

x

-

(http://play.google.com/store/apps/details?id=com.analyticsvidhya.android)



Analytics Vidhya

Learn everything about analytics

www.analyticsvidhya.com/blog/

(https://trainings.analyticsvidhya.com/courses/course-v1:AnalyticsVidhya+DS101+2018T2/about?utm_source=AVtopblogBanner)

(https://datahack.analyticsvidhya.com/contest/american-



www.analyticsvidhya.com/blog/category/business-analytics/

Reply.

Data Exploration

www.analyticsvidhya.com/blog/author/sunil-ray/, JANUARY 10, 2016

(https://trainings.analyticsvidhya.com/courses/course-v1:AnalyticsVidhya+CVDL101+CVDL101_T1/about?utm_source=CV101AVBlogBanner&utm_medium=Stickybanner2utm_campaign=CV101banner)

rescue.

I can confidently say this, because I've been through such situations, a lot.

I have been a Business Analytics professional for close to three years now. In my initial days, one of my mentor suggested me to spend significant time on exploration and analyzing data. Following his advice has served me well.

I've created this tutorial to help you understand the underlying techniques of data exploration. As always, I've tried my best to explain these concepts in the simplest manner. For better understanding, I've taken up few examples to demonstrate the complicated concepts.

^

Subscribe!



(uploads/2016/01/de.jpg)

(<https://trainings.analyticsvidhya.com/courses/course-v1:AnalyticsVidhya+DS101+2018T2/about?>



ration

required ?

missing value ?

Treatment

(https://trainings.analyticsvidhya.com/courses/course-v1:AnalyticsVidhya+CVDL101+CVDL101_T1/about?

v1:AnalyticsVidhya+CVDL101+CVDL101_T1/about?

How to detect outlier ?

utm_source=CV101AVBlogBanner&utm_medium=Stickybanner2utm_campaign=CV101banner)

How to remove outlier ?

4. The Art of Feature Engineering

- What is Feature Engineering ?
- What is the process of Feature Engineering ?
- What is Variable Transformation ?
- When should we use variable transformation ?
- What are the common methods of variable transformation ?
- What is feature variable creation and its benefits ?

Let's get started.

Subscribe!

1. Steps of Data Exploration and Preparation

Remember the quality of your inputs decide the quality of your output. So, once you have got your business problem, you will spend a lot of time and efforts here. With my personal estimate, data exploration will take up to 70% of your total project time.



Next, you need to clean and prepare your data for building your predictive model:

([https://trainings.analyticsvidhya.com/courses/course-v1:AnalyticsVidhya+DS101+2018T2/about?](https://trainings.analyticsvidhya.com/courses/course-v1:AnalyticsVidhya+DS101+2018T2/about?utm_source=AnalyticsVidhya&utm_medium=BlogBanner&utm_campaign=CV101banner)

7 multiple times before we come up with our refined model.



Let's understand this step more clearly by taking an example. (output) variables. Next, identify the data type and category of the

Let's understand this step more clearly by taking an example. ([https://trainings.analyticsvidhya.com/courses/course-](https://trainings.analyticsvidhya.com/courses/course-v1:AnalyticsVidhya+CVDL101+CVDL101_T1/about?utm_source=AnalyticsVidhya&utm_medium=BlogBanner&utm_campaign=CV101banner)

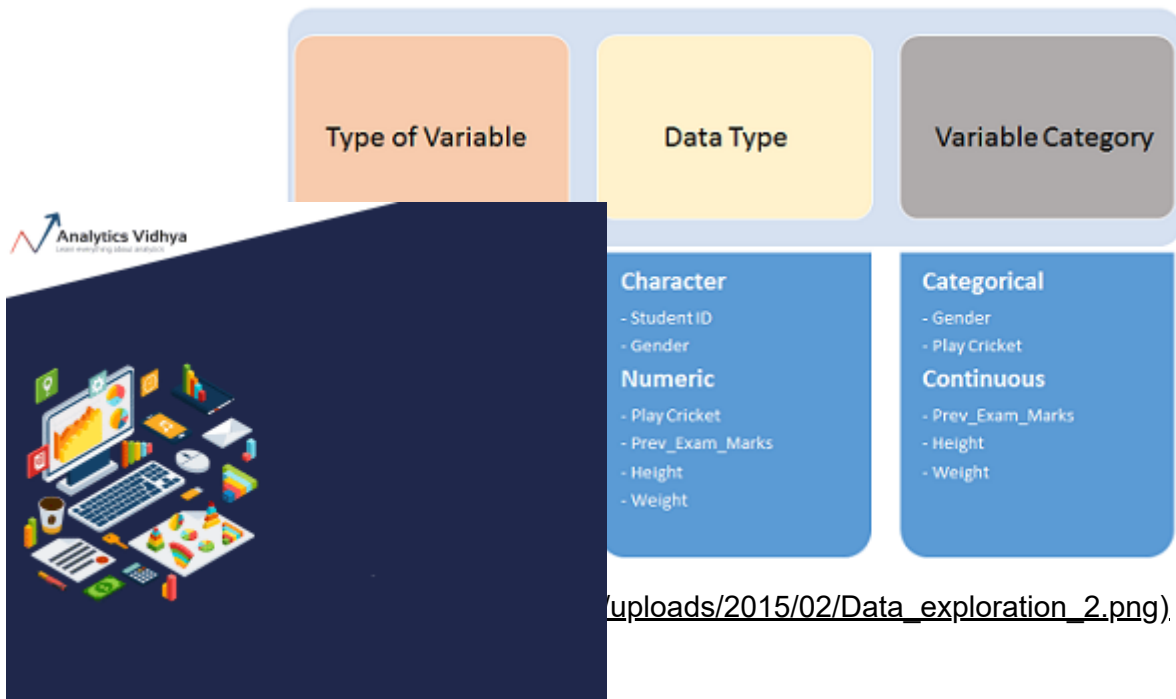
[v1:AnalyticsVidhya+CVDL101+CVDL101_T1/about?](https://trainings.analyticsvidhya.com/courses/course-v1:AnalyticsVidhya+CVDL101+CVDL101_T1/about?utm_source=AnalyticsVidhya&utm_medium=BlogBanner&utm_campaign=CV101banner)

Example:- Suppose, we want to predict, whether the students will play cricket or not (refer below data set). Here you need to identify predictor variables, target variable, data type of variables and category of variables.

