression models applied over different data sets, relative evaluation metrics should be used. A relative metric is one that is relative to some specific reference. The use of relative metrics becomes even more relevant when the target variable in the data sets upon which a specific regression model is applied differ greatly in magnitude, as it affects the magnitude of the error.

The *mean absolute error* (MAE) is the most commonly used metric for evaluating the predictive regression model. *MAE* measures the average magnitude of the errors in a set of

predictions without their direction. It calculates the average of the absolute differences between predicted and actual values of the target variable by giving equal weight to their individual differences.

$$\text{MAE} = \frac{1}{n} \sum_{i=1}^n |y_i - \hat{y}_i| \quad (2.17)$$

where  $y_i$  is the observed value, and  $\hat{y}_i$  is the value predicted by the model.

*Mean squared error* (MSE) is another commonly used metric that calculates the average squared difference between the predicted value of the target variable and its actual value (Bishop, 2006; Kuhn & Johnson, 2013).

$$\text{MSE} = \frac{1}{n} \sum_{i=0}^n (y_i - \hat{y}_i)^2 \quad (2.18)$$

The square root of *MSE* is also used as a metric. Known as *root mean square error* (RMSE), the advantage of this metric is that it has the same unit as the target variable, which is the main reason for using it instead of *MSE*.

$$\text{RMSE} = \sqrt{\sum_{i=1}^{i=n} \frac{1}{n} (y_i - \hat{y}_i)^2} \quad (2.19)$$

The correlation between the predicted and actual values of the target variable is an important measure and thus needs to be reported. The standard metric used for calculating the linear correlation between the prediction and actual value is *correlation coefficient* ( $r$ ). The value of  $r$  ranges from -1 to 1. A positive value of  $r$  indicates a positive correlation, whereas a negative value indicates negative correlation. A value of  $r$  that is closer to 0 indicates a smaller correlation; a perfect correlation is when  $r$  is either -1 or +1.

$$r = \frac{\sum_{i=1}^n (\hat{y}_i - \bar{\hat{y}}_i)(y_i - \bar{y}_i)}{\sqrt{\sum_{i=1}^n ((\hat{y}_i - \bar{\hat{y}}_i)^2 \sum_{i=1}^n (y_i - \bar{y}_i)^2)}} \quad (2.20)$$

The proportion of variability of the target variable in the data captured by the regression model is also important to measure. For this purpose, the *coefficient of determination* ( $R^2$ ) metric is commonly used. The value of  $R^2$  ranges from 0 to 1. The value 1 implies a perfect correlation, while 0 implies no correlation at all. It is important to note that even if the value of  $R^2$  is close to 1, i.e. nearly a perfect correlation, the predicted values may still be far from the actual values because  $R^2$  does not account for the bias between predicted and actual values of the target variable.

$$R^2 = \left( \frac{\sum_{i=1}^n ((\hat{y}_i - \bar{y}_i)(y_i - \bar{y}_i))}{\sqrt{\sum_{i=1}^n ((\hat{y}_i - \bar{y}_i)^2 \sum_{i=1}^n (y_i - \bar{y}_i)^2)}} \right)^2 \quad (2.21)$$

For the special case of predictive regression, i.e. time series forecasting, the commonly used metric for evaluating the performance of forecasting models is *mean absolute percentage error* (MAPE). If, for example, sales of a specific product are to be forecasted based on its previous sales values, then the *MAPE* of the forecast model can be calculated as below:

$$MAPE = \frac{100}{n} \sum_{t=1}^n \left| \frac{A_t - F_t}{A_t} \right| \quad (2.22)$$

where  $A_t$  is actual sales, and  $F_t$  is forecasted sales.

This section introduced the major concepts used in data mining from the data-driven AI perspective. A study of these concepts and techniques provides a brief understanding of how they have been used in the case studies presented in this thesis. In Section 2.1.1, we discussed the types of data, data quality and data pre-processing techniques, which are the preliminary steps for data analysis. Section 2.1.2 explained the data mining tasks, elaborating first on descriptive tasks and then focusing on insights about the datasets, data exploration, summary statistics, clustering and association analysis. Second, it discussed predictive tasks, in which

the value of the target attribute is predicted based on the values of the independent attributes. Section 2.1.3 discussed different data mining techniques applied in this thesis, and Section 2.1.4 described the data sampling and model evaluation techniques that were used. In the next section, we will cover the background of the F&A industry and of textile management and discuss the related literature on the context and characteristics of the F&A industry from the perspective of supply chain management.

## **2.2 Fashion, apparel and textile management**

Textile management is an interdisciplinary research field and could broadly be defined as the study of all issues relating to management functions in the context of the textile industry. This means that research within textile management takes inspiration from multiple fields of research, such as business administration, industrial economy and fashion. The aim of this chapter is to outline some of the F&A concepts used in the articles that relate to the field of textile management. The following section describes the concepts of textiles management and provides an overview of the traditional F&A supply chain, including characteristics of the F&A industry and challenges, particularly those related to customer behaviour and sales forecasting.

### **2.2.1 Textiles, clothing, apparel, garments and fashion**

Since, as mentioned above, textile management is a multidisciplinary subject that is highly context-dependent, it is important to start with definitions of the core concept, i.e. textiles. Textiles are man-made materials that consist of fibres (natural or man-made) (Humphries, 2009). Textiles are today used in a wide variety of products and industries, such as car tires and aeroplanes as well as for biotechnological uses. The textile industry thus expands far beyond clothes and apparel, which are the focus of this thesis, and produces both consumer goods and industry products.

As explained in the introduction, this study mainly focuses on the consumer industry, specifically the fashion and apparel (F&A) industry. There are a number of widely used concepts with similar meanings, such as clothes and apparel. “Clothes” is defined as the things that you wear to cover, protect or decorate your body (Cambridge Online English



Dictionary, 2021). Thus, clothes may be made of textiles, but might also be made of other materials. “Garment” is defined by the same source as a piece of clothing and “apparel” as a particular type of clothes that are sold in a shop. All three concepts are often used interchangeably in the literature (Kaiser, 2014), as a self-contained noun as well as in combination with other words (for example, garment industry). In relation to retail, however, apparel seems to be the most appropriate.

Fashion, on the other hand, is a more complex concept. It may be defined as “an expression that is widely accepted by a group of people over time” (Fernie & Sparks, 1998). The term “fashion” could thus relate to many different things that people use to adorn their bodies and to lead a stylish life, such as cars, electronic devices, music, clothes, glasses, house decorative materials and jewellery. In relation to clothes and apparel, fashion thus refers to the symbolic and intangible values of clothes rather than to the functional values described above (Kaiser, 2014; Kawamura, 2005).

The fashion and apparel (F&A) industry manufactures a wide variety of fashion garments, using a range of textiles and with many design and style variations (Le Bon, 2014). The other related terms used for the F&A industry are garment industry and clothing industry. However, the products manufactured in the clothing or garment industries may not necessarily have the fashion element, as fashion is a symbolic and intangible add-on to the clothes or garments that fulfil people’s clothing requirements (Kawamura, 2005). Within this thesis, although variations of the concept were used in early articles, later articles have adopted the term “F&A industry” to imply that the examined data sets are related to retail. For consistency, the same term is also used when discussing the supply chain below.

## **2.2.2 F&A industry from a supply chain perspective**

*Supply chain management* (SCM), as defined by Cooper et al. (1997), is the combination of key business processes, from end-user through original suppliers, that provide products, services and information and add value for customers and other stakeholders. The supply chain in the F&A industry involves processes, known as SCM processes, that transform the raw materials into finished products and deliver them to the end-customers. In the literature, the terms *customer* and *consumer* are interchangeably used; however, their definitions differ. The *customer* is simply someone who purchases the product or service, while the *consumer* is the actual user of that product or service (Baines et al., 2013). In this thesis, datasets that are

collected within the F&A industry primarily focus on the customer, i.e. the buyer of the products, and they do not really reveal who the consumer is. Thus, “customer” is the most relevant term when discussing these data sets. The SCM processes can be defined by the associated business processes, a structured and measured set of activities that produce a specific product or output for a particular market or customer (Cooper et al., 1997). The supply chain is comprised of a *value chain*, which is a *means* by which business actions transform inputs into valuable output. The value chain is defined as the range of activities that add value in each SCM process and deliver a quality product to the end-customer (Walters, 2002). The value of any product or service is the result of its ability to fulfil its customers’ demand. Traditionally, the value chain in an industry begins with its core competencies and assets and moves to various raw materials and inputs, to a product offering, and finally to the end-customers via various marketing channels. In today’s context, the industries are increasingly becoming customer-centric and so is their value chain. The idea of the customer as the end of the value chain has been challenged by the reversed value chain, which is driven by the customer priorities and needs (Walters, 2002). Thus, this value chain starts with identifying customer needs, then identifies the appropriate marketing channels, services and products, and then the inputs, raw materials and core competencies needed to fulfil them.

The global fashion and apparel supply chain is a complicated system of various actors and is usually laden with complex decision-making (Stevenson et al., 2005). It deals with a wide range of raw materials, including fibre, yarn, fabric, dyeing and other chemicals and the related processes. The latter can be broadly divided into four stages, as shown in Figure 2.9: design, fabric production, apparel production and distribution of products. The supply chain has previously operated on a push model in which brand owners or retailers (buyers) provide information to manufacturers, such as the design or technical specifications of the fabric and garment to be produced, the volume of products and the sizes in which the garment is to be produced. Fabric and garment manufacturers follow the instructions to create samples, which are then converted into finished fabric and garments after the buyer approves them. The finished garments are then delivered to a wholesaler or retailer. If it is wholesaler, there is a new actor who acts as intermediary between customer and wholesaler. The garments are sold through one or more channels in the context of the retailer, such as brick-and-mortar stores, web-shops (e-commerce), department stores and multi-brand retail. Finally, finished apparel is delivered to the customers through the distribution process via wholesalers and retailers. Designers, who are usually employed by retailers, are in charge of creating collections based

on market and trend analysis. In most cases, retailers do not own any manufacturing facilities, but they play an important role in attempting to bring products to market.

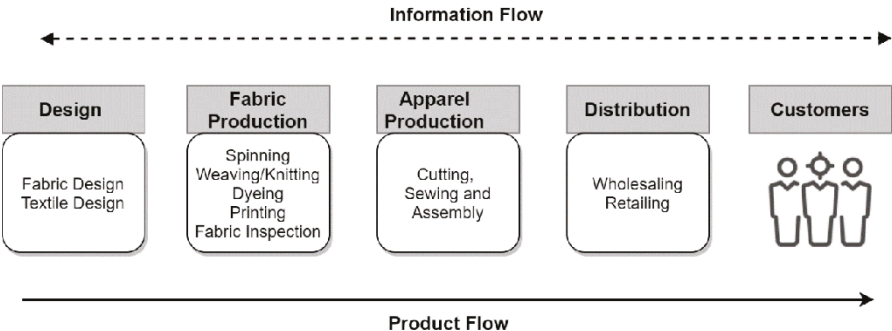


Figure 2.9: The traditional F&A supply chain

As we can observe from Figure 2.9, the process of transformation of fibres into the final F&A product includes many entities, which results in a long and complex supply chain and entails several manufacturing activities. In this supply chain, the F&A retailers are considered to be the main intermediate actor, as they drive the whole supply chain flow, starting from when they order the F&A items from the raw material providers to when they deliver them to the end-customers. Recently, the F&A supply chain has been undergoing numerous transformations due to the dynamic conditions of the F&A industry, including new market trends, large product variations (Vaagen & Wallace, 2008) and varying customer preferences (De Felice, 2012).

These transformations render customer demands unpredictable and uncertain (Alturki et al., 2012). Certain undesired characteristics of the market, such as product imitation and counterfeits, lead to erosion of the competitive advantage of the original brands. This, in turn, compels F&A brands to regularly introduce innovative products, owing to which the lifecycle of the F&A items is increasingly becoming shorter (Barnes & Lea-Greenwood, 2010). As the F&A industry is experiencing these varying market features, traditional supply chains cannot meet its needs. Instead, the industry should adopt new, responsive, demand-driven strategies that focus on enhancing product availability and innovation and the customer experience (Lam & Postle, 2006). A fully integrated and multitier supply chain network is crucial to facilitating the real-time sharing of inventory and demand-related information among its key segments, thereby reducing the operating costs and improving customer satisfaction. Moreover, this approach to information sharing and focus on demand-driven, responsive

strategies could reduce the risks associated with uncertain demand and inventory levels, as these methods enable real-time visibility across the supply chain and allow companies to react quickly to uncertain market situations (Budd et al., 2012).

As discussed in the above context, improved supply chain performance is considered critical to the F&A industry. Two widely used models to improve the flexibility and performance of supply chains operating under volatile demand situations are *Lean* and *Agile* models (Battista & Schiraldi, 2013). The term “*Lean*” refers to the ability of a supply chain to strictly fulfil customer requirements, such as order completion, accuracy and delivery time, by using more accurate demand forecasts and real-time information exchange (Bruce et al., 2004). The term “*Agile*” refers to the ability of a supply chain to perform well in an environment where the product demand is highly uncertain (Christopher et al., 2004). Depending on the specific market environment in which companies operate, hybrid approaches may be adopted in which one or more of several models might be used.

Overall, faced with numerous challenges to effectively manage supply chains, the F&A companies frequently deal with complex decision-making problems that include sales forecast, customer behaviour analysis and formulation of quick response strategies. The goal of profit maximisation and customer satisfaction depends largely on how effectively F&A companies manage these decision problems, as they constitute the success factors of F&A supply chain management (Thomassey, 2010).

### **2.2.3 Retailing in the F&A supply chain**

As discussed above, *retailers* are the main actors between production and the customer in the downstream supply chain network; sales, distribution, assortment and marketing constitute some of the key decision areas. All the activities directly related to the sale of goods and services to the ultimate end-customer for personal and non-business use are referred to as *retailing* (Baines et al., 2013). The F&A retailers make orders to the suppliers and provide the customers with the goods manufactured by the upstream companies. Hence, the retailers are considered to be the main intermediaries between suppliers and customers, and they play a vital role in all the major F&A supply chain operations. In order to effectively cater to the consumers’ specific demands, inventory planning and the design of marketing strategies are vital for F&A retailers (Bhardwaj & Fairhurst, 2010; Choi, 2016). Digitalisation offers retailers new ways to add value to the value chain and customers. The digital traces that the

customers leave when they select and purchase products from the retailers provide transaction data that, through use of various analytical tools, may be used by retailers to derive insights into customer behaviour.

Transaction data thus offers F&A retailers new opportunities to address the challenges, such as sales forecasting, customer behaviour analysis and formulation of quick response strategies, that arise from the complex nature of the F&A supply chain, especially regarding difficult-to-manage sourcing and manufacturing processes (Thomassey, 2010)

### ***2.2.3.1 Analysing customer purchasing behaviour***

Analysis of the customer's purchasing behaviour helps marketing professionals determine their customers' demands and needs. Understanding this process is beneficial to businesses because it allows them to better tailor their marketing initiatives to previous marketing efforts that have successfully influenced customers to buy. The F&A retail industry is undergoing a rapid transitory phase due to various factors. These include new advances in manufacturing technology, such as Industry 4.0, global trade reforms and volatile consumer buying behaviour (Grieco et al., 2017; Taplin, 2014; Wang & Ha-Brookshire, 2018). One strategy is for F&A retailers to design effective and strategic marketing and promotional campaigns to increase their customer retention rate, thereby improving their competitiveness in the market. Another strategy is to strive for high-quality products and services in order to improve the customer buying experience and thereby encourage the development of brand loyalty (Dahana et al., 2019).

Rapidly evolving customer preferences have led F&A retail brands to implement fast-fashion business models in order to serve their customers' ever-changing appetite for new fashion styles (Bhardwaj & Fairhurst, 2010). Popular brands, such as Zara and H&M, introduce new styles within as short as three weeks' time (Jin & Bennur, 2015). Unlike in the traditional F&A business context, today's advanced technology, such as 3D scanners, IoT, AI, big data analytics and cloud-based database technology, provides F&A companies with the unique ability to map out their customers' specific needs very quickly and further necessitates improvement of the supply chain operations in order to satisfy these needs (Miklosik & Evans, 2020). However, the complexity of the decision-making involved in these operations is increased by the nature of the data that is generated through customers' shopping transactions.

Today's multi-channel customer is always connected, and today's researcher talks about the consumer decision-making journey rather than process in order to illustrate the complexity that it displays (Lynch & Barnes, 2020). Within a single purchase, the customer interacts not only with the retailer but with numerous other channels. Tracing customer data, therefore, becomes a complex task for the retailers. F&A retail companies cannot sustain their profit by relying on traditional decision models that are not capable of handling the complexities that have resulted from the new technologies and their continuing and significant influence on customers' shopping behaviour (Gupta et al., 2018). Moreover, they are compelled to evolve their strategic operational and management processes from time to time to mirror their customers' varied shopping traits and market needs.

One important aspect of a customer buying behaviour analysis is the "customer churn", which refers to irregular or uncertain purchase behaviour that is magnified by disloyalty towards or dissatisfaction with retail brands (Burez & Van den Poel, 2009). Customer churn is caused by a variety of factors that need to be clearly identified, as it could cause profit losses to the F&A retailers. The customer churn behaviour can affect the business growth of the F&A retail industry; therefore, detecting it in advance becomes a critical decision problem. The major challenge facing F&A retailers is to minimise customer churn rates that impede business growth and harm customer loyalty.

F&A retailers need to sustain their business profit by offering high-quality products and services to their customers so as to retain loyalty and trust (Edelstein, 2001). However, this comes with the need to shift from the traditional marketing approach, such as mass marketing, as it is incompatible with a modern retail market that has been revolutionised by advanced digital technologies. Such transformations provide many avenues for F&A retailers to make use of AI and digital technologies to study their customers' behaviour and effectively predict churn-related traits.

### ***2.2.3.2 Sales forecasting in the F&A industry***

Generally, the decision-making in F&A retailing starts with the allocation of the budget for procurement or other operations. After designers select the items that should be included in the collection, sales forecasting allow managers to initiate the purchase or production process (Wong & Guo, 2010). F&A retail operates in a competitive market in which efficient inventory control and management play a crucial role in increasing the business profit. Accurate sales forecasting is, therefore, essential in order to be successful in such a

competitive market environment. If the sales forecasts are less accurate, situations such as under-stocking or overstocking may arise, and these could have a strong negative influence on the profitability of the company (Agrawal & Schorling, 1996; Xia & Wong, 2014). In addition, it has been observed that the F&A industry operates with long supply chains, which lead to orders being placed before the level of demand for the products is clearly identified (Huang et al., 2017; Lee et al., 1997). Despite its relevance, sales forecasting is a complex decision problem, as product sales are strongly dependent on the individual tastes of consumers, which tend to vary significantly (Bhardwaj & Fairhurst, 2010; Yu et al., 2011). Moreover, the lifecycle of fashion products is quite short, which leads to a lack of historical information about the products (Choi et al., 2012). Besides, fashion products have a wide range of collections that vary in size, colour, style, etc., leading to variable SKUs (stock keeping units) (Liu et al., 2013). F&A products are typically weather specific; hence, the sales exhibit seasonal patterns. The F&A items sold in one season are less likely to be sold in the next, so sales campaigns and promotions, such as end-of-season sales, discount offers, etc., are used to ensure that stocked items are sold in the same season. Another important factor that affects F&A sales are fashion trends, which often last a few weeks; therefore, new styles and designs are introduced very frequently. This leads to the short lifecycle of the F&A products. Figure 2.10 depicts the difference between the product lifecycle in the F&A market and that in other markets after the products are launched.

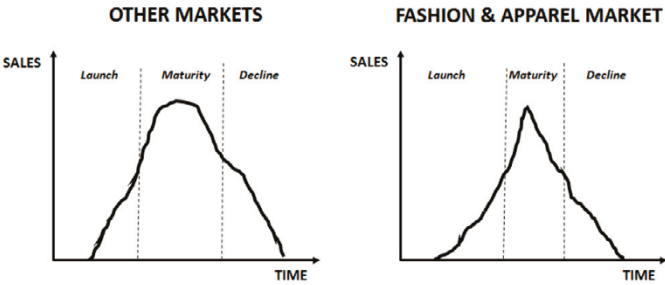


Figure 2.10: Product lifecycle in the F&A industry compared to other markets (Martino et al., 2017)

Given the impact of various exogenous factors on F&A sales, F&A retailers often operate under various constraints. For example, as the sales of F&A products are price sensitive, meaning customers are normally attracted to less expensive items, the F&A retailers often need to reduce the production cost of F&A products. This is the reason why most of the F&A companies outsource manufacturing to developing countries with relatively low wages

(Bhardwaj & Fairhurst, 2010). Due to the globalised nature of the F&A supply chain, factors such as lead time, inventory and natural calamities can be disruptive. Under such constraints, accurate sales forecasting plays a crucial role in allowing F&A retailers to be more efficient and responsive while making both strategic and operational decisions (Armstrong, 2001). Keeping the right inventory in the stock, timely replenishment of stock and scheduling order fulfilment are some of the key decisions that depend on accurate sales forecasts. The losses that F&A companies incur could be attributed to less accurate or flawed sales forecasts, as these lead to overstocking or stock-outs (Beheshti-Kashi et al., 2015).

### **2.3 Data-driven decision-making**

Data-driven decision-making refers to the decision-making process that is based on data analysis rather than solely on human intuition (Provost & Fawcett, 2013). For example, F&A retailers could simply use their experience of what customers purchased during previous summers to decide which products to advertise to them rather than actually analysing historical customer data to find out what they frequently bought during past summer months.

There are many studies of how data-driven decision-making benefits firms and affects their performance. Economist Erik Brynjolfsson and his colleagues (2011) developed a measure of data-driven decision-making that evaluates how strongly firms use data and analytical tools to make decisions (Provost & Fawcett, 2013). They found that the data-driven firms are more productive. They also found a causal relationship between data-driven decision-making and higher return on equity, return on assets, asset utilisation and market value.

In many industries, growing digitalisation is driving the datafication process. F&A retailers are increasingly collecting customer data to create customised products and services. Using customer data and effectively analysing it could enable F&A retailers to design their products in line with the customer choices and to make improved decisions. Customer data from other valuable sources, such as social media, emails and blogs, could be valuable for the F&A retailers' understanding of customer sentiments about all aspects of their preferences, allowing the retailers to effectively manage promotional campaigns, optimise their pricing strategies, develop new products and services and enhance customer loyalty. Moreover, digitalisation offers great opportunities for the F&A industry to improve the speed and flexibility of the SCM processes by adopting data-driven decision-making. Digitally connected supply chain partners can enhance communication between them and further improve the flexibility of their business operations.



### 3 Research methodology

*This chapter describes the methodological choices used in this thesis to answer the research questions. It discusses the research strategy, process, datasets, methods and evaluation procedures and metrics used to carry out this research.*

#### 3.1 Research strategy

This thesis investigates some important challenges in the F&A industry to which data analysis can provide value. To achieve this goal, empirical research targeting problems and utilising data sets from the fashion and apparel industry is conducted. The research work carried out in this thesis is multidisciplinary, and it draws knowledge from three distinct fields: textile engineering, management and computer science. In this work, knowledge from the textile engineering and management fields is used to understand the different stages of the textile apparel supply chain in current practice and the challenges facing the industry. Knowledge from the computer science field is used to investigate the potential of AI techniques to produce improved decision support.

The main scientific approach for evaluating algorithms and methods used in data-driven AI is *controlled experiments*, i.e. an empirical quantitative analysis targeting one or more properties of the algorithms or methods under investigation. Specifically, as described in Chapter 2, the approach of empirically evaluating predictive performance is standard in machine learning, with objective evaluation criteria and error metrics used for predictive performance. In particular, standard test protocols, like cross-validation, are employed to obtain unbiased estimations of performance on novel data. If possible, statistical inference is employed to rule out formulated null hypotheses, which typically state that the choice of method has no effect on the predictive performance.

In empirical machine learning, the main method to establish general properties of novel algorithms or methods is to compare the suggested approach to state-of-the-art alternatives in controlled experiments, normally using a large number of publicly available data sets for the benchmarking. Another frequently used option is performance of a *case study* (typically in the form of a proof-of-concept), during which a specific, often real-world, problem is

investigated. While such application studies are less general, they are often highly relevant for the industry, and they provide excellent starting points for future research.

The research approach used for this thesis was both qualitative (literature review) and quantitative (controlled experiments, standard metrics and protocols). Most importantly, all empirical investigations target real-world challenges from the F&A industry using real-world data sets. With this in mind, the results are often in the form of proof-of-concepts and the published papers, consequently, in the form of application papers.

To summarise, in this thesis we present a set of experimental studies that exploit F&A data sources and use data mining techniques to improve decision support. In an experimental study, a test (or a series of tests) is performed by adjusting some factors that influence the output. While performing the experimental studies in this thesis, factors such as choice of data mining techniques and model parameters are controlled. Each of our studies is guided by the formulated research question and objective. According to the research question, we perform data collection, data pre-processing, feature selection and model construction. The data mining techniques are applied on the pre-processed data and their default parameter values are used to run the experimental analysis. As a result, we generate the model and evaluate its performance using relevant evaluation metrics. The overall steps of the experimental design applied in this thesis are illustrated in Figure 3.1.

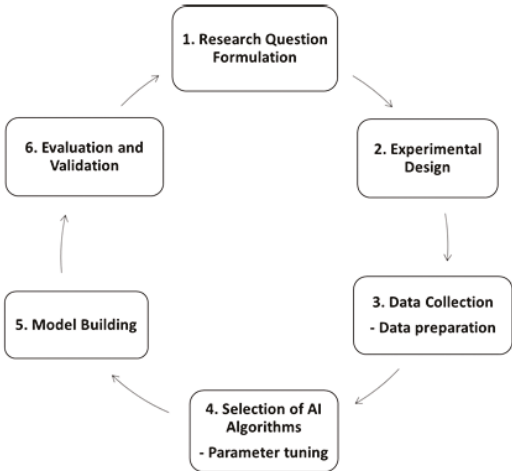


Figure 3.1: Experimental design of the thesis

### 3.2 Structure of PhD programme

The PhD project commenced in September 2017 under the framework of SMDTex (Sustainable Management and Design for textiles), a joint doctorate programme funded by the European Commission. The research plan was designed to conduct the research at three different universities: École Nationale Supérieure des Arts et Industries Textiles (ENSAIT)/The University of Lille—Sciences et Technologies in France (18 months, i.e. from Sept. 2017–Feb. 2018); University of Borås in Sweden (18 months, i.e. from Mar. 2018–Aug. 2020); and Soochow University in China (from Sept. 2020–July 2021).

The three universities specialise in different domains of textile automation, management and engineering. Therefore, the research was expected to be multidisciplinary and to contribute to the domains of management and engineering. More specifically, the aim of the research was to exploit F&A retail data using AI techniques to improve the decision-making process in the F&A retail industry.

This thesis is built on five papers that were published at different mobilities during the span of four years. Together they present three different case studies that were implemented at different mobilities.

During the first mobility, a broad understanding of AI in the F&A industry was achieved by performing an extensive literature study that focuses mainly on understanding the range of data sources and AI techniques applied at the different stages of the F&A supply chain. Further, to work on the F&A retailers' specific problem, *customer transaction* data were acquired from the European retailers. A preliminary study was performed to analyse customer purchase behaviour in F&A retail through use of an exploratory approach. Consequently, the first two articles (Paper I and Paper II) were produced during the first mobility.

Work during the second mobility continued the focus on the retailer perspective. In this mobility, *sales* and *social media* data were exploited and a SMBF (social media-based forecasting) model was implemented (Paper IV). Another, very different, data set that described *campaigns* using the attributes of different campaigns and promotions was also explored, and AI techniques were used to predict the success and profit of future campaigns (Paper V).

Due to the COVID-19 pandemic, the third mobility was conducted digitally from the University of Borås. During the third mobility, an extension of the work of the first mobility on *customer transaction* data was done by increasing the data size up to four years of data and

substantially widening the scope of the analyses. This study led to useful insights about the customers in the F&A retail industry, and as a result, a churn prediction model was implemented (Paper III). The overall research process and outcomes during the mobilities are outlined in Table 3.1.

Table 3.1: PhD program by mobility

	MOBILITY 1	MOBILITY 2	MOBILITY 3
	(ENSAIT, University of Lille, France)	(University of Borås, Sweden)	(Soochow University, China)
	Sept. 2017-Feb. 2019	March 2019-Aug. 2020	Sept. 2020-Aug. 2021
Research Outcome	<ul style="list-style-type: none"> <li>• A literature review was performed to acquire the F&amp;A domain knowledge from an AI perspective.</li> <li>• Investigated the different AI techniques that could be applied to deal with F&amp;A retail problems.</li> <li>• Data was collected from European retailer for the study of consumer behaviour in F&amp;A retail.</li> <li>• F&amp;A consumer behaviour was studied using 14 months of consumer purchase data.</li> </ul>	<ul style="list-style-type: none"> <li>• Social media (Twitter) and sales data were acquired from a European retailer.</li> <li>• Development of the SMBF forecasting model.</li> <li>• Campaign Data acquisition from a European retailer.</li> <li>• Studied the components of campaigns and developed predictive model to predict the campaign success and profit.</li> </ul>	<ul style="list-style-type: none"> <li>• Extended study of consumer behaviour with four years of data.</li> <li>• Development of churn prediction model using AI techniques.</li> </ul>
	Published Papers I and II	Published Papers IV and V	Published Paper III

### 3.3 Datasets

This section describes the datasets used in this thesis. The research studies in this thesis are broadly carried out through use of customer transaction data, sales data, Twitter data and campaign data from the F&A retail industry. These data sets were acquired from the F&A retailers, except for those utilised in the first study (Paper I).

Paper I is a literature review study, and the data is obtained from two extensive databases, *Scopus* and *Web of Science*. The detailed process of article review and selection is discussed in Paper I. To address the RQ1 of the thesis, i.e. to identify the key data-driven AI applications in the F&A industry, 149 total research articles were collected and thoroughly reviewed.

To answer RQ2 of the thesis, several data sources were explored, such as customer transaction data from F&A companies, social media data, apparel sales data and campaign data.

- **Customer transaction data**

For the case studies in Papers II and III, F&A customer transaction data were collected from a European (Italian) F&A company. The data consist of the purchase information of over 3.2 million customers. Paper II analyses and models the customer behaviour using fourteen months of historical customer data, whereas Paper III studies customer behaviour by applying advanced DM techniques to customer data that spans a four-year period (from 2013 to 2016). The data set used in Paper III contains over 28,726,819 transactions and 3,246,866 unique customers. Each row of this dataset consists of only a customer ID, the purchase date and the total revenue. Specifically, any information related to the purchased products, as well as demographic detail about the customer, is missing. However, since customer IDs are, of course, unique and remain with a specific customer, it is possible to follow the purchase history of customers over the years. The aim is consequently to exploit this very large but also rather limited data set in order to study customer behaviour.

- **Twitter data**

In Paper IV, the Twitter data of an Italian F&A retailer are obtained from the Twitter platform using Twitter API. Twitter is the sole data source, and the tweets represent the customers' opinions of the products they purchased from the F&A retailer. These tweets were retrieved and stored in the standard format. The time span chosen for the tweet collection is six months, which also corresponds to the apparel sales data used in Paper IV.

- **Apparel sales data**

In Paper IV, Twitter data are used to forecast garment sales based on apparel sales data spanning six months. The apparel sales data were provided by the Italian F&A retailer and consist of daily information on the sales of F&A products and the purchased quantity.

- **Campaign data**

The campaign data used in Paper V were obtained from a Swedish F&A retailer. This data contains information about over 826 campaigns and their 28 attributes. The data

attributes broadly consist of attributes explaining campaign types and promotions in the form of discounts, add-ons, requirements, etc.

A summary of the data used in this thesis is presented in Table 3.2.

Table 3.2: A summary of data

<b>Paper</b>	<b>Data acquisition</b>	<b>Description of the data</b>
<b>I</b>	Articles extracted from Scopus and Web of Science databases	Articles from journals and conferences since 1989 (total: 1019 Articles)
<b>II</b>	F&A retail data (14 months)	5,770,844 transactions and 1,020,923 distinct customers
<b>III</b>	F&A retail data (4 years)	28,726,819 transactions and 3,246,566 unique customers
<b>IV</b>	Twitter and apparel sales data (F&A retailer)	Six months of sales data (quantity sold and time), Twitter data attributes
<b>V</b>	Campaign data from F&A retailer	826 campaigns with 28 attributes

### 3.4 Methods

In this thesis, we employed a mixed methodological approach to perform the research studies presented in the appended papers. The methods used in each of the papers can be categorised as either qualitative or quantitative.

For Paper I, the review of existing literature on data-driven AI applications in the F&A industry is conducted using a mixed methodological approach to fulfil the purpose of the thesis research, i.e. to describe the current status of research on applications of data-driven AI in the F&A industry in terms of data and techniques. For the purpose of identifying existing research on the applications of data-driven AI in the domain of the F&A industry, the use of an exploratory approach was deemed appropriate to answer RQ1. The systematic literature review (SLR) method (Thorpe & Holt, 2007) used in Paper I is a structured method that enables the study of a vast number of existing research on the chosen research topic, i.e. the identification of the key data-driven AI applications in the F&A industry. The SLR method is

widely used to study literature on a specific research topic; additionally, it follows a systematic approach to address the goals of the literature review (Thorpe & Holt, 2007). Therefore, the SLR method is selected in order to make the research findings reproducible and transparent.

For Papers II–V, a quantitative method is used in which case studies are conducted using controlled experiments. Using this approach, data mining techniques are applied on the respective datasets used in these papers and are further evaluated following standard testing protocols, as described in Section 2.1.4. Table 3.3 summarises the tasks, their aims and the applied techniques in Papers II–V.

In Paper II, the aim is to identify the segments of customers in the customer transaction data and to analyse customers' purchase behaviour. The study of customer buying behaviour is important to the F&A retailers' ability to measure their business performance. The task in Paper II is descriptive in nature, i.e. the segments of valuable customers are to be explored from the customer data. The RFM analysis method, as discussed in Section 2.1.2, is applied to group the customers by their RFM values, which are dependent upon the identification of profitable customer segments.

In Paper III, which is the extended work of Paper II, we used four years of customer transaction data to predict customer churn behaviour. Customer churn is the tendency of customers to halt future purchasing. We used advanced AI techniques, i.e. RF and GB models, to model customer behaviour and predict churn. RF and GB are advanced AI techniques that are ensemble DT models. These models are selected to predict the customer churn because of their ability to improve generalisation performance by combining predictions from individual DTs.

In Paper IV, Twitter data are used to improve the apparel sales forecast. We used the text mining method to analyse customer sentiments from the Twitter data. Detection of the sentiment in the tweet is a predictive classification task for which the DM technique, i.e. the NB model, is used. Further, the NB model is integrated with the fuzzy time-series forecasting model, which maps the tweets' uncertainty and improves sales forecasting. The fuzzy time-series forecasting model is based on the fuzzy sets, which is discussed in detail in Chapter 2. The sentiments expressed by the customers are subjective; therefore, the fuzzy set method is selected to parse useful sentiment from the customer tweets based on a rule-based inferencing approach.

The predictive models, viz. the DT and RF models, are utilised in Paper V, the purpose of which was to analyse different attributes of campaign variables and to, based upon these attributes, predict the success and profitability of the campaigns. Studying historical campaign data to identify the factors that influence the success of campaigns is a valuable decision task that includes both the predictive regression and classification tasks. For these tasks, the DT and RF predictive models were selected because of their ability to improve generalisation performance.

Table 3.3: Summary of tasks and techniques by research paper

Research Paper	Research Aim	Tasks	Techniques
II	Customer acquisition, retention and prediction of future shopping behaviour	Exploratory	Transition matrix
III		Exploratory and predictive (classification)	Random forest (RF) and Gradient boost (GB)
IV	Improving sales forecasting using social media data	Predictive (classification and forecasting)	Naive Bayes (NB), Exponential forecasting, and Fuzzy sets
V	Improving campaign design by analysing campaign variables	Exploratory and predictive (classification and regression)	Decision tress (DT) and RF

### 3.5 Model evaluation

When evaluating the predictive performance of the AI techniques, several metrics are used to broadly assess the analytical results generated by the models. Detailed explanations of these evaluation metrics are given in Section 2.1.4.

Papers II, III and V make use of these evaluation metrics, i.e. accuracy, ROC and AUC. The accuracy of predictive models is calculated as the degree to which the observations are classified correctly. It is the proportion of correctly classified instances among the total number of instances. These papers also use other metrics, such as confusion matrix, ROC and AUC, as these are the standard visual ways of assessing the performance of classification models (see Section 2.1.4).

In Paper IV, the evaluation metric used for predictive regression, i.e. time series forecasting, is mean absolute percentage error (MAPE). MAPE is used as a performance measure to evaluate the accuracy of the forecasting model in terms of predicting the apparel sales. The



metric RMSE, which is the squared root of the average squared distance between the observed and predicted values of the target variable, is used in Paper V to evaluate the performance of the predictive regression models (DT and RF).

The evaluation metrics used for the regression, classification and forecasting tasks performed in Papers II, III, IV and V, respectively, are shown in Table 3.4

Table 3.4: The evaluation metrics used in the appended papers

<b>Research Paper</b>	Regression (RMSE, MAPE)	Classification (CA, Confusion matrix, ROC, AUC)
II		×
III		×
IV	×	
V	×	×

## 4 Summary of the appended papers

This chapter summarises the appended papers, focusing on the key results of and contributions made by each study. It first outlines the research contributions of the thesis author and then introduces the context in which the appended papers address the research questions outlined in Chapter 1.

This thesis is the result of five research publications. The links between the attached research papers and the research questions are presented in Table 4.1, and the scope of each study is outlined in Table 4.2.

Table 4.1: Categorisation of the appended papers according to the research questions

Paper	RQ1	RQ2
I	×	
II		×
III		×
IV		×
V		×

Table 4.2. Scope of the appended papers

Article	Scope	Purpose	Data	Task	Techniques
I	A comprehensive review of AI methods applied in the F&A industry	To identify the trends and gaps in the use of different types of datasets and data-driven AI applications in the F&A industry	Scopus, Web of Science		Systematic literature review
II	Customer behaviour analysis	To segment the customers based on their recency and to predict future purchase behaviour	Customer transaction data of European F&A retail	Exploratory	RFM segmentation, transition matrix
III	Churn prediction	To segment and predict customer churn for the upcoming year	Customer transaction data	Exploratory Analysis and Predictive (Classification)	Transition matrix, RF, GB
IV	Impact of fashion consumer social media data on sales	To develop a fuzzy time-series forecasting model that uses both historical sales data and social media data	Twitter data about the product and sales data (Italian/European)	Predictive (Classification and Forecasting)	Text mining, fuzzy, and forecasting (exponential model)
V	Study of the promotional campaigns of F&A retailers	To develop predictive models for estimating the outcomes of the campaigns	Campaign data	Exploratory and Predictive (Classification and Regression)	RF, DT

The research studies presented in this thesis were carried out at three separate institutions with the collaboration of several authors. Table 4.3 briefly describes the works of the different authors in each of the appended publications.

Table 4.3: Contributions of the authors in the appended papers

<b>Paper</b>	<b>First Author</b>	<b>Co-authors</b>	<b>Contribution</b>
<b>I</b>	Chandadevi Giri (CG)	Sheenam Jain (SJ), Pascal Brunaix (PB) and Xianyi Zeng (XZ)	CG and SJ are equally contributing authors, and they together conducted a review study. CG was responsible for categorising the articles from an AI perspective, while SJ was responsible for categorising the work from the F&A supply chain perspective. The study has been developed and approved in accordance with PB and XZ.
<b>II</b>	Chandadevi Giri	Sebastien Thomassey (ST) and Xianyi Zeng	The study was designed and completed with the support of ST and XZ. CG was responsible for data collection, experimental analysis and the structuring of the manuscript. ST and XZ helped to interpret the findings and refine the manuscript.
<b>III</b>	Chandadevi Giri	Ulf Johansson (UJ)	The study was designed and carried out under the supervision of UJ. All authors contributed to finalising the methodology. CG performed experimental analysis and wrote the manuscript. UJ contributed to the interpretation and refinement of the results of the manuscript.
<b>IV</b>	Chandadevi Giri	Sebastien Thomassey and Xianyi Zeng	The study was designed and conducted with the support of ST and XZ. CG collected data, carried out the experimental analysis and wrote the manuscript. ST and XZ helped to interpret the findings and to refine the manuscript.
<b>V</b>	Chandadevi Giri	Ulf Johansson and Tuwe Löfström (TL)	The study was designed and conducted with the support of UJ and TL. All authors contributed to the finalisation of the methodology. CG performed experimental analysis and wrote the manuscript. UJ and TL helped to interpret the findings and refine the manuscript.

#### **4.1 Paper I: A detailed review of artificial intelligence applied in the fashion and apparel industry**

This section includes two parts. The first, Section 4.1.1, summarises the research from the paper that focuses on classifying and analysing articles from an applied AI techniques perspective in the F&A supply chain. The second, Section 4.1.2, is the extended part of 4.1.1, and it analyses the extracted 149 articles based on the data used in the F&A supply chain.

#### 4.1.1 Summary of Paper I

**Purpose:** In the past few decades, AI has played an important role in transforming the F&A industry. However, research in this area is scattered and mainly concentrated at one stage in the supply chain. Therefore, it is difficult to understand the work done in the distinct domains of the F&A industry. Furthermore, there is a need to more closely examine the data-driven AI methods applied at various supply chain levels to optimise business operations. However, with the dissemination of AI technology, the complexity of business operations has increased, making it imperative for the various entities (industrial experts, academic researchers, and managers) to have interdisciplinary expertise. The purpose of this work is to conduct a literature review and analyse the current trends and applications of data-driven AI in the F&A industry.

**Methodology:** The literature review study conducted in Paper I is based on the systematic literature review (SLR) research method. An SLR methodology was chosen to make the research more rational, transparent and reproducible (Thorpe & Holt, 2007). The complete research workflow includes article retrieval, article selection, information extraction, article classification, analysis and findings. The article retrieval resulted in 1019 articles. This was followed by article selection, which included five steps and led to the selection of 149 articles deemed to be within the scope of this study. Hereafter, information was extracted by thoroughly examining the articles aligned with the research questions. These articles were classified on the basis of F&A supply chain stages and AI techniques.

