# DATA EXPLORATION

**BANCO SANTANDER**

Machine Learning and innovation

## 1. Steps of Data Exploration and Preparation

Below are the steps involved to understand, clean and prepare your data for building your predictive model:
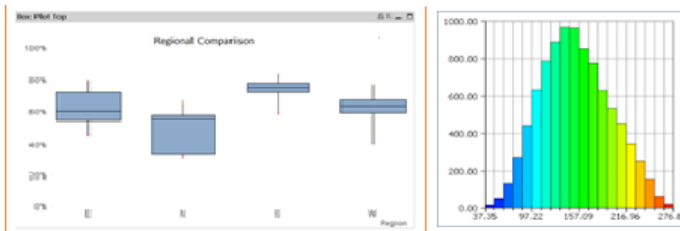
Variable Identification

Identify **Predictor** (Input) and **Target** (output) variables. Next, identify the data type and category of the variables.

Univariate Analysis

We explore variables one by one. Method to perform uni-variate analysis will depend on whether the variable type is categorical or continuous.

**Continuous Variables:-** In case of continuous variables, we need to understand the central tendency and spread of the variable. These are measured using various statistical metrics visualization methods as shown below:

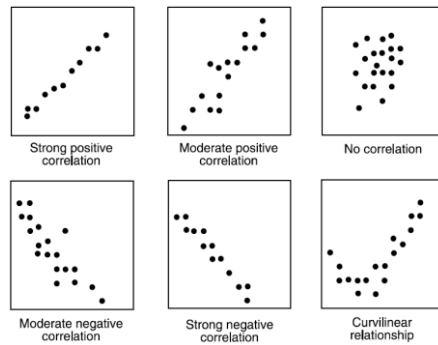| Central Tendency | Measure of Dispersion | Visualization Methods |
|---|---|---|
| Mean | Range | Histogram |
| Median | Quartile | Box Plot |
| Mode | IQR | |
| Min | Variance | |
| Max | Standard Deviation | |
| | Skewness and Kurtosis | |

**Note:** Univariate analysis is also used to highlight missing and outlier values. In the upcoming part of this series, we will look at methods to handle missing and outlier values. To know more about these methods, you can refer course [descriptive statistics from Udacity](descriptive statistics from Udacity).

**Categorical Variables:-** For categorical variables, we'll use frequency table to understand distribution of each category. We can also read as percentage of values under each category. It can be be measured using two metrics, **Count** and **Count%** against each category. Bar chart can be used as visualization.

Bi-variate Analysis

Bi-variate Analysis finds out the relationship between two variables. Here, we look for association and disassociation between variables at a pre-defined significance level. We can perform bi-variate analysis for any combination of categorical and continuous variables. The combination can be: Categorical & Categorical, Categorical & Continuous and Continuous & Continuous.

- **Continuous & Continuous:** While doing bi-variate analysis between two continuous variables, we should look at scatter plot. It is a nifty way to find out the relationship between two variables. The pattern of scatter plot indicates the relationship between variables. The relationship can be linear or non-linear.

Scatter plot shows the relationship between two variable but does not indicates the strength of relationship amongst them. To find the strength of the relationship, we use Correlation. Correlation varies between -1 and +1.
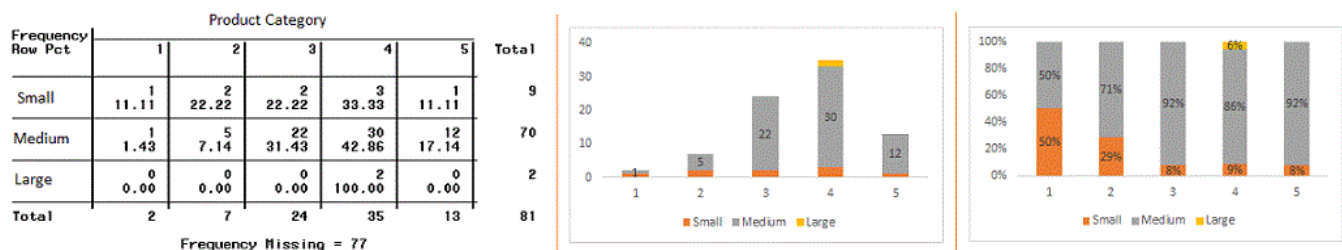
- -1: perfect negative linear correlation

- +1:perfect positive linear correlation and

- 0: No correlation

Correlation can be derived using following formula:

$$\text{Correlation} = \text{Covariance}(X,Y) / \text{SQRT}(\text{Var}(X) * \text{Var}(Y))$$

Categorical & Categorical: To find the relationship between two categorical variables, we can use following methods:

- **Two-way table:** We can start analyzing the relationship by creating a two-way table of count and count%. The rows represents the category of one variable and the columns represent the categories of the other variable. We show count or count% of observations available in each combination of row and column categories.