Student_ID	Gender	Prev_Exam_Marks	Height (cm)	Weight Caregory (kgs)	Play Cricket
S001	M	65	178	61	1
S002	F	75	174	56	0
S003	M	45	163	62	1
S004	M	57	175	70	0
S005	F	59	162	67	0

(https://www.analyticsvidhya.com/wp-content/uploads/2015/02/Data_exploration_11.png)Below, the variables have been defined in different category:

Subscribe!



([https://www.analyticsvidhya.com/wp-content/uploads/2015/02/Data_exploration_2.png](#)).

([https://trainings.analyticsvidhya.com/courses/course-v1:AnalyticsVidhya+DS101+2018T2/about?](https://trainings.analyticsvidhya.com/courses/course-v1:AnalyticsVidhya+DS101+2018T2/about?utm_source=CV101+AVBlogBanner&utm_medium=Stickybanner2utm_campaign=CV101banner)

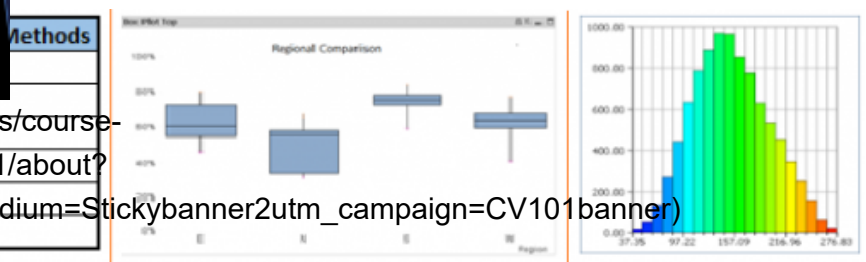


([https://trainings.analyticsvidhya.com/courses/course-v1:AnalyticsVidhya+CVDL101+CVDL101_T1/about?](https://trainings.analyticsvidhya.com/courses/course-v1:AnalyticsVidhya+CVDL101+CVDL101_T1/about?utm_source=CV101+AVBlogBanner&utm_medium=Stickybanner2utm_campaign=CV101banner)

[utm_source=CV101+AVBlogBanner&utm_medium=Stickybanner2utm_campaign=CV101banner](#))

one. Method to perform uni-variate analysis will depend on continuous. Let's look at these methods and statistical measures individually:

For continuous variables, we need to understand the central tendency and using various statistical metrics visualization methods as shown



(https://www.analyticsvidhya.com/wp-content/uploads/2015/02/Data_exploration_31.png) **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](#) (<https://www.udacity.com/course/ud827>).

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 measured using two metrics, **Count** and **Count%** against each category. Bar chart can be used as visualization.

Subscribe!

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, & Continuous. Different methods are used to tackle these



([https://trainings.analyticsvidhya.com/courses/course-v1:AnalyticsVidhya+DS101+2018T2/about?](https://trainings.analyticsvidhya.com/courses/course-v1:AnalyticsVidhya+DS101+2018T2/about?utm_source=CV101AVBlogBanner&utm_medium=StickyBanner&utm_campaign=CV101Banner)



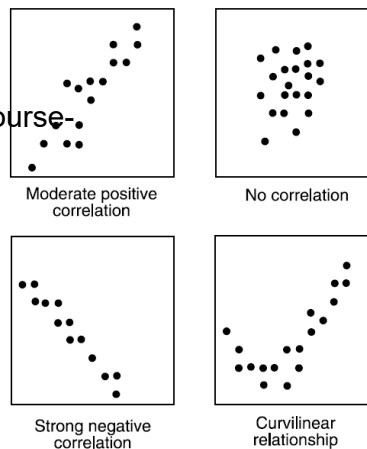
([https://trainings.analyticsvidhya.com/courses/course-v1:AnalyticsVidhya+CVDL101+CVDL101_T1/about?](https://trainings.analyticsvidhya.com/courses/course-v1:AnalyticsVidhya+CVDL101+CVDL101_T1/about?utm_source=CV101AVBlogBanner&utm_medium=StickyBanner&utm_campaign=CV101Banner)

(https://www.analyticsvidhya.com/wp-content/uploads/2015/02/Data_exploration_4.png)
 utm_source=CV101AVBlogBanner&utm_medium=StickyBanner&utm_campaign=CV101Banner)

- 0: No correlation

in detail:

For bi-variate analysis between two continuous variables, we should look at the relationship between two variables. The pattern of scatter plots varies. The relationship can be linear or non-linear.



(https://www.analyticsvidhya.com/wp-content/uploads/2015/02/Data_exploration_4.png) Scatter plot shows the relationship between two variables. It does not indicate the strength of relationship amongst them. To measure the strength of relationship, we use correlation. Correlation varies between -1 and +1.

Correlation can be derived using following formula:

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

Various tools have function or functionality to identify correlation between variables. In Excel, function CORREL() is used to return the correlation between two variables and SAS uses procedure PROC CORR to identify the correlation. These function returns Pearson Correlation value to identify the relationship between two variables:

Subscribe!

X	65	72	78	65	72	70	65	68
Y	72	69	79	69	84	75	60	73

Metrics	Formula	Value
Covariance (X, Y)	$COVAR(56:16, 57:17)$	18.77
Standard Deviation (X)		18.48
Standard Deviation (Y)		45.23
Correlation Coefficient		0.65

Analytics Vidhya
LEARN. GROW. SHINE. ABOUT ANALYTICS



(https://www.analyticsvidhya.com/wp-content/uploads/2015/02/Data_exploration_51.png)

relationship(0.65) between two variables X and Y.

relationship between two categorical variables, we can use following

Methods.