**Findings:** This paper gives a comprehensive review of research on data-driven AI applications in the F&A industry. The overall trend of the research articles on applications of AI reflect that the majority of the studies, 56% of the total research, were conducted in the last decade. It was found that ML and expert systems were the most applied AI techniques in the F&A industry, with a significant focus on apparel, fabric production and distribution and the least focus on the design stage. A large number of studies focused on F&A manufacturing and production, whereas relatively few focused on applications for the design stage of the F&A supply chain. There was also a lack of research on the application of advanced AI techniques, such as ensemble learning and deep learning. It was found that genetic algorithm (GA), artificial neural network (ANN) and fuzzy logic are the most widely used methods for prediction (e.g. yarn, fabric properties, and colour), evaluation (e.g. quality inspection and defect detection) and forecasting (e.g. trend analysis, supplier selection, and demand forecasting).

#### 4.1.2 Extended summary of Paper I

**Purpose:** The purpose of this extended analysis was to review the 149 articles (87 journals, 62 conferences) selected in Paper I and to further categorise and analyse them based on the data types (semi-structured, structured and unstructured) used in the F&A supply chain stages.

**Methodology:** The articles were studied to extract information on the different kinds of data used in the F&A supply chain stages.

**Findings:** The distribution of classification of articles by data types in the F&A supply chain is shown in Figure 4.1. It can be observed that most of the data used in the F&A industry were structured, and this accounts for 79% of the journal articles and 67% of the conference articles. On the other hand, few articles use semi-structured data; less than 5% of the total reviewed articles utilised semi-structured data. Last, about 17–31% of the articles are found to make use of unstructured data. The above distribution is further categorised by F&A supply chain stages (fabric production, apparel production, design and distribution), as illustrated in Figure 4.2.

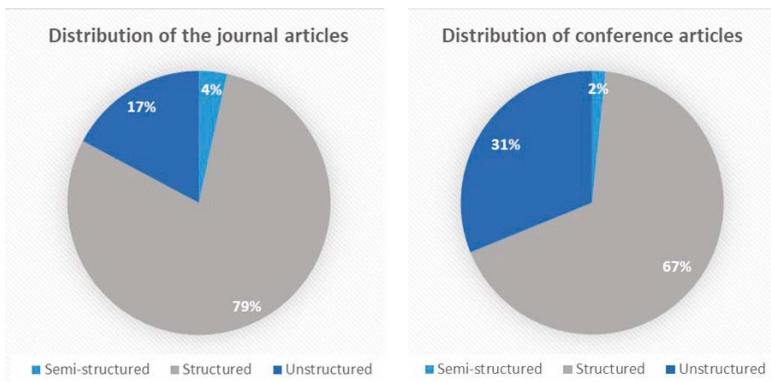


Figure 4.1: Overall distribution of articles by data types

The most used data in all supply chain stages was structured data, while unstructured data is the second most used data type. The semi-structured data type is the least used in both journal and conference articles for the apparel design and distribution stages. The categorisation of the different data types used in the 149 articles per the F&A supply chain stages can be seen in Table 4.4.

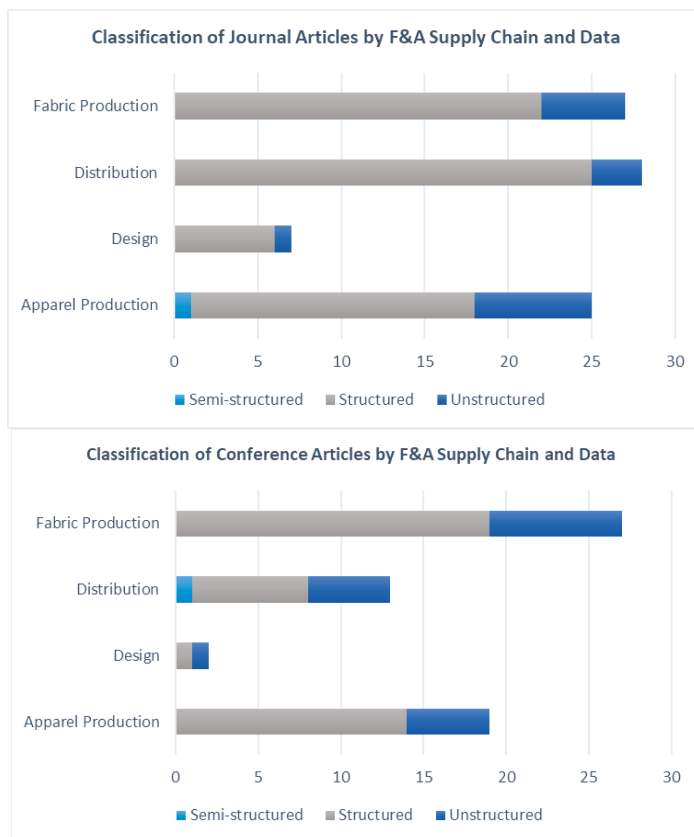


Figure 4.2: Distribution of articles by data and F&A supply chain stages

Table 4.4 : Categorisation of the data used in the articles per the F&A supply chain stages

Supply Chain Stage	Semi-Structured	Structured	Unstructured
Apparel Production	A database containing anthropometric measurements, demographic profiles and 3D body scans of samples of the studied population	Fabric data (properties, categories, types, colour, etc.), apparel production data from production plant with information about sewing machines and apparel produced, production allocation data (to allocate the resources in production)	3D scan of the body (anthropometric data), fashion colour image data, textile clothing image data, fabric defect and stitching defect images
		Apparel pattern data (type of garment, size and quantity)	Pattern data (images taken of left side of body, right/left sleeve, left overlap) and body posture 3D scan
		Manufacturing data with material information	Textile pattern image (square, triangle, circle, horizontal/vertical line, flower, curve, leaves, etc.)

		Production workstation data with information about the spreading, cutting, sewing, finishing and quality inspection workstations	
		Fabric tensile properties data and mechanical properties of fabrics	
<b>Design</b>		Clothing style and pattern data to design 3D clothes, design data (cutting pattern, outline curve and coordinates), and rating data used for design recommendation	
<b>Fabric Production</b>	Social media data (Instagram)	Anthropometric numerical data	
		Fabric properties data (fibre thickness, weft count, warp count, tightness, weight, wool content), physical and chemical properties of fabric data	Fabric image data with yarn /fibre defects (horizontal/vertical missing yarn, colour yarn, spot, hole, stitching defect) and defect-free
		Production scheduling data (list of jobs in production inline, machine effectiveness, type of product, and colour)	Fabric images pattern defects (stains, holes, texture properties like a shadow, different edges, folds, etc.)
		Fabric quality data (trash code, trash area, trash particle count, length, uniformity, short fibre index, strength, elongation, maturity, white level and yellow level)	Woven fabric images (plain, twill, stain), and fabric image defects (fabric flaws caused by material texture)
		Fibre data coming from industrial processes (fibre diameter, hauteur, fibre length distribution, short fibre content, yarn count, twist, draft ratio, and spinning speed), and yarn data (characteristics of fibre)	
		Friction and lubrication data of machines where fabric is produced, and production data (machines and materials)	
		Weaving data (warp-related stop when one or more warp ends break), filling-related stop (filling yarn fails to pass through the warp successfully) and other stops (all other reasons, including operator intervention, stop for maintenance action, power outage, etc.)	
		Spinning mills data (for predicting yarn strength)	
		Dyeing order data (for fabric roll selection from inventory, picking, machine loading and scheduling), dyeing data (class of dyes, group of dyes, dye, process, and machines), cotton allocation data (yarn quality and yarn product) and cotton allocation list (brands in production line)	
		Scheduling jobs data (machine loads, fabric roll selection and sequencing operations, SKU)	
<b>Distribution</b>	Social media data (Instagram)	Sales data (product, quantity, date sold, price, store), inventory data (batch size, cost, speed, manufacturers, and suppliers)	Blogs, online reviews, news, magazines, clothing article from the affiliate network
		Retail data with information about the product and sales	3D scan from garment try-on
		Retail data with information about suppliers and products, and retailer catalogue data with information about the product	Street photos (enriched with outfits presented on different body shapes of individuals belonging to various age groups)

## 4.2 Paper II: Customer analytics in the fashion retail industry

**Purpose:** Over the last few years, the F&A retail market has faced significant challenges due to increasingly varying customer demands. This paper explores how customer analytics in the F&A retail industry can improve decision-making through use of statistical and AI techniques to analyse customer behaviour. By investing in these strategies and in data-driven AI techniques, the F&A retail industry could benefit from increased revenue, improved profits and higher customer satisfaction rates, thereby sustaining growth in volatile markets. This study aims to provide an overview of customer analytics in the age of information and to study customer behaviour using an exploratory approach.

**Methodology:** First, customer analytics strategy, scope and methods within the scope of the F&A industry are conceptualised and discussed, after which the exploratory case study on its application is conducted using the sales data. Customer *Recency* value is computed to map the customers' recent purchase behaviour. Two segments were created for the six-month interval and customers were divided based on their *Recency* values into four classes: *Inactive*, *Less active*, *Active* and *Highly active*. Transition matrices were used to measure the probabilities of customer movement from one class to another and revenue predictions for the next two six-month intervals.

**Findings:** The study finds that customer analytics is dependent on the companies' business goals and the problem they are attempting to solve. It could be based on behavioural or longitudinal social network analysis, or it could be focused on predicting customer behaviour at an individual level without considering any other information. Data can be collected and utilised to understand F&A customers from social networks (data from tweets, blogs, brand events, etc.), internal data (CRM, ERP, e-commerce, etc.) and external data (cookies and plug-ins). Collected data can be pre-processed using exploratory, descriptive and predictive techniques, depending on the level of analysis to be performed. Customer lifetime value (CLV), RFM modelling and churn prediction are some of the types of analysis used to analyse customer behaviour. The results from the case studies show that 75% of the inactive customers remain inactive in the next six months, while over 8% of the inactive customers become highly active. Over 10% of the newly acquired customers will remain highly active,



5% will remain active and about 69% will become inactive. It is clear that inactive customers do not generate revenue, and therefore acquisition of new customers is necessary for the generation of high revenue. The exploratory study using the statistical approaches presented in this paper can be used to understand the future behaviour of current customers. In the F&A industry, segmentation like this would help us identify which customers generate the most value for the company and determine how to build stable brand loyalty.

### **4.3 Paper III: Data-driven business understanding in the fashion and apparel industry**

**Purpose:** In retailing, data analytics is widely used as a means of gaining customer insights. Most retail datasets include personal or, at the very least, demographic information about customers. However, in some datasets, the only information available is that of the purchases made and their links to customers. In this application paper, an analysis was performed that used a very large real-world data set from the fashion retailing industry. The most interesting aspect of the data set, in spite of it only containing the actual purchases connected to the customer ID, was the fact that it covered four years in total. With this in mind, the primary goal of this study was to determine what an F&A company can learn about their business and their customers from such a data set. The purpose of this research study is:

- to analyse F&A sales data, with very limited information about the customers, by using a combination of descriptive and exploratory methods.
- to predict and understand customer churn using predictive modelling.

**Methodology:** The research methodology for this study is briefly described below:

- *Data set used*—The customer transaction data were collected from the Italian F&A retailer between 2013 and 2016. There are a total of 28,726,819 transactions and 3,246,866 unique customers in the data set.
- *Exploratory study*—Two different analyses were carried out with the overall goal of exploiting this very large but also quite limited data set in order to understand customer behaviour. In the first exploratory phase, aggregated customer data was analysed to understand the overall dynamics of the customers. Following that, an RFM analysis, which segmented the customers into 3 x 3 cubes, was performed for each

year. In the final section of the exploratory study, transition matrices were created in order to analyse the customer flows between RFM cells from one year to the next.

- *Customer churn modelling*—In the second part of the study, predictive models were generated to predict whether customers would remain loyal or churn, i.e. not make a purchase the following year. In the first predictive experiment, which targeted all active customers, the models were trained on inputs from 2013–2014 and on targets from 2015, while the evaluation was done with inputs from 2014–2015 and with targets from 2016. The second predictive experiment focused on the most important customers who had remained loyal for the first three years, providing churn predictions for the fourth year. Since the churners represent a minority class in this situation, under-sampling was used to create a balanced training set.
- *Evaluation*—The predictive task performed in this study was classification, so the models were evaluated using classification metrics. These quantitative results were, however, complemented by a qualitative analysis of the produced rule sets.

**Findings:** Many interesting patterns were discovered and analysed using a combination of descriptive methods, traditional RFM analysis and predictive modelling. A key insight from the exploratory study was that nearly half of all customers churn every year, indicating a very low level of customer loyalty. However, the churn is offset by the addition of a comparable number of new customers. On the other hand, customers that remain loyal for a couple of years are rather unlikely to churn in the future. As a result, churn prediction and churn prevention should be a top priority. In this study, trained predictive models were found to be sufficiently accurate, with a reasonable trade-off between precision and recall. When creating interpretable models, in the form of rule sets, to predict which of the loyal customers were most likely to churn the following year, a few rules stood out. In particular, few purchases and/or a decreasing number of purchases were discovered to be key indicators.

#### **4.4 Paper IV: Exploitation of social network data for forecasting garment sales**

**Purpose:** In today's digital world, customers can express and share their thoughts about their different shopping experiences through social media. Knowing consumer experiences in a consumer and market research framework is essential for the businesses to enhance the products and services they provide. The collection and use of such information have been

significant challenges for the F&A industry. Traditionally, companies depend on sales forecasting models to predict their potential product sales, which in turn allow them to plan their business operations efficiently. The fashion industry faces operating difficulties arising from uncertainties in product demand and customer preference. In this context, traditional forecasting models that consider the influence of marketing campaigns and promotions on retail purchases have not been able to capture the effects of social media. The purpose of this research is to study the influence of social media on sales by exploiting Twitter data to potentially improve sales forecasting.

**Methodology:** The complete research workflow for this study includes the following key steps: data collection, data pre-processing, correlation analysis, SMBF modelling and performance evaluation.

- *Data sets used*—Twitter and sales data were collected from the Italian F&A retailer over a period of six months.
- *Data pre-processing and sentiment extraction*— In this step, the text mining technique was used to clean the tweets. These cleaned tweets were used to extract sentiments using advanced NRC lexicon data and the NB classifier. Classification results fell into three sentiment indices, viz., positive ( $S_p$ ), negative ( $S_n$ ) and neutral ( $S_{ne}$ ).
- *Correlation analysis*—Extracted sentiments were aggregated weekly, and the correlation between the sentiments of the current week and the sales of the following week was computed. This step was performed to evaluate the correlation between tweets and sales.
- *SMBF (social media-based forecasting) model*—After investigating the relationship between tweet sentiment and sales, the SMBF model was developed to forecast apparel sales based on the historical sales data and social media sentiments. To this end, the FSIS model was developed to define the uncertainty of sentiments in terms of the corrective coefficient ( $V_s$ ), which describes the sales variance factor with actual sales. This was implemented with the exponential forecast (EF) model to estimate the final forecast of the SMBF model.
- *Performance validation*—The performance of the SMBF model was evaluated using the forecasting metric MAPE.

**Findings:** Of the classified tweets, over 47% were classified as positive, 27% as negative and 28% as neutral. The results of the classification demonstrate that customers have a rather

positive view of the selected fashion brand, as only 27 % of the overall tweets were negative. It was found that there is a correlation between sales and consumer tweets, which implies the influence of customer tweets on sales. The SMBF model, which is a fuzzy time-series forecasting model, was developed based on both the historical sales data and social media data. The SMBF model was assessed, its performance was compared with the EF model and SMBF was found to outperform EF. The results from this paper highlight that social media data help to improve the forecasting of garment sales and that the proposed model could be easily integrated with any time-series forecasting technique. The findings of this paper contribute to addressing the complex decision problem of forecasting fashion product sales by using social media data, which can provide decision-makers in the F&A industry with the valuable practical perspectives needed to optimise their market strategies.

#### **4.5 Paper V: Predictive modelling of campaigns to quantify performance in the fashion retail industry**

**Purpose:** The popularity of promotional campaigns has experienced growth in parallel with the explosion of customer shopping interactions across multi-channels. By offering attractive discounts or services, such as free delivery, free returns and gifts, promotional campaigns help fashion retailers retain customers. However, promotional campaigns are expensive for the retailers, especially if the conversion rates, and consequently the revenues, are low. The aim of this study is to use predictive modelling to identify the profitability and success rate of campaigns and to present a data-driven campaign prediction model, or simulation engine, for the F&A industry. With the proposed setup, profit and activation levels are to be predicted based on the different properties of a campaign. The individual components of prior campaigns can be freely combined into new promotions and campaigns, i.e. the simulation engine is not restricted to previous campaigns.

**Methodology:** The complete research framework developed for the study includes data collection, exploratory study, campaign modelling and model evaluation.

- **Data set used:** Campaign data with information about the campaign types, discount features and add-ons were collected from a Swedish fashion retailer.
- **Exploratory data analysis:** The exploratory study of the campaign types is discussed in detail in the paper (see Section IV.B in Paper V). The distribution of the type of

order features shows that more than 50% of the orders had the promotion *firstline*. Distribution of average profit per order through the campaign illustrates that the average profit is right-skewed, with most orders having a positive profit, whereas the distribution of the activation rate is highly skewed towards the left. It was found that the response rate of customers is typically less than 2%. Based on activation and profit, three class labels were defined: *unsuccessful*, *successful* and *highly successful*. These were used as target variables for a classification model.

- **Modelling:** Two ML models, DTs and RF, were applied, and their performance was compared.
- **Evaluation:** The profitability prediction model was evaluated using the regression metrics RMSE, MAE and  $R^2$ . The campaign success prediction model was evaluated using the classification metrics AUC and accuracy and by analysing the confusion matrix.

**Findings:** While the opaque RF model outperformed the interpretable decision and regression trees in terms of accuracy, some general and potentially valuable insights were derived by analysing the interpretable models. In particular, a high discount on the first item leads to highly profitable campaigns, and combining the discount with a free gift also leads to a highly successful campaign. Thus, this case study has demonstrated that data-driven methods can be used to understand and ultimately optimise campaigns and promotions. Such models could be used to simulate campaigns before they are executed, potentially providing retailers with a sophisticated campaign planning tool.

## 5 Conclusions

This chapter aims to conclude the RQs based on the results discussed in Chapter 4.

### 5.1 Concluding RQ1

*What are the widely used data types and key data-driven AI applications in the F&A industry?*

In this thesis, the findings from the literature study, which was aimed at exploring and identifying the different types of data and key data-driven AI applications in the F&A industry, suggest that structured data is the predominant data type in all supply chain stages of the F&A industry. Unstructured data is the second most used data type, while semi-structured data is the least used data type in the F&A industry. The use of semi-structured data mainly exists in the apparel design and distribution stages of the supply chain.

There is an increasing amount of research on the application of data-driven AI in the production and distribution phases of the value chain. The main activities in these phases for which AI has been employed include planning and scheduling fabric and apparel production as well as logistic networking, in which production floor assembly, cutting, inventory allocation and construction of supply networks are the major decision problems. While these stages have employed AI significantly, the design stage has received less attention. One of the reasons for this could be that the dynamic data related to customer specifications are neither generated nor utilised due to the lack of data sharing from other supply chain stages. Of the identified data-driven AI applications applied in the F&A industry, machine learning and expert systems are found to be the most utilised techniques. As presented in Paper 1, the categorisation of AI applications according to the supply chain stages in the F&A industry could be valuable for F&A firms, as it could help them decide on an appropriate AI approach for decision support in various operational levels in the supply chain stages.

### 5.2 Concluding RQ2

The conclusion of the three sub-research questions of RQ2 will be presented first, followed by the overall conclusion of RQ2.

### *How can data-driven AI techniques support and improve F&A retailers' decision-making?*

The first research sub-question is:

- a) How can fashion customer transaction data be utilised to analyse and predict customer behaviour?*

This thesis has presented a number of different examples in which transaction data can be utilised by F&A retailers as a basis for decision support. While the main technique is supervised learning, the applications vary depending on both the exact data and the problem formulation. A key insight from the empirical work of the thesis is therefore that data-driven approaches can support many different data mining tasks by using the supervised learning paradigm on processed transaction data. Specifically, it has been demonstrated that the predictive approach can be combined with both descriptive techniques and more exploratory methods in order to gain a better understanding of the business and customers. With these approaches, the F&A firms could, for instance, design better marketing and promotion strategies as well as more fully understand customer churn.

The second sub-question is:

- b) How can social media data be utilised to improve sales forecasts?*

As demonstrated, augmenting customer data with sentiments from social media can improve the quality of sales forecasting, especially for short-term forecasts. A specific example is the impact of consumer sentiments harvested in the current week on the product sales of the next week. Improved sales forecast accuracy, as demonstrated in the thesis results, could be instrumental in the F&A retail industry's struggle to reduce overproduction and thereby reduce waste and production costs. The F&A industry could improve the quality of their products and services by focusing more closely on the current trends and customer sentiments about them that can be extracted from social media platforms.

- c) How can historical campaign data be utilised to design successful campaigns and promotions?*

In the empirical work, it has been shown that data analysis, in the form of predictive modelling, can be applied to planning and simulation of promotions and campaigns, a task that normally does not utilise data-driven approaches. Applied predictive techniques could enable F&A retailers to analyse the effectiveness of the campaigns prior to their strategic

execution. The identification of the key attributes that influence customers to make a purchase from a campaign or promotion could also be achieved by using an approach based on the techniques presented in this thesis.

### **Overall conclusion of RQ2:**

*How can data-driven AI techniques support and improve F&A retailers' decision-making?*

The data-driven projects presented in this thesis span a number of applications in F&A retailing, thus showing the overall benefit of utilising data-driven AI for decision support in that domain. One important lesson is that it is often necessary to combine not only different data sources but also techniques in order to provide a deeper understanding of the targeted aspect. This was demonstrated by combining transactional data with sentiment data in one study but proved even more critical when analysing processed sales data using predictive, descriptive and explorative techniques. The result was a complex and quite comprehensive description of a company's business and customers over time, despite utilising a data set that described the customers only by their orders (in the form of date of purchase and money spent). Another general result is that while opaque predictive models typically outperform interpretable alternatives, general and valuable insights can be found by inspecting and analysing comprehensible models.

Reviewing the data mining tasks that have been addressed, churn prediction and sales forecasting have, of course, been frequently targeted with data-driven approaches. The novelty in these examples therefore lies mainly in the kind of data used for the actual modelling. Normally, churn prediction utilises personalised data, containing at the very least some demographic data. However, this research has shown that application of RFM-modelling as a feature engineering step could be used to obtain reasonable predictive performance and discover interesting patterns from just the sales data. Regarding sales forecasting, the main contribution is the addition of sentiment data to improve short-term predictions. The campaign and promotion optimisation example, finally, is more novel in itself, and the presented study shows the benefit of using the data-driven approach. Specifically, the generated predictive model makes it possible to simulate, i.e. predict the activation level and profit of, combinations of promotions and campaigns that have not actually been used.

All in all, this thesis has provided a number of examples that demonstrate how data-driven AI techniques can support and improve F&A retailers' decision-making.



**Overall conclusion of the thesis:**

In line with the purpose of this thesis, a series of studies have demonstrated that the different types of data generated by F&A retailers can be processed and transformed into valuable business and customer insights through use of a number of data-driven AI techniques, including both descriptive and predictive alternatives. On a high level, the F&A industry can use data-driven AI techniques to improve production, as well as supply chain operations, by using data from many different sources and in extremely varied formats. The data-driven projects used as case studies or proofs-of-concept in this thesis were typically undertaken by combining different data sources and techniques in novel ways, with the overall purpose of improving decision support in the F&A industry.

## 6 Discussion and future work

This section briefly considers both the importance and some potential limitations of the research conducted in this thesis. In addition, a number of suggestions for future work is presented.

In this thesis, different AI techniques and varying types of data have been explored for improved decision support in the F&A retail industry. As described in Chapter 3, the studies performed should be regarded as proofs-of-concept, demonstrating the potential of data-driven approaches utilising AI techniques. With this in mind, the main contribution of the thesis is the combined effect of all the studies, which shows the versatility of data-driven AI, rather than each and every method, algorithm and micro-technique used in the studies. Consequently, this thesis should not be regarded as a set of best-practices recipes that could be utilised off-the-shelf by an F&A retailer. Having said that, some approaches, like combining sales data and sentiment analyses for short-term forecasting as well as utilising information from past campaigns to build a simulation engine for campaigns and promotions, could probably be successfully employed by a number of different retailers with a minimum of fine-tuning. However, as always when it comes to data-driven AI, two necessary requirements for success are a deep understanding of a business problem that should be tackled utilising data analytics, and, of course, access to high-quality data.

In the first case study (Paper II and Paper III), the transaction data that were used to study customer behaviour had very limited information about the customers and orders. Demographic data (gender, age, country, region, etc.) and product-related attributes (size, style, colour, etc.) were both missing. In this respect, it could be interesting to study the impact of these attributes on the performance of the applied AI techniques.

The second case study (Paper IV) was focused on weekly forecasting using only an F&A retail brand's historical sales data. In the future, a short-term forecasting model, as per the product categories, could produce interesting contributions for the product-level forecasting. One important suggestion would be to integrate the FSIS model with other forecasting models, such as ARIMA, GARCH, etc. Another interesting approach for future research could be the exploitation of product image data to enhance the performance of the FSIS forecast model presented in this thesis.

In the third case study (Paper V), we demonstrated how advanced AI techniques (DT and RF) can predict the performance of campaigns. This was performed by studying campaign data, in which the main attributes were discounts and add-on features. In this context, the combination of enriched campaign data with more meta-attributes, including the customer demographics, purchase type and product category, and advanced AI techniques could be an interesting alternative worth exploring in order to further improve predictions of campaign success.

Finally, it must be noted that research looking at how the results from data-driven approaches are actually used by retailers is generally lacking. Which strategies are used on top of the predictive models, e.g. forecasting reduced sales the next week due to some negative sentiments in social media? How do marketers benefit from access to interpretable models describing, for instance, customers likely to churn? Are promotional campaigns created *in silico* received as expected or is a final “human touch” necessary? These and many similar questions should be interesting for both marketing and computer science scholars to study since it is *how* the AI is used that will ultimately determine success rather than the exact AI technique. What is obvious, however, is that data-driven AI will be pervasive in F&A retailing, as well as in general retailing and in many other application areas.

## References

- Aggarwal, C. C. (2017). *Outlier analysis*. Springer International Publishing.
- Agrawal, D., & Schorling, C. (1996). Market share forecasting: An empirical comparison of artificial neural networks and multinomial logit model. *Journal of Retailing*, 72(4), 383–407.
- Alturki, A., Bandara, W., & Gable, G. G. (2012). Design science research and the core of information systems. In *International Conference on Design Science Research in Information Systems* (pp. 309-327). Springer, Berlin, Heidelberg.
- Armstrong, J. S. (Ed.). (2001). *30 principles of forecasting*. Springer.
- Baines, P., Fill C., & Page K. (2013). *Essentials of marketing*. Oxford University Press.
- Ballings, M., & Van den Poel, D. (2012). Customer event history for churn prediction: How long is long enough? *Expert Systems with Applications*, 39(18), 13517–13522.
- Barnes, L., & Lea-Greenwood, G. (2010). Fast fashion in the retail store environment. *International Journal of Retail and Distribution Management*, 38(10), 760–772.
- Barnett, V. (1978). *Outliers in statistical data*. John Wiley & Sons.
- Battista, C., & Massimiliano, M. S. (2013). The logistic maturity model: Application to a fashion company. *International Journal of Engineering Business Management*, 5, 5-29.
- Beg, I., & Ashraf, S. (2009). Similarity measures for fuzzy sets. *Applied and Computational Mathematics*, 8(2), 192–202.
- Beheshti-Kashi, S. et al. (2015). A survey on retail sales forecasting and prediction in fashion markets. *Systems Science & Control Engineering: An Open Access Journal*, 3(1), 154–161.
- Berry, M. J. A., & Linoff, G. S. (2004). *Data mining techniques: For marketing, sales, and customer relationship management*. John Wiley & Sons.
- Bhardwaj, V., & Fairhurst, A. (2010). Fast fashion: Response to changes in the fashion industry. *The International Review of Retail, Distribution and Consumer Research*, 20(1), 165–173.
- Bishop, C. M. (2006). *Pattern recognition and machine learning*. Springer.
- Le Bon, C. (2014). *Fashion marketing: Influencing consumer choice and loyalty with fashion products*. Business Expert Press.
- Bradlow, E. T., Gangwar, M., Kopalle, P., & Voleti, S. (2017). The role of big data and predictive analytics in retailing. *Journal of Retailing*, 93(1), 79–95.
- Breiman, L. (2001). Random forests. *Machine Learning*, 45(1), 5–32.
- Bruce, M., Daly, L., & Towers, N. (2004). Lean or agile: A solution for supply chain management in the textiles and clothing industry? *International Journal of Operations and Production Management*, 24(1–2), 151–170.
- Brynjolfsson, E., Hitt, L. M., & Kim, H. H. (2011). Strength in numbers: How does data-driven decisionmaking affect firm performance? *SSRN* (1819486).

- Budd, J., Knizek, C., & Tevelson, B. (2012). The demand-driven supply chain: Making it work and delivering results. *Boston Consulting Group*.
- Burez, J., & Van den Poel, D. (2009). Handling class imbalance in customer churn prediction. *Expert Systems with Applications*, 36(3), 4626–4636.
- Cambridge University. (2021). *Online English dictionary*. Retrieved August 23, 2021, from URL <https://dictionary.cambridge.org/dictionary/english/clothes>.
- Choi, T. M. (2016). Information systems for the fashion and apparel industry. *Information Systems for the Fashion and Apparel Industry*. Woodhead Publishing.
- Choi, T. M., Hui, C. L., Ng, S. F., & Yong, Y. (2012). Color trend forecasting of fashionable products with very few historical data. *IEEE Transactions on Systems, Man and Cybernetics Part C: Applications and Reviews*, 42(6), 1003–1010.
- Christopher, M., Lawson, R., & Peck, H. (2004). Creating agile supply chains in the fashion industry. *International Journal of Retail & Distribution Management*, 32(8), 367–376.
- Cooper, M. C., Lambert, D. M., & Pagh, J. D. (1997). Supply chain management: More than a new name for logistics. *The International Journal of Logistics Management*, 8(1), 1–14.
- Craven, M., & Shavlik, J. (1999). *Citeseer Rule extraction: Where do we go from here*. University of Wisconsin Machine Learning Research Group working Paper, 99.
- Dahana, W. D., Yukihiro, M., & Morisada, M. (2019). Linking lifestyle to customer lifetime value: An exploratory study in an online fashion retail market. *Journal of Business Research*, 99, 319–331.
- Dubois, D., Prade, H. M., & Henri. (1980). *Fuzzy sets and systems: Theory and applications*. Academic Press.
- Edelstein, H. (2001). Building profitable customer relationships with data mining. In *Customer relationship management* (pp. 339–51) Vieweg+ Teubner Verlag, Wiesbaden.
- De Felice, F., Gnoni, M. G., & Petrillo, A. (2012). A multi-criteria approach for sustainable mass customisation in the fashion supply chain. *International Journal of Mass Customisation*, 4(3–4), 220–238.
- Ferraris, A., Mazzoleni, A., Devalle, A., & Couturier, J. (2019). Big data analytics capabilities and knowledge management: Impact on firm performance. *Management Decision* 57(8), 1923–1936.
- Fitzgerald, M., Kruschwitz, N., Bonnet, D., & Welch, M. (2013). Embracing digital technology: A new strategic imperative. *MIT Sloan Management Review*, 55(2), 1.
- Friedman, J. H. (2001). Greedy function approximation: A gradient boosting machine. *The Annals of Statistics*, 29(5), 1189–1232.
- Gottwald, S. (2010). An early approach toward graded identity and graded membership in set theory. *Fuzzy Sets and Systems*, 161(18), 2369–2379.
- Grabher, G., Ibert, O., & Flohr, S. (2008). The neglected king: The customer in the new knowledge ecology of innovation. *Economic Geography*, 84(3), 253–280.
- Grieco, A., Caricato, P., Gianfreda, D., Pesce, M., Rigon, V., Tregnagli, L., & Voglino, A. (2017). An industry 4.0 case study in fashion manufacturing. *Procedia Manufacturing*,

11, 871-877.

- Griva, A., Bardaki, C., Pramatari, K., & Doukidis, G. (2021). Factors affecting customer analytics: Evidence from three retail cases. *Information Systems Frontiers*, 1–24.
- Griva, A., Bardaki, C., Pramatari, K., & Papakiriakopoulos, D. (2018). Retail business analytics: Customer visit segmentation using market basket data. *Expert Systems with Applications*, 100, 1–16.
- Gupta, S., Kar, A. K., Baabdullah, A., & Al-Khowaiter, W. A. A. (2018). Big data with cognitive computing: A review for the future. *International Journal of Information Management*, 42, 78–89.
- Haenlein, M., & Kaplan, A. (2019). A brief history of artificial intelligence: On the past, present, and future of artificial intelligence. *California Management Review*, 61(4), 5–14.
- Hand, D. J. (2007). Principles of data mining. In *Drug safety* (pp. 621–622).
- Huang, H., Huang, H., & Liu, Q. (2017). Intelligent retail forecasting system for new clothing products considering stock-out. *Fibres and Textiles in Eastern Europe*, 25(0), 10–16.
- Hüllermeier, E. (2005). Fuzzy methods in machine learning and data mining: Status and prospects. *Fuzzy Sets and Systems*, 156(3), 387–406.
- Humphries, M. (2009). *Fabric glossary*. Pearson/Prentice Hall.
- Jin, B., & Bennur, S. (2015). Does the importance of apparel product attributes differ by country? Testing Kano's theory of attractive quality in four countries. *Clothing and Textiles Research Journal*, 33(1), 35–50.
- Kaplan, A., & Haenlein, M. (2019). Siri, Siri, in my hand: Who's the fairest in the land? On the interpretations, illustrations, and implications of artificial intelligence. *Business Horizons*, 62(1), 15–25.
- Kawamura, Y. (2005). *Fashion-ology: An introduction to fashion studies: Dress, Body, Culture*. Oxford, UK: Berg Publishers.
- Kuhn, M., & Johnson, K. (2013). *Applied predictive modeling*. Springer.
- Lam, J. K. C., & Postle, R. (2006). Textile and apparel supply chain management in Hong Kong. *International Journal of Clothing Science and Technology*, 18(4), 265–277.
- Lee, H. L., Padmanabhan, V., & Whang, S. (1997). Information distortion in a supply chain: The bullwhip effect. *Management Science*, 43(4), 546–558.
- Liu, N., Ren, S., Choi, T. M., Hui, C. L., & Ng, S. F. (2013). Sales forecasting for fashion retailing service industry: A review. *Mathematical Problems in Engineering* 2013.
- Martino, G., Iannone, R., Fera, M., Miranda, S., & Riemma, S. (2017). Fashion retailing: A framework for supply chain optimization. *Uncertain Supply Chain Management*, 5(3), 243–272.
- Miklosik, A., & Evans, N. (2020). Impact of big data and machine learning on digital transformation in marketing: A literature review. *IEEE Access*, 8, 101284–101292.
- OECD Publishing. (2019). *Going digital: Shaping policies, improving lives*. Organisation for Economic Co-operation and Development OECD.

- Olson, D. L. (2007). Data mining in business services. *Service Business*, *1*(3), 181–193.
- Pagani, M., & Pardo, C. (2017). The impact of digital technology on relationships in a business network. *Industrial Marketing Management*, *67*, 185–192.
- Pang-Ning Tan, Steinbach, M., Karpatne, A., & Kumar, V. (2016). *Introduction to data mining* (2nd ed.). Pearson.
- Pedregosa, F., Varoquaux, G., Gramfort, A., Michel, V., Thirion, B., Grisel, O., ... & Duchesnay, E. (2011). Scikit-learn: Machine learning in Python. *The Journal of machine Learning research*, *12*, 2825-2830.
- Pedrycz, W. (1994). Why triangular membership functions? *Fuzzy Sets and Systems*, *64*(1), 21–30.
- Peng, Y., Kou, G., Shi, Y., & Chen, Z. (2008). A descriptive framework for the field of data mining and knowledge discovery. In *International journal of information technology and decision making* *7*(04), 639–682. World Scientific Publishing Company.
- Piatetsky-Shapiro, G., & Parker, G. (2011). Lesson: Data mining and knowledge discovery: An introduction. *Introduction to Data Mining*.
- Provost, F., & Fawcett, T. (2013). Data science and its relationship to big data and data-driven decision making. *Big Data*, *1*(1), 51–59.
- Pyle, D. (1999). *Data preparation for data mining*. Morgan kaufmann.
- Ramaswamy, V., & Ozcan, K. (2016). Brand value co-creation in a digitalized world: An integrative framework and research implications. *International Journal of Research in Marketing*, *33*(1), 93–106.
- Russell, S., & Norvig, P. (2009). *Artificial intelligence: A modern approach* (3rd ed.). Prentice Hall Press.
- Salgado, C. M., Azevedo, C., Proença, H., & Vieira, S. M. (2016). Noise versus outliers. In *Secondary Analysis of Electronic Health Records*, 163–183.
- Schallmo, D. R. A., & Williams, C. A. (2018). *Digital transformation now!* Springer International Publishing.
- Schallmo, D., Williams, C. A., & Boardman, L. (2020). Digital transformation of business models—best practice, enablers, and roadmap. *Digital Disruptive Innovation*, 119-138.
- Taplin, I. M. (2014). Global commodity chains and fast fashion: How the apparel industry continues to re-invent itself. *Competition & Change*, *18*(3), 246–264.
- Thamm, A., Gramlich, M., & Borek, A. (2020). *The ultimate data and AI guide: 150 FAQs about artificial intelligence, machine learning and data*. Data AI Press.
- Thomassey, S. (2010). Sales forecasts in clothing industry: The key success factor of the supply chain management. *International Journal of Production Economics*, *128*(2), 470–483.
- Thorpe, R., & Holt, R. (2007). In *The SAGE dictionary of qualitative management research*. SAGE Publications Ltd.
- Thrall, J. H., Li, X., Li, Q., Cruz, C., Do, S., Dreyer, K., & Brink, J. (2018). Artificial intelligence and machine learning in radiology: Opportunities, challenges, pitfalls, and

- criteria for success. *Journal of the American College of Radiology*, 15(3), 504–508.
- Vaagen, H., & Wallace, S. W. (2008). Product variety arising from hedging in the fashion supply chains. *International Journal of Production Economics*, 114(2), 431–455.
- Viertl, R. (2011). Statistical methods for non-precise data. In *International encyclopedia of statistical science*, 1442–1444.
- Walters, D. (2002). *Operations strategy: A value chain approach*. Macmillan International Higher Education.
- Wang, B., & Ha-Brookshire, J. E. (2018). Exploration of digital competency requirements within the fashion supply chain with an anticipation of industry 4.0. *International Journal of Fashion Design, Technology and Education*, 11(3), 333–342.
- Williams, Graham. "Descriptive and predictive analytics." In *Data Mining with Rattle and R*, pp. 171-177. Springer, New York, NY, 2011.
- Wong, W. K., & Guo, Z. X. (2010). A hybrid intelligent model for medium-term sales forecasting in fashion retail supply chains using extreme learning machine and harmony search algorithm. *International Journal of Production Economics*, 128(2), 614–624.
- Xia, M., & Wong, W. K. (2014). A seasonal discrete grey forecasting model for fashion retailing. *Knowledge-Based Systems*, 57, 119–126.
- Xiong, G., Zhu, F., Liu, X., Dong, X., Huang, W., Chen, S., & Zhao, K. (2015). Cyber-Physical-Social system in intelligent transportation. *IEEE/CAA Journal of Automatica Sinica*, 2(3), 320–333.
- Xu, X., & Hua, Q. (2017). Industrial big data analysis in smart factory: Current status and research strategies. *IEEE Access*, 5, 17543–17551.
- Yu, Y., Choi, T.-M., & Hui, C.-L. (2011). An intelligent fast sales forecasting model for fashion products. *Expert Systems with Applications*, 38(6), 7373–7379.
- Zadeh, L. A. (1965). Fuzzy sets. *Information and Control*, 8(3), 338–353.
- Zadeh, L. A. (1975). The concept of a linguistic variable and its application to approximate reasoning—I. *Information Sciences*, 8(3), 199–249.
- Zou, X., Wong, W. K., & Mo, D. (2019). Fashion meets AI technology. In *Advances in intelligent systems and computing* (pp. 255–267). Springer Verlag.
- 
-



## **Appended Papers**



## Paper I

**Giri, C., Jain, S., Zeng, X., & Bruniaux, P.** (2019). A detailed review of artificial intelligence applied in the fashion and apparel industry. *IEEE Access*, 7, 95376-95396.

Available at: <https://ieeexplore.ieee.org/document/8763948>

DOI: [10.1109/ACCESS.2019.2928979](https://doi.org/10.1109/ACCESS.2019.2928979)



# A Detailed Review of Artificial Intelligence Applied in the Fashion and Apparel Industry

CHANDADEVI GIRI<sup>1,2,3,4</sup>, SHEENAM JAIN<sup>1,2,3,4</sup>, XIANYI ZENG<sup>1,4</sup>, AND PASCAL BRUNIAUX<sup>1,4</sup>

<sup>1</sup>Laboratoire de Génie et Matériaux Textiles (GEMTEX), ENSAIT, F-59000 Lille, France

<sup>2</sup>The Swedish School of Textiles, University of Borås, S-50190 Borås, Sweden

<sup>3</sup>College of Textile and Clothing Engineering, Soochow University, Suzhou 215168, China

<sup>4</sup>Automatique, Génie informatique, Traitement du Signal et des Images, Université Lille Nord de France, F-59000 Lille, France

Corresponding authors: Chandadevi Giri (chanda.giri2@gmail.com) and Sheenam Jain (sheenam.jain21@gmail.com), both contributed to this work equally.

This work was supported by the European Commission through the Framework of Erasmus Mundus Joint Doctorate Program-SMDTex.

