

Big Mart Sales: Practice Problem

Introduction

This is an extensive exploratory tutorial on the Big Mart Sales challenge. It's a regression practice problem wherein we have to predict sales product-wise and store-wise. This tutorial is intended for the beginners who want to learn how to solve a regression problem in R. We will follow the following table of content.

1. Problem Statement
2. Hypothesis Generation
3. Data Structure and Content
4. Exploratory Data Analysis
 - Univariate Analysis
 - Bivariate Analysis
5. Missing Value Treatment
6. Feature Engineering
7. Encoding Categorical Variables
 - Label Encoding
 - One Hot Encoding
8. PreProcessing Data
9. Modeling
 - Linear Regression
 - Lasso Regression
 - Ridge Regression
 - RandomForest
 - XGBoost
10. End Notes

Problem Statement

The data scientists at BigMart have collected 2013 sales data for 1559 products across 10 stores in different cities. Also, certain attributes of each product and store have been defined. The aim is to build a predictive model and find out the sales of each product at a particular store.

Using this model, BigMart will try to understand the properties of products and stores which play a key role in increasing sales.

Hypothesis Generation

This is a very important stage in any machine learning process. It basically involves brainstorming and coming up with as many ideas as possible about what could affect the target variable. It helps us in exploring the data at hand more efficiently and effectively. Hypothesis Generation should be done before seeing the data or else we will end up with biased hypotheses. Following are some of the hypotheses based on the problem statement.

1. Sales are higher during weekends.
2. Higher sales during morning and late evening.
3. Higher sales during end of the year.
4. Store size affects the sales.
5. Location of the store affects the sales.
6. Items with more shelf space sell more.

You can come up with more hypotheses of your own, the more the better. Let's begin exploring the dataset and try to find interesting patterns.

We'll first load the required packages.

[Hide](#)

```
library(data.table)
library(dplyr)
library(ggplot2)
library(caret)
library(corrplot)
library(e1071)
library(xgboost)
library(cowplot)
```

We use data.table's fread() to speed up reading in the datasets.

[Hide](#)

```
train = fread("Train_UWu5bXk.csv")
test = fread("Test_u94Q5KV.csv")
submission = fread("SampleSubmission_Tmn039y.csv")
```

Data Structure and Content

It's good to quickly check the dimensions of our data.

[Hide](#)

```
dim(train);dim(test)
```

```
[1] 8523  12
[1] 5681  11
```

train dataset has 8523 rows and 12 features and test has 5681 rows and 11 columns. train has 1 extra column which is the target variable. We will predict this target variable for the test dataset later in this tutorial.

Feature names of train and test datasets

[Hide](#)

```
names(train)
```

[1] "Item_Identifier"	"Item_Weight"	"Item_Fat_Content"	"Item_Visibility"
[5] "Item_Type"	"Item_MRP"	"Outlet_Identifier"	"Outlet_Establishment_Year"
[9] "Outlet_Size"	"Outlet_Location_Type"	"Outlet_Type"	"Item_Outlet_Sales"

Hide

```
names(test)
```

[1] "Item_Identifier"	"Item_Weight"	"Item_Fat_Content"	"Item_Visibility"
[5] "Item_Type"	"Item_MRP"	"Outlet_Identifier"	"Outlet_Establishment_Year"
[9] "Outlet_Size"	"Outlet_Location_Type"	"Outlet_Type"	

Item_Outlet_Sales is present in train but not in test dataset because this is the target variable that we have to predict.

Structure of train and test datasets

Hide

```
str(train)
```

```
Classes 'data.table' and 'data.frame': 8523 obs. of 12 variables:
 $ Item_Identifier      : chr  "FDA15" "DRC01" "FDN15" "FDX07" ...
 $ Item_Weight          : num  9.3 5.92 17.5 19.2 8.93 ...
 $ Item_Fat_Content     : chr  "Low Fat" "Regular" "Low Fat" "Regular" ...
 $ Item_Visibility      : num  0.016 0.0193 0.0168 0 0 ...
 $ Item_Type           : chr  "Dairy" "Soft Drinks" "Meat" "Fruits and Vegetables" ...
 $ Item_MRP            : num  249.8 48.3 141.6 182.1 53.9 ...
 $ Outlet_Identifier    : chr  "OUT049" "OUT018" "OUT049" "OUT010" ...
 $ Outlet_Establishment_Year: int  1999 2009 1999 1998 1987 2009 1987 1985 2002 2007 ...
 $ Outlet_Size         : chr  "Medium" "Medium" "Medium" "" ...
 $ Outlet_Location_Type : chr  "Tier 1" "Tier 3" "Tier 1" "Tier 3" ...
 $ Outlet_Type         : chr  "Supermarket Type1" "Supermarket Type2" "Supermarket Type1"
 "Grocery Store" ...
 $ Item_Outlet_Sales    : num  3735 443 2097 732 995 ...
 - attr(*, ".internal.selfref")=<externalptr>
```

Hide

```
str(test)
```

```
Classes 'data.table' and 'data.frame': 5681 obs. of 11 variables:
 $ Item_Identifier      : chr  "FDW58" "FDW14" "NCN55" "FDQ58" ...
 $ Item_Weight         : num  20.75 8.3 14.6 7.32 NA ...
 $ Item_Fat_Content     : chr  "Low Fat" "reg" "Low Fat" "Low Fat" ...
 $ Item_Visibility     : num  0.00756 0.03843 0.09957 0.01539 0.1186 ...
 $ Item_Type           : chr  "Snack Foods" "Dairy" "Others" "Snack Foods" ...
 $ Item_MRP            : num  107.9 87.3 241.8 155 234.2 ...
 $ Outlet_Identifier    : chr  "OUT049" "OUT017" "OUT010" "OUT017" ...
 $ Outlet_Establishment_Year: int  1999 2007 1998 2007 1985 1997 2009 1985 2002 2007 ...
 $ Outlet_Size         : chr  "Medium" "" "" "" ...
 $ Outlet_Location_Type : chr  "Tier 1" "Tier 2" "Tier 3" "Tier 2" ...
 $ Outlet_Type         : chr  "Supermarket Type1" "Supermarket Type1" "Grocery Store" "Supermarket Type1" ...
 - attr(*, ".internal.selfref")=externalptr>
```

To explore data in any data science competition, it is advisable to append test data to the train data. So, we will combine both train and test to carry out data visualization, feature engineering, one-hot encoding, and label encoding. Later we would split this combined data back to train and test datasets.

[Hide](#)

```
test[,Item_Outlet_Sales := NA]
combi = rbind(train, test) # combining train and test datasets
dim(combi)
```

```
[1] 14204    12
```

Exploratory Data Analysis

Univariate Analysis

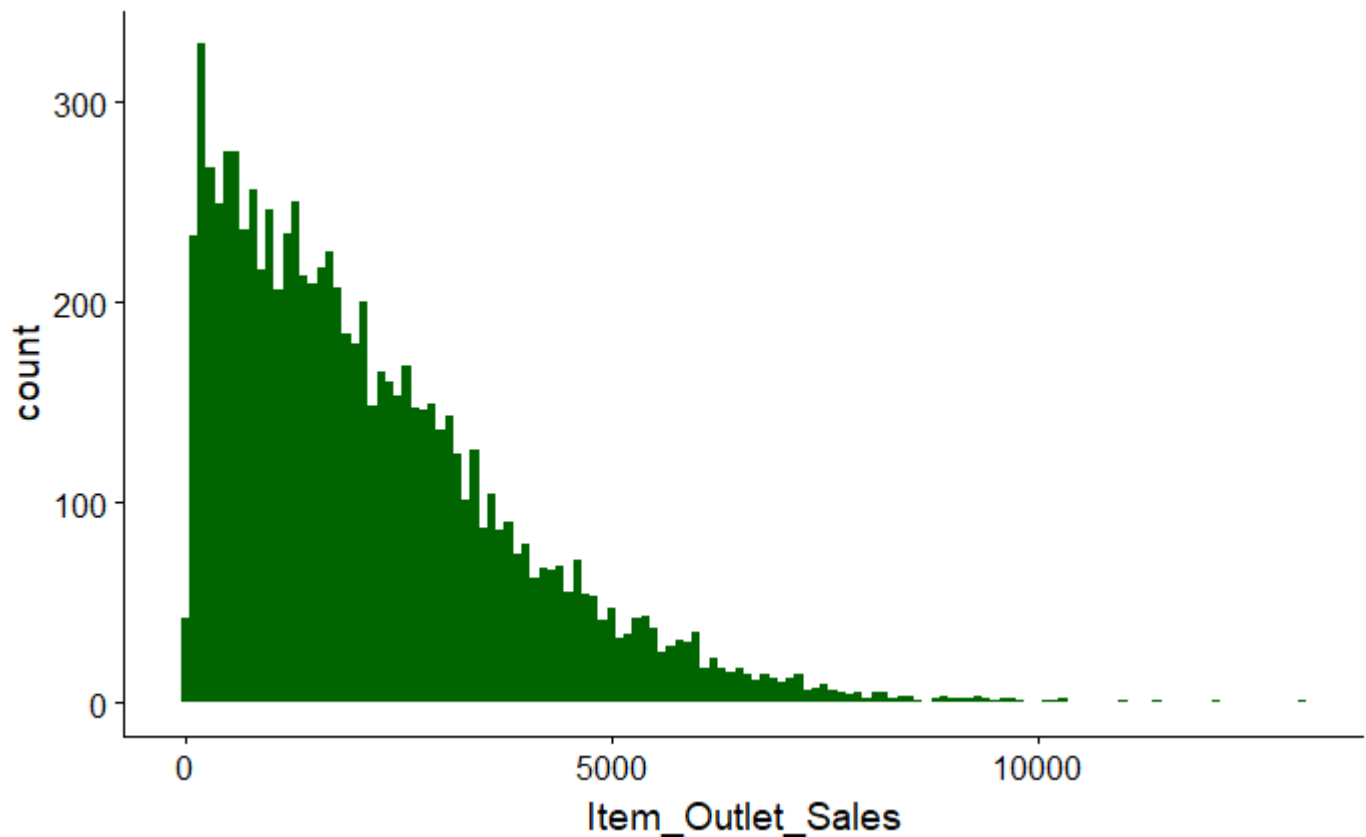
We will start off by plotting and exploring all the individual variables to gain some insights. In this tutorial ggplot2 package has been used to generate all the plots.