(https://trainings.analyticsvidhya.com/courses/course-

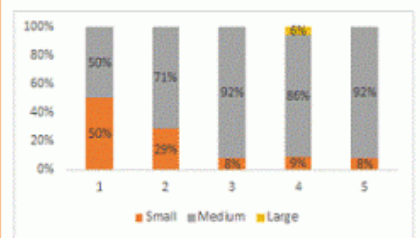
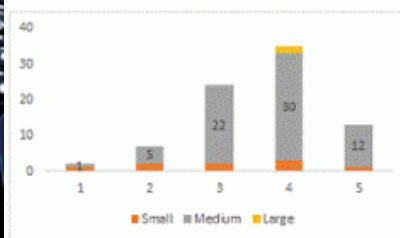
v1:AnalyticsVidhya+DS101+2018T2/about?

Analytics Vidhya
LEARN. GROW. SHINE. ABOUT ANALYTICS



ing the relationship by creating a two-way table of count and category of one variable and the columns represent the categories of count% of observations available in each combination of row and

is more of a visual form of Two-way table.



(https://trainings.analyticsvidhya.com/courses/course-

(https://www.analyticsvidhya.com/wp-content/uploads/2015/02/Data_exploration_6.gif)

V1:AnalyticsVidhya+CVDL101+CVDL101-T1/about?

utm_source=CV101AVBlogBanner&utm_medium=Stickybanner2utm_campaign=CV101banner)

- **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.

Subscribe!

Probability less than 0.05: It indicates that the relationship between the variables is significant at 95% confidence. The chi-square test statistic for a test of independence of two categorical variables is found by:

$$X^2 = \sum (O - E)^2 / E$$

content/uploads/2015/02/Data_exploration_7.png)where O is the observed frequency and E is the expected frequency under the null hypothesis and computed

$$= \frac{\text{row total} \times \text{column total}}{\text{sample size}}$$

content/uploads/2015/02/Data_exploration_8.png).

Expected count for product category 1 to be of small size is 0.22. It is 0.22 times the column total for Product category (2) then dividing by 0.22 is conducted for each cell. Statistical Measures used to analyze



(https://trainings.analyticsvidhya.com/courses/course-v1:AnalyticsVidhya+DS101+2018T2/about?utm_source=AnalyticsVidhya+DS101+2018T2+Banner&utm_medium=Banner&utm_campaign=AnalyticsVidhya+DS101+2018T2+Banner)



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

utm_source=CV101AVBlogBanner&utm_medium=StickyBanner2utm_campaign=CV101banner)

$$z = \frac{|x_1 - x_2|}{\sqrt{\frac{s_1^2}{n_1} + \frac{s_2^2}{n_2}}}$$

(https://www.analyticsvidhya.com/wp-content/uploads/2015/02/ztestformula1.jpg)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.

Subscribe!

$$t = \frac{\bar{X}_1 - \bar{X}_2}{\sqrt{S^2 \left(\frac{1}{N_1} + \frac{1}{N_2} \right)}}$$

$$S^2 = \frac{(N_1 - 1)S_1^2 + (N_2 - 1)S_2^2}{N_1 + N_2 - 2}$$

where:

- \bar{X}_1, \bar{X}_2 : Averages
- S_1^2, S_2^2 : Variances
- N_1, N_2 : Counts
- t : has t distribution with $N_1 + N_2 - 2$ degree of freedom



(<https://content/uploads/2015/02/ttest.png>)

page of more than two groups is statistically different.

of five different exercises. For this, we recruit 20 men and (20 groups). Their weights are recorded after a few weeks. We need to know if the weight gain on them is significantly different or not. This can be done by using a t-test on each.

Three stages of Data Exploration, Variable Identification, Uni-Variate and Bi-Variate analysis. We also looked at various statistical and visual methods to identify the variables. [https://trainings.analyticsvidhya.com/courses/course-v1:AnalyticsVidhya+DS101+2018T2/about?](https://trainings.analyticsvidhya.com/courses/course-v1:AnalyticsVidhya+DS101+2018T2/about?utm_source=AnalyticsVidhya+DS101+2018T2+T1/about?utm_medium=Stickypbanner2&utm_campaign=CV101banner)



ing values Treatment. More importantly, we will also look at why identifying them is necessary.

Missing data in the training data set can reduce the power / fit of a model or can lead to a biased model because we have not analysed the behavior and relationship with other variables correctly. It can lead to wrong prediction or classification. [https://trainings.analyticsvidhya.com/courses/course-v1:AnalyticsVidhya+CVDL101+CVDL101_T1/about?](https://trainings.analyticsvidhya.com/courses/course-v1:AnalyticsVidhya+CVDL101+CVDL101_T1/about?utm_source=CV101AVBlogBanner&utm_medium=Stickypbanner2&utm_campaign=CV101banner)

Missing data in the training data set can reduce the power / fit of a model or can lead to a biased model because we have not analysed the behavior and relationship with other variables correctly. It can lead to wrong prediction or classification.

Subscribe!

Name	Weight	Gender	Play Cricket/ Not
Mr. Amit	58	M	Y
Mr. Anil	61	M	Y
Miss Swati	58	F	N
Miss Richa	55		Y
Mr. Steve	55	M	N

Name	Weight	Gender	Play Cricket/ Not
Mr. Amit	58	M	Y
Mr. Anil	61	M	Y
Miss Swati	58	F	N
Miss Richa	55	F	Y
Mr. Steve	55	M	N
Miss Reena	64	F	Y
Miss Rashmi	57	F	Y
Mr. Kunal	57	M	N

Gender	#Students	#Play Cricket	%Play Cricket
F	4	3	75%
M	4	2	50%

Analytics Vidhya
Learn something about analytics



[https://trainings.analyticsvidhya.com/courses/course-v1:AnalyticsVidhya+DS101+2018T2/about?](https://trainings.analyticsvidhya.com/courses/course-v1:AnalyticsVidhya+DS101+2018T2/about?utm_source=CV10+AVBlogBanner&utm_medium=StickyBanner&utm_campaign=CV10Banner)

Analytics Vidhya
Learn something about analytics



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

- **Missing completely at random:** This is a case when the probability of missing variable is same for all observations. For example: respondents of data collection process decide that they will declare their earning after tossing a fair coin. If an head occurs, respondent declares his / her earnings & vice versa. Here each observation has equal chance of missing value.
- **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. For example: We are collecting data for age and female has higher missing value compare to male.
- **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. For example: In a medical study, if a particular diagnostic causes discomfort, then there is higher chance of drop out from the study. This missing value is not at random unless we have included "discomfort" as an input variable for all patients.

