

Veer Pal Singh Data Scientist March 20, 2022

Table of contents

01	Introduction	
02	Datasets	
03	Data Preparation	
04	Feature Engineering	
05	EDA Analysis	
06	Model Building	
07	Dashboard	
08	Model Selection	
09	Visualization	
10	Recommendations	

Market Mix Model: Budget Optimization

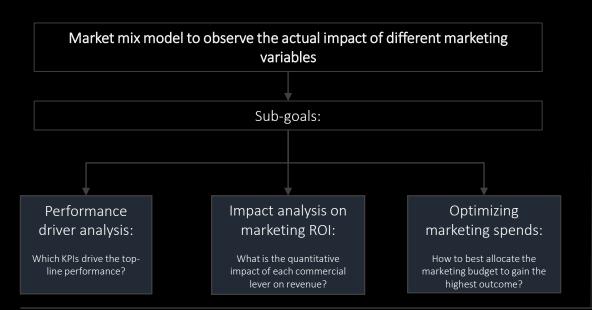
ElecKart is an e-commerce firm based out of Ontario, Canada specializing in electronic products. Over the last one year, they had spent a significant amount of money on marketing. Occasionally, they had also offered big-ticket promotions (similar to the Big Billion Day). They are about to create a marketing budget for the next year, which includes spending on commercials, online campaigns, and pricing & promotion strategies. The CFO feels that the money spent over the last 12 months on marketing was not sufficiently impactful, and, that they can either cut on the budget or reallocate it optimally across marketing levers to improve the revenue response using Market Mix Model.

Imagine that you are a part of the marketing team working on budget optimization. You need to develop a market mix model to observe the actual impact of different marketing variables over the last year.

Using your understanding of the model, you have to recommend the optimal budget allocation for different marketing levers for the next year.

Business Objective

To develop a market mix model for ElecKart (an e-commerce firm based out of Ontario, Canada) for 3 product sub-categories - Camera Accessory, Gaming Accessory, and Home Audio - to observe the actual impact of various marketing variables over one year (July 2015 to June 2016) and recommend the optimal budget allocation for different marketing levers for the next year.



Methodology

The approach for this project has been designed to follow the CRISP-DM Framework. The various stages of the framework are represented below in a sequential flow:



Codes: https://github.com/veer2701/Market-Mix-Model/tree/main

Description of Datasets

The data was collected from multiple sources with the following information:

- Main Consumer file with order details had features such as order date, order ID, gross merchandise value, number of units sold for specific products, etc.
- Media Investment file with monthly spend on various advertising channels the amount invested in each advertising medium for the past year
- Sale Calendar file showing dates from the past year when there was a promotional offer
- Monthly NPS score file showing net promotion score and company stock value for last year
- Stock Index of the company on a monthly basis
- Product details with information like product category and vertical (camera accessory, home audio, and gaming accessory)
- Weather file having detailed weather reports from last year in the state of Ontario, Canada

Product, price, promotion, and place form the four Ps of the marketing mix. These are the key factors that are involved in introducing a product or service to the public.



Data Preparation and Data Cleaning

Data preparation is the process of preparing raw data so that it is suitable for further processing and analysis. Key steps include collecting, cleaning, and labeling raw data into a form suitable for machine learning (ML) algorithms and then exploring and visualizing the data.

Handling Incorrect values in some columns

- Imputing "\N" value in delivery days &delivery days by 0
- Treating incorrect GMV values (where gmv > product_mrp units) by imputing the faulty MRP values with GMV/units
- Handling Negative values for product_procurement_sla, deliverybdaysdelivery daysays by dropping them
- Handling large values(0.3%) for product_procurement_sla by dropping them