Target Variable

Since our target variable is continuous, we can visualise it by plotting its histogram.

[Hide](#)

```
ggplot(train) + geom_histogram(aes(train$Item_Outlet_Sales), binwidth = 100, fill = "darkgreen")
+
  xlab("Item_Outlet_Sales")
```



As you can see, it is a right skewed variable and would need some data transformation to treat its skewness. We will come back to it later.

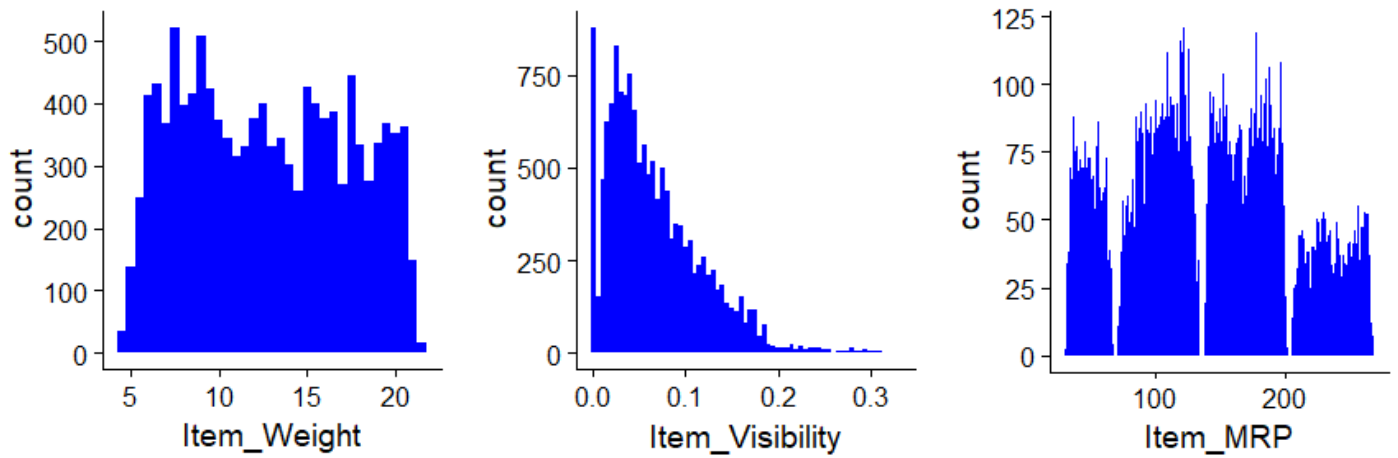
Independent Variables (numeric variables)

Now let's check the numeric independent variables. We'll again use the histograms for visualizations because that will help us in visualizing the distribution of individual variables.

Hide

```
p1 = ggplot(combi) + geom_histogram(aes(Item_Weight), binwidth = 0.5, fill = "blue")
p2 = ggplot(combi) + geom_histogram(aes(Item_Visibility), binwidth = 0.005, fill = "blue")
p3 = ggplot(combi) + geom_histogram(aes(Item_MRP), binwidth = 1, fill = "blue")
plot_grid(p1, p2, p3, nrow = 1) # plot_grid() from cowplot package
```

Removed 2439 rows containing non-finite values (stat_bin).



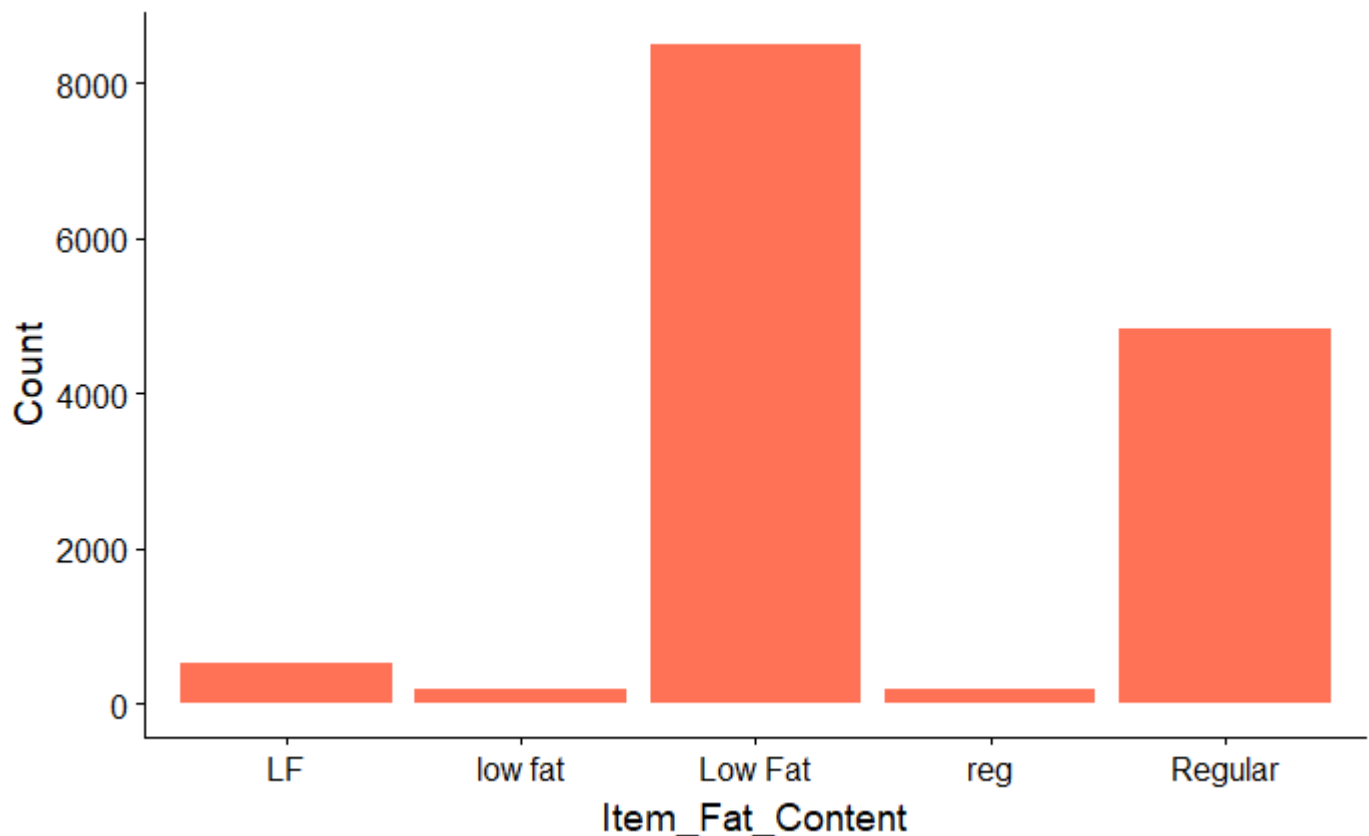
As you can see, there is no clear pattern in Item_Weight and Item_MRP. However, Item_Visibility is right-skewed and should be transformed to curb its skewness.

Independent Variables (categorical variables)

Now we'll try to explore and gain some insights from the categorical variables. A categorical variable or feature can have only a finite set of values. Let's first plot Item_Fat_Content.

Hide

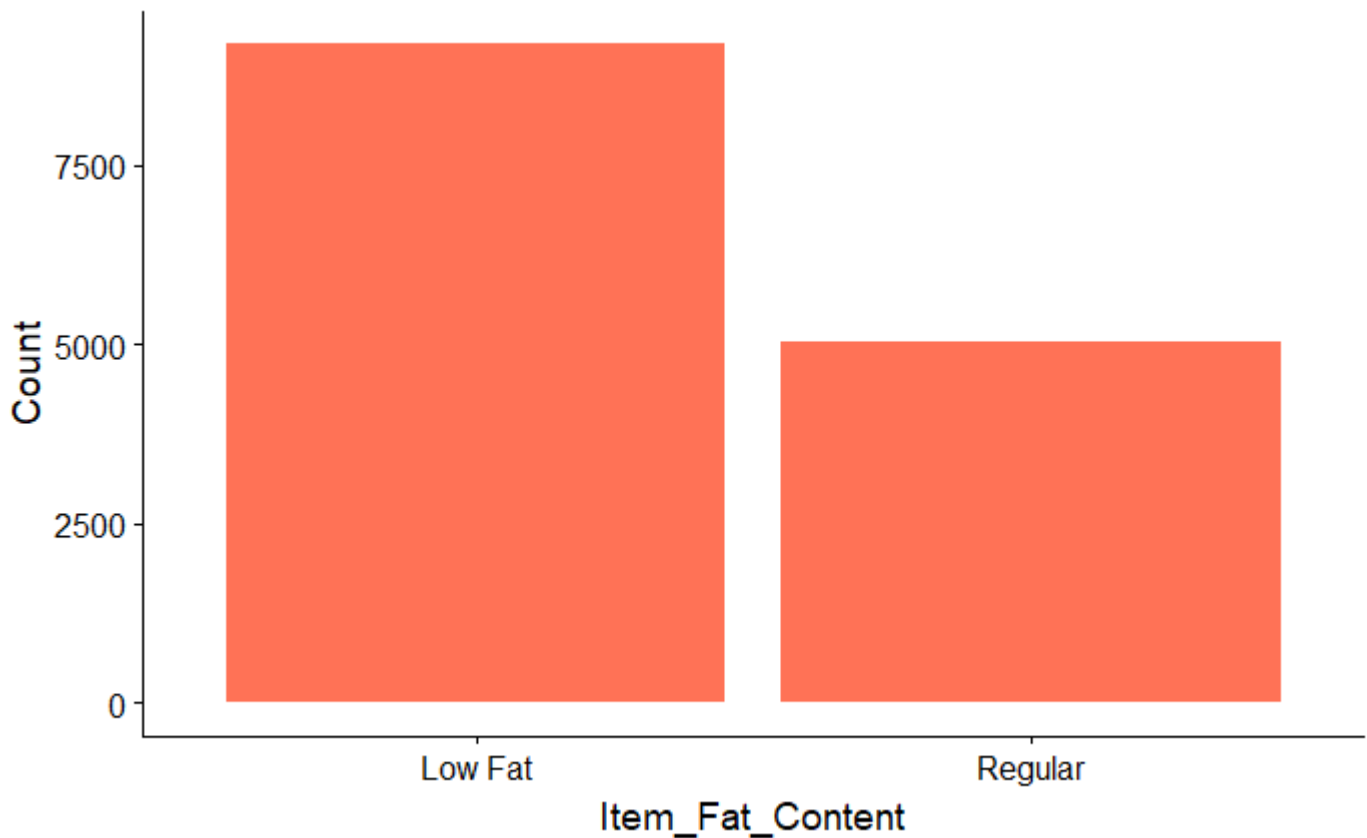
```
ggplot(combi %>% group_by(Item_Fat_Content) %>% summarise(Count = n())) +
  geom_bar(aes(Item_Fat_Content, Count), stat = "identity", fill = "coral1")
```



In the figure above, 'LF', 'low fat', and 'Low Fat' are the same category and can be combined into one. Similarly we can be done for 'reg' and 'Regular' into one. After making these corrections we'll plot the same figure again.

