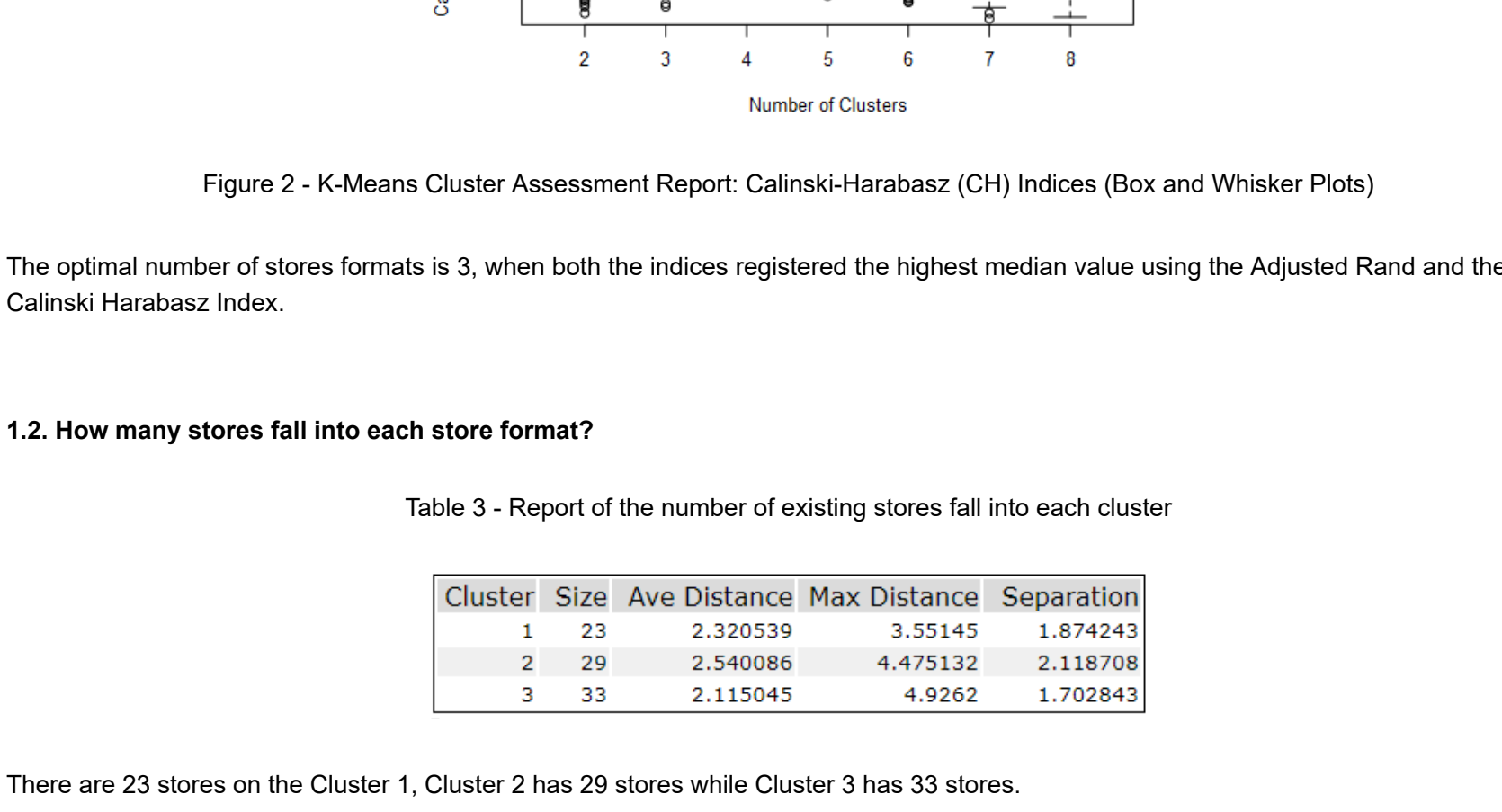
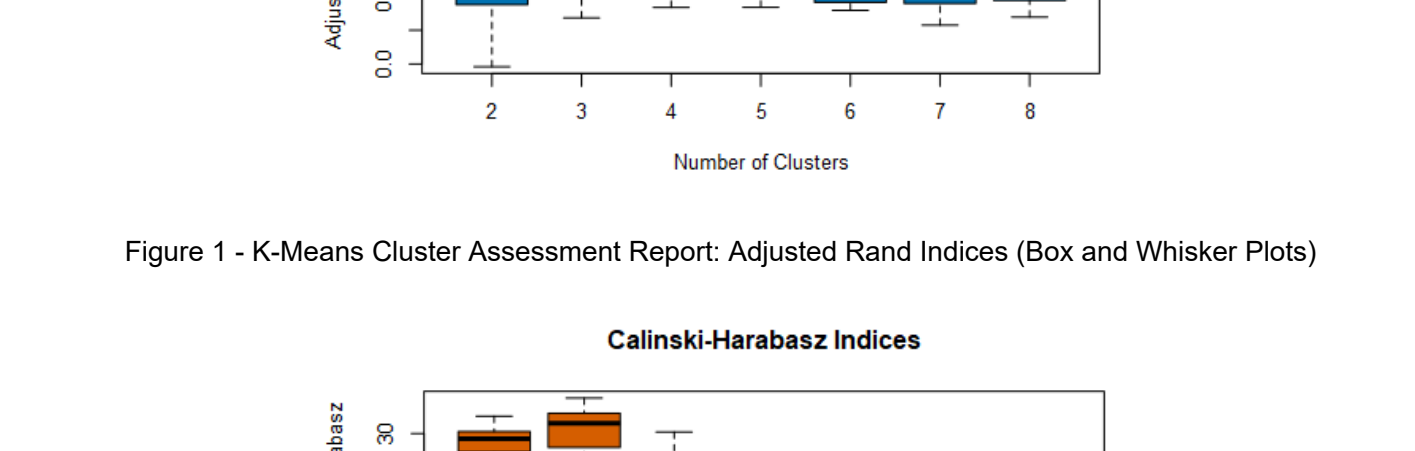
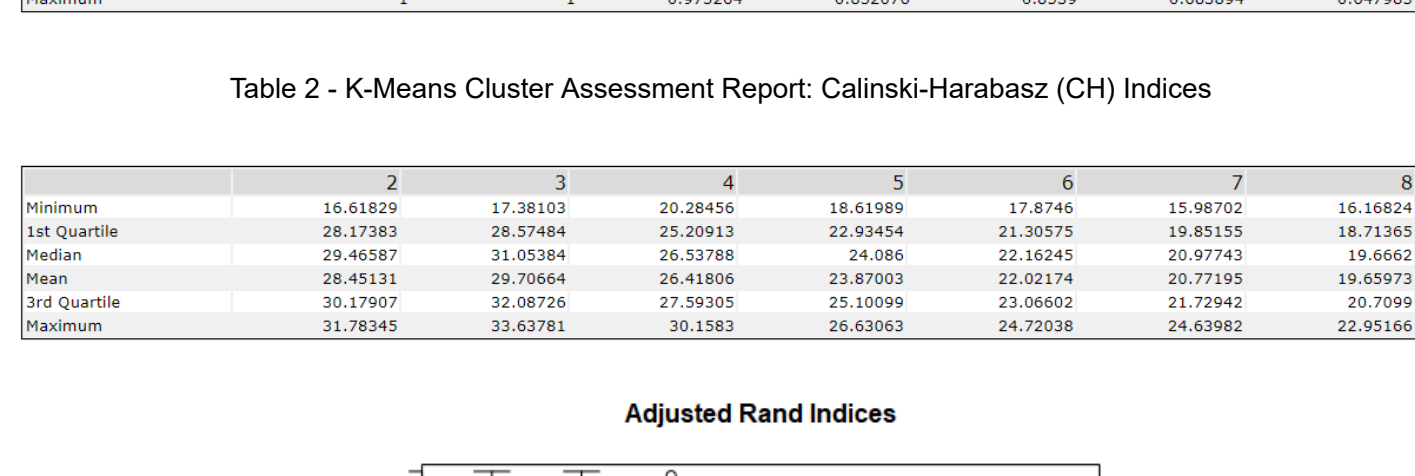


