

# PreProcessing Data

## What is Data PreProcessing?

In simple words, pre-processing refers to the transformations applied to your data before feeding it to the algorithm. It involves further cleaning of data, data transformation, data scaling and many more things.

For our data, we will deal with the skewness and scale the numerical variables

## Removing 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 its distribution as close to normal distribution as possible. In our data, variables `Item_Visibility` and `price_per_unit_wt` are highly skewed. So, we will treat their skewness with the help of log transformation.

```
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.

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

Splitting the combined data *combi* back to train and test set.

```
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 for test dataset
```

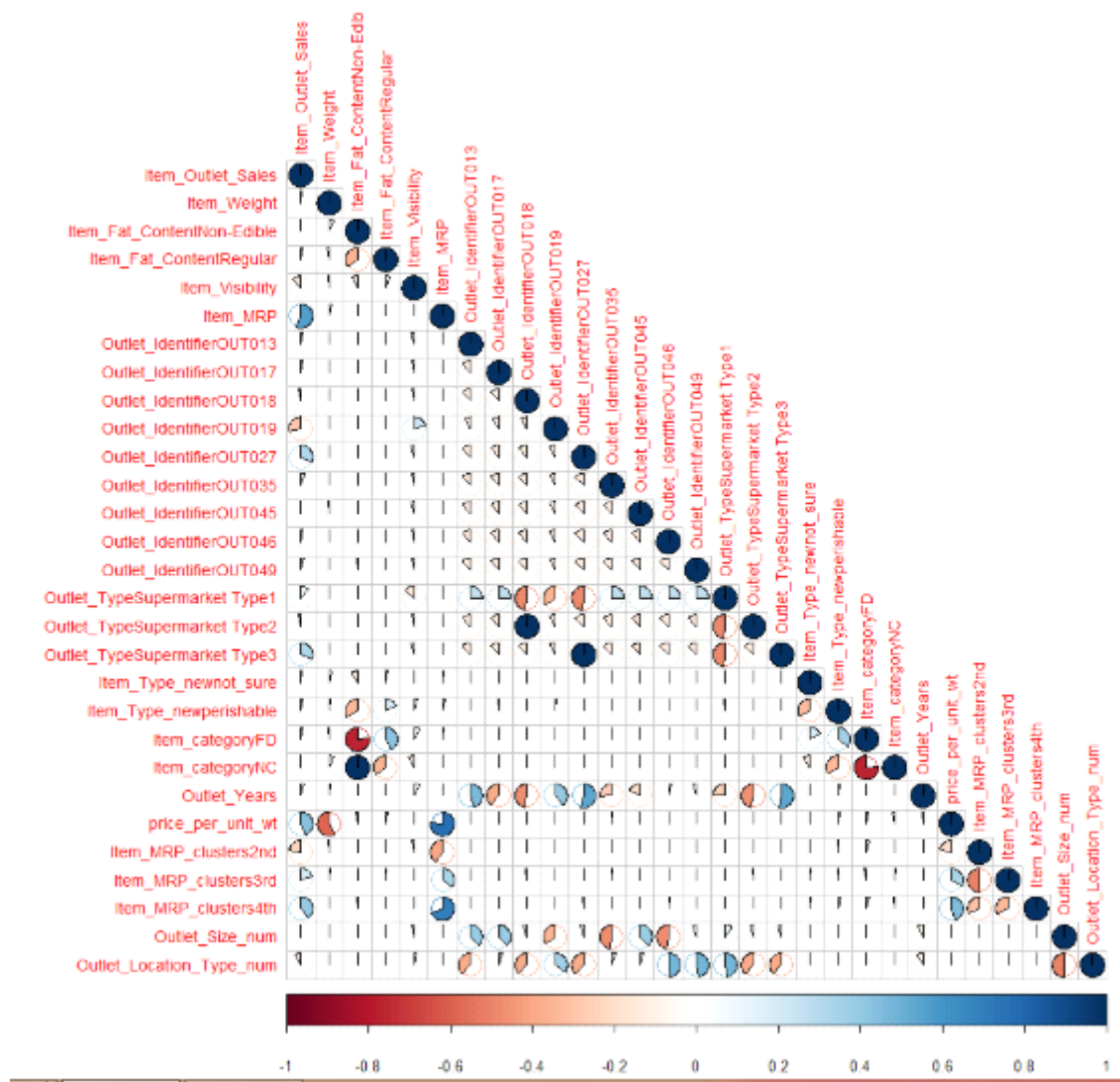
## Correlated Variables

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

1. negative correlation:  $< 0$  and  $\geq -1$
2. positive correlation:  $> 0$  and  $\leq 1$
3. no correlation:  $0$

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

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



The correlation plot above shows correlation between all the possible pairs of variables in out data. The correlation between any two variables is represented by a pie. A blueish pie indicates positive correlation and reddish pie indicates negative correlation. The magnitude of the correlation is denoted by the area covered by the pie.

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.

>