[Hide](#)

```
combi$Item_Fat_Content[combi$Item_Fat_Content == "LF"] = "Low Fat"  
combi$Item_Fat_Content[combi$Item_Fat_Content == "low fat"] = "Low Fat"  
combi$Item_Fat_Content[combi$Item_Fat_Content == "reg"] = "Regular"  
ggplot(combi %>% group_by(Item_Fat_Content) %>% summarise(Count = n())) +  
  geom_bar(aes(Item_Fat_Content, Count), stat = "identity", fill = "coral1")
```



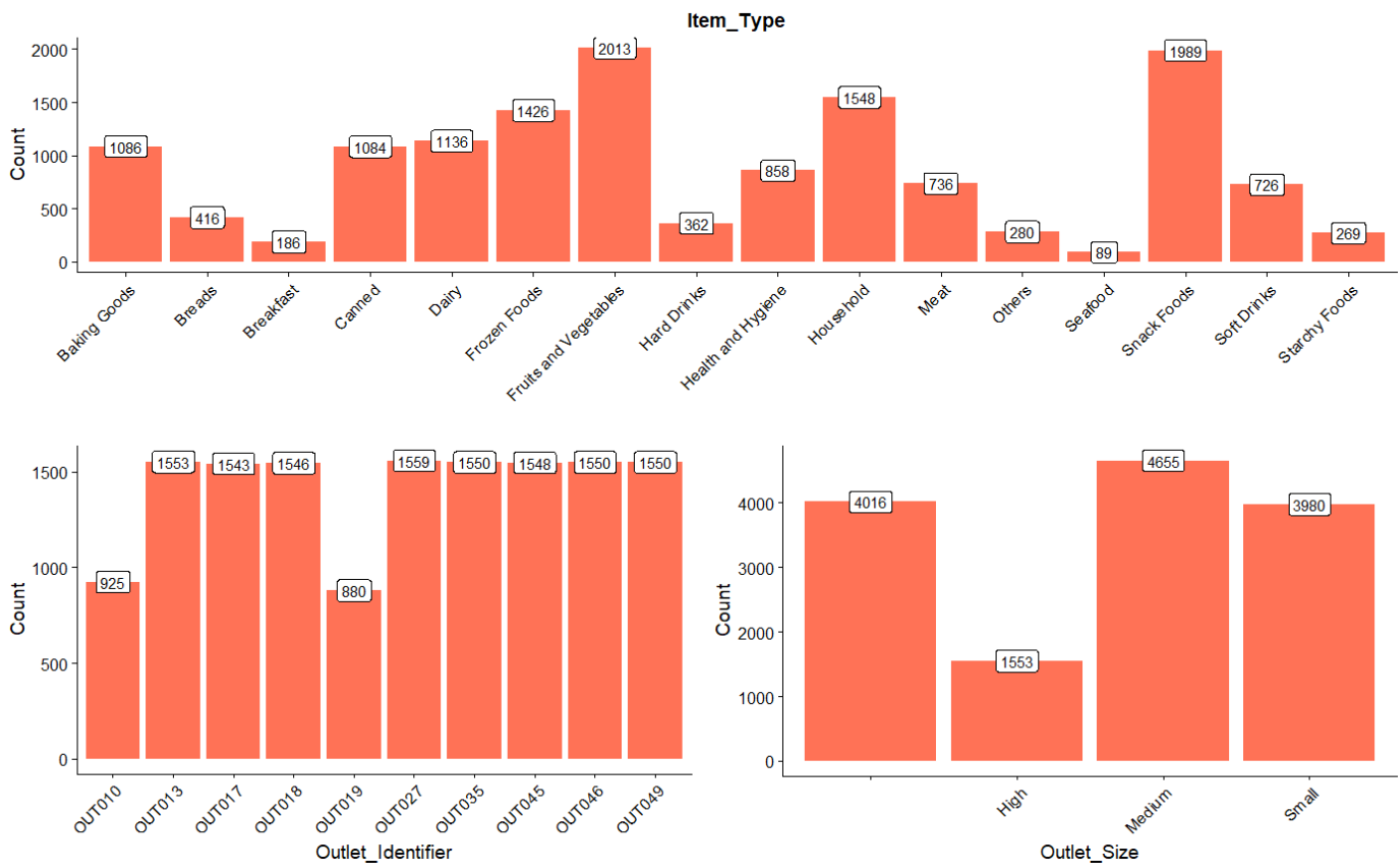
Now let's check the other categorical variables.

[Hide](#)

```

p4 = ggplot(combi %>% group_by(Item_Type) %>% summarise(Count = n())) +
  geom_bar(aes(Item_Type, Count), stat = "identity", fill = "coral1") +
  xlab("") +
  geom_label(aes(Item_Type, Count, label = Count), vjust = 0.5) +
  theme(axis.text.x = element_text(angle = 45, hjust = 1))+
  ggtitle("Item_Type")
p5 = ggplot(combi %>% group_by(Outlet_Identifier) %>% summarise(Count = n())) +
  geom_bar(aes(Outlet_Identifier, Count), stat = "identity", fill = "coral1") +
  geom_label(aes(Outlet_Identifier, Count, label = Count), vjust = 0.5) +
  theme(axis.text.x = element_text(angle = 45, hjust = 1))
p6 = ggplot(combi %>% group_by(Outlet_Size) %>% summarise(Count = n())) +
  geom_bar(aes(Outlet_Size, Count), stat = "identity", fill = "coral1") +
  geom_label(aes(Outlet_Size, Count, label = Count), vjust = 0.5) +
  theme(axis.text.x = element_text(angle = 45, hjust = 1))
second_row = plot_grid(p5, p6, nrow = 1)
plot_grid(p4, second_row, ncol = 1)

```



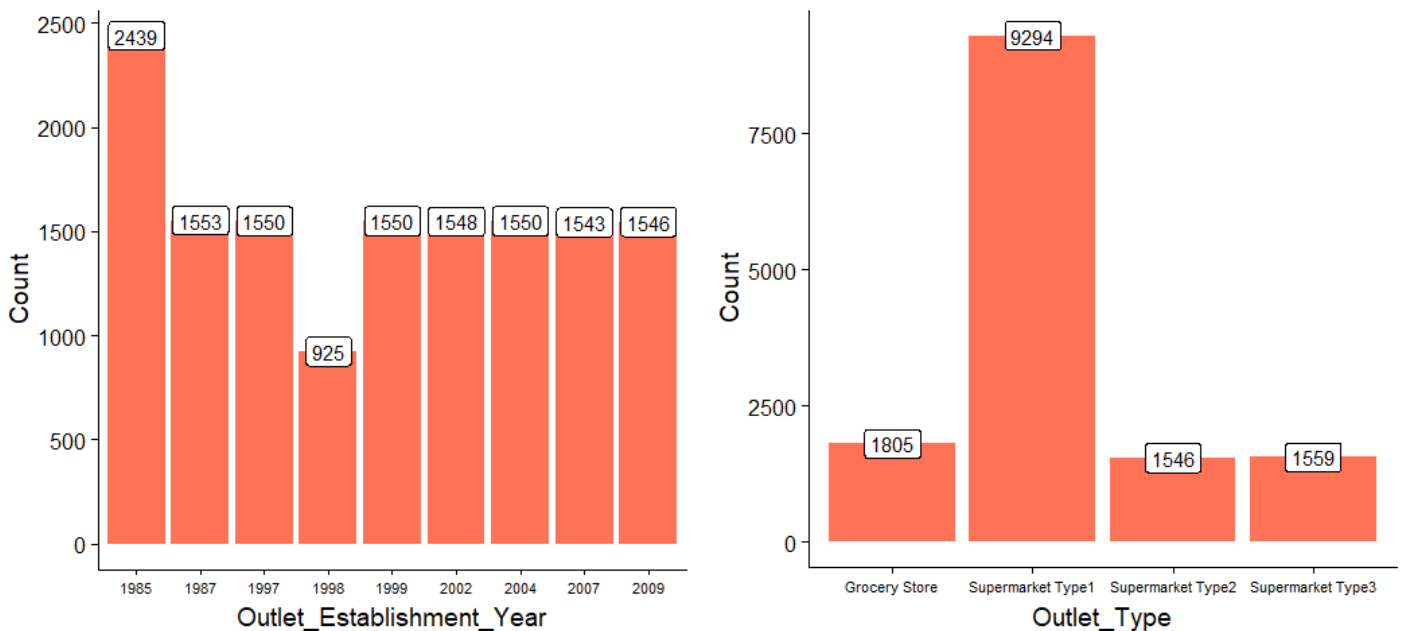
In Outlet_Size's plot, for 4016 observations, Outlet_Size is blank or missing. We will check for this in the bivariate analysis to substitute the missing values in the Outlet_Size.

Hide


```

p7 = ggplot(combi %>% group_by(Outlet_Establishment_Year) %>% summarise(Count = n())) +
  geom_bar(aes(factor(Outlet_Establishment_Year), Count), stat = "identity", fill = "coral1") +
  geom_label(aes(factor(Outlet_Establishment_Year), Count, label = Count), vjust = 0.5) +
  xlab("Outlet_Establishment_Year") +
  theme(axis.text.x = element_text(size = 8.5))
p8 = ggplot(combi %>% group_by(Outlet_Type) %>% summarise(Count = n())) +
  geom_bar(aes(Outlet_Type, Count), stat = "identity", fill = "coral1") +
  geom_label(aes(factor(Outlet_Type), Count, label = Count), vjust = 0.5) +
  theme(axis.text.x = element_text(size = 8.5))
plot_grid(p7, p8, ncol = 2)

```



There are lesser number of observations in the data for the outlets established in the year 1998 as compared to the other years. Supermarket Type 1 seems to be the most popular category of Outlet_Type.