<u>De-Duplication ahead and dropped them of</u> Data

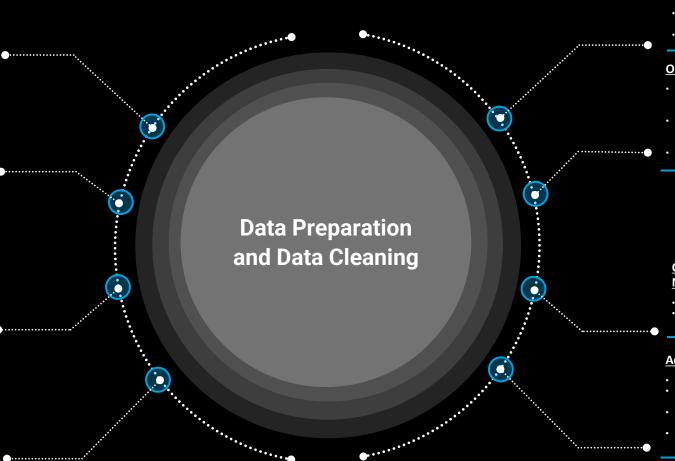
 After converting all column values to lower case, we see that there are around 99283 (6.33%) rows that are duplicates. We went De-Duplication ahead and dropped them

Treating Null values and Whitespaces

- Initially there weren't any NULL values in the dataframe.
 However, there were quite a few Whitespaces present in some of the columns in the dataframe
- We first converted these whitespaces to NaNs and the dropped these values

Selecting One Year Data

 Selecting1 Year Data from July, 2015 – June, 2016. In the process, 592 records were dropped



Dropping Insignificant columns

- Dropping Columns with Single Unique Value (as it doesn't add any information to the analysis)
- Dropping some of the 'ld' Columns which are insignificant to the analysis

Outlier Treatment

- Since we have already deleted some records on erroneous grounds, in order that we don't lose any further data, we chose not to delete outlier values
- For the variables 'SLA', 'deliverybdays', 'deliverybdays', 'gmv', 'product_mrp', 'list_price' where outliers are present, we CAPPED the values above 99 percentile to the value corresponding to 99 percentile
- Thus the outliers couldn't affect the predictive model while at the same time there was enough data to build a generalizable

Converting Categorical Attributes to Numerical Form

- · Binary encoding for categorical variable with 2 levels
- One Hot Encoding for categorical variable with multiple levels by creating dummy variables

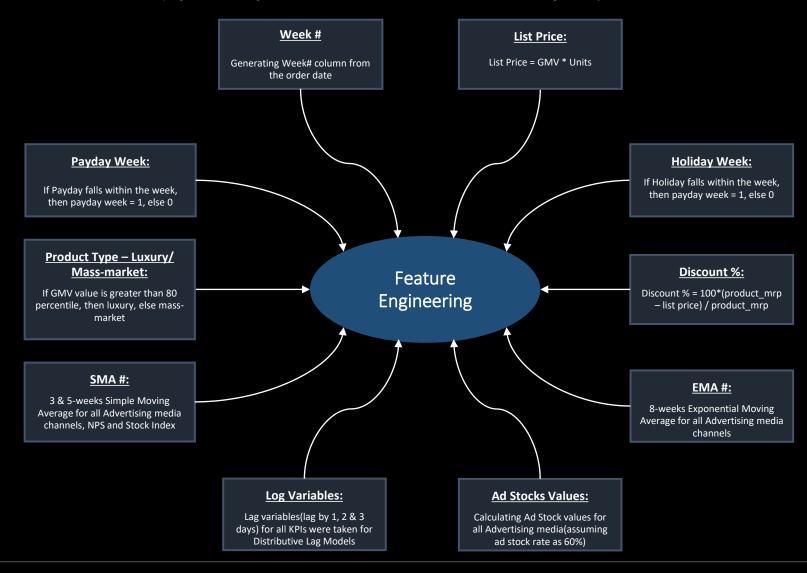
Additional Data Preparation for Model Building

- Merging Order dataset with all other secondary dataframes
- Extracting 3 separate dataframes for 3 product subcategories camera accessory, home audio and gaming accessory
- Roll Up daily Order Data to Weekly Level by aggregating the numeric variables based on Week
- Scaling and dividing the master dataframes into train and test datasets for all 3 product subcategories

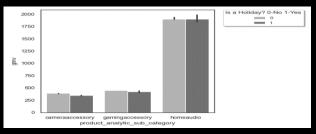
Codes: https://github.com/veer2701/Market-Mix-Model/tree/main

Feature Engineering: Creation of New KPI's

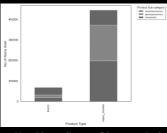
Marketing KPIs are measurable metrics your company can track to gauge performance over channel-specific marketing activities. It proves the Return on Investment (ROI) about the effectiveness of marketing campaigns and strategies so the business can double down on what's working and adjust for what's not.