Subscribe!

- **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. For example: People with higher or lower income are likely to provide non-response to their earning.



Missing values ?

Deletion and Pair Wise Deletion.

observations where any of the variable is missing. Simplicity is one method, but this method reduces the power of model because it

on analysis with all cases in which the variables of interest are good is, it keeps as many cases available for analysis. One of the cases different sample size for different variables.

(https://training.analyticsvidhya.com/wp-content/uploads/2015/02/Data_Exploration_2_2.png)
v1:AnalyticsVidhya+DS101+2018T2/about?



Manpower	Sales
25	343
.	280
33	332
.	272
25	.
29	326
26	259
32	297

Pair wise deletion

Gender	Manpower	Sales
M	25	343
F	.	280
M	33	332
M	.	272
F	25	.
M	29	326
.	26	259
M	32	297

(https://training.analyticsvidhya.com/wp-content/uploads/2015/02/Data_Exploration_2_2.png)

the nature of missing data is “**Missing completely at random**”

(https://training.analyticsvidhya.com/wp-content/uploads/2015/02/Data_Exploration_2_2.png)

v1:AnalyticsVidhya+CVDL101+CVDL101_T1/about?

utmsource=AnalyticsVidhya&utm_medium=organic&utm_campaign=Data_Exploration_2_2

Mean / Mode / Median Imputation: Imputation is a method to fill in the missing values with estimated ones. The objective is to employ known relationships that can be identified in the valid values of the data set to assist in estimating the missing values. Mean / Mode / Median imputation is one of the most frequently used methods. 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. Like in above table, variable “**Manpower**” is missing so we take average of all non missing values of “**Manpower**” (28.33) and then replace missing value with it.
- **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. For “**Male**”, we will replace missing values of manpower with 29.75 and for “**Female**” with 25!

Subscribe!

3. 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



([https://trainings.analyticsvidhya.com/courses/course-v1:AnalyticsVidhya+DS101+2018152/about?](https://trainings.analyticsvidhya.com/courses/course-v1:AnalyticsVidhya+DS101+2018152/about?utm_source=CV101AVBlogBanner&utm_medium=Stickybanner&utm_campaign=CV101Banner)



([https://trainings.analyticsvidhya.com/courses/course-v1:AnalyticsVidhya+DS101+2018152/about?](https://trainings.analyticsvidhya.com/courses/course-v1:AnalyticsVidhya+DS101+2018152/about?utm_source=CV101AVBlogBanner&utm_medium=Stickybanner&utm_campaign=CV101Banner)

After dealing with missing values, the next task is to deal with outliers. Often, we tend to neglect outliers while building models. This is a discouraging practice. Outliers tend to make your data skewed and reduces accuracy. Let's learn more about outlier treatment.

3. Techniques of Outlier Detection and Treatment

What is an Outlier?

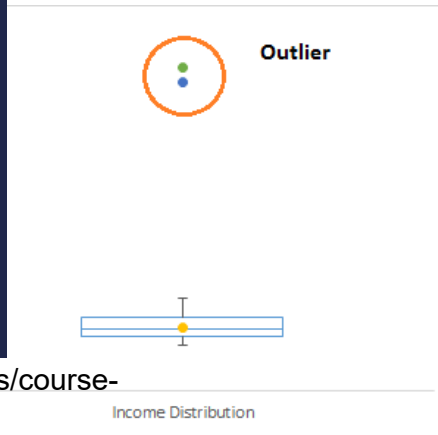
Subscribe!

Outlier is a commonly used terminology by analysts and data scientists as it needs close attention else it can result in wildly wrong estimations. Simply speaking, Outlier is an observation that appears far away and diverges from an overall pattern in a sample.

Let's take an example, we do customer profiling and find out that the average annual income of customers is \$4 million. These two customers have an annual income of \$4 and \$4.2 million. These two customers are outliers from the population. These two observations will be seen as Outliers.



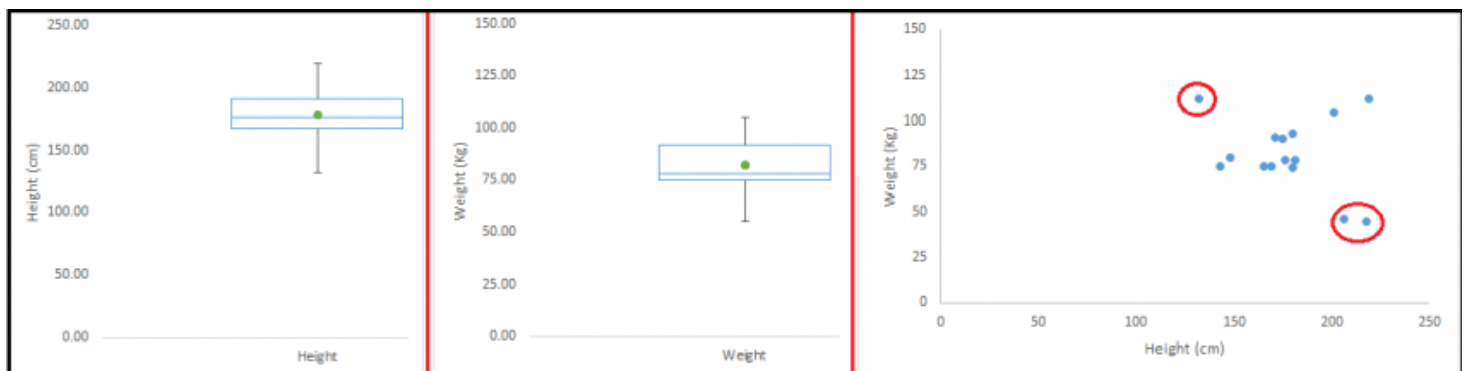
(<https://trainings.analyticsvidhya.com/courses/course-v1:AnalyticsVidhya+DS101+2018T2/about?>



(<https://www.analyticsvidhya.com/wp-content/uploads/2015/02/Outlier.png>)

Univariate and Multivariate. Above, we have discussed the example of univariate distribution. When we look at the distribution of a single variable, it is univariate. In order to find them, you have to look at distributions in multi-

Let's take an example, 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 outliers in univariate distribution. 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.