Bivariate Analysis

After looking at every feature individually, let's now explore them again with respect to the target variable. Here we will make use of scatter plots for continuous or numeric variables and violin plots for the categorical variables.

```
train = combi[1:nrow(train)]
```

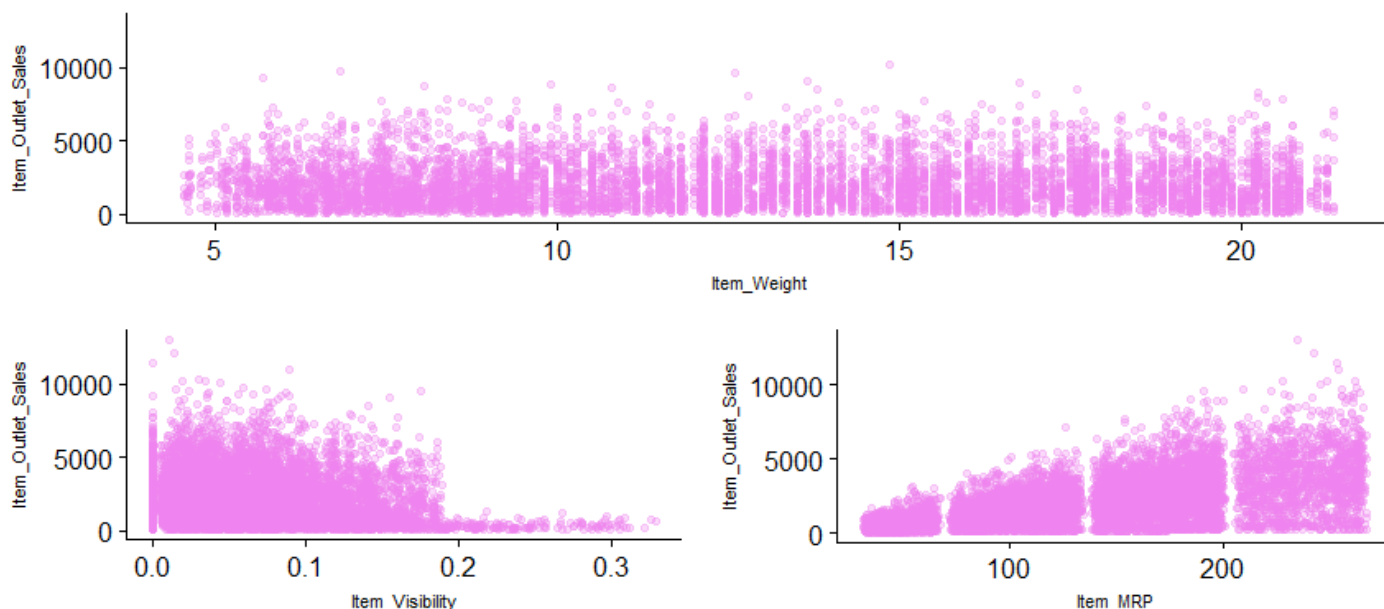
[Hide](#)
[Hide](#)

```

p9 = ggplot(train) + geom_point(aes(Item_Weight, Item_Outlet_Sales), colour = "violet", alpha = 0.3) +
  theme(axis.title = element_text(size = 8.5))
p10 = ggplot(train) + geom_point(aes(Item_Visibility, Item_Outlet_Sales), colour = "violet", alpha = 0.3) +
  theme(axis.title = element_text(size = 8.5))
p11 = ggplot(train) + geom_point(aes(Item_MRP, Item_Outlet_Sales), colour = "violet", alpha = 0.3) +
  theme(axis.title = element_text(size = 8.5))
second_row_2 = plot_grid(p10, p11, ncol = 2)
plot_grid(p9, second_row_2, nrow = 2)

```

Removed 1463 rows containing missing values (geom_point).



Item_Outlet_Sales is spread well across the entire range of the Item_Weight without any obvious pattern. In the Item_Visibility vs Item_Outlet_Sales, there is a string of points at Item_Visibility = 0.0 which seems strange as item visibility cannot be completely zero. We will take note of this issue and deal with it in the later stages.

In the third plot of Item_MRP vs Item_Outlet_Sales, we can clearly see 4 segments of prices that can be used in feature engineering to create a new variable.

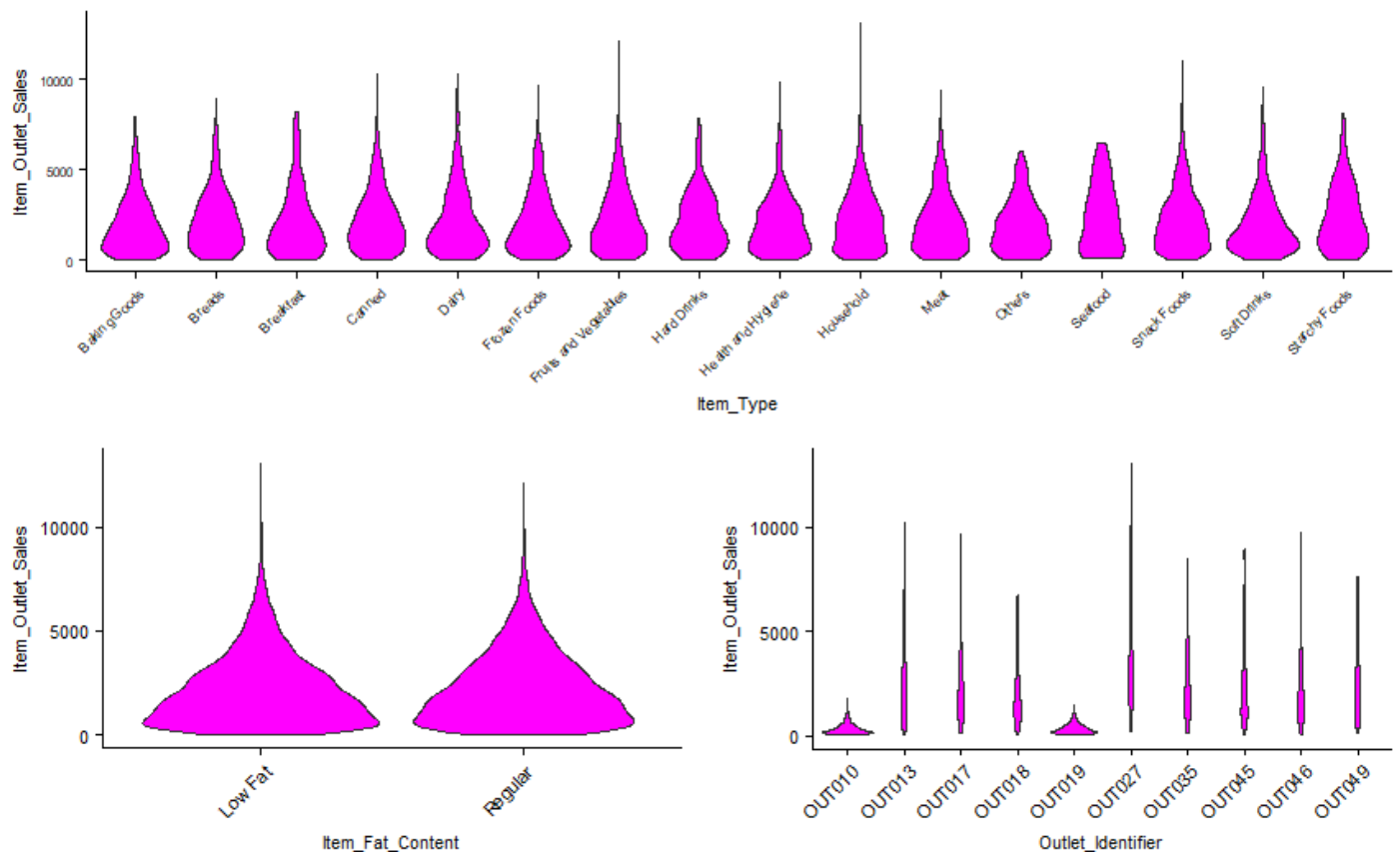
Now we'll visualise the categorical variables with respect to Item_Outlet_Sales. We will try to check the distribution of the target variable across all the categories of each of the categorical variable. We can do this using boxplots but instead we'll use the violin plots as they show the full distribution of the data. The horizontal width of a violin plot at a particular level indicates the concentration of data at that level.

Hide

```

p12 = ggplot(train) + geom_violin(aes(Item_Type, Item_Outlet_Sales), fill = "magenta") +
  theme(axis.text.x = element_text(angle = 45, hjust = 1),
        axis.text = element_text(size = 6),
        axis.title = element_text(size = 8.5))
p13 = ggplot(train) + geom_violin(aes(Item_Fat_Content, Item_Outlet_Sales), fill = "magenta") +
  theme(axis.text.x = element_text(angle = 45, hjust = 1),
        axis.text = element_text(size = 8),
        axis.title = element_text(size = 8.5))
p14 = ggplot(train) + geom_violin(aes(Outlet_Identifier, Item_Outlet_Sales), fill = "magenta") +
  theme(axis.text.x = element_text(angle = 45, hjust = 1),
        axis.text = element_text(size = 8),
        axis.title = element_text(size = 8.5))
second_row_3 = plot_grid(p13, p14, ncol = 2)
plot_grid(p12, second_row_3, ncol = 1)

```



Distribution of Item_Outlet_Sales across the categories of Item_Type is not very distinct and same is the case with Item_Fat_Content. However, the distribution for OUT010 and OUT019 categories of Outlet_Identifier are quite similar and very much different from the rest of the categories of Outlet_Identifier.