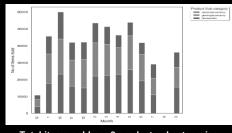
EDA Analysis & Visualization:



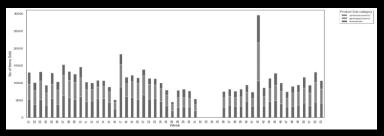
Average Revenue from Holiday/Non-holiday days for the 3 product subcategories



No of items(Luxury/Mass-market) sold per 3 product subcategories



Total items sold per 3 product subcategories per Month



Total items sold per 3 product subcategories per Week

Top 10 Product Verticals which brought the Maximum Revenue for 3 product sub-categories

	product_analytic_sub_category	$product_analytic_vertical$	gmv	product_in_category
0	homeaudio	homeaudiospeaker	1.873206e+08	homeaudiospeaker in homeaudio
1	cameraaccessory	lens	1.085308e+08	lens in cameraaccessory
2	gamingaccessory	gamepad	6.187440e+07	gamepad in gamingaccessory
3	gamingaccessory	gamingheadset	3.199049e+07	gamingheadset in gamingaccessory
4	cameraaccessory	binoculars	2.658427e+07	binoculars in cameraaccessory
5	gamingaccessory	gamingmouse	2.632837e+07	gamingmouse in gamingaccessory
6	cameraaccessory	camerabattery	2.356174e+07	camerabattery in cameraaccessory
7	cameraaccessory	camerabag	2.249499e+07	camerabag in cameraaccessory
8	cameraaccessory	flash	2.228150e+07	flash in cameraaccessory
9	homeaudio	fmradio	2.222170e+07	fmradio in homeaudio

Home Audio Speaker under Home Audio segment brought the largest revenue followed by Camera Lens under Camera Accessory & Gamepad under Gaming Accessory

Top 10 Product Verticals with most no of sales for 3 product sub-categories

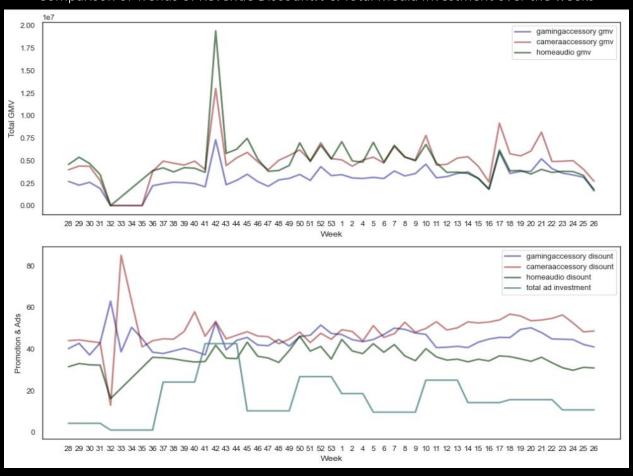
	product_analytic_sub_category	product_analytic_vertical	units	product_in_category
0	homeaudio	homeaudiospeaker	76581	homeaudiospeaker in homeaudio
1	gamingaccessory	gamingheadset	59928	gamingheadset in gamingaccessory
2	gamingaccessory	gamepad	52437	gamepad in gamingaccessory
3	cameraaccessory	flash	47808	flash in cameraaccessory
4	gamingaccessory	gamingmouse	35470	gamingmouse in gamingaccessory
5	cameraaccessory	camerabattery	35107	camerabattery in cameraaccessory
6	cameraaccessory	lens	32350	lens in cameraaccessory
7	cameraaccessory	cameratripod	31220	cameratripod in cameraaccessory
8	homeaudio	fmradio	24681	fmradio in homeaudio
9	cameraaccessory	camerabag	15842	camerabag in cameraaccessory

Home Audio Speaker under Home Audio segment had the most no of sales followed by Gaming Headset & Gamepad under Gaming Accessory

EDA Analysis & Visualization (contd.):

Time series analysis helps organizations understand the underlying causes of trends or systemic patterns over time. Using data visualizations, business users can see seasonal trends and dig deeper into why these trends occur. With modern analytics platforms, these visualizations can go far beyond line graphs.