**ABSTRACT** The enormous impact of artificial intelligence has been realized in transforming the fashion and apparel industry in the past decades. However, the research in this domain is scattered and mainly focuses on one of the stages of the supply chain. Due to this, it is difficult to comprehend the work conducted in the distinct domain of the fashion and apparel industry. Therefore, this paper aims to study the impact and the significance of artificial intelligence in the fashion and apparel industry in the last decades throughout the supply chain. Following this objective, we performed a systematic literature review of research articles (journal and conference) associated with artificial intelligence in the fashion and apparel industry. Articles were retrieved from two popular databases “Scopus” and “Web of Science” and the article screening was completed in five phases resulting in 149 articles. This was followed by article categorization which was grounded on the proposed taxonomy and was completed in two steps. First, the research articles were categorized according to the artificial intelligence methods applied such as machine learning, expert systems, decision support system, optimization, and image recognition and computer vision. Second, the articles were categorized based on supply chain stages targeted such as design, fabric production, apparel production, and distribution. In addition, the supply chain stages were further classified based on business-to-business (B2B) and business-to-consumer (B2C) to give a broader outlook of the industry. As a result of the categorizations, research gaps were identified in the applications of AI techniques, at the supply chain stages and from a business (B2B/B2C) perspective. Based on these gaps, the future prospects of the AI in this domain are discussed. These can benefit the researchers in academics and industrial practitioners working in the domain of the fashion and apparel industry.

**INDEX TERMS** Artificial intelligence, big data analytics, machine learning, expert systems, fashion and apparel industry.

## I. INTRODUCTION

Fashion and apparel (F&A) industry is one of the largest economies contributing 38% to the Asia Pacific, 26% to Europe and 22% to North America [1]. According to Business of Fashion, (2019), F&A sales are projected to grow by 7.5% and 5.5% in the Asia Pacific and Europe respectively. F&A is also one of the largest waste producers globally [3] because of problems like overproduction and product returns.

The associate editor coordinating the review of this manuscript and approving it for publication was Zijian Zhang.

The principal reason behind this is the consumer’s dissatisfaction with the products offered by the industry in terms of size, color, and style. Hence, it is essential for the industry to become customer-centric for successfully regulating environment-friendly manufacturing practices. Consequently, it is important that the industry adopt sustainable production practices to alleviate waste production and management. One of the ways of achieving this can be by taking advantage of emerging Artificial Intelligence (AI) techniques for creating a sustainable digital supply chain [4].

In the past decades, AI has transformed many industries like health, transportation, and manufacturing due to its capability to solve problems using conventional mathematical models [5]–[7]. The application of AI has been recognized in the F&A industry at various stages such as apparel design, pattern making, forecasting sales production, supply chain management [8], [9].

With the emergence of globalization and digitalization, AI has gained attention to connect businesses globally. In the last decade, the F&A industry has utilized AI to a certain extent for improving supply chain processes like apparel production [10], fabric inspection [11], distribution [12]. This was important as the F&A industry is volatile and it is always challenging to quickly respond to change in trends and continuously evolving consumer's demands.

An additional impact of digitalization is noticed in consumer behavior in the F&A industry. The increase in awareness and advent of new offline and online mediums has changed the contemporary consumer's decision-making pattern, influenced by the various online and offline mediums [13]. It is, therefore, important to create digital platforms for efficient requirements elicitation and collection. This can be attained by utilizing the benefits accompanied by Information technology (IT), Artificial intelligence (AI) techniques, big data analytical tools and other current technologies [14].

Evidently, the F&A industry is one of the most dynamic industries with new data being generated every time a new garment is designed, produced and sold [15]. However, the industry still lacks the extensive adoption of AI methods. The industry is still using computational tools based on classical algorithms and modern AI techniques are confined to academic research. Hence, it is a requisite for the industry to adopt new AI techniques to have a competitive advantage and improve business profitability. To do this, it is indispensable to have a consolidated description of different AI techniques used in research to target various business problems in the F&A supply chain.

After scrutinizing the extant literature in this domain, we encountered a few review articles, where the focus was on either AI or supply chain in F&A. For instance, the review conducted in [8] shows categorization of research articles on the basis of four operation processes in the apparel industry: apparel design, manufacturing, retailing, and supply chain management. This study presented the limitations of academic research that hinders the application of AI methods at an industrial level and also found that the F&A industry received less recognition from AI research groups. The work represented in [16] was restricted to AI algorithms, "Decision support systems" and "Intelligent systems" in the textile and apparel supply chain. In addition, this study only considered journal articles for the review. In contrast to these two reviews, the review carried out in [17] focuses on "Data mining and Machine learning models" implemented in the textile industry. According to this study, classification techniques were applied more frequently as compared to clustering techniques.

Despite valuable contributions to the previous literature reviews, when observed, none of the reviews studied the overall impact of AI in the F&A industry. In addition, there is a need to have a broader outlook of AI techniques employed for improving business operations at different supply chain stages. Furthermore, no study emphasized on defining the F&A supply chain stages according to the business perspective. Every business is composed of Business-to-Business (B2B) and Business-to-Consumer (B2C) transactions. In a traditional business setting, every personnel involved in business operations has knowledge confined to a specific domain. However, with the proliferation of AI technology, the complexity of business operations has risen, making it important for individuals (industrial researchers, academic researchers, managers) to have interdisciplinary knowledge.

By segregating the supply chain operations into B2B and B2C, the purpose is to provide a roadmap for individuals who are willing to expand the horizon of their expertise to help the F&A industry in improving their business models and profitability.

The objectives of this study are threefold. First, to do an in-depth analysis of the ongoing trend of AI in the F&A industry over the last decades. For this, no time constraint was introduced while retrieving the articles from scientific databases. Second, to understand the exploitation of AI techniques employed at various F&A chain stages. This is to examine the industrial transformation from a technical perspective. Third, to comprehend the utilization of AI techniques with a business perspective in F&A supply chain. Hence, this paper addresses the following research questions:

RQ1. What is the impact of Artificial Intelligence on F&A Industry over the past decades?

RQ2. Where have the AI methods been applied in F&A supply chain?

RQ3. To what extent has research addressed the supply chain problems from a B2B and/or B2C perspective?

In this direction, this research aims to conduct a systematic and comprehensive literature review of AI methods applied in the F&A industry in the past decades. This study is viable for an independent researcher to understand AI trend in F&A irrespective of their domain. Another important attribute of this work is the consideration of all journal and conference publications, which is rarely found in other review studies.

The remaining article is organized as follows: Section II outlines the research framework for conducting the systematic literature review. Section III describes the steps involved in the article screening process. Section IV represents the taxonomy proposed for the classification of AI methods and F&A supply chain. This is followed by section V that discusses the analysis and findings of the review process. Section VI and VII present the research gaps identified, future implications, conclusion, and limitations.

## II. RESEARCH FRAMEWORK

In an attempt to answer the research questions, this study presents a systematic literature review (SLR) focusing on

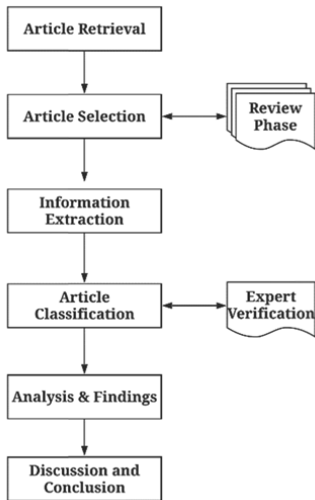


FIGURE 1. Systematic literature review: Research framework.

TABLE 1. Competencies of researchers.

Researcher	Competencies	
	Major	Minor
First Researcher	Artificial Intelligence, Data Science, Expert Systems, Machine Learning	Fashion, Textile, Supply Chain, Management
Second Researcher	Fashion and Apparel Supply Chain, Fashion Technology, Information Technology, Data analysis	Machine Learning, Artificial Intelligence
Expert Researcher	Significant Knowledge of Both Domains (AI and F&A)	N/A

artificial intelligence methods applied in the F&A industry. An SLR methodology was chosen to make the research more rational, transparent and reproducible [18].

Based on the research focus, the methodology adopted is shown in Figure 1. The review process commenced with collecting and preparing data from scientific databases. Subsequently, articles were selected in five phases (depicted in Figure 2), strictly adhering to the inclusion and exclusion criteria defined in table 4 and 5. Finally, the selected articles were considered for classification (described in section IV) and further analysis complying with the research questions. There were two researchers involved in the entire review process and one expert researcher for the validation of the classification process. The competencies of each researcher can be seen in the following Table 1.

### III. ARTICLE SCREENING PROCESS

The article screening process is presented in Figure 2. It is comprised of three steps, namely article retrieval, article selection, and information extraction.

TABLE 2. Synonyms of the targeted search keywords.

Artificial Intelligence (AI)	Fashion and Apparel (F&A)
Machine Learning	Fashion Industry
Deep Learning	Garment Industry
Data Mining	Apparel Industry
Artificial Intelligence	Clothing Industry
Data Analytics	Textile Industry
Expert Systems	
Knowledge Systems	
Intelligent Systems	
Decision Support Systems	
Data Management	

#### A. ARTICLE RETRIEVAL

This section discusses the steps involved in article retrieval, which is the initial part of the article selection process. The first step was to choose the databases to conduct the SLR. Two popular scientific databases, Scopus and Web of Science, were selected because of their popularity in academia. In addition, these databases index most of the journals and conference proceedings. Especially, most of the work in this research’s domain is also indexed in these two databases [19], [20].

This was followed by formulating the search string, which included all the synonyms related to artificial intelligence and the F&A industry (shown in Table 2). The final search string defined for both the databases are as follows:

##### 1) SEARCH STRING FOR SCOPUS

TITLE-ABS-KEY ( (“Machine learning” OR “deep learning” OR “data mining” OR “artificial intelligence” OR “data analytics” OR “expert system” OR “knowledge system” OR “intelligent system” OR “decision support system”) AND ( ( fashion OR garment\* OR apparel\* OR cloth\* OR textile\* ) industry\* ) ) AND ( LIMIT-TO (LANGUAGE, “ENGLISH” ) ) )

##### 2) SEARCH STRING FOR WEB OF SCIENCE

TS = (((“Machine learning” OR “deep learning” OR “data mining” OR “artificial intelligence” OR “data analytics” OR “expert system” OR “knowledge system” OR “intelligent system” OR “decision support system”) AND ( ( fashion OR garment\* OR apparel\* OR cloth\* OR textile\* ) industr\* ) )

Refined By: LANGUAGES: (ENGLISH)

where,

TITLE-ABS-KEY/ TS = Title, Abstract, and Keywords  
AND/ OR = Boolean operators to connect different keywords

\* = used for loose/approximate phrase

“ = used for exact phrase

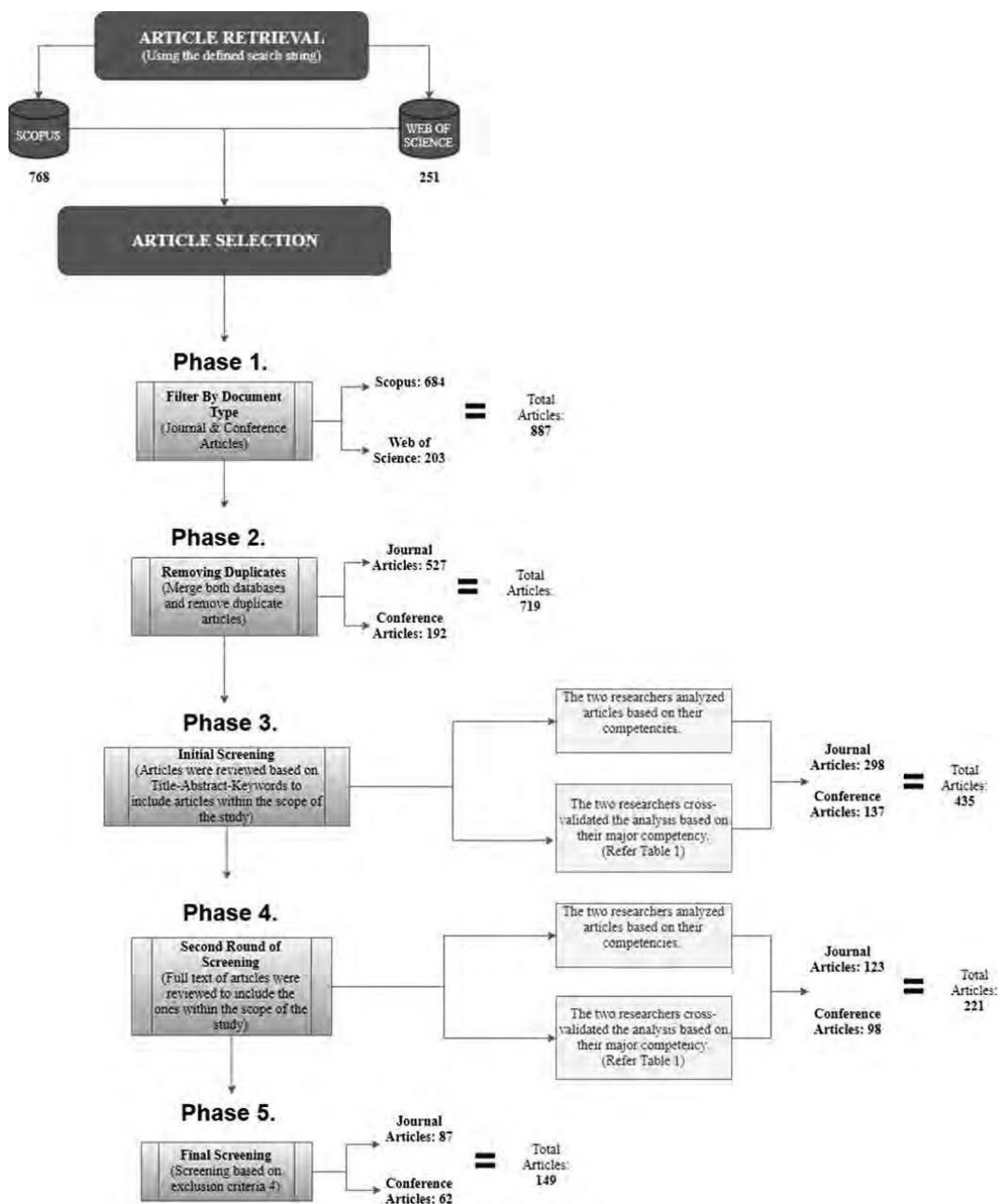


FIGURE 2. Article screening process.



TABLE 3. Extracted documents.

Scopus		Web of Science	
Document Type	Number of Articles	Document Type	Number of Articles
Conference Paper	319	Articles	138
Article	352	Conference Paper	65
Conference Review		Book and Book Chapter	48
	44		
Book and Book Chapter	32		
Review	13		
Article in Press	6		
Editorial	1		
Note	1		
Total	768		251

TABLE 4. Inclusion criteria.

Number	Criteria	Reason for Inclusion
1	No time Constraint	To understand the overall trend of AI in the F&A domain to answer RQ1
2	All Journal Articles and Conference Proceedings indexed in Scopus and Web of Science Databases	These databases index most of the journals and conference proceedings in the research field; To make sure that all articles relevant for addressing our three RQs are fetched
3	All countries and markets	To prevent biases while investigating the literature
4	All researches that applied AI techniques in F&A domain	To recognize all the stages in F&A supply chain where the implementation and execution of AI has been realized (to answer RQ2)
5	All researches conducted with a perspective of B2B, B2C or Both in F&A domain	To examine the extent with which business problems have been acknowledged using AI (to answer RQ3)

TABLE 5. Exclusion criteria.

Number	Criteria	Reason for Exclusion
1	Grey literature	To maintain the scientific reliability of the literature review
2	Non-English articles	To eliminate the misapprehension while scrutinizing the articles and avoid language barrier
3	Industries that are not F&A industry	As the study concentrated on Fashion and Apparel industry
4	Research discussing theoretical and/or conceptual frameworks	To ensure the empirical validation of the AI models that can be applied in F&A industry

The execution of these search strings on Scopus and Web of Science yielded 768 and 251 articles (total articles 1019) respectively. The different document types are shown in Table 3. In Scopus, the research articles were found from a time-period of 33 years (1989-2018). Whereas on Web of Science, the time period was of 18 years (1991-2017). It should be noted that no time constraint was applied while searching for articles as the aim was to study all the work done in the research domain, fulfilling the goal of RQ1. The article selection process was carried out using certain inclusion and exclusion criteria enumerated in Table 4 and 5.

**B. ARTICLE SELECTION**

This section describes the rigorous screening process employed by the two researchers involved in order to select the articles relevant to address the research questions. The screening included five phases as shown in Figure 2. In the ‘Phase 1’, the articles were filtered by document type in accordance with inclusion criteria 2 and exclusion criteria 1,

resulting into 684 articles in Scopus and 203 articles in Web of Science (total articles 887). In the ‘Phase 2’, the articles from both data sets were merged into one and redundant articles were eliminated, reducing the articles to 527 from journals and 192 from conference proceedings (total articles 719). In the ‘Phase 3’, initial screening was carried out by analyzing the “Title-Abstract-Keywords”, conforming to inclusion criteria 3 & 4 and exclusion criteria 3. The initial screening was conducted in two sub-phases. First, the two researchers analyzed articles according to their competencies. Second, the two researchers cross-validated the analysis based on their major competencies (refer to Table 1). At this stage, the number of articles decreased to 298 from journals and 137 from conference proceedings (total articles 435).

Similarly, considering the same inclusion and exclusion criteria in the ‘Phase 4’, the two researchers first studied and analyzed the “Full text” of the articles, and then cross-validated the analysis based on their major competencies (refer Table 1). While accessing the full texts, a few conference articles were encountered having published only abstracts. Such abstracts were excluded from the study. At this point, the number of remaining articles were 123 from journals and 98 from conference proceedings (total articles 221).

Lastly, in the ‘Phase 5’, the articles were scanned based on the exclusion criteria 4. The rationale was to include studies where the conceptual AI model was implemented and empirically validated. The final count of the articles was 87 from journals and 62 from conference proceedings (total articles 149). In all the phases, the articles were excluded based on the consensus between the two researchers.

The final 149 articles were considered for the classification based on supply chain stages, applied artificial intelligence techniques, and business perspective: B2B and B2C. To accomplish this, different stages in F&A supply chain and classes in AI (explained in detail in section IV) were defined. Further, the F&A supply chain stages were categorized into Business-to-Business (B2B) and Business-to-Customer (B2C). This classification was important to get a clear outlook of the different AI classes applied at the F&A supply chain stages to address the research questions as this would help to identify opportunities with AI to accomplish business-related problems in F&A industry.

**C. INFORMATION EXTRACTION**

This section discusses the process followed for extracting information and classifying the selected articles based on supply chain stages, artificial intelligence classes, B2B and B2C to address our research questions RQ1, RQ2, and RQ3. The 149 articles were thoroughly examined to extract the following information:

- 1) *Applied AI class and algorithm*
- 2) *Supply chain stage under study*
- 3) *Business perspective: B2B and B2C*
- 4) *Research gaps Identified*

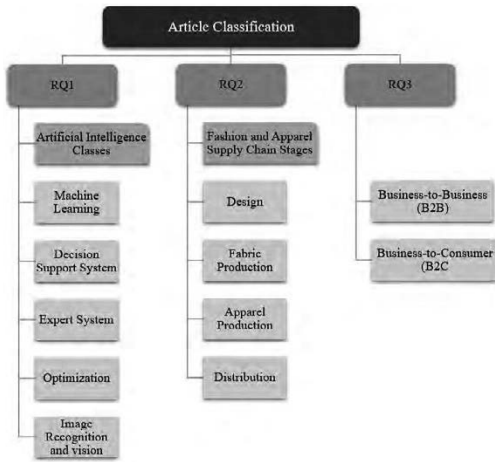


FIGURE 3. Article classification based on research questions.

The article classification conforming to the research questions is represented in Figure 3. As it can be seen, RQ1 is focused on understanding the overall trend of AI in the F&A industry. Hence, the focus of the screening process was limited to those articles discussing the implementation and execution of AI techniques in the F&A industry. To acknowledge RQ1, AI techniques were divided into five categories: Machine Learning, Decision Support System, Expert System, Optimization, and Image Recognition & Vision. The algorithms considered under each class are discussed in section IV.B. While extracting information, these classes were assigned to the articles.

RQ2 is aimed at identifying the various stages in the supply chain at which the AI method was employed. Hence, during the information extraction stage, the supply chain stage under study was recorded. To acknowledge RQ2, the supply chain stages were classified as Design, Fabric Production, Apparel Production, and Distribution. The processes considered under these stages are shown in Figure 4. While extracting information, the articles were assigned these supply chain stages.

RQ3 aims to understand the extent of business problems being a focus of research studies. To do this, the supply chain stages identified were further categorized from a business perspective into B2B and B2C (discussed in detail in section IV.A and Table 6). These classes were allocated to the research articles during information extraction.

This classification of research articles was verified with the help of an expert researcher actively involved in research related to artificial intelligence and F&A industry from the past two decades. The competency of the expert researcher is also mentioned in Table 1.

TABLE 6. B2B and B2C activities in the F&A industry.

B2B	B2C
Fashion Design	Fashion Design
Textile Design	Textile Design
Spinning	Dyeing & Printing
Weaving or Knitting	Cutting
Dyeing, Printing, Finishing & Inspection	Sewing & Assembly
Cutting	Finished Garment
Sewing & Assembly	Retailing
Finished Textile	E-commerce
Wholesaling	
Retailing	

IV. ADOPTED STRUCTURE FOR CLASSIFICATION OF ARTICLES

This section elucidates the structure of the fashion and apparel supply chain and attempts to cluster different supply chain stages into B2B and B2C. This is discussed in sub-section IV.A, which proposes a taxonomy to address RQ 2 and RQ 3 respectively. Similarly, AI techniques were assembled into five classes as explained in the sub-section IV.B to propose a taxonomy to address RQ 1.

A. PROPOSED TAXONOMY OF FASHION & APPAREL SUPPLY CHAIN STAGES

The fashion and apparel supply chain is a complex network of various actors designated worldwide. It deals with a diversity of raw materials: fiber, yarn, fabric, dyestuff, and other chemicals, and the related processes are broadly classified into four stages: design, fabric production, apparel production, and distribution as shown in figure 4. Traditionally, the supply chain follows a push system [21], where the brand owners or retailers (buyer) provide the manufacturers with information like the design or technical specification of the fabric and garment to be produced, the volume of the products, sizes in which the garment is to be produced. The fabric and garment producers follow the instructions to create samples, which upon approval by the buyer are converted into finished fabric and garment respectively. Usually, the finished fabrics are the raw material for the apparel production process. Finally, the finished garments are transported to a wholesaler or retailer. In the case of the wholesaler, there is another actor, which acts as a distributor between the consumer and the wholesaler. On the other hand, in the case of the retailer, the garments are sold through one or more channels, for instance, brick & mortar stores, web-shops (e-commerce), departmental stores, multi-brand retailers.

The designers employed by retailers are responsible for creating collections based on the current market and trend analysis. In most scenarios, retailers do not own any production house and play an important role to bring the products into the market. Hence, in the conventional supply chain, all the actors from design up to retailers/brand owners are

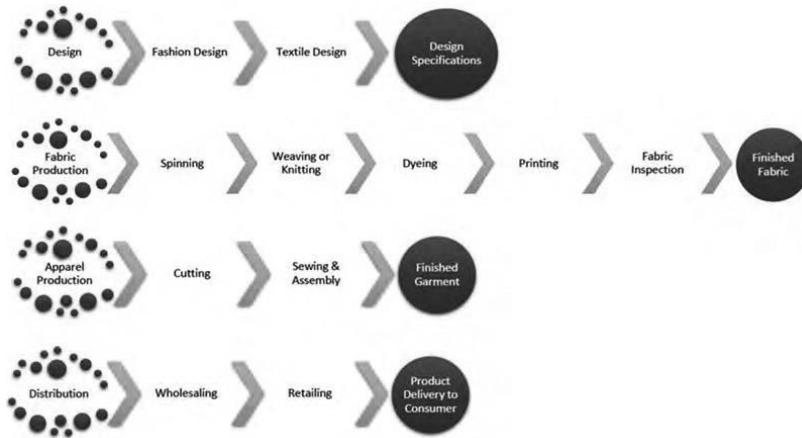


FIGURE 4. Stages in F&A supply chain.

considered as Business to business (B2B) as their primary customers are other businesses, while retailers are considered as business to consumer (B2C) as their primary customers are the end-users or consumers. However, in the past decade, with the advent of e-commerce, the definition of B2B and B2C has evolved [22]. Therefore, it has become important for the industry to adapt to this change and create new business strategies. It has also become vital to give a comprehensive demarcation between B2B and B2C, and how AI can help in combating problems at these segments.

#### 1) B2B (BUSINESS-TO-BUSINESS)

The F&A industry has a convoluted supply chain due to diverse product categories and their short lifecycle. The contemporary consumer has increased awareness and information related to the latest styles and designs [23]. Therefore, consumer buying behavior and engagement has changed. This is highly influenced by the proliferation of social media and internet [24], which is a widespread medium for dissemination of information related to the latest fashion trends, upcoming fashion weeks and popular celebrities. Due to this, the F&A supply chain has to rapidly change the collections to fulfill the growing consumer demand [25]. Hence, it is driven by a combination of business-to-business (B2B) and business-to-consumer (B2C) transactions. In any business, the number of B2B transactions is higher in comparison to that of B2C transactions [26], [27]. The main reason for this is that for every product there can be as many B2B transactions as there are sub-components or raw materials involved, while there will be only one B2C transaction.

Business to business (B2B) in the F&A industry is referred to as the commerce between two or more businesses. A B2B transaction, thus, will occur when a business demands raw material for the production process to manufacture the

product (e.g. garment manufacturer buying yarn), needs services for operational reasons (e.g. employing a third-party logistics service provider), re-sells goods and services produced by other businesses (e.g. a retailer buying products from manufacturer). The goal of a B2B transaction is to help their business stay profitable, competitive and successful. Table 2 shows the classification of the supply chain into B2B.

#### 2) B2C (BUSINESS-TO-CONSUMER)

B2C refers to the transactions conducted directly between a business and consumers who are the end-users of its product and/or services [28]. Behind a B2C transaction is a well-researched consumer regarding their options in order to find the best price and quality tradeoff. Traditionally, B2C referred to outlet shopping, however, with the rise of internet, smartphones and other mobile technology, a set of completely new B2C business channels have developed in the form of e-commerce, m-commerce, social media commerce or selling products and services over the internet [29]. It has become important for B2C companies to be omnipresent because of uncertainty in consumer behavior over different retail channels [30]. The success of a B2C model depends on the capacity of a business to evolve based on new technologies that are widely used by the consumers. Businesses that rely on B2C sales must maintain good relations with consumers to ensure their retention and loyalty.

In addition to this, another set of B2C transactions have evolved with the adoption of mass customization (MC) to fulfill growing consumer needs. In this, consumer interactions with the business increases and can even occur at the product development stage. Hence, fashion or textile design, dyeing, printing, cutting, sewing & assembly can all have customer involvement. Table 2 shows the classification of the supply chain into B2B and B2C.

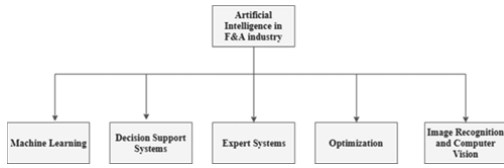


FIGURE 5. Classification of AI in the F&A industry.

## B. PROPOSED TAXONOMY OF APPLIED AI METHODS IN F&A SUPPLY CHAIN

Artificial intelligence has already proved its capability to solve the real world problems due to its heuristic characteristics of generalizing data. In the last three decades, the F&A industry has undergone a number of changes and AI has played a key role in this transformation. Currently, the F&A industry is equipped with advanced machines required at the various stages of apparel production, which has improved the overall efficiency of the industrial processes [9]. Application of AI is well explained and categorized by operating processes at a managerial level in the F&A [8]. However, these researches lack the categorization of the applied AI in the F&A supply chain. The study in [8] explains the research issues in the operating process of the apparel industry. This study found that AI research has 45% contribution to apparel manufacturing issues, approximately 9% to apparel forecasting and 4.2% to fashion recommendation.

Research in [17] is a comprehensive review of classification and clustering techniques utilized in the F&A industry and shows that classification algorithms have been used more than clustering. On the contrary, this research work does not talk about linear and non-linear predictive models. Moreover, this research does not convey about the current applications of computer vision and deep learning [31], customer analytics [32], optimization techniques [33], and big data analytics for digital manufacturing and customization [15], [34]. Taking into account the recent research conducted in the area of AI in F&A, this study categorizes the AI into five broad classes as shown in Figure 5.

### 1) MACHINE LEARNING

Machine Learning is a technical process by which the computers are trained to perform the assigned task without human intervention and learns from the patterns of the data itself. Mathematical models are built on historical data to predict and find hidden patterns to make a future decision [35]. Machine Learning can be classified as Supervised or Unsupervised learning.

**Supervised Learning-** is a parametric model and it has input (independent variables) and target variable (dependent variable) [36]. Supervised model performance can be improved by optimizing the model parameters through iterative processes [37]. Based on the research problem, it could be a classification or regression task and this relies on dependent variable whether it is categorical or numerical.

**Unsupervised Learning-** models have only input attributes or independent variables with the main task of grouping similar data points. This grouping the similar pattern data points is called clustering and this process creates their own labels [35].

Machine learning has been implemented in the F&A industry for sales prediction [38], trend analysis, color prediction [39], demand forecasting [40], fabric defect detection [41], predicting fabric behavior using mechanical properties [42].

### 2) DECISION SUPPORT SYSTEMS

The decision support system (DSS) is used in an organization at the commercial level for taking mid-level or high-level managerial decisions. It can be automatized or regulated by a human or blend of both. Few authors considered the decision support systems as a software tool whereas others considered it as a system that can be integrated with the business to make intelligent decisions [43]. Research in [44] states that DSS combines the mathematical model with conventional data retrieval methods; it is flexible and adapts to the organizational environment as per defined strategy. In the F&A industry, it is widely used to industrialize innumerable tasks by optimizing decision making process in the supply chain [45]. Decision support systems help the various actors in apparel manufacturing and production to choose appropriate process and resources to decrease the overall cost and enhance the performance of the apparel supply chain [46].

### 3) EXPERT SYSTEMS

In artificial intelligence, 'Expert system' is a system that makes a decision without human intervention [47]. It uses a reasoning approach to solve the complex problem, characterized by "if-then" rules. The first system was found around the 1970s and then gained popularity by 1980s [48]. They were considered the first popular software in the field of AI [36]. Expert systems are classified as Inference engine and knowledge base. 'Knowledge base' works on the principle of facts and rules, while 'inference engine' uses the rules to learn the facts and derive new facts [49]. In the F&A industry, it is applied in apparel manufacturing and production to select appropriate processes and equipment in order to generate minimal environmental pollution [50]. Furthermore, it has been applied for creating a recommendation engine in fashion retailing to improve the overall satisfaction of customers [51].

### 4) OPTIMIZATION

Artificial intelligence has the ability to solve complex problems and find numerous solutions by intelligent searching [36], [52]. Classical search algorithm starts with some random guess and this is improved using the iterative process. 'Hill climbing', 'Beam search' and 'Random optimization' are some of these methods [53]. Machine learning algorithms use 'search algorithms', which are based on optimization techniques. Simple exhaustive searches [52], [54]

are too slow and therefore ‘Heuristics’ approach is adapted to serve as a technique to find a solution. The limitation of the heuristics search approach is that it fails to work with smaller datasets [55]. An evolutionary algorithm is another form of optimization search, which starts with the initial guesses of the population permitting them to mutate, recombine and select the best one while discarding others. Popular Evolutionary algorithms are genetic algorithms (GA), gene expression programming and genetic programming [56], [57]. Distributed search method could be done using ‘swarm intelligent’ algorithms. GA is extensively used in the F&A industry to overcome the problems of scheduling and design layout of the apparel production [58], [59]. GA has the ability to respond to quick changes in the fashion industry. This algorithm has been used to improve the fitting services as well [60].

### 5) IMAGE RECOGNITION AND VISION

In Artificial intelligence, Computer vision is a scientific area, which trains a machine to achieve high-level interpretation of the images or videos. These images or videos can come from many sources such as the medical field, global sensing position, cameras [61], [62]. The principal tasks of computer vision algorithms are extraction, pre-processing, exploring the high dimensional data and creating supervised or unsupervised models [63]. Models use the concept of geometry, statistics, physics, and machine learning theory to get insights into the image understanding [64]. Object recognition, video tracking, motion estimation are some of the sub-areas in the field of computer vision [65]. Machine vision is applied in F&A to automate many industrial applications like inspection and process control [66]. Image recognition and vision is also popular for content-based image retrieval systems, virtual try-on and augmented reality in the F&A industry [67]–[69].

## V. ANALYSIS AND FINDINGS

According to our research framework, we selected 149 articles, which includes articles from both journals and conference proceedings. The articles were classified based on the taxonomy discussed in the previous section (refer section IV). This section discusses the result of our review addressing the three research questions in the form of distribution of articles. Section V.A and V.B correspond to RQ 1 by presenting the overall trend of AI in the F&A industry. Section V.C supports RQ 2 by showing articles that applied AI methods at various F&A supply chain stages. Section V.D correlates to RQ 3 by exhibiting articles by B2B and B2C. All the extracted information from the articles a) AI class, b) Methods, c) Supply chain stages and processes, d) B2B/B2C and the corresponding count of articles are consolidated in the form of Table 8 and 9 for journals and conferences respectively.

### A. THE OVERALL DISTRIBUTION OF ARTICLES OVER TIME

The overall trend of the articles published in three decades (1989 to 2018) is shown in figure 6. As can be observed, maximum research in the field of AI in F&A has been

Number of Research Papers

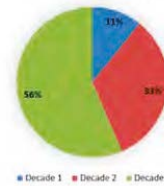


FIGURE 6. Decade-wise article representation of AI in F&A.

carried out in the last decade (2009–2018), which accounts for approximately 56% of the total articles reviewed. While in the first two decades, it was 11% and 33% respectively. Hence, even though AI methods were introduced long back in the 1950s [36] but their capability was realized much later in the last decade.

The detailed trend per year of articles published in journals and conference proceedings is depicted in Figure 7. As can be seen, the overall importance of this research domain has been equal in both journals and conferences. The only downward slope is visible for conference publications between the year 2013 and 2016 as highlighted in Figure 7. Apart from this, figure 8 shows the top 10 authors and institutions contributing to the literature of AI in the F&A industry.

### B. DISTRIBUTION OF ARTICLES BY APPLIED AI OVER TIME

In this section, a comprehensive result of the review in terms of the number of articles classified based on the AI classes defined in section IV.B is shown in figure 9. The maximum number of articles are published in the field of machine learning with a total contribution of 42%, followed by expert systems with 28%, and the least contribution of the other three AI classes.

The distribution of articles in journals and conferences since 1989 has been shown in figure 10. The AI class ‘Machine learning’ has been applied multiple times in journal articles since 1991. There are two peaks visible for journal articles in the year 2007 and 2017 with 4 and 7 articles published respectively. Whereas for conference articles, it has been applied since the year 2000 with three major peaks in the year 2010, 2012 and 2016 with 4, 3, and 7 articles respectively. On the other hand, the AI class ‘Expert system’ was widely used since 1994 in journal articles while in conference articles there has been no research after 2014. For AI class ‘decision support system’, there has been very little work in conference articles, while a gradual increase in presence is realized in journal articles since 2010. For ‘image recognition’, its presence is visible since 2009 in both journal and conference articles, being the least applied AI class (also shown in figure 10). In addition, it can be noticed that since the year 2017, its application has increased in journal articles. In contrast to other classes, for optimization, there were more articles in the conference as compared to journals.



FIGURE 7. Overall trend of AI in F&A since 1989.

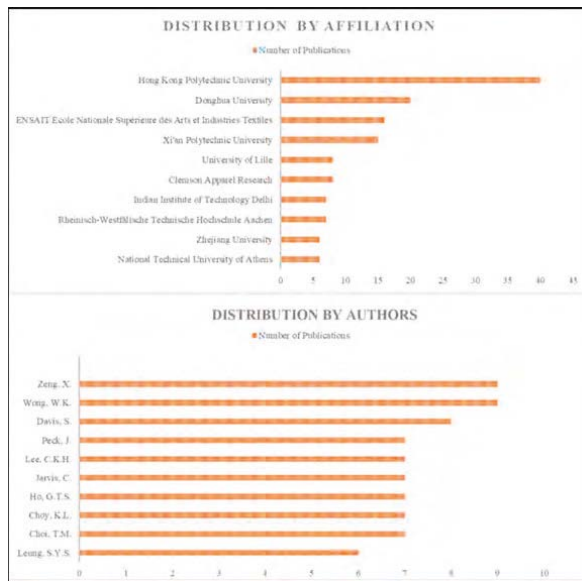


FIGURE 8. Top 10 affiliations and authors.

**C. DISTRIBUTION OF ARTICLES BY APPLICATION OF AI IN F&A SUPPLY CHAIN**

As shown in table 7, machine learning and expert systems have been applied widely at the four F&A supply chain stages, with least research in design. This is followed by optimization that has been applied with the least focus at the distribution stage.

The classification of articles by applied AI methods in F&A supply chain is shown in figure 11. In apparel production, all AI classes are applied in journal articles, while in conference articles image recognition and decision support system has been not used. In design, research has focused on three AI classes: optimization, machine learning, and expert systems for journal articles, whereas the expert system is not applied in conference articles. In distribution, majorly used

AI classes are machine learning, decision support systems, and expert systems in journal articles, while the focus has been on machine learning and image recognition in conference articles. In fabric production, all AI classes have been widely used with a major focus on machine learning and Expert systems. Additionally, there has been growing use of image recognition in fabric inspection, which is a process under fabric production (described in section IV.A).

**D. DISTRIBUTION OF ARTICLES BY B2B AND B2C OVER TIME**

As discussed in the previous sections the importance of B2B and B2C in F&A, we classified the articles based on their business focus. As can be seen in figure 12, research is focused on solving the issues related to B2B and there is

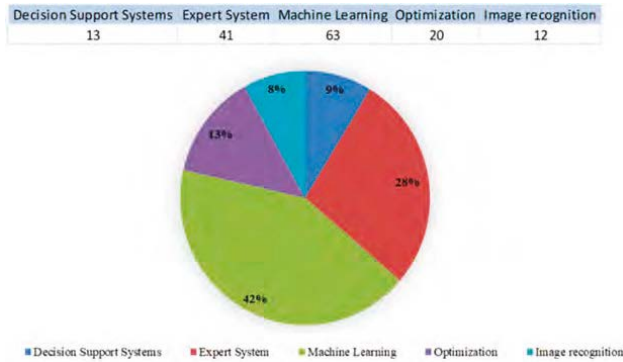


FIGURE 9. Total distribution of articles by AI methods applied.

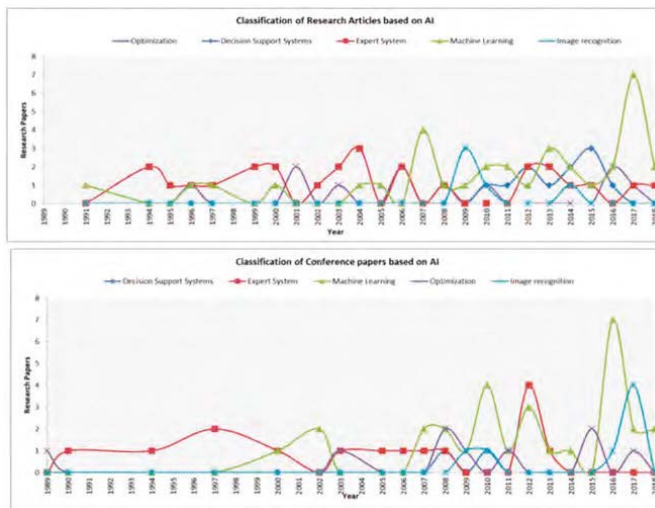


FIGURE 10. Distribution of articles by applied AI over time.

TABLE 7. Count of applied artificial articles in F&A supply chain.

	Decision support systems	Expert System	Machine Learning	Optimization	Image recognition	
Apparel Production	4	11	19	8	2	44
Design	0	4	2	3	0	9
Distribution	7	10	20	2	3	42
Fabric Production	2	16	22	7	7	54
	13	41	63	20	12	Total =149

little attention on B2C both in journal and conference articles. To get a clear picture, figure 13 shows the distribution of articles in three decades separately for B2B, B2C and both.

There were total 149 articles reviewed, out of which 122, 13 and 14 belonged to B2B, B2C and both (B2B/B2C)

respectively. If we consider B2B, as shown in figure 13, substantial research has been carried out in all three decades as compared to B2C. In the case of B2C, the total number of articles published in itself is low i.e. 13. Out of this, only two were published in the second decade and rest were

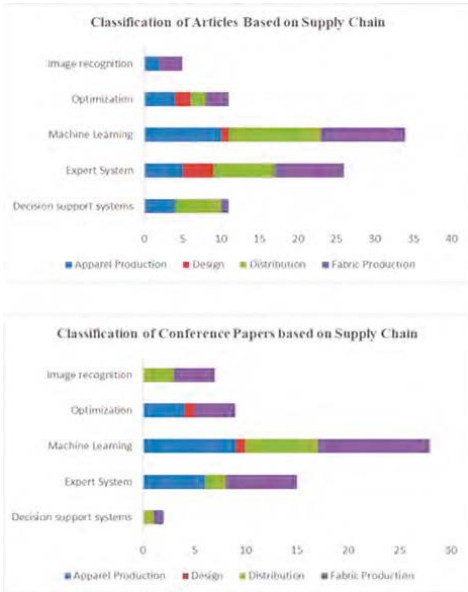


FIGURE 11. Distribution by supply chain processes.

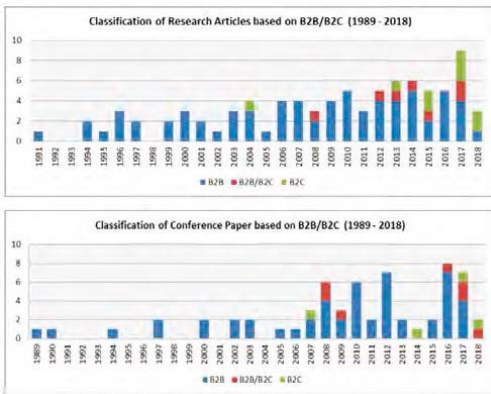


FIGURE 12. Distribution of articles by B2B and B2C.

published in decade 3. There were no articles published in decade 1. By both, it is meant that the focus of the study is both B2B and B2C. The total number of articles, in this case, is 14, out of which ten were published in the third decade and four in the second decade. There were no articles published in decade 1 for this as well. As a conclusion, we can say that little work has been done with the perspective of B2C and there is a need to attend to this gap.