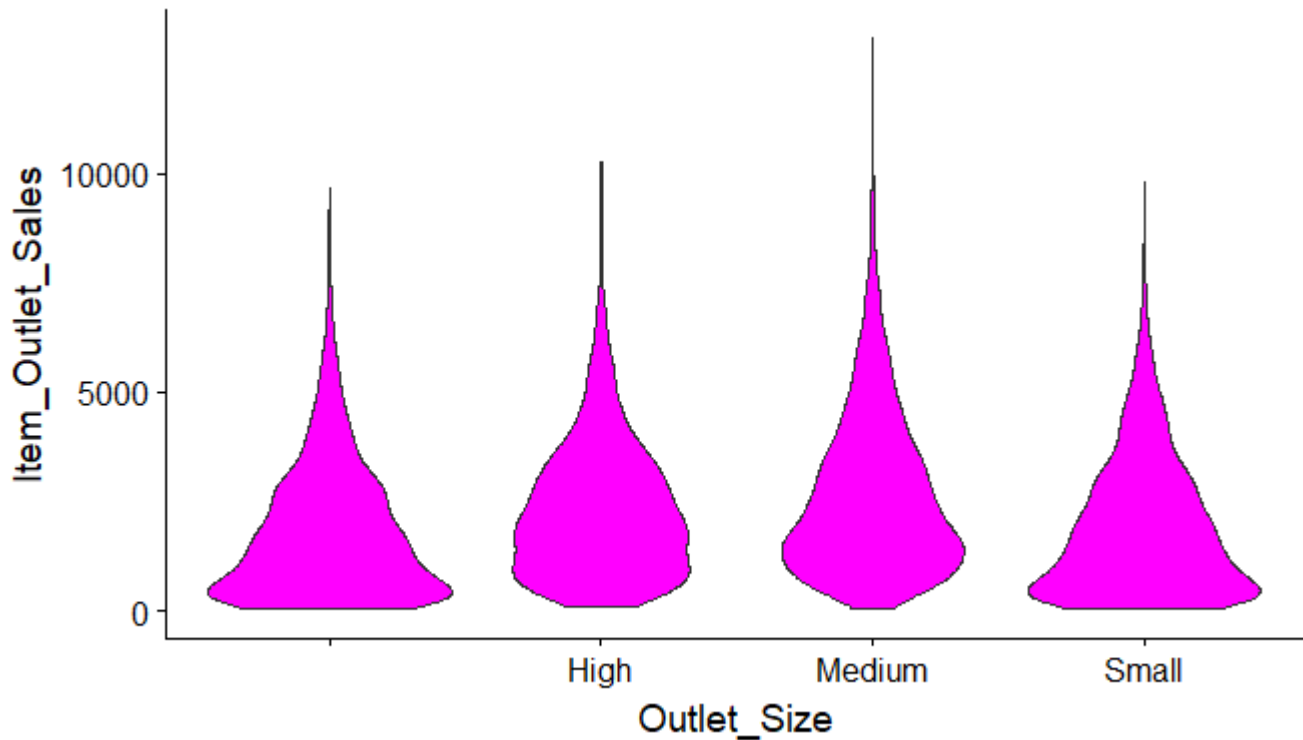
In the univariate analysis, we came to know about the empty values in Outlet_Size variable. Let's check the distribution of the target variable across Outlet_Size.

[Hide](#)

```

ggplot(train) + geom_violin(aes(Outlet_Size, Item_Outlet_Sales), fill = "magenta")

```

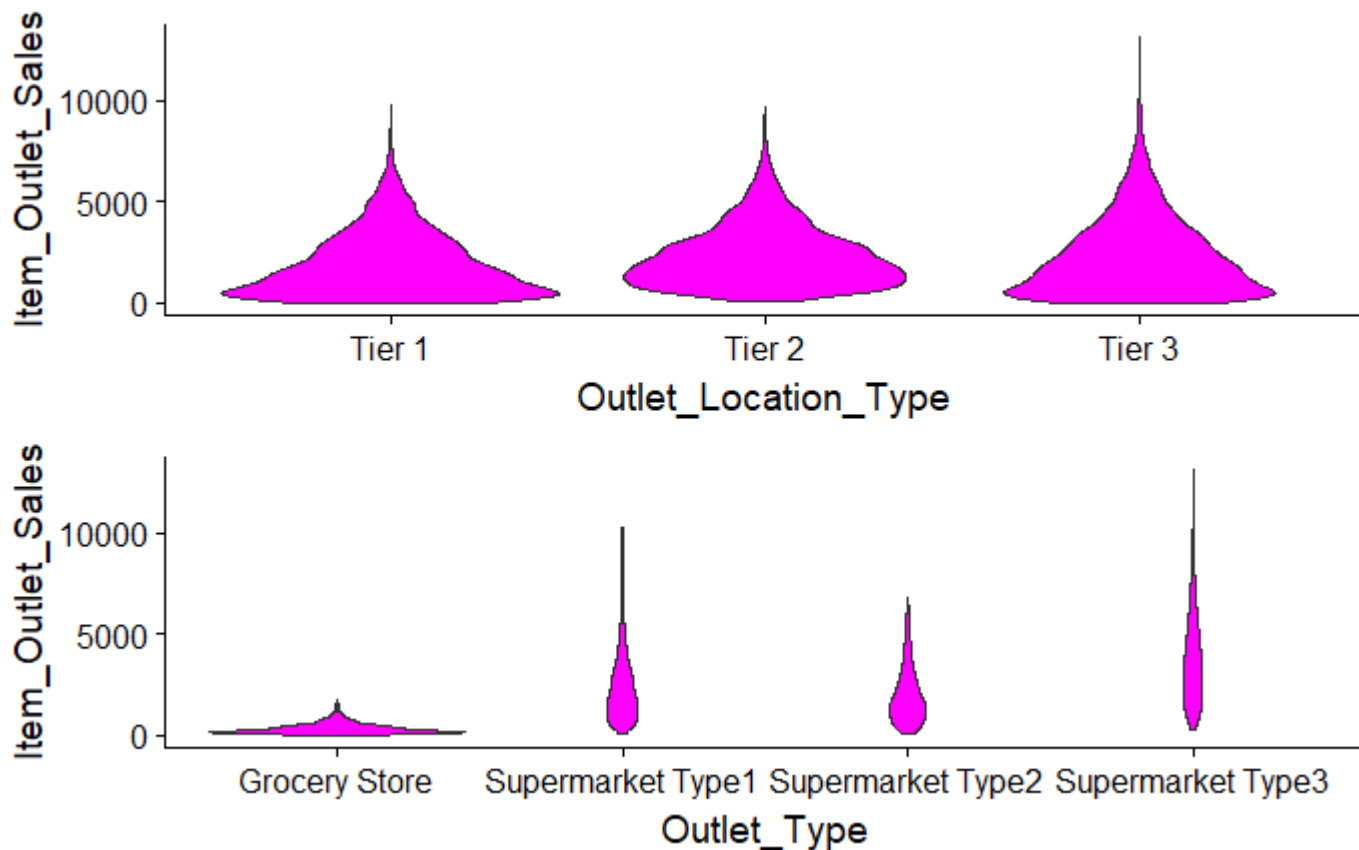


The distribution of 'Small' Outlet_Size is almost identical to the distribution of the blank category (first violin) of Outlet_Size. So, we can substitute the blanks in Outlet_Size with 'Small'. Please note that this is not the only way to impute missing values, but for the time being we will go ahead and impute the missing values with 'Small'.

Let's examine the remaining variables.

[Hide](#)

```
p15 = ggplot(train) + geom_violin(aes(Outlet_Location_Type, Item_Outlet_Sales), fill = "magenta"
)
p16 = ggplot(train) + geom_violin(aes(Outlet_Type, Item_Outlet_Sales), fill = "magenta")
plot_grid(p15, p16, ncol = 1)
```



Tier 1 and Tier 3 locations of Outlet_Location_Type look similar. In the Outlet_Type, Grocery Store has most of its data points around the lower sales values as compared to the other categories. These are the kind of insights that we can extract by visualizing our data. Hence, data visualization should be an important part of any kind data analysis.

Missing Value Treatment

Missing data can have a severe impact on building predictive models because the missing values might contain some vital information which could help in making better predictions. So, it becomes imperative to carry out missing data imputation. There are different methods to treat missing values based on the problem and the data. Some of the common techniques are as follows:

1. **Deletion of rows:** In train dataset, observations having missing values in any variable are deleted. The downside of this method is the loss of information and drop in prediction power of model.
2. **Mean/Median/Mode Imputation:** In case of continuous variable, missing values can be replaced with mean or median of all known values of that variable. For categorical variables, we can use mode of the given values to replace the missing values.
3. **Building Prediction Model:** We can even make a predictive model to impute missing data in a variable. Here we will treat the variable having missing data as the target variable and the other variables as predictors. We will divide our data into 2 datasets—one without any missing value for that variable and the other with missing values for that variable. The former set would be used as training set to build the predictive model and it would then be applied to the latter set to predict the missing values.

You can try the following code to quickly find missing values in a dataset.

Hide

```
colSums(is.na(combi))
```

	Item_Identifier	Item_Weight	Item_Fat_Content	Item_Vis
ibility	0	2439	0	
0	Item_Type	Item_MRP	Outlet_Identifier	Outlet_Establishme
nt_Year	0	0	0	
0	Outlet_Size	Outlet_Location_Type	Outlet_Type	Item_Outle
t_Sales	0	0	0	
5681				

As you can see above, we have missing values in Item_Weight and Item_Outlet_Sales. Missing data in Item_Outlet_Sales can be ignored since they belong to the test dataset. We'll now impute Item_Weight with mean weight based on the Item_Identifier variable.

Hide

```
missing_index = which(is.na(combi$Item_Weight))
for(i in missing_index){

  item = combi$Item_Identifier[i]
  combi$Item_Weight[i] = mean(combi$Item_Weight[combi$Item_Identifier == item], na.rm = T)

}
```

Replacing 0's in Item_Visibility variable

Hide

```
zero_index = which(combi$Item_Visibility == 0)
for(i in zero_index){

  item = combi$Item_Identifier[i]
  combi$Item_Visibility[i] = mean(combi$Item_Visibility[combi$Item_Identifier == item], na.rm = T)

}
```

Feature Engineering

Most of the times the given features in a dataset are not enough to give satisfactory predictions. In such cases, we have to create new features which might help in improving the model's performance. Let's try to create some new features for our dataset.

We can have a look at the Item_Type variable and classify the categories into perishable and non_perishable as per our understanding and make it into a new feature.

