

Retail Sales Prediction

By:

Azhar Ali

Mohd Taufique

Pushpam Raghuvanshi

AlmaBetter, Bangalore

Abstract:

Sales forecasting refers to the process of estimating demand for or sales of a particular product over a specific period of time.

Businesses use sales forecasts to determine what revenue they will be generating in a particular timespan to empower themselves with powerful and strategic business plans. Important decisions such as budgets, hiring, incentives, goals, acquisitions and various other growth plans are affected by the revenue the company is going to make in the coming months and for these plans to be as effective as they are planned to be it is important for these forecasts to also be as good.

The sales forecasts are also different from the sales-goals a company has. Sales-goals is what a company wants to happen to execute their future plans for the business. On the other hand sales forecasts are what is going to happen on the basis of past records, data, trends and various improvement measures taken.

The work here predicts the sales for a drug store chain in the European market for a time period of six weeks and compares the results of machine learning algorithms.

Problem Statement:

Rossmann operates over 3,000 drug stores in 7

European countries. Currently, Rossmann store managers are tasked with predicting their daily sales for up to six weeks in advance. Store sales are influenced by many factors, including promotions, competition, school and state holidays, seasonality, and locality. With thousands of individual managers predicting sales based on their unique circumstances, the accuracy of results can be quite varied.

You are provided with historical sales data for 1,115 Rossmann stores. The task is to forecast the "Sales" column for the test set. Note that some stores in the dataset were temporarily closed for refurbishment.

Introduction:

The interest for a product continues to change occasionally. No business can work on its monetary growth without assessing client interest and future demand of items precisely. Sales forecasting refers to the process of estimating demand for or sales of a particular product over a specific period of time.

For a good sales forecast, it is extremely important to get a good dataset as well. Forecasts heavily depend on the past records, trends and patterns observed for sales of a particular store. The variations could be due to a number of reason

In this Retail Sales Prediction, machine learning models are created that predict sales of these 1115 drug stores across the European market and compare the results of these models. In addition to this, an effort has been made to analyze and find all the features that are contributing to higher sales and the features which are leading to lower sales, so that improvement plans can be worked upon.

Approaches:

The approach followed here is to first check the sanctity of the data and then understand the features involved. The events followed were in our approach:

- **Understanding the business problem and the datasets**
- **Data cleaning and preprocessing-** finding null values and imputing them with appropriate values.
Converting categorical values into appropriate data types and merging the datasets provided to get a final dataset to work upon.
- **Exploratory data analysis-**of categorical and continuous variables against our target variable.
- **Data manipulation-**feature selection and engineering, feature scaling, outlier detection and treatment and encoding categorical features.
- **Modeling-** The baseline model- Linear Regression was chosen and other models like Ridge, Lasso, Decision Tree, Random forest implemented to maximize the performance. Hyperparameter & Cross Validation (CV) is also used to get the best model along with the best parameters.

- **Model Performance and Evaluation**
- **Store wise Sales Predictions**
- **Conclusion and Recommendations**

Understanding the Data:

First step involved is understanding the data and getting answers to some basic questions like; What is the data about? How many rows or observations are there in it? How many features are there in it? What are the data types? Are there any missing values? And anything that could be relevant and useful to our investigation. Let's just understand the dataset first and the features involved before proceeding further.

Our dataset consists of two csv files, the first consists of historical data with 1017209 rows or observations and 9 columns with no null values. The second dataset was supplementary information about the stores with 1115 rows and 10 columns and a lot of missing values in a few columns. The data types were of integer, float and object in nature.

Let's define the features involved:

- **Id** -an Id that represents a (Store, Date) tuple within the set
- **Store** -a unique Id for each store
- **Sales** -the turnover for any given day (Dependent Variable)
- **Customers** -the number of customers on a given day
- **Open** -an indicator for whether the store was open: 0 = closed, 1 = open
- **StateHoliday** -indicates a state holiday. Normally all stores, with few exceptions, are closed on state holidays. Note that all schools are closed on public holidays and weekends. a = public holiday, b = Easter holiday, c = Christmas, 0 = None

- **SchoolHoliday** -indicates if the (Store, Date) was affected by the closure of public schools
- **StoreType** -differentiates between 4 different store models: a, b, c, d
- **Assortment**- describes an assortment level: a = basic, b = extra, c = extended. An assortment strategy in retailing involves the number and type of products that stores display for purchase by consumers.
- **CompetitionDistance**- distance in meters to the nearest competitor store
- **CompetitionOpenSince[Month/Year]**- gives the approximate year and month of the time the nearest competitor was opened
- **Promo**- indicates whether a store is running a promo on that day
- **Promo2**- Promo2 is a continuing and consecutive promotion for some stores: 0 = store is not participating, 1 = store is participating
- **Promo2Since[Year/Week]** -describes the year and calendar week when the store started participating in Promo2
- **PromoInterval**- describes the consecutive intervals Promo2 is started, naming the months the promotion is started anew. E.g. "Feb,May,Aug,Nov" means each round starts in February, May, August, November of any given year for that store.

