

RETAIL IN DETAIL

ECE/CS 498 DSG Final Project

Spring 2020

Badrinarayanan R
br17
Industrial Engineering

Anirudh Sharma
ashar29
Information Management

Anunay Sharma
anunays2
Industrial Engineering

Abstract—Retail industry is a vast and major sector of the economy comprised of companies selling finished products to the customers. Around 66% of the U.S. gross domestic product (GDP) comes from retail consumption. The industry initially started with brick and mortar store retailers who are engaged in the sale of products from physical locations to make the customers purchase onsite. The last 7 years have seen the inception of e-tailers where products are promoted online and arrives at the doorstep of customers with quick turnaround time. Online purchases take place at the tap of a button on the mobile and make customers' life easier than ever. Although the internet era has made the retail segment more convenient, the retailers face day-to-day challenges across multiple domains which needs more than human predictions to be successfully managed, and provide an efficient and sustainable solution considering factors like time, money, and volume. The presence of statistics and machine learning models has been widespread and has made its entry very well in the field retail too. The retail data is an invaluable resource in the current period and is fed as an input to these analytical models to derive actionable insights for the retail stores. The scope of this project is concerned with the primary challenges in retail: customer segmentation, inventory, and sales planning, and promotion strategies. Each of these problems is broken down to understand the existing scenarios and are finally treated with a solution that has a solid foundation from the statistical standpoint

I. INTRODUCTION

Retailing is an integral part of modern society. Consumers highly depend on retail stores (both online and offline) for their various goods and service requirements. Earlier, goods and services were made available through the process of trading. But in present times trading is replaced by buying and selling goods which makes retail stores an important part of the supply chain and thus remains an attractive business sector for many. With a presence worth several trillion dollars, retail remains a significant contributor to the global economy. [1] With the recent technological advancements, analytics in retail has gained much popularity for its aid in helping the business to make key decisions. It enables the retailer to come up with standard methodologies that dissect the customer segment and product categories and can help in boosting its revenue. This project puts various concepts of retail analytics to use providing analytical insights that can be essential for making marketing, and procurement decisions. We aim to

leverage several machine learning techniques to analyze the retail business for a single retailer.

After careful understanding of the data, we have decided to answer 3 main questions that could prove to be beneficial for the retail store owner of this data set:

1) Customer Segmentation: Grouping customers into several clusters based on their purchase attributes and other relevant features thereby enabling the client to target a precise set of audiences for relevant promotions and offers (targeted marketing).

2) Sales forecasting: Analyzing the trends in the data to predict the sales of products to be sold in the upcoming weeks for the retail company

3) Affinity/Purchase/Basket Analysis: Exploring the customer purchase behavior on a day-to-day basis and to understand product affinity (i.e. Product X sells along with Y/Z)

There are precedent examples of various companies using these techniques successfully to enhance profits.

Infiniti Research, a leading market intelligence solutions provider located in London recently announced the completion of its latest success story on market segmentation analysis. A leading player in the money transfer market wanted to conduct a global agent satisfaction survey in a set of European Countries. With the help of a customer intelligence solution, the client was able to improve its agent management program by gaining an understanding of what the agents value. They were able to ensure a better satisfaction level and hence improve on sales and customer satisfaction. Sales forecasting is used by many e-commerce giants such as Amazon offers an automated solution called Amazon forecast which uses time-series data to predict future sales and product demand.

There will be an effort to cover all the problems proposed herein as much depth as possible and to provide insights that can be beneficial to the retail store owners.

II. EXPLORATORY ANALYSIS AND CHALLENGES

The dataset for this project has its source from Kaggle website as a repository. There has been no problem statement created nor any analysis done on this data to generate any insights. All the approaches and problems which will be discussed in this paper are self-defined and were not inspired

by any source. There are a total of 3 different data sources namely the customer data, product mapping data, and the transaction data for each customer who had visited the store. The transaction data records a total of 4 years of data for analysis and modeling.

As a first step, we perform exploratory data analysis to find patterns in the data. Fig. 1 shows the sales performance of different store types across different months since 2011. The e-shop dominates with twice the sales as compared to others. This is quite reasonable as this was around the time when customers were shifting to e-commerce platforms like Amazon and eBay. An interesting observation to keep in mind from the graph is the sudden decline in the sales data from 2014. This phenomenon is not temporary as it continues until the end of 2014.