Project 5: Combine Predictive Techniques

1. Determine Store Formats for Existing Stores

1.1. What is the optimal number of store formats? How did you arrive at that number?



The optimal number of stores formats is 3, when both the indices registered the highest median value using the Adjusted Rand and the Calinski Harabasz Index.

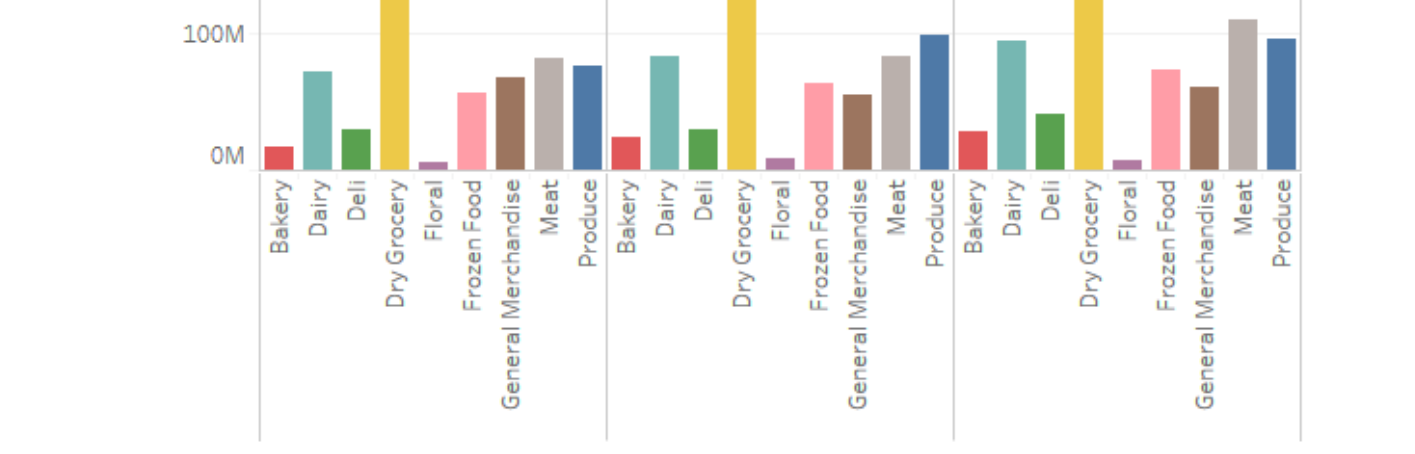
1.2. How many stores fall into each store format?

Table 3 - Report of the number of existing stores fall into each cluster

Cluster	Size	Ave Distance	Max Distance	Separation
1	23	2.320539	3.55145	1.874243
2	29	2.540086	4.475132	2.118708
3	33	2.115045	4.9262	1.702843

There are 23 stores on the Cluster 1, Cluster 2 has 29 stores while Cluster 3 has 33 stores.