(https://www.analyticsvidhya.com/wp-content/uploads/2015/02/Outlier_21.png)

Subscribe!

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:



(<https://trainings.analyticsvidhya.com/courses/course-v1:AnalyticsVidhya+DS401012018/about/>)



(<https://trainings.analyticsvidhya.com/courses/course-v1:AnalyticsVidhya+DS401012018/about/>)

• **Natural Outlier:** When an outlier is not artificial (due to error), it is a natural outlier. For instance: In my last assignment with one of the renowned insurance company, I noticed that the performance of top 50 financial advisors was far higher than rest of the population. Surprisingly, it was not due to any error. Hence, whenever we perform any data mining activity with advisors, we used to treat this segment separately.

more detail:

For example: As errors caused during data collection, recording, or entry can lead to outliers. For example: Annual income of a customer is \$100,000. Accidentally, the data entry person enters 1,000,000 in the figure. Now the income becomes \$1,000,000 which is 10 times the outlier value when compared with rest of the population.

Another common source of outliers is measurement error. This is caused when the measurement instrument used turns out to be faulty. For example: There are 10 weighing machines. 9 of them are correct. If a faulty machine is used, the weight measured by people on the faulty machine will be higher / lower than the rest measured on faulty machine can lead to outliers.

Another source of outliers is experimental error. For example: In a 100m sprint of 7 runners, one runner is concentrating on the 'Go' call which caused him to start late. Hence, his total run time can be an outlier.

Outliers are also found in self-reported measures that involves sensitive data. For example: A survey report the amount of alcohol that they consume. Only a fraction of the actual values might look like outliers because rest of the teens

When we perform data mining, we extract data from multiple sources. It is possible that data entry errors may lead to outliers in the dataset.

For example: Suppose we want to measure the height of athletes. By mistake, we include a few basketball players in the sample. This inclusion is likely to cause outliers in the dataset.

(<https://trainings.analyticsvidhya.com/courses/course-v1:AnalyticsVidhya+DS401012018/about/>)

utm_source=CV101AVBlogBanner&utm_medium=StickyBanner&utm_campaign=CV101Banner

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

Subscribe!

- 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.



([https://trainings.analyticsvidhya.com/courses/course-v1:AnalyticsVidhya+DS101+2018T2/about?](https://trainings.analyticsvidhya.com/courses/course-v1:AnalyticsVidhya+DS101+2018T2/about?utm_source=CV101AvBlogBanner&utm_medium=StickyBanner2utm_campaign=CV101banner)



([https://trainings.analyticsvidhya.com/courses/course-v1:AnalyticsVidhya+DS101+2018T2/about?](https://trainings.analyticsvidhya.com/courses/course-v1:AnalyticsVidhya+DS101+2018T2/about?utm_source=CV101AvBlogBanner&utm_medium=StickyBanner2utm_campaign=CV101banner)

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 $-1.5 \times IQR$ to $1.5 \times IQR$
- 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.

an example to check what happens to a data set with and without

With Outlier
4, 4, 5, 5, 5, 5, 6, 6, 6, 7, 7, 300
Mean = 30.00
Median = 5.50
Mode = 5.00
Standard Deviation = 85.03

([uploads/2015/02/Outlier_31.png](#)).

significantly different mean and standard deviation. In the first without the outlier, average soars to 30. This would change the

Subscribe!

- In SAS, we can use PROC Univariate, PROC SGPLOT. To identify outliers and influential observation, we also look at statistical measure like STUDENT, COOKD, RSTUDENT and others.

How to remove Outliers?



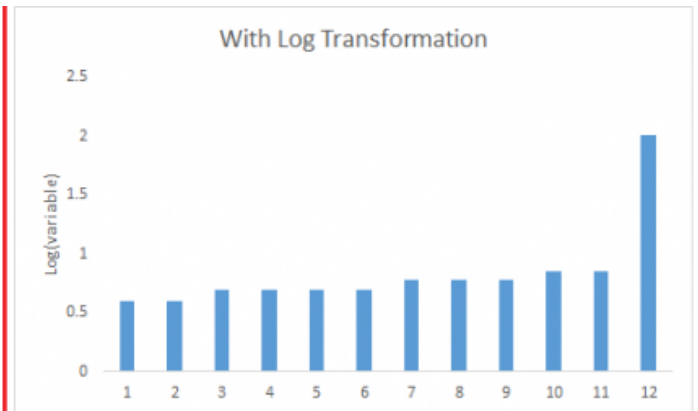
([https://trainings.analyticsvidhya.com/courses/course-v1:AnalyticsVidhya+DS101+2018T2/about?](https://trainings.analyticsvidhya.com/courses/course-v1:AnalyticsVidhya+DS101+2018T2/about?utm_source=CV101AVBlogBanner&utm_medium=Stickybanner2utm_campaign=CV101banner)



Similar to the methods of missing values like deleting observations, treating them as a separate group, imputing values and other statistical techniques used to deal with outliers:

Deleting values if it is due to data entry error, data processing error or other reasons. We can also use trimming at both ends to remove outliers.

Transforming variables can also eliminate outliers. Natural log of a value can be used. Binning is also a form of variable transformation. Decision trees are less sensitive due to binning of variable. We can also use the process of



([uploads/2015/02/Transformation_1.png](#))

Imputing. Like imputation of missing values (<https://www.analyticsvidhya.com/blog/2015/02/7-steps-data-exploration-preparation-building-model-part-2/>), 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.

Till here, we have learnt about steps of data exploration, missing value treatment and techniques of outlier detection and treatment. These 3 stages will make your raw data better in terms of information availability and accuracy. Let's now proceed to the final stage of data exploration. It is Feature Engineering.

Subscribe!

4. The Art of Feature Engineering



of extracting more information from existing data. You are not only making the data you already have more useful.

Let's take an example of foot fall in a shopping mall based on dates. If you try and use the date as a feature, you can extract meaningful insights from the data. This is because the foot fall is more on weekends than it is by the day of the week. Now this information about day of the week can be used to bring it out to make your model better.