Fig. 1. Sales trend across months for different store types

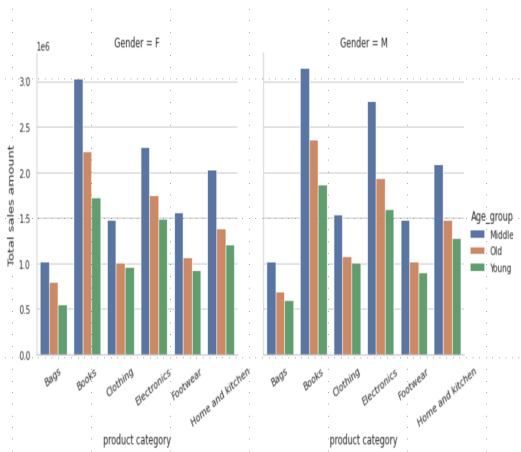


Fig. 2. Total sales amount for different product categories filtered on gender and age group

Fig. 2 indicates that middle-aged customers are the major contributors to the sales across all the product categories which is expected as people generally become capable of spending high amounts around this age. Another observation from this visualization is that books and bags segment are respectively

the highest and lowest sales categories for both the genders (Male and Female). From figure3 it can be seen that teleshops and e-shops contribute to more than 60% of the transactions in the retail store performance

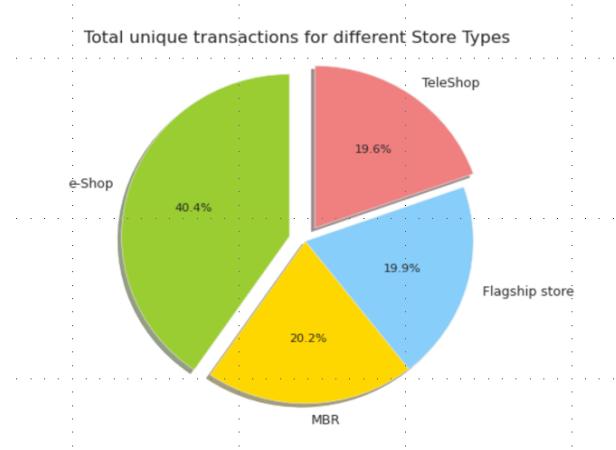


Fig. 3. Distribution of different store types

A. Roadblocks

One of the major roadblock in any data science problem is the data discrepancy. Some of them are mentioned as follows:

- As we can see, the data of 2014 seems to be very suspicious as there is 144% decrease in the overall sales.
- There are transactions where the quantity happens to be negative. These can be considered as returns by the customer but still there were no transactions found where the customer had purchased it well before.
- There is no information on the multiple purchases made by the customer in a single transaction. This would limit providing a product/SKU level recommendation

B. Measures

In order to make the profiling and modelling process unbiased, we have taken certain measures to transform the data and make it ready for further work:

- Restrict the data from 2011-2013
- Remove all the transactions having negative quantities
- The most significant metrics like sales, #transactions and quantity go through a process called as winsorization. It is a transformation done by limiting the extreme values of a data to reduce the effect of outliers

III. PROBLEM 1 - CUSTOMER SEGMENTATION

Customer segmentation will allow the retail stores to learn a great deal about their customers so that they can cater to their needs more efficiently. This will also allow tailoring their communication depending on the customer's life cycle and to prepare a better acquisition strategy. The method used for this problem is K-means clustering. The first step for this method was feature selection. Since there were only a few numerical columns in the data set, an iterative approach was implemented to select the germane variables for clustering.

A. Approach

The k-means algorithm works as follows:

- Step 1: Randomly choose k data points (seeds) to be the initial centroids i.e., cluster centers
- Step 2: Assign each data point to the closest centroid
- Step 3: Re-compute the centroids using the current cluster members
- Step 4: If a convergence criterion is not met, goes to step 2 and the algorithm continues until convergence is met.