Data Cleaning and Preprocessing:

Handling missing values is an important skill in the data analysis process. If there are very few missing values compared to the size of the dataset, we may choose to drop rows that have

missing values. Otherwise, it is better to replace them with appropriate values.

It is necessary to check and handle these values before feeding it to the models, so as to obtain good insights on what the data is trying to say and make great characterisation and predictions which will in turn help improve the business's growth.

The historical records dataset had no null values.

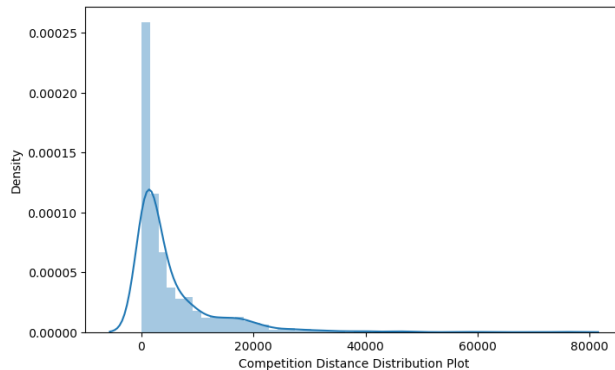
The dataset had a lot of nulls in the following columns:

- CompetitionOpenSinceMonth
- CompetitionOpenSinceYear
- Promo2SinceWeek
- Promo2SinceYear
- PromoInterval

Store	0
StoreType	0
Assortment	0
CompetitionDistance	3
CompetitionOpenSinceMonth	354
CompetitionOpenSinceYear	354
Promo2	0
Promo2SinceWeek	544
Promo2SinceYear	544
PromoInterval	544

Store	0
DayOfWeek	0
Sales	0
Customers	0
Open	0
Promo	0
StateHoliday	0
SchoolHoliday	0
dtype:	int64

- ‘CompetitionDistance’ - Competition Distance is the distance in meters to the nearest competitor store.
The Competition Distance distribution plot shows the distances at which generally the stores are opened



It seems like most of the values of the CompetitionDistance are towards the left and the distribution is skewed on the right. Median is more robust to outlier effect hence median was imputed in the null values.

Right skewed distributions occur when the long tail is on the right side of the distribution also called as positive skewed distribution which essentially suggests that there are positive outliers far along which influences the mean. It seems like most of the values of the CompetitionDistance in the column are in the lower range. Consequently, the longer tail in an asymmetrical distribution pulls the mean away from the most common values. The mean is greater than the median. The mean overestimates the most common values in the distribution and hence median is used in this case, it is more robust to outlier effect and hence median is used to impute the missing values in this feature.

- CompetitionOpenSinceMonth- gives the

approximate month of the time the nearest competitor was opened. The mode of the column is used to impute the missing values in the column as it gives the most occurring month.

- CompetitionOpenSinceYear-gives the approximate year of the time the nearest competitor was opened. The mode of the column is used to impute the missing values in the column as it gives the most occurring month.
- Promo2SinceWeek, Promo2SinceYear and PromoInterval are NaN wherever Promo2 is 0 or False as can be seen in the first look of the dataset. They are replaced with 0.

Lastly before proceeding further, the two datasets were merged on the common column of ‘Store’ to get everything together for the analysis.

Exploratory Data Analysis:

Exploratory data analysis is a crucial part of data analysis. It involves exploring and analyzing the dataset given to find out patterns, trends and conclusions to make better decisions related to the data, often using statistical graphics and other data visualization tools to summarize the results. The visualization tools involved in the investigation are python libraries- matplotlib and seaborn.

The goal here is to explore the relationships of different variables with ‘Sales’ to see what factors might be contributing to the high and low sales numbers.

Feature Exploration:

There are two kinds of features in the dataset: Categorical and Non Categorical Variables.

Categorical- A categorical variable is a variable that can take on one of a limited, and usually

fixed, number of possible values putting a particular category to the observation.

Non Categorical/Continuous- A non categorical or continuous variable is a variable whose value is obtained by measuring, i.e., one which can take on an uncountable set of values.

Both of them are analyzed separately. Categorical data is usually analyzed through count plots and barplots in accordance with the target variable and that is what is done here too. On the other hand Numeric or Continuous variables were analyzed through distribution plots, box plots and scatterplots to get useful insights.

