

Cross-Channel Predictive Analytics for Retail Distribution Decisions

by
James B. Coles
B.S., Electrical Engineering, University of California San Diego (2008)
M.S., Electrical Engineering, University of California San Diego (2009)

Submitted to the MIT Sloan School of Management and MIT Institute for Data, Systems and Society in partial fulfillment of the requirements for the degrees of

Master of Business Administration
and
Master of Science in Engineering Systems

in conjunction with the Leaders for Global Operations Program at the
MASSACHUSETTS INSTITUTE OF TECHNOLOGY
June 2017

© James B. Coles, MMXVII. All rights reserved.
The author hereby grants to MIT permission to reproduce and to distribute publicly paper and electronic copies of this thesis document in whole or in part in any medium now known or hereafter created.

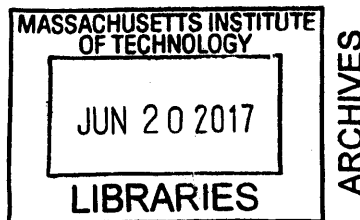
Author **Signature redacted**
MIT Sloan School of Management and MIT Institute for Data, Systems and Society
May 12, 2017

Certified by **Signature redacted**
Georgia Perakis, Thesis Supervisor
William F. Pounds Professor of Management Science
MIT Sloan School of Management

Certified by **Signature redacted**
Bruce Cameron, Thesis Supervisor
Director, MIT System Architecture Lab
MIT System Design and Management

Accepted by **Signature redacted**
Maura Herson
Director, MBA Program
MIT Sloan School of Management

Accepted by **Signature redacted**
John N. Tsitsiklis
Clarence J. Lebel Professor of Electrical Engineering
MIT Institute for Data, Systems and Society Graduate Officer



Cross-Channel Predictive Analytics for Retail Distribution

Decisions

by

James B. Coles

Submitted to the MIT Sloan School of Management and MIT Institute for Data, Systems
and Society
on May 12, 2017, in partial fulfillment of the
requirements for the degrees of
Master of Business Administration
and
Master of Science in Engineering Systems

Abstract

Distribution demand forecasting at Zara currently considers historical sales of products modified by expert knowledge inputs in an algorithm developed to calculate the shipment required to meet demand for the next sales period. In 2010, the introduction of Zara.com provided customers an additional channel to complete purchases and interact with the brand while providing Zara significant insight into changing customer preferences to supplement the expert knowledge of the Zara team.

This thesis investigates the utility of the data collected in the online sales channel for increasing the accuracy of the distribution demand forecasts. Two forecast types are considered: Initial Shipments for which no historical data exists, and Replenishment Shipments which have historical data. Forecasts are performed for both brick-and-mortar and e-commerce sales channels to demonstrate cross-channel utility of the data. The study presents a review of available datasets to identify those of potential interest and describes meaningful features engineered from raw datasets. By applying machine learning algorithms, significant features are identified and a predictive model is developed demonstrating significant WMAPE improvement for initial shipments to brick-and-mortar stores (0.23), moderate improvement for replenishment shipments to e-commerce (0.05) and limited improvement for replenishments to brick-and-mortar stores (<0.04). The results of this study demonstrate the potential for significant reduction of inventory requirements to maintain customer service levels and provides a baseline for future cross-channel forecasting work.

Thesis Supervisor: Georgia Perakis
Title: William F. Pounds Professor of Management Science

Thesis Supervisor: Bruce Cameron
Title: Director, MIT System Architecture Lab

Acknowledgments

The months I spent as part of the Zara team at the headquarters in A Coruña were a highlight of my experience in the LGO program at MIT. I'm filled with gratitude for the longstanding partnership between MIT and Inditex and the opportunity that the LGO program allowed me to engage in. This work is the result of the collective efforts of so many people.

Thank you Zara Distribución team for adopting me as part of the family. My project supervisors Ane Insausti Altuna, Jose Luis Goñi and Iván Escudero Rial were there from the beginning, introducing me to the world of fashion and providing advice whenever I needed it. Alberte Dapena Mora, you answered all of my questions twice and taught me everything I know about distribution and forecasting. Rocio, David, Lorena, Marcos, Leticia, Begoña, Susan, Ana, Pepa, Eva, Patricia, Laura, Carolina, Juan, Sergio, Egowitz, Miguel and Javier thank you all so much for sharing and teaching me about your fabulous work.

Thank you .com team, you introduced me to the cool new world inside Zara. Thank you to Fernando Talín Mariño and Rubén Botana Saavedra for being my .com guides. Jose Manuel Corredoira Corras and Javier Martinez Roldan, you taught me to use your databases, saved me when things made no sense, and made this project work. David, Juan Albuin, Juan Villacampa, Mateo, Olaia, Guillermo, and Miriam, thanks for helping to ease me in and sharing time with me to learn the inner workings of .com.

Thank you to my support at MIT. With regular check-ins, visit, ideas and support, my advisor Georgia Perakis not only kept me focused and on track but encouraged me to finish. Bruce Cameron provided the voice of experience and provided me with the perspective to consider the full context of my work. Lennart Baardman and Divya Singhvi, you helped keep me sane while lost in algorithms and gave me hope that there are good things to come.

Thank you Spanish familia. Alberte, Rocio, Javi, Mateo, Yunchi, Kyoko, Rita, Xiaowen, Aya, Mari, Joyce, Alex, Carlos, Emil, Sarah, Jasmine, Cat, Andres, Ana, Gabriel, Mateo, Javi and Ainoa. I miss you and our lunches. My card is sitting on my desk as I write.

Thank you family. Mom and Dad for visiting and leaving Mom with me. I will treasure our time in Galicia. Sharon, my wife for supporting me while I was away, exploring my new home with me, and loving me through fun times and hanger.

THIS PAGE INTENTIONALLY LEFT BLANK

Contents

1	Introduction	15
1.1	Project Motivation	16
1.1.1	Demand Forecasting	16
1.2	Project Approach	18
1.3	Thesis Contributions	18
1.4	Thesis Overview	19
2	Literature Review	21
2.1	Prior MIT Thesis Work	21
2.2	Retail Fashion Demand Forecasting	22
2.3	Data Analytics in E-commerce	23
2.4	Omni-Channel Demand Studies	23
2.5	MIT 12	23
3	Background	25
3.1	Zara Organization	25
3.1.1	Zara.com	26
3.1.2	Zara Distribution and Operations	26
3.2	Key Stakeholders	26
3.3	Product Classification	27
3.4	Distribution Demand Forecasting	28
3.4.1	Demand Forecasting Overview	29
3.4.2	Distribution Demand Forecast Types	30

3.4.3	Demand Forecast Calculations	30
3.4.4	Inventory Shipment Calculation	32
4	Methodology	35
4.1	Data Analysis	35
4.1.1	Business Objective	36
4.1.2	Data Collection	36
4.1.3	Feature Engineering	36
4.1.4	Modeling	36
4.1.5	Model Evaluation	37
4.1.6	Deployment	37
4.2	Data Sources	37
5	Data Collection	39
5.1	Data Identification	39
5.1.1	Process Research	39
5.1.2	Brick-and-Mortar Stores	40
5.1.3	Zara.com Experience	41
5.2	Dataset Availability	44
5.2.1	Distribution Datasets	45
5.2.2	E-Commerce Datasets	46
5.2.3	Supplemental Datasets	46
6	Feature Engineering	49
6.1	Data Filtering and Cleaning	49
6.1.1	Filtering Parameters	49
6.1.2	Data Cleaning	51
6.2	Feature Construction	51
6.2.1	Notation	52
6.2.2	Article Features	52
6.2.3	Location Features	53

6.2.4	Temporal Features	53
6.2.5	Website Structure Features	54
6.2.6	Customer Behavior Features	56
7	Modeling	59
7.1	Modeling Approach	59
7.2	Data Pre-Processing	60
7.2.1	Feature Reduction	61
7.3	Modeling Algorithms	64
7.3.1	Initial Feature Investigations	64
7.3.2	Algorithm Selection	64
7.3.3	Model Adjustments and Prediction Modifications	67
7.3.4	Model Construction and Testing	67
7.3.5	Coverage Measurement	67
8	Results Analysis	71
8.1	Prediction Error	71
8.1.1	Initial Products	72
8.1.2	Replenishment Products	74
8.1.3	OLS Linear Model Performance Comparison	78
8.2	Distribution Inventory Effects	79
8.2.1	Initial Products	79
8.2.2	Replenishment Products	80
9	Conclusions and Recommendations	85
9.1	Summary of Results	85
9.2	Recommendations	86
9.2.1	Implementation	86
9.2.2	Future Investigations	87
9.3	Conclusion	89

THIS PAGE INTENTIONALLY LEFT BLANK

List of Figures

3-1	An example of demand levels for a single MCC in a single location. The y-axis is normalized to present a percent of maximum article, location demand. . .	28
3-2	The current demand forecasting timeline for new and replenishment product shipments in the physical and e-commerce sales channels	29
4-1	The data analysis process adapted from CRISP-DM	35
5-1	An example website layout of the Zara.com storefront with navigation tree, search bar and filter options	42
5-2	An example product details page with coming soon indicators	43
5-3	A representative database map depicting the available databases and selected tables of datasets contained in the database. Connections are shown with directional arrows, no arrow means there is no direct connection.	45
7-1	The correlation matrix for the training set of a single replenishment feature set. Red highlighted entries represent cross-correlations greater than 0.7 and were flagged for possible removal.	63
7-2	Example tuning of the value of m to produce minimum out-of-bag error in the training dataset.	66
8-1	Forecasting error for baseline and updated models against actual V_m for a subset of test entries.	73
8-2	Forecasting error for baseline and updated models against actual V_m for a subset of test entries for brick-and-mortar stores.	75

8-3	Forecasting error for baseline and updated models against actual V_m for a subset of test entries for zara.com.	76
8-4	The distribution of forecast error ratios used to calculate the coverage required for brick-and-mortar stores.	82
8-5	The distribution of forecast error ratios used to calculate the coverage required for Zara.com fulfillment centers.	83

List of Tables

6.1	Summary of available data for initial shipment analysis.	51
6.2	Summary of available data for replenishment shipment analysis. MCC represents the number of articles available, Dates or cycles are the number of distinct forecast dates or periods, and Subfamilies is the number of subfamilies represented by the articles	51
7.1	Percentage of applicable initial shipment data entries after forecast type filtering for valid website and subscription data.	61
7.2	Percentage of applicable replenishment shipment data entries after forecast type filtering for valid website and subscription data.	61
7.3	Example frequency ratio calculation for Days New (Prior Week)	62
7.4	Initial shipment demand forecast WMAPE improvements with and without subscriptions	66
8.1	Initial shipment demand forecast WMAPE improvements with and without subscriptions	72
8.2	Initial shipment feature importance rankings for brick-and-mortar stores. . .	74
8.3	Replenishment shipment demand forecast WMAPE improvements with and without subscriptions	75
8.4	Initial shipment feature importance rankings for brick-and-mortar stores. . .	77
8.5	Initial shipment feature importance rankings for Zara.com forecasts.	78
8.6	Initial shipment demand forecast comparison using an OLS linear model and an random forest algorithm.	78

8.7	Replenishment shipment demand forecast comparison using an OLS linear model and an random forest algorithm.	79
8.8	Initial shipment relative coverage measurement and improvement across all stores in the selected market.	80
8.9	Initial shipment calculated average shipments and improvement across all stores in the selected market.	80
8.10	Replenishment shipment relative coverage measurement improvement across all stores in the selected market.	80
8.11	Replenishment shipment calculated average shipments improvement across all stores in the selected market.	81

Chapter 1

Introduction

Zara is the flagship brand of the Spanish headquartered Inditex group, the largest global fashion retailer with over 7,000 stores across 88 markets amongst 8 different brands. The Zara brand itself operates more than 2,000 stores and accounts for over 65% of net sales for the group. This footprint continues to grow; Zara introduced a net of 77 new stores in 2015 and is currently present in over 27 online markets through the e-commerce channel Zara.com [1].

Zara operates using a business model offering about 9,000 new articles every season delivered to stores twice a week. This model encourages frequent customer engagement with the brick-and-mortar and online stores to discover the latest designs available. This engagement is further encouraged by the fact that Zara does not utilize traditional marketing methods to advertise new products but relies on high visibility storefronts and online presence to showcase new products.

Zara customers are provided the latest fashion trends at affordable prices thanks to Zara's agile supply chain and inventory distribution processes. By maintaining a high level of awareness to the fast changing trends and demand preferences of its customers, Zara is able to deliver inventory to the right locations at the right time to fulfill customer needs while minimizing excess inventory costs. This ability is a critical component of success in an industry that operates on such a fast pace as fashion.

1.1 Project Motivation

The introduction of the e-commerce sales channel in 2010 at Zara.com created the opportunity for Zara to learn more about customers and to respond even faster to changing trends. As consumers increasingly embrace e-commerce to supply their fashion demand, the level of customer engagement at Zara.com has continued to grow. The datasets that are collected through the normal course of e-commerce operations provide the Zara.com team with great insight into the changing needs of their customers. By identifying the ways that this data can help to inform the processes across the entire organization, significant opportunities for organizational operations improvement can be identified.

The goal of this work is two-fold. First, it provides a critical initial look at the utility and availability of the datasets collected by the e-commerce group at Zara.com to the broader Zara organization. In the absence of traditional marketing campaigns, this data provides critical insights into Zara's customers and their preferences. An analysis of the datasets available to the various teams across the organization in an effort to increase information sharing and availability allowed a number of potential applications of the combined databases to be identified.

Additionally, this project describes a first investigation into the cross-channel effects between operations in the e-commerce and brick-and-mortar domains of Zara's retail operations. Focused on improving the customer demand forecast for articles across both channels, this effort combines the available e-commerce datasets with the existing demand forecasting process to create an updated model for predicting article sales.

1.1.1 Demand Forecasting

An essential strategic advantage in Zara's retail model is an accurate prediction of customer demand. In an ideal case, the retailer knows exactly how many units will be sold in a given location at a given time and have just enough inventory to fill those sales. In reality, retailers do not know exactly how many units will be sold, and must optimize the amount of inventory shipped based on a prediction of future sales and the certainty of that prediction. The generation of this prediction is called demand forecasting. Increasing the expected accuracy

of the forecast reduces the expected error of the number of units shipped versus the number of units actually demanded. By reducing this uncertainty, retailers can continue to meet customer demand with the same service level while minimizing the additional inventory sent to cover possible prediction error.

Traditional Demand Forecasting

For many traditional brick-and-mortar retailers including Zara, demand forecasting is typically performed using a historical average of sales data. This data is often supplemented by the knowledge of sales associates and management that interact with customers to develop a better understanding of customer trends and customer reactions to current products. This information is then combined through linear regression or another statistical analysis tools to develop a description of the past sales and the expected change in sales based on the historical trend modified by the expert knowledge of the store employees.

E-commerce retailers have additional tools to predict future sales of items. In addition to the traditional historical sales data used in brick-and-mortar stores, e-commerce websites monitor traffic through the pages of the website in an effort to improve the user experience. An additional benefit of this data is the insight it lends to the navigation of a customer through the store. This data is similar to what a traditional brick-and-mortar store would have if they tracked the movements of every customer from the time they entered until the time they exited the store. What e-commerce retailers gain in raw data, however, they lose in the personal knowledge and understanding of a customer developed through the interaction of employees in the brick-and-mortar stores.

Multi-Channel Demand Forecasting Potential

Multi-channel retailers operate in multiple retail channels simultaneously. For example, Zara operates both brick-and-mortar and e-commerce retail sales channels. The network of brick-and-mortar stores provides Zara with significant personal interaction with customers and allows for the commercial teams to develop a deep understanding of the preferences of those customers. Similarly, the e-commerce sales channel increases Zara's accessible market size by reaching customers not near a physical Zara location while allowing the company

to collect website traffic data to better understand customer preferences and trends. While each channel offers clear benefits alone, the potential benefit of combining the advantages to inform each other is an area of particular interest.