This exercise of bringing out information from data is known as feature engineering.

([https://trainings.analyticsvidhya.com/courses/course-v1:AnalyticsVidhya+DS101+2018T2/about?](https://trainings.analyticsvidhya.com/courses/course-v1:AnalyticsVidhya+DS101+2018T2/about?utm_source=AnalyticsVidhya+CVDL101+CVDL101_T1/about?utm_medium=Stickybanner2utm_campaign=CV101banner)



Feature Engineering ?

You have completed the first 5 steps in data exploration – Variable Selection (<https://www.analyticsvidhya.com/blog/2015/02/data-exploration-5-steps/>), Missing Value Imputation (<https://www.analyticsvidhya.com/blog/2015/02/7-steps-data-exploration/>), Outliers Detection (<https://www.analyticsvidhya.com/blog/2015/02/outliers-detection-treatment-dataset/>) and Outliers Treatment (<https://www.analyticsvidhya.com/blog/2015/02/outliers-detection-treatment-dataset/>). Feature engineering

• Variable transformation.
([https://trainings.analyticsvidhya.com/courses/course-v1:AnalyticsVidhya+CVDL101+CVDL101_T1/about?](https://trainings.analyticsvidhya.com/courses/course-v1:AnalyticsVidhya+CVDL101+CVDL101_T1/about?utm_source=AnalyticsVidhya+CVDL101+CVDL101_T1/about?utm_medium=Stickybanner2utm_campaign=CV101banner)

• Variable / Feature creation.
([https://trainings.analyticsvidhya.com/courses/course-v1:AnalyticsVidhya+CVDL101+CVDL101_T1/about?](https://trainings.analyticsvidhya.com/courses/course-v1:AnalyticsVidhya+CVDL101+CVDL101_T1/about?utm_source=AnalyticsVidhya+CVDL101+CVDL101_T1/about?utm_medium=Stickybanner2utm_campaign=CV101banner)
utm_source=CV101AVBlogBanner&utm_medium=Stickybanner2utm_campaign=CV101banner). 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.

Let's look at the situations when variable transformation is useful.

Subscribe!

When should we use Variable Transformation?

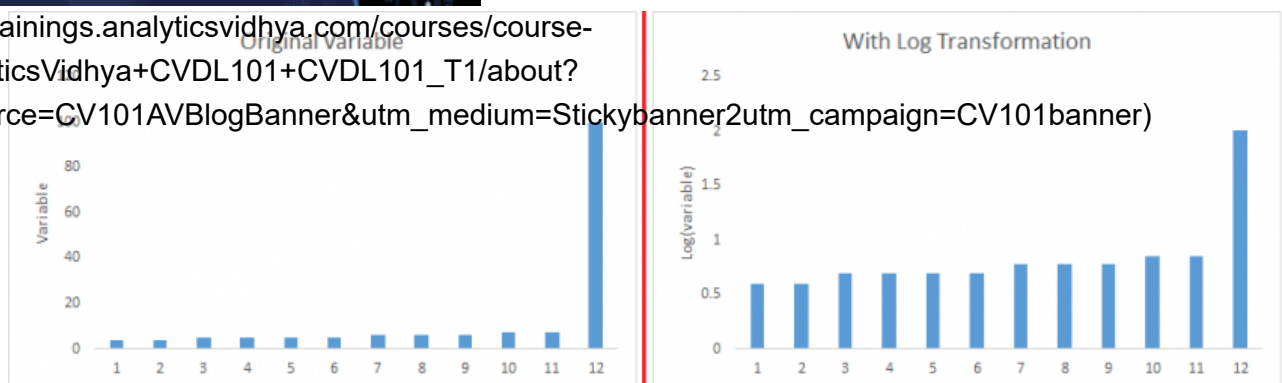
Below are the situations where variable transformation is a requisite:



([https://trainings.analyticsvidhya.com/courses/course-v1:AnalyticsVidhya+DS101+2018T2/about?](https://trainings.analyticsvidhya.com/courses/course-v1:AnalyticsVidhya+DS101+2018T2/about?utm_source=CV101AVBlogBanner&utm_medium=Stickybanner2utm_campaign=CV101banner)



(https://trainings.analyticsvidhya.com/courses/course-v1:AnalyticsVidhya+CVDL101+CVDL101_T1/about?utm_source=CV101AVBlogBanner&utm_medium=Stickybanner2utm_campaign=CV101banner)



(https://www.analyticsvidhya.com/wp-content/uploads/2015/03/Transformation_1.png)

- Variable Transformation is also done from an **implementation point of view** (Human involvement). Let's understand it more clearly. In one of my project on employee performance, I found that age has direct correlation with performance of the employee i.e. higher the age, better the performance.

Subscribe!

an implementation stand point, launching age based programme might present implementation challenge. However, categorizing the sales agents in three age group buckets of <30 years, 30-45 years and >45 and then formulating three different strategies for each group is a judicious approach. This categorization technique is known as Binning of Variables.



Variable Transformation?

Some variables are skewed. As discussed, some of them include square root, logarithmic, and many others. Let's look at these methods in this section.

One common transformation method used to change the shape of the distribution plot. It is generally used for reducing right skewness of the data and handling zero or negative values as well.

([https://trainings.analyticsvidhya.com/courses/course-](https://trainings.analyticsvidhya.com/courses/course-v1:AnalyticsVidhya+CV101+CV101_T1/about?utm_source=CV101AVBlogBanner&utm_medium=Stickybanner2utm_campaign=CV101banner)

[v1:AnalyticsVidhya+CV101+CV101_T1/about?](https://trainings.analyticsvidhya.com/courses/course-v1:AnalyticsVidhya+CV101+CV101_T1/about?utm_source=CV101AVBlogBanner&utm_medium=Stickybanner2utm_campaign=CV101banner)



The square and cube root of a variable has a sound effect on variable distribution. Square root is an arithmetic transformation. Cube root has its own advantage. It can handle zero. Square root can be applied to positive values including

negative values. It is performed on original values, percentile or frequency. The choice of transformation is based on business understanding. For example, we can categorize data into three groups: High, Average and Low. We can also perform co-variate transformation on more than one variables.