- **Stacked Column Chart:** This method is more of a visual form of Two-way table.



- **Chi-Square Test:** This test is used to derive the statistical significance of relationship between the variables. Also, it tests whether the evidence in the sample is strong enough to generalize that the relationship for a larger population as well. Chi-square is based on the difference between the expected and observed frequencies in one or more categories in the two-way table. It returns probability for the computed chi-square distribution with the degree of freedom.

Probability of 0: It indicates that both categorical variable are dependent

Probability of 1: It shows that both variables are independent.

Probability less than 0.05: It indicates that the relationship between the variables is significant at 95% confidence.

Statistical Measures used to analyze the power of relationship are:

- Cramer's V for Nominal Categorical Variable

- Mantel-Haenszed Chi-Square for ordinal categorical variable.

Categorical & Continuous: While exploring relation between categorical and continuous variables, we can draw box plots for each level of categorical variables. If levels are small in number, it will not show the statistical significance. To look at the statistical significance we can perform Z-test, T-test or ANOVA.

- **Z-Test/ T-Test**: Either test assess whether mean of two groups are statistically different from each other or not.

- If the probability of Z is small then the difference of two averages is more significant. The T-test is very similar to Z-test but it is used when number of observation for both categories is less than 30.

## 2. Missing Value Treatment

## Why my data has missing values?

- **Data Extraction**

- **Data Collection:** These errors occur at time of data collection and are harder to correct. They can be categorized in four types:

 I. Missing completely at random: This is a case when the probability of missing variable is same for all observations.

 II. Missing at random: This is a case when variable is missing at random and missing ratio varies for different values / level of other input variables.

 III. Missing that depends on unobserved predictors: This is a case when the missing values are not random and are related to the unobserved input variable.

 IV. Missing that depends on the missing value itself: This is a case when the probability of missing value is directly correlated with missing value itself.

## Which are the methods to treat missing values ?

Deletion It is of two types: List Wise Deletion and Pair Wise Deletion.

- In list wise deletion, we delete observations where any of the variable is missing. Simplicity is one of the major advantage of this method, but this method reduces the power of model because it reduces the sample size.

- In pair wise deletion, we perform analysis with all cases in which the variables of interest are present. Advantage of this method is, it keeps as many cases available for analysis. One of the disadvantage of this method, it uses different sample size for different variables.



| List wise deletion | | | | Pair wise deletion | | |
|---|---|---|---|---|---|---|
| Gender | Manpower | Sales | | Gender | Manpower | Sales |
| M | 25 | 343 | | M | 25 | 343 |
| F | . | 280 | | F | . | 280 |
| M | 33 | 332 | | M | 33 | 332 |
| M | . | 272 | | M | . | 272 |
| F | 25 | . | | F | 25 | . |
| M | 29 | 326 | | M | 29 | 326 |
| | 26 | 259 | | | 26 | 259 |
| M | 32 | 297 | | M | 32 | 297 |

*Deletion methods are used when the nature of missing data is "Missing completely at random" else non random missing values can bias the model output.*

**Mean/ Mode/ Median Imputation**: Imputation is a method to fill in the missing values with estimated ones.

It consists of replacing the missing data for a given attribute by the mean or median (quantitative attribute) or mode (qualitative attribute) of all known values of that variable. It can be of two types:

- **Generalized Imputation:** In this case, we calculate the mean or median for all non missing values of that variable then replace missing value with mean or median.

- **Similar case Imputation:** In this case, we calculate average for gender "**Male"** (29.75) and "**Female**" (25) individually of non missing values then replace the missing value based on gender.