To effectively use both sets of information together, the multi-channel retailer must understand the implications of each dataset on the other in a study of cross-channel demand effects. Through careful analysis of data between sales channels, retailers stand to better understand the customer needs while streamlining internal processes, and maximizing the number of customers served.

1.2 Project Approach

The proposed investigation was divided into two significant phases of Research and Analysis to address the two primary goals. The research phase included an investigation of the current demand forecasting and inventory distribution processes used by Zara’s distribution team. Additionally, research included a preliminary survey of the datasets currently available across the Zara organization and the level of access available to each stakeholder.

The analysis phase of the project focused on interpreting the data collected to inform company operations. This process collected the identified datasets, aggregated and transformed the raw data into meaningful features, and developed predictive models to explain the trends observed. These models were then tested to demonstrate applicability to both sales channels and demonstrate the utility of the datasets for predicting customer demand.

1.3 Thesis Contributions

The contributions of this thesis focus on applying machine learning techniques to the demand forecasting process for multi-channel retailers. In particular, the cross-channel effects between e-commerce and brick-and-mortar retail channels are demonstrated using the demand forecasting process. A random forest model is presented that incorporates cross-channel data to significantly improve demand forecast accuracy. Finally, a baseline feature set is defined for future study of cross-channel effects and demand forecasting.

Primarily, this work develops a random forest model to demonstrate the utility of e-commerce data to supplement traditional demand forecasting techniques. Two specific demand forecasting processes are defined corresponding to two distinct product types: new products with no historical sales data and replenishment products with at least one week of historical sales data. Each process is investigated comparing the performance of models incorporating e-commerce data with the current (baseline) forecasting model. This work shows that, with sufficient training data, a random forest model incorporating e-commerce data provides a significant advantage compared to current processes focused on historical sales. Absolute WMAPE improvement of 0.23 was observed for new article forecasts in brick-and-mortar stores, while improvements up to 0.06 were observed for replenishment shipments of existing articles in the e-commerce channel.

In support of the model development effort, a significant portion of the thesis work focused on data collection. Specifically, the identification of available datasets in the structure of a multi-channel retail operation and incorporating these datasets across database boundaries into a single, interpretable feature set. The result of this effort is a feature set that provides a framework for prioritized data collection in the e-commerce channel with a focus on optimizing utility for demand forecasting. The minimum required feature set is identified which will form the baseline for future studies involving additional datasets of interest.

1.4 Thesis Overview

This thesis is organized in chapters to facilitate ease of review according to the following outline:

Chapter 1 Introduces the project, the business motivation for the investigation performed, and the structure of the document.

Chapter 2 Reviews relevant literature as it applies to this investigation including previous work performed in partnership between Zara and MIT LGO.

Chapter 3 Develops the context of the project including a background of Zara and the fashion industry it operates in. The chapter also provides background descriptions of

the processes involved in the distribution operations for e-commerce and brick-and-mortar stores.

Chapter 4 Describes the detailed methodology used to identify and analyze the datasets in the study.

Chapter 5 Identifies the sources of data identified and analyzed in the investigation and describes their significance to the work.

Chapter 6 Documents the feature engineering process, transforming raw datasets into useful features to be interpreted by a model. The calculation of the features, rationale behind their use, and the resulting feature vectors and feature set is described in this chapter.

Chapter 7 Explains the modeling process used to create the predictive model demonstrating the utility of the feature sets produced. Additionally, the process of determining results as WMAPE rates and expected average inventory shipment reductions are described in this chapter.

Chapter 8 Presents the modeling results for the tests described in Chapter 7 and provides an interpretation of the meaning and cause of the results.

Chapter 9 Summarizes the findings of the investigation and the implications for Zara's operations. Additional areas of investigation and development as well as practical steps for implementation are described here.

Chapter 2

Literature Review

Significant prior research has been conducted in many of the fields related to this investigation. This work is the continuation of a longstanding partnership between the Massachusetts Institute of Technology (MIT) Leaders for Global Operations (LGO) program and Zara's distribution team representing the MIT 12 project. Together, Zara and LGO have investigated a number of operations research projects resulting in a number of tools still in use today. Additionally, the topics of demand forecasting in the fashion industry, data analytics, and omni-channel retailing have been the subjects of extensive study and form a solid foundation for the work described in this thesis.

2.1 Prior MIT Thesis Work

The research presented in this thesis marks the twelfth (MIT 12) LGO project conducted in collaboration with Zara. Prior projects performed in the partnership between the MIT LGO program and Zara's distribution team have spanned a number of operations optimization topics. Investigations by former students have spanned many aspects of operations within Zara's organization. The subject of demand forecasting has been one of particular interest spanning multiple prior theses. Both Initial Shipment and Replenishment demand forecasting processes have been carefully examined and continually optimized over the many years of partnership. In particular, prior thesis research performed by Correa [2], Garro [3], Garcia [4], and Kong [5] were essential resources forming the baseline processes and concepts

expanded upon in MIT 12.

2.2 Retail Fashion Demand Forecasting

Demand forecasting is a topic of significant academic interest, particularly with applications in the retail fashion industry. Indeed, the sheer number of reviews of demand forecasting models available testifies to the popularity of the subject; Choi et al.(2011) [6], Liu et al. (2013) [7], Nenni et al. (2013) [8], Thomassey (2014) [9], and Ren et al. (2016) [10] to name a few. Current literature emphasizes the complexity of the process given the traditionally long product development times, short product lifecycles, high product variety, and volatile demand profiles of the industry. In forecasting demand, retailers must identify the appropriate data aggregation level for the forecasting model used to overcome some of these effects. Additionally, retailers identify and collect relevant explanatory variables, or features, to explain demand variation in such (non-exhaustive) categories as: Item Features and Fashion Trends, Retailing Strategy, Marketing Strategy, Macro-economic Data, Calendar Data, Competition, and Weather Data [9].

A number of approaches to fashion demand forecasting have been proposed and studied over decades of research. Demand forecasting and inventory distribution optimization at Zara has been previously studied and implemented in research by Caro and Gallien over a longstanding research relationship [11]. Statistical methods are a popular approach, particularly the time-series analysis method autoregressive integrated moving average (ARIMA) and the seasonal variant (SARIMA). These methods are fast to perform, however cannot consider all of the relevant explanatory factors necessary. Artificial intelligence algorithms such as evolutionary neural networks (ENN) [12] have been explored, however these algorithms require long training times preclude them from significant practical use. Extreme learning machines are a special implementation of neural network called a single hidden-layer feed forward network and require less training time than ENN [13]. Recent research has been conducted in hybrid approaches, combining statistical analysis with artificial intelligence algorithms to improve forecasting performance with short prediction time [14]. Progress continues to be made in the fashion demand forecasting industry, decreasing the

required amount of time to train and predict using the models while increasing the accuracy of the forecasts produced.

2.3 Data Analytics in E-commerce

Data analytics in e-commerce has been widely studied as a field of significant interest to academics and companies. E-commerce firms collect a wealth of information about the traffic to and around their websites. Tools such as Google Analytics and Piwik allow website owners to track the digital footprints of customers through carefully designed data architectures [15]. This data can be used in a variety of applications from simple traffic analysis and marketing campaign monitoring to demand forecasting for retail pricing decisions using machine learning models [16].

2.4 Omni-Channel Demand Studies

There have been recent efforts to conduct research focused on the omni-channel effects on retail operations and customer demand. Omni-channel retailers offer products to customers via two or more sales channels, typically brick-and-mortar stores, e-commerce websites, and mobile apps. Indeed, studies have investigated correlations between trips to brick-and-mortar stores and online stores to determine that online shopping encourages brick-and-mortar shopping trips, but the relationship is not reciprocal [17]. As consumers increasingly engage these alternate channels, the distinction between channels is diminishing. Retailers are adapting their strategies to adjust to this shift and embrace the opportunity to understand customer engagement across channels better through data analytics [18].

2.5 MIT 12

In context of the significant established research on the subject of demand forecasting, this thesis seeks to form a bridge between existing models. As brick-and-mortar retailers continue to grow e-commerce sales, datasets collected by the retailer websites become signifi-

cant sources of insight into customer preferences. The established research investigates the current best practices for demand forecasting among fashion retailers and in e-commerce applications, however significant research has not been published exploring the application of data analysis in multi-channel fashion retail. Specifically, the cross-channel effects between brick-and-mortar and e-commerce demand indicator datasets is explored in this thesis, demonstrating a significant potential increase in demand forecast accuracy attributable to combining information from both sales channels.

Chapter 3

Background

Zara is well known in the fashion industry. With net sales in 2015 of more than €13M [1], Zara is one of the world's largest fashion retail brands competing with the likes of Hennes & Mauritz (H&M), Fast Retailing's Uniqlo brand, and Gap. Zara's highly vertically integrated and responsive supply chain combined with its ability to understand industry trends and preferences of its customers provide significant competitive advantages to the brand in the fast paced, highly variable world of fashion. The internal organization of Zara and the operations processes relevant to demand forecasting and e-commerce data collection are presented here to provide context for the rest of the investigation.

3.1 Zara Organization

The Zara brand consists of three departments : Women, Men, and Kids. Each department is further divided into several tipos de producto (English: product type). At the corporate headquarters, each department has groups of product manager teams assigned to each product type that oversee the commercial aspects of the stores for a particular market (e.g. Spain, China or the U.S.A). Additionally, the Distribució (English: Distribution) team coordinates and monitors distribution of inventory throughout Zara's network of stores.

3.1.1 Zara.com

Zara.com, the e-commerce sales channel for the Zara brand, was launched in late 2010 to offer more options for customers to engage with the brand. Since it was established, Zara.com has experienced impressive year over year growth and is currently present in over 27 markets worldwide [1]. Customer orders to Zara.com are fulfilled through warehouses called Fulfillment Centers (FC) located in or near the markets that they serve. Some markets are served by multiple FCs and some FCs fulfill orders from multiple markets.

Zara.com operates within the Zara organization. The commercial team within Zara.com is responsible for individual departments and product types, similar to the organization of the brick-and-mortar commercial teams. Distribution teams are responsible to coordinate with the Zara.com commercial team as another set of product managers in the same way they do for the brick-and-mortar stores. Additionally, the Zara.com includes engineers to monitor website traffic data collection and analysis as well as general information to maintain smooth operations.

3.1.2 Zara Distribution and Operations

The primary responsibility of the distribution team is to ensure that the number of items in each article that are available are allocated to the correct stores in the right size mixture. This responsibility also includes store operations processes for moving inventory between the sales floor and the stock room as well as inventory transfers between stores and returns. A key component to optimizing the distribution of inventory across this complex network is the demand forecast.

3.2 Key Stakeholders

A key motivation of this work is the collaboration opportunity between the e-commerce and brick-and-mortar store operations. As such, key stakeholders exist in each organization as well as across the wider company and represent critical interfaces for the work described.

Product Managers Zara's product managers are commercial experts responsible for man-

aging the products sold at each retail location including Zara.com storefronts. Product managers are intimately aware of customer preferences and provide location-specific expert product knowledge that is used by the Distribution team in demand forecasting.

Zara Distribución (Distribution) The Zara Distribución team is responsible for the distribution of all inventory across all sales channels for the company including brick-and-mortar stores and e-commerce websites. The Distribution team is the team most directly involved in the work described here.

Information Technology (IT) Engineers The IT Engineers are responsible for the implementation of the Inditex information infrastructure. These engineers develop and maintain databases and tools, including demand forecasting tools investigated in this work, used by teams across the organization. Additionally, IT engineers maintain and operate the Zara.com website, associated databases, and website traffic analytics.

3.3 Product Classification

Zara products are identified by multiple levels of classification:

Family The highest level of product classification is familia (English: family) and represents the most generic grouping of product. (e.g. Dresses or Shirts)

Subfamily The second level of product classification, called subfamilia (English: subfamily) associates the product and family to a specific product type (e.g. Woman Dresses or Basic Shirts). Subfamily categorization is subjectively assigned by the commercial buyers, each containing a varying number of articles at a given time. This decision does not involve the distribution demand forecasting team.

Article (MCC) Each product sold by Zara is identified with an article code (MCC) that serves as a unique identifier for a particular article. The article code is comprised of three features: Modelo, Calidad and Color (English: Model, Quality, and Color).

MCCt An individual item is known as an MCCt at Zara where t represents talla (English: size). The MCCt code adds specific size information to the article code and represents

the most basic product identification level.

While final product orders and shipments are placed on an MCCt level, the analysis performed in this work and demand forecasting algorithms focuses on the MCC level.

3.4 Distribution Demand Forecasting

Fashion products have short lifecycles with notoriously volatile demand profiles as shown in Figure 3-1. Zara's products are no exception. Articles can transition from introduction through sold out and discontinued in as little as four weeks. With such short lifetimes, it is critical to distribute inventory as accurately as possible to ensure maximum satisfaction of customer demand.

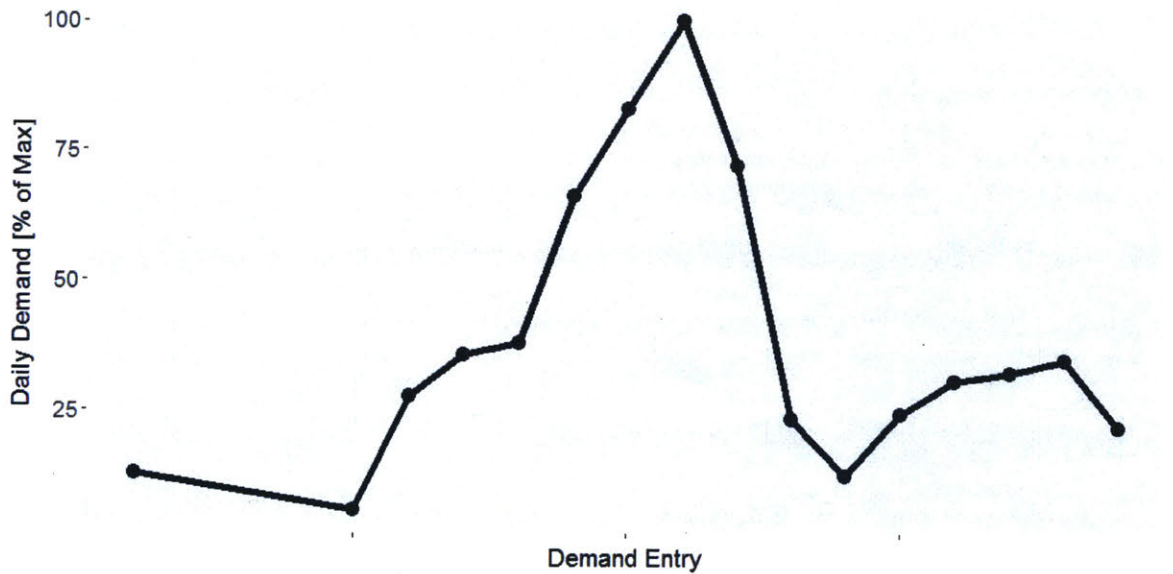


Figure 3-1: An example of demand levels for a single MCC in a single location. The y-axis is normalized to present a percent of maximum article, location demand.

It should be noted that demand forecasting exists in many areas of Zara's operations. Only distribution demand forecasting processes are considered in the scope of this work. Additional processes are applied by different groups for applications including product design and purchasing decisions made by the commercial buyers group and for forecasting demand in the clearance period.