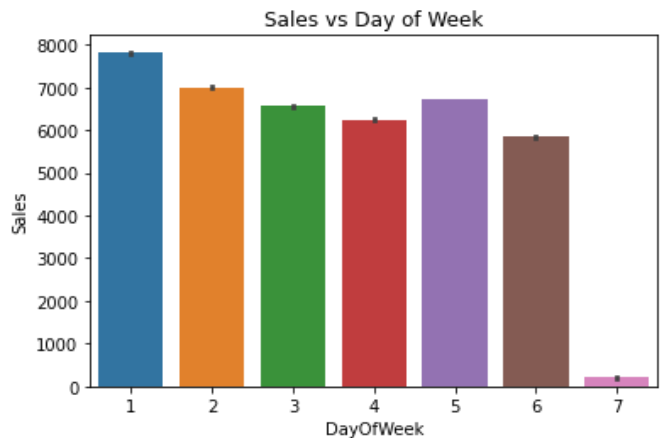
Hypotheses:

Just by observing the head of the dataset and understanding the features involved in it, the following hypotheses could be framed:

- There's a feature called "DayOfWeek" with the values 1-7 denoting each day of the week. There would be a week off probably Sunday when the stores would be closed and we would get low overall sales.
- Customers would have a positive correlation with Sales.
- The Store type and Assortment strategy involved would be having a certain effect on sales as well. Some premium high quality products would fetch more revenue.
- Promotion should be having a positive correlation with Sales.
- Some stores were closed due to refurbishment, those would generate 0 revenue for that time period.
- Stores are influenced by seasonality, probably before holidays sales would be high.

Categorical Columns Insights:

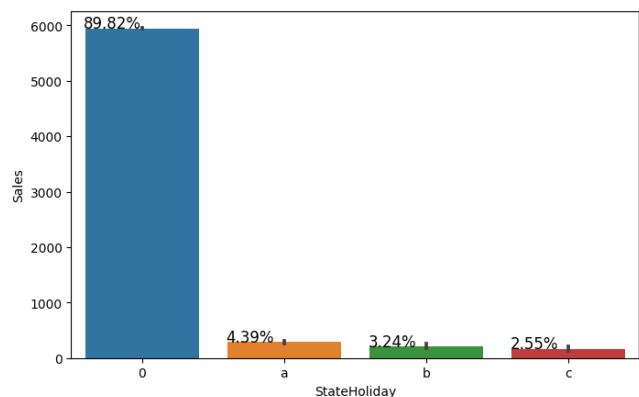
1. Total Sales on Weekdays



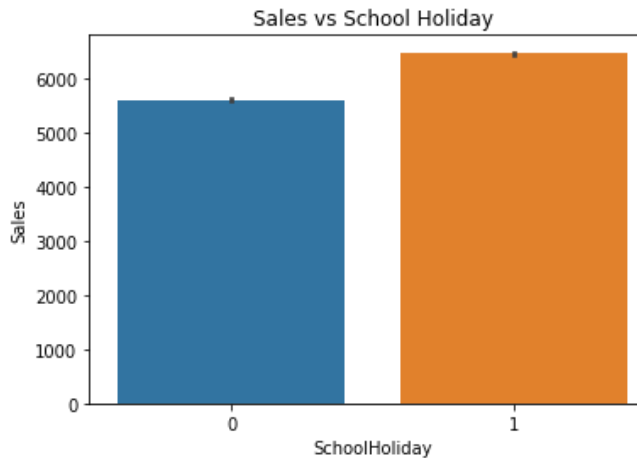
Here it can be deduced that there were more sales on Monday, probably because shops generally remain closed on Sundays which had the lowest sales in a week.

2. Sales VS State Holiday

Sales were low whenever there was a State Holiday indicating only a few stores were open on these days.



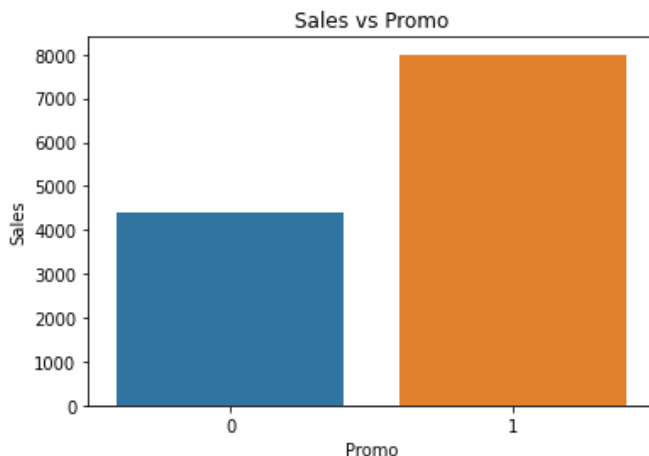
3. Sales vs School Holiday



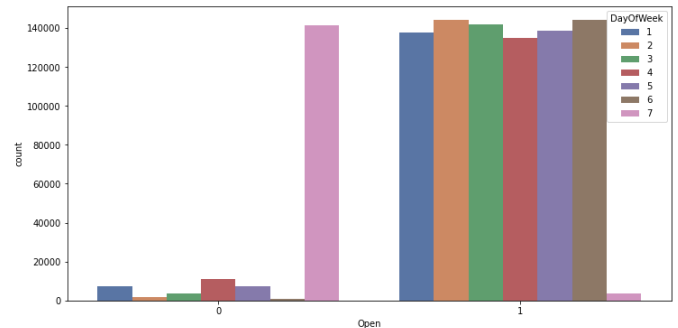
Sales were high on an average on School Holidays indicating School Holidays weren't compulsory by the law and comparatively more sales were recorded on holidays.