**Prediction Model**: Prediction model is one of the sophisticated method for handling missing data. Here, we create a predictive model to estimate values that will substitute the missing data. In this case, we divide our data set into two sets: One set with no missing values for the variable and another one with missing values. First data set become training data set of the model while second data set with missing values is test data set and variable with missing values is treated as target variable. Next, we create a model to predict target variable based on other attributes of the training data set and populate missing values of test data set.We can use regression, ANOVA, Logistic regression and various modeling technique to perform this. There are 2 drawbacks for this approach:

1. The model estimated values are usually more well-behaved than the true values

2. If there are no relationships with attributes in the data set and the attribute with missing values, then the model will not be precise for estimating missing values.

**KNN Imputation:** In this method of imputation, the missing values of an attribute are imputed using the given number of attributes that are most similar to the attribute whose values are missing. The similarity of two attributes is determined using a distance function. It is also known to have certain advantage & disadvantages.

Advantages:

- k-nearest neighbour can predict both qualitative & quantitative attributes

- Creation of predictive model for each attribute with missing data is not required

- Attributes with multiple missing values can be easily treated

- Correlation structure of the data is taken into consideration

Disadvantage:

- KNN algorithm is very time-consuming in analyzing large database. It searches through all the dataset looking for the most similar instances.

- Choice of k-value is very critical. Higher value of k would include attributes which are significantly different from what we need whereas lower value of k implies missing out of significant attributes.

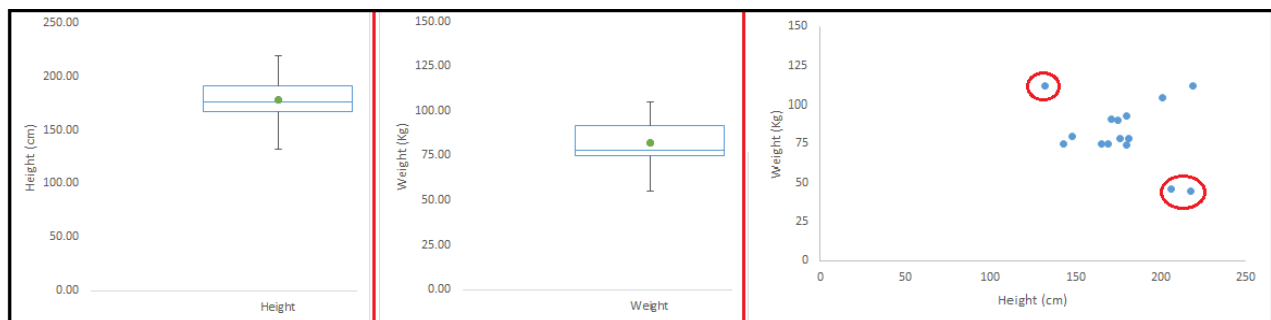3. Techniques of Outlier Detection and Treatment

## What is an Outlier?

Outlier is an observation that appears far away and diverges from an overall pattern in a sample.

## What are the types of Outliers?

Outlier can be of two types: **Univariate** and **Multivariate**.

Multi-variate outliers are outliers in an n-dimensional space. In order to find them, you have to look at distributions in multi-dimensions.

Let us understand this with an example. Let us say we are understanding the relationship between height and weight. Below, we have univariate and bivariate distribution for Height, Weight. Take a look at the box plot. We do not have any outlier (above and below 1.5*IQR, most common method). Now look at the scatter plot. Here, we have two values below and one above the average in a specific segment of weight and height.



## What causes Outliers?

Whenever we come across outliers, the ideal way to tackle them is to find out the reason of having these outliers. The method to deal with them would then depend on the reason of their occurrence. Causes of outliers can be classified in two broad categories:

1. **Artificial (Error) / Non-natural**

2. **Natural**.

Let's understand various types of outliers in more detail:

- Data Entry Errors: Human errors.

- Measurement Error.

- Experimental Error.

- Intentional Outlier.

- Data Processing Error: Whenever we perform data mining.

- Sampling error

- Natural Outlier: When an outlier is not artificial (due to error), it is a natural outlier.