Hide

```
# create a new feature 'Item_Type_new'
perishable = c("Breads", "Breakfast", "Dairy", "Fruits and Vegetables", "Meat", "Seafood")
non_perishable = c("Baking Goods", "Canned", "Frozen Foods", "Hard Drinks", "Health and Hygiene"
,
                  "Household", "Soft Drinks")
combi[,Item_Type_new := ifelse(Item_Type %in% perishable, "perishable",
                              ifelse(Item_Type %in% non_perishable, "non_perishable", "not_sur
e"))]
```

Let's compare Item_Type with the first 2 characters of Item_Identifier, i.e., 'DR', 'FD', and 'NC'. These identifiers most probably stand drinks, food, and non-consumable.

Hide

```
table(combi$Item_Type, substr(combi$Item_Identifier, 1, 2))
```

	DR	FD	NC
Baking Goods	0	1086	0
Breads	0	416	0
Breakfast	0	186	0
Canned	0	1084	0
Dairy	229	907	0
Frozen Foods	0	1426	0
Fruits and Vegetables	0	2013	0
Hard Drinks	362	0	0
Health and Hygiene	0	0	858
Household	0	0	1548
Meat	0	736	0
Others	0	0	280
Seafood	0	89	0
Snack Foods	0	1989	0
Soft Drinks	726	0	0
Starchy Foods	0	269	0

Based on the above table we can create a new feature. Let's call it Item_category.

Hide

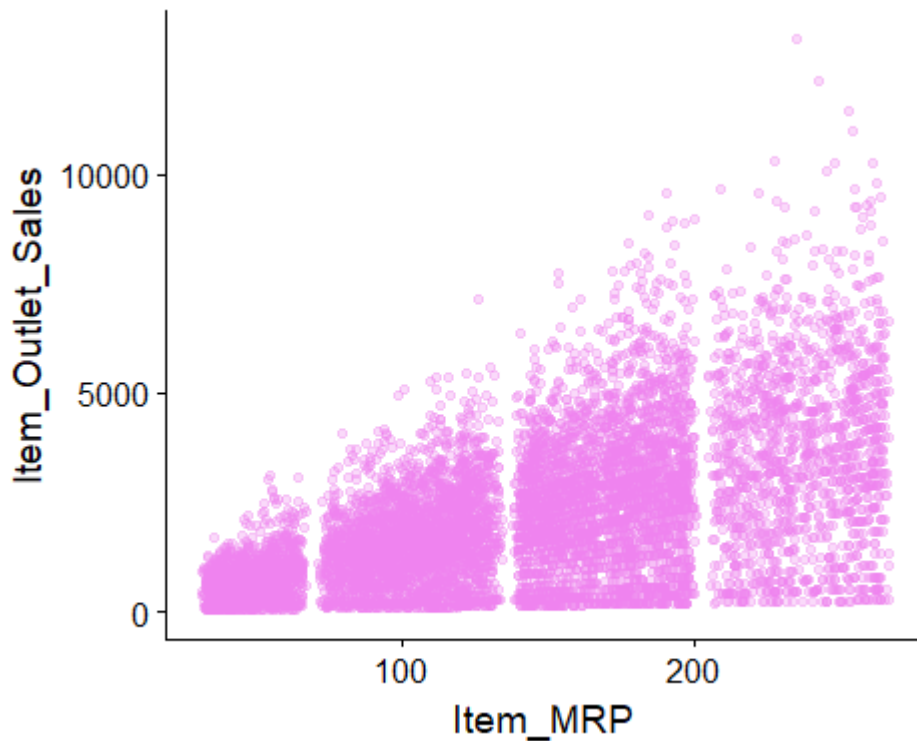
```
combi[,Item_category := substr(combi$Item_Identifier, 1, 2)]
```

We will also change the values of Item_Fat_Content wherever Item_category is 'NC' because non-consumable items cannot have any fat content. We will also create a couple of more features — Outlet_Years (years of operation) and price_per_unit_wt (price per unit weight).

Hide

```
combi$Item_Fat_Content[combi$Item_category == "NC"] = "Non-Edible"
# years of operation of outlets
combi[,Outlet_Years := 2013 - Outlet_Establishment_Year]
combi$Outlet_Establishment_Year = as.factor(combi$Outlet_Establishment_Year)
# Price per unit weight
combi[,price_per_unit_wt := Item_MRP/Item_Weight]
```

Earlier in the Item_MRP vs Item_Outlet_Sales plot, we saw Item_MRP was spread across in 4 chunks. We can use k Means clustering to create 4 groups using Item_MRP variable. K Means clustering requires prior knowledge of K i.e. no. of clusters we want to divide our data into. Here we will go ahead with K=4. Let's try it out.


[Hide](#)

```
Item_MRP_clusters = kmeans(combi$Item_MRP, centers = 4)
table(Item_MRP_clusters$cluster) # display no. of observations in each cluster
```

```
1    2    3    4
2556 4317 4931 2400
```

Let's use these clusters as a new independent feature in our data.

[Hide](#)

```
combi$Item_MRP_clusters = as.factor(Item_MRP_clusters$cluster)
```

Encoding Categorical Variables

In this stage, we will convert our categorical variables into numerical ones. We will use 2 techniques — Label Encoding and One Hot Encoding.

1. **Label encoding** simply means converting each category in a variable to a number. It is more suitable for ordinal variables — categorical variables with some order.
2. In **One hot encoding**, each category of a categorical variable is converted into a new binary column (1/0).

We will use both the encoding techniques.

Label encoding for the categorical variables

We will label encode Outlet_Size and Outlet_Location_Type as these are ordinal variables.

[Hide](#)

```
combi[,Outlet_Size_num := ifelse(Outlet_Size == "Small", 0,
                                ifelse(Outlet_Size == "Medium", 1, 2))]
combi[,Outlet_Location_Type_num := ifelse(Outlet_Location_Type == "Tier 3", 0,
                                          ifelse(Outlet_Location_Type == "Tier 2", 1, 2))]
# removing categorical variables after label encoding
combi[, c("Outlet_Size", "Outlet_Location_Type") := NULL]
```

One hot encoding for the categorical variable

[Hide](#)

```
ohe = dummyVars("~.", data = combi[, -c("Item_Identifier", "Outlet_Establishment_Year", "Item_Type")], fullRank = T)
ohe_df = data.table(predict(ohe, combi[, -c("Item_Identifier", "Outlet_Establishment_Year", "Item_Type")]))
combi = cbind(combi[, "Item_Identifier"], ohe_df)
```

PreProcessing Data

Before feeding our data into any model, it is a good practice to preprocess the data. We will do preprocessing on both independent variables and target variable

Checking Skewness

Skewness in variables is undesirable for predictive modeling. Some machine learning methods assume normally distributed data and a skewed variable can be transformed by taking its log, square root, or cube root so as to make the distribution of the skewed variable as close to normal distribution as possible.

[Hide](#)

```
skewness(combi$Item_Visibility); skewness(combi$price_per_unit_wt)
```

```
[1] 1.254046
[1] 1.304551
```

[Hide](#)

```
combi[,Item_Visibility := log(Item_Visibility + 1)] # log + 1 to avoid division by zero
combi[,price_per_unit_wt := log(price_per_unit_wt + 1)]
```

Scaling numeric predictors

Let's scale and center the numeric variables to make them have a mean of zero, standard deviation of one and scale of 0 to 1. Scaling and centering is required for linear regression models.

[Hide](#)

```
num_vars = which(sapply(combi, is.numeric)) # index of numeric features
num_vars_names = names(num_vars)
combi_numeric = combi[,setdiff(num_vars_names, "Item_Outlet_Sales"), with = F]
prep_num = preProcess(combi_numeric, method=c("center", "scale"))
combi_numeric_norm = predict(prepare_num, combi_numeric)
```

[Hide](#)

```
combi[,setdiff(num_vars_names, "Item_Outlet_Sales") := NULL] # removing numeric independent variables
combi = cbind(combi, combi_numeric_norm)
```

[Hide](#)

```
train = combi[1:nrow(train)]
test = combi[(nrow(train) + 1):nrow(combi)]
test[,Item_Outlet_Sales := NULL] # removing Item_Outlet_Sales as it contains only NA
```

Correlated Variables

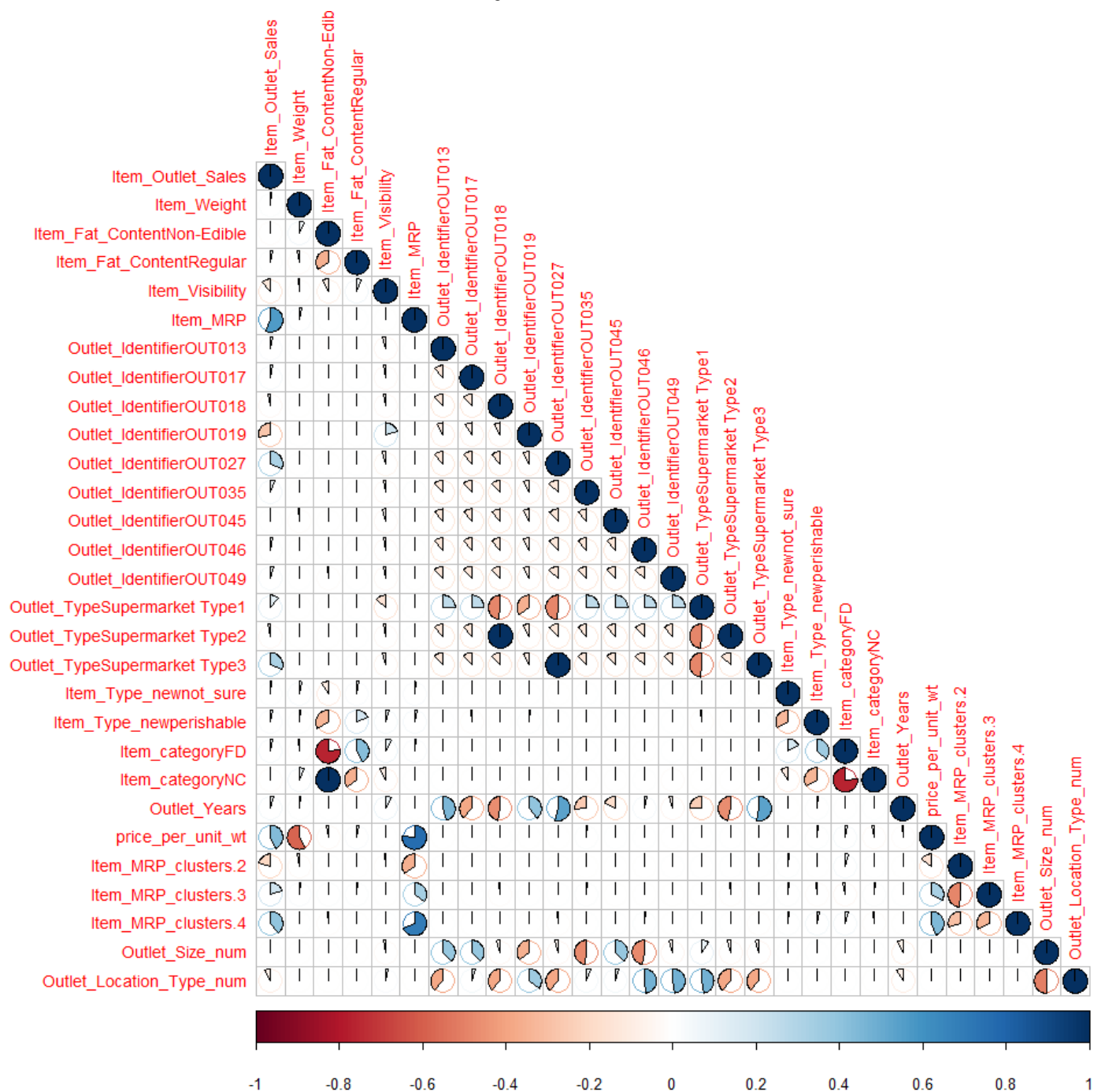
Let's examine the correlated features of train dataset. Correlation varies from -1 to 1.