1.3. Based on the results of the clustering model, what is one way that the clusters differ from one another?

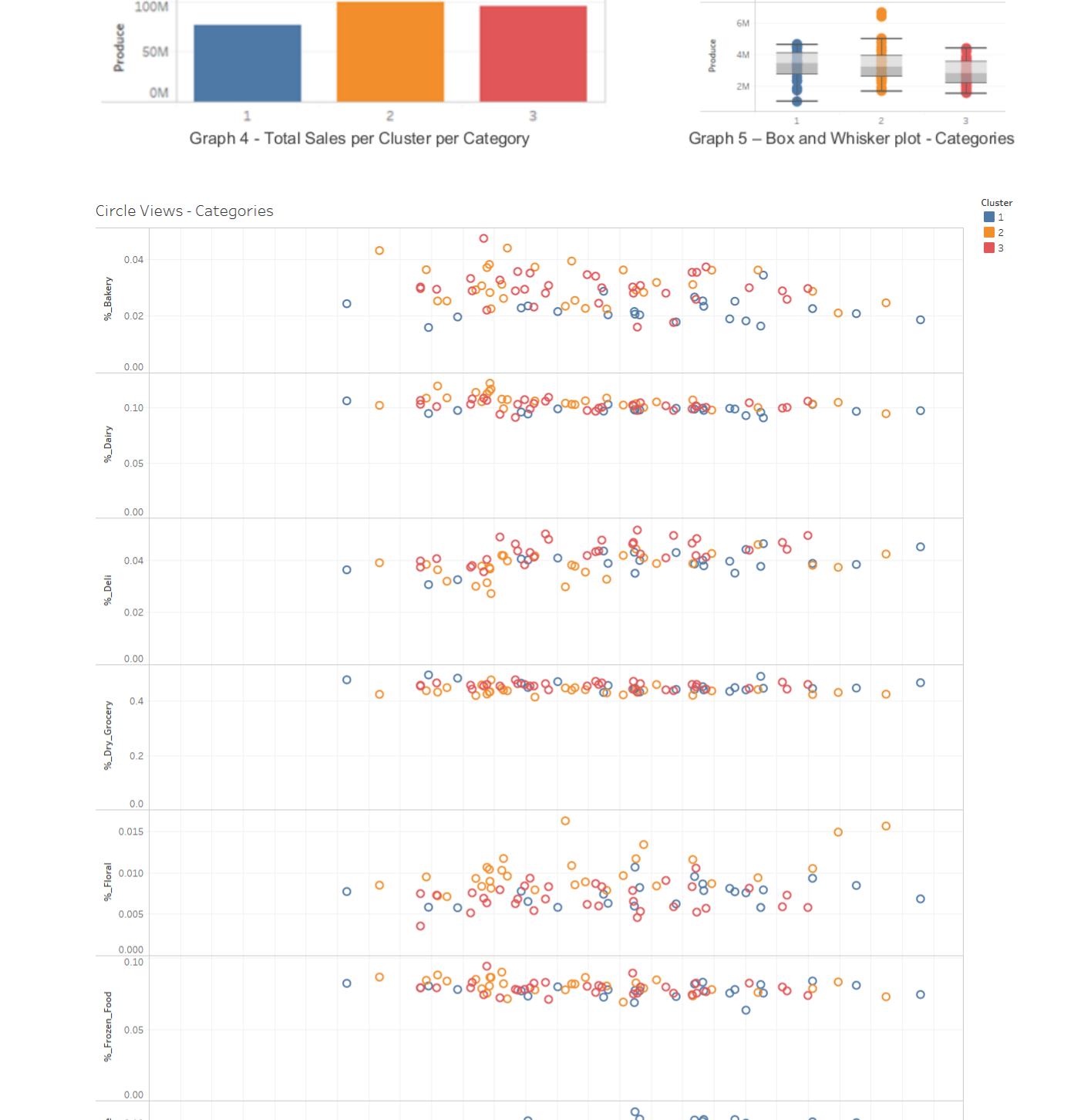


According to Graph 1, we could see that the Cluster 3 has the most total sales of them and has the most concentrated sales values, according to Graph 2, on the other hand, Cluster 1 has the lowest total sales and has the most sparse sales values.

Cluster 2 is on the middle ground between Cluster 1 and 3, both in total sales value and in spacing.



After analyzing above Graph 3, is remarkable the difference between the dry grocery category and the other values. One of the main factors of Cluster 3 has the highest sales values is because of its dry grocery category.