3.4.1 Demand Forecasting Overview

Zara's distribution demand forecasting supports the business need for quick and accurate short term demand forecasts. Every 3-4 days (one forecast period), a new distribution demand forecast is developed for each Zara store and new shipments are sent to restock the inventory required. Due to the volatility of fashion retail demand, forecasts are made for a short term from the expected date of shipment arrival.

Distribution demand forecasts are based on a weighted average of historical sales of the article being forecast in the store being considered. Two distinct distribution shipment types are considered in the forecasting, each with a specific process applied to it. The process applied depends on the type of distribution shipment to be sent: Initial Shipment or Replenishment Shipment. The shipment type is determined by stage of the product lifecycle that the particular (MCC, store) combination is in. Both distribution types are considered in this work.

The general timeline of the forecasting and delivery process for one market is outlined in Figure 3-2. Demand forecasts are generated for each forecast period, or twice per week. The lead time for items varies based on destination, however shipments are scheduled in groups based on location such that most articles arrive in stores on approximately the same day of the week for all stores receiving a shipment. Shipments arrive twice per week corresponding to the forecasting periods, once for weekend sales and once for weekday sales.

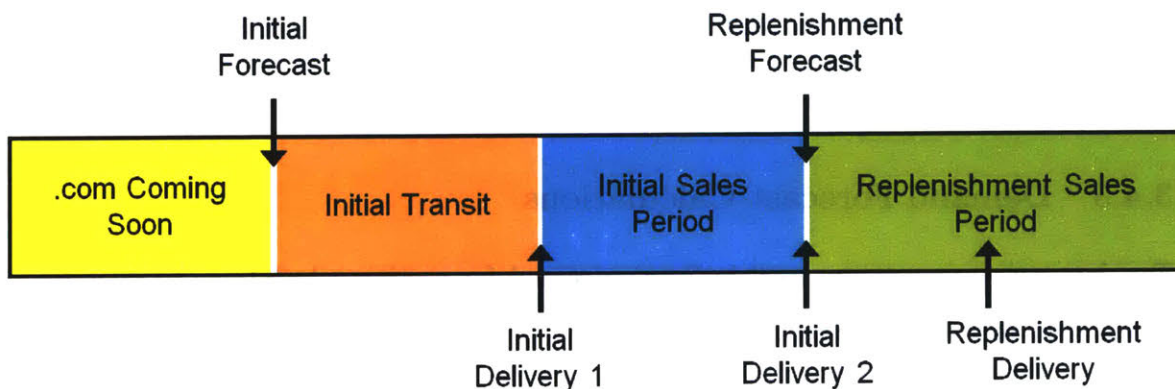


Figure 3-2: The current demand forecasting timeline for new and replenishment product shipments in the physical and e-commerce sales channels

3.4.2 Distribution Demand Forecast Types

Initial Shipment

Initial shipment forecasting is performed when there is insufficient prior sales data for the (MCC, store) combination. This occurs when an article is first introduced to a store and for subsequent forecasts made with less than one week of historical sales data, usually the first two forecast periods.

In this process, similar articles are subjectively defined through discussion with the commercial teams to identify articles that have been sent to the store that are expected to have similar demand profiles to the article in question. Each similar article is assigned a weight and the composite weighted historical demand of the similar articles for a selected time period, usually corresponding to the first two weeks of sales, is used to predict the new article performance. The similar articles are used as proxies for the article until one week of data has been collected and forecasting can be performed using the replenishment shipment process. When there is sales data for the (MCC, store), but for a period of less than one week, that sales data can be mixed in with the similar articles to inform the next shipment.

Replenishment Shipment

The replenishment shipment process describes all forecasts made with at least 1 week of (MCC, store) historical sales data available. In these forecasts, a weighted average of historical sales from the specific (MCC, store) over the previous week is used to determine the shipment quantity.

3.4.3 Demand Forecast Calculations

Two key measurements are used in Zara's demand forecasting calculations to describe an article's historical sales performance.

Días Posibles de Ventas (DPV) Días posibles de ventas (English: days of possible sale) for a given MCC (m) in a specified location (l) of interest represents the number of days that two conditions simultaneously hold true.

1. The store of interest (brick-and-mortar or website) is open for business and can conduct sales.
2. The article of interest has available stock in the sales channel of interest on the date of interest.

Although the time period used for calculation varies in practice, for the purposes of this study, this metric is calculated over the period of a single sales week ($N = 7days$) and is therefore calculated as:

$$DPV_{l,m} = \sum_{t=0}^N open_{l,t} AND(stock_{l,m,t} > 0)$$

Venta Media (Vm) The venta media (English: average sales) is the average sales per day of an article adjusted for days of possible sale. This metric is typically calculated for a given MCC in a specific store location or market aggregation level. The Vm approximates the untruncated customer demand for the article in that it considers the expected average sales rate of an article per day assuming sufficient stock. Demand measurements are therefore only truncated on the day that the stockout occurs and in most cases do not significantly influence the Vm value.

$$Vm_{l,m} = \frac{\sum_{t=0}^N sales_{l,m,t}}{DPV_{l,m}}$$

The Vm measurement is the primary input used to forecast the expected sales of articles in a particular store. For initial shipments, the Vm is calculated for the selected similar articles over a selected time period based on commercial and distribution experience. The weighted average of the Vm is calculated across the group and this is used as the representative article Vm for initial shipment calculations. Replenishment shipments simply use the previous 1 week of average sales in the (MCC, store).

Demand Forecast Accuracy

The accuracy of the predicted V_m is measured for each demand forecast produced by the distribution team. The standard error metric used to measure the error of the demand forecast is the weighted mean absolute percent error (WMAPE). The WMAPE is calculated over a group of articles of interest using the predicted and actual values of the V_m as:

$$WMAPE = \sum \frac{|actual|}{\sum |actual|} * \frac{|prediction - actual|}{|actual|} = \frac{\sum |prediction - actual|}{\sum |actual|}$$

Therefore, a high accuracy demand forecast will have a low WMAPE while a low accuracy forecast will have a high WMAPE.

Comparing WMAPE for the two forecast types using the baseline algorithm yields a higher WMAPE for initial shipments versus replenishment shipments. This is due to the difference in using actual sales data for replenishment calculations rather than similar articles as used in the initial shipment case. Similarly, comparing WMAPE by sales channels demonstrates a higher WMAPE for brick-and-mortar stores versus online FCs. The difference between channels can be explained by the volume of shipments sent to the different locations and the effect of regional demand pooling in the online FCs.

3.4.4 Inventory Shipment Calculation

The final step of the demand forecasting process is the calculation of the units to ship to each store. The objective stock level is the target level of inventory for a given (MCCT, store) combination based on the expected customer demand for the tuple. Objective stock is calculated by modifying the predicted V_m for the given shipment type by a factor representing the transit time and the safety stock (coverage) required for the desired service level of the article in a model. Additional modifications are made by a combination of manual adjustments and automated models to account for new articles, anticipated macro trends like holidays, and the size allocation distribution of the article. The calculated objective stock is then compared with the shipment maximum and minimum criteria. These limits ensure that no store receives a disproportionate number of units while meeting display criteria set by the commercial team to ensure that stores have at least enough inventory to populate a

full in-store display of the article. Therefore, the objective stock level for a given (MCCT, store) is expressed as:

$$Stock_{objective} = \begin{cases} Stock_{calculated}, & \text{if } Stock_{calculated} > \text{Display Minimum} \\ \text{Shipment Maximum}, & \text{if } Stock_{calculated} > \text{Shipment Maximum} \\ \text{Display Minimum}, & \text{otherwise} \end{cases}$$

where

$$Stock_{calculated} = Vm * transit * coverage * size$$

The total units shipped is calculated using the objective stock level and subtracting the total expected stock in the store at the time of shipment arrival.

$$Shipment = Stock_{objective} - Stock_{instore}$$

Finally, each store receives a rank for a subfamily based on its percentage of total weekly sales for the articles in the subfamily. Article shipments are then allocated to stores by subfamily ranking, with the highest ranked stores receiving the calculated shipment first, continuing until the total initial shipment inventory is allocated among the stores. The distribution team releases the proposed shipments to the appropriate commercial teams and modifies the values as needed after considering additional commercial criteria.

THIS PAGE INTENTIONALLY LEFT BLANK

Chapter 4

Methodology

This chapter introduces the specific methodologies applied to the investigation of the cross-channel demand forecasting capabilities at Zara. The methodologies employed were selected based upon the review of current processes as described in literature (see Chapter 2) in addition to processes currently in place within the distribution department at Zara.

4.1 Data Analysis

Successful implementation of a data analytics project requires a standardized process to guide the investigation. For the purposes of this project, the Cross Industry Standard Process for Data Mining (CRISP-DM) was adapted to define the process as shown in Figure 4-1. This approach ensured a standardized process that allowed the business objective to continually guide the development and iteration of the work.

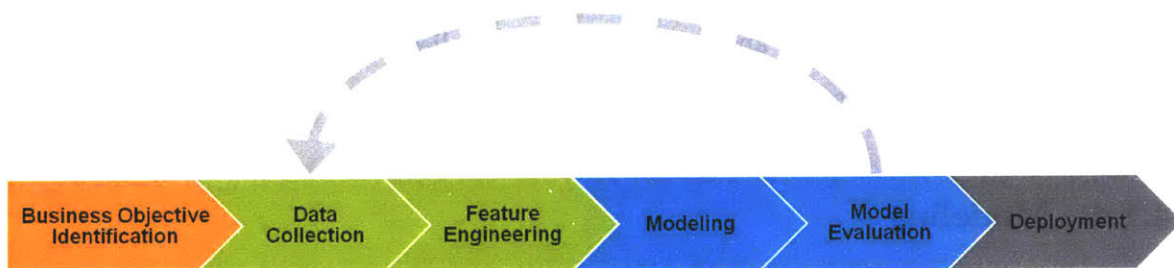


Figure 4-1: The data analysis process adapted from CRISP-DM

4.1.1 Business Objective

The process begins with business understanding, approached in this work as the identification of the business objective and problem formulation as discussed in Section 1.1. As described, the business objective is to provide a high service level to the customer by stocking sufficient inventory in a given location while minimizing operational costs. This work seeks to achieve this by producing a demand forecast with the minimum possible WMAPE for each demand forecast type as described in Section 3.4.

4.1.2 Data Collection

After clearly defining the business objective, datasets must be identified and collected in the Data Collection phase. Datasets are identified through a number of approaches including process research and interviews with experts in addition to an analysis of the datasets directly associated with the business objective. Once identified, datasets must be located and collected if available. If unavailable, datasets can be updated to incorporate specific data points of interest in future data collection efforts for inclusion in analysis at a later time.

4.1.3 Feature Engineering

Feature Engineering is the process of transforming raw datasets into useful features for interpretation by analytical models. In this work, this process is performed using a combination of industry knowledge collected through process analysis and training in combination with preliminary analysis of feature relationships. Additionally, all datasets are cleaned to eliminate erroneous values prior to analysis. Features are collected into feature vectors for each entry of interest, the full set of which comprise the feature space.

4.1.4 Modeling

Using the constructed feature set, the modeling phase seeks to build a descriptive model for a subset of feature vectors comprising a feature space called the training set. The training set is used to develop a model that accurately describes the observed behavior of the data.

Subsequently, the model developed using the training set is then applied to the remaining feature vectors that were held back in a blind test. This dataset represents unseen data and demonstrates the model's ability to predict the outcome of new data.

4.1.5 Model Evaluation

After models are built, trained, and tested around the feature set, model performance is evaluated to determine the usefulness of the prediction. Multiple model types as well as features that are included in the feature vectors can be adjusted in an iterative manner to converge on the best performing solution for the problem posed. Additionally, this phase must develop a metric that links the output of the model to the business objective identified in the Business Objective phase of the process.

4.1.6 Deployment

After the best performing model for the business objective of interest has been developed and identified, the model can continue to the deployment phase. In this phase, not directly addressed in this work, models are integrated into the standard business processes for testing in a real environment. Performance of the new process is carefully monitored in the initial implementation stages of this deployment and compared to the expected performance of the model as identified in the Model Evaluation phase of the process. Model performance should continually be checked against expectations and feature vectors should be updated to reflect changes in business processes and deficiencies in model performance.

4.2 Data Sources

Zara maintains a number of key databases storing data collected from various organizations within the firm. In particular, datasets were sourced from both the Zara.com and Distribución departments. The data collected describe the customer behavior on the Zara.com website, sales in both retail channels, product and store details and inventory information. Details regarding the individual datasets used in this analysis are discussed in Chapter 5.

THIS PAGE INTENTIONALLY LEFT BLANK

Chapter 5

Data Collection

Data collection represents a critical step in data mining and analysis. In this study, databases across the Zara organization were surveyed to identify the datasets of greatest potential value for demand forecasting. Through interviews with stakeholders and experts across organizations, a mapping of available databases was constructed to guide the data selection process and feature engineering phase of the project.

5.1 Data Identification

Zara maintains extensive databases describing the daily operations of the firm with high levels of detail and accuracy. To begin the data collection process, a preliminary investigation to identify potentially useful datasets was conducted. The proposed collection of datasets represented the data required to understand the customer preferences of both Zara.com and brick-and-mortar store customers.

5.1.1 Process Research

To clarify the data requirements that support Zara's demand forecasting, interviews were conducted throughout retail operations. The goal of these interviews was to develop an end-to-end concept of the process of identifying and fulfilling customer demand for a given MCC. For specific process descriptions regarding demand forecasting methods, see Section 3. The

identification of significant data points through general process understanding is discussed here. Interviews began with customers and customer-product interfaces in the brick-and-mortar stores and e-commerce channels to understand the process used by the commercial teams to determine which products to stock in each store. Process understanding was combined with hypotheses from other stakeholders to identify datasets that might provide additional customer preference insight to the Distribution demand forecasts.

5.1.2 Brick-and-Mortar Stores

In-store Purchases

In brick-and-mortar stores, customers enter the store without prior knowledge of the articles available but with a general sense of their personal fashion desires. Customer demand is therefore identified and fulfilled in two primary ways.

In the case of independent customers, the customer browses the sales floor of the store to identify products that match their preferences. Product preferences are determined by product features including type, style, material, color, fit and price. If a product is found to match the customer's preferences, the customer purchases the item and sales data is recorded in the sales database. If products exist but are not found in the store, the customer demand is left unrecorded.

Some customers engage with store associates to determine if available inventory matches their preferences. Store associates work with customers to determine the products the customer is interested in based on product features and comparison products. Through this process, the associate develops knowledge of the individual customer's preferences and matches those preferences with the available inventory in store. If products are found to match the preferences of the customer, the customer may purchase the articles. If an article is not found in the store, the associate may help the customer locate a suitable match at another store or on *Zara.com*. If no channels offer an article that matches the customer's preferences, no sale is recorded but the store associate remembers the preferences of the customer and communicates them to the regional Product Manager either directly or through the store manager during regular meetings.

Product Manager Forecasting

When forecasting demand, Product Managers combine sales data, personal experience with the region and knowledge of fashion trends, and reports from sales associates regarding customer preferences to inform their decisions. Additional information including store display configurations and layout are considered in discussions between the Product Manager and the store employees to inform expected customer purchases as well.

5.1.3 Zara.com Experience

At Zara.com, customers freely browse the storefront via one of three access methods: Zara's mobile application, mobile optimized website, and the desktop optimized website. All three methods directly access the specific interface developed by the Product Manager for the customer's region. The e-commerce experience does not give the customer access to store associates, thus each customer behaves as the independent customer in a brick-and-mortar store.

Website Navigation

The storefront allows the customer to browse available articles through the navigation tree, a tool that assigns categories to each product and allows customers to display products within the category of interest. Products may be assigned to multiple categories. Product filters can be applied to limit the articles displayed, for example selecting specific product subcategories, price limits, colors, and sizes. A special product category called Última Semana (English: New In) is placed at the top of the navigation tree. The New In category contains products released in the last two weeks and represents the most current trends.