Stopping/Convergence criterion:

- No (or minimal) re-assignments of data points to different clusters
- No (or minimal) change of centroids
- Minimal decrease in the sum of squared error(SSE)

$$SSE = \sum_{j=1}^k \sum_{x \in C_j} dist(x, m_j)^2 \quad (1)$$

$$m_j = 1/n_j \sum_{x \in C_j} x \quad (2)$$

where,

C_j is the j^{th} cluster,

m_j is the centroid of cluster C_j ,

n_j is the number of points in cluster C_j ,

$dist(x, m_j)$ is the distance between data point x and centroid m_j (generally Euclidean)

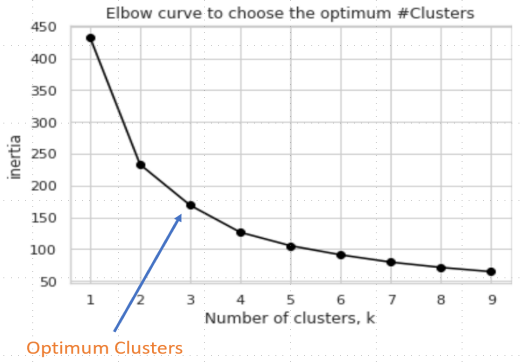


Fig. 4. Elbow curve to chose optimum number of clusters

The variables considered are Sales per Transaction (tells about the basket value), Age (Customer's age), and the number of transactions (Frequency of purchase), quantity, and quantity/transaction (basket size). The K means algorithm was run with different combinations of these variables and finally, we could see 3 well-distinguished clusters for the following set of variables : **Sales/transaction, Quantity, and the number of transactions**. The optimum number of clusters was chosen based on the statistical test of the elbow curve, which shows the within the sum of square distances for different clusters in fig. 4.

The elbow method runs k-means clustering on the data set for a range of values for k (from 1-10) and then for each value of k computes an average score for all clusters. The distortion score is computed as the sum of square distances from each point to its assigned center. As the drop in the inertia is not that significant after the 3rd cluster, the optimum number of clusters selected was 3.

B. Results

The fig. 5 and fig. 6 shows the cluster in a 3D space. While looking at it, although there is some overlap between the clusters, a decision boundary can be seen. After analyzing the data points in these clusters carefully it was inferred that the three clusters would be the bargain hunters/seasonal purchasers(BH)(lesser basket value and transactions), the High Spenders (HS) (with higher basket value) and the regular shoppers(RS) (more frequent purchases). This makes sense from a business aspect which can be a great insight for the retail store when running any promotions or offers.

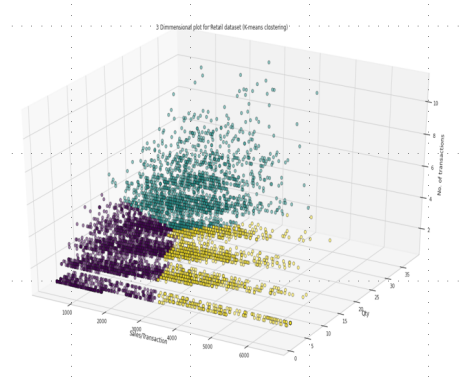


Fig. 5. 3-dimensional scatter plot for K-means

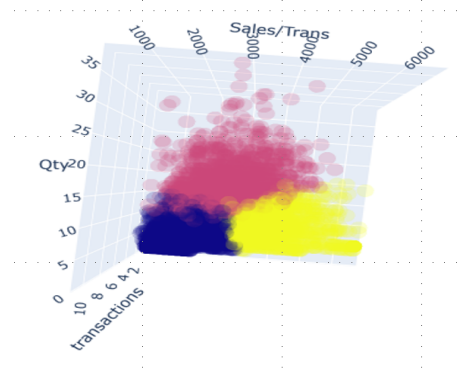


Fig. 6. Top view of the scatter plot for K-means

IV. PROBLEM 2 - SALES FORECASTING

What if the retail shop wants to be well prepared to handle a large crowd of customers during an occasion? It needs to have a very well defined strategy in setting up logistics and inventory planning. Prior to this, every company engages in

a process of estimating its future sales which can be very beneficial in telling the company how to efficiently manage its resources and inventory very well. According to research, companies with accurate sales are 10% more likely to grow their revenue year over year.