Comparison of Trends of Revenue Discount% & Total Media Investment over the weeks



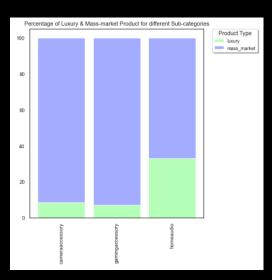
The following observations were noted from the above time series plots:

- For the week #42 (during Thanks giving), all the graphs show a steep rise. Revenue increased because of both higher discount% and increased Ad Investment.
- For the week 32(August), Revenue generated was the lowest from all 3 product subcategories. This can be observed as a direct relation to minimum amount of total investment in Ads. Discount was also lowest for all products apart from camera accessories. Post this dip in revenue, discount% was increased to bring about higher sales. This increase in Discount% was observed most in the case of gaming accessories. However, barring home audio products, the revenue from other products was seen to be constant for the next 3 weeks after which, the revenue started to pick up.
- In general the average discount % offered for home audio products is lesser compared to that of the other product subcategories.

EDA Analysis & Visualization:

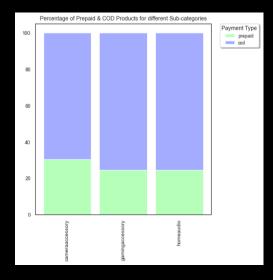
Finding the percentage of Luxury & Massmarket Products from 3 sub-categories

	$product_analytic_sub_category$	luxury	mass_market
0	cameraaccessory	18423	197478
1	gamingaccessory	13007	172869
2	homeaudio	36747	74314



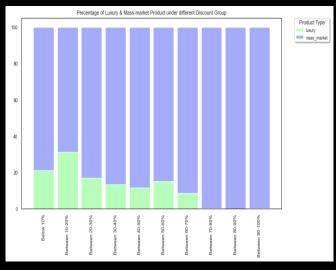
Percentage of luxury products under Home Audio is much more compared to the other sub categories. Finding the percentage of COD & Prepaid Products from 3 sub-categories

	product_analytic_sub_category	prepaid	cod
0	cameraaccessory	65462	150439
1	gamingaccessory	45273	140603
2	homeaudio	27123	83938



Percentage of prepaid payments under Camera Accessory was observed to be slightly more compared to that of the other sub categories. Finding the percentage of Luxury and Mass_market Products under different Discount groups

	Discount Bins	luxury	mass_market
0	Below 10%	9110	34093
1	Between 10-20%	14604	31972
2	Between 20-30%	10077	48907
3	Between 30-40%	8593	56255
4	Between 40-50%	7477	57617
5	Between 50-60%	13335	73808
6	Between 60-70%	4847	51506
7	Between 70-80%	68	41789
8	Between 80-90%	48	39781
9	Between 90-100%	18	8933



Percentage of luxury products were given a discount between 10-20%.

Description of Model Building

The primary objective of the case study being Revenue prediction and determination of important KPIs that influence the revenue growth, we have to build the following Linear Regression models:

1

Additive Model

Linear model is used to capture the current effect of several KPIs. This model assumes an additive relationship between the different KPIs. Hence their impacts are also additive towards the dependent Y variable. The equation can be represented as:

 $Y = \alpha + \beta 1At + \beta 2Pt + \beta 3Dt + \beta 4Qt + \beta 5Tt + \epsilon$

2

Multiplicative Model

Multiplicative model is used when there are interactions between the KPIs. To fit a multiplicative model, take logarithms of the data(on both sides of the model), then analyze the log data as before.