The customer also has the option to search for their desired article using keywords entered in the search field. The Zara.com team maintains an optimized search results algorithm to display the most relevant article results based on the keywords entered. This function seeks to partially serve the role of an in-store associate suggesting products based on a discussion of the customer's interests. Figure 5-1 shows an example of a typical Zara.com storefront.

Articles from the selected category or search results are displayed in a grid of either two

ÚLTIMA SEMANA MUJER

ÚLTIMA SEMANA

ABRIGOS
 CHAQUETAS
 BLAZERS
 VESTIDOS
 MONOS
 CAMISAS
 BODY
 PUNTO
 PANTALONES
 JEANS
 FALDAS
 GAMBETAS
 BUDASERAS
 ZAPATOS
 BOLSOS
 ACCESORIOS
 PERFUMES
 TARJETA REGALO
 JOIN LIFE
 MONDAY TO FRIDAY
 CASHMERE
 SPECIAL PRICES

#pearls TRENDS
 #stripes TRENDS
 #vichy TRENDS
 REBAJAS



NEW
CROP TOP ESTAMPADO FLORES
29.95 EUR



NEW
ABRIGO MANGA VOLUMINOSA
89.95 EUR



NEW
VESTIDO ESTAMPADO FLORAL
39.95 EUR



NEW
TOP ESTAMPADO FLORES
29.95 EUR



NEW
FALDA ESTAMPADA FLORES
29.95 EUR



NEW
FALDA ESTAMPADA FLORES
29.95 EUR



NEW
ZAPATO TACÓN PIEL AMARILLO
49.95 EUR



NEW
BLAZER MANGA ABULLOHADA
49.95 EUR



NEW
PANTALÓN ANCHO CON CINTURÓN
29.95 EUR



NEW
JERSEY ABERTURAS LATERALES
12.95 EUR



NEW
PANTALÓN ANCHO RAYAS
39.95 EUR



NEW
CAMISA RAYAS Y PARCHES
39.95 EUR

Figure 5-1: An example website layout of the Zara.com storefront with navigation tree, search bar and filter options

or six products per row as selected by the customer. Each article is represented by a photo chosen by the commercial team to best display the design aesthetic of the article. All of the articles that match the desired criteria are displayed in a single scrollable page. Articles with multiple color options representing multiple MCC have a separate entry for each color option. For more information or to make a purchase, the customer must select one of the options displayed in the results.

Article Details

Each MCC has a dedicated page with additional photographs and details about the article including a written article description, color, materials and care, sizing information, price, and product availability. The option to view a product Guía de Tallas (English: Size Guide) is provided for customers that desire more information on the fit of a particular article. If a product is produced in but not currently available in a given size, that size is greyed out on the page.

The Product Manager may decide to indicate that additional stock is coming with an envelope symbol next to the unavailable size. By clicking this symbol, the customer will have the option to enter an email address to be notified when the article is in stock. These subscriptions can be made for articles that were once in stock but are currently out of stock, or for articles that will be offered for the first time soon. Figure 5-2 shows an example product details page.

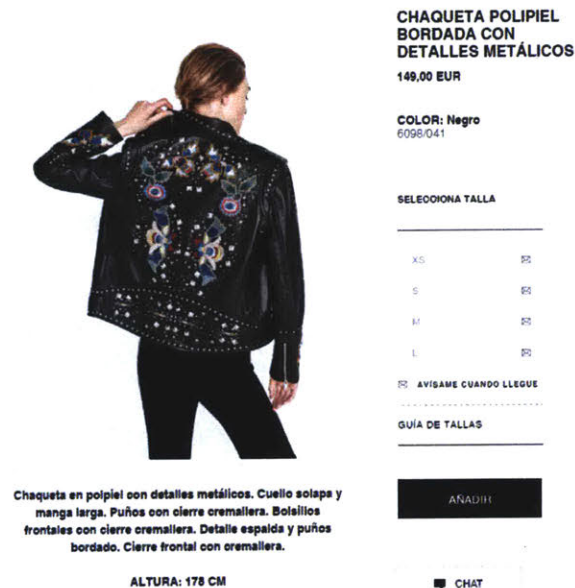


Figure 5-2: An example product details page with coming soon indicators

The product details page also provides customers with the option to check for article availability in a brick-and-mortar store. When the option is selected, the customer is asked to provide the size desired and location information (zipcode, address, city, etc.). The website will then provide the customer the current inventory information for the article at the closest stores to the location requested. The results include the level of inventory and store details including address and phone number.

Purchase Completion

If the customer wishes to purchase an item, they have the option to add the item in a given size to their cart and continue browsing for more items or continue to purchase only the selected item. The cart allows the customer to view the articles they have selected during

their session, adjust the number of units of each MCCt desired, remove an MCCt entirely from the cart, or to continue to complete the order.

Orders are completed by either logging into or creating a *Zara.com* account or by using the checkout as a guest option. If logged in, customer details including address and payment information from previous orders is pre-filled and the current order will be added to the customer's purchase history. When checking out as a guest, essential details are requested for payment processing and shipping information and order details are sent to the email address provided. When checking out as a guest, the order is not associated with any customer account or purchase history.

Product Manager Forecasting

Zara.com Product Managers use similar information as brick-and-mortar Product Managers for estimating customer demand. Generally, past sales of the same or similar articles are considered in addition to subscription information, the organization of the website, and information on trends in the specific region. As e-commerce sales can vary widely throughout the sales day, Product Managers also consider same day sales trends using real time sales data to catch very recent trends. Because *Zara.com* stocks a wide variety of products, the top performing articles receive particular attention to ensure adequate inventory is available to satisfy expected demand.

5.2 Dataset Availability

As previously described, *Zara* maintains vast collections of data to inform business operations. The datasets of interest were determined based upon the review of the Product Manager and customer experience processes outlined in Section 5.1. Interviews with stakeholders and the commercial team product experts yielded a preliminary set of datasets identified for collection. Once datasets of interest were identified, dataset availability was assessed using the database map, as shown in Figure 5-3. The datasets collected can be categorized by the organization that prepared them for this study: Distribution Datasets, E-commerce Datasets, and Supplemental Datasets.

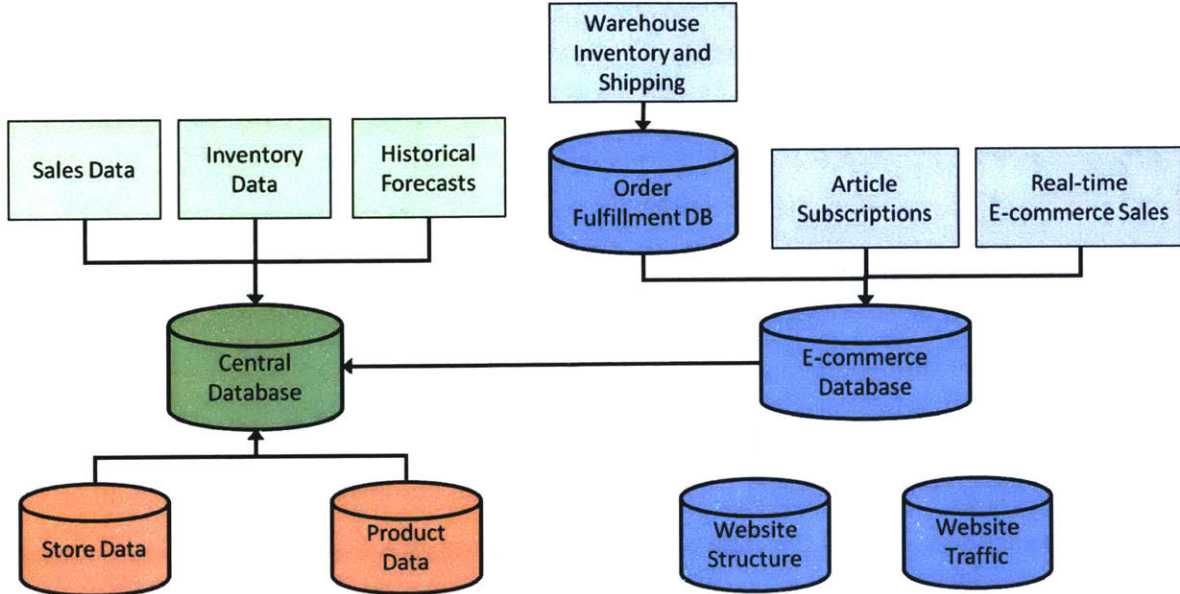


Figure 5-3: A representative database map depicting the available databases and selected tables of datasets contained in the database. Connections are shown with directional arrows, no arrow means there is no direct connection.

5.2.1 Distribution Datasets

The Zara Distribution team maintains datasets reporting sales and inventory data to inform store operations and distribution decisions. As discussed in Chapter 3, demand forecasting processes currently utilize the article sales data for each (MCC, location) 2-tuple. For replenishment products, demand is determined by same article past sales, while for initial product offerings demand is determined by similar product past sales. This data is used to produce a Vm prediction for the (MCC, location) in the current forecasting algorithm.

Inventory shipment data consists of shipment proposals (the output of Distribution team algorithms) and actual shipment volumes for each store. Additional inventory data records the daily inventory level for each MCCt in each location including distribution warehouses and in transit to stores. Inventory data is particularly useful in determining when customer demand measured by sales numbers is truncated due to stock-outs.

Finally, Distribution also uses store opening data to track the days that articles were available for sale. Store closures are not standardized across markets or regions, and therefore must be recorded individually for each store to ensure accurate demand forecasts based on

DPV and *V_m* calculations.

5.2.2 E-Commerce Datasets

In addition to the standard sales and inventory data that is used by Distribution, Zara.com collects data specific to website operations. Two of the most significant datasets include website structure and traffic data.

Website structure data is recorded to perform quality checking on the website storefronts across the Zara.com organization. The structure of the storefront is recorded as links connecting each page to another. The result of these links is a map of source and sink page nodes and the associated links within the website. Additionally, relative display position data of the links on each page is recorded to determine the layout of each page within the website. This layout information stores a snapshot of the positioning of each category in the navigation tree, as well as the positioning of each article within each of the navigation categories.

As described in Section 2.3 website traffic analysis is a common practice for internet businesses to better understand company performance. E-commerce businesses frequently use this information to measure customer acquisition and retention, and to inform marketing decisions. Zara.com collects anonymized website traffic information datasets that record the links that are selected on each page of the website and the order they are clicked to recreate a customer's path through the storefront. All website traffic information is time-stamped with the time and date of access.

5.2.3 Supplemental Datasets

Supplemental datasets are used in some Distribution demand forecasting and inventory movement calculations, but not as directly as the sales and inventory data discussed in Section 5.2.1. These datasets are typically used to group stores or articles for analysis or to extend the capabilities of other algorithms.

Store Data

Each retail sales location maintained by Zara has an associated descriptive set of data. Basic location information including street address, city, region, and postal code are all stored in a database referenced to the stores. Additionally, store configuration information regarding the sales floor and stockroom inventory capacities are recorded for inventory planning purposes.

Article Data

Each article has a number of data points associated with it. The product code represented by the MCCt records information on the model, quality or material, color, and size of the article. Additionally, articles are categorized into subfamilies as described in Section 3.1.

THIS PAGE INTENTIONALLY LEFT BLANK

Chapter 6

Feature Engineering

Feature Engineering is a critical step of data analysis. The principal objective of feature engineering is to transform raw datasets into useful inputs for a predictive model. In the context of this investigation, this includes data filtering to constrain volume and more accurately simulate operations, data cleaning to ensure accuracy, and various transformations to convert data into interpretable fields. Feature engineering was performed iteratively with modeling. New features were developed both based on opportunity when new ideas for potential performance improvement were identified, and on need when model validation performance was below expectations. In total, more than 45 features were constructed for analysis using raw datasets collected from all data sources listed in Chapter 5 to lend insight into customer preferences.

6.1 Data Filtering and Cleaning

6.1.1 Filtering Parameters

Initial filtering of the datasets was performed to accommodate processing platform and operational process constraints. Due to the large volume of data collected in both retail sales channels as described in Chapter 5, the analysis was focused on a specific subset of demand forecasts.

Market, Product and Period Restrictions

A single commercial market was selected for analysis. The market was selected based on criteria including total market size of both brick-and-mortar and e-commerce sales, the total number of store locations in the market, the number of e-commerce fulfillment centers serving the market, and sales growth trends of the market in both channels. The market selected was a top growth and sales area for e-commerce, providing a suitable e-commerce dataset for analysis of customer preferences. Additionally, the market was within the top markets globally for brick-and-mortar sales.

The sales period used for feature creation was selected based on the period that datasets were available from the Zara.com website traffic database at a sufficient level of detail. Datasets were restructured during the period of collection affecting comparability between datasets taken at different time periods. Therefore, a sales period of interest was defined to be a single sales quarter of a single sales season. This period of interest did not include any sales during the Rebajas (English: Clearance) period which might affect the significance of predictive variables in the analysis due to the effects that clearance pricing have on customer preferences and behavior. This period of interest covers a total of 19 inventory shipment cycles of demand forecasting.

As noted in Chapter 1, Zara introduces about 9,000 new articles each season. All articles are not distributed to all brick-and-mortar stores, but selectively based on expected demand. Zara.com offers all of the available articles since the e-commerce channel serves the entire market. The data associated with this extreme product variety was therefore limited to a subset of products believed to be representative of the greatest opportunity for demand forecasting improvement. This analysis focuses on data from 58 distinct subfamilies of products in Zara's Señora (English: Women's) department. The 58 subfamilies selected cover more than 2500 different MCCs.

A summary of available data points of interest post filtering is provided in Tables 6.1 and 6.2.

	Brick-and-Mortar	E-Commerce
MCC	1347	135
Dates	76	42
Subfamilies	55	38

Table 6.1: Summary of available data for initial shipment analysis.

	Brick-and-Mortar	E-Commerce
MCC	2636	2588
Cycles	19	19
Subfamilies	58	58

Table 6.2: Summary of available data for replenishment shipment analysis. MCC represents the number of articles available, Dates or cycles are the number of distinct forecast dates or periods, and Subfamilies is the number of subfamilies represented by the articles

6.1.2 Data Cleaning

All selected datasets were cleaned prior to analysis to avoid introducing biases due to missing values, false zeros, invalid entries and other detectable errors in the datasets. Feature vectors with missing values or zeros representing missing data for the datapoints required for feature construction were eliminated. Zeros were preserved and, in some cases, added where the existence of a zero value was known to be real instead of missing (e.g. zero sales for an item that is in stock). Invalid data entries include data points of an incorrect type, incomplete or incorrect information (e.g. invalid address, zipcode combinations) and were eliminated from the dataset to avoid confusion. All remaining values were reformatted to provide consistent data types to processing algorithms.

6.2 Feature Construction

Features were created using raw datasets described in Chapter 5. The goal of feature engineering is to construct features, using the raw datasets stored in the databases, that can be used to more effectively describe the state of the system to the model. Raw data points collected from the databases are transformed into feature vectors \mathbf{F} for use in predictive algorithms, where \mathbf{F} with n elements is defined as:

$$\mathbf{F} = [x_1, x_2, x_3, \dots, x_n]$$

Each feature vector is identified by a set of key values to distinguish between independent entries. Each key set describes a specific demand forecast type using the article (MC), color, time (representing the date the forecast is made), and location information as $[MC, C, time, store]$ called an entry. The full set of feature keys represents the range of the forecast in this investigation. The feature vectors associated with the full set of entries comprises the feature set and describes the full dataset used in the modeling activity.