Fig. 7 shows the weekly sales distribution of the retail store. This data exhibits different aspects of time series. The red boxes marked in the graph denote a pattern known as cycle which is referred in the later sections



Fig. 7. Weekly sales trend

A. Approach

Linear Regression

There are different ways to approach the problem of forecasting. One of the classic statistical way of predicting a real valued output is linear regression. It tries to model the relationship between a dependent variable and several independent variables to find the line that best fits the data with least squared error. The program for the linear regression is as follows:

$$\min \sum_{i=1}^n (y^{(i)} - w^T x^{(i)} - b)^2 \quad (3)$$

where,

$y^{(i)}$ is the dependent variable

$x^{(i)}$ is a matrix of independent variables

w and b are the coefficient and the intercept for the equation

There are four principal assumptions which justify the use of linear regression models for purposes of prediction:

- **Linearity:** The relationship between dependent and independent variables should be linear
- **homoscedasticity:** The variance of residuals should be same for any value of independent variable
- **No Auto correlation:** The error terms shouldn't be correlated to one another
- **Normality:** The errors should be normally distributed
- **Independence:** The variables are independent of each other

Time series Analysis

Time series analysis refers to the collection of data points analyzed at constant intervals to determine the behavior or pattern of the variables. This analysis helps us to move further towards the goal of forecasting and predictions. The unique thing that distinguishes linear regression from time series is the time dependency. The observations are dependent on one another.

An important criterion that needs to be satisfied to use the time series formulation would be stationarity. The assumption

states that the mean and variance of a time series are constant over time. Once they are ensured to be constant, it would be easier to solve the modeling problem. Most time-series data usually have at least one of these kinds of patterns: trend, seasonality, or cycles. [2]

Trend - The trend describes the general behavior of a time series. If a time series has a positive long term slope over time, it has an upward trend and if there is a negative slope, it has a downward trend.

Seasonality - A seasonal pattern is any kind of fluctuation in a time series that is caused by calendar related events. These events can be the time of year or the time of the day or the week. Seasonality always has fixed frequencies. The seasonal patterns start and end in the same period of a week. Consider the example of the occasion Black Friday and Cyber Monday. The sales at this period are meant to go up and this will be observed distinctly in the data.

Cycle - Cycles are defined as the rises and fall with non-fixed magnitudes that can last more than a calendar year. They are not repetitive. Usually, they result from external factors that make them much harder to predict. The time series forecasting leverages these patterns to produce reliable predictions. The fig. 8 below shows the split of the weekly sales data into different components

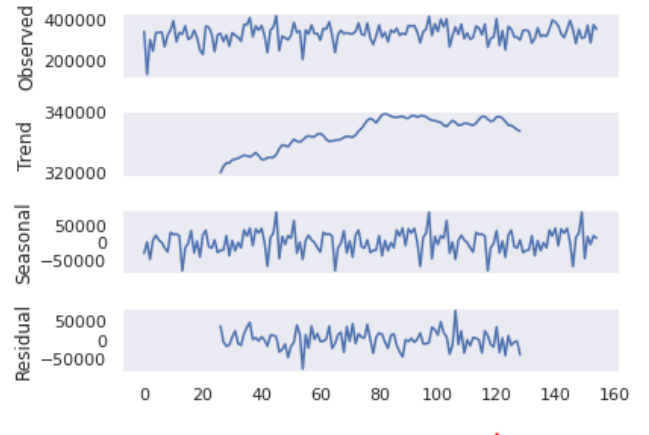


Fig. 8. Components of time series

B. Methodology

In this project, we started the modeling process with both the algorithms mentioned above. A special class of time series implemented is the ARIMA (Auto-Regressive Integrated Moving Average). These models take into account the lag terms of the variables along with the forecast errors as the only predictors for the forecasting.

our main objective is to forecast the sales value of the retail store. We would want to estimate the sales at a weekly level as the lower the granularity, the more accurate are the results. There are a total of 155 weeks of data from 2011-2013. This data follows a series of steps to achieving the final predictions explained in Fig 9:

- The dataset is split into a subset called as training and testing set. The training set is the set on which the whole algorithm trains by capturing all the variations in the past. The test set is the validation set here where we would get to understand how good our model is. As a rule of thumb, training and testing sets are split in the fashion of 80:20.
- Now this data is being tested with the assumptions of linear and time series analysis. Here, the assumptions of linear regression would not be met due to the presence of auto-correlation but still we force-fit the model to understand the baseline score through the most classic statistical procedure
- The ARIMA model has several parameters that go into the model:
AR term (p) tells about the number of lagged observations to be considered in the model.
MA (q) gives the number of lagged forecast errors to be taken into account.
Differencing term (d) is the number of differences between the observations to be taken to keep the time series stationary. All these parameters are derived from Autocorrelation(ACF) and Partial Auto-correlation (PACF) plots shown in fig. 10
- From the plots, we could see that the values of p and q should be ~ 3 or 4 and is found out by iterating with different combinations
- Various accuracy metrics are captured to understand the performance of both models.