4. Sales VS Promo

Promotion has a positive effect on Sales indicating high sales for stores with Promo=1. This can be clearly seen in the figure below.

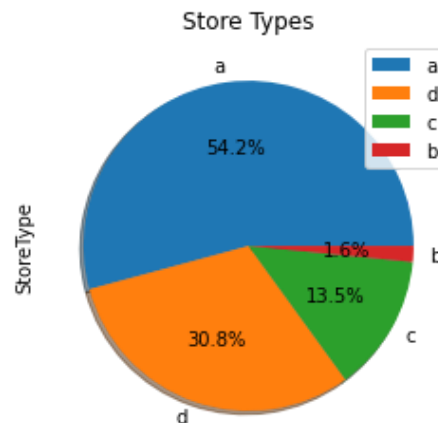


5. Store open/closed over day of week



This is a count plot of open shops according to the day of the week. It's clear that the number of shops open on Sundays were very less and hence low sales. Some shops were closed on weekdays as well accounting to the stores closed due to refurbishment or holidays.

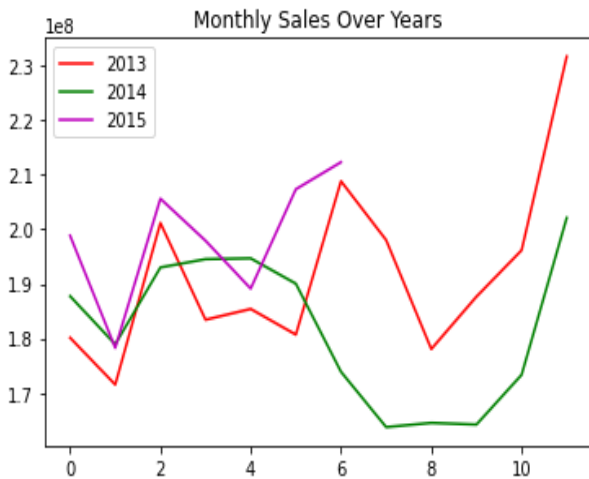
6. Proportion store Type



Through the pie chart for the various store types, it can be clearly observed that even though type a stores had the most proportion, type b stores were very low in quantity.

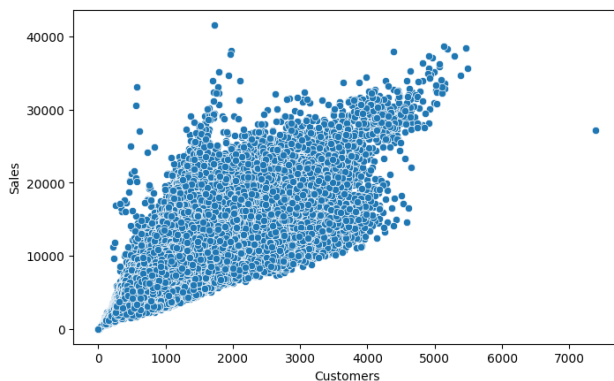
Continuous Features Insights:

1. Average Sales Over Year/Month



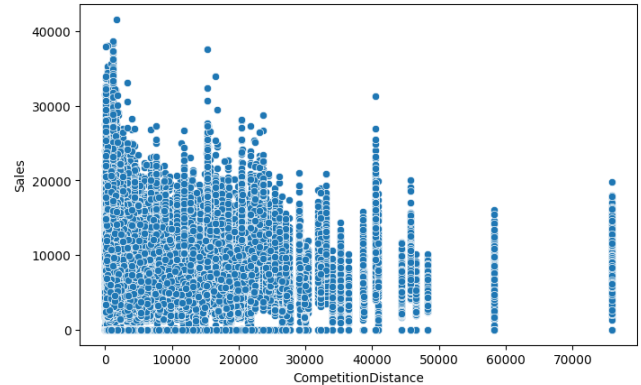
Here's a plot of Monthly Sales over the years. Sales rise up by the end of the year before the holidays. Sales for 2014 went down there for a couple months - July to September, indicating stores closed due to refurbishment.