 $Y = e^{\alpha} . X1^{\beta}1 . X2^{\beta}2 . X3^{\beta}3 . X4^{\beta}4 . X5^{\beta}5 + \epsilon \text{ and } InY = \alpha + \beta 1In(X1) + \beta 2In(X2) + \beta 3In(X3) + \beta 4In(X4) + \beta 5In(X5) + \epsilon'$

3

Koyck Model

Koyck model is used to capture the carry-over effect of different KPIs, ie.to model the current revenue figures based on the past figures of the KPIs. The Koyck tells us that the current revenue generated is not just influenced by the different independent attributes, but also because of the revenue generated over the last periods. Yt = α + β 1X1 + β 2X2 + β 3X3 + β 4X4 + β 5X5 + ϵ and Yt = α + μ Yt-1 + β 1X1 + β 2X2 + β 3X3 + β 4X4 + β 5X5 + ϵ

4

Distributive Lag Model (Additive)

This is a more generalizable distributed lag model that captures the carry-over effect of all the variables:

 $Yt = \alpha + \mu 1 Yt - 1 + \mu 2 Yt - 2 + \mu 3 Yt - 3 + \dots \bullet + \beta 1 X1t + \beta 1 X1t - 1 + \beta 1 X1t - 2 + \dots \bullet + \beta 2 X2t + \beta 2 X2t - 1 + \beta 2 X2t - 2 + \dots \bullet + \beta 3 X3t + \beta 3 X3t - 1 + \beta 3 X3t - 2 + \dots \bullet + \beta 4 X4t + \beta 4 X4t - 1 + \beta 4 X4t - 2 + \dots \bullet + \beta 5 X5t + \beta 5 X5t - 1 + \beta 5 X5t - 2 + \dots \bullet + \epsilon$

5

Distributive Lag Model (Multiplicative)

Distributive Lag Model(Multiplicative) will help us capture the interactions between current and carry over effects of the KPIs. $Yt = \alpha + \mu 1 \ln(Yt-1) + \mu 2 \ln(Yt-2) + \mu 3 \ln(Yt-3) + \dots \bullet + \beta 1 \ln(X1t-1) + \beta 1 \ln(X1t-2) + \dots \bullet + \beta 2 \ln(X2t-1) + \beta 2 \ln(X2t-1) + \beta 2 \ln(X2t-2) + \dots \bullet + \beta 3 \ln(X3t-1) + \beta 3 \ln(X3t-2) + \dots \bullet + \beta 4 \ln(X4t-1) + \beta 4 \ln(X4t-2) + \dots \bullet + \beta 5 \ln(X5t-1) + \beta 5 \ln(X5t-2) + \dots \bullet + \epsilon'$

Dashboard:

The following table contains the details of all models built, their accuracy scores and the top 5 KPIs returned by them:

Product Sub-category	Linear Regression Model	Cross Validation	R2 Score MSE	Score	Top 5 KPis
	Additive	No Yes	0.83 0 -0.8	0.17 1.08	product_vertical_lens, product_vertical_camerabattery, product_vertical_camerabag, product_vertical_camerabousing. Online marketing
	Multiplicative	No Yes	0.84	0.36	product_vertical_lens, product_vertical_camerabattery, is_mass_market, product_vertical_camerabatterycharger, TV
cameraaccessory	Koyck	No Yes	② 0.84 ② ③ 0.27 ④	0.16	product_vertical_lens, product_vertical_camerabag, product_vertical_camerahousing, product_vertical_camerabattery, Online marketing
	Distributive Lag Model (Additive)	No Yes	② 0.87 ② ③ 0.82 ②	0.12	product_vertical_lens, product_vertical_filter, product_vertical_camerabag, product_vertical_cameraremotecontrol, is_mass_market
	Distributive Lag Model (Multiplicative)	No Yes	0.77 0.82	0.5 0.18	is_mass_market, product_vertical_lens, product_vertical_cameraaccessory, product_vertical_camerabattery, product_vertical_cameratripod
		No	② 0.93 ②	0.05	
	Additive	Yes	0.51	0.49	product_vertical_gamepad, product_vertical_gamingheadset, is_mass_market, product_vertical_gamingaccessorykit, product_vertical_gamingmouse
	Multiplicative	No Yes	0.94	0.09	product_vertical_gamingheadset, is_mass_market, product_vertical_gamingmouse, product_vertical_gamepad, Online marketing_SMA_3
gamingaccessory	Koyck	No Yes	0.93	0.05	product_vertical_gamepad, product_vertical_gamingheadset, is_mass_market, product_vertical_gamingaccessorykit, product_vertical_gamingmouse
	Distributive Lag Model (Additive)	No Yes	0.87	0.1	$product_vertical_gamepad, product_vertical_gaming accessory kit, is_mass_market, product_vertical_motion controller, product_vertical_gaming keyboard$
	Distributive Lag Model (Multiplicative)	No Yes	0.93	0.11	product_vertical_gamepad, product_vertical_gamingmouse, is_mass_market, product_vertical_gamingkeyboard, is_cod
	Additive	No Yes	0.96	0.09	product_vertical_homeaudiospeaker, is_mass_market, Digital_SMA_I, product_vertical_fmradio, is_cod
homeaudio	Multiplicative	No Yes	0.86 O	0.34	product_vertical_homeaudiospeaker, is_mass_market, product_vertical_fmradio, Radio_Ad_Stock, Sponsorship
	Koyck	No Yes	0.96	0.09	product_vertical_homeaudiospeaker, is_mass_market, is_cod, NPS, Mean Temp
	Distributive Lag Model (Additive)	No Yes	0.42 0.53	1.39 0.47	product_vertical_homeaudiospeaker, product_vertical_karaokeplayer, is_mass_market, is_cod, product_vertical_fmradio
	Distributive Lag Model (Multiplicative)	No Yes	0.23 0.57	0.26	product_vertical_homeaudiospeaker, is_mass_market, product_vertical_fmradio, is_cod, product_vertical_voicerecorder

Model Selection:

- The criteria of choosing the model is based on the accuracy parameters R2 score & MSE score and the business relevance of the important attributes chosen by the model.
- Also we tried to choose models with cross validation because even though the ones without, sometimes give us good scores, they are not very dependable & generalizable, owing to limited dataset.
- By referring to the model dashboard, we finalize the following models for the 3 mentioned product subcategories Camera Accessory, Gaming Accessory and Home Audio:

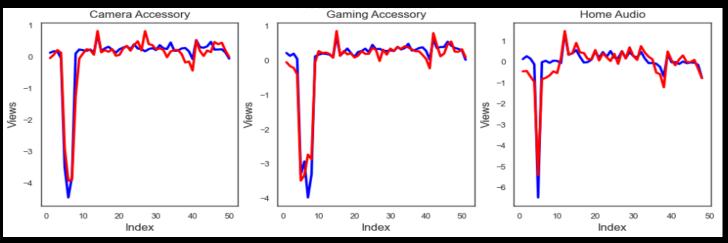
Product Sub-category	Linear Regression Model	R-square on Test Dataset	Mean Square Error	Top 5 KPIs
cameraaccessory	Multiplicative with CV	0.91	0.09	product_vertical_lens (0.181)
				product_vertical_camerabattery (0.160)
				is_mass_market (0.149)
				product_vertical_camerabatterycharger (0.121)
				TV (0.105)
gamingaccessory	Multiplicative with CV	0.94	0.06	product_vertical_gamingheadset (0.250)
				is_mass_market (0.234)
				$product_vertical_gaming mouse \ (\textbf{0.224})$
				product_vertical_gamepad (0.211)
				Online marketing_SMA_3 (0.157)
cameraaccessory	Multiplicative with CV	0.86	0.14	product_vertical_homeaudiospeaker (0.469)
				is_mass_market (0.289)
				product_vertical_fmradio (0.224)
				Radio_Ad_Stock (0.147)
				Sponsorship (0.121)

- We notice that all the 3 chosen models for the 3 subcategories are Multiplicative models.
- This fact tells us that there exists some interaction between the KPIs for all the 3 model.
- These models tell us about the growth of revenue vs the interactive growth of the KPIs

Model Validation:

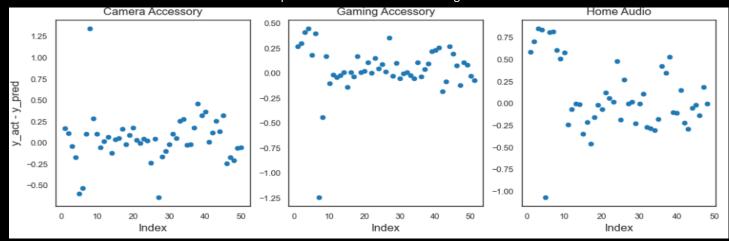
Actual vs. Predicted

Plotting the distribution of the error terms. The error terms follow a normal distribution with mean at 0 barring a few outlier values.



Error Terms

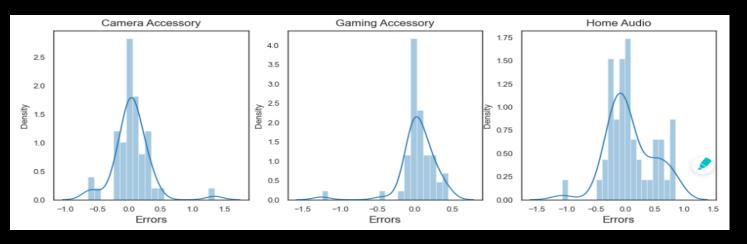
Drawing a scatter plot of the Error Terms to check the spread to ensure that the error terms have constant variance (homoscedasticity). The variance doesn't increase or decrease or follow a pattern as the error values change.



Model Validation (contd.):

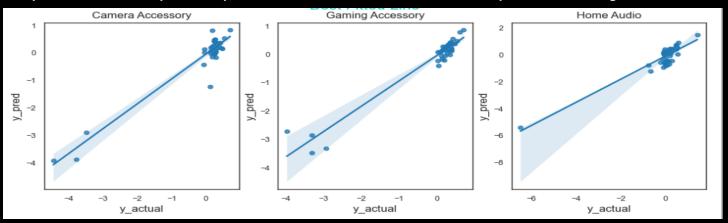
Distribution of Error Terms

Plotting the distribution of the error terms. The error terms follow a normal distribution with mean at 0 barring a few outlier values.



Best Fitted Line

Plotting a scatter plot with actual and predicted price values from the dataset to check the spread and drawing the best fitted line through it.



Equation of Best Fitted Line

Considering the top 5 KPIs from the models for our 3 product subcategories, we can see that the equation of our best fitted lines as follows:

Camera Accessory

Revenue = $0.0 + (0.181 \times product_vertical_lens) + (0.160 \times product_vertical_camerabattery) + (0.149 \times is_mass_market) + (0.121 \times product_vertical_camerabatterycharger) + (0.105 \times TV) + ...$

Gaming Accessory

Revenue = $0.0 + (0.250 \times product_vertical_gamingheadset) + (0.234 \times is_mass_market) + (0.224 \times product_vertical_gamingmouse) + (0.211 \times product_vertical_gamepad) + (0.157 \times Online marketing_SMA_3) + ...$

Home Audio

Revenue = $0.0 + (0.469 \times \text{product_vertical_homeaudiospeaker}) + (0.289 \times \text{is_mass_market}) + (0.224 \times \text{product_vertical_fmradio}) + (0.147 \times \text{Radio_Ad_Stock}) + (0.121 \times \text{Sponsorship}) + ...$

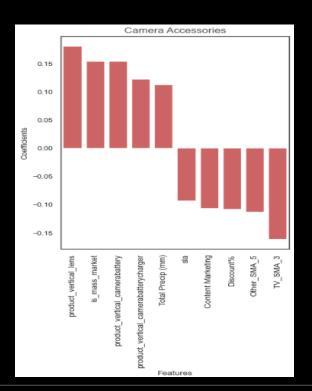
This equation implies how the revenue can grow with a unit growth in any of these independent KPIs with all other KPIs held constant

Recommendations:

Top 5 features that affect Each of the 3 Product Sub-categories (both positively and adversely) as per our Chosen Models'