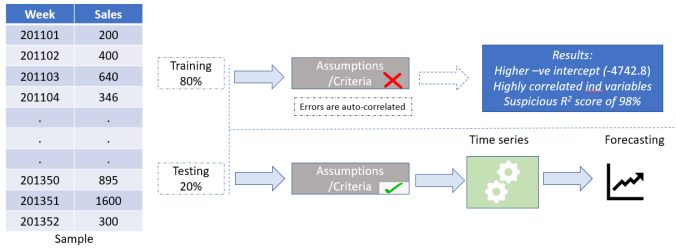


Fig. 9. Flowchart for Time series modelling

C. Forecasting Results

The linear regression outputs a linear equation with a very high -ve intercept of -4742.8. This value is quite inconsistent (as sales can't be negative) as this essentially signifies that without any of the lag variables effect, the model would predict this value as the sales for the future. Moreover, the R^2 yields a score of $\sim 98\%$ which looks to be very suspicious. Additionally, it is to be kept in mind that the data violates the assumption of auto-correlation in the case of linear regression.

When we look at the results of ARIMA model in fig. 11, we could see some interesting observations. The p-values of most of the variables used in the model are less than 0.05 and comes to be very significant. The AIC score (3047) tells about the amount of information lost by the model this score is the lowest among all the iterations picked by the model. Finally, the best model achieves a MAPE of 9.76%

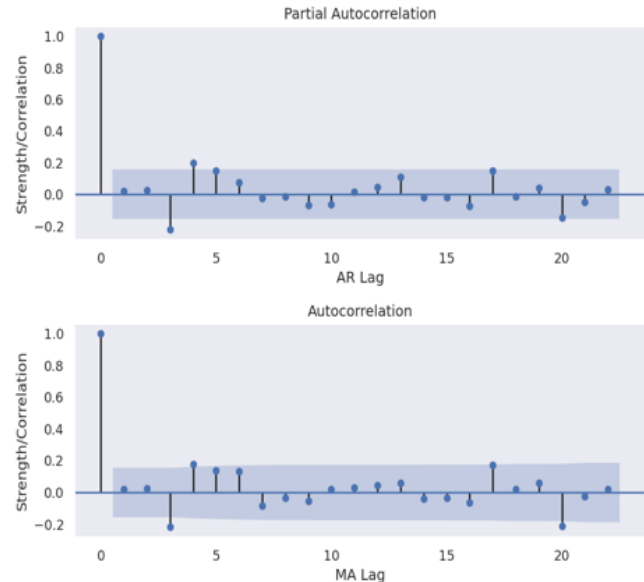


Fig. 10. PACF and ACF plots

ARIMA Model Results							
Dep. Variable:	total_amt	No. Observations:	126				
Model:	ARMA(3, 4)	Log Likelihood	-1514.835				
Method:	css-mle	S.D. of innovations	39394.206				
Date:	Thu, 07 May 2020	AIC	3047.669				
Time:	17:50:45	BIC	3073.196				
Sample:	0	HQIC	3058.040				
	coef	std err	z	P> z	[0.025	0.975]	
const	3.285e+05	3899.993	84.228	0.000	3.21e+05	3.36e+05	
ar.L1.total_amt	0.0595	0.105	0.566	0.572	-0.146	0.265	
ar.L2.total_amt	-0.2028	0.097	-2.092	0.039	-0.393	-0.013	
ar.L3.total_amt	-0.8597	0.106	-8.141	0.000	-1.067	-0.653	
ma.L1.total_amt	-0.0452	0.133	-0.339	0.735	-0.306	0.216	
ma.L2.total_amt	0.3551	0.130	2.728	0.007	0.100	0.610	
ma.L3.total_amt	0.6978	0.136	5.133	0.000	0.431	0.964	
ma.L4.total_amt	0.2199	0.109	2.017	0.046	0.006	0.434	
Roots							
	Real	Imaginary	Modulus	Frequency			
AR.1	0.4621	-0.8884j	1.0014	-0.1737			
AR.2	0.4621	+0.8884j	1.0014	-0.1737			
AR.3	-1.1600	-0.0000j	1.1600	-0.5000			
MA.1	0.4765	-0.8793j	1.0001	-0.1710			
MA.2	0.4765	+0.8793j	1.0001	-0.1710			
MA.3	-2.0631	-0.5385j	2.1322	-0.4594			
MA.4	-2.0631	+0.5385j	2.1322	-0.4594			