## What is the impact of Outliers on a dataset?

Outliers can drastically change the results of the data analysis and statistical modeling. There are numerous unfavourable impacts of outliers in the data set:

- It increases the error variance and reduces the power of statistical tests

- If the outliers are non-randomly distributed, they can decrease normality

- They can bias or influence estimates that may be of substantive interest

- They can also impact the basic assumption of Regression, ANOVA and other statistical model assumptions.

| Without Outlier | With Outlier |
|---|---|
| 4, 4, 5, 5, 5, 5, 6, 6, 6, 7, 7 | 4, 4, 5, 5, 5, 5, 6, 6, 6, 7, 7,300 |
| Mean = 5.45 | Mean = 30.00 |
| Median = 5.00 | Median = 5.50 |
| Mode = 5.00 | Mode = 5.00 |
| Standard Deviation = 1.04 | Standard Deviation = 85.03 |

## How to detect Outliers?

Most commonly used method to detect outliers is visualization. We use various visualization methods, like **Box-plot**, **Histogram**, **Scatter Plot** (above, we have used box plot and scatter plot for visualization). Some analysts also various thumb rules to detect outliers. Some of them are:

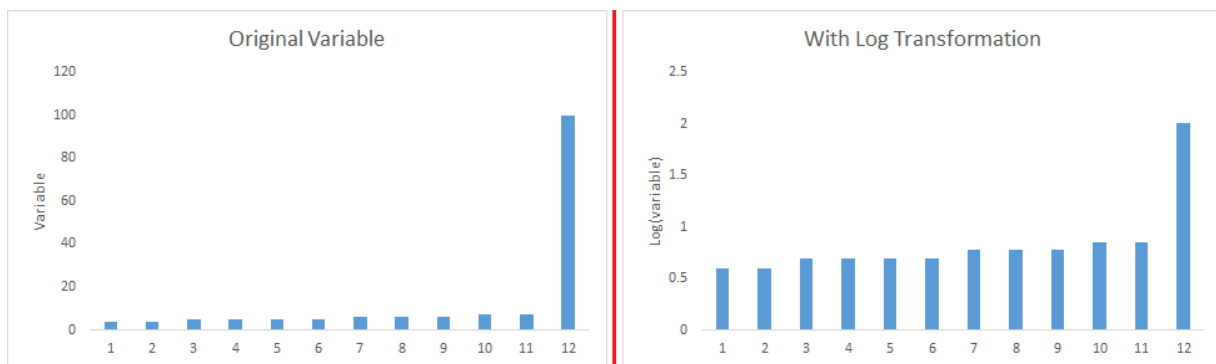- Any value, which is beyond the range of -1.5 x IQR to 1.5 x IQR

- Use capping methods. Any value which out of range of 5th and 95th percentile can be considered as outlier

- Data points, three or more standard deviation away from mean are considered outlier

- Outlier detection is merely a special case of the examination of data for influential data points and it also depends on the business understanding

- Bivariate and multivariate outliers are typically measured using either an index of influence or leverage, or distance. Popular indices such as Mahalanobis' distance and Cook's *D* are frequently used to detect outliers.

# How to remove Outliers?

Most of the ways to deal with outliers are similar to the methods of missing values like deleting observations, transforming them, binning them, treat them as a separate group, imputing values and other statistical methods. Here, we will discuss the common techniques used to deal with outliers:

**Deleting observations:** We delete outlier values if it is due to data entry error, data processing error or outlier observations are very small in numbers. We can also use trimming at both ends to remove outliers.

**Transforming and binning values:** Transforming variables can also eliminate outliers. Natural log of a value reduces the variation caused by extreme values. Binning is also a form of variable transformation. Decision Tree algorithm allows to deal with outliers well due to binning of variable. We can also use the process of assigning weights to different observations.



**Imputing:** Like imputation of missing values, we can also impute outliers. We can use mean, median, mode imputation methods. Before imputing values, we should analyse if it is natural outlier or artificial. If it is artificial, we can go with imputing values. We can also use statistical model to predict values of outlier observation and after that we can impute it with predicted values.

**Treat separately:** If there are significant number of outliers, we should treat them separately in the statistical model. One of the approach is to treat both groups as two different groups and build individual model for both groups and then combine the output.

## 4. The Art of Feature Engineering

### What is Feature Engineering?

Feature engineering is the science (and art) of extracting more information from existing data. You are not adding any new data here, but you are actually making the data you already have more useful.