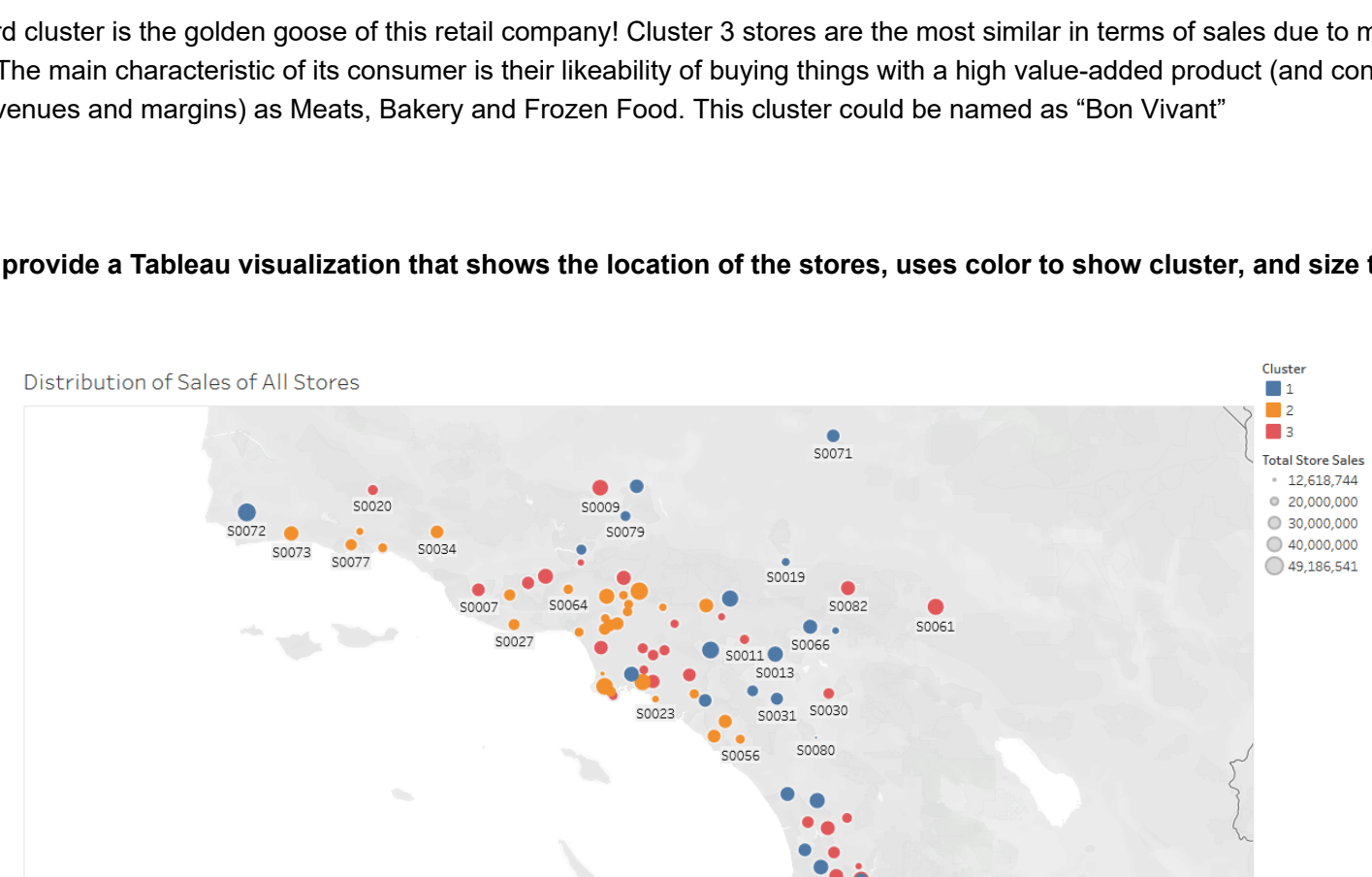


After analyzing Graph 4, Graph 5 and Graph 6, we could see that even though the dry grocery is the largest source of revenue for all clusters, we can see some specificities of each cluster, like the following:

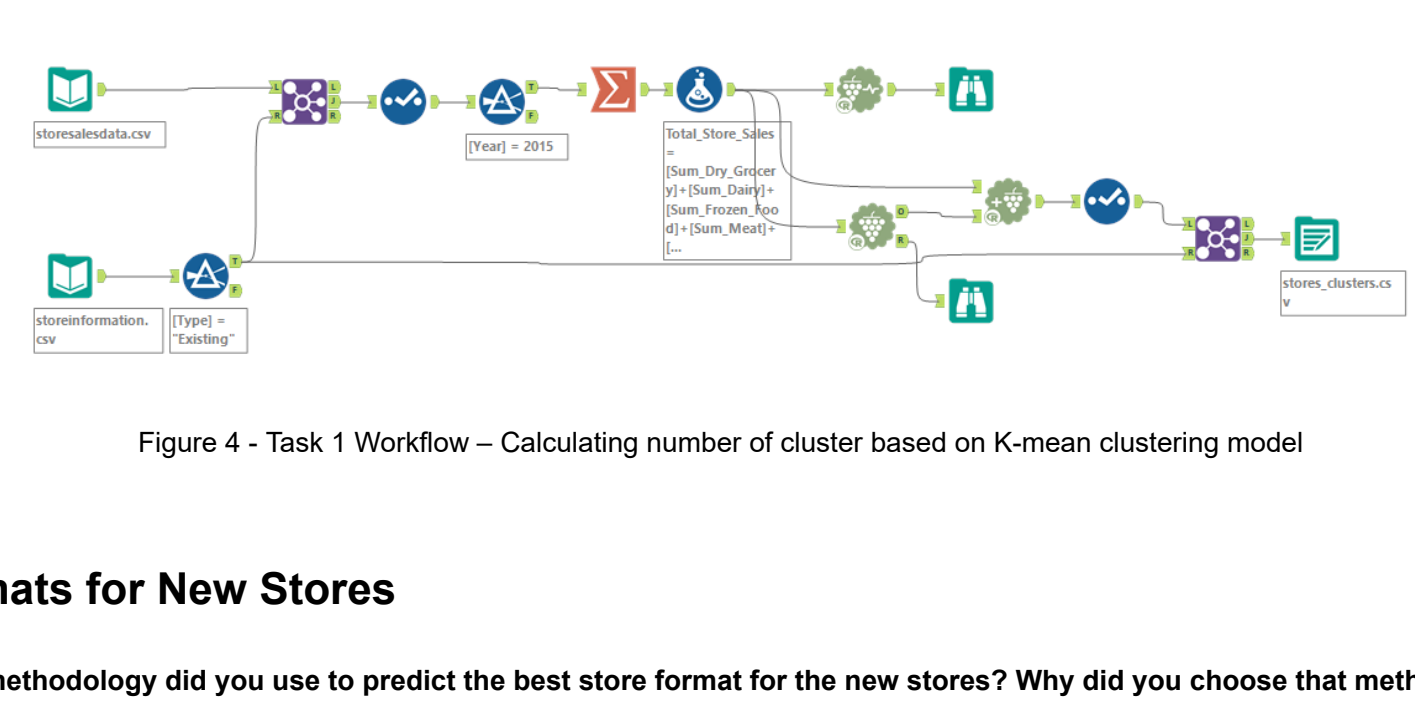
- Cluster 1:
 - Top 1 – General Merchandising
 - Top 2 – None
 - Top 3 – Bakery, Dairy, Deli, Dry Grocery, Floral, Frozen Food, Meat, Produce
- Cluster 2:
 - Top 1 – Floral, Produce
 - Top 2 – Bakery, Dairy, Deli, Dry Grocery, Frozen Food, Meat
 - Top 3 – General Merchandising
- Cluster 3:
 - Top 1 – Bakery, Dairy, Deli, Dry Grocery, Frozen Food, Meat
 - Top 2 – Floral, General Merchandising, Produce
 - Top 3 – None

- The first cluster has significant higher General Merchandising sales than other categories and the source of revenue of this cluster comes very diluted from each category, besides being the cluster with less effectiveness in sales. This cluster could be named as "Just Things";
- The second cluster has the most balanced incomes and the main characteristic of its consumer is their likeability of buying Produce and Floral things and not likely to buying General Merchandising. This cluster could be named as "Organics";
- The third cluster is the golden goose of this retail company! Cluster 3 stores are the most similar in terms of sales due to more compact range. The main characteristic of its consumer is their likeability of buying things with a high value-added product (and consequently high revenues and margins) as Meats, Bakery and Frozen Food. This cluster could be named as "Bon Vivant"

1.4. Please provide a Tableau visualization that shows the location of the stores, uses color to show cluster, and size to show total sales.



1.5. Workflow



2. Formats for New Stores

2.1. What methodology did you use to predict the best store format for the new stores? Why did you choose that methodology?

The following Model Comparison Report shows comparison matrix between Decision Tree, Forest Model and Boosted Model. Boosted Model is chosen despite having same accuracy as Forest Model due to higher F1 value.

Model Comparison Report				
Fit and error measures				
Model	Accuracy	F1	Accuracy_1	Accuracy_2
DT_Cluster	0.7059	0.7685	0.7500	1.0000
FM_Cluster	0.8235	0.8426	0.7500	1.0000
BM_Cluster	0.8235	0.8889	1.0000	1.0000

Model: model names in the current comparison.