## VI. DISCUSSION AND IMPLICATIONS

Based on the review conducted following the proposed research framework and addressing the three research questions, this study recognized a number of gaps with respect to applied AI in the F&A industry. These gaps in research provide the foundation to recommend future research direction in the F&A and AI domain.

### A. GAPS AND IMPLICATIONS

As discussed in section IV.A, most of the research related to AI in F&A industry has been carried out in the last decade (2009 to 2018), accounting to 56% of the total number of publications in the year 1989-2018. This shows that even though the AI methods existed since 1989, they have recently gained popularity in research related to F&A industrial problems. Although AI has left its footprint in research, it is still far from being implemented at the industrial level. One of the reasons for this is that researchers working in AI may lack expertise in F&A and at the same time, the professionals working in the industry may lack expertise in AI. In addition, industries are skeptical about the benefits of AI and big data analytics. Therefore, it is important that they look at the cost and benefit tradeoff to be able to exploit the full potential of AI.

In section V.B, we classified the articles based on their focus on B2B and B2C. The result demonstrates that most of the research work has been carried to target B2B business problems, accounting to 81% of the total number of publications. Whereas, little research has been conducted with a B2C perspective, accounting to approximately 8% of the total publications. This clearly illustrates that research needs to focus on B2C business problems. According to the State of Fashion (2019), two of the major industrial challenges faced will be rapidly changing consumer preferences and competition from online and omnichannel. Therefore, research related to the F&A industry needs to shift their focus to B2C, taking into account the importance of AI. AI can help to analyze consumer footprint omnichannel, which can help in creating personalized consumer database or profiles helping to improve business profitability and providing a competitive advantage.

In section V.C and V.D, we classified the articles based on applied AI method in the F&A supply chain. Most of the research articles fall under machine learning and expert systems class, which has been extensively applied at the supply chain stages: apparel production, fabric production, and distribution. Decision support systems, optimization and image recognition class have a more or less similar number of research articles published in this domain. Their application was seen to some extent at all three stages. However, least representation of these algorithms was observed at the design stage. It gives an impression that little focus has been given to design-related problems and hence, there is a huge scope of AI applications at this stage. For instance, AI methods can be used to create systems that can help fashion and product designers to capture consumer needs and preferences more





TABLE 8. AI methods used at various supply chain stages and processes in journal articles.

AI Class	Method/ Technique used	Department Targeted	Process Targeted	B2B/B2C	Journal Articles Count	Reference
<b>Machine Learning</b>	BP-ANN 9 back propagation artificial neural network; k-means clustering; sequential clustering; fuzzy logic; A two-level clustering method (SOM network+ K-means); Naïve Bayes; Support vector machine; Gene expression programming (GEP); FCM (fuzzy clustering using MSE); Non parametric regression forecasting; supervised clustering; K-Medoids ;CN2-SD; ANN regression; RFM modeling; Association rule; ELM (extreme learning machine); GA; fuzzy constraint logic system, fuzzy rules, fuzzy sets; feed-forward neural network, back-propagation algorithm; decision tree, classification, and regression tree, factor analysis; Regression; Treelerner; root mean square; Fuzzy Efficiency based Classifier System; Logistic regression; Bucket Brigade Algorithm; neural networks using the error back propagation mode, e-neuro fuzzy engine; DIDT technique (Top-Down Induction of Decision Trees ID3; data mining; text mining; semantic data analysis;	Apparel Manufacturing 10, Design 1, Distribution 14, Fabric Production 11	Cutting 2; Dyeing/ Printing/ Finishing/ Inspection 1; Finished Garments 2; Retailing 12; Sewing 4; Spinning 9;	B2B/B2C 4, B2B 25, B2C 5	34	[70] [38] [71] [72] [73] [74] [75] [76] [39] [77] [78] [79] [80] [81] [82] [83] [84] [85] [86] [41] [87] [87] [88] [89] [90] [91] [92] [93] [42] [94] [95] [96]
<b>Decision support system</b>	Fuzzy logic; fuzzy association rule mining (FARM); classification, regression, clustering and association analysis; Linear optimization with constraints; Fuzzy inference; Fuzzy aggregation; adaptive-network-based fuzzy inference system (ANFIS); analytic hierarchy process (AHP); TOPSIS;	Apparel Manufacturing 4, Distribution 6, Fabric Production 1,	Finished Garments 4; Retailing 6; Spinning 1;	B2B 9, B2C 2	11	[97] [98] [99] [100] [101] [12] [102] [103] [104] [105] [106]
<b>Expert System</b>	Association rules; ES named ES-EXITUS has been implemented using the SSM and the DMM; Fuzzy association rule mining; Fuzzy logic; clustering and probabilistic neural network (PNN); hybrid OLAP-association rule mining; Ontology, semantic web, multiple agents; Genetic Algorithm; gradient descent optimization, fuzzy sets; Chi-square test, Correspondence analysis; parametric cubic spline and bi-cubic surface patch, object-oriented technology for building the knowledge base; Linear programming, computer-based heuristic; Semantic network, heuristic rules; Bézier curve models evolutionary model; Sensitivity analysis, Cognitive mapping technique, cluster analysis; Normalization model; Programming language used Microsoft Visual C++ version 4.0, Rule-based expert system; Object-oriented representation technique; t-test, sensory evaluation;	Apparel Manufacturing 5, Design 4, Distribution 8, Fabric Production 9,	Cutting 1; Dyeing, Printing, Finishing, Inspection 5; Fashion Design 2; Finished Garments 1; Retailing 5; Sewing 2; Textile Design 2; Weaving or knitting 1; Wholesaling 1;	B2B/B2C 3, B2B 21, B2C 2	26	[107] [108] [109] [110] [111] [112] [113] [114] [115] [116] [117] [50] [118] [119] [120] [121] [122] [123] [124] [125] [126] [127] [128]
<b>Optimization</b>	Constraint and non-constraint optimization; simulation-based model; fuzzy rule optimization; Tabu-Bees algorithm; linear approximation; Evolutionary algorithms; Genetic algorithm; Morse function, topological analysis; Content-based filtering, wavelet decomposition using Haar transform collaborative filtering, vector correlation using the Pearson correlation coefficient; symbolic regression module; Multiple regression analysis, extrapolative forecasting and an adaptive Holt-Winters forecasting;	Apparel Manufacturing 4, Design 2, Distribution 2, Fabric Production 3,	Cutting 2; Fashion Design 1; Finished Garments; Retailing 2; Sewing 1; Spinning 1; Weaving or knitting 1; Wholesaling 1;	B2B 11	11	[129] [130] [131] [132] [46] [133] [134] [94] [135] [136]
<b>Vision</b>	ANN and image processing; K-means clustering, Naïve Bayesian, and a multi-layered perceptron (MLP); NN and GA; Back propagation neural network (NN);	Apparel Manufacturing 2, Fabric Production 3,	Finished Garments 2; Sewing 1; Spinning 1; Textile Design; Weaving or knitting 1;	B2B 5	5	[11] [137] [138] [139] [140]

TABLE 9. AI methods used at various supply chain stages and processes in conference articles.

AI Class	Methods/ Techniques used	Department Targeted	Process Targeted	B2B/B2C	Conference paper count	Reference
<b>Machine Learning</b>	SOA-based data mining framework, classification, ARIMA and KNN models, text mining, naive Bayes classifier, SOM neural network, EM cluster and ELM (extreme learning machine) same prediction, SVM and AdaBoost, artificial neural network, case-based reasoning, supervised learning, self-organizing maps, principal component analysis, type-2 fuzzy sets, clustering, correlation analysis, optimal bandwidth selection in kernel density, ontology, RDF, Multilayer Perceptron, J48 decision tree, k-nearest neighbor, classifier ripper, C4.5 and PART, neuro-fuzzy with subtractive clustering and genetic algorithm (ANFIS-GA) technique, Java, C/C++, Viswanathan-Bagchi algorithm, correlation, wavelet transform, neural network.	Fabric Production 11, Distribution 7, Apparel Manufacturing 9, Design 1,	Weaving or knitting 3, Retailing 8, Spinning 5, Dyeing, printing and finishing 2, Finished Garments 3, Sewing 3, Cutting 2, Textile Design 1, Fashion Design 1	B2B/B2C 2, B2B 23, B2C 3	<b>28</b>	[141] [142] [143] [144] [145] [146] [147] [148] [149] [150] [151] [152] [153] [154] [155] [156] [157] [158] [159] [160] [161] [162] [163] [164] [165] [166] [167] [168]
<b>Decision support system</b>	Self-adaptive genetic algorithm, genetic algorithm, top-down and bottom-up analysis, dynamic optimization algorithms, Maximum Principle of Pontryagin.	Distribution 1, Fabric Production 1	Finished Garments 1, Retailing 1	B2B 1, B2B/B2C 1	<b>2</b>	[169] [170]
<b>Expert System</b>	Rule-based, rough set theory, fuzzy, case-based reasoning, fuzzy logic, fuzzy logic sensory evaluation, Fuzzy neural network, Unsupervised learning, fuzzy clustering, genetic algorithm, approximate reasoning module, Rule-based System Shell, Metric-based Fuzzy Logic and artificial neural network, If-Then Rules for knowledge base, least-square regression analysis, linear regression, event series	Distribution 2, Fabric Production 7, Apparel manufacturing 6	Retailing 2, Spinning 3, Dyeing, Printing, Finishing and Inspection 2, Finished garment 1, Yam to Fabric 1, Cutting 1	B2B 14, B2B/B2C 1	<b>15</b>	[171] [172] [173] [174] [10] [175] [176] [177] [178] [179] [180] [181] [182] [183] [184]
<b>Optimization</b>	Stochastic descent, list algorithm, Evolutionary computing and genetic algorithm, Fuzzy set, geometric analysis method, Mirabit algorithm, apriori algorithm, Heuristic methods,	Fabric Production 4, Apparel manufacturing 4, Design 1	Spinning 3, Finished garments 1, Dyeing, printing and finishing 1, Cutting 3	B2B 8, B2B/B2C 1	<b>9</b>	[185] [59] [186] [187] [188] [189] [190] [191] [192]
<b>Vision</b>	Conditional Random Fields (CRF), Bayesian classification, CNN based classifier, computer vision, classification, consensus style centralizing auto-encoder (CSCAE), Gabor filter, Gaussian kernel, image processing using IMAQ, median filter, stereovision method	Fabric Production 4, Distribution 3	Weaving or knitting 1, Spinning 3, Retailing 3	B2B 5, B2C 1, B2B/B2C 1	<b>7</b>	[193] [194] [195] [196] [197] [198] [199]

and product designers. They can take its advantage in predictive analysis of future trends based on historical and real-time data. This can also prove promising in improving the existing recommendation engines, which currently rely on collaborative or content-based filtering. These engines can be improved by integrating with consumer data from social media and real-time trends from fashion blogs, magazines and other social networks like Pinterest, Instagram. Additionally, the performance of existing predictive models can be

enhanced with the help of advanced techniques like ensemble learning and transfer learning. An instance is the use of random forest instead of decision trees, the use of which has outperformed the classical model in terms of computational time [201], [202].

Another application area in F&A is mass customization, where machine learning can be used to reduce the lead times by creating a classifier that could be trained on the existing style database, enabling the product designers to

prepare the raw material inventory in advance. Pre-trained deep learning models using a library like Keras [203], inception model [204] along with big data analytics can be used to create co-design platforms with style recommendations helping consumers in co-designing garments.

One of the important application of image recognition and computer vision is to target the key consumer pain point of not having an appropriate size of the garment. This problem has been addressed with the help of virtual fitting tools. However, these tools are still at a nascent stage and can be highly improved with the help of these techniques.

As we have noticed that fuzzy techniques and genetic algorithms are exhaustively used for expert systems, decision support systems and optimization. These techniques can be combined with advanced AI techniques to enhance the computational ability of a machine learning algorithm. Similarly, if the classical forecasting model is combined with AI it can lead to better forecasting in terms of seasonality and trends.

It is evident from the review that the F&A industry still lacks an integrated platform for data sharing and communication amongst its stakeholders and consumers. An integrated platform has become a necessity to quickly respond to customer's growing needs and preferences and improving consumer satisfaction and loyalty. AI has a lot of potentials to create and maintain such a platform. This platform can help the industry to provide interactive communication, improved supply chain organization, hence, leading to a digitally connected supply chain. Another possibility is by merging AI techniques with blockchain technology, which can ensure security and transparency between consumers and various supply chain actors. The industry can be benefitted by integrating their business with cloud-based technologies like Microsoft Azure, Amazon web services, IBM Watson, etc., and parallel computing tools for big data analytics like Hadoop and Hive. If the fashion industry can successfully adopt the aforementioned AI techniques, it will be easier to integrate B2B and B2C leading to a sustainable business orientation of B2B2C. This will be fully achievable only when the F&A industry equip 'FAIR' data principle [205] in strategizing their business.

## VII. CONCLUSION

The aim of the study was to conduct a systematic literature review to address the three defined research questions (RQ1, RQ2, and RQ3). In line with the research framework, we retrieved 1019 articles published between 1989 and 2018, from two popular academic databases: Scopus and Web of Science. The article screening process was carried out in five phases (shown in figure 2), which resulted in 149 articles. To extract information from these articles and address our research questions, a taxonomy was proposed considering AI methods and F&A supply chain stages acknowledging RQ1 and RQ2 respectively. To acknowledge RQ3, F&A supply chain was further classified into B2B and B2C.

The research analysis says that most of the work in the field of F&A was carried out in the last decade (2009-2018) with

the most applied AI categories being "Machine Learning" and "Expert Systems". It was observed that the techniques most used in Machine Learning were predictive algorithms like regression and SVM, and whereas in the case of Expert Systems were "Artificial Neural Network", "Genetic Algorithm" and "Fuzzy Logic" for modeling F&A supply chain problems. Certainly, no application of algorithms like "deep learning" and "transfer learning" was realized. Further, very few research articles talked about "Big Data" in the field of F&A, which clearly states that the industry has not fully realized the potential of data analytics and AI.

This research found that F&A supply chain stages: "Apparel Production", "Fabric Production" and "Distribution" received maximum attention when applying AI techniques, whereas "Design" was the least focused. Additionally, a significant contribution was noted towards B2B problems compared to B2C. Hence, research needs to adopt a B2C perspective to be able to offer consumer-oriented solutions to the industry.

A comprehensive review reveals the research gap and implication presented in section VI stating that F&A needs a transformation in their supply chain by using AI techniques at the industrial level. With this, the industry can move towards a digital and sustainable supply chain. The implications and future directions proposed in this study can be beneficial for academic and industrial researchers, and industrial practitioners, who are willing to provide a substantial contribution to the subject area.

Regardless of this valuable work, there are some limitations to this study. First, the articles were searched in two databases, while there could be other databases that are relevant. Second, the publications apart from English were excluded. There could be valuable research available in other languages. Third, even though the review process was completed rigorously, it could still be prone to human error. Fourth, this study was restricted to researches with an empirical validation of the proposed AI models. There could be beneficial theoretical and conceptual frameworks in research, which were missed because of the exclusion criterion. Fifth, even though the synonyms considered for article retrieval were carefully decided, the study could have missed some articles due to different terminology being used. Seventh, grey literature was excluded from the review process, which also includes industrial reports. These reports can provide helpful insights and contribution in this domain.

## ACKNOWLEDGMENT

The authors would like to thank and appreciate the support of the expert researcher for helping us with validating the proposed taxonomy.

## REFERENCES

- [1] Lenzing, Statista. *Demand Share of Apparel Market Worldwide From 2005 to 2020, by Region*. Accessed: May 10, 2019. [Online]. Available: <https://www.statista.com/statistics/821457/demand-share-of-global-apparel-market-by-region/>

- [2] Statista. *Worldwide Forecasted Sales Growth in the Fashion Industry in 2019, by Region*. Accessed: May 10, 2019. [Online]. Available: <https://www.statista.com/statistics/802943/fashion-industry-sales-growth-worldwide-by-region/>
- [3] G. Sweeney, "It's the second dirtiest thing in the world—And you're wearing it," *AlterNet*, Aug. 2015.
- [4] X. Shang, X. Liu, G. Xiong, C. Cheng, Y. Ma, and T. R. Nyberg, "Social manufacturing cloud service platform for the mass customization in apparel industry," in *Proc. IEEE Int. Conf. Service Oper. Logistics, Informatics*, Jul. 2013, pp. 220–224.
- [5] J. H. Thrall, X. Li, Q. Li, C. Cruz, S. Do, K. Dreyer, and J. Brink, "Artificial intelligence and machine learning in radiology: Opportunities, challenges, pitfalls, and criteria for success," *J. Amer. College Radiol.*, vol. 15, no. 3, pp. 504–508, 2018.
- [6] G. Xiong, F. Zhu, X. Liu, X. Dong, W. Huang, S. Chen, and K. Zhao, "Cyber-physical-social system in intelligent transportation," *IEEE/CAA J. Autom. Sinica*, vol. 2, no. 3, pp. 320–333, Jul. 2015.
- [7] X. Xu and Q. Hua, "Industrial big data analysis in smart factory: Current status and research strategies," *IEEE Access*, vol. 5, pp. 17543–17551, 2017.
- [8] Z. X. Guo, W. K. Wong, S. Y. S. Leung, and M. Li, "Applications of artificial intelligence in the apparel industry: A review," *Text. Res. J.*, vol. 81, no. 18, pp. 1871–1892, 2011.
- [9] R. Nayak and R. Padhye, "Artificial intelligence and its application in the apparel industry," in *Automation in Garment Manufacturing*. Amsterdam, The Netherlands: Elsevier, 2018, pp. 109–138.
- [10] C. K. H. Lee, K. L. Choy, K. M. Y. Law, and G. T. S. Ho, "An intelligent system for production resources planning in Hong Kong garment industry," in *Proc. IEEE Int. Conf. Ind. Eng. Eng. Manage.*, Dec. 2012, pp. 889–893.
- [11] G. M. Nasira and P. Banumathi, "An intelligent system for automatic fabric inspection," *Asian J. Inf. Technol.*, vol. 13, no. 6, pp. 308–312, 2014.
- [12] M.-K. Chen, Y.-H. Wang, and T.-Y. Hung, "Establishing an order allocation decision support system via learning curve model for apparel logistics," *J. Ind. Prod. Eng.*, vol. 31, no. 5, pp. 274–285, 2014.
- [13] C. Giri, N. Harale, S. Thomassey, and X. Zeng, "Analysis of consumer emotions about fashion brands: An exploratory study," in *Proc. Data Sci. Knowl. Eng. Sens. Decis. Support*, 2018, pp. 1567–1574.
- [14] A. Acharya, S. K. Singh, V. Pereira, and P. Singh, "Big data, knowledge co-creation and decision making in fashion industry," *Int. J. Inf. Manage.*, vol. 42, pp. 90–101, Oct. 2018.
- [15] S. Jain, J. Bruniaux, X. Zeng, and P. Bruniaux, "Big data in fashion industry," *IOP Conf. Mater. Sci. Eng.*, vol. 254, no. 15, Oct. 2017, Art. no. 152005.
- [16] E. W. T. Ngai, S. Peng, P. Alexander, and K. K. L. Moon, "Decision support and intelligent systems in the textile and apparel supply chain: An academic review of research articles," *Expert Syst. Appl.*, vol. 41, no. 1, pp. 81–91, 2014.
- [17] P. Yildirim, D. Birant, and T. Alpyildiz, "Data mining and machine learning in textile industry," *Wiley Interdiscip. Rev. Data Mining Knowl. Discovery*, vol. 8, no. 1, p. e1228, Jan. 2018.
- [18] A. Booth, A. Sutton, and D. Papaioannou, *Systematic Approaches to a Successful Literature Review*. Newbury Park, CA, USA: Sage, 2016.
- [19] S. A. AlRyalat, L. W. Malkawi, and S. M. Momani, "Comparing bibliometric analysis using pubmed, scopus, and Web of science databases," *J. Vis. Express*, to be published.
- [20] A. Martín-Martín, E. Orduña-Malea, M. Thelwall, and E. D. López-Cózar, "Google scholar, Web of science, and scopus: A systematic comparison of citations in 252 subject categories," *J. Informetrics*, vol. 12, no. 4, pp. 1160–1177, Nov. 2018.
- [21] B. Mihm, "Fast fashion in a flat world: Global sourcing strategies," *IBER*, vol. 9, no. 6, Jun. 2010.
- [22] P. Sen, *B2B Marketing Evolution: Customers Demand B2C Experiences*. Accessed: May 10, 2019. [Online]. Available: <https://www.the-future-of-commerce.com/2019/03/14/b2b-marketing-evolution/>
- [23] J. M. Kozar and K. Y. H. Connell, "Barriers to socially responsible apparel purchasing behavior: Are consumers right?" in *Marketing Orientations in a Dynamic Business World*. Cham, Switzerland: Springer, 2017, pp. 79–85.
- [24] J. Manyika, M. Chui, B. Brown, J. Bughin, and R. Dobbs, *Big Data: The Next Frontier for Innovation, Competition, and Productivity*. McKinsey Global Institute, 2011.
- [25] V. Bhardwaj and A. Fairhurst, "Fast fashion: Response to changes in the fashion industry," *Int. Rev. Retail, Distrib. Consum. Res.*, vol. 20, no. 1, pp. 165–173, Feb. 2010.
- [26] R. L. Sandhusen, *Marketing*. Hauppauge, NY, USA: B.E.S. Publishing, 2008.
- [27] G. B. Shelly and H. J. Rosenblatt, *Systems Analysis and Design*. Borton, MA, USA: Course Technology, 2011.
- [28] M. Singh, "E-services and their role in B2C e-commerce," *Manag. Service Qual. Int. J.*, vol. 12, no. 6, pp. 434–446, Dec. 2002.
- [29] L. J. Anderson-Connell, P. V. Ulrich, and E. L. Brannon, "A consumer-driven model for mass customization in the apparel market," *J. Fashion Marketing Manag. Int. J.*, vol. 6, no. 3, pp. 240–258, Sep. 2002.
- [30] E. Brynjolfsson, Y. J. Hu, and M. S. Rahman, *Competing in the Age of Omnichannel Retailing*, vol. 54, no. 4. Cambridge, MA, USA: MIT, 2013.
- [31] L. B. Hales, M. L. Hales, and D. T. Collins, "Improving real-time expert control systems through deep data mining of plant data and global plant-wide energy monitoring and analysis," in *Proc. SME Annu. Meeting Exhib. SME, Meeting Preprints*, 2012, pp. 414–417.
- [32] C. Murray, "Retailers' disconnect with shoppers is costing them," A Forrester Consulting Thought Leadership Paper Commissioned by Cognizant, Tech. Rep., 2017. [Online]. Available: <https://www.cognizant.com/whitepapers/retailers-disconnect-with-shoppers-is-costing-them-codex2693.pdf>
- [33] Z.-H. Hu and X.-K. Yu, "Optimization of fast-fashion apparel transshipment among retailers," *Text. Res. J.*, vol. 84, no. 20, pp. 2127–2139, 2014.
- [34] M. Y. Wong, Y. Zhou, and H. Xu, "Big data in fashion industry: Color cycle mining from runway data," in *Proc. AMCIS, Surfing IT Innov. Wave-22nd Amer. Conf. Inf. Syst.*, 2016, pp. 1–10.
- [35] C. Bishop, *Pattern Recognition and Machine Learning*. Springer, 2006.
- [36] R. Stuart and P. Norvig, *Artificial Intelligence: A Modern Approach*. Upper Saddle River, NJ, USA: Prentice-Hall, 2009.
- [37] M. Mohri, A. Rostamizadeh, and A. Talwalkar, *Foundations of Machine Learning*. Cambridge, MA, USA: MIT Press, 2012.
- [38] S. V. Kumar and S. Poonkuzhali, "Improving the sales of garments by forecasting market trends using data mining techniques," *Int. J. Pure Appl. Math.*, vol. 119, no. 7, pp. 797–805, 2018.
- [39] S.-W. Hsiao, C.-H. Lee, R.-Q. Chen, and C.-H. Yen, "An intelligent system for fashion colour prediction based on fuzzy C-means and gray theory," *Color Res. Appl.*, vol. 42, no. 2, pp. 273–285, Apr. 2017.
- [40] Brahmadeep and S. Thomassey, "Intelligent demand forecasting systems for fast fashion," in *Information Systems for the Fashion and Apparel Industry*. Woodhead Publishing, 2016, pp. 145–161. [Online]. Available: [https://www.researchgate.net/publication/301345185\\_Intelligent\\_demand\\_forecasting\\_systems\\_for\\_fast\\_fashion](https://www.researchgate.net/publication/301345185_Intelligent_demand_forecasting_systems_for_fast_fashion)
- [41] A. Ghosh, T. Guha, R. B. Bhar, and S. Das, "Pattern classification of fabric defects using support vector machines," *Int. J. Clothing Sci. Technol.*, vol. 23, nos. 2–3, pp. 142–151, 2011.
- [42] D. Z. Pavlinić and J. Geršak, "Design of the system for prediction of fabric behaviour in garment manufacturing processes," *Int. J. Cloth. Sci. Technol.*, vol. 16, nos. 1–2, pp. 252–261, 2004.
- [43] P. G. Keen, "Decision support systems: A research perspective," in *Decision Support Systems: Issues and Challenges: Proceedings of an International Task Force Meeting*, 1980, pp. 23–44.
- [44] R. H. Sprague, Jr., "A framework for the development of decision support systems," *MIS Quart.*, vol. 4, no. 4, pp. 1–26, 1980.
- [45] Y. Tu and E. H. H. Yeung, "Integrated maintenance management system in a textile company," *Int. J. Adv. Manuf. Technol.*, vol. 13, no. 6, pp. 453–462, 1997.
- [46] W. K. Wong and S. Y. S. Leung, "Genetic optimization of fabric utilization in apparel manufacturing," *Int. J. Prod. Econ.*, vol. 114, no. 1, pp. 376–387, 2008.
- [47] P. Jackson, *Introduction to Expert Systems*, vol. 2. Reading, MA, USA: Addison-Wesley, 1990.
- [48] C. T. Leondes, *Expert Systems: The Technology of Knowledge Management and Decision Making for the 21st Century*. New York, NY, USA: Academic, 2002.
- [49] N. N. Stella and A. O. Chuks, "Expert system: A catalyst in educational development in Nigeria," 2011. [Online]. Available: <http://hrmars.com/admin/pics/261.pdf>
- [50] K. Metaxiotis, "RECO: An expert system for the reduction of environmental cost in the textile industry," *Inf. Manage. Comput. Secur.*, vol. 12, no. 3, pp. 218–227, 2004.

- [51] W. K. Wong, X. H. Zeng, and W. M. R. Au, "A decision support tool for apparel coordination through integrating the knowledge-based attribute evaluation expert system and the T-S fuzzy neural network," *Expert Syst. Appl.*, vol. 36, no. 2, pp. 2377–2390, 2009.
- [52] G. F. Luger and W. A. Stubblefield, *Artificial Intelligence: Structures and Strategies for Complex Problem Solving*. San Francisco, CA, USA: Benjamin Cummings, 1993.
- [53] D. Poole, A. Mackworth, and R. Goebel, *Computational Intelligence: A Logical Approach*. New York, NY, USA: Oxford Univ. Press, 1998.
- [54] N. J. Nilsson, *Artificial Intelligence: A New Synthesis*. San Mateo, CA, USA: Morgan Kaufmann, 1998.
- [55] G. Tecuci, "Artificial intelligence," *Wiley Interdiscip. Rev. Comput. Statist.*, vol. 4, no. 2, pp. 168–180, Mar./Apr. 2012.
- [56] J. H. Holland, *Adaptation in Natural and Artificial Systems: An Introductory Analysis With Applications to Biology, Control, and Artificial Intelligence*. Cambridge, MA, USA: MIT Press, 1992.
- [57] J. H. Holland, *Genetic Programming Complex Adaptive Systems Genetic Programming On the Programming of Computers by Means of Natural Selection*. Cambridge, MA, USA: MIT Press, 1992.
- [58] R. Guruprasad and B. K. Behera, "Genetic algorithms and its application to textiles," *Textile Asia*, vol. 40, nos. 4–5, pp. 35–38, 2009.
- [59] Z.-J. Lu, Q. Xiang, Y.-M. Wu, and J. Gu, "Application of support vector machine and genetic algorithm optimization for quality prediction within complex industrial process," in *Proc. IEEE Int. Conf. Ind. Inform. (INDIN)*, Jul. 2015, pp. 98–103.
- [60] P. C. L. Hui, K. C. C. Chan, K. W. Yeung, and F. S. F. Ng, "Application of artificial neural networks to the prediction of sewing performance of fabrics," *Int. J. Cloth. Sci. Technol.*, vol. 19, no. 5, pp. 291–318, Oct. 2007.
- [61] D. H. Ballard and C. M. Brown, *Computer Vision*. Upper Saddle River, NJ, USA: Prentice-Hall, 1982.
- [62] M. Sonka, V. Hlavac, and R. Boyle, *Image Processing: Analysis and Machine Vision*. Boston, MA, USA: Thompson Learning, 2008.
- [63] B. Jahne, *Computer Vision and Applications: A Guide for Students and Practitioners*. Amsterdam, The Netherlands: Elsevier, 2000.
- [64] D. Forsyth and J. Ponce, *Computer Vision: A Modern Approach*. Upper Saddle River, NJ, USA: Prentice-Hall, 2003.
- [65] T. Morris, *Computer Vision and Image Processing—Tim Morris—Macmillan International Higher Education*. Palgrave Macmillan. Accessed: May 10, 2019. [Online]. Available: <https://www.macmillanihe.com/page/detail/computer-vision-and-image-processing-tim-morris?sf1=barcode&st1=9780333994511>
- [66] C. Steger, M. Ulrich, and C. Wiedemann, *Machine Vision Algorithms and Applications*. Hoboken, NJ, USA: Wiley, 2018.
- [67] C.-F. J. Kuo, C.-L. Lee, and C.-Y. Shih, "Image database of printed fabric with repeating dot patterns part (I)—Image archiving," *Text. Res. J.*, vol. 87, no. 17, pp. 2089–2105, Oct. 2017.
- [68] M. Yuan, I. R. Khan, F. Farbiz, S. Yao, A. Niswar, and M.-H. Foo, "A mixed reality virtual clothes try-on system," *IEEE Trans. Multimedia*, vol. 15, no. 8, pp. 1958–1968, Dec. 2013.
- [69] G. A. Cushen and M. S. Nixon, "Markerless real-time garment retexturing from monocular 3D reconstruction," in *Proc. IEEE Int. Conf. Signal Image Process. Appl. (ICSIPA)*, Nov. 2011, pp. 88–93.
- [70] Y. Zhao, J. Song, A. Montazeri, M. M. Gupta, Y. Lin, C. Wang, and W. J. Zhang, "Mining affective words to capture customer's affective response to apparel products," *Text. Res. J.*, vol. 88, no. 12, pp. 1426–1436, Jun. 2018.
- [71] K. Liu, J. Wang, E. Kamalha, V. Li, and X. Zeng, "Construction of a prediction model for body dimensions used in garment pattern making based on anthropometric data learning," *J. Textile Inst.*, vol. 108, no. 12, pp. 2107–2114, Apr. 2017.
- [72] A. F. Tehrani and D. Ahrens, "Modified sequential k-means clustering by utilizing response: A case study for fashion products," *Expert Syst.*, vol. 34, no. 6, Dec. 2017, Art. no. e12226.
- [73] B.-R. Li, Y. Wang, and K.-S. Wang, "A novel method for the evaluation of fashion product design based on data mining," *Adv. Manuf.*, vol. 5, no. 4, pp. 370–376, Dec. 2017.
- [74] M. Hamad, S. Thomassey, and P. Bruniaux, "A new sizing system based on 3D shape descriptor for morphology clustering," *Comput. Ind. Eng.*, vol. 113, pp. 683–692, Nov. 2017.
- [75] K. Liu, X. Zeng, P. Bruniaux, J. Wang, E. Kamalha, and X. Tao, "Fit evaluation of virtual garment try-on by learning from digital pressure data," *Knowl.-Based Syst.*, vol. 133, pp. 174–182, Oct. 2017.
- [76] A. Fallahpour, K. Y. Wong, E. U. Olugu, and S. N. Musa, "A predictive integrated genetic-based model for supplier evaluation and selection," *Int. J. Fuzzy Syst.*, vol. 19, no. 4, pp. 1041–1057, Aug. 2017.
- [77] K. J. Ferreira, B. H. A. Lee, and D. Simchi-Levi, "Analytics for an online retailer: Demand forecasting and price optimization," *Manuf. Service Oper. Manage.*, vol. 18, no. 1, pp. 69–88, 2016.
- [78] A. F. Tehrani and D. Ahrens, "Supervised regression clustering: A case study for fashion products," *Int. J. Bus. Anal.*, vol. 3, no. 4, pp. 21–40, 2016.
- [79] P. Q. Brito, C. Soares, S. Almeida, A. Monte, and M. Buyoet, "Customer segmentation in a large database of an online customized fashion business," *Robot. Comput. Integr. Manuf.*, vol. 36, pp. 93–100, Dec. 2015.
- [80] V. Mozafari and P. Payvandy, "Application of data mining technique in predicting worsted spun yarn quality," *J. Text. Inst.*, vol. 105, no. 1, pp. 100–108, 2014.
- [81] C.-C. Chen, "RFID-based intelligent shopping environment: A comprehensive evaluation framework with neural computing approach," *Neural Comput. Appl.*, vol. 25, nos. 7–8, pp. 1685–1697, 2014.
- [82] S. M. Darwish, "Soft computing applied to the build of textile defects inspection system," *IET Comput. Vis.*, vol. 7, no. 5, pp. 373–381, Oct. 2013.
- [83] C.-Y. Tsai and C.-H. Hsu, "Developing standard elderly aged female size charts based on anthropometric data to improve manufacturing using artificial neural network-based data mining," *Theor. Issues Ergon. Sci.*, vol. 14, no. 3, pp. 258–272, 2013.
- [84] J.-F. Liang, J.-M. Liang, and J.-P. Wang, "Empirical study on B/C apparel consumption behavior based on data mining technology," *J. Donghua Univ. (English Ed.)*, vol. 30, no. 6, pp. 530–536, 2013.
- [85] H. L. Viktor, I. Peña, and E. Paquet, "Who are our clients: Consumer segmentation through explorative data mining," *Int. J. Data Mining, Model. Manage.*, vol. 4, no. 3, pp. 286–308, 2012.
- [86] Y. Yu, T.-M. Choi, and C.-L. Hui, "An intelligent fast sales forecasting model for fashion products," *Expert Syst. Appl.*, vol. 38, no. 6, pp. 7373–7379, 2011.
- [87] B. Karthikeyan and L. M. Sztandera, "Analysis of tactile perceptions of textile materials using artificial intelligence techniques: Part 1: Forward engineering," *Int. J. Clothing Sci. Technol.*, vol. 22, nos. 2–3, pp. 187–201, 2010.
- [88] C.-H. Hsu, "Data mining to improve industrial standards and enhance production and marketing: An empirical study in apparel industry," *Expert Syst. Appl.*, vol. 36, no. 3, pp. 4185–4191, 2009.
- [89] C. Altinoz, "Supplier selection for industry: A fuzzy rule-based scoring approach with a focus on usability," *Int. J. Integr. Supply Manage.*, vol. 4, nos. 3–4, pp. 303–321, 2008.
- [90] A. E. Amin, A. S. El-Geheini, I. A. El-Hawary, and R. A. El-Beali, "Detecting the fault from spectrograms by using genetic algorithm techniques," *Autex Res. J.*, vol. 7, no. 2, pp. 80–88, 2007.
- [91] P. N. Koustoumpardis, J. S. Fourkiotis, and N. A. Aspragathos, "Intelligent evaluation of fabrics' extensibility from robotized tensile test," *Int. J. Clothing Sci. Technol.*, vol. 19, no. 2, pp. 80–98, 2007.
- [92] C.-H. Hsu and M.-J. J. Wang, "Using innovative technology to establish sizing systems," *Int. J. Innov. Learn.*, vol. 2, no. 3, pp. 233–245, 2005.
- [93] I. Soufflet, M. Calonnier, and C. Dacremont, "A comparison between industrial experts' and novices' haptic perceptual organization: A tool to identify descriptors of the handle of fabrics," *Food Qual. Preference*, vol. 15, nos. 7–8, pp. 689–699, 2004.
- [94] K.-Y. Jung, Y.-J. Na, and J.-H. Lee, "Creating user-adapted design recommender system through collaborative filtering and content based filtering," in *Proc. Portuguese Conf. Artif. Intell. Berlin, Germany: Springer*, 2003, pp. 204–208.
- [95] S. Sette and L. Boullart, "An implementation of genetic algorithms for rule based machine learning," *Eng. Appl. Artif. Intell.*, vol. 13, no. 4, pp. 381–390, 2000.
- [96] C. K. Park, D. H. Lee, and T. J. Kang, "A new evaluation of seam pucker and its applications," *Int. J. Clothing Sci. Technol.*, vol. 9, nos. 2–3, pp. 252–255, 1997.
- [97] A. Aksoy and N. Öztürk, "Design of an intelligent decision support system for global outsourcing decisions in the apparel industry," *J. Text. Inst.*, vol. 107, no. 10, pp. 1322–1335, 2016.
- [98] C. K. H. Lee, Y. K. Tse, G. T. S. Ho, and K. L. Choy, "Fuzzy association rule mining for fashion product development," *Ind. Manage. Data Syst.*, vol. 115, no. 2, pp. 383–399, 2015.