1. negative correlation: < 0 and ≥ -1
2. positive correlation: > 0 and ≤ 1
3. no correlation: 0

It is not desirable to have correlated features if we are using linear regressions.

[Hide](#)

```
cor_train = cor(train[, -c("Item_Identifier")])
corrplot(cor_train, method = "pie", type = "lower", tl.cex = 0.9)
```



Variables `price_per_unit_wt` and `Item_Weight` are highly correlated as the former one was created from the latter. Similarly `price_per_unit_wt` and `Item_MRP` are highly correlated for the same reason.

Modeling

Finally we have arrived at most interesting stage of the whole process — predictive modeling. We will start off with the simpler models and gradually move on to more sophisticated models. We will build the models using...

- Linear Regression

- Lasso Regression
- Ridge Regression
- RandomForest
- XGBoost

We will use **5-fold cross validation** in all the models we are going to build. Basically cross validation gives an idea how well a model generalizes to unseen data.

Linear Regression

Linear regression is the simplest and most widely used statistical technique for predictive modeling. Given below is the linear regression equation:

$$Y = \theta_1 X_1 + \theta_2 X_2 + \dots + \theta_n X_n$$

where X_1, X_2, \dots, X_n are the independent variables, Y is the target variable and all thetas are the coefficients. Magnitude of a coefficient wrt to the other coefficients determines the importance of the corresponding independent variable.

For a good linear regression model, the data should satisfy a few assumptions which can be found in this article (<https://www.analyticsvidhya.com/blog/2016/07/deeper-regression-analysis-assumptions-plots-solutions/>) in detail. One of these assumptions is that of absence of multicollinearity, i.e, the independent variables should be correlated. However, as per the correlation plot above, we have a few highly correlated independent variables in our data. This issue of multicollinearity can be dealt with regularization.

For the time being, let's build our linear regression model with all the variables.

Hide

```
set.seed(1234)
my_control = trainControl(method="cv", number=5)
linear_reg_mod = train(x = train[, -c("Item_Identifier", "Item_Outlet_Sales")], y = train$Item_Outlet_Sales,
                       method='glmnet', trControl= my_control)
print("5- fold cross validation scores:")
```

```
[1] "5- fold cross validation scores:"
```

Hide

```
print(round(linear_reg_mod$resample$RMSE, 2))
```

```
[1] 1159.97 1121.63 1089.28 1115.11 1162.17
```

Making predictions on the test dataset

Hide

```
#submission$Item_Outlet_Sales = (predict(linear_reg_mod, test[, -c("Item_Identifier"))))^2
submission$Item_Outlet_Sales = predict(linear_reg_mod, test[, -c("Item_Identifier")])
write.csv(submission, "Linear_Reg_submit_2_21_Apr_18.csv", row.names = F)
```

Leaderboard score: 1202.26

Regularised regression models can handle the correlated independent variables well and helps in overcoming overfitting. **Ridge** penalty shrinks the coefficients of correlated predictors towards each other while the **Lasso** tends to pick one of a pair of correlated features and discard the other. The tuning parameter **lambda** controls the strength of the penalty.

Lasso Regression

[Hide](#)

```
set.seed(1235)
my_control = trainControl(method="cv", number=5)
Grid = expand.grid(alpha = 1, lambda = seq(0.001, 0.1, by = 0.0002))

lasso_linear_reg_mod = train(x = train[, -c("Item_Identifier", "Item_Outlet_Sales")], y = train
$Item_Outlet_Sales,
                             method='glmnet', trControl= my_control, tuneGrid = Grid)
```

Leaderboard score: 1202.26

Ridge Regression

[Hide](#)

```
set.seed(1236)
my_control = trainControl(method="cv", number=5)
Grid = expand.grid(alpha = 0, lambda = seq(0.001, 0.1, by = 0.0002))

ridge_linear_reg_mod = train(x = train[, -c("Item_Identifier", "Item_Outlet_Sales")], y = train
$Item_Outlet_Sales,
                             method='glmnet', trControl= my_control, tuneGrid = Grid)
```

Leaderboard score: 1219.08

RandomForest

RandomForest is a tree based bootstrapping algorithm wherein a certain no. of weak learners (decision trees) are combined to make a powerful prediction model. For every individual learner, a random sample of rows and a few randomly chosen variables are used to build a decision tree model. Final prediction can be a function of all the predictions made by the individual learners. In case of regression problem, the final prediction can be mean of all the predictions. For detailed explanation visit this article (<https://www.analyticsvidhya.com/blog/2016/04/comprehensive-tutorial-tree-based-modeling-scratch-in-python/>).

We will now build a RandomForest model with 400 trees. The other tuning parameters used here are `mtry` — no. of predictor variables randomly sampled at each split, and `min.node.size` — minimum size of terminal nodes (setting this number large causes smaller trees and reduces overfitting).

Hide

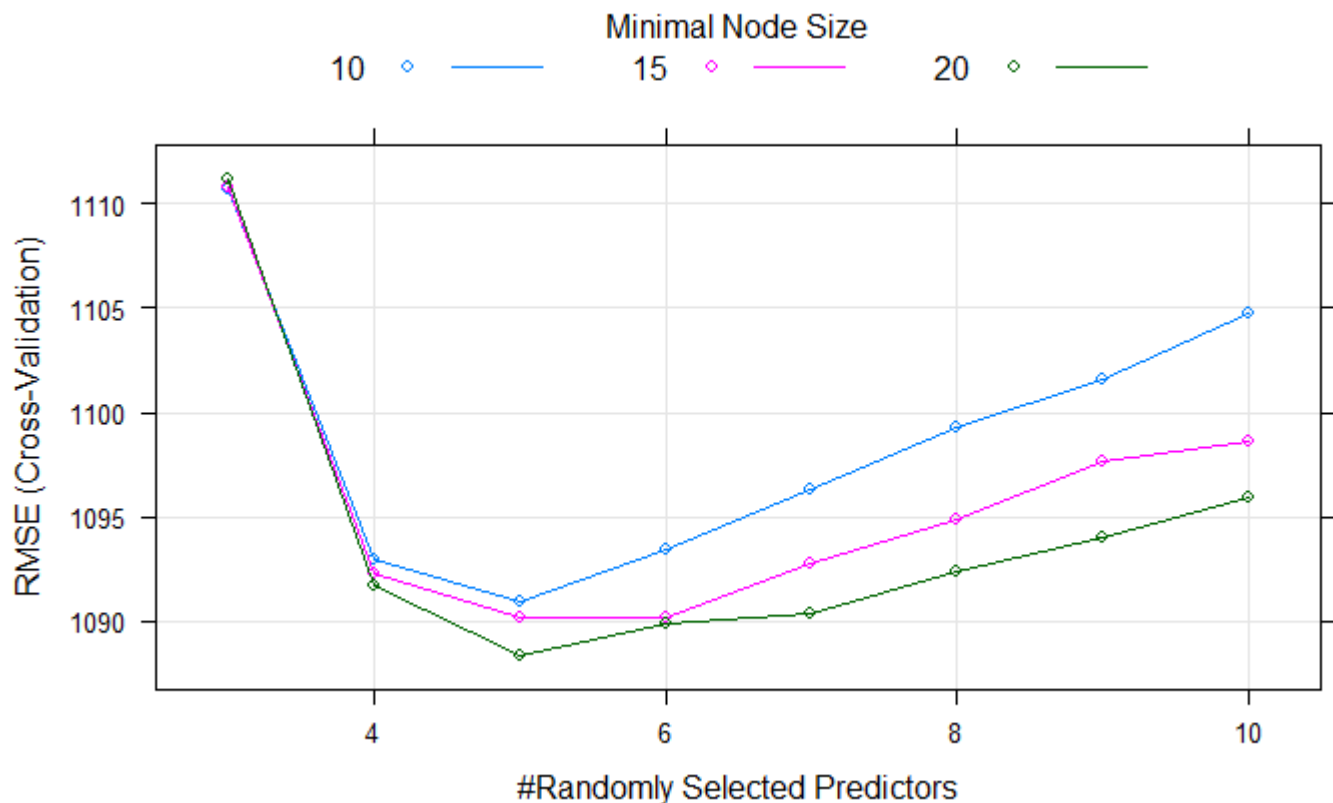
```
set.seed(1237)
my_control = trainControl(method="cv", number=5)

tgrid = expand.grid(
  .mtry = c(3:10),
  .splitrule = "variance",
  .min.node.size = c(10,15,20)
)

rf_mod = train(x = train[, -c("Item_Identifier", "Item_Outlet_Sales")],
  y = train$Item_Outlet_Sales,
  method='ranger',
  trControl= my_control,
  tuneGrid = tgrid,
  num.trees = 400,
  importance = "permutation")
```

Hide

```
plot(rf_mod)
```



As per the plot shown above, the best score is achieved at `mtry = 5` and `min.node.size = 20`.