Accuracy: overall accuracy; number of correct predictions of all classes divided by total sample number.

Accuracy_[class name]: accuracy of Class [class name] is defined as the number of cases that are correctly predicted to be Class [class name] divided by the total number of cases that actually belong to Class [class name]; this measure is also known as recall.

AUC: area under the ROC curve; only available for two-class classification.

F1: F1 score, 2 * precision * recall / (precision + recall). The precision measure is the percentage of actual members of a class that were predicted to be in that class divided by the total number of cases predicted to be in that class. In situations where there are three or more classes, average precision and average recall values across classes are used to calculate the F1 score.

Confusion matrix of BM_Cluster			
	Actual_1	Actual_2	Actual_3
Predicted_1	4	0	1
Predicted_2	0	4	2
Predicted_3	0	0	6

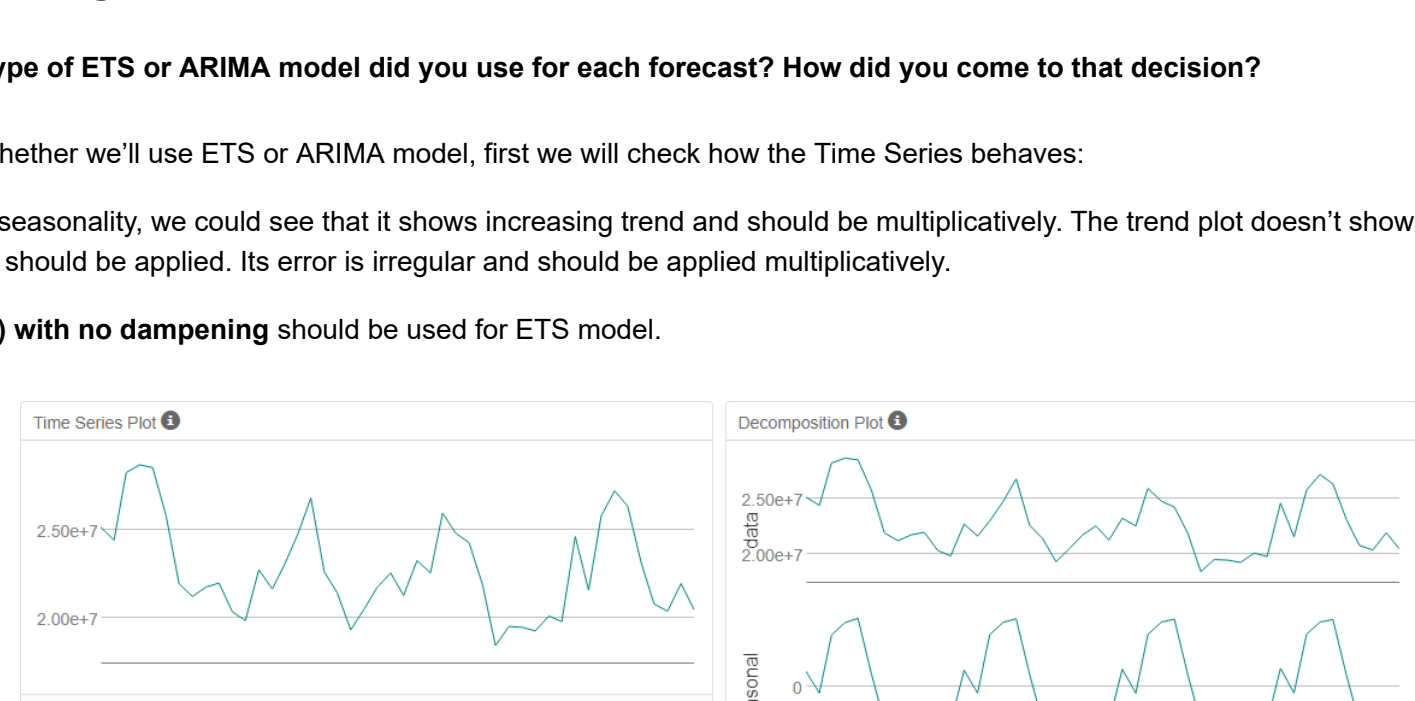
Confusion matrix of DT_Cluster			
	Actual_1	Actual_2	Actual_3
Predicted_1	3	0	2
Predicted_2	0	0	2
Predicted_3	1	0	5

Confusion matrix of FM_Cluster			
	Actual_1	Actual_2	Actual_3
Predicted_1	3	0	1
Predicted_2	0	4	1
Predicted_3	1	0	7

2.2. What format do each of the 10 new stores fall into?

Store Number	Segment
S0086	2
S0087	1
S0088	3
S0089	2
S0090	2
S0091	1
S0092	2
S0093	1
S0094	2
S0095	2