Fig. 11. ARIMA - Summary

D. Seasoned ARIMA (SARIMA) - Results

Although the ARIMA model performs well with the obtained parameters, there is always room to make them better by considering even more intricate factors into account. SARIMA model [3] takes into the seasonal terms to provide better results from the existing model. The seasonal components are the counterparts to ARIMA parameters p , d and q are denoted by P , D and Q . A series of iterations using `pdm arima` was done to find the out the best seasonal parameters (fig. 12). The outputs are shown in the figure below

The results of SARIMA suggests that the model was able to capture more variation than the previous models (fig. 13). The seasonal component plays a major role in the retail sales as there are several offers available during different seasons and occasions. The MAPE of the model gives 10.7%, slightly higher than ARIMA but the AIC (1817) and BIC(1828) of the


```

Performing stepwise search to minimize aic
Fit ARIMA(1,1,1)x(0,1,1,52) [intercept=True]; AIC=nan, BIC=nan, Time=nan seconds
Fit ARIMA(0,1,0)x(0,1,0,52) [intercept=True]; AIC=1874.011, BIC=1878.592, Time=0.366 seconds
Fit ARIMA(1,1,0)x(1,1,0,52) [intercept=True]; AIC=1824.251, BIC=1833.413, Time=3.562 seconds
Fit ARIMA(0,1,1)x(0,1,1,52) [intercept=True]; AIC=nan, BIC=nan, Time=nan seconds
Fit ARIMA(0,1,0)x(0,1,0,52) [intercept=False]; AIC=1885.167, BIC=1887.457, Time=0.296 seconds
Fit ARIMA(1,1,0)x(0,1,0,52) [intercept=True]; AIC=1836.528, BIC=1843.399, Time=0.450 seconds
Fit ARIMA(1,1,0)x(2,1,0,52) [intercept=True]; AIC=nan, BIC=nan, Time=nan seconds
Fit ARIMA(1,1,0)x(1,1,1,52) [intercept=True]; AIC=nan, BIC=nan, Time=nan seconds
Fit ARIMA(1,1,0)x(0,1,1,52) [intercept=True]; AIC=nan, BIC=nan, Time=nan seconds
Fit ARIMA(2,1,0)x(2,1,1,52) [intercept=True]; AIC=nan, BIC=nan, Time=nan seconds
Fit ARIMA(0,1,0)x(1,1,0,52) [intercept=True]; AIC=1854.698, BIC=1861.479, Time=3.085 seconds
Fit ARIMA(2,1,0)x(1,1,0,52) [intercept=True]; AIC=1824.047, BIC=1835.499, Time=4.926 seconds
Fit ARIMA(2,1,0)x(0,1,0,52) [intercept=True]; AIC=1839.002, BIC=1848.164, Time=1.133 seconds
Fit ARIMA(2,1,0)x(2,1,0,52) [intercept=True]; AIC=nan, BIC=nan, Time=nan seconds
Fit ARIMA(2,1,0)x(1,1,1,52) [intercept=True]; AIC=nan, BIC=nan, Time=nan seconds
Fit ARIMA(2,1,0)x(0,1,1,52) [intercept=True]; AIC=nan, BIC=nan, Time=nan seconds
Fit ARIMA(2,1,0)x(2,1,1,52) [intercept=True]; AIC=nan, BIC=nan, Time=nan seconds
Fit ARIMA(3,1,0)x(1,1,0,52) [intercept=True]; AIC=1818.201, BIC=1831.944, Time=0.691 seconds
Fit ARIMA(3,1,0)x(0,1,0,52) [intercept=True]; AIC=1834.289, BIC=1845.741, Time=1.787 seconds
Fit ARIMA(3,1,0)x(2,1,0,52) [intercept=True]; AIC=nan, BIC=nan, Time=nan seconds
Fit ARIMA(3,1,0)x(1,1,1,52) [intercept=True]; AIC=nan, BIC=nan, Time=nan seconds
Fit ARIMA(3,1,0)x(0,1,1,52) [intercept=True]; AIC=nan, BIC=nan, Time=nan seconds
Fit ARIMA(3,1,0)x(2,1,1,52) [intercept=True]; AIC=nan, BIC=nan, Time=nan seconds
Fit ARIMA(3,1,1)x(1,1,0,52) [intercept=True]; AIC=1819.950, BIC=1835.984, Time=13.449 seconds
Fit ARIMA(2,1,1)x(1,1,0,52) [intercept=True]; AIC=1824.203, BIC=1837.946, Time=12.377 seconds
Total fit time: 50.209 seconds