- [99] Z.-H. Hu, C. Wei, and X.-K. Yu, "Apparel distribution with uncertain try-on time by evolutionary algorithm," *Int. J. Clothing Sci. Technol.*, vol. 27, no. 1, pp. 75–90, 2015.
- [100] A. S. Kumar and V. M. Bhaskaran, "Multidirectional product decision support system for user activation assessment in textile industry using collaborative filtering methods," *Int. J. Appl. Eng. Res.*, vol. 10, no. 20, pp. 17205–17209, 2015.
- [101] C. K. H. Lee, K. L. Choy, K. M. Y. Law, and G. T. S. Ho, "Application of intelligent data management in resource allocation for effective operation of manufacturing systems," *J. Manuf. Syst.*, vol. 33, no. 3, pp. 412–422, 2014.
- [102] D. Nakandala, P. Samaranyake, and H. C. W. Lau, "A fuzzy-based decision support model for monitoring on-time delivery performance: A textile industry case study," *Eur. J. Oper. Res.*, vol. 225, no. 3, pp. 507–517, 2013.
- [103] B. Rabenasolo and X. Zeng, "A risk-based multi-criteria decision support system for sustainable development in the textile supply chain," in *Handbook on Decision Making* (Intelligent Systems Reference Library), vol. 33, J. Lu, L. C. Jain, and G. Zhang, Eds. Berlin, Germany: Springer, 2012.
- [104] A. Aksoy, N. Ozturk, and E. Sucky, "A decision support system for demand forecasting in the clothing industry," *Int. J. Clothing Sci. Technol.*, vol. 24, no. 4, pp. 221–236, 2012.
- [105] D. Kumar, J. Singh, and O. P. Singh, "A decision support system for supplier selection for Indian textile industry using analytic hierarchy process based on fuzzy simulation," *Int. J. Bus. Perform. Supply Chain Model.*, vol. 3, no. 4, pp. 364–382, 2011.
- [106] A. Majumdar, R. Mangla, and A. Gupta, "Developing a decision support system software for cotton fibre grading and selection," *Indian J. Fibre Text. Res.*, vol. 35, no. 3, pp. 195–200, 2010.
- [107] S. Chakraborty and K. Prasad, "A quality function deployment-based expert system for cotton fibre selection," *J. Inst. Eng. E.*, vol. 99, no. 1, pp. 43–53, Jun. 2018.
- [108] V. Istrat and N. Lalić, "Association rules as a decision making model in the textile industry," *Fibres Text. Eastern Eur.*, vol. 25, no. 4, pp. 8–14, 2017.
- [109] K. Santiago-Santiago, A. L. Laureano-Cruces, J. M. A. Antuñano-Barranco, O. Domínguez-Pérez, and E. Sarmiento-Bustos, "An expert system to improve the functioning of the clothing industry," *Int. J. Clothing Sci. Technol.*, vol. 27, no. 1, pp. 99–128, 2015.
- [110] C. K. H. Lee, G. T. S. Ho, K. L. Choy, and G. K. H. Pang, "A RFID-based recursive process mining system for quality assurance in the garment industry," *Int. J. Prod. Res.*, vol. 52, no. 14, pp. 4216–4238, 2014.
- [111] S. Jamal, H. Esmaili, and S. E. Maryam, "Developing a hybrid intelligent model for constructing a size recommendation expert system in textile industries," *Int. J. Clothing Sci. Technol.*, vol. 25, no. 5, pp. 338–349, 2013.
- [112] C. K. H. Lee, K. L. Choy, G. T. S. Ho, K. S. Chin, K. M. Y. Law, and Y. K. Tse, "A hybrid OLAP-association rule mining based quality management system for extracting defect patterns in the garment industry," *Expert Syst. Appl.*, vol. 40, no. 7, pp. 2435–2446, 2013.
- [113] Y. Yu, C.-L. Hui, and T.-M. Choi, "An empirical study of intelligent expert systems on forecasting of fashion color trend," *Expert Syst. Appl.*, vol. 39, no. 4, pp. 4383–4389, 2012.
- [114] W. K. Wong, S. Y. S. Leung, Z. X. Guo, X. H. Zeng, and P. Y. Mok, "Intelligent product cross-selling system with radio frequency identification technology for retailing," *Int. J. Prod. Econ.*, vol. 135, no. 1, pp. 308–319, 2012.
- [115] W.-S. Lo, T.-P. Hong, and R. Jeng, "A framework of E-SCM multi-agent systems in the fashion industry," *Int. J. Prod. Econ.*, vol. 114, no. 2, pp. 594–614, 2008.
- [116] C. Rong-Chang, L. Chih-Chang, and L. Shiu-Shiun, "An automatic decision support system based on genetic algorithm for global apparel manufacturing," *Int. J. Soft Comput.*, vol. 1, no. 1, pp. 17–21, 2006.
- [117] T. W. Lau, P. C. L. Hui, F. S. F. Ng, and K. C. C. Chan, "A new fuzzy approach to improve fashion product development," *Comput. Ind.*, vol. 57, no. 1, pp. 82–92, 2006.
- [118] I. P. Tatsiopoulos, S. T. Ponis, and E. A. Hadzilias, "An e-releaser of production orders in the extended enterprise," *Prod. Planning Control*, vol. 15, no. 2, pp. 119–132, 2004.
- [119] Y. Liu and Z.-F. Geng, "Three-dimensional garment computer aided intelligent design," *J. Ind. Text.*, vol. 33, no. 1, pp. 43–54, 2003.
- [120] W. D. Cooper and C. Saydam, "A general decision support systems approach to the port scheduling problem for pressure beck operations," *J. Text. Inst.*, vol. 94, nos. 1–2, pp. 1–11, 2003.
- [121] C. Saydam and W. D. Cooper, "A decision support system for scheduling jobs on multi-port dyeing machines," *Int. J. Oper. Prod. Manage.*, vol. 22, nos. 9–10, pp. 1054–1065, 2002.
- [122] C. Jarvis, S. Davis, B. Kernodle, W. Masuchun, and S. AngannaMuthu, "Decision support for balanced inventory flow replenishment system," *Int. J. Clothing Sci. Technol.*, vol. 12, no. 6, pp. 100–102, 2000.
- [123] R. Convert, L. Schacher, and P. Viallier, "An expert system for the dyeing recipes determination," *J. Intell. Manuf.*, vol. 11, no. 2, pp. 145–155, 2000.
- [124] C. Eckert, I. A. N. Kelly, and M. Stacey, "Interactive generative systems for conceptual design: An empirical perspective," *Artif. Intell. Eng. Des. Anal. Manuf.*, vol. 13, no. 4, pp. 303–320, 1999.
- [125] C. A. B. e Costa, L. Ensslin, E. C. Corrêa, and J.-C. Vansnick, "Decision support systems in action: Integrated application in a multicriteria decision aid process," *Eur. J. Oper. Res.*, vol. 113, no. 2, pp. 315–335, 1999.
- [126] L. J. Anderson, C. Warfield, M. Barry, and C. Emery, "Toward a national model of an electronic decision support system for domestic sourcing of textiles and apparel," *Clothing Text. Res. J.*, vol. 15, no. 2, pp. 65–75, 1997.
- [127] C. K. Park, D. H. Lee, and T. Kang, "Knowledge-base construction of a garment manufacturing expert system," *Int. J. Clothing Sci. Technol.*, vol. 8, no. 5, pp. 11–28, 1996.
- [128] E. K. Bye and M. R. DeLong, "A visual sensory evaluation of the results of two pattern grading methods," *Clothing Text. Res. J.*, vol. 12, no. 4, pp. 1–7, 1994.
- [129] G. Martino, R. Iannone, M. Fera, S. Miranda, and S. Riemma, "Fashion retailing: A framework for supply chain optimization," *Uncertain Supply Chain Manage.*, vol. 5, no. 3, pp. 243–272, 2017.
- [130] C. K. H. Lee, K. L. Choy, G. T. S. Ho, and C. H. Y. Lam, "A slippery genetic algorithm-based process mining system for achieving better quality assurance in the garment industry," *Expert Syst. Appl.*, vol. 46, pp. 236–248, Mar. 2016.
- [131] G. Martino, B. Yuce, R. Iannone, and M. S. Packianather, "Optimisation of the replenishment problem in the fashion retail industry using Tabu-Bees algorithm," *IFAC-PapersOnLine*, vol. 49, no. 12, pp. 1685–1690, 2016.
- [132] F. Caro and J. Gallien, "Inventory management of a fast-fashion retail network," *Oper. Res.*, vol. 58, no. 2, pp. 257–273, 2010.
- [133] S.-S. Li, R.-C. Chen, and C.-C. Lin, "A genetic algorithm-based decision support system for allocating international apparel demand," *WSEAS Trans. Inf. Sci. Appl.*, vol. 3, no. 7, pp. 1294–1299, 2006.
- [134] N. Werghi, Y. Xiao, and J. P. Siebert, "A functional-based segmentation of human body scans in arbitrary postures," *IEEE Trans. Man, Cybern., B, Cybern.*, vol. 36, no. 1, pp. 153–165, Feb. 2006.
- [135] S. Sette and L. Boullart, "Genetic programming: Principles and applications," *Eng. Appl. Artif. Intell.*, vol. 14, no. 6, pp. 727–736, 2001.
- [136] N. I. Karacapılıdis and C. P. Pappis, "Production planning and control in textile industry: A case study," *Comput. Ind.*, vol. 30, no. 2, pp. 127–144, 1996.
- [137] Y. Shin, Y. Kim, and E. Y. Kim, "Automatic textile image annotation by predicting emotional concepts from visual features," *Image Vis. Comput.*, vol. 28, no. 3, pp. 526–537, 2010.
- [138] C. W. M. Yuen, W. K. Wong, S. Q. Qian, L. K. Chan, and E. H. K. Fung, "A hybrid model using genetic algorithm and neural network for classifying garment defects," *Expert Syst. Appl.*, vol. 36, no. 2, pp. 2037–2047, 2009.
- [139] W. K. Wong, C. W. M. Yuen, D. D. Fan, L. K. Chan, and E. H. K. Fung, "Stitching defect detection and classification using wavelet transform and BP neural network," *Expert Syst. Appl.*, vol. 36, no. 2, pp. 3845–3856, Mar. 2009.
- [140] C. W. M. Yuen, W. K. Wong, S. Q. Qian, D. D. Fan, L. K. Chan, and E. H. K. Fung, "Fabric stitching inspection using segmented window technique and BP neural network," *Text. Res. J.*, vol. 79, no. 1, pp. 24–35, 2009.
- [141] Q. Lu, Z.-J. Lyu, Q. Xiang, Y. Zhou, and J. Bao, "Research on data mining service and its application case in complex industrial process," in *Proc. IEEE Int. Conf. Autom. Sci. Eng.*, Aug. 2017, pp. 1124–1129.
- [142] D. Banerjee, N. Ganguly, S. Sural, and K. S. Rao, "One for the road: Recommending male street attire," in *Proc. Pacific-Asia Conf. Knowl. Discovery Data Mining*, 2018, pp. 571–582.

- [143] D. Venkataraman, N. Vinay, T. V. Vamsi, S. P. Boppudi, R. Y. Reddy, and P. Balasubramanian, "Yarn price prediction using advanced analytics model," in *Proc. IEEE Int. Conf. Comput. Intell. Comput. Res. (ICICIC)*, Dec. 2016, pp. 1–8.
- [144] C. Fiarni, H. Maharani, and R. Pratama, "Sentiment analysis system for Indonesia online retail shop review using hierarchy Naive Bayes technique," in *Proc. 4th Int. Conf. Inf. Commun. Technol. (ICoICT)*, May 2016, pp. 1–6.
- [145] W. Lin and Z. Miao, "The construction of fabric grading model based on SOM algorithm," in *Proc. Int. Conf. Intell. Transp., Big Data Smart City (ICITBS)*, Dec. 2015, pp. 841–843.
- [146] M. Wong, Y. Zhou, and H. Xu, "Big data in fashion industry: Color cycle mining from runway data," in *Proc. 22nd Amer. Conf. Inf. Syst.*, 2016, pp. 1–10.
- [147] A. P. Martins, P. Bruniiaux, and S. Thomassey, "Cluster-based sales forecasting of fast fashion using linguistic variables and ELM," in *Proc. 12th Int. FLINS Conf. Uncertainty Modeling Knowl. Eng. Decis. Making*, 2016, pp. 978–983.
- [148] D. Siegmund, A. Kuijper, and A. Braun, "Stereo-image normalization of voluminous objects improves textile defect recognition," in *Proc. Int. Symp. Vis. Comput. Cham, Switzerland: Springer*, 2016, pp. 181–192.
- [149] J. B. Pérez, A. G. Arrieta, A. H. Encinas, and A. Q. Dios, "Textile engineering and case based reasoning," in *Proc. 13th Int. Conf. Distrib. Comput. Artif. Intell.*, vol. 474, 2016, pp. 423–431.
- [150] D. P. P. Mesquita, A. N. A. Neto, J. F. de Queiroz Neto, J. P. P. Gomes, and L. R. Rodrigues, "Using robust extreme learning machines to predict cotton yarn strength and hairiness," in *Proc. 24th Eur. Symp. Artif. Neural Netw. (ESANN)*, 2016, pp. 65–70.
- [151] C. D. Kreyenhagen, T. I. Aleshin, J. E. Bouchard, A. M. I. Wise, and R. K. Zalegowski, "Using supervised learning to classify clothing brand styles," in *Proc. IEEE Syst. Inf. Eng. Design Symp. (SIEDS)*, Apr. 2014, pp. 239–243.
- [152] A. Kitipong, W. Rueangsirasak, and R. Chairsricharoen, "Classification system for traditional textile: Case study of the batik," in *Proc. 13th Int. Symp. Commun. Inf. Technol., Commun. Inf. Technol. New Life Style Beyond Cloud (ISCIT)*, Sep. 2013, pp. 767–771.
- [153] P. F. Li, J. Wang, H. H. Zhang, and J. F. Jing, "Automatic woven fabric classification based on support vector machine," in *Proc. Int. Conf. Autom. Control Artif. Intell. (ACAI)*, no. 598, Mar. 2012, pp. 581–584.
- [154] N. A. Khalifa, S. M. Darwish, and M. A. El-Iskandarani, "Automated textile defects recognition system using computer vision and interval type-2 fuzzy logic," in *Proc. 1st Int. Conf. Innov. Eng. Syst. (ICIES)*, Dec. 2012, pp. 148–153.
- [155] E. Yesil, M. Kaya, and S. Siradag, "Fuzzy forecast combiner design for fast fashion demand forecasting," in *Proc. Int. Symp. Innov. Intell. Syst. Appl.*, Jul. 2012, pp. 1–5.
- [156] X. G. Wang, X. Z. Li, and Y. Li, "Application of cluster algorithm in clothes shape classifying," *Appl. Mech. Mater.*, vols. 55–57, pp. 1058–1062, May 2011.
- [157] X. Chen, Y. Luo, and F. Zhu, "The application of data mining in FFE of the fashion product development," in *Proc. Int. Symp. Comput. Intell. Design (ISCID)*, vol. 1, Oct. 2010, pp. 215–217.
- [158] L. Z. Fan and Q. Qiao, "Research on a fashion knowledge management platform for women garment development," in *Proc. Int. Conf. Manage. Service Sci. (MASS)*, Aug. 2010, pp. 1–3.
- [159] Y. Li, V. E. Kuzmichev, Y. Luo, and X. Wang, "Garment industry oriented clothes shape classifying by cluster," in *Proc. 2nd Int. Conf. Ind. Mechatronics Automat. (ICIMA)*, vol. 2, May 2010, pp. 472–475.
- [160] M. Selvanayaki, M. S. Vijaya, K. S. Jamuna, and S. Karpagavalli, "Supervised learning approach for predicting the quality of cotton using WEKA," in *Proc. Int. Conf. Bus. Admin. Inf. Process.*, vol. 70, 2010, pp. 382–384.
- [161] I. Peña, H. L. Viktor, and E. Paquet, "Explorative data mining for the sizing of population groups," in *Proc. 1st Int. Conf. Knowl. Discovery Inf. Retr. (KDIR)*, 2009, pp. 152–159.
- [162] L. Koehl, X. Zeng, M. Camargo, C. Fonteix, and F. Delmotte, "Analysis and identification of fashion oriented industrial products using fuzzy logic techniques," in *Proc. 3rd Int. Conf. Intell. Syst. Knowl. Eng. (ISKE)*, Nov. 2008, pp. 465–470.
- [163] J. De Armas, C. León, G. Miranda, and C. Segura, "Remote service to solve the two-dimensional cutting stock problem: An application to the canary islands costume," in *Proc. 2nd Int. Conf. Complex. Intell. Softw. Intensive Syst. (CISIS)*, Mar. 2008, pp. 971–976.
- [164] S. Ierace, R. Pinto, and S. Cavaliere, "Application of neural networks to condition based maintenance: A case study in the textile industry," *IFAC Proc. Vols.*, vol. 40, no. 3, pp. 147–152, 2007.
- [165] N. Y. Kim, Y. Shin, and E. Y. Kim, "Emotion-based textile indexing using neural networks," in *Proc. Int. Symp. Consum. Electron. (ISCE)*, 2007, pp. 349–357.
- [166] S. Da Silva Camargo and P. M. Engel, "MiRABIT: A new algorithm for mining association rules," in *Proc. Int. Conf. Chilean Comput. Sci. Soc. (SCCC)*, Nov. 2002, pp. 162–166.
- [167] Y. Kita and N. Kita, "A model-driven method of estimating the state of clothes for manipulating it," in *Proc. IEEE Workshop Appl. Comput. Vis.*, Dec. 2002, pp. 63–69.
- [168] T. Sano and H. Yamamoto, "Intelligent CAD system for Japanese kimono," in *Proc. 26th Annu. Conf. IEEE Ind. Electron. Soc. (IECON)*, vol. 1, Oct. 2000, pp. 942–947.
- [169] J. Shao, Z. Zhao, J. Wang, and P. Song, "Production management and decision system oriented to the textile enterprise based on multi-agent," in *Proc. Int. Conf. Comput. Design Appl. (ICDDA)*, vol. 4, Jun. 2010, pp. V445–V449.
- [170] Y. Chen and L.-C. Wang, "Fuzzy logic of matching sense on fashion image," in *Proc. 4th Int. Conf. Natural Comput. (ICNC)*, vol. 7, Oct. 2008, pp. 85–89.
- [171] M. Min, "A rule based expert system for analysis of mobile sales data on fashion market," in *Proc. Int. Conf. Inf. Sci. Appl. (ICISA)*, Jun. 2013, pp. 1–2.
- [172] Z.-J. Lv, Q. Xiang, and J.-G. Yang, "A novel data mining method on quality control within spinning process," *Appl. Mech. Mater.*, vol. 224, pp. 87–92, Nov. 2012.
- [173] Z. Xue, X. Zeng, L. Koehl, and Y. Chen, "Development of a fuzzy inclusion measure for investigating relations between visual and tactile perception of textile products," in *Proc. IEEE Int. Conf. Multisensor Fusion Integr. Intell. Syst.*, Sep. 2012, pp. 108–113.
- [174] B. S. Villanueva and M. Sánchez-Marré, "Case-based reasoning applied to textile industry processes," in *Proc. Int. Conf. Case-Based Reasoning*, Berlin, Germany: Springer, 2012, pp. 428–442.
- [175] L.-C. Wang, Y. Chen, and Y. Wang, "Formalization of fashion sensory data based on fuzzy set theory," in *Proc. 4th Int. Conf. Natural Comput. (ICNC)*, vol. 7, Oct. 2008, pp. 80–84.
- [176] J. L. Su, Z. Ouyang, and Y. Chen, "Design of agile infrastructure for manufacturing system with FNN based Web-enabled technology solutions," in *Proc. 7th Int. Conf. Intell. Syst. Design Appl. (ISDA)*, Oct. 2007, pp. 79–83.
- [177] R.-C. Chen, S.-S. Li, C.-C. Lin, and T.-S. Chen, "Application of self-adaptive genetic algorithm on allocating international demand to global production facilities," in *Proc. World Congr. Intell. Control Automat. (WCICA)*, vol. 2, Jun. 2006, pp. 7152–7156.
- [178] R.-C. Chen, S.-S. Li, C.-C. Lin, and C.-C. Feng, "A GA-based global decision support system for garment production," in *Proc. Int. Conf. Neural Netw. Brain (ICNNB)*, vol. 2, Oct. 2005, pp. 805–809.
- [179] M. H. F. Zarandi and M. Esmailian, "A systematic fuzzy modeling for scheduling of textile manufacturing system," in *Proc. Annu. Conf. North Amer. Fuzzy Inf. Process. Soc. (NAFIPS)*, Jul. 2003, pp. 359–364.
- [180] G. S. Loo, B. C. P. Tang, and L. Janczewski, "An adaptable human-agent collaboration information system in manufacturing (HACISM)," in *Proc. Int. Workshop Database Expert Syst. Appl. (DEXA)*, Sep. 2000, pp. 445–449.
- [181] G. Gonçalves, J. B. De Sousa, and F. L. Pereira, "The development of an integrated decision support system for a textile company: A systems engineering approach," in *Proc. Eur. Control Conf. (ECC)*, Jul. 1997, pp. 3543–3548.
- [182] G. Brewka and C. Habel, "KI-97: Advances in artificial intelligence," in *Proc. 21st Annu. German Conf. Artif. Intell.*, vol. 21, Freiburg, Germany: Springer, Sep. 1997.
- [183] D. A. Xerokostas and K. G. Aravossis, "A solution procedure for solving the nonrectangular cutting stock problem of the clothing industry using expert knowledge within an intelligent CAD environment," in *Proc. Int. Conf. Database Expert Syst. Appl.*, Berlin, Germany: Springer, 1994, pp. 214–224.
- [184] P. T. Kleeman, T. A. Calogero, C. L. Cromwell, J. D. Nan, P. F. Sanborn, and P. B. Wright, "A computer integrated system for autonomous management of production (illustrated by applications in textile manufacturing)," in *Proc. Rensselaer's 2nd Int. Conf. Comput. Integr. Manuf.*, May 1990, pp. 324–330.



- [185] N. Klement, C. Silva, and O. Gibaru, "Solving a discrete lot sizing and scheduling problem with unrelated parallel machines and sequence dependent setup using a generic decision support tool," in *Proc. IFIP Int. Conf. Adv. Prod. Manage. Syst.*, vol. 513, 2017, pp. 459–466.
- [186] A. Uber, Jr., P. J. de Freitas Filho, and R. A. Silveira, "E-HIPS: An extension of the framework HIPS for stagger of distributed process in production systems based on multiagent systems and memetic algorithms," in *Proc. Mex. Int. Conf. Artif. Intell.* Cham, Switzerland: Springer, 2015, pp. 413–430.
- [187] P. Y. Mok, "Intelligent apparel production planning for optimizing manual operations using fuzzy set theory and evolutionary algorithms," in *Proc. 5th Int. Workshop Genetic Evol. Fuzzy Syst. Symp. Comput. Intell. (IEEE GEFS SSCI)*, Apr. 2011, pp. 103–110.
- [188] L. S. Admuthé and S. D. Apte, "Computational model using ANFIS and GA: Application for textile spinning process," in *Proc. 2nd IEEE Int. Conf. Comput. Sci. Inf. Technol. (ICCSIT)*, Aug. 2009, pp. 110–114.
- [189] L. Suyi and Z. Leduo, "Textile pattern generation technique based on quasi-regular pattern theory and their transform," in *Proc. Pacific-Asia Workshop Comput. Intell. Ind. Appl. (PACIA)*, vol. 2, Dec. 2008, pp. 264–266.
- [190] D. Wen, G. Liu, Y. Zhou, and S. Jin, "The study of enhancing authenticity of cloth surface intersection processing," in *Proc. Pacific-Asia Workshop Comput. Intell. Ind. Appl. (PACIA)*, vol. 1, Dec. 2008, pp. 977–981.
- [191] S. Vorasitchai and S. Madarassmi, "Improvements on layout of garment patterns for efficient fabric consumption," in *Proc. IEEE Int. Symp. Circuits Syst.*, vol. 4, May 2003, pp. IV552–IV555.
- [192] J. Távora and H. Coelho, "A path planner for the cutting of nested irregular layouts," in *Proc. Portuguese Conf. Artif. Intell.*, in Lecture Notes in Computer Science, vol. 390, 1989, pp. 246–256.
- [193] B. Kalra, K. Srivastava, and M. Prateek, "Computer vision based personalized clothing assistance system: A proposed model," in *Proc. 2nd Int. Conf. Next Gener. Comput. Technol. (NGCT)*, Oct. 2016, pp. 341–346.
- [194] M. T. Habib, S. B. Shuvo, M. S. Uddin, and F. Ahmed, "Automated textile defect classification by Bayesian classifier based on statistical features," in *Proc. Int. Workshop Comput. Intell. (IWCI)*, Dec. 2016, pp. 101–105.
- [195] Y.-I. Ha, S. Kwon, M. Cha, and J. Joo, "Fashion conversation data on instagram," in *Proc. 11th Int. Conf. Web Social Media (ICWSM)*, 2017, pp. 418–427.
- [196] D. Siegmund, T. Samartzidis, B. Fu, A. Braun, and A. Kuijper, "Fiber defect detection of inhomogeneous voluminous textiles," in *Proc. Mex. Conf. Pattern Recognit.*, 2017, pp. 278–287.
- [197] S. Jiang, M. Shao, C. Jia, and Y. Fu, "Consensus style centralizing auto-encoder for weak style classification," in *Proc. AAAI*, 2016, pp. 1223–1229.
- [198] L. Fan and G. Jiang, "Optimized Gabor filter parameters for uniform texture flaw detection," in *Proc. IEEE Int. Conf. Intell. Syst. Knowl. Eng. (ISKE)*, Nov. 2010, pp. 173–176.
- [199] V. Carvalho, F. Soares, and R. Vasconcelos, "Artificial intelligence and image processing based techniques: A tool for yarns parameterization and fabrics prediction," in *Proc. IEEE Conf. Emerg. Technol. Factory Automat. (ETFA)*, Sep. 2009, pp. 1–4.
- [200] I. Amed, A. Balchandani, M. Beltrami, A. Berg, S. Hedrich, and F. Rölkens, "The state of fashion 2019: A year of awakening," McKinsey & Company, Nov. 2018. [Online]. Available: <https://www.mckinsey.com/industries/retail/our-insights/the-state-of-fashion-2019-a-year-of-awakening>
- [201] T. G. Dietterich, "Ensemble learning," *Handbook Brain Theory Neural Netw.*, vol. 2, pp. 110–125, Mar. 2002.
- [202] C. Zhang and Y. Ma, *Ensemble Machine Learning: Methods and Applications*. New York, NY, USA: Springer, 2012.
- [203] K. Ramasubramanian and A. Singh, "Deep learning using keras and tensorflow," in *Machine Learning With R*. Berkeley, CA, USA: Apress, 2019, pp. 667–688.
- [204] C. Szegedy, V. Vanhoucke, S. Ioffe, J. Shlens, and Z. Wojna, "Rethinking the inception architecture for computer vision," in *Proc. IEEE Conf. Comput. Vis. Pattern Recognit. (CVPR)*, Jun. 2016, pp. 2818–2826.
- [205] M. Boeckhout, G. A. Ziehluis, and A. L. Bredenoord, "The FAIR guiding principles for data stewardship: Fair enough?" *Eur. J. Hum. Genet.*, vol. 26, no. 7, pp. 931–936, May 2018.



**CHANDADEVI GIRI** received the M.Sc. degree in big data systems from National Research University, Moscow, in 2015. She is currently pursuing the Ph.D. degree through the framework of Erasmus Mundus Joint Doctorate Program-Sustainable Management and Design of Textile (SMDTex). She was a Senior Analyst for a multinational IT consultant company Cognizant Technology Solutions, Mumbai, for three years. She has been an Erasmus Scholar, since 2017. Her research was supported by the European Commission. She is affiliated with three universities including the ENSAIT, Lille University of Science and Technology, France, the University of Borås, Sweden, and Soochow University, China. Her current research interests include artificial intelligence, recommendation systems, big data analytics, predictive modeling, machine learning, and expert systems.



**SHEENAM JAIN** received the B.Tech. degree from BIET, India, in 2014, and the master's degree in fashion technology from NIFT, New Delhi, India, in 2016. She is currently pursuing the Ph.D. degree through the framework of Erasmus Mundus Joint Doctorate Program-Sustainable Management and Design of Textile (SMDTex) affiliated by three universities, including ENSAIT/GEMTEX, the Lille University of Science and Technology, France, the University of Borås, Sweden, and Soochow University, China. She is also an Erasmus Scholar. Her research was supported by the European Commission. Her current research interests include big data, artificial intelligence, recommendation systems, apparel industry, digitalization, and the mass customization of apparel.



**XIANYI ZENG** received the B.Eng. degree from Tsinghua University, Beijing, China, in 1986, and the Ph.D. degree from the Centre d'Automatique, Université des Sciences et Technologies de Lille, Villeneuve-d'Ascq, France, in 1992. He is currently a Professor with the Ecole Nationale Supérieure des Arts et Industries Textiles (ENSAIT), Roubaix, France, and the Director of the GEMTEX Laboratory, Roubaix. His research interests include intelligent decision support systems for fashion and material design, and the modeling and analysis of human perception and cognition on industrial products and their integration into virtual products.



**PASCAL BRUNIAUX** was born in Denain, France, in 1959. He received the Ph.D. degree in automatics from the Lille University of Science and Technology, Villeneuve d'Ascq, France, in 1988. In 2007, he became a Full Professor in computer, automatic, and signal processing. He has published over 17 scientific book chapters, one scientific book, and 42 papers in reviewed international journals. He has presented over 97 papers at international conferences. He is a Supervisor of 22 Ph.D. students. His research interests include the modeling of textile structures, the modeling and virtualization of the human being, the analysis of textile comfort and consumer well-being, supervised and unsupervised classification of 3-D morphologies, the analysis of the process of customization of clothing, and setting up the virtual tailor.



## Paper II

**Giri, C.,** Thomassey, S., & Zeng, X. (2019). Customer analytics in fashion retail industry. In *Functional Textiles and Clothing* (pp. 349-361). Springer, Singapore

Available at: [https://link.springer.com/chapter/10.1007%2F978-981-13-7721-1\\_27](https://link.springer.com/chapter/10.1007%2F978-981-13-7721-1_27)

DOI: [https://doi.org/10.1007/978-981-13-7721-1\\_27](https://doi.org/10.1007/978-981-13-7721-1_27)



# Customer Analytics in Fashion Retail Industry



Chandadevi Giri, Sebastien Thomassey and Xianyi Zeng

**Abstract** This paper aims to give an overview of customer analytics in fashion retail industry in the era of big data. Fashion retail industry has been facing significant challenges since last few years due to rapidly varying customer demands. Nowadays, customers are much more informed and connected because of social media and other channels on the Internet. They demand more personalized services, and perception is not sufficient to understand our customers. Therefore, we need data to understand our customers and meet their expectation. We will discuss how customer analytics can create value in fashion retail industry, strategies and methodology to examine the consumer data. Employing and investing in these methods and technologies, industry will benefit from improved revenues, improve in sales, higher customer retention rates and thereby it will sustain in the uncertain markets. Segments are created using recency value of the customers, and their future behavior is predicted using transition matrix.

**Keywords** Customer analytics · Big data · Segmentation · Consumer behavior · Fashion retail industry

## 1 Introduction

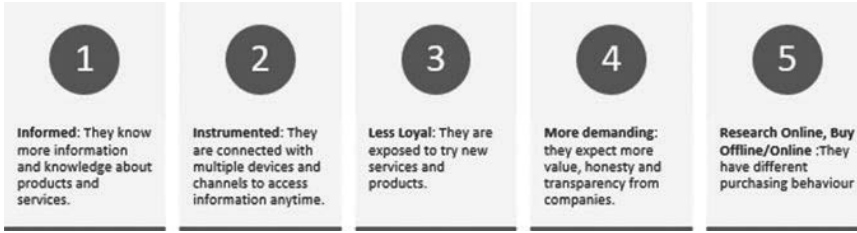
The present study is part of a Ph.D. project in sustainable design and management of textiles, focusing on e-commerce and consumer analytics to understand consumer behavior, and it is necessary to understand data apart from the perception of the consumer. The contemporary fashion industry is challenging, and integrating customer analytics with its business models will enable it to achieve business goals and high

---

C. Giri (✉) · S. Thomassey · X. Zeng  
Gemtex, Ensait, 2 Allée Louise et Victor Champier, 59056 Roubaix, France  
e-mail: chandadevi.giri@ensait.fr

C. Giri  
University of Borås, 501 90 Borås, Sweden

Soochow, College of Textile and Clothing Engineering, Suzhou 21506, China



**Fig. 1** Features of new customers

revenue growth. To illustrate how customer analytics can be effective, we have created customer segments on fashion apparel retail data to study customer behavior and the revenue generated by them. This will be discussed in detail in Sect. 5, customer data analysis.

Understanding consumer preferences has been a major challenge for fashion retailers. As per the market research of Statista [1], the revenue generated in the year 2017 from fashion globally is accounted for US\$406,476 million globally and it is estimated to increase by 11.6% by 2022. In the fashion industry, the largest market segment is “Clothing” that contributes a market volume of US\$272,599 million in 2017. Revenue generated in China itself contributed around US\$164,219 million in 2017. Also, it is estimated that large no. of users will be buying fashion products online by 2022, and the clothing industry will have the maximum market share. Currently, China, the USA, and UK are the major players in the fashion industry, and they are expected to grow in the future. It is important for the retailers to invest in big data analytics to understand customer preferences in real time. Big data provides an opportunity to understand customers in more precise way, and this leads to the emergence of new analytics area which is termed as “customer analytics.” “Customer analytics refers to the collection, management, analysis, and strategic leverage of a firm’s granular data about the behavior(s) of its customers” [2]. A new type of consumers has evolved. The evolution of customer behavior can be seen in Fig. 1. Therefore, it is essential to study the customer insights from the generated data to understand the fashion market trends.

## 2 Research Problem and Goal

The premise of the research problem emerges from the fact that fashion retailers are continuously facing many challenges in terms of predicting customer behavior in real time and adopting new strategies to fulfill customer demands. In line with this

premise, the goal of this research was to study the following:

1. How to classify customers based on their purchasing behavior evolving over time?
2. How to predict future customer movement and revenue generated by them?

The goal of the customer analytics is to address the above problems. It is important to understand whether consumer is creating value to the business, or whether they are satisfied with the services provided by the retail company, and their preferences in order to take appropriate actions to improve the services and products. Customer analytics can help retail companies to retain their customer base, increase revenue growth, and to predict consumer behavior, and eventually, the objective of value creation by each consumer could be achieved. Therefore, it is imperative to segment different group of the consumers as per their preferences and the value they create.

### 3 Customer Analytics

This section presents an overview of customer analytics strategy, scope, and methodology as shown in Fig. 2.

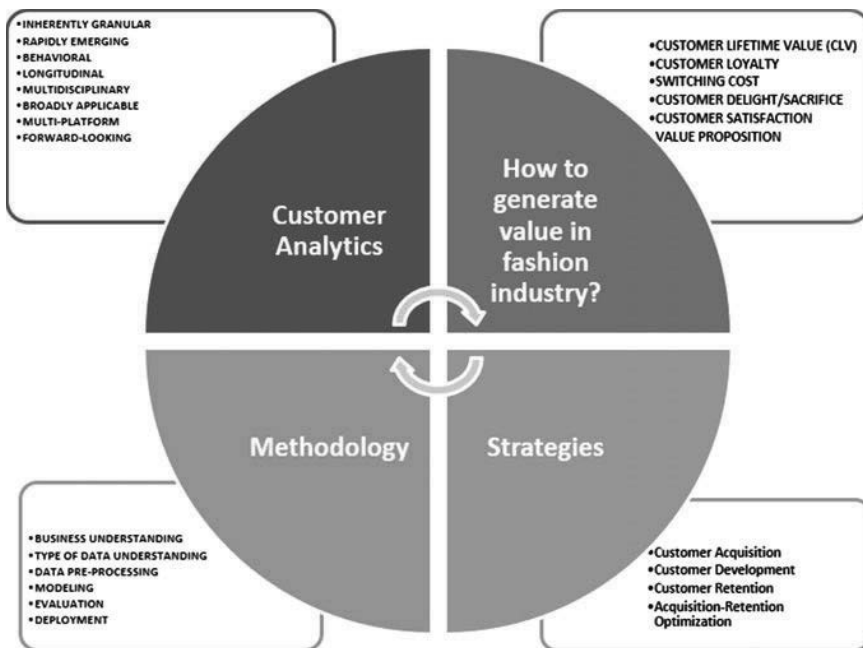


Fig. 2 Customer analytics, strategy, scope, and methodology

### ***3.1 Feature of Customer Analytics***

Company's business goals can define how customer analytics can be incorporated into their business intelligence. It could be focused on the prediction of customer behavior at an individual level without considering any other information, or it could be used by merging information from different systems to analyze, or it could be based on behavioral or longitudinal social network analysis. In other words, it depends on the business goals of the companies and which problem they are trying to address.

### ***3.2 Value Generation in Fashion Retail Industry***

Fashion retail industry is complicated, and it is quite difficult to understand the consumer choices toward the product. We can generate a high business value by identifying consumer lifetime value which can in turn help the fashion industry to achieve their profit goals. Customer Life Time Value (CLV) is the predicted value that businesses will derive from their entire relationship with their customers [3]. Well-known machine learning (ML) and probabilistic models such as Bayesian Inferences, Moving Averages, Regressions, and Pareto/NBD (Negative Binomial Distribution) can be used to predict CLV [4]. Customer segmentation can be done considering the demographic, geographic, behavioral factors using K-means clustering [5]. Association analysis can also be used for building recommendation system [6].

### ***3.3 Strategy***

All industries are consumer oriented, and consumers are crucial for their business success. The main strategy for fashion industry is to expand their customer base. As mentioned in the introduction, the dynamic nature of customers' buying pattern drives retailers to improve their strategy for customer acquisition and retention.

### ***3.4 Customer Analytics Methodologies***

Methodology for customer analytics often depends on the business problems that we are trying to address. Fashion social network data can be collected from tweets, boards, blogs to understand the hot topics, current trends, brands, events, criticism, etc. in fashion industry [7]. Internal sources are ERP, CRM, fashion e-commerce, etc., and external sources are cookies, plug-ins Adobe flash, etc. [8]. Data can be pre-processed to get into structured form; then, analytical method such as descriptive and predictive can be applied, and models can be evaluated. According to PwC and



**Table 1** Customer analytics with ML methods and statistics

Analysis	ML methods and statistics
Future profitability	Neural networks
CLV	Statistics
Potential CLV	Multi-regression
CLV profiling	Supervised clustering
Churn	Decision trees/Classification
RFM profiling	Decision trees
Acquisition modeling	Neural networks
Response analysis and modeling	Neural networks
Response index	Statistics
ROI	Structured procedures
Campaigns	Regression and structured procedures

SAP retailer survey [9], 39% of retailers ranked “Ability to turn customer data into intelligent and actionable insight” one of their greatest challenges. There is a huge gap between the big data and fashion retail industry. Retailers are more concerned with the data collection. After collecting data, they do not know what to do with such a huge and highly complex customer data. There is a lack of systems for tracking their minute-wise inventory. Therefore, understanding business, data, and customers is very important but retailers have to invest on analytics for creating valuable insights from the chunks of data that could benefit the industry to improve their products and services. Fashion industry is in a greater need to use advanced business intelligence tools, data analytics platforms, big data tools for capturing and processing data. Table 1 lists different machine learning method which can be used for customer analytics to predict profitability, life cycle of the consumers, loyalty, and campaigns.

## 4 Customer Data Analysis

For this study, we used dataset from the apparel industry, which spans from 2015-10-01 to 2016-12-01. To create segments, we used “Recency” value of each customer as the main indicator. Segments are created based on the recency value for six months, and each segment is further classified as “Inactive,” “Less Active,” “Active,” “Highly Active,” and “New Customers.” Given the 14-month time span of our dataset, we created two segments: one is 01-06-2016 to 01-12-2016, and the other is from 01-12-2015 to 01-06-2016. Using these two segments, we computed the transition matrix which is used for predicting the behavior and movement of customers within the five categories of each segment in the next one year. Based on the no. of the customers in each segment, revenue and cumulative revenue are calculated and predicted for next one year. No discount rate is considered for this study, but full price only.

### 4.1 Methodology

Transition matrix is used to compute the likelihood of the future occurrence depending on the current state. Let us assume that fashion retailers have  $X_n$  customers where  $X$  represents the total number of customers in a given segment at current state  $n$ ,  $X \in S$ , where  $S = \{\text{Inactive, Less Active, Highly Active, Active, New Customers}\}$ , then the probability  $P_{nm}$  of the customer in next state  $m$  will be given as

$$P_{nm} = P(X_{t+1} = m | X_t = n) \tag{1}$$

Thus, the transition matrix is conditioned on the present state, and the previous and the upcoming conditions are independent. In our study, we predicted the number of customers in each segment in an organization and their transition to the next state. Suppose the states are 1, 2, 3, ...  $r$ , where  $r$  represents the row, then the transition matrix for the different segments can be represented in matrix form as shown in Eq. 2. Thus, the probability of the customers in a segment in state  $m$  conditioned on state  $n$  can be represented by Eq. 3, and sum of all probabilities in a row will be equal to 1.

$$P(S) = \begin{bmatrix} S_{11} & S_{12} & S_{13} & \dots & S_{1r} \\ S_{21} & S_{22} & S_{23} & \dots & S_{2r} \\ \dots & \dots & \dots & \dots & \dots \\ \dots & \dots & \dots & \dots & \dots \\ S_{r1} & S_{r2} & S_{r3} & \dots & S_{rr} \end{bmatrix} \tag{2}$$

$$\sum_{m=1}^r P(S_{nm}) = \sum_{m=1}^r P(X_{t+1} = m | X_t = n) = 1 \tag{3}$$

### 4.2 Descriptive Analysis of the Data

Dataset is comprised of 5,770,844 customer transaction data with 1,020,923 distinct customers. For the customer analysis, we considered three variables: ‘‘Customer ID,’’ ‘‘Purchase Amount,’’ and ‘‘Date of Purchase,’’ see Table 2. The statistics show that average spending by each customer is 38.60 units, and maximum spending for each transaction is 199.90 units. The maximum time lapse with last transactions is 426.77 days (approx. 14 months).

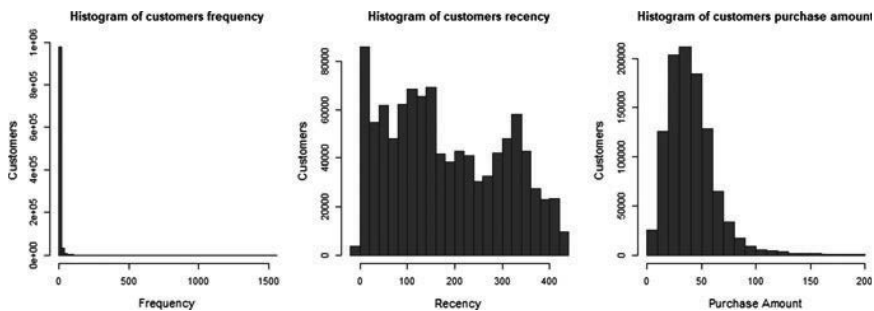
We have calculated three variables for customer data, i.e. recency, frequency, and monetary Value (RFM). Recency is the no. of the days lapsed between customers’ recent transaction date and the last transaction date. The bigger the recency value, the less active the customers are. Frequency is defined as the no. of the purchases made by the customers in a given period of time, and monetary value is the amount of money customers spent in each transaction. RFM is calculated for each customer. We can see from Table 3 that the average recency of a customer is 180 days, and

**Table 2** Summary statistics of customer data

Purchase_amount	Date_of_purchase	Days_since
Min.: -4.08	Min.: 2015-10-01	Min.: -0.2292
1st Qu.: 19.95	1st Qu.: 2016-01-10	1st Qu.: 120.7708
Med Median: 34.95	Median: 2016-04-22	Median: 222.7708
Mean: 38.60	Mean: 2016-04-25	Mean: 219.4961
3rd Qu.: 53.15	3rd Qu.: 2016-08-02	3rd Qu.: 325.7708
Max.: 199.90	Max.: 2016-12-01	Max.: 426.7708

**Table 3** Summary statistics of recency, frequency, and monetary (amount) for customer

Recency	Frequency	Amount
Min.: -0.2292	Min.: 1.00	Min.: 0.00
1st Qu.: 80.7708	1st Qu.: 2.00	1st Qu.: 25.82
Median: 155.7708	Median: 3.00	Median: 37.84
Mean: 180.4872	Mean: 5.65	Mean: 40.93
3rd Qu.: 288.7708	3rd Qu.: 6.00	3rd Qu.: 51.35
Max.: 426.7708	Max.: 1555.00	Max.: 199.90

**Fig. 3** Histogram of recency, frequency, and purchase amount

maximum recency is 426 days. The average spent amount per transaction is 40.93 units while the maximum spent amount for each transaction is 199.90 units. According to the frequency, customers had made at least one purchase and on an average five purchases. Histogram of RFM values is depicted in Fig. 3.

### 4.3 Segmentation Based on Recency

Two segments were created for the six-month interval based on recency values. In each segment, classes were assigned as Inactive, Less Active, Active, and Highly

**Table 4** Segmentation class based on recency value

Recency (days)	Customer class
>180	Inactive
≤180 and >90	Less active
≤90 and >50	Active
≤50	Highly active

**Table 5** Average recency, frequency, and purchase amount for each class within Segment 1

	Segment 1	Recency	First_purchase	Frequency	Amount
1	Inactive	296.47006	317.42518	3.286801	40.69390
2	Less active	131.86588	217.59175	5.252555	34.16796
3	Active	69.98107	195.52441	6.867103	44.75028
4	Highly active	21.90170	301.30038	15.947734	43.73544
5	New customers	23.62203	25.01742	3.271787	62.07015

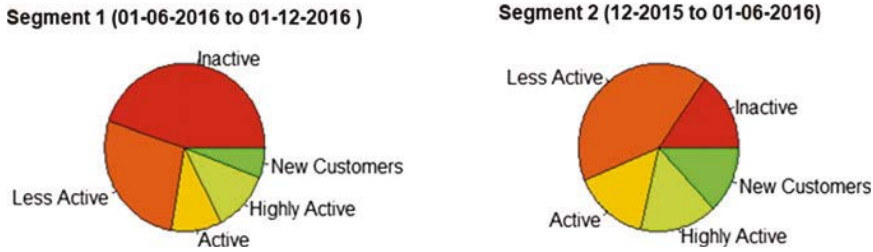
**Table 6** Average recency, frequency, and purchase amount for each class within Segment 2

	Segment 2	Recency	First_purchase	Frequency	Amount
1	Inactive	212.44278	214.87399	2.388128	52.89656
2	Less active	135.97081	150.84615	3.346653	33.82017
3	Active	68.74415	117.36808	4.974818	43.37889
4	Highly active	26.68699	175.74538	12.261362	39.70714
5	New customers	24.02350	25.32776	2.695572	43.46406

Active customers based on the recency values as shown in Table 4. The class “New Customer” is calculated as Customer in segment = “Highly Active” and first purchase  $\leq 50$ . This will help us to identify the customers who were not present in the Segment 2 but are newly added in the Segment 1. Highly Active customers are those whose recency is less than or equal to 50 days while inactive customers are those whose recency is more than 180 days. Recency, frequency, and average purchase amount for Segments 1 and 2 can be seen in Tables 5 and 6.