2.3. Workflow



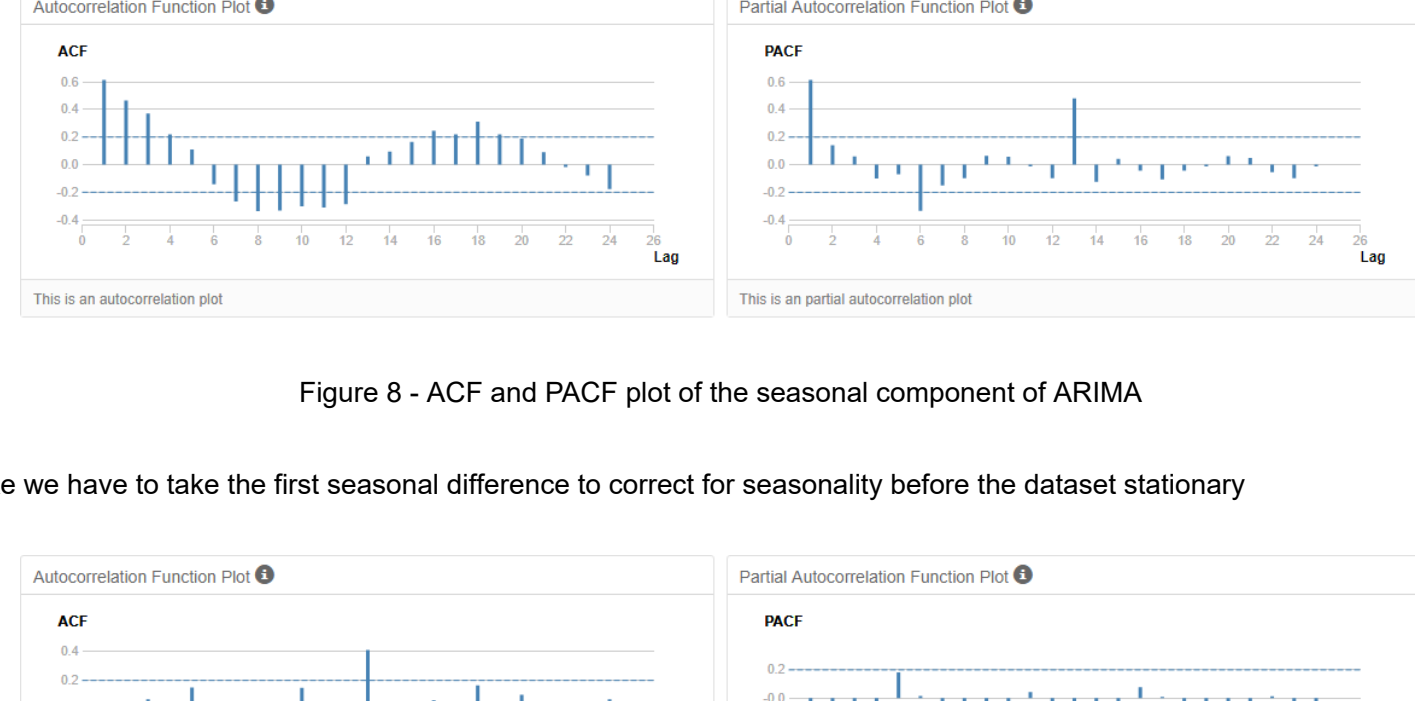
3. Predicting Produce Sales

3.1. What type of ETS or ARIMA model did you use for each forecast? How did you come to that decision?

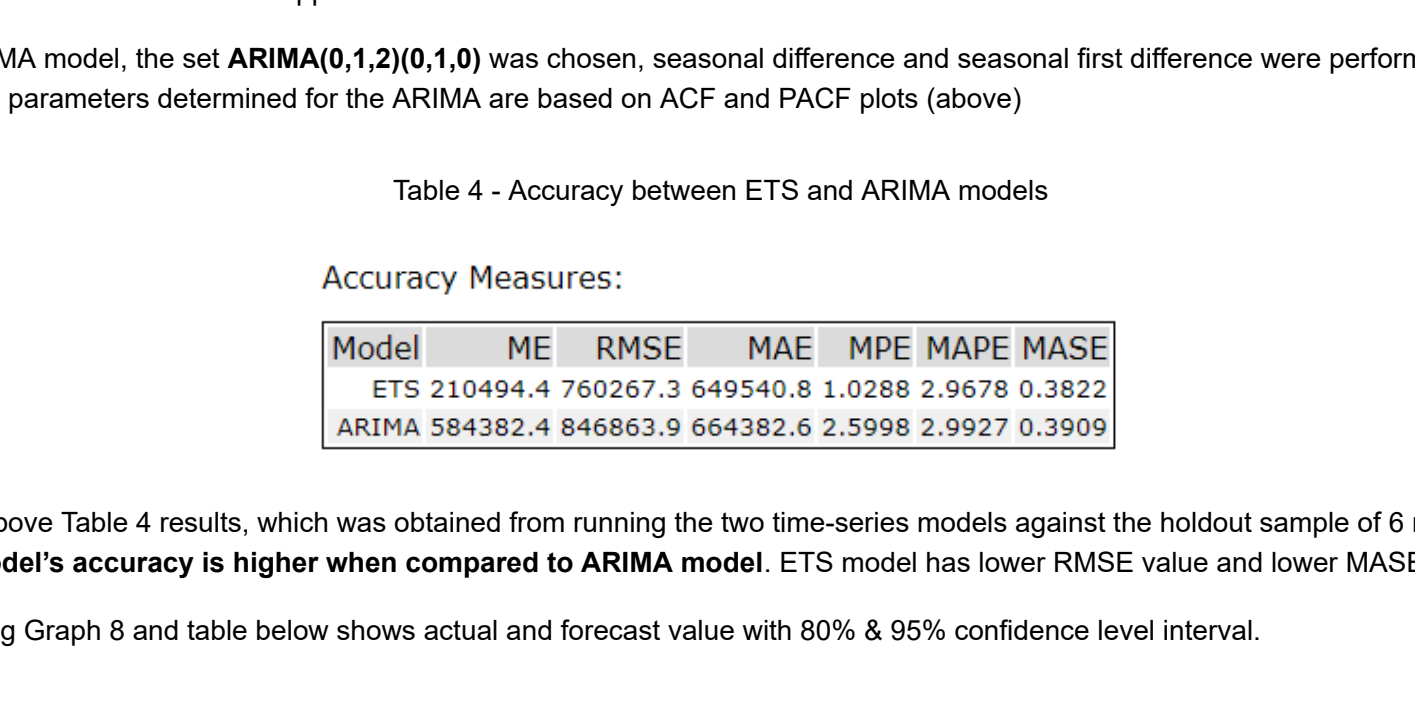
To decide whether we'll use ETS or ARIMA model, first we will check how the Time Series behaves:

Looking for seasonality, we could see that it shows increasing trend and should be multiplicatively. The trend plot doesn't show any trending, and nothing should be applied. Its error is irregular and should be applied multiplicatively.