What is Feature / Variable Creation & its Benefits?

([https://trainings.analyticsvidhya.com/courses/course-](https://trainings.analyticsvidhya.com/courses/course-v1:AnalyticsVidhya+CV101+CV101_T1/about?utm_source=CV101AVBlogBanner&utm_medium=Stickybanner2utm_campaign=CV101banner)

[v1:AnalyticsVidhya+CV101+CV101_T1/about?](https://trainings.analyticsvidhya.com/courses/course-v1:AnalyticsVidhya+CV101+CV101_T1/about?utm_source=CV101AVBlogBanner&utm_medium=Stickybanner2utm_campaign=CV101banner)

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

Subscribe!

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. Let's look at it through "**Titanic – Kaggle competition**" (<https://www.kaggle.com/c/titanic-gettingStarted/data>). In this data set, variable age has missing



used the salutation (Master, Mr, Miss, Mrs) of name as a new variable to create? Honestly, this depends on business curiosity and the set of hypothesis he might have about the of variables, binning variables and other methods of variable to create new variables.

The most common application of dummy variable is to convert categorical variables. Dummy variables are also called Indicator Variables. It is a predictor in statistical models. Categorical variable can take 'gender'. We can produce two variables, namely, "**Var_Male**" with values 1 (Male) and 0 (No Male). We can produce two variables, namely, "**Var_Female**" with values 1 (Female) and 0 (No Female).

(<https://www.analyticsvidhya.com/blog/2013/11/simple-manipulations-extract-data/>) which can be applied to your data.



Gender	Var_Male	Var_Female
Male	1	0
Female	0	1
Female	0	1
Male	1	0
Female	0	1
Male	1	0
Male	1	0

(<https://www.analyticsvidhya.com/uploads/2015/03/Dummy.png>).

Information / creation ideas

(<https://www.analyticsvidhya.com/blog/2013/11/simple-manipulations-extract-data/>) which can be applied to your data.

utm_source=CV101AVBlogBanner&utm_medium=Stickybanner2utm_campaign=CV101banner)

End Notes

As mentioned in the beginning, quality and efforts invested in data exploration differentiates a good model from a bad model.

This ends our guide on data exploration and preparation. In this comprehensive guide, we looked at the seven steps of data exploration in detail. The aim of this series was to provide an in depth and step by step guide to an extremely important process in data science.

Subscribe!

Personally, I enjoyed writing this guide and would love to learn from your feedback. Did you find this guide useful? I would appreciate your suggestions/feedback. Please feel free to ask your questions through comments below.

If you like what you just read & want to continue your analytics learning, [subscribe to our newsletter](#) ([https://mailchi.mp/140700020/analyticsvidhya](#)), [follow us on LinkedIn](#) ([https://www.linkedin.com/company/analyticsvidhya](#)) or like our [facebook page](#) ([https://www.facebook.com/analyticsvidhya](#)).



Analytics Vidhya's Android APP



[Share this](#) ([https://trainings.analyticsvidhya.com/courses/course-v1:AnalyticsVidhya+DS101+2018T2/about?](#)



[like this](#) ([https://trainings.analyticsvidhya.com/courses/course-v1:AnalyticsVidhya+CVDL101+CVDL101_T1/about?](#)

utm_source=CV101AVBlogBanner&utm_medium=Stickybanner2utm_campaign=CV101banner)

TAGS : [BIVARIATE ANALYSIS \(HTTPS://WWW.ANALYTICSVIDHYA.COM/BLOG/TAG/BIVARIATE-ANALYSIS/\)](#), [DATA EXPLORATION](#)

([HTTPS://WWW.ANALYTICSVIDHYA.COM/BLOG/TAG/DATA-EXPLORATION/](#)), [DUMMY VARIABLES](#)

([HTTPS://WWW.ANALYTICSVIDHYA.COM/BLOG/TAG/DUMMY-VARIABLES/](#)), [FEATURE ENGINEERING](#)

([HTTPS://WWW.ANALYTICSVIDHYA.COM/BLOG/TAG/FEATURE-ENGINEERING/](#)), [KNN IMPUTATION](#)

([HTTPS://WWW.ANALYTICSVIDHYA.COM/BLOG/TAG/KNN-IMPUTATION/](#)), [MEDIAN IMPUTATION](#)

([HTTPS://WWW.ANALYTICSVIDHYA.COM/BLOG/TAG/MEDIAN-IMPUTATION/](#)), [MISSING VALUE](#)

([HTTPS://WWW.ANALYTICSVIDHYA.COM/BLOG/TAG/MISSING-VALUE/](#)), [MISSING VALUE IMPUTATION](#)

([HTTPS://WWW.ANALYTICSVIDHYA.COM/BLOG/TAG/MISSING-VALUE-IMPUTATION/](#)), [ONE HOT ENCODING](#)

([HTTPS://WWW.ANALYTICSVIDHYA.COM/BLOG/TAG/ONE-HOT-ENCODING/](#)), [OUTLIER REMOVAL](#)

Subscribe!

([HTTPS://WWW.ANALYTICSVIDHYA.COM/BLOG/TAG/OUTLIER-REMOVAL/](https://www.analyticsvidhya.com/blog/tag/outlier-removal/)), [UNIVARIATE ANALYSIS](https://www.analyticsvidhya.com/blog/tag/univariate-analysis/)

([HTTPS://WWW.ANALYTICSVIDHYA.COM/BLOG/TAG/UNIVARIATE-ANALYSIS/](https://www.analyticsvidhya.com/blog/tag/univariate-analysis/))



NEXT ARTICLE

perfectly captures the growth of Data Science

2016/01/20-powerful-images-perfectly-captures-growth-data-science/)

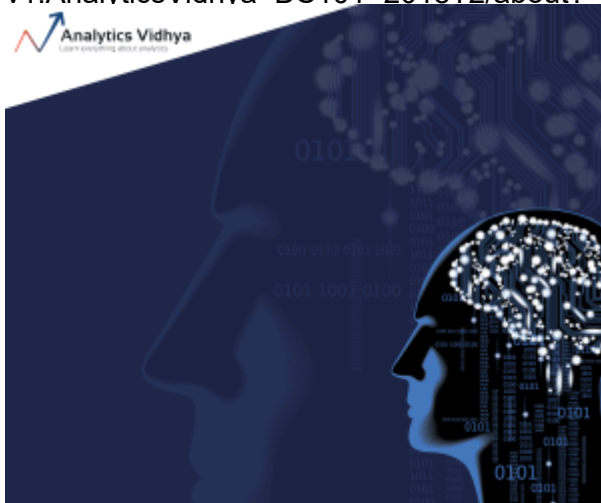
...