2. Scatterplot - Sales and Customer



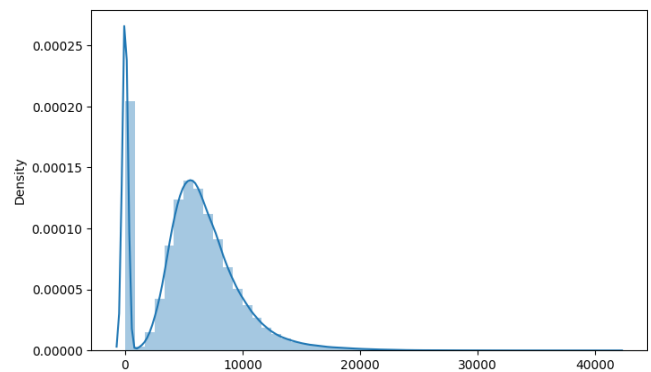
Sales and Customer scatter plot showed a direct positive relation between them with a few outliers.

3. Scatterplot - Sales and Competition Distance



From the above scatter plot it can be observed that mostly the competitor stores weren't that far from each other and the stores densely located near each other saw more sales. This could indicate competition between busy locations vs remote locations.

4. Sales Distribution Plot



Here's a distribution plot of the Sales column. The drop in sales indicates the 0 sales accounting to the stores temporarily closed due to refurbishment.

Correlation:

Correlation is a statistical term used to measure the degree in which two variables move in relation to each other. A perfect positive correlation means that the correlation coefficient is exactly 1. This implies that as one variable moves, either up or down, the other moves in the same direction. A perfect negative correlation means that two variables move in opposite directions, while a zero correlation implies no linear relationship at all.

Store	1	8.5e-06	0.0051	0.024	4.7e-05	5.8e-05	0.00064	0.00029	0.0015	2.3e-05	0.0014	-0.026	-0.037	-0.00022	0.0085	0.011	0.0085
DayOfWeek	-8.5e-06	1	-0.46	-0.39	-0.53	-0.39	-0.21	0.0019	-0.0054	0.0051	-0.0049	-2.5e-05	2.6e-06	-2.2e-05	0.00017	0.00021	0.00017
Sales	-0.0051	-0.46	1	0.89	0.68	0.45	0.085	0.024	0.049	-0.012	0.053	-0.019	0.023	0.0087	-0.091	-0.044	-0.091
Customers	0.024	-0.39	0.89	1	0.62	0.32	0.072	-0.0012	0.038	-0.0045	0.041	-0.1	-0.024	0.0065	-0.15	-0.098	-0.15
Open	-4.7e-05	-0.53	0.68	0.62	1	0.3	0.086	-0.001	-0.0068	0.033	0.0048	0.008	0.0014	0.0026	-0.0083	-0.0074	-0.0083
Promo	-5.8e-05	-0.39	0.45	0.32	0.3	1	0.067	0.024	-0.012	-0.11	0.00057	0.00014	-7.7e-06	0.00013	-0.00098	-0.0012	-0.00098
SchoolHoliday	-0.00064	-0.21	0.085	0.072	0.086	0.067	1	-0.037	0.1	0.031	0.071	-0.0037	0.00012	0.0018	-0.0069	-0.0067	-0.0069
Year	-0.00029	0.0019	0.024	-0.0012	-0.001	0.024	-0.037	1	-0.27	-0.0025	-0.26	0.00071	-4.2e-05	0.00066	-0.005	-0.0061	-0.005
Month	-0.0015	-0.0054	0.049	0.038	0.0068	-0.012	0.1	-0.27	1	0.012	0.012	0.0036	-0.00022	0.0033	-0.025	-0.031	-0.025
Day	-2.3e-05	0.0051	-0.012	-0.0045	0.033	-0.11	0.031	-0.0025	0.012	1	0.07	4.9e-05	-1.2e-06	4.7e-05	0.00035	-0.00044	0.00035
WeekOfYear	-0.0014	-0.0049	0.053	0.041	0.0048	0.00057	0.071	-0.26	0.012	0.07	1	0.0035	-0.00021	0.0032	-0.025	-0.03	-0.025
CompetitionDistance	-0.026	-2.5e-05	-0.019	-0.1	0.008	0.00014	-0.0037	0.00071	0.0036	4.9e-05	0.0035	1	-0.049	0.02	-0.14	-0.12	-0.14
CompetitionOpenSinceMonth	-0.037	2.6e-06	-0.023	-0.024	0.0014	-7.7e-06	0.00012	4.2e-05	-0.00022	-1.2e-06	-0.00021	-0.049	1	0.058	0.022	0.02	0.022
CompetitionOpenSinceYear	-0.00022	-2.2e-05	0.0087	0.0065	0.0026	0.00013	0.0018	0.00066	0.0033	4.7e-05	0.0032	0.02	0.058	1	-0.023	-0.027	-0.023
Promo2	-0.0085	0.00017	-0.091	-0.15	-0.0083	-0.00098	-0.0069	-0.005	-0.025	-0.00035	-0.025	-0.14	0.022	-0.023	1	0.76	1
Promo2SinceWeek	-0.011	0.00021	-0.044	-0.098	-0.0074	-0.0012	-0.0067	-0.0061	-0.031	-0.00044	-0.03	-0.12	0.02	-0.027	0.76	1	0.76
Promo2SinceYear	-0.0085	0.00017	-0.091	-0.15	-0.0083	-0.00098	-0.0069	-0.005	-0.025	-0.00035	-0.025	-0.14	0.022	-0.023	1	0.76	1