### What is the process of Feature Engineering ?

You perform feature engineering once you have completed the first 5 steps in data exploration – Variable Identification, Univariate, Bivariate Analysis, Missing Values Imputation and Outliers Treatment. Feature engineering itself can be divided in 2 steps:

- Variable transformation.

- Variable / Feature creation.

These two techniques are vital in data exploration and have a remarkable impact on the power of prediction. Let's understand each of this step in more details.

### What is Variable Transformation?

In data modelling, transformation refers to the replacement of a variable by a function. For instance, replacing a variable x by the square / cube root or logarithm x is a transformation. In other words, transformation is a process that changes the distribution or relationship of a variable with others.

### When should we use Variable Transformation?

Below are the situations where variable transformation is a requisite:

- When we want to **change the scale** of a variable or standardize the values of a variable for better understanding. While this transformation is a must if you have data in different scales, this transformation does not change the shape of the variable distribution

- When we can **transform complex non-linear relationships into linear relationships**. Existence of a linear relationship between variables is easier to comprehend compared to a non-linear or curved relation. Transformation helps us to convert a non-linear relation into linear relation. Scatter plot can be used to find the relationship between two continuous variables. These transformations also improve the prediction. Log transformation is one of the commonly used transformation technique used in these situations.

**Symmetric distribution is preferred over skewed distribution** as it is easier to interpret and generate inferences. Some modeling techniques requires normal distribution of variables. So, whenever we have a skewed distribution, we can use transformations which reduce skewness. For right skewed distribution, we take square / cube root or logarithm of variable and for left skewed, we take square / cube or exponential of variables.

*Variable Transformation is also done from an **implementation point of view.***

## What are the common methods of Variable Transformation?

There are various methods used to transform variables. As discussed, some of them include square root, cube root, logarithmic, binning, reciprocal and many others. Let's look at these methods in detail by highlighting the pros and cons of these transformation methods.

- **Logarithm:** Log of a variable is a common transformation method used to change the shape of distribution of the variable on a distribution plot. It is generally used for reducing right skewness of variables. Though, It can't be applied to zero or negative values as well.

- **Square / Cube root:** The square and cube root of a variable has a sound effect on variable distribution. However, it is not as significant as logarithmic transformation. Cube root has its own advantage. It can be applied to negative values including zero. Square root can be applied to positive values including zero.

- **Binning:** It is used to categorize variables. It is performed on original values, percentile or frequency. Decision of categorization technique is based on business understanding. For example, we can categorize income in three categories, namely: High, Average and Low. We can also perform co-variate binning which depends on the value of more than one variables.

## What is Feature / Variable Creation & its Benefits?

Feature / Variable creation is a process to generate a new variables / features based on existing variable(s). For example, say, we have date(dd-mm-yy) as an input variable in a data set. We can generate new variables like day, month, year, week, weekday that may have better relationship with target variable. This step is used to highlight the hidden relationship in a variable:

| Emp_Code | Gender | Date | New_Day | New_Month | New_Year |
|----------|--------|-----------|---------|-----------|----------|
| A001 | Male | 21-Sep-11 | 21 | 9 | 2011 |
| A002 | Female | 27-Feb-13 | 27 | 2 | 2013 |
| A003 | Female | 14-Nov-12 | 14 | 11 | 2012 |
| A004 | Male | 07-Apr-13 | 7 | 4 | 2013 |
| A005 | Female | 21-Jan-11 | 21 | 1 | 2011 |
| A006 | Male | 26-Apr-13 | 26 | 4 | 2013 |
| A007 | Male | 15-Mar-12 | 15 | 3 | 2012 |

There are various techniques to create new features. Let's look at the some of the commonly used methods:

- **Creating derived variables:** This refers to creating new variables from existing variable(s) using set of functions or different methods.

- **Creating dummy variables:** One of the most common application of dummy variable is to convert categorical variable into numerical variables. Dummy variables are also called Indicator Variables. It is useful to take categorical variable as a predictor in statistical models. Categorical variable can take values 0 and 1.

| Emp_Code | Gender | Var_Male | Var_Female |
|----------|--------|----------|------------|
| A001 | Male | 1 | 0 |
| A002 | Female | 0 | 1 |
| A003 | Female | 0 | 1 |
| A004 | Male | 1 | 0 |
| A005 | Female | 0 | 1 |
| A006 | Male | 1 | 0 |
| A007 | Male | 1 | 0 |