In Fig. 4, we have depicted how new customers who were absent in Segment 2 are now evolved in the Segment 1. As the logic behind our segmentation is to identify the similarity between the customers in different segments and understanding their behaviors after 6 months, we grouped them together according to their recency criteria for each segment. From the Segment 1, it is evident that new customers have been acquired that were not present in the Segment 2. Those who were new in Segment 2 transferred into another category “Inactive” and “Less Active” in Segment 1. Same can be observed for the highly active consumer, no. of highly active customers reduced in the segmentation 1, which are current customers.

The concept of transition matrix is employed to identify the probabilities of customers changing their segments. In other words, it is important to measure the likeli-



**Fig. 4** Number of customers of each class in Segment 1 and Segment 2

**Table 7** Transition matrix as per class of customers in Segment 1 and Segment 2

	Inactive less	Active	Active	Highly active	New customers
Inactive	83,293	12,661	4869	8986	0
Less active	208,564	48,854	15,118	24,620	0
Active	64,694	19,288	8322	16,226	0
Highly active	37,687	25,518	13,197	36,721	0
New customers	64,861	14,182	5430	9628	0

hood of inactive customers from the Segment 2 to moving to the Less Active, Active, and Highly Active classes of the Segment 1. Therefore, transition matrix enables us to find the number of customers in each class of the Segment 1. The probabilities from transition matrix are crucial in directing the future marketing campaigns and in targeting the potential customers. Transition matrix showing the probabilities of customers evolving in the next segment (Segment 1) is depicted in Fig. 4.

We have created transition matrix for the two segments to understand how the customers have changed their status from Segment 2 to Segment 1 in Table 7. The classes shown horizontally are Segment 1, and the one vertically is Segment 2. Based on vertical classes, we will see the movement of the customer in the horizontal classes. For example, in the first row, the customers who were inactive in segment 2, 83,293 customers remain inactive, 12,661 becomes less active, 8986 becomes active, and 4869 becomes highly active in Segment 1. So, we can say that very less inactive customers joined the groups: “Active” and “Highly Active.” Now, if we see what happened to the new customers in Segment 2, the fifth row of the transition table, 64,861 customers became inactive and 14,182 customers became less active, while only few remained in the groups: “Active” and “Highly Active.”

After we divided the rows in Table 7 by the sum of the customers for a given class, we get the transition probabilities as shown in Table 8. From Table 8, we can interpret that 75% customers who were inactive in Segment 2 will remain inactive in Segment 1, while about 8% inactive customers in Segment 2 will become highly active in Segment 1. For new customers in Segment 2, about 10% of them will remain

**Table 8** Transition matrix by probability

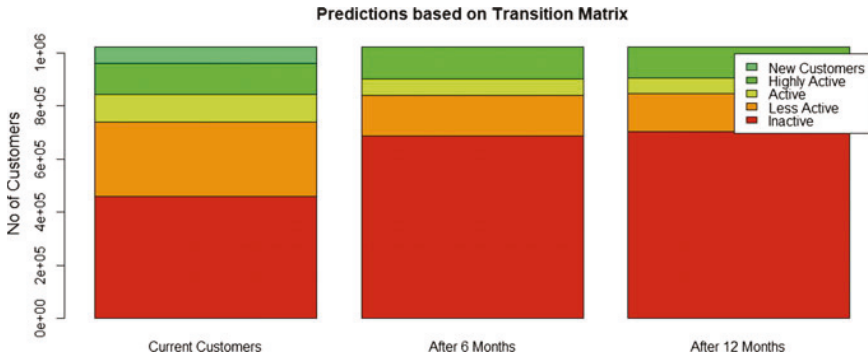
	Inactive	Less active	Active	Highly active	New customers
Inactive	0.75852617	0.11530020	0.04434063	0.08183300	0.00000000
Less active	0.70186703	0.16440523	0.05087563	0.08285210	0.00000000
Active	0.59609325	0.17772045	0.07667926	0.14950705	0.00000000
Highly active	0.33315064	0.22557747	0.11666063	0.32461126	0.00000000
New customers	0.68927004	0.15071041	0.05770396	0.10231560	0.00000000

highly active, 5% will remain active, and about 69% will remain inactive in Segment 1. Likewise, we can interpret the results for other classes too.

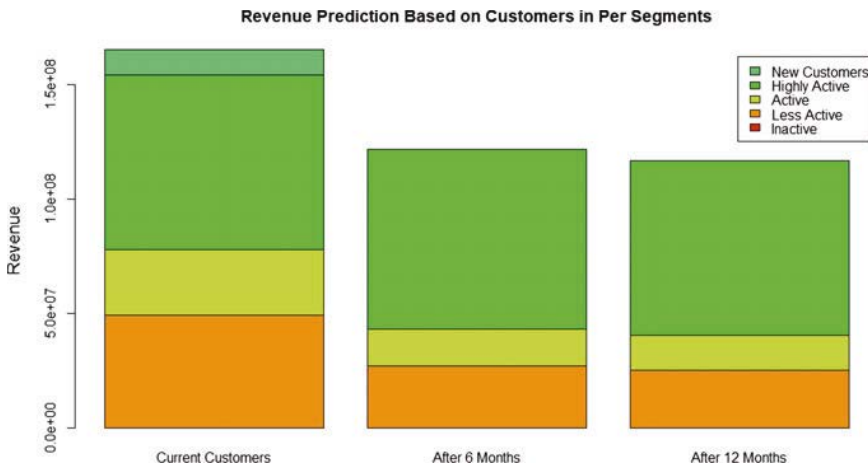
#### ***4.4 Prediction Based on Transition Matrix for the Next 6 Months and 12 Months***

It is often necessary to predict the no. of new customers in the next time period as they can potentially contribute to the revenue growth. Based on the prediction, managers design their marketing campaigns targeting new customers and the customers that are less likely to evolve, or in other words who are not exhibiting any movement to the active segments can be removed from the target group. Based on the results of transition matrix, we predicted the no. of the customers for the next 6 and 12 months and the revenue generated by them as shown in Figs. 5 and 6. Table 9 shows the predicted value of total number of customers in each class of the segments after 6 and 12 months. We can perceive that the number of the inactive users will be increasing in next one year, which means that new customers, active customers, and highly active customers from Segment 1 (current customers) will move to either Inactive or Less Active group of the same segment. As we have seen that the probability of “Inactive” customers shifting to “Active” group is quite less. New customers will be evolving in the future for 6 and 12 months of prediction. However, Fig. 5 shows the probability of the customer’s transition to other groups after 6 and 12 months. We can see from Table 9 that the no. of the customers in “Inactive” and “Less Active” class is increasing, while in “Active” group is decreasing. No. of “Highly Active” users remain approximately the same.

Intuitively, if the number of customers becoming active in the Segment 1 is higher, it means that the revenue generated by them will be higher. The future revenue or the revenue that will be generated in the next segment or in future segments can be predicted by looking at the classes of segments to which they belong. Table 10 shows revenue generated by each class in a segment currently and after 6 and 12 months. From Table 10, it is evident that the inactive customers cannot generate revenue as indicated by values “0,” whereas new customers are significant to high revenue generation. However, we could not predict the number of new customers and revenue



**Fig. 5** Prediction using transition matrix for no. of customers for each class for 6 and 12 months



**Fig. 6** Prediction of revenue generated by each class after 6 and 12 months

**Table 9** Number of customers in each class after 6 and 12 based on Segment 1

	Current customers	After 6 months	After 12 months
Inactive	459,099	687,299.59	704,646.68
Less active	279,777	153,078.56	142,277.96
Active	103,584	59,794.53	56,935.07
Highly active	118,107	120,750.32	117,063.29
New customers	60,356	0.00	0.00

**Table 10** Revenue generated by each class after 6 and 12 months based on Segment 1

	Current customers	After 6 months	After 12 months
Inactive	0	0	0
Less active	49,552,927	27,112,632	25,199,675
Active	26,369,141	16,376,269	15,593,133
Highly active	76,394,171	78,103,929	75,719,075
New customers	10,850,470	0	0

generated by them in the future time periods (6 and 12 months) because the Customer IDs of the new customers will be unique and different from the Customer IDs in the dataset we used for the study.

## 5 Results and Conclusion

Segment is created on 14-month customer data, and five different classes were defined based on the calculated recency value. As a result, we could identify five different customer classes purchase behavior over time, including revenues, accordingly. Transition matrix is used to predict the number of the customers in each segment class and revenue generated by them for the next 6 and 12 months. This kind of segmentation in a fashion apparel industry would help us to identify which segment of customers generates high value to the organization and how they can be retained for a long period. Besides, we can also analyze consumer behavior in detail by studying their purchasing behavior. As this segmentation is created by considering only one parameter “recency” value of the customer, it could be further improved by including other parameters: “Frequency” and “Monetary Value.” We will consider machine learning methods for predicting the customer’s behavior in our future work.

The fashion industry is dynamic and sensitive to quick changes in the customer behavior with the seasons and trends. Customer analytics will help the fashion industry to quickly respond to changing customer preferences. By applying customer analytics, we can analyze the buying behavior of the consumers and their preferences. Furthermore, it will also benefit supply chain and inventory management and it will be easy for the retailers to make decisions based on real-time tracking systems, reducing losses and helping the company to operate more environmentally friendly. Customer Analytics in fashion retail industry will help to customize the profiles of their consumer, enhance the personalized recommendation services, more loyalty programs, will give the opportunity to know their customers better than before and will help the business to create value from it. Thus, focusing on customer analytics in the era of big data, industries could be benefitted more than ever in the history.

**Acknowledgements** We are grateful to SMDTex—Sustainable Design and Management of Textiles, Erasmus Mundus Joint Doctoral commission for creating good research environment to take



out this research and their consistent support. We would like to show our deep gratitude to Mr. Giuseppe Craparotta, Senior Data Scientist, Evo Pricing for helping us with the data. Without him, it would have been difficult for us to execute this research work. We would also like to thank you Evo Pricing and their innovative team for their contribution to the research.

## References

1. Statista—The Statistics Portal for Market Data, Market Research and Market Studies. (n.d.). Retrieved from <https://www.statista.com> 15 Oct 2017
2. Home—Wharton CAI. (n.d.). Retrieved from <http://wcai.wharton.upenn.edu/> 16 Oct 2017
3. Strickland, P.J.: *Data Science and Analytics for Ordinary People*. S.I.: Lulu com (2015)
4. Wilkinson, J.W., Trinh, G., Lee, R., Brown, N.: Can the negative binomial distribution predict industrial purchases? *J. Bus. Ind. Market.* **31**(4), 543–552 (2016). <https://doi.org/10.1108/jbim-05-2014-0105>
5. Li, Z.: Research on customer segmentation in retailing based on clustering model. In: 2011 International Conference on Computer Science and Service System (CSSS) (2011). <https://doi.org/10.1109/csss.2011.5974496>
6. Tan, P., Steinbach, M., Kumar, V.: *Introduction to Data Mining*. Pearson, Nueva Delhi (India) (2016)
7. Data is Fashionable [Web log post]. (n.d.) (2016)
8. Curto, J., Braulio, N.: *Customer Analytics*. Editorial UOC (2015)
9. Verhoef, P., Kooge, E., Walk, N.: *Creating Value with Big Data Analytics*. Routledge (2016)



## Paper III

**Giri, C., Johansson, U.** (2021, September). Data-driven Business Understanding in the Fashion and Apparel Industry. *In The 18th International Conference on Modeling Decisions for Artificial Intelligence (MDAI 2021)*, Umeå, Sweden.

Available at: <http://www.mdai.cat/mdai2021/proceedings.MDAI2021.usb.pdf>



# Data-driven Business Understanding in the Fashion and Apparel Industry\*

Chandadevi Giri<sup>1</sup> and Ulf Johansson<sup>2</sup>

<sup>1</sup> Swedish School of Textiles, University of Borås, Sweden,  
`chandadevi.giri@hb.se`

<sup>2</sup> Dept. of Computing, Jönköping University, Sweden  
`ulf.johansson@ju.se`

**Abstract.** Data analytics is pervasive in retailing as a key tool to gain customer insights. Often, the data sets used are large, but also rich, i.e., they contain specific information, including demographic details, about individual customers. Typical usage of the analytics include personalized recommendations, churn prediction and estimating customer lifetime value. In this application paper, an investigation is carried out using a very large real-world data set from the fashion retailing industry, containing only limited information. Specifically, while the purchases can be connected to individual customers, there is no additional information available about the customers. With this in mind, the main purpose is to discover what the company can learn about their business and their customers as a group, based on the available data. The exploratory analysis uses data from four years, where each year has more than 1 million customers and 6 million transactions. Using traditional RFM (*Recency, Frequency and Monetary*) analysis, including looking at the transitions between different segments between two years, some interesting patterns can be observed. As an example, more than half of the customers are replaced each year. In a second experiment, the possibility to predict which of the customers that are the most likely to not make a purchase the next year is examined. Interestingly enough, while the two algorithms evaluated obtained very similar f-measures; the random forest had a substantially higher precision, while the gradient boosting showed clearly better recall. In the last experiment, targeting only the customers that have remained loyal for at least three years, rule sets describing patterns and trends that are strong indicators for churn or not are inspected and analyzed.

**Keywords:** RFM modeling · Churn prediction · Fashion and apparel.

## 1 Introduction

In retailing, customer retention is integral to the business growth. Actually, as suggested by Zhao [15], even a marginal increase in customer retention can lead

---

\* This research work is conducted under the framework of the SMDTex project (2017) funded by European Commission. We are grateful to the Evo Pricing company for providing the data for this research.

to as high as 90 percent increase in profit. Fashion retailers fiercely compete in the market to maximize customer retention by developing advanced marketing campaign strategies, and by offering unique products and services to their customers, thereby promoting customers' loyalty towards the brands, see e.g., [5].

Most retailing datasets contain personal, or at the very least demographic data, about the customers. In some datasets, however, the only information available is the purchases made, and a possibility to link all purchases made by the same customer. The overall purpose of this study is to investigate one such dataset from the fashion business, combining exploratory and predictive techniques in order to obtain customer insights. To the best of our knowledge, this is the first studied carried out on such a large dataset, covering only purchases, in the fashion and apparel retailing industry. For this analysis, the main method employed is RFM modeling. While RFM has been used frequently for exploratory analysis, it is rarely used as attributes in predictive modeling.

RFM modeling is one of the most popular methods for market segmentation to study consumer behavior. It quantifies customer behaviour in terms of Recency (R), Frequency (F), and Monetary value (M). Although customer churn management has been studied in the fashion industry, a majority of these studies focus on the segmentation aspects of customer data, providing only implicit knowledge about the customers' loyalty. One obvious reason for the lack of studies directly targeting churn, using predictive modeling, in the fashion industry is the problem of defining churning in this context. Since the most common definition is that a churning customer will not make a purchase in the coming period, which could be for instance six months or one year, the connection to RFM, in particular recency and frequency, is obvious.

## 2 Background

### 2.1 RFM analysis

RFM, i.e., *Recency, Frequency and Monetary* are numerical indicators describing the purchase history of consumers for a specified interval e.g., monthly, quarterly, or yearly [8]. Recency denotes the time elapsed since the last transaction made by the customer, frequency represents the number of transactions made by the customer and monetary indicates the total spending within the particular time interval. Customers with low recency, high frequency, and high monetary are of course the most active and profitable see e.g., [1]. While RFM analysis is widely applied for targeted market segmentation and for studying consumer behaviour, it is rarely combined with predictive modeling.

### 2.2 Churn prediction

Churn prediction constitutes the crucial element in Customer Relationship Management (CRM) strategy of any retail industry, where the main focus is retaining

loyal customers and providing them with the satisfactory services and products [6]. Hence, customer churn prediction is a crucial decision problem faced by the retail industries as it focuses on predicting the customers that are likely to churn in the near future. Identifying the factors that drive customer churn is one of the complex business problems in the F&A retail industry. Customers that are turning away from the fashion brands are a greater risk for the fashion retailers, and therefore, development of prediction models based on historical customer data and sales data is extremely important to identify the risky, loyal and valuable customers for the fashion retailers.

### 2.3 Related work

The application of data-driven techniques utilizing statistical or machine learning techniques for customer churn prediction has been investigated in several industries. One example is the Telecom industry, where churning means terminating a subscription for cellular services. Specifically, Sharma, Kumar and Panigrahi [9] proposed utilizing neural network models for predicting customer churn in that setting. In a similar case, but for landline telecommunications, Huang et al. [7] used a set of modeling techniques to predict customer churn, while identifying some key features that drive customer churn. In a study by Burez and Van Den Poel [3] random forests and logistic regression were applied to customer churn prediction in the financial service industry. Here, the the inherent class imbalance problem was handled by using different sampling schemes, specifically investigating the benefit of random and under-sampling.

The problem of customer segmentation based on loyalty or lifetime values has been heavily investigated in the existing literature. Looking specifically at the fashion retail industry, Dachyar et al [4] applied K-means to customer transaction data from an Indonesian brand to create customer segments based on customer lifetime value. By combining RFM scores with a customer lifetime value model, Yoseph and Heikkila [14] developed a modified regression model to produce customer segments using point-of-sales data from a medium-sized fashion retailer in Kuwait.

Many churn studies focusing on improved predictive performance exist. As an example, Zhao et al. [15] applied novel balanced and weighted random forests when predicting customer churn in the Chinese banking industry. Similarly, Tsai and Lu [10] proposed a hybrid prediction model combining standard neural networks and self-organizing maps to improve model accuracies for customer churn prediction. In yet another study, Xia et al. [13] used support vector machines to predict customer churn with very high accuracy. In a comparative study, by Vafeiadis et al. [11] it was found that Adaboost with support vector models as base classifiers were the most accurate. Finally, Verbecke et al. [12] obtained increased predictive performance by combining traditional CRM data with social media data on the customer level.

### 3 Method

Consumer transaction data were obtained from a European fashion apparel company during the four years 2013 - 2016. In total, the data set contains 28,726,819 transactions and 3,246,866 unique customers. Here it must be noted that each transaction contains only the customer ID, the date of the purchase and the total revenue of that transaction. Specifically, no information about the products purchased, nor about the customers, is available. However, since customer IDs of course are unique and remain with a specific customer, it is possible to follow the purchase history of customers over the years.

With the overall purpose of exploiting this very large but also rather limited data set in order to understand customer behavior, two different analyses were performed. In the first, exploratory, part, we start by looking at some customer data on the aggregated level, specifically identifying the proportion of all customers that remain active between two years, and how many new customers that make their first purchase. After that, the customers are, based on yearly data, divided into a RFM cube of  $3 \times 3$ , which is used to give a better understanding of different customer segments. In the last part of the exploratory investigation, transition matrices showing how customers move between RFM cells from one year to the next are created and analyzed.

In this study, the RFM analysis is carried out on a yearly basis, i.e., we could consider it as taking place on January 1, the next year. Consequently the recency (R) value for a customer is the number of days that has elapsed since the last purchase, the frequency value (F) is the number of purchases (transactions) the previous year, and the monetary value (M) is the total revenue generated by that customer the previous year. When creating the  $3 \times 3$  cube, each attribute was binned to contain an equal number of instances. As expected, customers that were inactive during the year are not part of the RFM cube. When describing the cells, the convention is to use the three numbers representing R, F and M in that order. In addition, the highest value (here three) represents the most active or valuable customers, i.e., the cell 333 contains the customers with the most recent purchases, the highest numbers of purchases and the largest total monetary value. It must be noted that since these values are highly correlated, the number of customers in each cell will vary substantially. As an example, we would expect very few customers in cells with  $F=3$  and  $M=1$ . In the second, predictive, analysis, random forest [2] and gradient boosting models, each consisting of 300 base classifiers, were trained to predict whether customers will remain loyal or churn, i.e., not make a purchase the next year. Here, RFM values and scores are used as input variables. In the first predictive experiment, predictions are made for the following year based on two years history, and it includes all customers that were active the last year. More specifically, the models are trained using the years 2013-2014 as inputs and 2015 as target, but the evaluation is done using 2014-2015 as inputs and 2016 as the target. In the experimentation, different feature sets combining the RFM scores (i.e., 1-3 for each factor), the RFM values (i.e., the actual number of purchases etc.) and in-cell relative values, were tried. The in-cell relative values are between 0-1, and represents where in



the cell that specific customer is for each factor. It is calculated as the (actual value-lowest value in cell) / (highest value in cell - lowest value in cell). As an example, if a customer’s monetary value is 50, and the lowest and highest monetary values in that cell are 20 and 120, respectively, the in-cell relative value would be  $(50-20)/(120-20)=0.3$ . All-in-all the following five different feature sets were evaluated:

- **Scores:** RFM scores
- **Values:** RFM values
- **Scores+:** RFM scores plus in-cell relative values
- **Values+:** RFM values plus in-cell relative values
- **All:** RFM scores, RFM values plus in-cell relative values

The second experiment, targets only the very important customers that have remained loyal for the first three years, producing churn predictions for the fourth and final year. In this setting, it should be noted that the churners represent the minority class, so undersampling was used producing a balanced training set. In addition, since actionable and interpretable rules were prioritized, the feature set used contained only F-scores and M-scores from all three previous years, plus the R score from the previous year. Furthermore, an interpretable model (a rule set) was evaluated together with the opaque models, with the purpose of finding valuable (individual) rules. In this experiment, a 70/30 randomized split into training and test set was used. The number of test instances is 124083.

## 4 Results

Table 1 presents the descriptive analytics of the four year data in which total transactions, unique customers, new customers, and churning customers are shown. Interestingly, almost more than 50 % of the total customers in each year are churning, which implies the low degree of customer loyalty. However, the churn is offset by the addition of approximately the similar number of new customers.

**Table 1.** Yearly data

Year	Transactions	Customers	New	Churning
2013	6420939	1538104	N/A	868975
2014	7424685	1519641	850512	804059
2015	7444853	1390248	674666	730719
2016	6916972	1242253	582724	N/A

Figure 1 shows the customer distribution as per the RFM score value. Out of 27 bins represented by 27 bars in each RFM score column for each year, the bins: 233, 332, 323 and 333 represent the highly active and valuable customers.

Here, the picture is very polarized; while two of the three most populated bins, i.e., 111 and 211, contain low value customers, the bin 333 also holds a large number of customers.

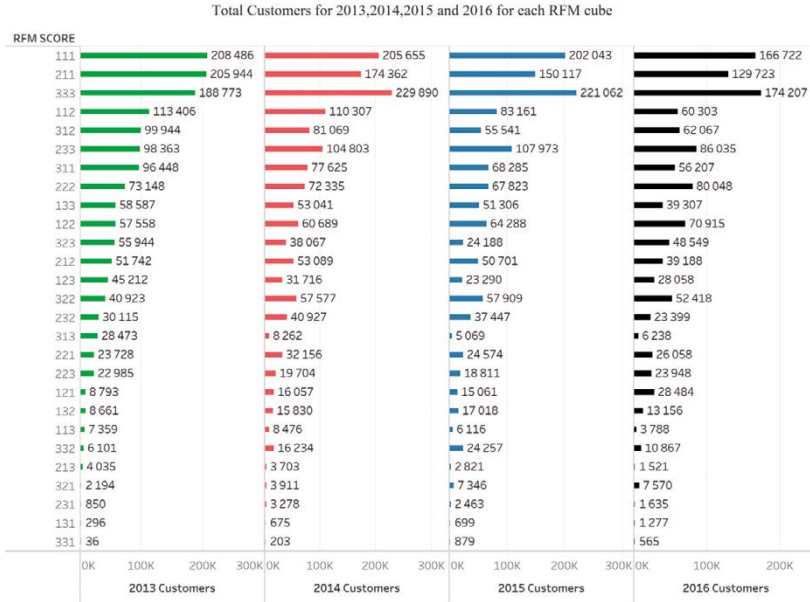


Fig. 1. Customer distribution over cells

Figure 2 illustrates the average frequency value for each bin. Here we see that the purchasing frequency of customers vary from single orders to more than 10 per year. Customers in the bins with the largest F-values typically have more than 10 purchases per year, making them very active. The fact that some of the customers make so many orders should be reassuring for the company since it is a clear indication that satisfied customers will return for more purchases.

Figure 3 represents the mean monetary values of the customers in each cell. As expected, Bin 333 represents the highest monetary value, but it should be noted that there are a substantial amount of very valuable customers (i.e., high M value) that still have relatively low F-values in bins 113, 123, 213, 223. These customers apparently make fewer but often larger purchases. Obviously, the two groups of valuable customers (frequent and infrequent purchasers) should be treated differently in the company's CRM. Figure 4 illustrates the proportion of old and new customers in each RFM cell. The general pattern is that new customers actually constitute a majority in most cells. The obvious exception is 333 where more than two customers out of three remain from last year. An interesting group of cells is 113, 213, 123 and 223, i.e., highly valuable customers

with not that many purchases; in these cells almost half of the customers are actually new, showing that some new customers make rather few but large orders.

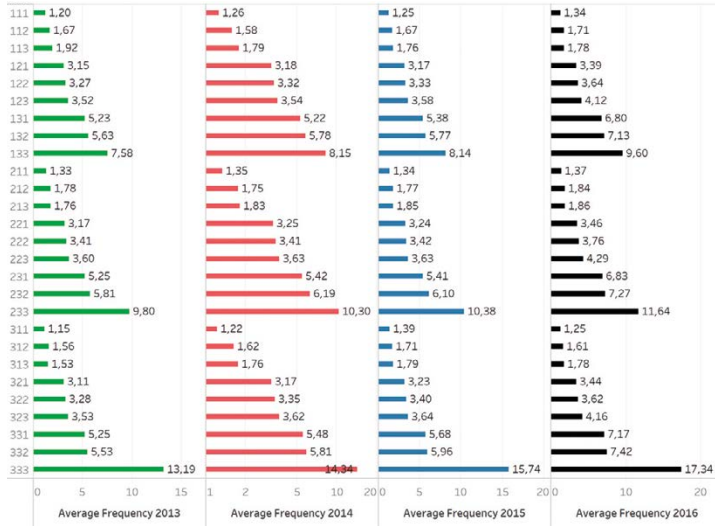


Fig. 2. Mean frequencies values in cells

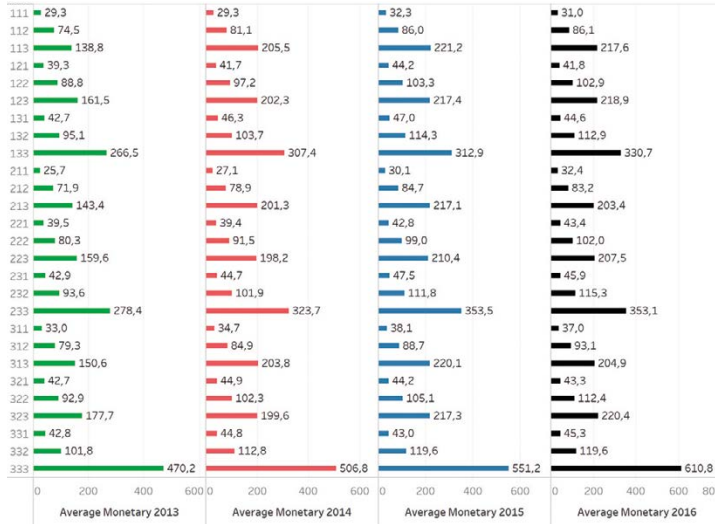


Fig. 3. Mean monetary values (€) in cells



Fig. 4. New customers

Figure 5 in a similar way depicts the proportion of churners and loyal customers in each bin.

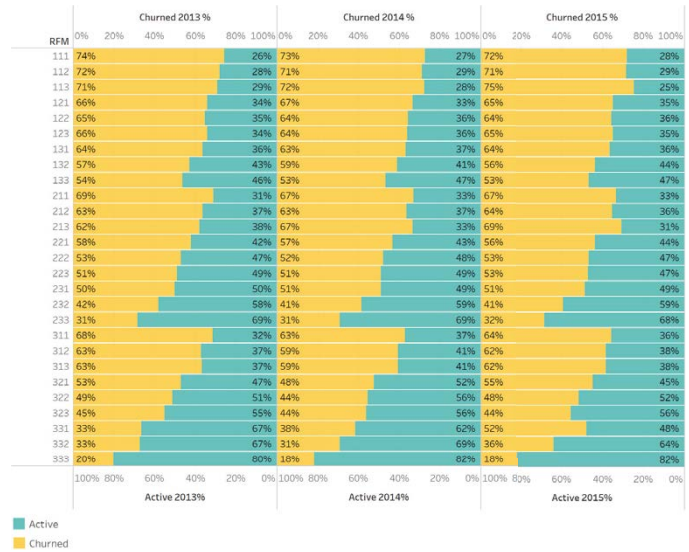


Fig. 5. Churners

This is, of course, a very important view for the company, showing for instance that, on the very general level, valuable customers tend to remain loyal. Specifically, cells representing high F and M values (i.e., bins 233, 332 and 333) exhibit less customer churn and consequently, these customers have a tendency to remain active in the next year. With this in mind, it becomes imperative to target customers in cells like 232, 222, 322 and 223, where the churn rates are approximately 50% trying to move them to cells with highest retention rates.



Fig. 6. Largest flows of customers between years

Figure 6 shows the largest flow of remaining customers between the RFM cells. Here it must be noted that while the results are aggregated over all years, the variation between years is fairly small. As seen in the graph, five cells receive most of the loyal customers; 111, 211, 221, 232, 233 and 333. Starting with some specific numbers, it is of course good to see that 60% of the customers in 333 that stay loyal also remain in either that cell or 233, i.e., making frequent purchases resulting in high M-values. On another level, a pattern that the company must be aware of is the fact that customers making frequent purchases, but with low M-values (i.e., cells 131, 121, 231, 221, 331 and 321) generally stay in these cells. It would be clearly beneficial for the company if these customers could be moved to cells with higher M-values; both for the immediate profit and for the lower churn rates in these cells. Table 2 below presents the results from the first predictive experiment. The main result is that independent of the feature set used, the models obtain rather high F-measures, providing a reasonable balance

between precision and recall. It is interesting to note that the two algorithms utilize different ways to get very similar f-measures; the random forest shows a relatively high precision, while the gradient boosting has higher recall. Overall, this shows that it is indeed possible to predict churners, and that the simplest feature set, i.e., using only RFM scores, is sufficient for that purpose.

**Table 2.** Experiment 1 - predictive performance

Model	Features	Acc	Prec	Recall	AUC	F-measure
RF300	Scores	0.684	0.755	0.680	0.720	0.715
RF300	Values	0.683	0.755	0.678	0.730	0.715
RF300	Scores+	0.684	0.753	0.681	0.732	0.715
RF300	Values+	0.681	0.742	0.680	0.731	0.710
RF300	All	0.681	0.744	0.680	0.728	0.710
GB300	Scores	0.684	0.680	0.754	0.741	0.715
GB300	Values	0.685	0.685	0.742	0.744	0.713
GB300	Scores+	0.684	0.683	0.743	0.743	0.712
GB300	Values+	0.686	0.685	0.743	0.745	0.713
GB300	All	0.686	0.686	0.741	0.745	0.712

Table 3 below shows the overall results from the second predictive experiment, i.e., when looking only at customers that have remained loyal for three years. Again, the different techniques obtain similar predictive performance, although the opaque ensemble models have a small edge. From a high-level perspective, the precision is rather good, but the recall somewhat disappointing. While the predictive performance give a clear indication of the difficulty of the problem, it is not obvious how a retention strategy utilizing this should be devised.

**Table 3.** Experiment 2 - predictive performance

Model	Acc	Prec	Recall	AUC	F-measure
RF300	0.651	0.759	0.479	0.723	0.588
GB300	0.650	0.764	0.478	0.733	0.588
Rule set	0.649	0.751	0.477	0.729	0.583

With this in mind, a further analysis of the produced rule set was undertaken. While the rule set contains in total 454 rules, most rules apply to only a handful of test instances. In Table 4 and Table 5 below, we show five sample rules for the churn class and three for the loyal class. These rules both cover a large number of test instances and exhibit a fairly high accuracy. The obvious usage of these rules is twofold; since they together cover approximately 60% of the test data

set, they could of course be used for prediction, but they could also be used for inspection and analysis, showing strong indicators for whether customers are likely to churn or not.

Starting with the rules that identify likely churners, we see that frequency scores, in particular the last two years, are very important. The first two rules and the last given below, all have  $F=1$  and  $F_1=1$ , i.e., covering customers in the lowest bin for frequency the last two years. It must be noted that these rules also have a very high accuracy, further strengthening the importance of this pattern. Furthermore, the other two rules also include conditions related to low frequencies. Rule 4 is slightly different, covering customers in the medium  $F$  bin last year. Here, however, the  $R$  score is 1, so a possible interpretation is that these customers made a few purchases in the beginning of the year, but after that became inactive, which of course is consistent with a churning behavior.

**Table 4.** Experiment 2 - Rules for churners

Rule	Inst. cov.	Perc.	Acc.
$M = 1 \wedge R > 1 \wedge F_2 < 3 \wedge F_1 = 1 \wedge F = 1$	8280	6.7	0.746
$M = 1 \wedge R = 1 \wedge F_2 = 1 \wedge F_1 = 1 \wedge F = 1$	5497	4.4	0.781
$R = 2 \wedge F_2 = 1 \wedge F_1 = 2 \wedge F = 1$	5346	4.3	0.730
$R = 1 \wedge F_1 = 1 \wedge F = 2$	4390	3.4	0.712
$M > 1 \wedge R = 1 \wedge F_2 < 3 \wedge F_1 = 1 \wedge F = 1$	2629	2.1	0.762

Looking at the rules selecting customers that will remain loyal, there is actually one rule that while covering more than a quarter of all customers, is also extremely accurate. This rule says that customers in RFM cell 333 that also were in the  $F=3$  bin the year before, are very unlikely to churn.

**Table 5.** Experiment 2 - Rules for loyal customers

Rule	Inst. cov.	Perc.	Acc.
$M = 3 \wedge R = 3 \wedge F_1 = 3 \wedge F = 3$	32624	26.3	0.845
$M = 3 \wedge R = 2 \wedge F_2 = 3 \wedge F_1 = 3 \wedge F = 3$	9097	7.0	0.711
$M = 3 \wedge R = 3 \wedge F_1 = 2 \wedge F = 3$	5655	4.6	0.694

## 5 Concluding remarks

We have in this application paper demonstrated how a large data set, containing only purchase histories from the Fashion & Apparel industry, can be utilized in order to obtain a business understanding. Using a combination of traditional RFM analysis and predictive modeling, many interesting patterns were discovered and analyzed. A key insight is of course that while 50% of the customers

are replaced by new ones every year, the loyal customers tend to be most valuable. Consequently, churn prediction and prevention should be a prioritized goal. Here, trained predictive models were deemed to be sufficiently accurate, exhibiting a reasonable trade-off between precision and recall. When inducing rule sets predicting which of the customers that have remained loyal for three years that were the most likely to churn the following year, a few rules highlighting in particular few purchases and/or a decreasing number of purchases, were found to be key indicators.

## References

1. Ballings, M., Van den Poel, D.: Customer event history for churn prediction: How long is long enough? *Expert Systems with Applications* **39**(18), 13517–13522 (2012)
2. Breiman, L.: Random forests. *Machine Learning* **45**(1), 5–32 (2001)
3. Burez, J., Van den Poel, D.: Crm at a pay-tv company: Using analytical models to reduce customer attrition by targeted marketing for subscription services. *Expert Systems with Applications* **32**(2), 277–288 (2007)
4. Dachyar, M., Esperanca, F., Nurcahyo, R.: Loyalty improvement of indonesian local brand fashion customer based on customer lifetime value (clv) segmentation. In: *IOP Conference Series: Materials Science and Engineering*. vol. 598, p. 012116. IOP Publishing (2019)
5. Dahana, W.D., Miwa, Y., Morisada, M.: Linking lifestyle to customer lifetime value: An exploratory study in an online fashion retail market. *Journal of Business Research* **99**, 319–331 (2019)
6. Edelstein, H.: Building profitable customer relationships with data mining. In: *Customer Relationship Management*, pp. 339–351. Springer (2000)
7. Huang, B., Kechadi, M.T., Buckley, B.: Customer churn prediction in telecommunications. *Expert Systems with Applications* **39**(1), 1414–1425 (2012)
8. Reinartz, W., Kumar, V.: The mismanagement of customer loyalty. *Harvard business review* **80**(7), 86–94 (2002)
9. Sharma, A., Panigrahi, D.P.K.: Article: A neural network based approach for predicting customer churn in cellular network services. *International Journal of Computer Applications* **27**(11), 26–31 (August 2011), full text available
10. Tsai, C.F., Lu, Y.H.: Customer churn prediction by hybrid neural networks. *Expert Systems with Applications* **36**(10), 12547–12553 (2009)
11. Vafeiadis, T., Diamantaras, K.I., Sarigiannidis, G., Chatzisavvas, K.C.: A comparison of machine learning techniques for customer churn prediction. *Simulation Modelling Practice and Theory* **55**, 1–9 (2015)
12. Verbeke, W., Martens, D., Baesens, B.: Social network analysis for customer churn prediction. *Applied Soft Computing* **14**, 431–446 (2014)
13. Xia, G.e., Jin, W.d.: Model of customer churn prediction on support vector machine. *Systems Engineering-Theory & Practice* **28**(1), 71–77 (2008)
14. Yoseph, F., Heikkila, M.: Segmenting retail customers with an enhanced rfm and a hybrid regression/clustering method. In: *2018 International Conference on Machine Learning and Data Engineering (iCMLDE)*. pp. 108–116 (2018). <https://doi.org/10.1109/iCMLDE.2018.00029>
15. Zhao, Y., Li, B., Li, X., Liu, W., Ren, S.: Customer churn prediction using improved one-class support vector machine. In: *International Conference on Advanced Data Mining and Applications*. pp. 300–306. Springer (2005)



## Paper IV

**Giri, C.,** Thomassey, S., & Zeng, X. (2019). Exploitation of social network data for forecasting garment sales. *International Journal of Computational Intelligence Systems*, 12(2), 1423-1435.

Available at: <https://www.atlantis-press.com/journals/ijcis/125922606>

DOI: <https://doi.org/10.2991/ijcis.d.191109.001>



Special Issue

# Exploitation of Social Network Data for Forecasting Garment Sales

Chandadevi Giri<sup>1,2,3,4,\*</sup>, Sebastien Thomassey<sup>1</sup>, Xianyi Zeng<sup>1</sup><sup>1</sup>Laboratoire de Génie et Matériaux Textiles (GEMTEX), ENSAIT, F-59000 Lille, France<sup>2</sup>The Swedish School of Textiles, University of Borås, S-50190 Borås, Sweden<sup>3</sup>College of Textile and Clothing Engineering, Soochow University, Suzhou 215168, China<sup>4</sup>Automatique, Génie informatique, Traitement du Signal et des Images, Université Lille Nord de France, F-59000 Lille, France

## ARTICLE INFO

### Article History

Received 24 Apr 2019

Accepted 26 Sep 2019

### Keywords

Social Media Data

Forecasting

Naïve Bayes

Sentiment analysis

Fuzzy forecasting model

## ABSTRACT

Growing use of social media such as Twitter, Instagram, Facebook, etc., by consumers leads to the vast repository of consumer generated data. Collecting and exploiting these data has been a great challenge for clothing industry. This paper aims to study the impact of Twitter on garment sales. In this direction, we have collected tweets and sales data for one of the popular apparel brands for 6 months from April 2018 – September 2018. Lexicon Approach was used to classify Tweets by sentence using Naïve Bayes model applying enhanced version of Lexicon dictionary. Sentiments were extracted from consumer tweets, which was used to map the uncertainty in forecasting model. The results from this study indicate that there is a correlation between the apparel sales and consumer tweets for an apparel brand. “Social Media Based Forecasting (SMBF)” is designed which is a fuzzy time series forecasting model to forecast sales using historical sales data and social media data. SMBF was evaluated and its performance was compared with Exponential Forecasting (EF) model. SMBF model outperforms the EF model. The result from this study demonstrated that social media data helps to improve the forecasting of garment sales and this model could be easily integrated to any time series forecasting model.

© 2019 The Authors. Published by Atlantis Press SARL.

This is an open access article distributed under the CC BY-NC 4.0 license (<http://creativecommons.org/licenses/by-nc/4.0/>).

## 1. INTRODUCTION

Due to the advent of social media, people, in general, have, like never before, been able to express and share their opinions about their various shopping experiences. In a consumer and market research framework, knowing about consumer experiences is crucial for the companies to improve their products and services they provide. Companies are now in a fierce competition to collect information about the experiences and opinions of their customers from social media as this information provides them with the valuable insights, which can be used to improve the satisfaction of their customers. Traditionally, companies relied on sales forecasting models to estimate their future product sales, which in turn help them in the strategic and effective management of their business operations. However, embedding consumer opinions into sales forecast models and studying their influence on the sales is increasingly becoming a new research trend [1]. There is a shift from the study of effects of marketing and advertisement on product sales to the study of social media influence. Social media is considered to be the digital word-of-mouth platform, which significantly influences consumers shopping behavior and thus, the product sales [2]. Many studies have explored the relationship between consumer sentiments on social media and the product sales [3]. Fashion industry is

grappling with numerous challenges of understanding their customers' opinions and adjusting their business models accordingly [4]. The recent trend of fast fashion is evident to the fact that consumers are rapidly evolving in terms of their changing preferences for new products, designs and features, and also the short product life cycle. Many industries such as television, telecommunication, banks, etc., tap social media to extract customer opinions, and adjust their operations based on the insights derived from social media [5]. However, there is a dearth of research on models in the fashion industry that study the factors such as social media and its effects on product sales. Fashion industry faces operational challenges emanating from the uncertainty in product demand and consumer choices. In this context, traditional forecast models that take into account the effects of marketing campaigns and advertisements on product sales fail to capture the effects of social media. Given the need for real time sales forecast in order to be able to manage various logistics related problems such as inventory, fashion industry needs to consider studying the effects of social media on the sales.