Leaderboard score: 1157.25

XGBoost Model

XGBoost is a fast and efficient algorithm and has been used to by the winners of many data science competitions. It's a boosting algorithm and you may refer this article (<https://www.analyticsvidhya.com/blog/2015/11/quick-introduction-boosting-algorithms-machine-learning/>) to know more about boosting. XGBoost works only with numeric variables and we have already done that. There are many tuning parameters in XGBoost which can be broadly classified into General Parameters, Booster Parameters and Task Parameters.

- General parameters refers to which booster we are using to do boosting. The commonly used are tree or linear model
- Booster parameters depends on which booster you have chosen
- Learning Task parameters that decides on the learning scenario, for example, regression tasks may use different parameters with ranking tasks.

Let's have a look at the parameters that we are going to use in our model.

1. **eta**: It is also known as the learning rate or the shrinkage factor. It actually shrinks the feature weights to make the boosting process more conservative. The range is 0 to 1. Low eta value means model is more robust to overfitting.
2. **gamma**: The range is 0 to ∞ . Larger the gamma more conservative the algorithm is.
3. **max_depth**: We can specify maximum depth of a tree using this parameter.
4. **subsample**: It is the proportion of rows that the model will randomly select to grow trees.
5. **colsample_bytree**: It is the ratio of variables randomly chosen for build each tree in the model.

[Hide](#)

```
param_list = list(  
  
    objective = "reg:linear",  
    eta=0.01,  
    gamma = 1,  
    max_depth=6,  
    subsample=0.8,  
    colsample_bytree=0.5  
)
```

[Hide](#)

```
dtrain = xgb.DMatrix(data = as.matrix(train[,-c("Item_Identifier", "Item_Outlet_Sales")])), label  
= train$Item_Outlet_Sales)  
dtest = xgb.DMatrix(data = as.matrix(test[,-c("Item_Identifier")]))
```

[Hide](#)

```
set.seed(112)
xgbcv = xgb.cv(params = param_list,
               data = dtrain,
               nrounds = 1000,
               nfold = 5,
               print_every_n = 10,
               early_stopping_rounds = 30,
               maximize = F)
```



```
[1] train-rmse:2746.726611+7.335351 test-rmse:2746.741455+30.997620
Multiple eval metrics are present. Will use test_rmse for early stopping.
Will train until test_rmse hasn't improved in 30 rounds.

[11] train-rmse:2537.407764+5.897640 test-rmse:2539.186767+31.033561
[21] train-rmse:2349.479150+4.506506 test-rmse:2353.149951+32.307349
[31] train-rmse:2182.710596+4.820427 test-rmse:2188.130664+31.642597
[41] train-rmse:2033.751758+3.185745 test-rmse:2040.717822+31.669930
[51] train-rmse:1901.117578+2.685327 test-rmse:1910.519629+30.535673
[61] train-rmse:1784.638452+2.785332 test-rmse:1796.658545+28.845284
[71] train-rmse:1680.817847+2.442780 test-rmse:1695.664795+28.498337
[81] train-rmse:1590.710376+2.355847 test-rmse:1608.133936+28.083114
[91] train-rmse:1511.284009+1.938146 test-rmse:1531.907446+27.753792
[101] train-rmse:1441.667432+1.669035 test-rmse:1465.339282+29.325773
[111] train-rmse:1381.855786+2.373448 test-rmse:1408.609692+27.946078
[121] train-rmse:1329.842700+2.424946 test-rmse:1359.539600+27.418424
[131] train-rmse:1283.422485+2.876224 test-rmse:1316.416382+26.402673
[141] train-rmse:1243.852637+2.540695 test-rmse:1280.471533+25.467234
[151] train-rmse:1210.363989+2.557588 test-rmse:1250.352270+25.341492
[161] train-rmse:1180.807373+2.830153 test-rmse:1224.257568+24.374810
[171] train-rmse:1155.128735+2.460500 test-rmse:1202.033203+23.914640
[181] train-rmse:1133.309326+2.638564 test-rmse:1183.441016+23.106944
[191] train-rmse:1114.457251+2.588460 test-rmse:1167.850073+22.394809
[201] train-rmse:1098.201440+3.044061 test-rmse:1154.797119+21.841495
[211] train-rmse:1083.706030+3.528924 test-rmse:1143.217773+21.357599
[221] train-rmse:1071.321094+3.570172 test-rmse:1134.130224+20.882285
[231] train-rmse:1060.626660+3.701624 test-rmse:1126.470825+20.581249
[241] train-rmse:1051.289917+3.347872 test-rmse:1119.966406+20.567948
[251] train-rmse:1043.086231+3.365346 test-rmse:1114.744653+20.188325
[261] train-rmse:1035.549756+3.017828 test-rmse:1110.146753+19.898095
[271] train-rmse:1028.656812+3.034549 test-rmse:1106.056055+19.461418
[281] train-rmse:1022.382519+3.010057 test-rmse:1102.918799+19.266198
[291] train-rmse:1016.968725+2.859337 test-rmse:1100.386548+19.148728
[301] train-rmse:1011.780688+2.829672 test-rmse:1098.218945+18.923820
[311] train-rmse:1007.186755+2.809837 test-rmse:1096.416211+18.866274
[321] train-rmse:1002.984875+2.724723 test-rmse:1095.036182+18.856279
[331] train-rmse:999.086707+2.610513 test-rmse:1093.826953+18.811409
[341] train-rmse:995.230347+2.489401 test-rmse:1092.823755+18.730222
[351] train-rmse:991.635303+2.562694 test-rmse:1091.962793+18.797019
[361] train-rmse:988.453491+2.598685 test-rmse:1091.367261+18.735331
[371] train-rmse:985.262891+2.528210 test-rmse:1090.897730+18.676589
[381] train-rmse:982.168066+2.491125 test-rmse:1090.334131+18.727500
[391] train-rmse:979.388269+2.496498 test-rmse:1090.044995+18.728524
[401] train-rmse:976.569849+2.411764 test-rmse:1089.951758+18.669657
[411] train-rmse:973.715760+2.416732 test-rmse:1089.936035+18.802668
[421] train-rmse:971.002002+2.453059 test-rmse:1089.889087+18.826629
[431] train-rmse:968.387219+2.438929 test-rmse:1089.820801+18.791873
[441] train-rmse:965.961414+2.681547 test-rmse:1089.849414+18.777218
[451] train-rmse:963.600769+2.733766 test-rmse:1089.912305+18.756355
Stopping. Best iteration:
[430] train-rmse:968.672278+2.411026 test-rmse:1089.806958+18.821152
```

[Hide](#)

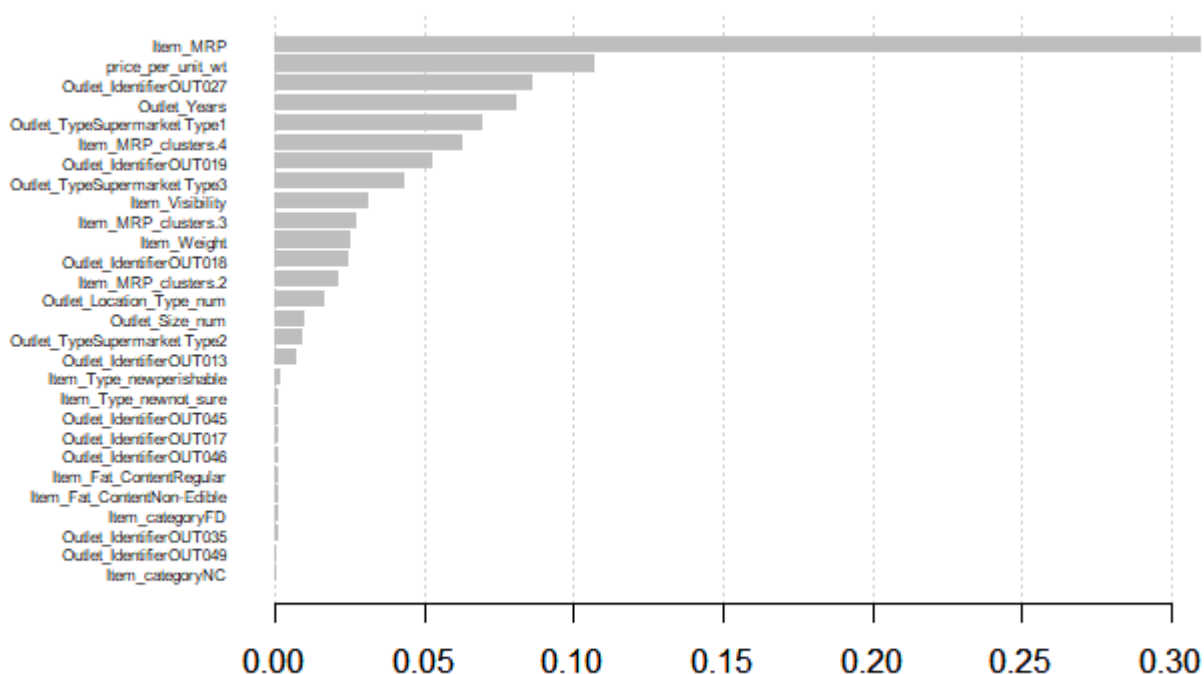
```
xgb_model = xgb.train(data = dtrain, params = param_list, nrounds = 430)
```

Leaderboard score: 1154.70

Variable Importance

[Hide](#)

```
var_imp = xgb.importance(feature_names = setdiff(names(train), c("Item_Identifier", "Item_Outlet_Sales")),
                        model = xgb_model)
xgb.plot.importance(var_imp)
```



As expected Item_MRP is the most important variable in predicting the target variable. New features created by us, like price_per_unit_wt, Outlet_Years, Item_MRP_Clusters, are also among the top most important variables. This is why feature engineering plays such a crucial role in predictive modeling.

End Notes

After trying and testing 5 different algorithms, the best score on the public leaderboard is achieved by XGBoost (1154.70), followed by RandomForest (1157.25). I hope this tutorial was helpful for you in understanding how a machine learning competition is approached and what are the steps one should go through to build a robust model. So, do replicate this analysis and let us know if you face any issues.