- Day of the week has a negative correlation indicating low sales as the weekends, and promo, customers and open has positive correlation.
- State Holiday has a negative correlation suggesting that stores are mostly closed on state holidays indicating low sales.
- CompetitionDistance showing negative correlation suggests that as the distance increases sales reduce, which was also observed through the scatterplot earlier.
- There's multicollinearity involved in the dataset as well. The features telling the same story like Promo2, Promo2 since week and year are showing multicollinearity.

Removing Multicollinearity:

- Removing the features which are having $VIF > 10$ because it will affect & interpret the result.
- $VIF \leq 10$ is usually preferred as this can easily explain the variance of 90% i.e, R-square becomes 90%. ($VIF = 1/(1-R^2)$).

	variables	VIF
0	Store	3.627575
1	DayOfWeek	4.513547
2	Customers	4.339790
3	Promo	1.946273
4	StateHoliday	1.003985
5	SchoolHoliday	1.247951
6	Day	3.847661
7	StoreType	1.916142
8	Assortment	2.049916
9	CompetitionDistance	1.532510
10	CompetitionOpenSinceMonth	6.574554
11	Promo2	4.752803
12	Promo2SinceWeek	3.737566

In this table we can see that all features are having $VIF < 10$. That looks pretty fine.

Data Manipulation:

Data manipulation involves manipulating and changing our dataset before feeding it to various regression machine learning models. This involves keeping important features, outlier treatment, feature scaling and creating dummy variables if necessary.

Feature Engineering:

- Some stores were closed due to refurbishment and some on account of week off or holidays. Those stores on those dates generated zero sales and hence removing the rows was important to avoid confusion by the algorithms and then removing the feature altogether because it wasn't providing any value in prediction of the sales.
- There were features that like Competition Open since Month and Year. It was combined to count the total months since the nearest competition was opened.
- Promo2SinceWeek, Promo2SinceYear indicated promotion 2 opened since week and year. These features were combined to count the total months since promotion 2 is run.
- PromoInterval has values [0, 'Feb, May, Aug, Nov', 'Jan, Apr, Jul, Oct', 'Mar, Jun, Sept, Dec'] Changing these to dummy variables.
- StoreType and Assortment have values [a,b,c,d] and [a,b,c] respectively. Changing these to Numerical values [0,1,2,3] and [0,1,2] respectively.

Outlier Detection:

In statistics, an outlier is a data point that differs significantly from other observations. Outliers can occur by chance in any distribution, but they often indicate either measurement error or that the population has a heavy-tailed distribution.

Z-score is a statistical measure that tells you how far a data point is from the rest of the dataset. In a more technical term, Z-score tells how many standard deviations away a given observation is from the mean.

$$z = (x - \text{mean}) / \text{standard deviation}$$

More than 3 standard deviations was considered as an outlier. Exploring the outliers dataframe, some important insights were generated:

- The data points with sales value higher than 28000 are very low and hence they can be considered as outliers.
- The outliers had day of the week as 7 i.e. Sunday and the store type for those observations were 'b'.
- Other outliers had promotion running on that day.
- Being open 24*7 along with all kinds of assortments available is probably the reason why it had higher average sales than any other store type.
- It can be well established that the outliers are showing this behavior for the stores with promotion = 1 and store type B. It would not be wise to treat them because the reasons behind this behavior seem fair.

- If the outliers are a valid occurrence it would be wise not to treat them by deleting or manipulating them especially when we have established the ups and downs of the target variable in relation to the other features. It is well established that there is seasonality involved and no linear relationship is possible to fit. For these kinds of dataset tree based machine learning algorithms are used which are robust to outlier effect.

Feature Scaling:

Feature Scaling is a technique to standardize the independent features present in the data in a fixed range. It is done to prevent biased nature of machine learning algorithms towards features with greater values and scale. The two techniques are:

Normalization: is a scaling technique in which values are shifted and rescaled so that they end up ranging between 0 and 1. It is also known as Min-Max scaling. [0,1]

$$X' = \frac{X - X_{min}}{X_{max} - X_{min}}$$

Standardization: is another scaling technique where the values are centered around the mean with a unit standard deviation. This means that the mean of the attribute becomes zero and the resultant distribution has a unit standard deviation. [-1,1]

$$X' = \frac{X - \mu}{\sigma}$$

Normalization of the continuous variables was done further.