Motivated by the aforementioned problem facing garment industry, this research, aims to forecast garment sales by exploiting twitter data and investigating its influence of twitter as a social media on garment sales. Focusing into line, this study presents a hybrid sales forecast model “Social Media Based Forecasting (SMBF)” which is a fuzzy time series forecasting model that maps

\* Corresponding author. Email: [chandadevi.giri@ensait.fr](mailto:chandadevi.giri@ensait.fr)

historical information of the sales and the “Impact of sentiments from social media” (Fuzzy Sentiment Impact on Sales [FSIS] model) on garment sales. This research uses real twitter data and historical sales data for an Italian fashion brand, and applied developed model to examine its forecast performance. We believe that this paper is the first endeavor to use sales forecast model on social media data in a fashion industry. In this study, we propose a sales forecast model that is combined with sentiment analysis of social media platform, and we investigate if the performance of garment sales forecast can be improved using this model. The results from this study contribute to addressing overarching research question of forecasting fashion product sales by exploiting social media data, and it can provide decision makers in fashion industry with significant practical insights to optimize their business strategies. Moreover, we show that given the high performance of the model presented in this paper, it is possible to enable sustainable garment production as it would prevent the risk of over production of garments.

Detailed discussion of the existing literature is explained in Section 2. Experimental design and methodology is discussed in Section 3. This is followed by experimental results in Section 4, illustrating the classification results, correlation results and the model performance. Sections 5 and 6 explain the results, discussion and conclusion.

## 2. LITERATURE REVIEW

This section discusses existing literature related to sales forecasting in a fashion industry, in general, and the use of social media data based sales forecasting for various products. Maria E. Nenni, in paper [6], presents an excellent review of different existing sales forecasting methods used in fashion industry, and explores their suitability to forecast demands and sales with different characteristics. It is highlighted that inaccuracy or low performance of forecast models could be the result of unavailability of historical data, short selling and high level of uncertainty related to consumer preferences and satisfaction. Limitations of traditional sales forecast models such as auto-regressive (AR) model, extreme learning machine (ELM) model and artificial neural network (ANN) model are highlighted in the work [7] by studying the performance of newly developed ELM model with adaptive metrics of input and comparing its accuracy with the traditional ones. In the study carried out by [8], more accurate models based on fuzzy logic, data mining and neural network are used to improve the precision of sales forecasting in a clothing industry. The first study of consumer opinions in terms of online ratings given by consumers is studied [9], wherein Bass model is developed to forecast movie sales on box office. This study was the first approach to investigate the relationship between online reviews as a proxy for spread of consumers’ word-of-mouth and its influence on sales forecast. In a study done by [10], a linear regression model is used in combination with sentiment classifier to forecast box office sales, and the study (concluded) concludes that sentiments from Twitter (highly influence) have high influence on movie sales on box office. To forecast fluctuations in Amazon’s daily products sales, natural language processing is used in [11] to extract sentiments from online product reviews, and the relationship between sentiments and the product sales (have) has been studied. In a similar study done by [12], AR sentiment aware model has been applied to forecast rev-

enue on box office by extracting sentiments from online blogs and their influence on sales was found significant. Another domain, such as earthquake forecasting, where twitter data is used to predict damages caused by earthquakes is studied in [13]. In this study, machine learning models such as naïve Bayes (NB) and support vector machine (SVM) are used to classify tweets related to earthquake incidents and further embedded with spatial smoothing and regression models to estimate the loss due to earthquake. Very interesting application of social media analytics is presented in the study by [14], where authors used deep learning methods to forecast disease outbreak from twitter data with high accuracy. Popularity prediction model based on deep learning is developed in [15], where they extracted social media data to predict overall public interest and social trend. Another interesting industry example of social media based prediction is presented in [16], where machine learning models are applied to derive sentiments from news and social media and used to predict stock prices. Study done by [17] presents a graph theory based convolutional neural network model to forecast propaganda in a political system and also to predict election outcomes in South Africa and Kenya by extracting sentiments from tweets. Due to the growing trend of fast fashion, product life cycle is shrinking rapidly, and therefore, traditional long-term sales forecast models are irrelevant and are not very useful. Short-term demand and replenishment forecasting in a fashion industry using deep learning methods based on real historical industry data is conducted in [18]. Moreover, the influence of consumer sentiments, derived from a big data stream of online customer reviews generated in an e-commerce industry, has been studied in [19] by applying efficient methods to visualize short term demand forecast. This study has emphasized that online consumer sentiments are the promising factor to forecast short-term demand at a product level. Consumer sentiments derived from social media platform such as twitter, are studied in [20] to gauge the current fashion market trend and to predict the consumer perception toward brands. Study in [21] provides a comprehensive discussion on the usefulness of social media data for the operational activities of the industries. It is observed that social media has proved to be a significant factor for understanding of (various social outcomes) consumer behavior and it has a high influence on today’s competitive digital market. In this respect, we aim to study the consumer sentiments from their tweets about fashion brand and to explore their influence on garment sales, which is the first study of its kind in a fashion industry domain.

## 3. EXPERIMENTAL DESIGN AND METHODS

This section presents the schema of all the steps involved in the SMBF model implementation. It is divided into sub-sections briefing different stages involved in the development of model.

Subsequently, research framework (adopted in) this study is explained in Section 3.1. In the following sections, the topics covered are as follows: data collection steps in Section 3.2, data preprocessing in Section 3.3, sentiment extraction from fashion brand tweets in Section 3.4, Exponential Forecast model explanation in Section 3.5, modeling uncertainty of tweet sentiments using fuzzy inference in Section 3.6, stepwise implementation of SMBF model in Section 3.7, and finally, forecasting performance measuring indicator is explained in Section 3.8.

### 3.1. Research Framework

Research outline is shown in Figure 1, which includes sales prediction exploiting social media data and sales data in the following steps:

- I. Data Collection—To accomplish the task, Twitter and sales data spanning six months were collected from Italian fashion brand.

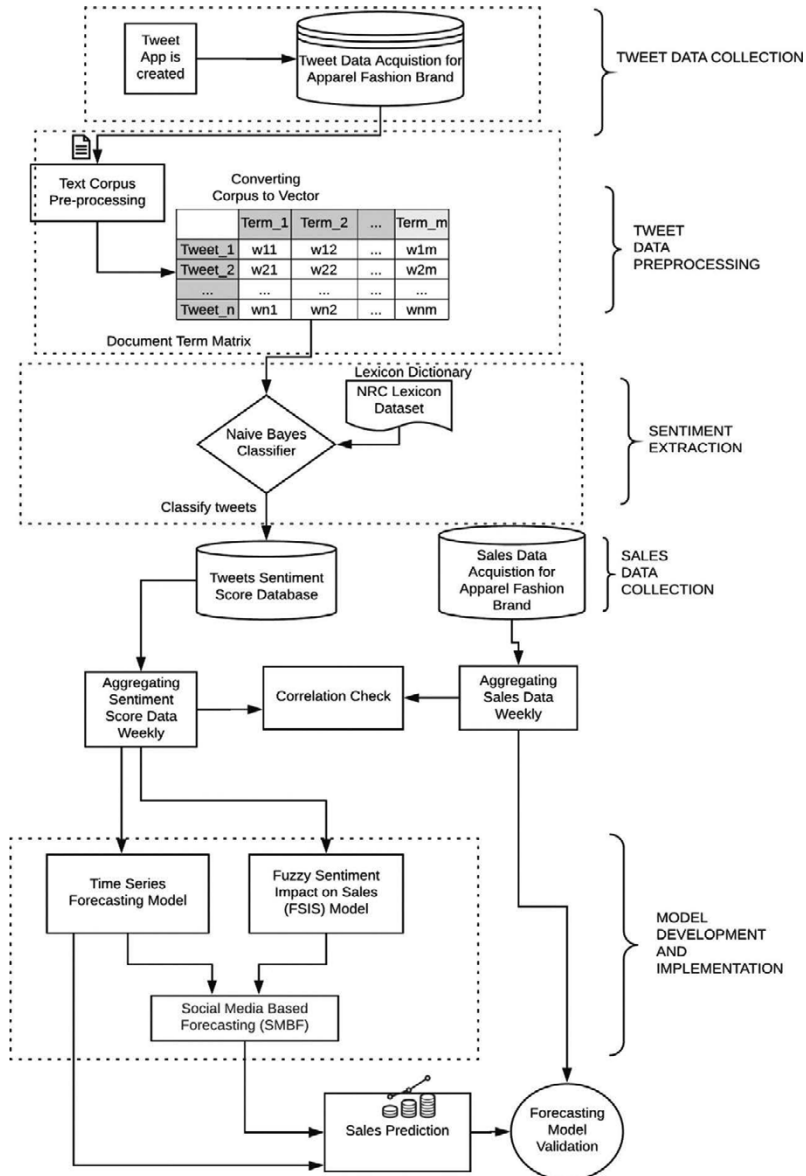


Figure 1 | Research framework.

**Table 1** | Attributes of Tweet data.

Serial No.	Attribute	Type	Description
1	Text	String	This is the text of the Tweet.
2	Favorited	Boolean	Specifies if the Tweet has been liked by the authenticated user
3	Favorite_count	Integer	Specifies the count of the Tweet that has been liked by the users
4	ReplyToSN	String	If the represented Tweet is a reply, this field will contain the screen name of the original Tweet's author
5	Created	String	This field states creation of Tweet at UTC time
6	Truncated	Boolean	If the original tweet exceeds the limit of 140 characters, the attribute "text" will be truncated and it will be represented by ellipsis "..."
7	ReplyToSID	String	If the specified Tweet is a reply, this field will contain the string representation of the original Tweet's ID.
8	Id	Integer	It represents a unique identifier of a tweet
9	ReplyToUID	Integer	If the tweet is a reply, this attribute will give the integer representation of the original Tweet's author ID.
10	StatusSource	String	It represents source of the Tweet created via Web, Android, iphone, etc.
11	ScreenName	String	This field gives the "screen name" of the user
12	RetweetCount	Integer	Number of times particular Tweet has been retweeted
13	Place	Places	User's location
14	Retweeted	Boolean	It indicates if the Tweet has been retweeted
15	Longitude/latitude	Coordinates	Indicates user's geographical location

- II. Data Preprocessing and Sentiment extraction—In this step, text mining technique is used for cleaning tweets. These cleaned tweets were then used for sentiment extraction using an enhanced NRC lexicon data and NB classifier. Classification results into three sentiment indicators Positive ( $S_p$ ), Negative ( $S_n$ ) and Neutral ( $S_{ne}$ ).
- III. Examine correlation—Extracted sentiments from earlier stage were aggregated weekly and the relationship between the sentiments of the current week and sales of the upcoming week was examined. This step was carried out to investigate relationship between tweets and sales.
- IV. SMBF Model—After exploring the relationship between tweet sentiments and sales, this model is designed to forecast garment sales based on historical sales and sentiments from social media. To attain this, FSIS model was built to identify the uncertainty of sentiments in terms of corrective coefficient ( $V_s$ ), which represents the sales variation factor with actual sales and this is implemented on Exponential Forecast (EF) to predict the final forecast from SMBF model.
- V. Performance Validation—The model performance was validated using forecasting measure, Mean Absolute Performance Error (MAPE).

### 3.2. Data Collection

Data used in this work was collected for one of the Fashion apparel brands, spanning time period of six months from April 2018 and September 2018. This study uses two data forms, first, Twitter data, which were collected using tweet API by creating tweet app to get authentication key [22]. Open source tool "Rstudio" and its package "TwitterR" were used to request API to Twitter using OAuth access [20]. Semantic keywords with hashtag "Brand name" were applied on tweet retrieval string. Collected tweets attributes is shown in the Table 1 and total number of tweets collected for six months interval was 838 as shown in Table 2.

**Table 2** | Collected Tweets.

Month	Collected Tweets
April	142
May	166
June	97
July	152
Aug	143
Sep	138
<b>Total</b>	<b>838</b>

In the next step, sales data for the same interval of time as the Twitter data was collected, which accounts for the daily sales of garments.

### 3.3. Twitter Data Preprocessing

Twitter data collected in previous step was in an unstructured form as it contained many duplicates, urls, punctuation, etc. It was necessary to clean and normalize tweet for analysis. As the first step of preprocessing, manual screening of tweets were performed to confirm whether or not it belongs to the fashion apparel brand that we targeted. While exploring the tweets, synonyms of selected brand name were found in some tweets. Moreover, some of the tweets were also related to politics, games and some other issues, and therefore such tweets were removed. Secondly, as only text is required for sentiment analysis, first attribute shown in the Table 1 was chosen. Text mining was applied to remove duplicates, whitespaces, hashtags, urls, stop words, retweet entities and numbers [20]. As a result, the total count of tweets reduced to 313, which is 37.35 % of the actual collected data.

### 3.4. Algorithm Formalization for Sentiment Extraction (Or Mining?) from Tweets

This section discusses the methodology applied for the classification of tweets. Lexical approach was used to extract sentiments from

tweets [23]. NRC lexicon dictionary [24] was used, which contains 5636 words labeled with “Positive” and “Negative.” For the tweet classification task, a probabilistic classifier “NB” algorithm [25] was used as it is the effective classification model for text classification. Moreover, it is highly scalable and it allocates the class labels of a finite set to feature vectors [23]. The algorithmic representation of NB is shown in Eq. (1).

$$P(X|Y) = \frac{P(Y|X) P(X)}{P(Y)} \tag{1}$$

where,

- $P(X|Y)$  is the posterior probability of a target class given attribute.
- $P(Y|X)$  is the likelihood, a probability of attribute given class.
- $P(X)$  is the prior probability of class.
- $P(Y)$  is the prior probability of predictor.

For the mathematical formulation of this research task Twitter dataset  $T$  is considered, which is a set of tweets  $t$  and can be represented as  $T = \{t_1, t_2, t_3 \dots \dots t_n\}$ . It was converted into corpus using text mining as explained in the Section 3.3. This corpus is a structured form of text, and is applicable for mathematical operations. Hence, each tweet could be considered as a set of words, i.e.,  $t = \{w_1, w_2, w_3 \dots \dots w_n\}$ . Similarly, Lexical labeled dictionary or dataset  $L$  [24] is composed of two attributes: words ( $W'$ ) and sentiment class ( $C$ ), and it is represented as  $L = \{W', C\}$ , where,  $W' = w'_1, w'_2, \dots w'_n$  is the set of words in lexicon dictionary and  $C = \{c_1, c_2, c_3 \dots \dots c_n\}$  is labeled sentiment class in lexicon dictionary. The sample of the lexicon dataset is shown in Table 3. Lexicon dictionary is labeled with two classes such as “Positive” and “Negative.”

For each given tweet  $t$  in Twitter dataset  $T$ , NB computes posterior probability of classes of sentiments,  $c \in C$ , where  $C = \{positive, negative\}$  and assigns the predicted class  $\hat{C}$ , which has the maximum posterior probability to the word of a given tweet. Accordingly, NB algorithm is represented in Eq. (1) can be rewritten as shown in Eq. (2).

$$\hat{C} = \underset{c \in C}{\operatorname{argmax}} P(c|t) = \underset{c \in C}{\operatorname{argmax}} \frac{P(t|c) P(c)}{P(t)} \tag{2}$$

Further, from Eq. (2), denominator is dropped as  $P(t)$  will remain same for all sentiment class and it will have an impact on “argmax” [26]. Therefore, it can be transformed as shown in Eq. (3).

$$\hat{C} = \underset{c \in C}{\operatorname{argmax}} P(t|c) P(c) = \underset{c \in C}{\operatorname{argmax}} P(w_1, w_2, w_3 \dots \dots w_n|c) P(c) \tag{3}$$

**Table 3** Example of the lexicon dictionary dataset.

Words ( $W'$ )	Sentiment Class ( $C$ )
Good	Positive
Sad	Negative
Honest	Positive

where,

$$P(w_1, w_2, w_3 \dots w_n|c) = P(w_1|c) \times P(w_2|c) \dots \times P(w_n|c)$$

Therefore, Eq. (3) can be rewritten as below

$$\hat{C} = \underset{c \in C}{\operatorname{argmax}} P(c) \prod_{w \in t} P(w|c) \tag{4}$$

This model follows the assumptions: position of the features (words) are not taken into account and the feature probabilities  $P(w_n|C_n)$  are independent of the given class  $c \in C$  [26]. “Prior” and “Likelihood” probability values are calculated for the occurrence of the most probable class and score is assigned.

Furthermore, formalization of the above equations was implemented for the ease of applying it in programming language illustrated in the Figure 2. It is assumed that the algorithm will calculate the maximum posterior probability for the occurrence of a word in a given tweet that is present in a lexical dataset. Therefore,  $P(c)$  is calculated, which is a prior probability of sentiment class and it can be calculated using the below Eq. (5) and likelihood using Eq. (6).

$$\text{prior } [c] = P(c) = \frac{Nc}{N} \tag{5}$$

where,

$N$  = Total words in lexicon dictionary

$Nc$  = counts of words from lexicon dictionary for each class  $c$

$$\begin{aligned} \text{Likelihood} &= \prod_{w \in t} P(w|c) = P(w_1, w_2, w_3 \dots \dots w_n|c) \tag{6} \\ &= \frac{\text{count}(w, c)}{\sum_{w' \in W} \text{count}(w', c)} \end{aligned}$$

where,

$(w, c)$  = number of occurrences of  $w$  in a Tweet  $t$

$(w', c)$  = occurrence of  $w'$  in lexicon dictionary  $L$  labelled with class  $c$

If the tweet does not have the identical word that exists in the lexical dictionary, then the likelihood in Eq. (6) will be zero as shown in Eq. (7) and this cannot be conditioned, therefore, Laplace smoothing [26] is introduced by adding 1 illustrated in Eq. (8).

$$\prod_{w \in t} P(w|c) = \frac{\text{count}(w, c)}{\sum_{w' \in W} \text{count}(w', c)} = 0 \tag{7}$$

$$\prod_{w \in t} P(w|c) = \frac{\text{count}(w, c) + 1}{\sum_{w' \in W} \text{count}(w', c) + 1} \tag{8}$$

While applying this equation in programming, a problematic condition could ascend when count of terms increases. As we know that each probability value of term  $P(w|c)$  will vary between 0 and 1, and if we multiply them, the product result will start approaching towards zero value. This leads to the problem in the programming language where representation of extremely small number is given as “double” or “long double” type that relies on standard data type “floating,” which is eventually fixed number and its precision reduces with miniscule numbers. The common way to combat this problem is to write the probabilities in “log-probabilities.” This transformation is considered due to number of reasons [27];

- The value of probability lies between 0 and 1, whereas log-probability can lie between  $-\infty$  and 0.
- The log-probabilities are easily comparable between smaller and larger numbers, that argmax does in probabilities for example  $m < n$ , then  $\ln(m) < \ln(n)$ .
- Arithmetic works well, i.e., logarithm of the multiplication of the variable is the addition of the logarithms, i.e.,

$$\ln(m.n) = \ln(m) + \ln(n)$$

Logarithmic representation of the Eq. (4) is represented below comprehensively with the final Eq. (9).

$$\hat{C} = \underset{c \in C}{\operatorname{argmax}} P(c) \prod_{w \in t} P(w|c)$$

$$\hat{C} = \ln \left[ \prod_{w \in t} P(w|c) \right] + \ln P(c)$$

$$\hat{C} = \left[ \sum_{w=1}^{w=n} \ln P(w|c) \right] + \ln P(c)$$

$$\hat{C} = \ln \left[ \frac{\operatorname{count}(w, c) + 1}{\sum_{w' \in W} \operatorname{count}(w', c) + 1} \right] + \left[ \frac{Nc}{N} \right] \quad (9)$$

Algorithm presented in Figure 2 assigns score to a word according to its frequency that matches with the words in a lexical dataset. The final score is the absolute of addition of logarithmic calculation of prior and likelihood as shown in Eq. (10). So, the algorithm assigns the class score for all words in a tweet. And, we are interested to know if the given tweet, which is as a whole sentence or statement,

is whether positive or negative. Therefore, a new metrics “Best fit” is used for categorizing the tweets as “Positive” or “Negative.”

$$\operatorname{score}[c] = \operatorname{abs}(\operatorname{prior}[c] + \operatorname{likelihood}[w, c]) \quad (10)$$

“Best fit” metrics for a given tweet is calculated as the ratio of the scores of positive and negative class as illustrated in Eq. (11). Classes for a given tweet will be assigned based of the conditions of score shown in Table 4. That is, if the value of “Best fit” score is greater than 1, then the tweet is positive; if it is less than 1 then it will be categorized as “Negative” tweet; and if the score is equal to 1, then it will be categorized as “Neutral” tweet.

$$\operatorname{Best\_Fit} = \frac{\operatorname{Total\ score\ of\ Positive\ words\ in\ a\ tweet}}{\operatorname{Total\ Score\ of\ Negative\ words\ in\ a\ tweet}} \quad (11)$$

### 3.5. Forecasting Model

This research study uses Exponential forecasting (EF) as a base model for forecasting the sales quantity. Exponential smoothing is a well-known forecasting method and was introduced by [28] and [29]. It can be used as an alternative to group of ARIMA models [30]. EF model equation is shown in the Eq. (12):

$$F_{t+1} = \alpha A_t + (1 - \alpha) F_t \quad (12)$$

where,

$F_{t+1}$  = Forecast for the time  $t + 1$

**Table 4** | Best fit criteria.

Best Fit	Tweet
Greater than 1	Positive
Less than 1	Negative
Equal to 1	Neutral

**PSEUDOCODE:**

**INPUT :**

1.  $T = \{t_1, t_2, t_3, \dots, t_n\}$  #Tweet Cleaned Dataset where,  $t = \{w_1, w_2, w_3 \dots \dots w_n\}$ , set of words in a tweet

2.  $L = \{W, C\}$  #Classes of Sentiments, where,  $C = \{c_1, c_2, c_3 \dots \dots c_n\}$  and  $W = w_1, w_2, \dots, w_n$

**PROCESS:**

#Calculate

1.  $\operatorname{prior}[c] = \ln \frac{Nc}{N}$

# N, Nc = Total words in dataset and counts of words from Dataset in each class c (NRC dataset)

2.  $\operatorname{likelihood}[w, c] = \ln \left[ \frac{\operatorname{count}(w, c) + 1}{\sum_{w \in W} \operatorname{count}(w, c) + 1} \right]$

3.  $\operatorname{score}[c] = \operatorname{abs}(\operatorname{prior}[c] + \operatorname{likelihood}[w, c])$

For each word  $w$  in  $L$  ( $L = \text{Lexicon of dataset } D$ ),  $\operatorname{count}(w, c) = \text{number of occurrences of } w \text{ in Tweets}$  and  $\operatorname{count}(w, c) = \text{number of occurrences of } w \text{ in } L$ .

**OUTPUT :** Predicted Class  $c$  for the words  $w$  in a tweet  $t$  with a score

**Figure 2** | Algorithm for assigning score to the word class.



$\alpha$  = Smoothing constant  $0 \leq \alpha \leq 1$

$A_t$  = Actual data

$F_t$  = Forecast data for time  $t$

Sales data of fashion garment industry was aggregated weekly and an EF was performed. “Corrective Coefficient ( $V_s$ )” was calculated, which represents the variation in sales as shown in Eq. (13).  $V_s$  is used as an output variable for FSIS, mapping the impact of tweets optimizing the predicting ability of the forecast model.

$$\text{Corrective Coefficient } (V_s) = \frac{\text{Actual Sales} - \text{EF Sales Forecast}}{\text{EF Sales Forecast}} \quad (13)$$

### 3.6. Modelling Tweet’s Sentiment Uncertainty for Predicting Sales Using Fuzzy Theory

This section discusses, in detail, the process involved in modelling tweet sentiments. This section is divided into two sub-sections briefing the details of model design and defines linguistic variables for analyzing the impact of Tweet sentiments. Section 3.6.1 outlines the architecture of FSIS model and section 3.6.2 explains the fuzzy linguistic variables for the FSIS model.

The uncertainty is represented by fuzzy sets in classical theory and each element in a set hold a membership degree [31–33]. Fuzzy set can be denoted as a pair  $(U, m)$ , where  $U$  is a set and  $m$  is a membership function of a fuzzy set, which assigns the value between 0 and 1 and it can be denoted as  $U \rightarrow [0, 1]$  [34]. Therefore, for a fuzzy set  $P = (U, m)$ , membership function  $m = \mu(A)$  maps the function between 0 to 1. For a finite set  $U = \{a_1, \dots, a_n\}$ , it can be represented as  $\{m(a_1)/a_1, \dots, m(a_n)/a_n\}$ . Let  $a \in U$ , then  $\alpha$  will not be considered in the fuzzy set  $(U, m)$  if  $m(a) = 0$ , will be fully considered if  $m(a) = 1$  and it will be partly considered if  $0 < m(a) < 1$  [35]. It follows the mathematical operations of complement, intersection and union. For the fuzzy set  $P$  and  $Q$ ,  $P, Q \subseteq U$ ,  $u \in U$ , complement operation can be represented as  $\mu_p(u) = 1 - \mu_p(u)$ , intersection operation can be represented as  $\mu_{P \cap Q}(u) = \min\{\mu_p(u), \mu_Q(u)\}$  and Union of two fuzzy sets can be shown as  $\mu_{P \cup Q}(u) = \max\{\mu_p(u), \mu_Q(u)\}$ . Fuzzy set has “Commutative,” “Associative,” “Distributive,” “Identity” and “Transitive” properties. Triangular membership [36] is popular to solve engineering optimization problems due to its ability to provide solution instantly. Degree of Membership function for a fuzzy set  $P$  can be illustrated in Fig. 3 and Eq. (14) defined by the lower and upper limit,  $a$  and  $b$  where,  $a < b < c$ .

$$\mu_{P(x)} = \begin{cases} 0 & x \leq a \\ \frac{x-a}{b-a} & a < x \leq b \\ \frac{c-x}{c-b} & b < x < c \\ 0 & x \geq c \end{cases} \quad (14)$$

#### 3.6.1. FSIS architecture

FSIS architecture is defined using classical “Mamdani” Fuzzy inference proposed by [32,37], the first control system base, fashioning it by if-then rules, which can be gained from human knowledge. This work was inspired by [31] and realized its application of fuzzy inference for complex decision mechanism. Mamdani inference maps the input attributes to the output using fuzzy logic inference which involves steps as follows;

- (a) First step involves defining the input and output linguistic variables followed by defining membership function for each linguistic variables
- (b) Fuzzification—In this step, input parameters are mapped to appropriate linguistic variables based on pre-defined membership function that can be viewed as fuzzy set. Membership functions are linked with the weight factors that regulate the influence for each defined rule.
- (c) Knowledge base—It has a database of fuzzy set and defined rules. “Facts” are symbolized through linguistic variables and reasoning is followed by defined if-then rules.
- (d) Inference Engine—Rule inference is managed by inference engine, where human knowledge can be combined without any difficulty using linguistic rules.
- (e) Defuzzification—Aggregation of all fuzzy sets to give the single output to draw a conclusion. Centroid method is used for obtaining the crisp number as shown in Eq. (15).

$$x^* = \frac{\int \mu_p(x) . x dx}{\int \mu_p(x) . dx} \quad (15)$$

Model design of FSIS involves above steps and this is explained in Section 3.6.2 with input parameters as sentiments and output parameter as corrective coefficient as explained in the Section 3.5, and corresponding defined membership functions.

#### 3.6.2. Defining universal and fuzzy set for sentiments for FSIS

For formalizing and modeling twitter data with garment sales performance, we used intelligent fuzzy technique as discussed in the Section 3.6.1. FSIS constitutes three fuzzy sentiment input variables  $I$  as  $S_p$ (Positive),  $S_n$ (Negative),  $S_{ne}$ (Neutral) and Output  $O$

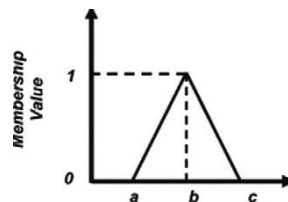


Figure 3 | Triangular membership.

as  $V_s$  (Corrective coefficient). Triangular membership function is assigned to input variables in the range between [0 1] and output variable in the range [-25 25]. Aggregated weekly sentiment data were used for all three input variables to feed into model. The key objective in modelling sentiments, viz.,  $S_p$  (Positive);  $S_n$  (Negative); and  $S_{ne}$  (Neutral) retrieved from tweets with fuzzy inference is to capture human emotions on social media weekly and to forecast the sales variation based on these consumer emotions. Human emotions are uncertain and they change from time to time. So, for example, maybe in one week, sentiment "Positive" could be high and it could lead to a positive effect on sales of upcoming week. Similarly, if it is "Negative," then it could lead to a negative impact on sales of upcoming week. The designs of the linguistic variables for sentiments are taken as "Low," "Medium" and "High." Thus,  $I$  is a set of sentiments, i.e.,  $Input(I) = \{S_p, S_n, S_{ne}\}$  and membership sets for  $I$  is defined in three linguistic variables as "Low," "Medium" and "High" for each input variables and can be written as  $S_p = \{Low, Medium, High\}$ ,  $S_n = \{Low, Medium, High\}$ ,  $S_{ne} = \{Low, Medium, High\}$ , and for output variable  $Output(O) = \{V_s\}$ , membership sets is defined as  $V_s = \{Low, Medium, High\}$ , where,  $V_s$  is a corrective coefficient and it is calculated as a difference between actual sales and EF sales as shown in Eq. (13).

Triangular membership function for the input variable in set  $I$  is defined between limit [0 1] and the output variable in the limit [-25 25], which is a corrective coefficient in percentage of sales and it will predict sales variation from -25% to 25%. (Table 5 describes) the membership variables and values for the  $I$  and  $O$  respectively are shown in Table 5. Figures 4 and 5 depict the triangular membership function for variables  $I$  and  $O$ .

Calculation of the membership function value for the Input variable " $S_p = \{Low, Medium, High\}$ " is illustrated below:

$$\mu(S_p(Low)) = \frac{0.5 - x}{0.5} \quad 0 \leq x \leq 0.5$$

Table 5 | Membership function for input variable I (Input) and O (Output).

Input/ Membership	Low	Medium	High
$(S_p)_w$	[0 0 0.5]	[0 0.5 1]	[0.5 1 1]
$(S_n)_w$	[0 0 0.5]	[0 0.5 1]	[0.5 1 1]
$(S_{ne})_w$	[0 0 0.5]	[0 0.5 1]	[0.5 1 1]

Output/ Membership	Low	Medium	High
$(V_s)_{w+1}$	[-25 -25 0]	[-25 0 25]	[0 25 25]

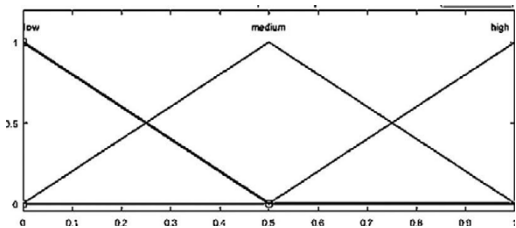


Figure 4 | Triangular membership for input variables ( $S_p, S_n, S_{ne}$ ).

$$\mu(S_p(Medium)) = \frac{x}{0.5} \quad 0 \leq x \leq 0.5 \quad \text{and} \quad \frac{1-x}{0.5} \quad 0.5 \leq x \leq 1$$

$$\mu(S_p(High)) = \frac{x-0.5}{0.5} \quad 0.5 \leq x \leq 1$$

Similarly, membership function values is calculated for  $S_n$  and  $S_{ne}$ . Besides, the membership function value for corrective coefficient " $V_s = \{Low, Medium, High\}$ " is calculated as below:

$$\mu(V_s(Low)) = \frac{0 - x}{0 - (-25)} \quad -25 \leq x \leq 0$$

$$\mu(S_p(Medium)) = \frac{x - (-25)}{0 - (-25)} \quad -25 \leq x \leq 0 \quad \text{and}$$

$$\frac{25 - x}{25 - 0} \quad 0 \leq x \leq 25$$

$$\mu(S_p(High)) = \frac{x - 25}{25 - 0} \quad 0 \leq x \leq 25$$

Considering the impacts of the weekly sentiments, fourteen conditional rules are defined by taking into account human knowledge as shown in Table 5. These rules are defined as

"Rule i : IF condition i THEN action I"

To demonstrate the functioning mechanism of rules consider first row of the Table 6, which states that if  $S_p = High$ ,  $S_n = low$  and  $S_{ne} = Medium$ , then  $V_s = High$ , i.e., if the aggregated weekly positive sentiment is high, negative sentiment is low, and the neutral sentiments is medium, then the sales performance will be high and its values will lie between 0 % to 25 % range, meaning there is an incremental variation in sales than the previous week. Similarly, expected sales variation decision will be taken by the model considering the inference of input variables and defined rules. The final score of the output, i.e., defuzzification of the fuzzy set and the crisp number is calculated by the centroid method using Eq. (15). Architecture of FSIS is depicted in Figure 6.

### 3.7. SMBF Model

SMBF model is a time series model, which takes into account the impact of sentiments on sales using FSIS. Parameters of the FSIS are optimized to improve the forecast ability of the model and it is implemented on the exponential time series forecasting and the formula its forecast estimate is given in Eq. (16), where  $V_s$  is corrective coefficient, calculated using Eq. (13) and  $EF$  is sales volume forecasted by the EF model as explained in section 3.5.

$$SMBF = (1 + V_s) \times EF \tag{16}$$

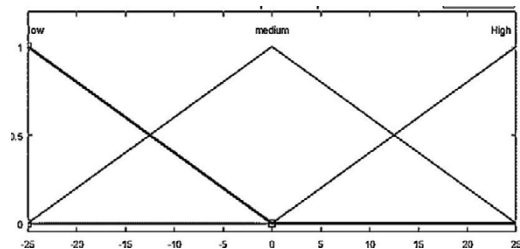


Figure 5 | Triangular membership for output variables ( $V_s$ ).

where,

SMBF = sales Forecast from fuzzy time series model

$V_s$  = corrective coefficient predicted by the FSIS model

EF = sales forecasted by Exponential Model

FSIS model parameters, i.e., “Input,” “Output” and “Rules” defined in Section 3.6.2 were optimized using the local optimization

method known as “pattern search method” [38], which is useful for the fast convergence. The model was tuned for 500 iterations, and as a result, model was optimized. Optimized parameters for “Input” and “Output” variables are shown in Table 7 and optimized rules are shown in the Table 8.

### 3.8. Validation Method and Performance Measure (MAPE)

Mean Absolute Percentage Error (MAPE) is a performance metrics for evaluating the performance of the forecasting models and it is denoted as in Eq. (17).

$$MAPE = \frac{100\%}{n} \sum_{t=1}^n \left| \frac{A_t - F_t}{A_t} \right| \quad (17)$$

where,

$A_t$  = actual sales at time  $t$

$F_t$  = forecasted sales at time  $t$

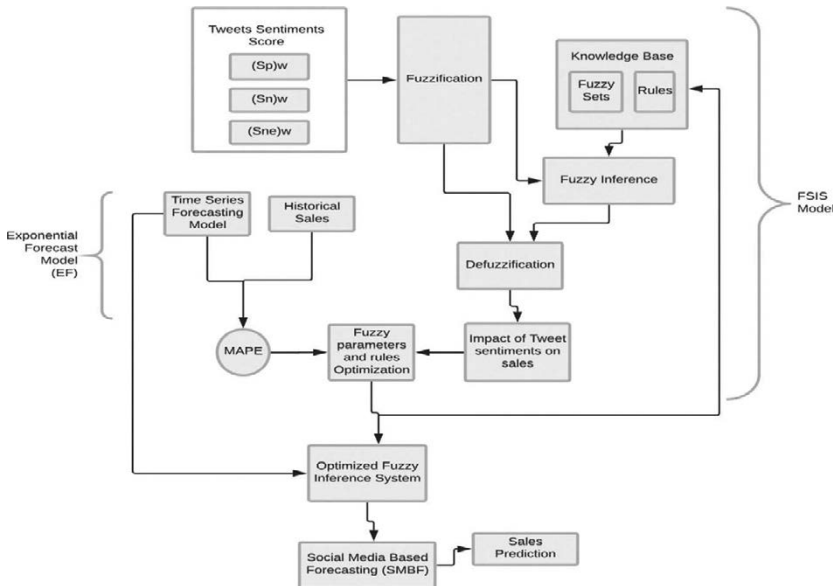
**Table 6** | Rules defined for Fuzzy Sentiment Impact on Sales (FSIS).

if	And	And	Then
$(S_p)_w$	$(S_n)_w$	$(S_{ne})_w$	$(V_s)_{w+1}$
High	Low	Medium	High
High	Medium	Low	High
Low	High	Medium	Low
Medium	High	Low	Medium
Low	High	Low	Low
Medium	Low	High	Medium
Low	Medium	High	Medium
Low	Low	Medium	High
Medium	Low	Medium	High
Low	Medium	Medium	Low
Medium	Medium	Low	Medium
Medium	Low	Low	Medium
Low	Low	High	Medium
High	Low	Low	High

**Table 7** | Optimized model parameter of FSIS.

Input/ Membership	Low	Medium	High
$(S_p)_w$	[0 0 0.5]	[0 0.875 1]	[0.75 1 1]
$(S_n)_w$	[0 0 0.5]	[0 1 1]	[0.125 0.75 1]
$(S_{ne})_w$	[0 0 0.5]	[0.125 0.375 1]	[1 1 1]
Output/ Membership	Low	Medium	High
$(V_s)_{w+1}$	[-25 -23.5 20.5]	[-15 -14.375 -14.25]	[22.5 25 25]

FSIS, Fuzzy Sentiment Impact on Sales.



**Figure 6** | Social Media Based Forecasting (SMBF) model.

**Table 8** Optimized rules defined for FSIS.

$(S_p)_w$	$(S_n)_w$	$(S_{ne})_w$	$(Vs)_{w+1}$
High	Medium	Medium	High
High	Medium	High	High
High	High	-	High
-	High	Medium	Low
High	High	High	Low
Medium	Low	High	Medium
High	High	Medium	Medium
High	Low	Medium	High
Medium	Low	Medium	High
Low	-	Low	Low
-	Medium	Low	Medium
High	High	Low	Medium
Low	Low	High	Medium
Low	High	Low	High

FSIS, Fuzzy Sentiment Impact on Sales.

### 4. EXPERIMENTAL RESULTS

This section presents the results of all the models discussed in previous section. This section is divided into three subsections: Section 4.1 summarizes the result of NB classifier, which is described in Section 3.4. It is used for extracting sentiments from the tweets and assigning sentiment score to a word in a tweet that matches with the lexical dictionary. Based on sentiment score, tweets are classified as “Positive,” “Negative” and “Neutral” as explained in Section 3.4. Further, in Section 4.2, computational mechanism of correlation between weekly aggregated Twitter sentiments for the current week and sales of the next week is discussed. Lastly, Section 4.3 evaluates the forecast performance of SMBF and EF models.

#### 4.1. Tweet Classification

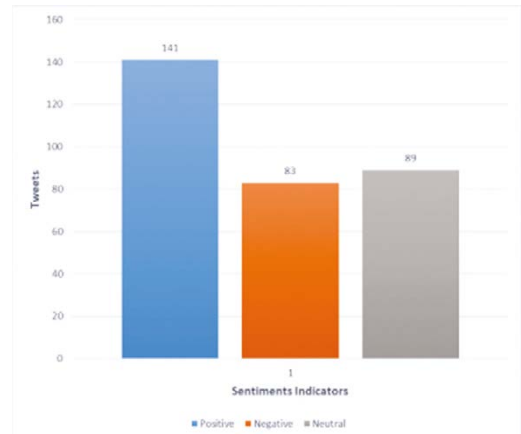
Tweets were assigned sentiment class as explained in Section 3.4. Logarithmic values of “Prior” and “Likelihood” were calculated of the lexicon dataset [24] and it is shown in the Table 9. The total count of the words labeled “Positive” is 2312 and “Negative” is 3324 in lexical dataset. (To see in detail) Detailed mechanism (of how the) of algorithm is explained in section 3.4. (Let us take one) Consider, for example, in a tweet, “style made genuine matching lining,” the word “genuine” exists in the lexical database as “Positive,” therefore, the score for this positive word in a tweet will be calculated using Eq. (10). As there is only one positive word that matches with the words in lexical dictionary and not with the negative word in a tweet, the score for the “Positive” class will be calculated by using logarithmic value in Table 9 as “ $0.891062 + (1 * 7.74586823) = 8.636929$ ”; and for the “Negative” class, it is “ $0.528006 + (0 * 8.1088924156) = 0.528006$ ,” and the final score for the best fit will be calculated using Eq. (11) and the result is shown in Table 10. In Table 10, “5” indicates the tweet number in a Tweeter data T. In this case, as the score is more than 1, the assigned class for this tweet is “Positive.” Similarly, all tweets were assigned scores and classified as “Positive,” “Negative” and “Neutral” based on their “Best Fit” score. Overall classification of cleaned tweets is illustrated in the Figure 7, and it can be seen that about 45% tweets were classified as “Positive”; 27% tweets were classified as “Negative” and approximately 28% of tweets were classified as “Neutral.” The classification results indicate that the chosen fashion brand has quite a positive

**Table 9** Log prior and log likelihood of lexical dictionary.

NRC Lexicon	Count of Words (Nc)	log(Nc/N)	log likelihood (w,c)
Positive	2312	0.891062	7.74586823
Negative	3324	0.528006	8.1088924156
Total (N)	5636		8.636929873

**Table 10** Result interpretation for a tweet has one positive word.

POS	NEG	POS/NEG	BEST FIT
5 8.63692987301857	0.528005717043232	16.357644612987	positive



**Figure 7** Classification of tweets for fashion brand.

outlook from customer perspectives as only 27% of total tweets were negative.

#### 4.2. Correlation Test

The main aim of this study was to investigate if the customer sentiments from tweets collected in the current week have any influence on the upcoming week’s sales. For this, assumption was made and it was verified using Pearson correlation [39]. Following assumptions were made to check the correlation between sentiments indicators  $S_p$ ,  $S_n$ ,  $S_{ne}$  and sales volume  $S_v$ , where  $S_p$ = aggregated “Positive” tweets for a week “w”,  $S_n$ = aggregated “Negative” tweets for a week “w”,  $S_{ne}$ = aggregated Neutral tweets for a week “w”,  $S_v$ = Change in volume of sales for week  $w + 1$ .

- If the correlation between  $S_p$  and  $S_v$  is positive, then the sales will increase.
- If the correlation between the  $S_n$  and  $S_v$  is negative the sales will be decrease.
- If the correlation between the  $S_{ne}$  and  $S_v$  is positive the sales will be increase.