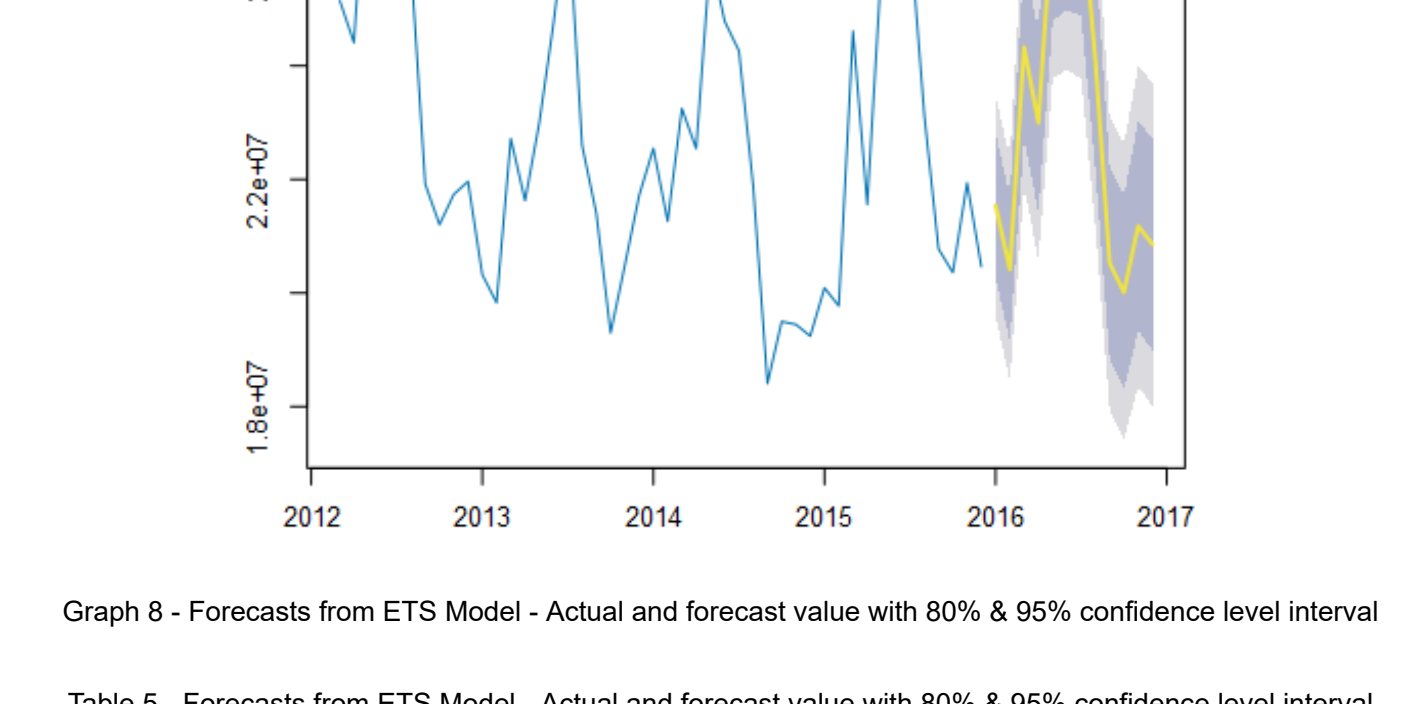
ETS(M,N,M) with no dampening should be used for ETS model.



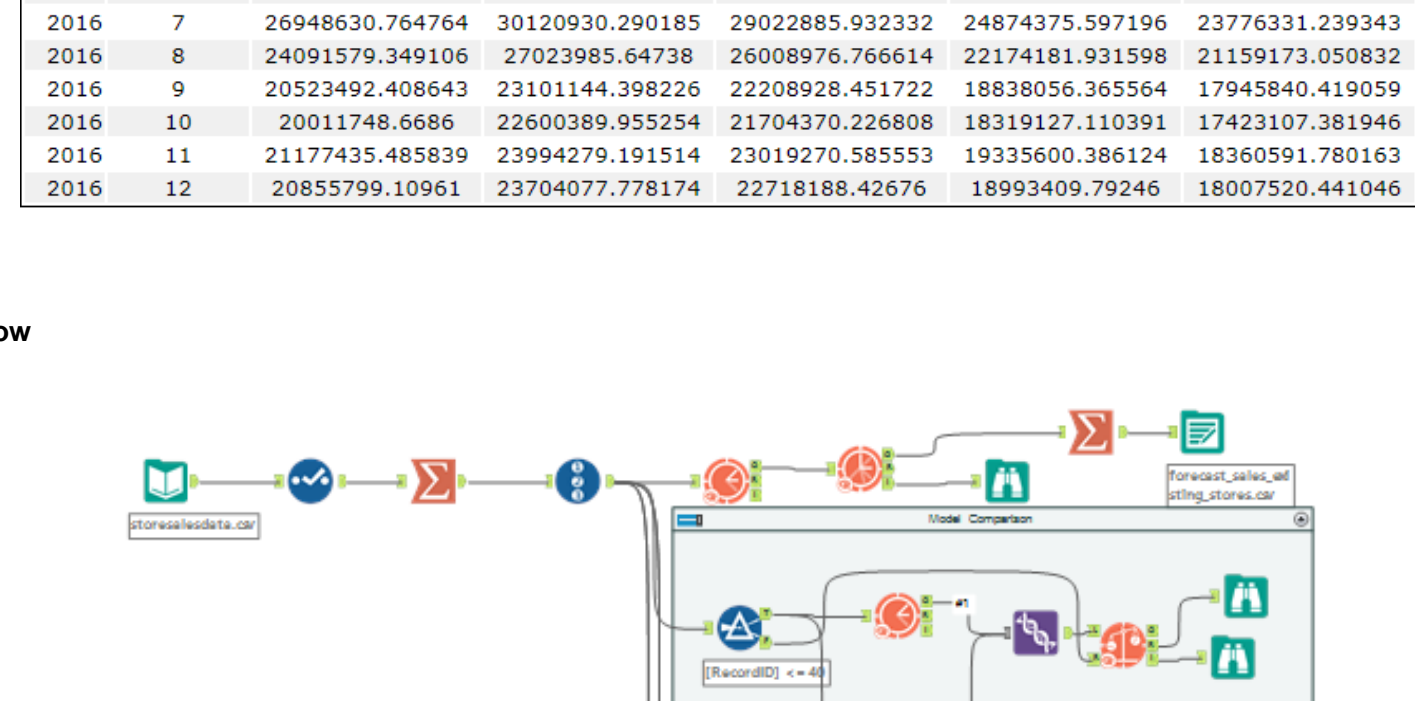
Because of the seasonality on these series, we need to differentiate our Time Series in order to Stationarize the series, as following.