```

Fig. 12. Parameter Tuning - SARIMA

model is greatly reduced. Also, the p-values the variables are within the standard significance level of 0.05. Thus, SARIMA model gives best results for the retail sales data and the forecast is shown in fig. 14

Statespace Model Results						
Dep. Variable:	total_amt			No. Observations:	126	
Model:	SARIMAX(3, 1, 0)x(1, 1, 0, 52)			Log Likelihood	-903.551	
Date:	Fri, 08 May 2020			AIC	1817.101	
Time:	00:06:35			BIC	1828.554	
Sample:	0			HQIC	1821.665	
	- 126					
Covariance Type:	opg					
	coef	std err	z	P> z	[0.025	0.975]
ar.L1	-0.6290	0.131	-4.811	0.000	-0.885	-0.373
ar.L2	-0.2707	0.079	-3.423	0.001	-0.426	-0.116
ar.L3	-0.2975	0.103	-2.891	0.004	-0.499	-0.096
ar.S.L52	-0.3747	0.033	-11.350	0.000	-0.439	-0.310
sigma2	2.663e+09	2.11e+11	1.26e+20	0.000	2.66e+09	2.66e+09
Ljung-Box (Q):	27.54			Jarque-Bera (JB):	96.07	
Prob(Q):	0.93			Prob(JB):	0.00	
Heteroskedasticity (H):	1.13			Skew:	-1.60	
Prob(H) (two-sided):	0.78			Kurtosis:	7.62	

Fig. 13. SARIMA - Summary

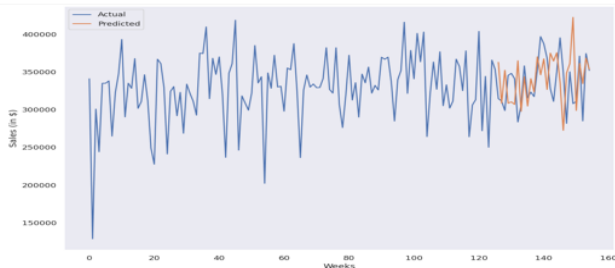


Fig. 14. SARIMA Forecast

V. CUSTOMER PRODUCT AFFINITY

Till now, two main aspects of retail marketing were covered. The former deals with the dissection of customers and the

latter help with sales prediction to enable better logistics and inventory management. A missing piece between these would be the type of products suggested to the segment of customers so that there is a higher chance of up-selling and cross-selling products. This is very similar to MBA (Market Basket Analysis), a widely used technique to identify the best possible combination of products or services frequently bought by customers. With the given constraints of data, the best possible solution was to apply advanced data analysis techniques and mining to provide close to accurate recommendations.

A. Approach

The approach to this problem would be to understand the purchase of different sub-categories by each customer of the retail store. A new line of analysis proposed from our end involves looking at the top3 sub-category purchases by each customer along with the chain of subcategories purchased by them. To get this recommendation more accurate, we leverage the customer segmentation to provide recommendations to different groups. The final step would be to provide offers on the set of preferred sub-categories to specific customers. Fig. 15 clearly explains this

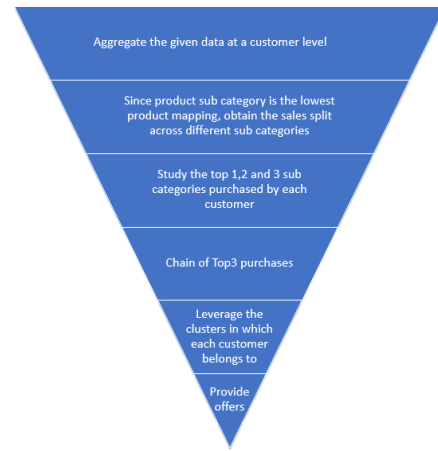


Fig. 15. Product Affinity Methodology

B. Results

The top 2 tables in fig. 16 shows the purchase analysis for the segment Regular Shoppers(RS). It is seen that the top 3 and the chain corresponds to women, men, and kids. This can help us achieve in deciding with a certain confidence that the Regular shoppers are more inclined towards the apparels section. Similarly, for the bottom 2 tables, the chain is more in line with the top2 and 3 purchases which shows the personal appliances and academic categories. This behavior is exhibited both by Bargain Hunters(BH) and High Spenders(HS) and providing such suggestions can have high sales conversions eventually leading to higher profits for the retail store. It is noted that the customer segmentation has had a bigger impact in grouping the customers and narrow down the analysis among the clusters