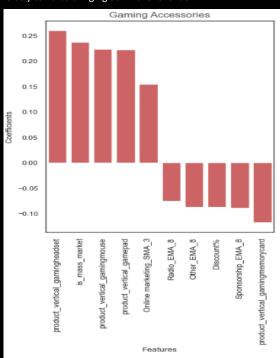
Camera Accessory

- Company should promote `Lens`, `Camera Batteries` & `Camera Battery Chargers` as they fetch the highest revenue.
- Advertisement spends on TV has a positive impact on revenue. One unit of TV spend can boost the revenue by 0.105 units. Content Marketing spends on the other hand impacts negatively.
- Mass-market` products are better contributors to the increased revenue in comparison to the Luxury products.
- Higher percentage of Discounts in general given for this sub category works adversely towards bringing down the revenue.



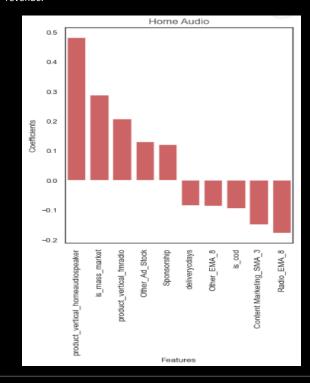
Gaming Accessory

- Company should promote `Gaming Headset`, `Gaming Mouse` & `Gamepad` as they fetch the highest revenue. On the contrary, `Gaming Memory Cards` results in loss.
- Advertisement spends on Online Marketing, Radio & Others have a positive cumulative impact on revenue. Sponsorship spends on the other hand has a negative cumulative effect.
- 'Mass-market' products are better contributors to the increased revenue in comparison to the Luxury products.
- Higher percentage of Discounts in general given for this sub category works adversely towards bringing down the revenue



Home Audio

- Company should promote 'Home Audio Speakers' & 'FM Radios as they fetch the highest revenue.
- `Mass-market` products are better contributors to the increased revenue in comparison to the Luxury products.
- Radio Ad Stock (carry over effect of Radio Advertisement) spends helps to boost the revenue to a significant extent.
- Advertisement spends on Sponsorship has a positive impact on revenue. Content Marketing spends on the other hand impacts negatively.
- COD payments in general for this sub category are bad in bringing down the revenue.



Recommendations (contd.):

General Recommendations:

- Most of the sales take place when Discount% is between 50-60%. However, that doesn't necessarily help in boosting the revenue. EDA shows that an average discount% between 10-20% is the most profitable for the company specially among luxury items.
- > In general most of the Home Audio items sold are luxury items and hence, customers prefer to use COD instead of paying upfront.
- During festive time(e.g. Thanksgiving) more investment is made on Advertisement and good promotional offers were rolled out. This usually boosts the revenue. However just providing discounts without properly advertising for it on several media channels doesn't help. We have seen that for the weeks 32 35(August), revenue generated was the lowest from all 3 product subcategories even though median discount% was raised after the initial drought. In fact, this dip in revenue can be observed as a direct relation to minimum amount of total investment in Ads during the given timeframe.

> Camera Accessories – Recommendation based on elasticity of KPI's

- > Focus on ads spend on Digital Means
- > Weekly sales have better impact
- > Procurement of products should be taken seriously

➤ Game Accessories – Recommendation based on elasticity of KPI's

- > Focus on digital ad stocks
- > As SEM investments are least effective, so further do not add revenue
- > Procurement of products should be taken seriously

➤ Home Audio – Recommendation based on elasticity of KPI's

- ➤ Need more focus on cash on delivery customers
- ➤ More advertisements through Television channels
- > Launch marketing schemes towards colder regions, it might help the sales

Challenges Faced:

- > Deciding on number of derived KPI's needed which are important as well as logical.
- > Arriving at a base dataset because many iterations were performed as the model results were leading to perfectly fitting model.
- > To solve this problem, correlation concept was used to remove correlated variables
- > Selection of significant variables which are finally decided by correlation matrix and non zero variance
- ➤ Identifying the functions/ packages to perform ad stock and creating Lag Variables
- > Limited availability of the sample implementation of Kyock, Lag and Multiplicative models in Python
- > Decision on removal of variables from analysis

Thank You!