6.2.1 Notation

Common notation for feature construction equations is noted below.

T = The total number of time periods applicable or desired for the entry of interest.

t = A single time period where $t = 0$ represents the date that the demand forecast is prepared.

C_t = The total number of categories that an article is found on the Zara.com website at time t .

c = A single category that an article can be found on the Zara.com website.

S = The total number of manufactured sizes for a given article.

s = A single size of a given article.

6.2.2 Article Features

Article features describe some of the features of an article that a Product Manager would consider in evaluating expected demand for an article.

Subfamily The Subfamily feature is directly pulled from raw article subfamily data as a subfamily code.

Colors The total number of distinct color options available for a given MC as calculated on the day prior to forecasting across the entire market of interest.

Price The price of the MCC pulled directly from the article information database in units of the local market currency for relative article comparison.

Total Stores The total number of distinct stores that the MCC combination has been shipped to since introduction inclusive of the key store of interest through the day before the forecast is made.

6.2.3 Location Features

Region The region of interest for the given entry. This region corresponds to the region of the physical region of the store or e-commerce FC location, the region of website traffic and the region e-commerce sales were shipped to for order fulfillment.

City The city of interest for a given entry. This city corresponds to the physical city of the store or e-commerce FC location.

Rank The internal sales rank of the store of interest in the department of the MCC of interest.

Mall A binary indicator extracted from store address and description datasets indicating whether the store of interest is located within a shopping mall complex or not.

Floor Capacity The number of units of inventory, irrespective of type, that can be stored on the sales floor of the store of interest.

Stockroom Capacity The number of units of inventory irrespective of type that can be stored in the stockroom of the store of interest.

Total Capacity The total inventory capacity of the store of interest calculated as $Capacity_{floor} + Capacity_{stock}$.

6.2.4 Temporal Features

Weekday The day of the week that the demand forecast is produced. This feature also reflects the day of the week that new inventory will first be available for sale assuming a constant transit time between forecast date and receipt at the store.

Initial Date The initial offering date of the MCC in the store of interest.

Initial E-Commerce Date The initial offering date of the MCC on Zara.com for the given market.

E-Commerce Lead The number of days between the initial offering dates of the MCC on Zara.com and in the store of interest calculated as $InitialDate_{store} - InitialDate_{com}$. This number represents the number of days an article was available for sale online before it was available in the brick-and-mortar store.

Article Age The number of days between the initial offering date of the MCC in the store of interest and the demand forecasting date calculated as $ForecastDate - InitialDate_{store}$. This number represents the total number of days an article has been available for sale in the store of interest, and is therefore only available for replenishment shipments.

6.2.5 Website Structure Features

Average Best Position (Cumulative and Prior Week) The average over the specified time period of the daily best position of the MC across all of the category pages found on the website. Time periods available include cumulative since the product was first introduced on Zara.com, and for the previous week beginning the day before forecasting is performed (only available for replenishment shipments). Calculated as:

$$\frac{1}{T} * \sum_{-(T+1)}^{t=-1} \min\{position_0, position_1, \dots, position_{C_t}\}$$

Average Worst Position (Cumulative and Prior Week) The average over the specified time period of the daily worst position of the MC across all of the category pages found on the website. Time periods available include cumulative since the product was first introduced on Zara.com, and for the previous week beginning the day before forecasting is performed (only available for replenishment shipments). Calculated as:

$$\frac{1}{T} * \sum_{-(T+1)}^{t=-1} \max\{position_0, position_1, \dots, position_{C_t}\}$$

Average Mean Position (Cumulative and Prior Week) The average over the specified time period of the daily mean position of the MC across all of the category pages found on the website. Time periods available include cumulative since the product was first introduced on Zara.com, and for the previous week beginning the day before forecasting is performed (only available for replenishment shipments). Calculated as:

$$\frac{1}{T} * \sum_{-(T+1)}^{t=-1} \left[\frac{1}{C_t} * \sum_{c=0}^{C_t} position_t \right]$$

Average Categories (Cumulative and Prior Week) The average number of categories that that MC can be found under on the Zara.com website calculated over the specified time of interest. Time periods available include cumulative since the product was first introduced on Zara.com, and for the previous week beginning the day before forecasting is performed (only available for replenishment shipments). Calculated as:

$$\frac{1}{T} * \sum_{-(T+1)}^{t=-1} C_t$$

Days New (Cumulative and Prior Week) The number of days that the MC was found in the New In category on Zara.com during the time period of interest. Time periods available include cumulative since the product was first introduced on Zara.com, and for the previous week beginning the day before forecasting is performed (only available for replenishment shipments).

Average In-stock (Cumulative and Prior Week) The average percentage of sizes of an MC available in stock on Zara.com over the time period of interest. Time periods available include cumulative since the product was first introduced on Zara.com, and for the previous week beginning the day before forecasting is performed (only available for replenishment shipments). Calculated as:

$$\frac{1}{T} * \sum_{-(T+1)}^{t=-1} \left[\frac{\sum_{s=0}^S instock_s}{S} \right]$$

6.2.6 Customer Behavior Features

Venta Media (V_m) The predicted average sales output from the Distribution team demand forecasting algorithms. This feature incorporates same-store sales information from prior weeks, adjusted for inventory levels and days of sale as described in Section 3.4. Using this output as an input feature for the model allows the proposed algorithm to operate in cooperation with current processes.

Prior E-Commerce Sales E-Commerce sales through Zara.com recorded during the time period of interest. Aggregated at the MC and MCC levels for each region of interest. Time periods available include cumulative since the product was first introduced on Zara.com and for the previous week beginning the day before forecasting is performed.

Views The number of times the article details page was viewed for an MC of interest across all available pages for the MC on Zara.com over the time period of interest for the region of interest. Time periods available include cumulative since the product was first introduced on Zara.com and for the previous week beginning the day before forecasting is performed.

Add to Cart The total number of times an MC of interest was added to a customer cart on Zara.com over the time period of interest for the region of interest. Time periods available include cumulative since the product was first introduced on Zara.com and for the previous week beginning the day before forecasting is performed.

Remove from Cart The total number of times an MC of interest was removed from a customer cart on Zara.com over the time period of interest for the region of interest. Time periods available include cumulative since the product was first introduced on Zara.com and for the previous week beginning the day before forecasting is performed.

Net Cart The net result of cart additions and removals for a given MC in the region of interest over the time period of interest calculated as $Cart_{add} - Cart_{remove}$. Time periods available include cumulative since the product was first introduced on Zara.com and for the previous week beginning the day before forecasting is performed.

Check In-Store The total number of times that the Check Availability In-Store option was selected on the article details page for a given MC by a customer in the region of interest over the time period of interest. Time periods available include cumulative since the product was first introduced on Zara.com and for the previous week beginning the day before forecasting is performed.

Size Guide The total number of times that the Size Guide option was selected on the article details page for a given MC by a customer in the region of interest over the time period of interest. Time periods available include cumulative since the product was first introduced on Zara.com and for the previous week beginning the day before forecasting is performed.

Subscriptions The total number of customer subscriptions made for a currently unavailable article on Zara.com over the time period of interest. A subscription indicates a customer request to be alerted when the article comes in stock but does not necessarily signal an intent to purchase. Subscription data is aggregated at the MC and MCC article levels and at the region and market location levels. Time periods available include cumulative since the product was first introduced on Zara.com and for the previous week beginning the day before forecasting is performed. Subscription data only exists for those articles that have a subscription option on the Zara.com website.

THIS PAGE INTENTIONALLY LEFT BLANK

Chapter 7

Modeling

The feature datasets described in Chapter 6 represent the inputs used to develop predictive models for the actual customer demand. Each feature vector $\mathbf{F}_{MCC,t,l}$ is paired with an associated actual demand value measured by the actual average sales per day and expressed as $VmReal_{MCC,t,l}$. Therefore the modeling effort seeks to describe the relationship in each $\{\mathbf{F}_{MCC,t,l}, VmReal_{MCC,t,l}\}$ set by identifying a function $f(\mathbf{F})$ assuming no prior knowledge to describe the underlying relationship between the feature vectors and associated outcomes such that, for all MCC, t, l in the set:

$$f(\mathbf{F}_{MCC,t,l}) = VmReal_{MCC,t,l}$$

7.1 Modeling Approach

A standard modeling workflow was applied to the demand forecasting model investigation. In the workflow, potential modeling algorithms are applied to a subset of the entries in the Feature Set. Each algorithm is tuned to optimize the performance of the method for comparison against other potential modeling algorithms. The best performing model is then used on a new dataset to quantify the model's performance on unseen data.

When developing a model to describe a predictive relationship, the dataset that the model is developed (trained) on should be distinct from the dataset that the model's performance is ultimately reported (tested) against to the extent possible. To ensure this separation, Feature

Sets were separated into three subsets describing the steps in the modeling workflow.

Training Data is data used for developing the initial descriptive model for a given algorithm. When the model is trained, the algorithm minimizes the error between the predicted and actual values of the dependent variable (here *VmReal*) for each entry.

Validation/Development Data is data that is used to evaluate intermediate performance of the trained models on unseen data. The validation dataset is used to tune parameters of the modeling algorithms being applied to maximize performance, and to compare the performance of different algorithms to each other to select the final candidate for the modeling effort.

Test Data is the set of entries that is isolated from all model development until the final stage performance evaluation. Test data represents new data that will be presented to the model in the deployment phase of the project and therefore must not be used to optimize but only to report model performance.

7.2 Data Pre-Processing

Prior to modeling, the feature set was pre-processed to prepare the data to match a simulation of the standard distribution operations process. For the demand forecasting process, pre-processing was performed during feature engineering by adjusting the values in each feature vector to represent only information available at the time when the demand forecast was made (t). For most features this required all data only include values through $t - 1$ as described in Section 6.2.

After the feature vectors were properly prepared to simulate the information available for the demand forecast type, entries without associated website features were eliminated from the Feature Set. Additionally, the feature set was split into subscription and normal feature sets. All entries that had associated subscription information (including 0 where applicable) for unavailable MCCt on the Zara.com website were included in the subscription dataset. All feature vectors with website data including those with subscriptions were included in the standard feature set but subscription data was ignored.

The combination of pre-processing filters results in reduced datasets. In Initial Shipment forecasting, website data must be available for a given item at least 7 days in advance of the first date of sales for the same MCC in the store of interest. This condition is particularly restrictive for Zara.com forecasting since it only can be true for those MCC that are sold by another FC or posted on the website as coming soon at least 7 days in advance of arrival at the e-commerce FC of interest. Similarly, the data for each Replenishment Shipment entry must be available at least 4 days prior to sale in the store of interest. A summary of data entries available for each demand forecasting type is shown in Tables 7.1 and 7.2. The % discarded indicates the number of entries that did not have either website or subscription data available as of the forecasting date.

Initial Shipment: Applicable Data Entries		
Channel	Website Data % Discarded	Subscriptions % Discarded
Zara.com	89.40%	99.03%
Brick-and-Mortar	75.93%	99.03%

Table 7.1: Percentage of applicable initial shipment data entries after forecast type filtering for valid website and subscription data.

Replenishment Shipment: Applicable Data Entries		
Channel	Website Data % Discarded	Subscriptions % Discarded
Zara.com	9.03%	91.61%
Brick-and-Mortar	14.14%	89.95%

Table 7.2: Percentage of applicable replenishment shipment data entries after forecast type filtering for valid website and subscription data.

7.2.1 Feature Reduction

Prior to modeling, the feature set was analyzed for significant features looking at feature variance and collinearity. Features that did not pass initial tests based on these measures were marked for removal consideration from final tested models. Marking these features allowed for feature prioritization identifying the features that are most likely to introduce noise or provide limited value in model prediction. It is advantageous to remove features that introduce unnecessary complexity or noise to a model to enable faster computation with limited system resources while preserving maximum model performance.

Zero and Near-Zero Variance

Analysis of feature variance is used to evaluate the relative amounts of information contained in each of the features. Features with zero or near-zero variances are less likely to contain useful information for the modeling algorithm and can therefore often be removed from the feature vectors in the feature set. Features with zero variance are those that have a single constant value in each feature vector and therefore offer no additional information to a model. These features make the feature vector unnecessarily large and can easily be removed from the dataset. There were no zero variance features found in the final feature set.

A less straightforward case of questionable information value is the case of a near-zero variance feature. Near-zero variance features contain only a limited number of unique values relative to the total number of samples, and the relative frequency of the most common value is much higher than that of the second most frequent value. The threshold values for each condition were based on literature research and began with a unique value percentage cutoff of 10% with a frequency ratio of 20:1 [19].

The features marked as near-zero variance from the feature set included Days New (Prior Week) and Average In-Stock. These variables were flagged for consideration, however were not removed from models immediately. The Days New variable was very frequently zero for replenishment shipments, especially in the e-commerce channel. Similarly, the Average In-Stock was most often 1 because most items were almost fully stocked at Zara.com during a given time period. Frequency ratios for these variables were re-run using a binary condition to determine whether sufficient variance existed to justify using the features in the models. All features were found to satisfy the criteria, however were kept for further consideration in the modeling phase. Table 7.3 shows an example frequency table for the values in the Days New variable of one test.

Days New (Prior Week)								
Value	0	1	2	3	4	5	6	7
Frequency	889654	40698	7206	13219	17844	31468	34279	16332

Table 7.3: Example frequency ratio calculation for Days New (Prior Week)

7.3 Modeling Algorithms

Models described in this section were built with the R 3.3.2 programming language [20] using the randomForest package [21], party package for regression trees using ctree [22], and the caret package for training and testing models [23]. Necessary data pre-processing not performed natively in a database was performed in R made extensive use of the dplyr [24] and tidyr [25] packages for data manipulation.

7.3.1 Initial Feature Investigations

A number of algorithms exist for use in explaining the relationships between pairs of predictive feature vectors and associated outcomes. When the underlying pattern of the relationship is well understood, traditional linear or non-linear regression techniques are excellent choices. These techniques are preferred for the simplicity of implementation and interpretation. Additionally, these algorithms are not computationally intensive and can therefore be run quickly against large datasets.

Feature relationships to associated outcomes were investigated to identify linear or simple non-linear relationships in the training dataset where available. No significant linear relationships were found between any single feature or unclustered combination of features in the feature set and the VmReal measurements from initial testing.

7.3.2 Algorithm Selection

Ordinary Least Squares Linear Regression

Ordinary least squares (OLS) linear regressions were run on each dataset to test the performance of a linear model using e-commerce data. Tests were run on the OLS linear model and subsequent models in parallel to determine the tradeoff between model complexity and performance. Linear models were tuned using 10 fold cross-validation with 3 repeats to minimize bias. Feature sets were trimmed using results from parallel tests as described in Section 8.1. Final trimmed feature sets were used to determine OLS linear model performance.

Random Forest

Due to the non-linear and uncertain nature of the underlying relationships between the features and actual demand, an algorithmic modeling approach was adopted. In particular, a random forest algorithm was used to assess applicability of the feature set to prediction of actual demand in both channels under both initial and replenishment product shipments.

The random forest algorithm was identified as the best potential candidate balancing available computing resources, training time, and predictor strength requirements. Although a strong predictor, a critical drawback of using the random forest algorithm in this context is the lack of interpretability as compared to linear regression and simple decision trees.

At a high level, the random forest algorithm is an ensemble predictor method in which a number of decision trees are generated based upon random subsets selected with replacement from the training data provided. With every subset selected, the remaining data is called the out-of-bag (oob) data and is used for internal algorithm error calculation. Similarly, as each tree is grown a random subset consisting of m features sampled from the set of M features in the feature set (s.t. $m \ll M$) is selected at each decision node and the best split is used in each case. Each tree is then grown to the full extent possible without pruning. The ensemble of trees is used as a collective to produce the predicted output for the regression as the average over the forest. [26]

In the random forest algorithm, two parameters can be tuned to optimize performance. The first parameter is the number of trees (ntree) that are generated that constitute the ensemble predictor. By varying ntree, the number of decision trees generated is varied. The value of ntree has been shown to minimally affect model accuracy in previous studies with a high ntree not contributing to overfitting errors [27]. Therefore, an ntree value of 101 was used for modeling purposes in this study after limited experimentation of values ranging from 51 through 501.