One hot encoding:

For categorical variables where no such ordinal relationship exists, the integer encoding is not enough. We have categorical data integers encoded with us, but assuming a natural order and allowing this data to the model may result in

poor performance.

Many of the features such as DayofWeek, StoreType and Assortments were categorical in nature and had to be one hot encoded to continue.

Modeling:

Factors affecting in choosing the model:

Determining which algorithm to use depends on many factors like the problem statement and the kind of output you want, type and size of the data, the available computational time, number of features, and observations in the data, to name a few.

Train-Test Split:

In machine learning, train/test split splits the data randomly, as there's no dependence from one observation to the other. That's not the case with time series data. Here, it's important to use values at the rear of the dataset for testing and everything else for training.

The 20% of the datasets were kept as a testing set and the rest of the historical data i.e. 80% were used in the training set.

Models Implemented:

Linear Regression:

A baseline is a simple model that provides reasonable results on a task and does not require much expertise and time to build. It is the most easy to interpret. Regression models a target prediction value based on independent variables. It is mostly used for finding out the relationship between variables and forecasting. Different regression models differ based on – the kind of relationship between dependent and independent variables they are considering, and the number of independent variables getting used.

Regularization:

- Regularization is one of the ways to improve our model to work on unseen data by ignoring the less important features.
- It avoids overfitting by adding a penalty to the model with high variance, thereby shrinking the beta coefficients to zero.

Lasso Regularization (L1)

- It stands for Least Absolute Shrinkage and Selection Operator
- It adds L1 the penalty
- L1 is the sum of the absolute value of the beta coefficients

Ridge Regularization (L2)

- It adds L2 as the penalty
- L2 is the sum of the square of the magnitude of beta coefficients

Decision Tree:

Decision Tree is a Supervised learning technique that can be used for both Classification and Regression problems. It is a tree-structured classifier, where internal nodes represent the features of a dataset, branches represent the decision rules and each leaf node represents the outcome.

In a Decision tree, there are two nodes, which are the Decision Node and Leaf Node. Decision nodes are used to make any decision and have multiple branches, whereas Leaf nodes are the output of those decisions and do not contain any further branches.

Random Forest:

Random forests are an ensemble learning method for classification and regression that operates by constructing a multitude of decision trees at training time. For regression tasks, the output of the random forest is the average of the results given by most trees.

In simple terms, random forest builds multiple decision trees and merges them together to get a more accurate and stable prediction.

Random Forest Regressor results were much better than our baseline model with a test R^2 of 0.8686.

Random Forest Hyperparameters:

- `max_depth`- The `max_depth` of a tree in Random Forest is defined as the longest path between the root node and the leaf node
- `min_sample_split`- a parameter that tells the decision tree in a random forest the minimum required number of observations in any given node in order to split it. The default value of the `min_sample_split` is assigned to 2. This means that if any terminal node has more than two observations and is not a pure node, we can split it further into subnodes.
- `max_leaf_nodes`- This hyperparameter sets a condition on the splitting of the nodes in the tree and hence restricts the growth of the tree. If after splitting we have more terminal nodes than the specified number of terminal nodes, it will stop the splitting and the tree will not grow further.

- **min_samples_leaf**- This Random Forest hyperparameter specifies the minimum number of samples that should be present in the leaf node after splitting a node.
- **n_estimators**- the number of tree
- **max_sample** (bootstrap sample)-The max_samples hyperparameter determines what fraction of the original dataset is given to any individual tree.
- **max_features**- This resembles the number of maximum features provided to each tree in a random forest.

Grid search cv searches on hyper parameters to fit and score various models and get the best estimator.

Random Forest Hyperparameter Tuned Model :

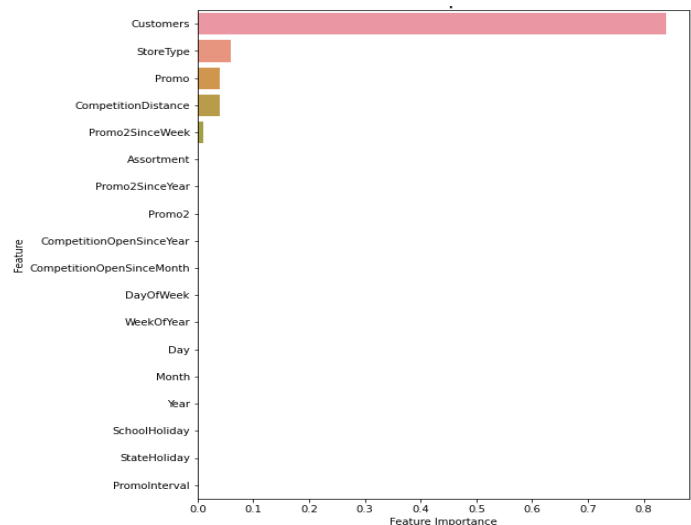
The maximum R^2 was seen in the Random Forest with hyperparameter and Cross validation (CV) tuned models with the value 0.9727. This indicates that all the trends and patterns that could be captured by these models without overfitting were done and the maximum level of performance achievable by the model was achieved.