Top1	Sales/ Transactions	Top2	Sales/ Transactions	Top3	Sales/ Transactions
Women	2,789.82	Mens	2,577.28	Women	2,619.63
Mens	2,736.77	Women	2,571.26	Mens	2,634.96
Kids	2,601.20	Kids	2,795.82	Kids	2,613.44
Non-Fiction	2,595.62	Cameras	2,563.27	Non-Fiction	2,670.04
Mobiles	2,639.31	Bath	2,767.47	Mobiles	2,523.96

Top1	Sales/ Transactions	Top2	Sales/ Transactions	Top3	Sales/ Transactions
Women	2,789.82	Academic	2562.113	Personal Appliances	2525.163
Mens	2,736.77	Mens	2675.798	Academic	2590.627
Kids	2,601.20	Women	2530.436	Non-Fiction	2522.863
Tools	2,595.62	Kids	2514.252	Women	2629.5
DIY	2,639.31	Personal Appliances	2454.903	Mens	2473.734

Top 3 Chain	#Transactions	Sales/Transactions
Women-Kids-Mens	49	133,343.67
Kids-Mens-Women	52	118,348.82
Women-Mens-Cameras	25	76,968.76
Women-Cameras-Personal Appliances	29	92,019.98
Women-Mens-Academic	27	64,980.73

Top 3 Chain	#Transactions	Sales/Transactions
Women-Personal Appliances-Non-Fiction	141	353211.04
Mens-Academic-Personal Appliances	118	295118.98
Kids-Academic-Personal Appliances	66	166722.95
Tools-Academic-Personal Appliances	40	97484.205
Academic-Personal Appliances-Non-Fiction	42	93890.745

Fig. 16. Left - Cluster Regular Shoppers and Right - Cluster Bargain Hunters and High Spenders

VI. DISCUSSION

To summarize, there are a total of three problems addressed in this report. The customer segments are obtained as a result of the statistical K-means algorithm and yield 3 categories: Bargain Hunters(BH), High Spenders(HS), and Regular Shoppers (RS). This looks very relevant as we could see there are around ~41% transactions where the basket value is higher than the average basket value for the given period. Similarly, out of the total customers, around ~33% of them have visited quite often than the average frequency of customer visits. For a retail store, managing inventory and logistics play a very vital role, and in the movement of goods. Hence, forecasting techniques help the retail stores to be well prepared for a huge inflow of customers during certain periods. The final problem suggests the kind of products to be promoted to customers. A brief look at the sales of different sub-categories completes this analysis. The products suggested to the regular shoppers have the highest sales share in the retail store while the ones suggested to the spenders and bargain hunters have the lowest share. This proves that the latter ones are very seasonal and are quickly grabbed by the customers to take advantage of the offers.

VII. MEMBER CONTRIBUTIONS

- **Badrinarayanan R** – Studied and implemented different time series modelling by reading up about different techniques in the domain of retail forecasting – **28%**
- **Anirudh Sharma** – Worked on the cluster methodologies by understanding the type of data available and different kind of distance measures to be considered for the problem – **28%**
- **Anunay Sharma** – Discussed and carried out the product subcategory recommendation analysis through several logical formulations – **28%**
- **As a team** – Exploratory Data Analysis to get a sense of the data - **5%** each

VIII. ACKNOWLEDGMENT

We would like to thank the TA's and the Professor for guiding us and providing suggestions on framing different problem statements during the initial review session We

REFERENCES

- [1] N. T. Karim *et al.*, "Customer and target individual face analysis for retail analytics." *2018 International Workshop on Advanced Image Technology (IWAIT), Advanced Image Technology (IWAIT), 2018 International Workshop on*, pp. 1 – 4, 2018.
- [2] [Online]. Available: <https://towardsdatascience.com/an-overview-of-time-series-forecasting-models-a2fa7a358fcb>
- [3] T.-M. Choi, Y. Yu, and K.-F. Au, "A hybrid sarima wavelet transform method for sales forecasting." *Decision Support Systems*, vol. 51, no. 1, pp. 130 – 140, 2011.