The second parameter is m , the number of features randomly selected and evaluated at each decision node in each tree. The value of m is held constant for each tree generated in the forest. Two effects related to the selection of m strongly contribute to the error rate of a random forest: the correlation between any of the constituent trees and the predictive

strength of each of those trees. By reducing the value of m , both the correlation between trees and the strength of the trees is reduced. Therefore, the value of m can be optimized to appropriately balance these two effects for minimum prediction error in the model. Values of m were tuned on each training dataset and evaluated on the OOB error of each test to identify the appropriate value. An example of an m tuning experiment is shown in Figure 7-2.

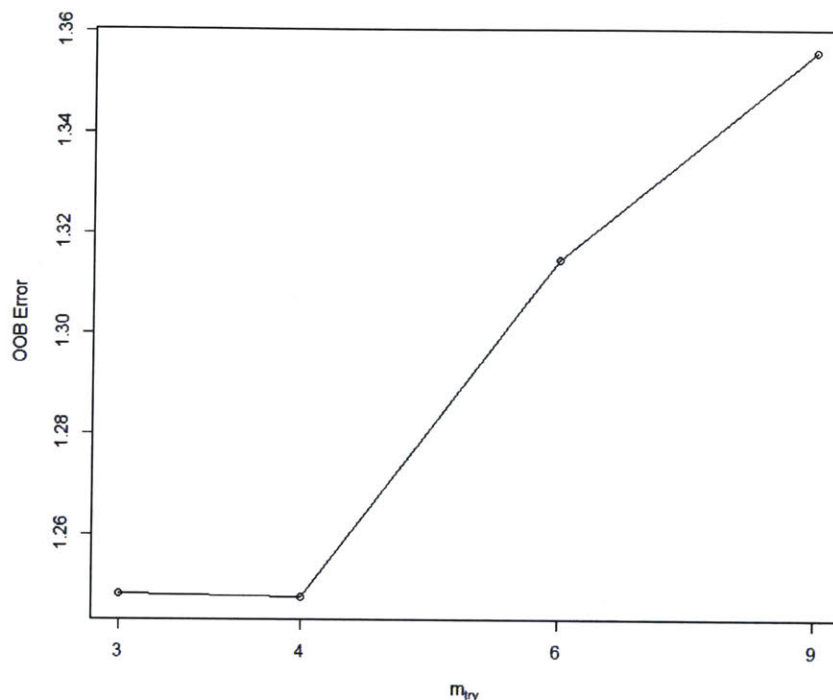


Figure 7-2: Example tuning of the value of m to produce minimum out-of-bag error in the training dataset.

The values of m used for each test are summarized in Table 7.4.

Channel	Initial m Value	Replen m Value
Zara.com	9	6
Zara.com w/ subs	10	10
Brick-and-Mortar	12	10
Brick-and-Mortar w/ subs	9	11

Table 7.4: Initial shipment demand forecast WMAPE improvements with and without subscriptions

7.3.3 Model Adjustments and Prediction Modifications

The predictions made by the models incorporating e-commerce features were determined to consistently over-estimate customer demand. This error was observed most frequently in very low (near-zero or zero) demand entries and most frequently with brick-and-mortar store replenishment shipments. Upon investigation, the overestimation error was determined to be a result of rounding error and a heuristic was introduced to the model. Since V_m is calculated as the average sales over a period of one week, the minimum non-zero real value for the actual V_m is 0.14 or $\frac{1}{7}$. Therefore, after the model made a prediction of actual sales, any predictions with $V_m < 0.14$ were rounded to zero. This heuristic effectively minimized the over-estimation error observed.

7.3.4 Model Construction and Testing

The random forest was trained on 70% of the dataset for each forecast type available, internally tested and tuned using oob error estimates. A validation set of 15% of data was used to test variable significance levels and to eliminate variables as needed to improve performance. Although the random forest algorithm is a relatively robust predictor even with highly correlated features [27], feature removal was investigated for possible performance improvements. Additional features were engineered and considered iteratively. When validation performance was unacceptable, new features were included in the feature set, existing features were removed, and the resulting feature set was tested for model performance improvement.

Following successful model validation, the final model was tested on the unseen 15% of data available in the dataset. WMAPE rates were recorded for each test and test results were compared to the baseline demand forecasting model forecast and results were interpreted as described in Chapter 8.

7.3.5 Coverage Measurement

Coverage is a measurement used to determine the objective stock level of an article in a given store as described in Section 3.4. The coverage value required for each (MCCT, store)

is determined by measuring the historical demand for an article compared to the forecasted demand of that article in every store for the same period. This measurement yields the forecast error as:

$$Error_{forecast} = \frac{VmReal}{VmForecast}$$

Thus a forecast error of greater than 1 represents a scenario where actual demand exceeded forecast demand while a forecast error of less than 1 represents a scenario where actual demand was less than the forecast demand. Finally, a forecast error of exactly 1 represents that the forecast and the actual demand were exactly the same.

Taken over a period of time equivalent to the period used for creating the demand forecast, a distribution of forecast errors is produced to describe the demand forecast performance. A service level (p) is selected by the distribution team through discussion with the commercial teams involved. Given p , the error level from the distribution is selected corresponding to the value that covers $p\%$ of the distribution. For a perfectly performing demand forecast, this error level would be 1. A $p\%$ error level greater than 1 indicates that demand forecast is biased below the required service level and therefore the inventory level must be increased. Similarly, if the error level is less than one, this indicates a consistently over-optimistic demand forecast and inventory shipped will be reduced.

By multiplying the identified forecast error by the lead time between forecasting article demand and arrival in stores, a coverage level is determined for each (MCCt, store). The coverage level calculated is applied to the demand forecast to produce the inventory shipment for each (MCCt, store) combination such that $Shipment \propto Coverage * Vm$. Thus, coverage is an appropriate proxy for reduction in inventory attributable to demand forecast accuracy improvements.

In this study, forecast error ratios, Δ coverage, and expected shipments were calculated across all (MCC, store) combinations present in the test dataset based on model and baseline demand forecasts. p was selected to be 90%, representing a reasonable service level for the fashion industry. In this implementation, cases when $Vm_{baseline} = 0$ were handled according to the value of $VmReal$. In cases where $VmBaseline = 0$ and $VmReal = 0$ then the forecast error was set to 1. If, however, $VmBaseline = 0$ and $VmReal > 0$ then the

forecast error was set to ∞ . This implementation was selected to reward correctly predicted 0 cases while penalizing under-estimation. This modification was not necessary with the updated model because 0 sales were never predicted for products with positive e-commerce data.

THIS PAGE INTENTIONALLY LEFT BLANK

Chapter 8

Results Analysis

Modeling results were analyzed for validation and final test runs and summarized for the final models. All Brick-and-mortar results shown use models that were built using both e-commerce and brick-and-mortar feature sets to predict demand for the specified channel of interest. E-commerce demand forecast models only considered e-commerce features and historical sales for forecasting. Analysis was performed considering reduction in WMAPE for modeling and forecasting accuracy measurement. Additionally, a translation of WMAPE improvement to coverage and expected inventory shipments was made to present results with an operational improvement focus. The results and findings of this analysis are summarized in the following sections.

8.1 Prediction Error

Weighted Mean Absolute Percent Error

WMAPE and associated significant features for the random forest model incorporating e-commerce features (Model) and for the current demand forecasting model (Baseline) were measured for each forecast type of interest: initial shipments to Zara.com and to brick-and-mortar stores, and replenishment shipments to Zara.com and to brick-and-mortar stores. The reduction in WMAPE was recorded for each case as the improvement provided by including e-commerce features ($\Delta WMAPE$).

Feature Importance

Significant features of interest were identified using a standard random forest variable importance calculation metric. Specifically, the percent increase in mean square error (MSE) metric was used to rank the feature significance in the regression. While building the random forest, the algorithm records the accuracy of the prediction using the oob samples for each tree grown. Next, each feature selected in the nodes is permuted one at a time in the oob samples presented to the tree. The error measured is then averaged to produce a measure of the MSE increase as a result of randomizing the value of each feature in the ensemble. [28]

It is important to note that the measurement of variable importance in this manner does not accurately yield the effective error introduced if a feature was completely removed from the feature set. When a feature is completely removed, other variables may be used instead to provide the same information, especially if there is significant correlation between them. Therefore, the measure of feature significance for highly correlated features was not measured. Instead, relative importance of features was measured using the minimum set of significant features required to maintain predictive performance.

8.1.1 Initial Products

WMAPE

In the initial shipments, the incorporation of e-commerce features provides a significant advantage to brick-and-mortar store demand forecasting in particular. Results from the modeling efforts for the brick-and-mortar as well as Zara.com channels are presented in Table 8.1. A plot of the error ($predicted - actual$) scaled by the mean error for the model type ($\frac{1}{N} \sum |predicted - actual|$) for the baseline and model for a subset of test entries analyzed are shown in Figure 8-1.

Initial Shipment: WMAPE Performance	
Channel	$\Delta WMAPE$
Brick-and-Mortar	-0.2269
Brick-and-Mortar w/ subs	-0.1985

Table 8.1: Initial shipment demand forecast WMAPE improvements with and without subscriptions

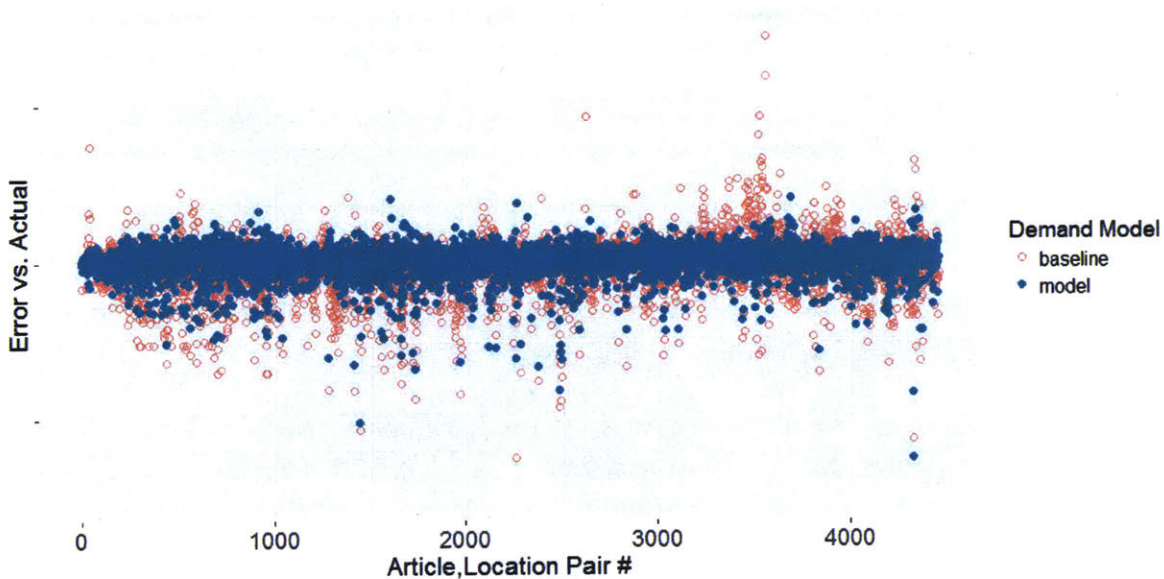


Figure 8-1: Forecasting error for baseline and updated models against actual V_m for a subset of test entries.

Results from initial product offerings at brick-and-mortar stores indicate a potential WMAPE reduction of 0.2269. A significant portion of the high WMAPE level in Zara's current baseline demand forecast can be attributed to the fact that no same store prior sales data exists for initial article shipments. Instead, surrogate similar articles are used as described in Section 3.4. For this reason using the same article and same region information from the e-commerce dataset provides the model with more accurate demand information to forecast with.

Unfortunately, the limited number of data entries (see Table 7.1) in the Zara.com applicable feature set post filtering was insufficient to build a robust random forest model as evidenced in the poor forecast results. This dataset limit results from the restriction that e-commerce data must be available at least one day prior to the forecasting date. In the Zara.com cases, this restricts the dataset to only those cases of articles available for subscription prior to availability for sale, or which were listed for another reason on the website prior to being available for sale.

Similarly, the entries available for the initial products shipped to brick-and-mortar stores were relatively limited as well (see Table 7.1). While the results do indicate significant

potential improvement, additional testing is necessary to confirm the findings with a reliable model.

Feature Importance

Feature importance measurements were run on both the full set of features available to the random forest model as well as a restricted set of uncorrelated features. The results of the feature significance tests are outlined in Table 8.2.

Initial Shipment: B&M Feature Importance		
Rank	Feature Name	% MSE Increase
1	Subfamily	0.2197
2	Vm	0.1786
3	Net Cart	0.1465
4	Avg. Best Position	0.1305
5	Ecommerce Lead	0.1251
6	Avg. Categories	0.1151
7	Rank	0.0946
8	Prior E-commerce Sales MCC	0.0872
9	Price	0.0740
10	Region	0.0637
11	Colors	0.0623
12	Weekday	0.0586
13	Days New	0.0577

Table 8.2: Initial shipment feature importance rankings for brick-and-mortar stores.

The results demonstrate that in the initial shipment case the subfamily of the article, baseline demand forecast prediction, website behavior and structure information are all significant features of the model. It is significant to note that while subfamily classifications provide the model with additional similar article information to use for predicting demand, the Vm feature is very significant to the model accuracy as well.

8.1.2 Replenishment Products

In the replenishment shipments, the incorporation of e-commerce features provides a measurable decrease in forecast error in both the brick-and-mortar and Zara.com sales channels. Results for the brick-and-mortar and Zara.com channels are summarized in Table 8.3. Scaled

forecast error comparison plots for the brick-and-mortar and e-commerce replenishment forecasts are shown in Figure 8-2 and 8-3 respectively with the same y-axis scale for comparison.

Replenishment Shipment: WMAPE Performance	
Channel	$\Delta WMAPE$
Zara.com	-0.0456
Zara.com w/ subs	-0.0719
Brick-and-Mortar	-0.0008
Brick-and-Mortar w/ subs	-0.0396

Table 8.3: Replenishment shipment demand forecast WMAPE improvements with and without subscriptions

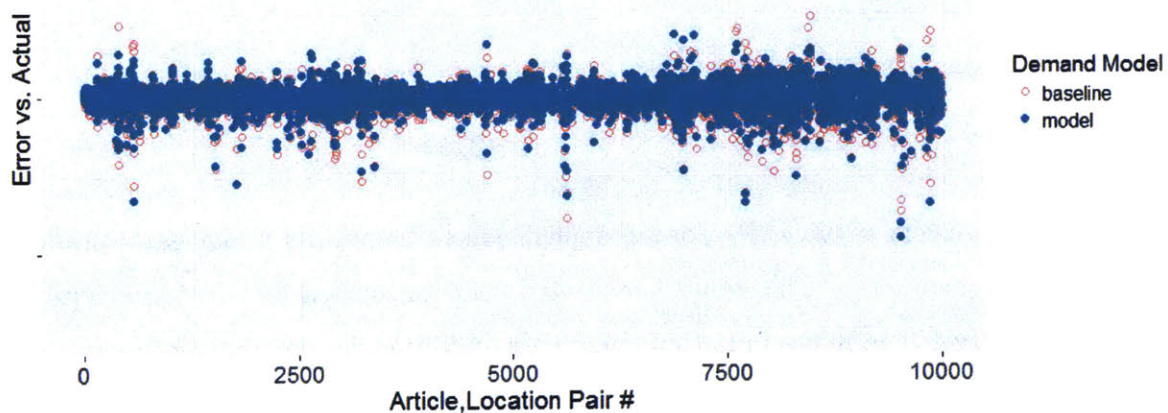


Figure 8-2: Forecasting error for baseline and updated models against actual V_m for a subset of test entries for brick-and-mortar stores.