Actual vs Predicted Values for Random Forest Hyperparameter & Cross-validation (CV) Tuned Model:

The following data frame or table shows the actual and predicted value by our best model which is a tuned random forest model using hyperparameters and cross validation.

	actual	predicted
0	6792	6277.666667
1	11585	11038.933333
2	11843	11386.666667
3	11961	10756.933333
4	4657	4456.000000
...
168863	9680	9094.066667
168864	4252	4434.333333
168865	2581	2466.466667
168866	3757	3371.000000
168867	3540	3652.933333
168868 rows × 2 columns		

Feature Importance:



The most important features in predicting the Sales were Customers, CompetitionDistance, StoreType and Promo. The value shown above is rounded up to 2 digit or decimal places only.

Model Performance and Evaluation:

ML Model predict sales for stores which are open and when there is some sales because there is no sales when store is closed.

	Linear Regression (OLS)	Lasso Regression (L1)	Ridge Regression (L2)	Decision Tree	Random Forest	Random Forest with hyperparameter tuning & CV
RMSE	1519.132	1520.620	1519.849	1433.262	1126.660	<u>513.893</u>
MAPE	15.902	15.920	15.908	15.478	12.442	<u>4.993</u>
R2	0.761	0.7607	0.7609	0.7609	0.8686	<u>0.9727</u>

Evaluation Metrics:

- Mean Absolute Error(MAE)- MAE is a very simple metric which calculates the mean of absolute difference between actual and predicted values.
- Mean Squared Error(MSE)- Mean squared error states the mean of the squared difference between actual and predicted value.
- Root Mean Squared Error(RMSE)- It is a simple square root of mean squared error.
- R Squared (R^2)- R2 score is a metric that tells the performance of your model, not the loss in an absolute sense that how well did your model perform. Hence, R2 squared is also known as Coefficient of Determination or sometimes also known as Goodness of fit. Its value ranges from 0 to 1. It can be negative if the model is performing worse than the base.
- Adjusted R Squared- The disadvantage of the R2 score is while adding new features in data the R2 score starts increasing or remains constant but it

never decreases because It assumes that while adding more data variance of data increases. Adjusted R^2 is adjusted for this disadvantage and shows the real value.

Conclusion :

The main objective of sales forecasting is to paint an accurate picture of expected sales. Sales teams aim to either hit their expected target or exceed it.

When the sales forecast is accurate, operations go smoothly and future planning for the company's growth is done efficiently.

Upon having this analysis it can be established that given the dataset, the model developed is able to explain 97.27 % of the variations and is able to predict the sales values in a good range.

Some important insights to draw from the analysis includes:

- There were more sales on Monday, probably because shops generally remain closed on Sundays which had the lowest sales in a week.
- The positive effect of promotion on Customers and Sales is observable.
- Most stores have competition distance within the lower range and had more sales than stores far away, probably indicating competition in busy locations vs remote locations.
- Store type 'b' though being few in number had the highest sales average. The reasons include all three kinds of assortments specially assortment level b which is only available at type b stores and being open on Sundays as well.

- The outliers in the dataset showed justifiable behavior. The outliers were either of store type b or had promotion going on which increased sales.
- Random Forest Tuned Model gave the best results. This indicates that all the trends and patterns that could be captured by these models without overfitting were done and the maximum level of performance achievable by the model was achieved.
- Approx. 50% stores are of type 'a'. There are very few stores of type 'b'.
- Store type 'b' has the highest sales and all other store types 'a', 'c', 'd' have nearly equal sales.
- December records the highest monthly sales. This may be due to Christmas and New Year.
- Sales is more when promos/offers are running on stores.

Recommendations:

- More stores should be encouraged for promotion as promo is directly correlated with Sales columns (Target Variable)
- Store type 'b' should be increased in number.
- There is seasonality involved. Hence, the stores should be encouraged to promote and take advantage of the holidays.

Challenges:

- The major challenge involved is the computational time and RAM needed to work upon such a large dataset.
- Understanding the meaning of a few Features/Columns.

- Handling Large amounts of Sales Data.
- Dealing with Null values, as there are many Null values.
- Understanding the business model of Retail Sales that how they work.
- Dealing with Categorical columns to make them numerical for use in ML model building.
- Also, forming different graphs to show insights from the dataset and to summarize the information and communicate the results and trends to the reader successfully.

References:

- AlmaBetter Resources
- Kaggle
- Analytics Vidhya Blogs
- Towards Data Science Blogs
- Github
- Scikit- Learn Org