Correlation values, as shown in the Table 11, indicate that there is a moderate positive correlation between  $S_{ne}$  and  $S_v$  that is 32%, weak

positive correlation between  $S_p$  and  $S_v$  with a value of 22% and negative correlation between  $S_n$  and  $S_v$  with  $-13\%$ . This result clearly shows that if there is more number of positive tweets and neutral tweets, it will have positive effect on sales, whereas if there is more number of negative tweets, it will have negative effect on sales.

### 4.3. SMBF Model Performance

In this research FSIS is combined and implemented on EF as explained in Section 3.7 and final forecast was achieved by SMBF using Eq. (16). Model performance of SMBF was evaluated by calculating MAPE values as explained in Section 3.8. MAPE values for EF model and SMBF model is shown in the Table 12. This clearly shows that adding the impact of the sentiments by FSIS model on EF model, exhibit the forecasting capability of model. It can be observed that the performance of the SMBF is better than EF, which solely relies on the historical sales information, whereas SMBF model counts on both historical sales data and social media data. This result illustrates that by adding the features of social media in our forecasting model, the performance of the model can be improved significantly.

The forecasted sales by EF and SMBF models were plotted to compare it with actual sales as shown in Figure 8, which clearly shows that forecast achieved by

SMBF model is significantly better than EF. The prediction of the SMBF model was very close to the actual as shown in Figure 8. As the collected data was during summer season, there is a peak in sales for week 27 and 34 while there is a slight drop in week 36. This could be attributed to the promotional or sales discount effect on the sales. These factors are not considered in this study. However, model SMBF tried to capture the trend to some extent.

## 5. DISCUSSION

This study is conducted on real tweets and sales data of fashion brand for spring and summer assortment. When collected tweets were analyzed after cleaning, the number of tweets was reduced to only 37.5% of the actual tweets in the collected data, which indicates that this brand is in establishing phase on social media. Data pre-processing of tweets was done in two ways: manually; and using text mining to ensure that the collected tweets belonged to the chosen brand. As the enhanced version of lexicon dataset (NRC lexicon) was used, NB classifier classified tweets quite satisfactorily. Most of the tweets were classified as “Positive” and “Neutral,” which represents the popularity of the brand amongst its consumers. This research investigated and established the relationship between the social media data and sales data. Although, the relationship between the social media data and sales data was moderate and weak, we can conclude that there is a relationship between tweet data and fashion apparel sales data, and this indicates that tweet data can be used to forecast garment sales. To deal with the fuzziness of human sentiments and fashion apparel sales, FSIS model was designed, which analyzes the impact of sentiments on sales. FSIS model was combined with EF model to form the SMBF model. The model parameters for SMBF model were optimized using pattern search method used for local optimization reducing the errors between the actual and predicted sales of the model. SMBF model performance is illustrated in Figure 8 and Table 12. SMBF outperforms the EF. There are two peaks in sales for week 28 and week 35, which are summer weeks and model has tried to capture the sales peak during summer to some extent. This shows that the garment sales can be influenced by the social media and if

Table 11 Correlation table.

Variables	Correlation
$S_p$ and $S_v$	0.22
$S_n$ and $S_v$	-0.13
$S_{ne}$ and $S_v$	0.32

Table 12 Model evaluation.

Forecasting Model	MAPE
EF	14.7911567205165
SMBF	10.1143453396686

EF, Exponential Forecast; SMBF, Social Media Based Forecasting.

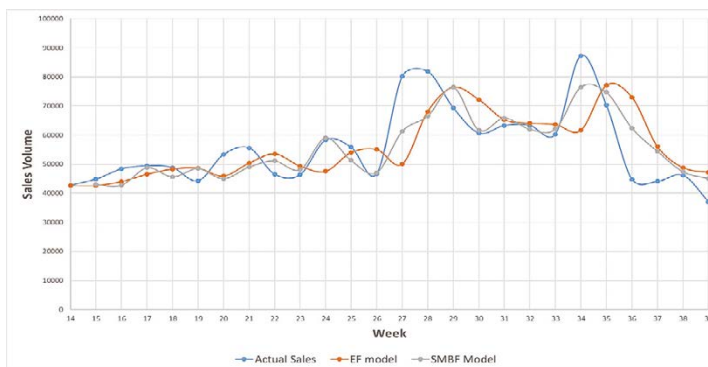


Figure 8 Forecasted sales by Exponential Forecast (EF) model and Social Media Based Forecasting (SMBF) model with actual sales.

its behavior is modelled in existing time series model, we can considerably improve the performance of the traditional forecasting of garment sales.

## 6. CONCLUSION AND FUTURE WORK

Taking into account SMBF model performance and results, we conclude that this model could be effective and beneficial to the fashion garment industry. This research proves that social media data can be used for forecasting sales volume. As the forecasting model will be sensitive to social media data, it is the responsibility of a company to increase its visibility on social media in terms of its services and advertisement.

This will improve the quality of collected social media data. The limitation of this study, firstly, lies in the fact that classification of tweets could be improved by adding more words to the dictionary for the class “Positive” and “Negative.” Secondly, as for the presented model, only “Twitter” as a social media data platform was considered, and if we combine other social media platforms such as “Facebook” and “Instagram,” model performance could be enhanced further.

The experiment in this study was focused on weekly forecast, and therefore, overall sales data of a brand was used for analysis. We would like to extend this work in future to create a short term forecasting model to forecast according to the product category. In this experiment, we combined FSIS model with EF model to build SMBF model for forecasting garment sales. In future, we will incorporate other forecasting models to enhance this study further.

## CONFLICT OF INTEREST

There is no conflict of interest.

## ACKNOWLEDGMENTS

This research work is conducted under the framework of SMDTex-Sustainable Management and Design in Textiles. We are grateful to the company, Evo Pricing, for providing us with the data to carry out this research work.

## REFERENCES

- [1] Z.P. Fan, Y.J. Che, Z.Y. Chen, Product sales forecasting using online reviews and historical sales data: a method combining the Bass model and sentiment analysis, *J. Bus. Res.* 74 (2017), 90–100.
- [2] Q. Ye, R. Law, B. Gu, The impact of online user reviews on hotel room sales, *Int. J. Hosp. Manage.* 28 (2009), 180–182.
- [3] N. Amblee, T. Bui, The impact of electronic-word-of-mouth on digital microproducts: an empirical investigation of Amazon shorts, in 2007 40th Annual Hawaii International Conference on System Sciences (HICSS07), 2007.
- [4] C. Giri, S. Thomassey, X. Zeng, Customer analytics in fashion retail industry, in: A. Majumdar, D. Gupta, S. Gupta (Eds.), *Functional Textiles and Clothing*, Springer, Singapore, 2019, pp. 349–361.
- [5] S. Yu, S. Kak, A Survey of Prediction Using Social Media, 2012. <https://arxiv.org/abs/1203.1647>.
- [6] M.E. Nenni, L. Giustiniano, L. Pirolò, Demand forecasting in the fashion industry: a review, *Int. J. Eng. Bus. Manage.* 5 (2013).
- [7] M. Xia, Y. Zhang, L. Weng, X. Ye, Fashion retailing forecasting based on extreme learning machine with adaptive metrics of inputs, *Knowl. Based Syst.* 36 (2012), 253–259.
- [8] S. Thomassey, Sales forecasts in clothing industry: the key success factor of the supply chain management, *Int. J. Prod. Econ.* 128 (2010), 470–483.
- [9] V. Setyani, Y.-Q. Zhu, A.N. Hidayanto, P.I. Sandhyaduhita, B. Hsiao, Exploring the psychological mechanisms from personalized advertisements to urge to buy impulsively on social media, *Int. J. Inf. Manage.* 48 (2019), 96–107.
- [10] S. Asur, B.A. Huberman, Predicting the future with social media, in 2010 IEEE/WIC/ACM International Conference on Web Intelligence and Intelligent Agent Technology, Toronto, 2010, pp. 492–499.
- [11] N. Archak, A. Ghose, P.G. Ipeirotis, Deriving the pricing power of product features by mining consumer reviews, *Manage. Sci.* 57 (2011), 1485–1509.
- [12] Y. Liu, X. Huang, A. An, X. Yu, ARSA: a sentiment-aware model for predicting sales performance using blogs, in SIGIR '07 Proceedings of the 30th Annual International ACM SIGIR Conference on Research and Development in Information Retrieval, Amsterdam, 2007.
- [13] M. Mendoza, B. Poblete, I. Valderrama, Nowcasting earthquake damages with Twitter, *EPJ Data Sci.* 8 (2019).
- [14] O. Şerban, N. Thapen, B. Maginnis, C. Hankin, V. Foot, Real-time processing of social media with SENTINEL: a syndromic surveillance system incorporating deep learning for health classification, *Inf. Process. Manage.* 56 (2019), 1166–1184.
- [15] G. Chen, Q. Kong, N. Xu, W. Mao, NPP: a neural popularity prediction model for social media content, *Neurocomputing.* 333 (2019), 221–230.
- [16] Z. Wang, S.-B. Ho, Z. Lin, Stock market prediction analysis by incorporating social and news opinion and sentiment, in 2018 IEEE International Conference on Data Mining Workshops (ICDMW), Singapore, 2018, pp. 1375–1380.
- [17] R. Vijayan, G. Mohler, Forecasting retweet count during elections using graph convolution neural networks, in 2018 IEEE 5th International Conference on Data Science and Advanced Analytics (DSAA), Turin, 2018, pp. 256–262.
- [18] W. Du, S.Y.S. Leung, C.K. Kwong, A multiobjective optimization-based neural network model for short-term replenishment forecasting in fashion industry, *Neurocomputing.* 151 (2015), 342–353.
- [19] E.W.K. See-To, E.W.T. Ngai, Customer reviews for demand distribution and sales nowcasting: a big data approach, *Ann. Oper. Res.* 270 (2018), 415–431.
- [20] C. Giri, N. Harale, S. Thomassey, X. Zeng, Analysis of consumer emotions about fashion brands: an exploratory study, in Conference on Data Science and Knowledge Engineering for Sensing Decision Support, Belfast, 2018, pp. 1567–1574.
- [21] S. Ren, H.-L. Chan, T. Siqin, Demand forecasting in retail operations for fashionable products: methods, practices, and real case study, *Ann. Oper. Res.* (2019).
- [22] J. Gentry, Twitter client for R, n.d. <http://geoffjentry.hexdump.org/twitteR.pdf>.

- [23] S.J. Russell, P. Norvig, E. Davis, *Artificial Intelligence A Modern Approach*, Pearson Education India, 2015.
- [24] M. Saif, NRC Emotion Lexicon. <https://saifmohammad.com/WebPages/NRC-Emotion-Lexicon.htm>.
- [25] M.E. Maron, Automatic indexing: an experimental inquiry, *J. ACM*. 8 (1961), 404–417.
- [26] C.D. Manning, P. Raghavan, H. Schütze, *Introduction to Information Retrieval*, Cambridge University Press, Cambridge, 2008.
- [27] S. Russell, P. Norvig, *Artificial Intelligence—A Modern Approach*, third ed., Pearson Education, 2013. <https://www.pearson.com/us/higher-education/program/Russell-Artificial-Intelligence-A-Modern-Approach-3rd-Edition/PGM156683.html>
- [28] R. Brown, *Exponential Smoothing for Predicting Demand*, Arthur D. Little, Cambridge, 1956.
- [29] C.C. Holt, Forecasting seasonals and trends by exponentially weighted moving averages, *Int. J. Forecast.* 20 (2004), 5–10.
- [30] A. Pankratz, *Forecasting with Univariate Box-Jenkins Models: Concepts and Cases*, John Wiley & Sons, New York, 2009.
- [31] L.A. Zadeh, Outline of a new approach to the analysis of complex systems, *IEEE Trans. Syst. Man Cybern. SMC-3* (1973), 28–44.
- [32] S. Gottwald, An early approach toward graded identity and graded membership in set theory, *Fuzzy Sets Syst.* 161 (2010), 2369–2379.
- [33] L.A. Zadeh, The concept of a linguistic variable and its application to approximate reasoning—I, *Inf. Sci. (Ny)*. 8 (1975), 199–249.
- [34] D. Dubois, H. Prade, *Fuzzy Sets and Systems Theory and Applications*, Academic Press, New York, 1980.
- [35] I. Beg, S. Ashraf, Similarity measures for fuzzy sets, *Appl. Comput. Math.* 8 (2009), 192–202.
- [36] W. Pedrycz, Why triangular membership functions?, *Fuzzy Sets Syst.* 64 (1994), 21–30.
- [37] E.H. Mamdani, S. Assilian, An experiment in linguistic synthesis with a fuzzy logic controller, *Int. J. Man. Mach. Stud.* 7 (1975), 1–13.
- [38] C. Audet, J.E. Dennis, Analysis of generalized pattern searches, *SIAM J. Optim.* 13 (2002), 889–903.
- [39] K. Yeager, *LibGuides: SPSS Tutorials: Pearson Correlation*. <https://www.semanticscholar.org/paper/LibGuides%3A-SPSS-Tutorials%3A-Pearson-Correlation-Yeager/7c8f5092aa5a6b963f5e71abb4021344296eedc8#related-papers>.





## Paper V

**Giri, C., Johansson, U., & Löfström, T.** (2019, December). Predictive Modeling of Campaigns to Quantify Performance in Fashion Retail Industry. In *2019 IEEE International Conference on Big Data (Big Data)* (pp. 2267-2273). IEEE.

Available at: <https://ieeexplore.ieee.org/document/9005492>

DOI: 10.1109/BigData47090.2019.9005492



# Predictive Modeling of Campaigns to Quantify Performance in Fashion Retail Industry

Chandadevi Giri

Email: [chandadevi.giri@hb.se](mailto:chandadevi.giri@hb.se)

*Dep. of Business Administration and  
Textile Management,  
University of Borås,  
Automation, Computer Engineering  
Signal & Image Analysis GEMTEX,  
ENSAIT*

University of Lille

*Dep. of Textile Engineering  
College of Textile and Clothing  
Engineering  
Soochow University*

Ulf Johansson

Email: [ulf.johansson@ju.se](mailto:ulf.johansson@ju.se)

*Dept. of Computer Science and  
Informatics, Jönköping University*

Tuwe Löfström

Email: [tuwe.lofstrom@ju.se](mailto:tuwe.lofstrom@ju.se)

*Dept. of Computer Science and  
Informatics, Jönköping University,*

**Abstract**— Managing campaigns and promotions effectively is vital for the fashion retail industry. While retailers invest a lot of money in campaigns, customer retention is often very low. At innovative retailers, data-driven methods, aimed at understanding and ultimately optimizing campaigns are introduced. In this application paper, machine learning techniques are employed to analyze data about campaigns and promotions from a leading Swedish e-retailer. More specifically, predictive modeling is used to forecast the profitability and activation of campaigns using different kinds of promotions. In the empirical investigation, regression models are generated to estimate the profitability, and classification models are used to predict the overall success of the campaigns. In both cases, random forests are compared to individual tree models. As expected, the more complex ensembles are more accurate, but the usage of interpretable tree models makes it possible to analyze the underlying relationships, simply by inspecting the trees. In conclusion, the accuracy of the predictive models must be deemed high enough to make these data-driven methods attractive.

**Keywords**—Fashion retail, Campaign Prediction, Machine Learning, Predictive Modeling, Decision Trees, Random Forest

## I. INTRODUCTION

Promotional campaigns have gained popularity with the growing multi-channel consumer interaction. With this in mind, it is not surprising to see the importance put on understanding promotional campaigns in the marketing literature. Both quantitative and qualitative works have been conducted in this domain, with applications in different retail industries. The background of this research is briefly outlined in [1], [2]. It has been observed from many studies that retailers segment their customers and offer different promotions in personalized ways to achieve maximum retention. Some papers e.g., [3][4], investigate scenarios of targeting faithful customers, while [5] suggests that long-term customers should be offered lower prices than newer. Dynamic pricing has been adopted by the industry that uses data from their databases to optimize the prices of the product based on market demands [6]. Several online fashion retailers offer promotions based on subscriptions and referrals [7], [8]. Conditional promotions are also popular, for example, a consumer will be offered a 1€ discount on the purchase of a minimum of 5€. This was analyzed in [9], where it was concluded that conditional campaigns were successful in

driving unplanned shopping. Technology has also enabled RFID-based systems which could help to understand the important elements in retailing [10], [11]. This has allowed the industry to store information about the consumer buying patterns and behavior in a better way. Many studies have been conducted to model consumer behavior that relies on the pricing factor [12][13]. Artificial intelligence [14] and machine learning [15] have gained importance in the fashion and clothing industry mainly focusing on the manufacturing sector and for improving supply chain management. However, the industry could benefit from using data analytics in several other scenarios. In this paper, we present a data-driven campaign prediction model to identify the profitability and success rate of campaigns.

## II. PROBLEM STATEMENT AND RESEARCH AIM

Attracting and retaining customers is challenging for most industries. However, for the fashion retailers, it is even harder, since there are many influencing factors for customers in the digital world. So, continuously working for retaining customers is vital for these companies, to gain or keep a competitive advantage. The promotional campaigns help fashion retailers to retain customers by offering attractive discounts or services like free delivery, free returns and gifts. Promotional campaigns are, however, very expensive for the retailers, in particular if the conversion rates, and consequently the revenues, are low. Therefore, fashion retailers need to have an appropriate and well-defined strategy for designing their campaigns. In practice, it is fair to say that the fashion retail industry still struggles to even identify the parameters driving customer attention. In this paper, we propose a data-driven method to identify the profitability and ultimately success rate of different campaigns. The suggested method uses predictive modeling where two different models are created:

1. **Profitability:** This regression model predicts the average profit from a campaign.
2. **Success:** This classification model predicts the overall success of a campaign

Fig. 1 below describes the business problem of fashion retailers and states some possible solutions investigated here. In the empirical study, we present a case study on 826 promotional campaigns from a leading Swedish fashion retailer. In more detail, we study campaign data features and use machine learning to model *profitability* and *success*. The results of this study are promising enough to argue that fashion retailers could use this or similar data-driven methods to optimize campaigns. This could also help them significantly in deciding the discount levels (e.g., 10%, 20%, etc.) and included services, e.g., free delivery, free return, etc.

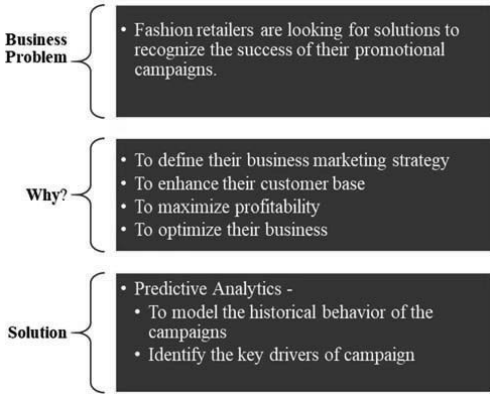


Fig. 1. Overview of Research Problem and Objective

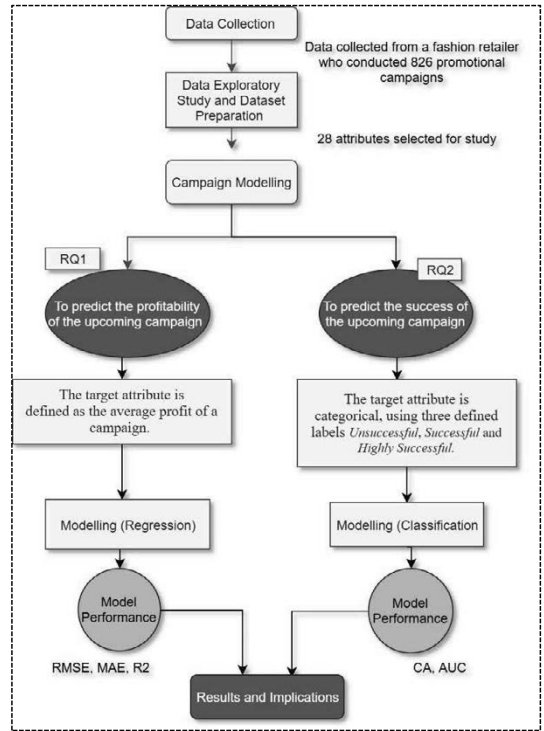


Fig. 2. Research Framework

The research framework used in this study is explained in Fig. 2, which outlines the research methodology of the work, describing the steps from data collection, data analysis and modeling. A detailed description of these steps is given in the following sections.

The succeeding sections are organized like this: First section III describes the methodology and the model evaluation metrics used in the research work. Section IV briefly discusses the steps employed to formulate the models, with the subsequent subsections discussing data collection, exploratory study, creating new features and developing models. Section V discusses the prediction results of the models. Section VI concludes this study.

### III. PREDICTIVE MODELING AND EVALUATION

This section describes the regression and classification models and the evaluation metrics used in this work.

#### A. Predictive Models

Two different kinds of predictive models are used for modeling the campaign behavior in the experimentation described below:

**Decision trees** is a tree-like representation of a model that is formed by the set of rules splitting at the nodes, and making

predictions in the leaves. The most popular decision tree models are CART[16] and C4.5/C5.0 [17]. The main advantage of decision trees over other machine learning modeling techniques is their interpretability, thus allowing explanations of individual predictions and inspection of overall underlying relationships.

**Random forest** [18], is an ensemble learning technique that can be used to develop both classification and regression models. While random forests consist of decision trees (or regression trees), these base models (called *random trees*) differ slightly from their standard counterparts. More specifically, in order to introduce the necessary diversity, bootstrap sampling is used, and the attributes available when optimizing a split are restricted to a randomized subset of the attributes.

#### B. Model Evaluation Metrics

In the empirical investigation, several evaluation metrics are used.

##### Regression

- RMSE is defined as the square root of the averaged squared errors.
- MAE is the mean absolute error of the predictions.
- $R^2$  is the proportion of the variance in the explained variable explained by the model.

### Classification

- Classification accuracy indicates the proportion of accurately classified examples.
- Area under ROC curve is a measure of the ordering capability of the classifier.

### IV. METHOD

This section briefs the steps followed in conducting the research aligned with the campaign prediction as explained in the research framework (Fig. 2). This section is divided into three parts, beginning with the data collection, then exploring the campaign data attributes followed by modeling.

TABLE I. CAMPAIGN DATA ATTRIBUTES

No.	Swedish	Description
1	check	True when campaign type was check
2	firstLine	True when campaign type was first line
3	firstAndSecondLine	True when campaign type was first and second line
4	combDiscount	True when campaign type was combination
5	allOrder	True when campaign type was all order
6	ladder	True when campaign type was ladder
7	threeForTwo	True when campaign type was three for two
8	discCheck	The percent discount when payed by internal check. Can only occur together with check campaigns
9	discFirst	The percent discount on the first item. Firstline, firstandsecondline, combDiscount, allorder and ladder must have a value for DisFirst
10	discRest	The percent discount on the second item. CombDiscount, and ladder must have a value for DiscRest. Firstandsecondline can have a DisSecond.
11	marketingDiscount	The average amount discount received per order (in SEK) for the entire campaign
12	freeGift	True when offered a free gift
13	payedGift	True when offered a payed gift
14	freeShipping	True when offered free shipping
15	freeExressShipping	True when offered free express shipping
16	freeReturn	True when offered free return
17	reqSale	True when the campaign required sale items
18	reqReducedPrice	True when the campaign required items with reduced price
19	reqBrandSelection	True when the campaign required items from specific brands
20	reqValue	True when the campaign required a minimum value
21	req#Items	True when the campaign required a minimum number of items
22	reqTime	True when the campaign required to be used within a limited time
23	reqOrdinaryPrice	True when the campaign required items with ordinary prices
24	reqRedOrdPrice	True when the campaign required items with reduced or ordinary prices
25	grossDemand	The average demand created for a campaign
26	NumRecipients	The number of recipients exposed to a campaign
27	NumOrders	The number of orders resulting from a campaign
28	profit	The average profit of the order, the target value for the regression modelling

### A. Data Collection

Data from 826 unique conducted campaigns were collected from a Swedish fashion retailer. Every campaign is described using the 28 features described in Table I. The first seven attributes represent the campaign types, which are mutually exclusive. The following four attributes (8-11) are the discount(s) offered in the campaign. The attributes 12-16 are addons, offering some additional services. The attributes 17-24 are all different requirements on the order to allow the campaign offerings to apply. Attribute 25, finally, is the average demand created for a campaign. The three attributes at the end (26-28) are either a target attribute or used to define a target attribute.

### B. Exploratory Data Analysis and Data Preparation

This section discusses the data attributes in more detail while visualizing their distributions.

1. Campaign types: There are 7 campaign types (see attributes 1-7 in Table I). The frequency distribution of the order type is shown in the Fig. 3. Almost 50 % of all campaigns are *first line* campaigns, followed by *all order* and *check* campaigns, with more than 100 campaigns each. There are also 26 unclassified campaigns, lacking campaign types.
2. The discount attributes of the campaigns are attributes 8-10, with attribute 11 representing the average discount for the entire order. Check discounts (*discCheck*) are only combined with *check* campaigns, whereas discounts on the first item (*discFirst*) are included in all types of campaigns except the *check* and *threeForTwo* campaigns. The attribute *discRest* are used in campaigns offering different discount levels, like *firstAndSecondLine*, *combination* and *ladder*. Distribution of the first discount attribute is shown in Fig. 3 below. As expected, since most campaigns include a discount on the first item, the number of campaigns with *discFirst* is much higher, with a majority of campaigns having discounts with 30% or more on the first item. For the same reason, most campaigns do not include check discounts or more than one discount, so the majority of campaigns have 0% on *discCheck* and *discRest*. The campaigns that do include these discounts offer discounts ranging from 10% to 50%.

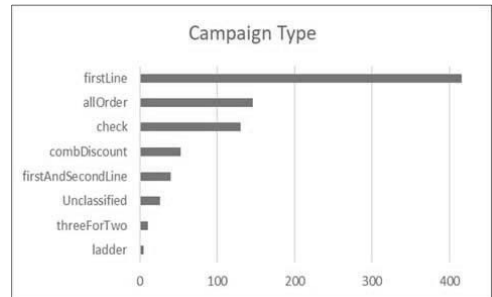


Fig.3. Distribution of order type features

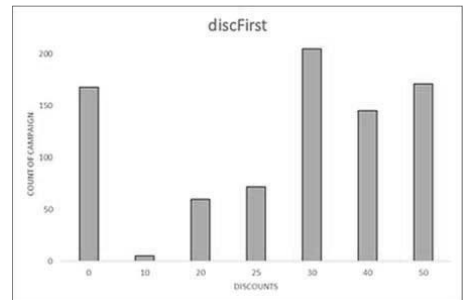


Fig. 4 Distributions of DiscFirst features

- Addon attributes (12-16) offer some extra service to the customers, like free shipping or a free gift.
- Requirement attributes (17-24) limits the applicability of the campaign to only apply if you meet the requirements.
- The gross demand attribute (25) is the average demand created for a campaign.
- Fig. 5 shows the distribution of average profit per order within campaign. The average profit is rightly skewed, with most orders having a positive profit. The profit numbers are in Swedish currency (SEK) and the majority of the data points lie between 0 and 200.

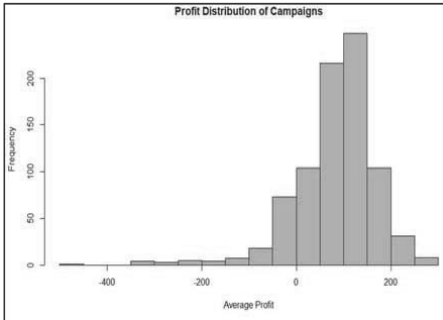


Fig.5. Average profit distribution of campaigns

For the regression experiment (RQ1) the target attribute is the average profit, as described in Table I. For the classification (RQ2) a new feature called *activation* was first created using

$$Activation = \frac{Total\ number\ of\ orders\ received}{Total\ number\ of\ receipts} \quad (1)$$

Looking at the distribution in Fig. 6 below, it can be observed that it is heavily left-skewed. Most of the data points fall between 0 and 0.02, i.e., the response rate is typically less than 2%.

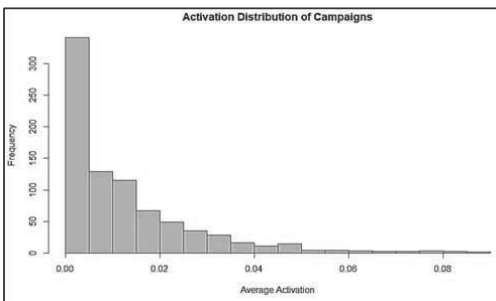


Fig.6. Activation distribution

Now, the classification target for *success* was defined based on both *profit* and *activation* as shown in Fig. 7. Three labels {*Unsuccessful*, *Successful*, *Highly Successful*} were introduced using the criteria in Table II. The obvious logic is that a truly successful campaign should have both high activation and profit.

TABLE II CAMPAIGN PERFORMANCE DEFINING CRITERIA

Model	Values (A=Activation and P=Profit)	Campaign Performance
LA and LP	$A \leq 0.02 \ \& \ P < 0$	Unsuccessful
LA and HP	$A \leq 0.02 \ \& \ P > 0$	Successful
HA and HP	$A > 0.02 \ \& \ P > 0$	Highly Successful

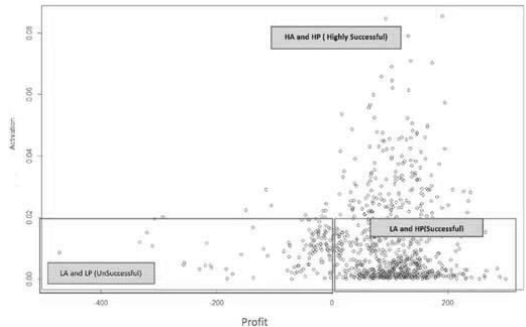


Fig.7. Labels of "Campaign Performance" using Profit and Activation

Out of 826 campaigns, 115 campaigns are classified as *Unsuccessful*, 546 campaigns as *Successful* and 165 campaigns as *Highly Successful*.

### C. Modeling

To predict the profit of the campaigns, regression trees and random forest were applied. Here the features 1 to 25 in Table I are used as input variables with the target *profit*. Standard 10-fold cross-validation was used for the evaluation. Default model parameters were used for the decision trees and the random forest, i.e., the forests consist of 100 trees.

To predict the success of the campaign, classification trees and random forest were applied. Again, the inputs consist of attributes 1-25 in Table I, with *success* as target. Learning parameters were the same as in the regression experiment.

## V. EXPERIMENTAL RESULTS AND DISCUSSION

This section presents and discuss the results.

### A. Profitability Model Evaluation

Table III shows the performance of the regression models predicting profitability. First of all, looking at the error metrics, it should be noted that the predictions are quite accurate, on average. Actually, as seen by the MAEs, the mean prediction error is lower than 40 SEK. As expected, the random forest outperformed the regression trees, but, again looking at MAE, the differences in actual numbers are not that large. Also from looking at R2 values, it is obvious that the models are fairly accurate, explaining more than 50% of the relationship between inputs and target.

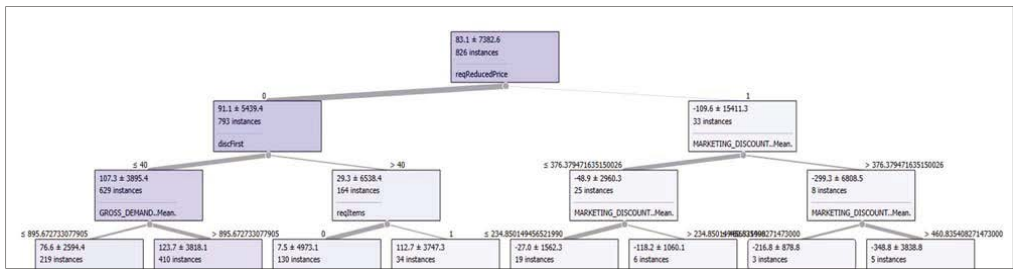


Fig. 8. Regression Tree of Profitability Model

TABLE III PREDICTION RESULTS ON TEST DATA FOR PROFITABILITY MODEL

Model	RMSE	MAE	R2
Regression Trees	57.458	38.531	0.553
Random Forest	45.783	32.164	0.716

Fig. 8 shows the top four levels of the regression tree built for the profitability task. The depth of the tree is limited for ease of understanding. A few general insights that can be drawn from the tree are that required reduced price results in negative profit, just as discounts higher than 40%, while a required minimum number of items results in increased profit.

### B. Campaign Success Model Evaluation

The results for the Campaign Success modelling explained in section IV.C are presented in Table IV. As can be expected, the random forest performs better than the classification tree, even if the difference is not that large.

TABLE IV CLASSIFICATION RESULTS ON CAMPAIGN TEST DATA

Model	AUC	CA
Classification Trees	0.802	0.741
Random Forest	0.843	0.768

Confusion matrices were created for both models to illustrate classifier predictions and is shown in Fig. 9. It can be seen that Random Forest has predicted all target class labels better than Classification trees with prediction accuracy of 87.2% for *Successful*, 52.7% for *Highly Successful* and 76.5% for *Unsuccessful*. It is worth noting that both models are fairly good at identifying the unsuccessful campaigns and that most misclassifications done for the unsuccessful campaigns were classified as successful rather than highly successful.

To demonstrate the model results, ROC curves are plotted for both the models with the target class instances as shown in Fig. 10. This represents the plot of ‘True Positive’ (Sensitivity) with ‘False Positive’ (Specificity), and it could be observed that the ROC curve for the class label *Unsuccessful* has better accuracy for both models than other two class labels, corresponding to the reflection done above. Furthermore, both models have a very steep beginning of the curve, indicating that they both can sort out most of the unsuccessful campaigns rather accurately and quickly.

Classification Tree is plotted to get full insights into the decision rules created by the model (see Fig. 11). The depth of the tree is kept to 4 levels for ease of interpretation. A successful campaign accounts for 66.1 % of the total campaigns. The insights from the tree can be summarized as: High discount on the first item (> 40%) together with a free gift identifies the majority of the unsuccessful campaigns; If it is not a *firstLine* campaign and the discount on the first item is high (> 40%), then the campaign is unsuccessful; Required reduced price will most often result in unsuccessful campaigns; Inclusion of a free gift seems to be a distinguishing mark for highly successful campaigns.

### VI. CONCLUDING REMARKS

In this paper, we have proposed and evaluated a data-driven solution for fashion retailers to model the success of promotional campaigns. In more detail, two different situations were modeled, (i) to predict the profitability of the promotional campaigns and (ii) to classify the campaign success level. To achieve this, we collected data from one of the leading Swedish fashion retailers, where data attributes describe discount features, order type features, promotional features, number of recipients, number of orders received,

Classification Tree						Random Forest					
		Predicted						Predicted			
		Highly Successful	Successful	Unsuccessful	Σ			Highly Successful	Successful	Unsuccessful	Σ
Actual	Highly Successful	74	90	1	165	Actual	Highly Successful	87	78	0	165
	Successful	51	466	29	546		Successful	49	476	21	546
	Unsuccessful	5	29	81	115		Unsuccessful	4	23	88	115
Σ		130	585	111	826	Σ		140	577	109	826

Fig.9. Confusion Matrix with the number of instances

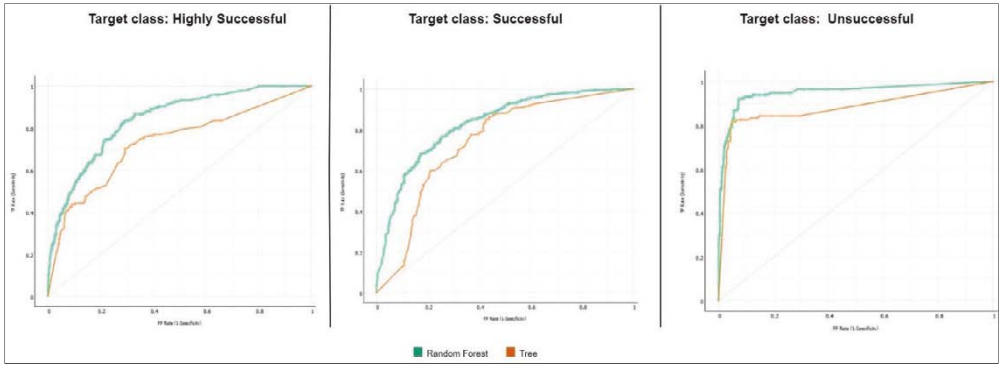


Fig.10. ROC Curve

and profit gained. For the modeling, both tree models and random forest were used. As expected, the random forest models outperformed the tree models, for both classification and regression, regarding the accuracy, but by analyzing the regression and classification trees some valuable insights were learned. Specifically, a high discount on the first item leads to highly profitable campaigns. Besides, if this was combined with a free gift and marketing discounts, it would lead to highly successful campaigns.

Thus, this case study demonstrated that data-driven methods can be used to understand, and ultimately optimize campaigns and promotions. Obviously, such models could be used to simulate campaigns prior to going live, potentially giving the retailers a sophisticated tool for campaign planning.

VII. ACKNOWLEDGMENT

Chandadevi Giri’s work is carried out under the framework of Sustainable Management and Design for textiles (SMDTex). Ulf Johansson and Tuve Löfström are part of the Data-Driven Innovation: Algorithms, Platforms and Ecosystems project, funded by the Knowledge foundation (Grant number: 20160035). The authors are grateful to the Swedish fashion retailer for supplying the data.

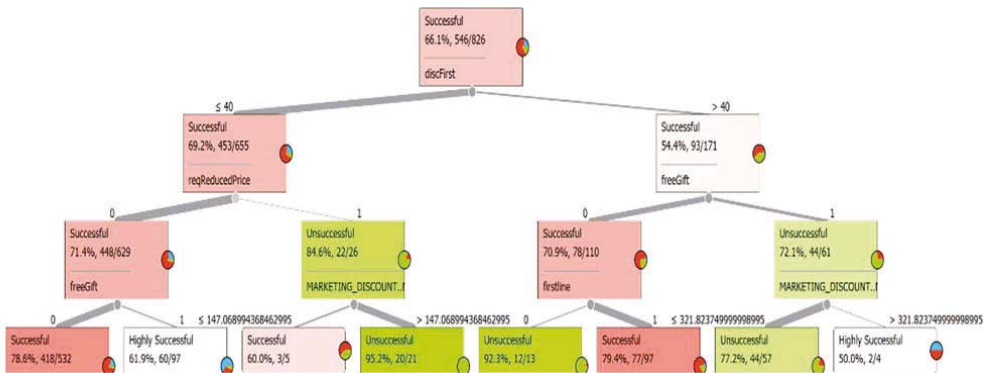


Fig. 11 Classification Tree for Campaign Success Model



REFERENCES

- [1] D. Grewal, K. L. Ailawadi, D. Gauri, K. Hall, P. Kopalle, and J. R. Robertson, "Innovations in Retail Pricing and Promotions," *J. Retail.*, vol. 87, pp. S43–S52, Jul. 2011.
- [2] K. L. Ailawadi, J. P. Beauchamp, N. Donthu, D. K. Gauri, and V. Shankar, "Communication and Promotion Decisions in Retailing: A Review and Directions for Future Research," *J. Retail.*, vol. 85, no. 1, pp. 42–55, 2009.
- [3] F. M. Feinberg, A. Krishna, and Z. J. Zhang, "Do we care what others Get? A Behaviorist Approach to Targeted Promotions," *J. Mark. Res.*, vol. 39, no. 3, pp. 277–291, Aug. 2002.
- [4] C. Giri, S. Thomassey, and X. Zeng, "Customer Analytics in Fashion Retail Industry," in *Functional Textiles and Clothing*, Singapore: Springer Singapore, 2019, pp. 349–361.
- [5] D. Grewal, D. M. Hardesty, and G. R. Iyer, "The effects of buyer identification and purchase timing on consumers' perceptions of trust, price fairness, and repurchase intentions," *J. Interact. Mark.*, vol. 18, no. 4, pp. 87–100, Jan. 2004.
- [6] T. T. Nagle and J. E. (College teacher) Hogan, *The strategy and tactics of pricing : a guide to growing more profitably*, 4th ed. Upper Saddle River, N.J. : Pearson/Prentice Hall, 2006.
- [7] M. J. Barone and T. Roy, "Does Exclusivity Always Pay Off? Exclusive Price Promotions and Consumer Response," *J. Mark.*, vol. 74, no. 2, pp. 121–132, Mar. 2010.
- [8] G. Ryu and L. Feick, "A Penny for Your Thoughts: Referral Reward Programs and Referral Likelihood," *Journal of Marketing*, vol. 71. Sage Publications, Inc., pp. 84–94.
- [9] L. Lee and D. Ariely, "Shopping Goals, Goal Concreteness, and Conditional Promotions," *J. Consum. Res.*, vol. 33, no. 1, pp. 60–70, Jun. 2006.
- [10] S. Ganesan, M. George, S. Jap, R. W. Palmatier, and B. Weitz, "Supply Chain Management and Retailer Performance: Emerging Trends, Issues, and Implications for Research and Practice," *J. Retail.*, vol. 85, no. 1, pp. 84–94, Mar. 2009.
- [11] S. K. Hui, P. S. Fader, and E. T. Bradlow, "Path Data in Marketing: An Integrative Framework and Prospectus for Model Building," *Mark. Sci.*, vol. 28, no. 2, pp. 320–335, Mar. 2009.
- [12] J. Lindsey-Mullikin and D. Grewal, "Imperfect Information: The Persistence of Price Dispersion on the Web," *J. Acad. Mark. Sci.*, vol. 34, no. 2, pp. 236–243, Apr. 2006.
- [13] S.-F. S. Chen, K. B. Monroe, and Y.-C. Lou, "The effects of framing price promotion messages on consumers' perceptions and purchase intentions," *J. Retail.*, vol. 74, no. 3, pp. 353–372, Sep. 1998.
- [14] R. Nayak and R. Padhye, *Artificial intelligence and its application in the apparel industry*. Elsevier Ltd, 2017.
- [15] C. M. Bishop and C. M., *Pattern recognition and machine learning*. Springer, 2006.
- [16] L. Breiman, J. Friedman, R. Olshen, C. S.- Group, and undefined 1984, "Classification and regression trees. Wadsworth Int."
- [17] J. R. (John R. Quinlan, *C4.5 : programs for machine learning*. Morgan Kaufmann Publishers, 1993.
- [18] L. Breiman, "Random Forests," *Mach. Learn.*, vol. 45, no. 1, pp. 5–32, 2001.