In both sales channels, the reduction in WMAPE for articles with subscriptions is very promising for replenishment shipments. For the Zara.com sales channel, the WMAPE is reduced by 0.0456 without considering subscriptions and 0.0719 when including subscription information. Similarly, the brick-and-mortar channel observes a WMAPE reduction of 0.0396 with subscriptions however no significant improvement for articles without subscriptions. The lack of improvement is likely due to the low expected sales volume. This would explain the improvement observed for articles with subscriptions, since they have approximately twice the expected average sales as those without due to popularity.

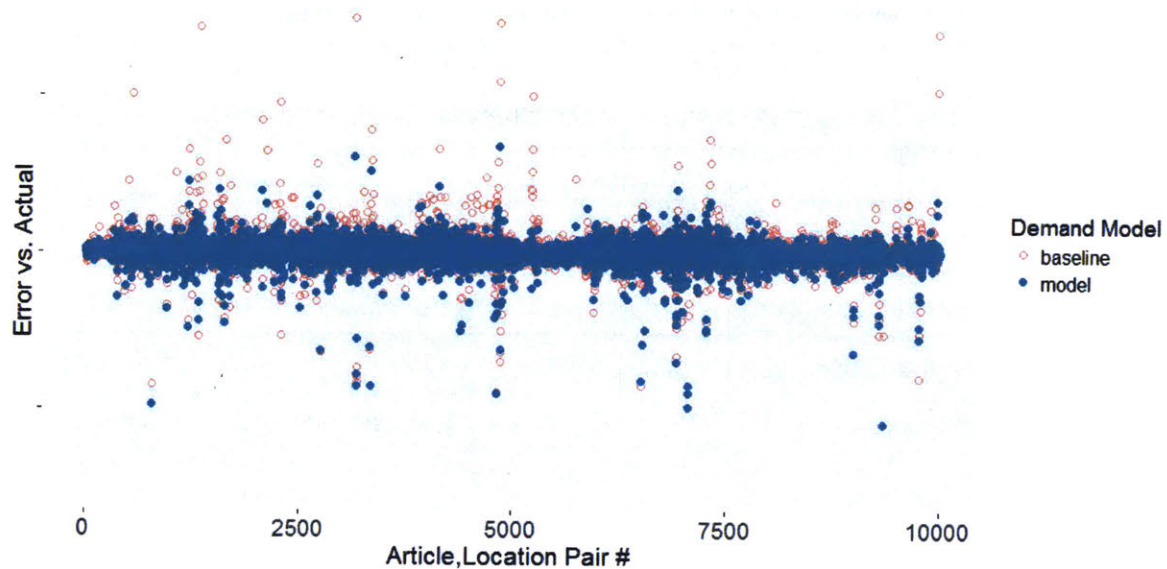


Figure 8-3: Forecasting error for baseline and updated models against actual Vm for a subset of test entries for zara.com.

The improvement in WMAPE for the replenishment shipments is not as drastic as in the initial shipment case. This result is expected because the baseline WMAPE levels are significantly lower than in the initial shipment case. Additionally the higher accuracy of the Zara.com forecast compared to the brick-and-mortar store forecast mirrors the accuracy of the baseline model. The effect is a result of the sales (and correspondingly shipment) volume to each store. In the Zara.com case, the demand forecast represents the pooled demand of the entire market and thus the volumes are much higher. Conversely, with many different retail locations and comparatively low sales and shipment volumes, the relative variance of the brick-and-mortar forecasts is expected to be greater than that of Zara.com.

In both cases, subscription information offers additional information useful for the prediction of future sales, however more strongly in the case of Zara.com. This is likely due to the fact that customers that subscribe on the website will be notified of coming stock with a link to the website in order to complete the sale.

Feature Importance

Feature importance measurements were run on both the full set of features available to the random forest model as well as a restricted set of uncorrelated features. The results of the feature significance tests are outlined in Table 8.4 and 8.5.

Replenishment Shipment: B&M Feature Importance				
E-commerce Data Only			w/ Subscriptions	
Rank	Feature Name	% MSE Increase	Feature Name	% MSE Increase
1	Vm	0.4018	Vm	0.5846
2	Subfamily	0.0544	Subfamily	0.1389
3	Stores	0.0506	Net Cart	0.0951
4	Net Cart	0.0410	Article Age	0.0824
5	Article Age	0.0492	Rank	0.0823
6	Article Age (Ecomm)	0.0479	Article Age (Ecomm)	0.0796
7	Ecomm Sales (MCC)	0.0273	Stores	0.0659
8	Avg. Best Position	0.0243	Ecomm Sales (MCC)	0.0538
9	Region	0.0193	Region	0.0467
10	Rank	0.0180	Avg. Best Position	0.0367
11	Price	0.0155	Total Capacity	0.0312
12	Avg. Categories	0.0101	Avg. Instock (MCC)	0.0311
13	Total Capacity	0.0072	Price	0.0252
14	Avg. Instock (MCC)	0.0056	Subscriptions (MCC)	0.0252
15	Days New	0.0048	Avg. Categories	0.0198
16	Mall	0.0013	Days New	0.0175
17			Mall	0.0043

Table 8.4: Initial shipment feature importance rankings for brick-and-mortar stores.

The results demonstrate that in the replenishment shipment case the baseline demand forecast prediction, subfamily of the article, website behavior and structure information are all significant features. The prior subscriptions to the article are particularly significant in the e-commerce channel since conversion to sale of subscriptions happens primarily on the website. Also of note is the significantly greater importance of the Vm output from the baseline. This indicates that the predictions of the baseline forecast are only slightly modified by the addition of the e-commerce data in the model. The performance of the baseline forecast is therefore an excellent feature to build off of.

Replenishment Shipment: Zara.com Feature Importance				
E-commerce Data Only			w/ Subscriptions	
Rank	Feature Name	% MSE Increase	Feature Name	% MSE Increase
1	Vm	8.775	Vm	19.41
2	Net Cart	2.162	Subs (MCC)	7.208
3	Subfamily	1.077	Net Cart	4.358
4	Article Age Ecomm	0.5210	Subfamily	3.153
5	Price	0.4343	Article Age Ecomm	1.120
6	Avg. Best Position	0.4284	Avg. Best Position	1.010
7	Region	0.3044	Region	1.005
8	Avg. Instock (MCC)	0.2655	Avg. Instock (MCC)	0.8135
9	Avg. Categories	0.2353	Price	0.6248
10	Days New	0.2020	Avg. Categories	0.4581
11			Days New	0.0806

Table 8.5: Initial shipment feature importance rankings for Zara.com forecasts.

8.1.3 OLS Linear Model Performance Comparison

An OLS linear model was also run on each forecast type considered. A comparison of results between the random forest model and the OLS linear model is shown in Table 8.6 for initial shipments and Table 8.7 for replenishment shipments.

Initial Shipment: Δ WMAPE Comparison		
Channel	Linear	Random Forest
Zara.com	-0.0757	-0.1759
Zara.com w/ subs	0.0993	0.0714
Brick-and-Mortar	-0.1231	-0.2269
Brick-and-Mortar w/ subs	-0.2035	-0.1985

Table 8.6: Initial shipment demand forecast comparison using an OLS linear model and a random forest algorithm.

While the linear model performed surprisingly well after tuning for initial products in the brick-and-mortar stores, overall performance of the linear model did not compare well with the random forest. The significant difference in results suggests that an OLS linear model cannot adequately describe the relationship between the e-commerce datasets and the demand for products in the .com and brick-and-mortar sales channels.

Replenishment Shipment: Δ WMAPE Comparison		
Channel	Linear	Random Forest
Zara.com	-0.0025	-0.0467
Zara.com w/ subs	-0.0067	-0.0576
Brick-and-Mortar	0.0416	-0.0008
Brick-and-Mortar w/ subs	0.0287	-0.0396

Table 8.7: Replenishment shipment demand forecast comparison using an OLS linear model and an random forest algorithm.

8.2 Distribution Inventory Effects

For most forecast types, e-commerce data provides significant potential for improving demand forecast accuracy in distribution. While WMAPE improvements are excellent indicators of potential performance benefits, the translation of five points of WMAPE reduction into tangible benefits for the firm are not obvious for company stakeholders. Translating the WMAPE reduction to a familiar inventory metric helps to more clearly communicate the potential advantage provided by the data. Here, the change in coverage required to maintain the desired service level was used to translate demand forecast accuracy effects into inventory reduction.

8.2.1 Initial Products

As described in Section 3.4, proposed shipments are determined by either the calculated shipment based on the demand forecast and coverage level required, or by the minimum shipment quantity to meet commercial display requirements. For the initial shipment case, coverage is not an effective proxy for the reduction in required inventory. Measuring the average coverage level required for the initial shipment dataset confirmed that the baseline prediction was too frequently 0, resulting in a non-meaningful (∞) coverage output for the samples of interest. The anticipated coverage required for the forecast made using the model is summarized in Table 8.8.

Initial: Coverage			
Channel	Baseline	Model	% Δ Coverage
Brick-and-Mortar	∞	2.034	NM
Brick-and-Mortar w/ subs	∞	1.973	NM

Table 8.8: Initial shipment relative coverage measurement and improvement across all stores in the selected market.

Initial: Average Shipment			
Channel	Baseline	Model	% Δ Avg. Shipment
Brick-and-Mortar	∞	10.79	NM
Brick-and-Mortar w/ subs	∞	12.14	NM

Table 8.9: Initial shipment calculated average shipments and improvement across all stores in the selected market.

8.2.2 Replenishment Products

The change in required inventory due to coverage for replenishment shipments is substantial. The results of the coverage change and resulting average shipment sent in each channel are summarized in Tables 8.10 and 8.11 respectively.

Replenishment: Coverage	
Channel	Δ Coverage
Brick-and-Mortar	-34.46%
Brick-and-Mortar w/ subs	-42.64%
Zara.com	-25.49%
Zara.com w/ subs	-18.28%

Table 8.10: Replenishment shipment relative coverage measurement improvement across all stores in the selected market.

The results demonstrate a significant change in both the coverage required to maintain a 90% service level and the expected average articles shipped as a result. The results show that in the best case, the model reduces the average shipment by 34.58% to the brick-and-mortar stores and by 27.57% to the Zara.com FC's. This change in inventory required represents a significant potential savings while maintaining a constant service level. The coverage differences suggest that the baseline forecast is underestimating demand relative to the e-commerce model.

Examining the distribution of forecast errors used to determine the relative coverage multiple for each forecast type lends insight into the effects underlying the performance dif-

Replenishment: Average Shipment	
Channel	Δ Avg. Shipment
Brick-and-Mortar	-26.06%
Brick-and-Mortar w/ subs	-34.58%
Zara.com	-27.57%
Zara.com w/ subs	-23.55%

Table 8.11: Replenishment shipment calculated average shipments improvement across all stores in the selected market.

ferences. The distributions of the forecast error ratios for the brick-and-mortar and Zara.com channels are shown in Figures 8-4 and 8-5 respectively.

In the distribution of the brick-and-mortar store forecast error ratios, there is a significant difference between the baseline and model distributions. The peak in the baseline distribution occurs at 1 due to the frequency at which the baseline predicts 0 sales for an (MCC, store) combination. The peaks that occur near zero represent very small demand forecasts from the baseline and proposed models. Because the proposed model always has information from the e-commerce dataset, it never predicts 0 but comes very close and has a significant peak at very small ratio values. It is also worth noting the thicker tail of the baseline distribution, thus while the baseline model is more accurately predicting the (MCC, store) with 0 sales while significantly underestimating many others, the proposed model appears to smooth the error out over the distribution and fall-off more rapidly.

The distributions of the Zara.com forecast error ratios demonstrate the expected distribution patterns. The peak of both distributions is close to 1, with the model being more concentrated and the baseline demonstrating a slight overestimation of demand. On the other end of the distribution, the thick tail of the baseline distribution shows significant under-estimation of demand in many cases while the model tail falls off relatively quickly. This difference in the tails is the likely explanation for the significant difference in coverage required in the two cases.

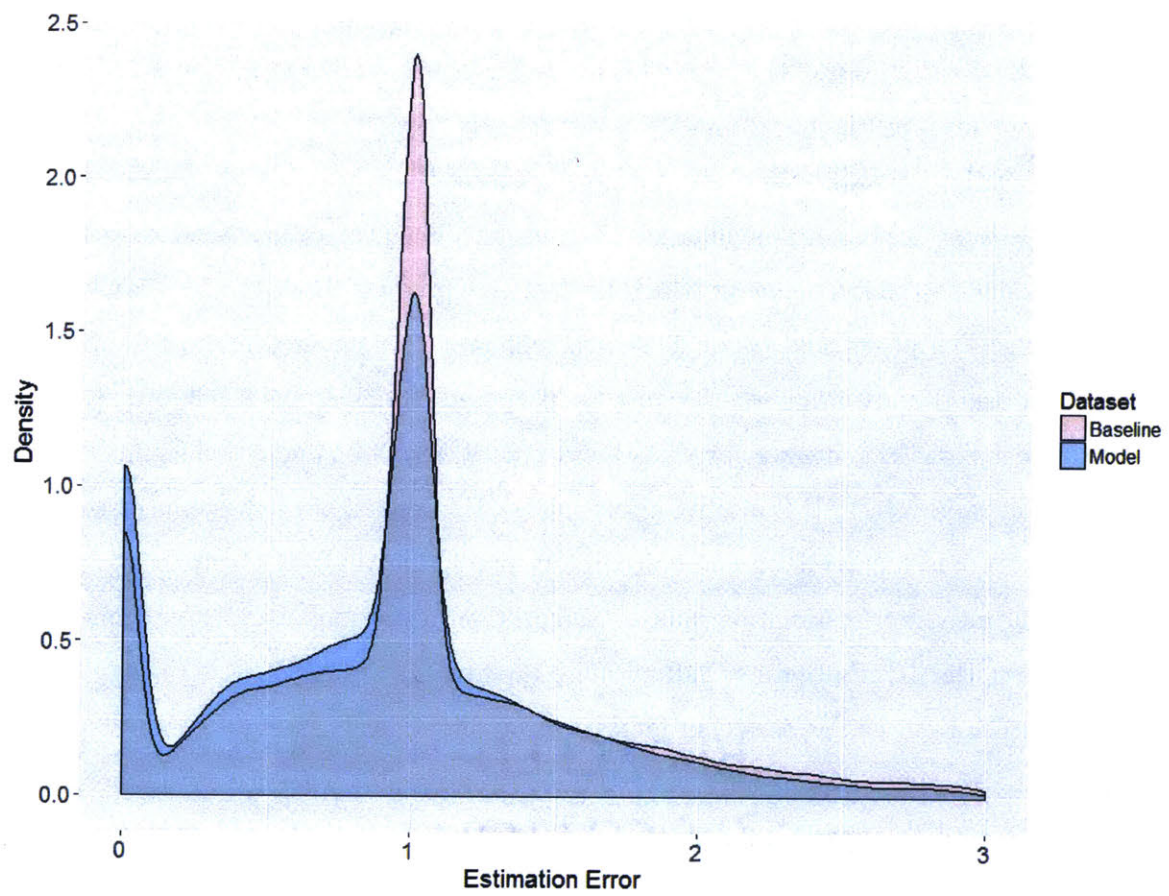


Figure 8-4: The distribution of forecast error ratios used to calculate the coverage required for brick-and-mortar stores.