First, we need to look at the seasonal differencing component, to allow us to account for the value as observed in the same season one year earlier, as figure below.



So, looks like we have to take the first seasonal difference to correct for seasonality before the dataset stationary



After plotting the first seasonal difference, we can see that the series has stationarized. We can see this through our ACF and PACF plots, the serial correlational has now disappeared.

For the ARIMA model, the set ARIMA(0,1,2)(0,1,0) was chosen, seasonal difference and seasonal first difference were performed. There is a lag-2. The parameters determined for the ARIMA are based on ACF and PACF plots (above)

Table 4 - Accuracy between ETS and ARIMA models

Accuracy Measures:						
Model	ME	RMSE	MAE	MPE	MAPE	MASE
ETS	210494.4	760267.3	649540.8	1.0288	2.9678	0.3822
ARIMA	584382.4	846863.9	664382.6	2.5998	2.9927	0.3909

Based on Table 4 results, which was obtained from running the two time-series models against the holdout sample of 6 months data, the ETS model's accuracy is higher when compared to ARIMA model. ETS model has lower RMSE value and lower MASE value.

The following Graph 8 and table below shows actual and forecast value with 80% & 95% confidence level interval.

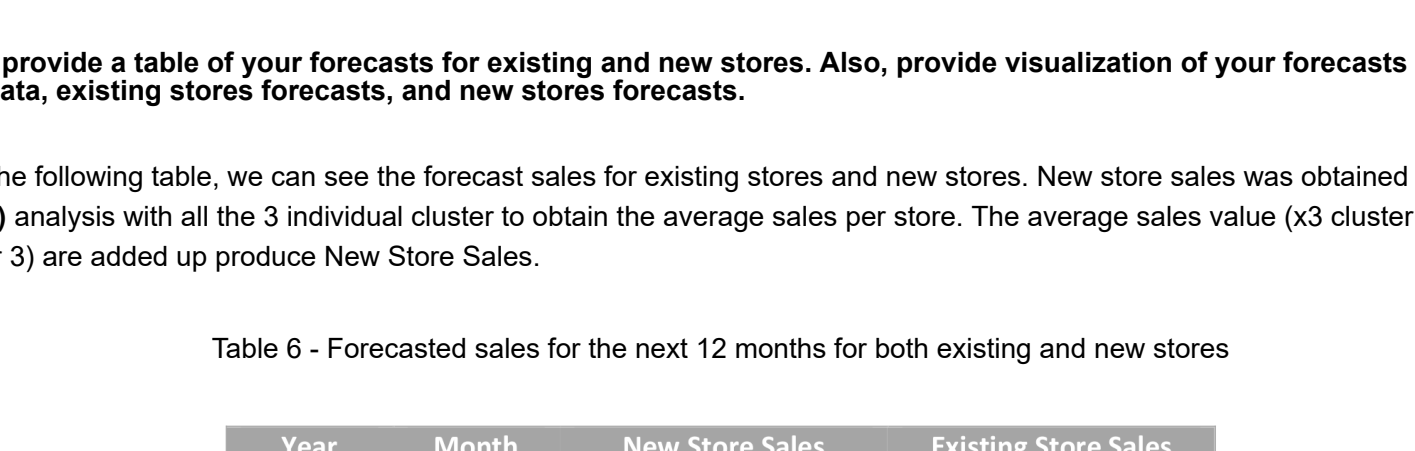


Table 5 - Forecasts from ETS Model - Actual and forecast value with 80% & 95% confidence level interval

Period	Sub_Period	forecast_high_95	forecast_high_80	forecast_low_80	forecast_low_95
2016	1	21539936.007499	23479964.557336	22808452.492932	20271419.522066
2016	2	20413770.60136	22357792.702597	21684898.329698	19142642.873021
2016	3	24325953.097628	26761721.213559	25918616.262307	22732369.932948
2016	4	22993466.348585	25403233.826166	24569128.696553	21417604.087517
2016	5	26691951.419156	29608731.673669	28599131.515834	24784771.322478
2016	6	26989964.010552	3005322.497686	28994294.191682	24985633.829422
2016	7	26948630.764754	30120930.290185	29022885.932332	24874375.597196
2016	8	24091579.349106	27023985.64738	26008976.766614	22174181.931598
2016	9	20523492.408643	2310144.398226	22208928.451722	18838056.365564
2016	10	2001748.6586	22600389.95254	21794370.226808	18319127.110391
2016	11	21177435.485839	23994279.101514	23019270.585553	19325600.386124
2016	12	20855799.10961	23704077.778174	22718188.42676	18993409.79246
	Total				

3.2. Workflow



3.3. Please provide a table of your forecasts for existing and new stores. Also, provide visualization of your forecasts that includes historical data, existing stores forecasts, and new stores forecasts.

Looking at the following table, we can see the forecast sales for existing stores and new stores. New store sales was obtained by using ETS(M,N,M) analysis with all the 3 individual cluster to obtain the average sales per store. The average sales value (x3 cluster 1, x6 cluster 2, x1 cluster 3) are added up produce New Store Sales.

Table 6 - Forecasted sales for the next 12 months for both existing and new stores

Year	Month	New Store Sales	Existing Store Sales
2016	1	\$ 2,587,450.85	\$ 21,539,936.01
2016	2	\$ 2,477,352.89	\$ 20,413,770.60
2016	3	\$ 2,913,185.24	\$ 24,325,953.10
2016	4	\$ 2,775,745.61	\$ 22,993,466.35
2016	5	\$ 3,150,866.84	\$ 26,691,951.42
2016	6	\$ 3,188,922.00	\$ 26,989,964.01
2016	7	\$ 3,214,745.65	\$ 26,948,630.76
2016	8	\$ 2,866,348.66	\$ 24,091,579.35
2016	9	\$ 2,538,726.85	\$ 20,523,492.41
2016	10	\$ 2,488,148.29	\$ 20,011,748.67
2016	11	\$ 2,595,270.39	\$ 21,177,435.49
2016	12	\$ 2,573,396.63	\$ 20,855,799.11
	Total	\$ 33,370,159.89	\$ 276,563,727.27

Using these predictive techniques, the customer can minimize investment risk and know what the expected profit values should be.