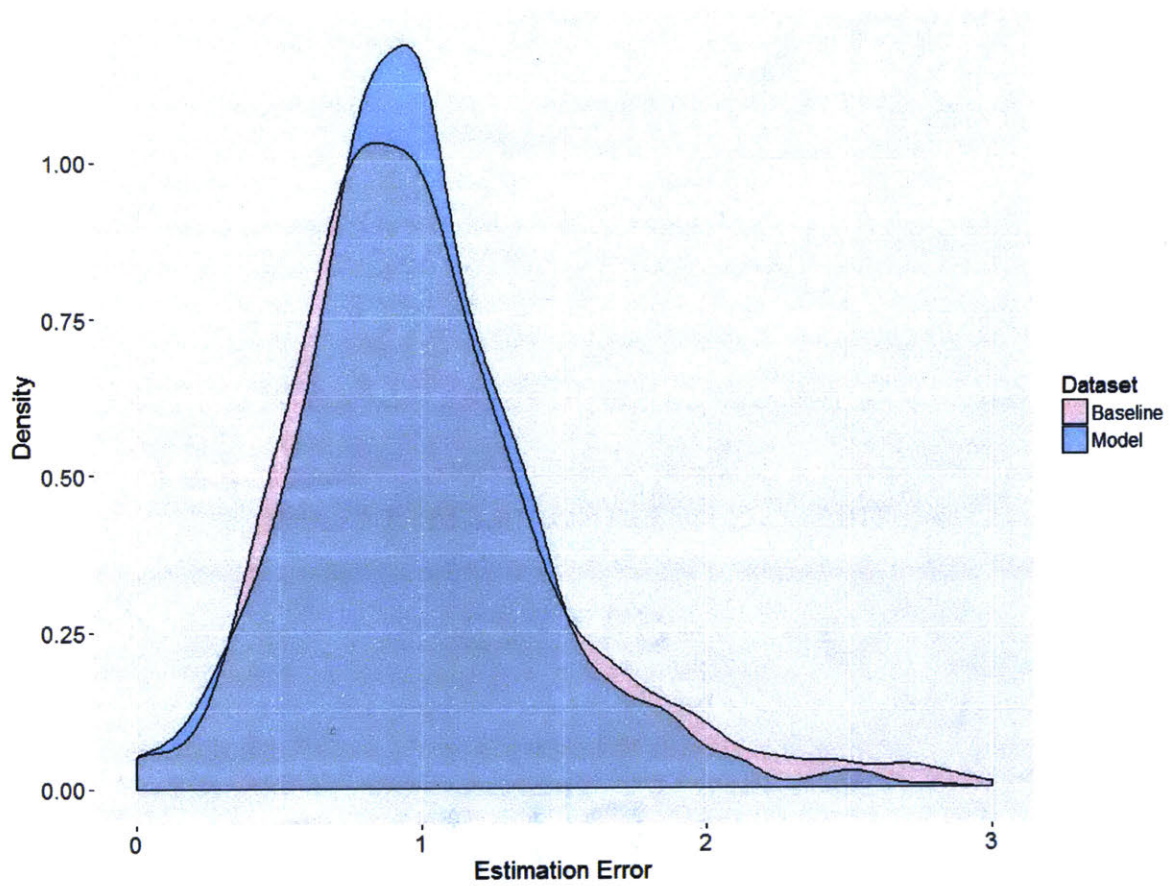


Figure 8-5: The distribution of forecast error ratios used to calculate the coverage required for Zara.com fulfillment centers.

THIS PAGE INTENTIONALLY LEFT BLANK

Chapter 9

Conclusions and Recommendations

9.1 Summary of Results

The results of this study show that incorporating e-commerce data collected by Zara.com into the distribution demand forecasting model can significantly improve forecast accuracy.

In the initial product forecasts, the modified model provides a significant increase in demand forecast accuracy as demonstrated in the WMAPE improvement results for brick-and-mortar stores. While dataset size restricted the modeling efforts for the Zara.com channel, the results from brick-and-mortar stores indicate that a similar result is possible for e-commerce as well. The ranking of significant features in the model indicates that some of the most significant indicators for demand forecasting are the subfamily of the article and the net cart activity or another correlated feature of website traffic modified by the positioning on the website and the length of time the article has been available to view. This data is available for a limited selection of items not yet available for sale on Zara.com and this study suggests it can be leveraged to improve forecast accuracy for initial shipments to e-commerce FCs as well.

The performance of the modified algorithm in replenishment shipment forecasts also demonstrates the utility of e-commerce datasets for predicting the demand of both the Zara.com and brick-and-mortar sales channels. The modest reductions of WMAPE in both cases represent significant inventory effects for the replenishment shipments to the stores and FCs. The significant variables analysis reveals that the baseline forecast is already a strong

predictor of future demand. The most influential e-commerce features used to adjust the prediction of the baseline include customer behavior, website structure and article features.

9.2 Recommendations

The results presented indicate significant predictive value in the datasets currently collected by Zara.com. In response, several next steps will enable Zara to realize immediate value from the work while continuing to pursue additional insight from the datasets in future investigations.

9.2.1 Implementation

E-commerce datasets have been demonstrated to provide a significant predictive value for distribution. In response, systematic and continuous collection of the significant features identified in this study is recommended. The datasets used in this study are all currently collected in a raw form by the Zara.com team. The raw datasets exist in separate databases and must currently be transformed and aggregated into meaningful features for model prediction. The most significant features identified in this study should be collected, aggregated, and migrated to a central shared database for use in both e-commerce operations and planning as well as distribution algorithm processing. These variables should include at minimum: the net cart activity or another correlated feature for high intent customer traffic, website structure and article position including the total number of days an article has been available on the website.

To extend the applicability of the prediction model to the e-commerce channel initial product offerings, research into operational process modifications in Zara.com can be conducted. While the analysis in this study demonstrates significant predictive value in the data collected online prior to sales of an article, this value must be balanced with the cannibalistic effect on the sales of currently available items. An additional study into the cannibalization effects of offering a preview of a product on the website prior to

9.2.2 Future Investigations

The work described represents a first step at integration of datasets across channels of a multi-channel retailer. While the results demonstrate significant utility in demand forecasting accuracy, additional work must be done to facilitate increased implementation efficiency and to better understand the dynamics between the available sales channels.

Additional Datasets

A significant drawback to the use of a random forest algorithm for evaluation in this study is the requirement for relatively large training datasets for each forecast type. This limitation was seen most clearly in the inability to predict initial shipments for the e-commerce channel due to lack of sufficient cases of website data prior to sales availability. Continued study of the applicability of the model should be pursued with additional data collected in subsequent seasons. Continued data collection will offer significantly more data for analysis than was available for this study in some cases. Models should be re-run to confirm the findings described here and to better understand the forecast improvement contributors.

Operational process modifications present an additional option for increased dataset availability. A significant limiting factor for the e-commerce initial shipment forecasts was the lack of products available for subscription prior to availability for sale on the website. Offering products on the website between 7 and 14 days prior to sale would allow customer preferences to be revealed through the customer behavior even without sales. Current e-commerce commercial processes allow this only for a small subset of items. Not only will this provide the capability to predict sales for e-commerce initial shipments, but it will provide significantly more data for the other distribution forecasts as well. It is important to note that this operational change must be considered in balance with the cannibalistic effects of revealing a forthcoming product on current product sales. If the lost sales to existing products outweighs the benefits gained in demand forecast accuracy then the process modification should not be pursued.

Additional data may also be collected at finer levels of aggregation. Website behavior data in particular was limited to regional aggregation while customers across different sub-regions

or cities in a given region may exhibit significantly different article preferences and demand profiles. Aggregation of website traffic data at a finer resolution should be investigated for significant forecasting accuracy effects. If sufficient traffic location data exists at this resolution, the algorithm might be able to better tailor the model prediction to specific stores rather than across entire regions as presented here.

The e-commerce features investigated in this study represent the most significant features currently collected by Zara.com for distribution optimization. Additional datasets can be collected and analyzed in future investigations to provide greater insight into the dynamic and highly variable customer preferences in the fashion industry. For example, the influence of social media, news, and fashion blogs on customer article preferences was postulated in this study without data to test the hypothesis. A future study incorporating this information into meaningful features may reveal a significant improvement and enable even greater future demand accuracy than the features currently used.

Algorithm Simplification

While the random forest model used in this study is a strong predictor, the training requirements and limited interpretability of the model are not ideal for implementation in distribution operations and widespread organizational acceptance. Although the OLS linear model did not perform as well as the random forest, the comparable performance for some forecasts indicate that simplified algorithms could be designed to better suit the forecasting problem. Simplified algorithms exist that could provide more interpretable and more efficient training and prediction times for practical implementation. Through further algorithm optimization, a model based on a regression tree or clustered linear regression may be able to approach or exceed the performance of the random forest used here with much faster runtimes. A particular benefit of these algorithms is the interpretability of the underlying models and ability to communicate the algorithm to the commercial teams and management for organizational buy-in and feedback.

Pilot Project

The actual operational impact of the WMAPE improvement is only partially demonstrated by the coverage and average inventory shipment changes measured here. For example, gains due to captured lost sales from stores that received no inventory (and therefore had no actual demand information) are not represented here and may represent a significant operational improvement. Therefore, a pilot test is recommended prior to full scale implementation of the modified algorithm. This pilot test will allow the distribution team to observe whether the reduction in required coverage suggested by this study still results in the anticipated service level for customers.

9.3 Conclusion

The work of the MIT 12 project demonstrates a first investigation of the utility of multi-channel data analytics for distribution decision making at Zara. By integrating e-commerce datasets collected by Zara.com with the current baseline distribution demand forecasting model, the specific use case of cross-channel demand forecasting was investigated. The results show promising improvements in demand forecast accuracy for both initial and replenishment shipments and suggest that implementation of a modified forecasting model will reduce inventory coverage requirements significantly while maintaining customer service levels. The implementation of an updated model and the corresponding processes encourages continued collaboration between the brick-and-mortar and e-commerce organizations within Zara. The work presented represents the latest effort in a partnership between MIT and Zara in operations research. Demand forecasting is only a the first of many potential applications of cross-channel datasets for operational optimization. The MIT 12 project lays the groundwork for continued data analysis in the continued partnership between MIT and Zara in the future.

THIS PAGE INTENTIONALLY LEFT BLANK

Bibliography

- [1] Inditex. Annual Report 2015. page 318, 2015.
- [2] Juan R Correa. Optimization of a Fast-Response Distribution Network. (2000), 2007.
- [3] Andres Garro. New product demand forecasting and distribution optimization: a case study at Zara. page 194, 2011.
- [4] Jose M Garcia. Demand Forecasting at Zara: A Look at Seasonality, Product Lifecycle and Cannibalization. 2014.
- [5] Evelyne L Kong, Ecole Centrale Paris, Georgia Perakis, Thesis Supervisor, Bruce Cameron, Thesis Supervisor, and Maura Herson. Cannibalization Effects of Products in Zara ' s Stores and Demand Forecasting by. 2015.
- [6] T. M. Choi, C. L. Hui, and Y. Yu. Intelligent time series fast forecasting for fashion sales: A research agenda. In *2011 International Conference on Machine Learning and Cybernetics*, volume 3, pages 1010–1014, July 2011.
- [7] Na Liu, Shuyun Ren, Tsan Ming Choi, Chi Leung Hui, and Sau Fun Ng. Sales forecasting for fashion retailing service industry: A review. *Mathematical Problems in Engineering*, 2013, 2013.
- [8] Maria Elena Nenni, Luca Giustiniano, and Luca Pirolo. Demand forecasting in the fashion industry: A review. *International Journal of Engineering Business Management*, 5(SPL.ISSUE), 2013.
- [9] Sébastien Thomassey. Sales Forecasting ni Apparel and Fashion Industry: A Review. In *Intelligent Fashion Forecasting Systems: Models and Applications*, chapter 2, pages 9–27. 2014.
- [10] Shuyun Ren, Hau-Ling Chan, and Pratibha Ram. A Comparative Study on Fashion Demand Forecasting Models with Multiple Sources of Uncertainty. *Annals of Operations Research*, pages 1–21, 2016.
- [11] Felipe Caro and Jeremie Gallien. Inventory Management of a Fast-Fashion Retail Network. *Operations Research*, 58(2):257–273, 2010.
- [12] Kin Fan Au, Tsan Ming Choi, and Yong Yu. Fashion retail forecasting by evolutionary neural networks. *International Journal of Production Economics*, 114(2):615–630, 2008.

- [13] Zhan-Li Sun, Tsan-Ming Choi, Kin-Fan Au, and Yong Yu. Sales forecasting using extreme learning machine with applications in fashion retailing. *Decision Support Systems*, 46(1):411–419, 2008.
- [14] Yong Yu, Tsan Ming Choi, and Chi Leung Hui. An intelligent fast sales forecasting model for fashion products. *Expert Systems with Applications*, 38(6):7373–7379, 2011.
- [15] María del Mar Roldán García, José García-Nieto, and José F. Aldana-Montes. An ontology-based data integration approach for web analytics in e-commerce. *Expert Systems with Applications*, 63:20–34, 2016.
- [16] Kris Johnson, Bin Hong, and David Simchi-levi. Analytics for an Online Retailer : Demand Forecasting and Price Optimization. *Manufacturing & Service Operations Management*, (2012):1–33, 2013.
- [17] Yiwei Zhou and Xiaokun (Cara) Wang. Explore the relationship between online shopping and shopping trips: An analysis with the 2009 NHTS data. *Transportation Research Part A: Policy and Practice*, 70(June):1–9, 2014.
- [18] Erik Brynjolfsson, Yu Jeffrey Hu, and Mohammad S Rahman. Competing in the Age of Omnichannel Retailing. *MIT Sloan Management Review*, 54(4):23–29, 2013.
- [19] Max Kuhn and Kjell Johnson. *Applied Predictive Modeling*. Springer, 2013.
- [20] R Core Team. *R: A Language and Environment for Statistical Computing*. R Foundation for Statistical Computing, Vienna, Austria, 2016.
- [21] Andy Liaw and Matthew Wiener. Classification and regression by randomforest. *R News*, 2(3):18–22, 2002.
- [22] Torsten Hothorn, Kurt Hornik, and Achim Zeileis. Unbiased Recursive Partitioning: A Conditional Inference Framework. *Journal of Computational and Graphical Statistics*, 15(3):651—674, 2006.
- [23] Max Kuhn. Contributions from Jed Wing, Steve Weston, Andre Williams, Chris Keefer, Allan Engelhardt, Tony Cooper, Zachary Mayer, Brenton Kenkel, the R Core Team, Michael Benesty, Reynald Lescarbeau, Andrew Ziem, Luca Scrucca, Yuan Tang, Can Candan, and Tyler Hunt. *caret: Classification and Regression Training*, 2016. R package version 6.0-73.
- [24] Hadley Wickham and Romain Francois. *dplyr: A Grammar of Data Manipulation*, 2016. R package version 0.5.0.
- [25] Hadley Wickham. *tidyr: Easily Tidy Data with 'spread()' and 'gather()' Functions*, 2017. R package version 0.6.1.
- [26] Leo Breiman. Randomforest2001. pages 1–33, 2001.
- [27] Leo Breiman. Statistical Modeling: The Two Cultures. *Statistical Science*, 16(3):199–215, 2001.

- [28] Trevor Hastie, Robert Tibshirani, and Jerome Friedman. The Elements of Statistical Learning. *Elements*, 1:337–387, 2009.