PREVIOUS ARTICLE

to Become a Data Scientist in 2016

m/blog/2016/01/ultimate-plan-data-scientist-2016/)

(<https://trainings.analyticsvidhya.com/courses/course-v1:AnalyticsVidhya+DS101+2018T2/about?>



(<https://www.analyticsvidhya.com/blog/author/sunil-ray/>).

[Analyticsvidhya.Com/Blog/Author/Sunil-Ray/](https://www.analyticsvidhya.com/blog/author/sunil-ray/))

...e professional with deep experience in the Indian Insurance industry. I have worked for various multi-national Insurance companies in last 7 years.

(https://trainings.analyticsvidhya.com/courses/course-v1:AnalyticsVidhya+CVDL101+CVDL101_T1/about?

utm_source=CV101AVBlogBanner&utm_medium=Stickybanner2utm_campaign=CV101banner)

RELATED ARTICLES

TAVISH SRIVASTAVA ([HTTPS://WWW.A...](https://www.a...))



(<https://www.analyticsvidhya.com/>

KUNAL JAIN ([HTTPS://WWW.ANALYTIC...](https://www.analytic...))



(<https://www.analyticsvidhya.com/>

ANALYTICS VIDHYA CONTENT TEAM ([H...](https://www.a...))



(<https://www.analyticsvidhya.com/>
Subscribe!

[time-series-data-r/](#)

Exploration of Time Series Data in R

(<https://www.analyticsvidhya.com/blog/2015/08/time-series-data-in-r/>)



based way to learn Machine Learning

([https://trainings.analyticsvidhya.com/courses/course-v1:AnalyticsVidhya+DS101+2018T2/about?](https://trainings.analyticsvidhya.com/courses/course-v1:AnalyticsVidhya+DS101+2018T2/about?utm_source=CV101AVBlogBanner&utm_medium=Stickybanner2utm_campaign=CV101banner))



Really useful and comprehensive, thanks!
([https://trainings.analyticsvidhya.com/courses/course-v1:AnalyticsVidhya+CVDL101+CVDL101_T1/about?](https://trainings.analyticsvidhya.com/courses/course-v1:AnalyticsVidhya+CVDL101+CVDL101_T1/about?utm_source=CV101AVBlogBanner&utm_medium=Stickybanner2utm_campaign=CV101banner))

utm_source=CV101AVBlogBanner&utm_medium=Stickybanner2utm_campaign=CV101banner)



BAGUINE BAZONGO

January 11, 2016 at 5:52 am (<https://www.analyticsvidhya.com/blog/2016/01/guide-data-exploration/#comment-103489>)

Hi Ray,

I would like to thank you very much for this useful post

I took more than 30 statistical courses but your post has summarized them for me

Now all things are clear about EDA

I'm member of the John Hopkins University Data Scientists (Coursera) Group

Best,



NANDU KULKARNI

[review-pg-diploma-in-data-analytics-by-upgrad-iiit-b/](#)

Course Review – PG Diploma in Data Analytics by UpGrad & IIIT-

RAY (<https://www.analyticsvidhya.com/blog/2015/08/time-series-data-in-r/>)



(<https://www.analyticsvidhya.com/blog/2014/10/introduction-sas-macro/>)

Introduction to SAS Macros

(<https://www.analyticsvidhya.com/blog/2014/10/introduction-sas-macro/>)

macro/)

[entrepreneurs-big-data-analytics-data-science/](#)

Top Datapreneurs who made data science what it is today.

(<https://www.analyticsvidhya.com/blog/2015/08/time-series-data-in-r/>)



(<https://www.analyticsvidhya.com/blog/2014/04/tricky-base-sas-interview-questions-part-ii/>)

Tricky Base SAS interview questions : Part-II

(<https://www.analyticsvidhya.com/blog/2014/04/tricky-base-sas-interview-questions-part-ii/>)

et a prompt response from the author. We request you to post **discussion portal** (<https://discuss.analyticsvidhya.com/>) to get your

Reply.

[analyticsvidhya.com/blog/2016/01/guide-data-exploration/#comment-103484](https://www.analyticsvidhya.com/blog/2016/01/guide-data-exploration/#comment-103484))

Reply.

Subscribe!

January 11, 2016 at 6:45 am (<https://www.analyticsvidhya.com/blog/2016/01/guide-data-exploration/#comment-103490>).

Excellent series of blog posts. Thanks and keep up the good work!



[Reply](#)

[analyticsvidhya.com/blog/2016/01/guide-data-exploration/#comment-103491](https://www.analyticsvidhya.com/blog/2016/01/guide-data-exploration/#comment-103491))

certainly a good refresher. Keep writing!

[Reply](#)

[analyticsvidhya.com/blog/2016/01/guide-data-exploration/#comment-103494](https://www.analyticsvidhya.com/blog/2016/01/guide-data-exploration/#comment-103494))

([https://trainings.analyticsvidhya.com/courses/course-v1:AnalyticsVidhya+CV101+CV101_T1/about?](https://trainings.analyticsvidhya.com/courses/course-v1:AnalyticsVidhya+CV101+CV101_T1/about?utm_source=CV101-AVBlogBanner&utm_medium=StickyBanner&utm_campaign=CV101Banner)

KARTHIKEYAN SANKARAN ([HTTP://WWW.TWITTER.COM/KARTHIKONBI](http://www.twitter.com/karthikonbi))

[Reply](#)



[analyticsvidhya.com/blog/2016/01/guide-data-exploration/#comment-103510](https://www.analyticsvidhya.com/blog/2016/01/guide-data-exploration/#comment-103510))

ts of Machine Learning. The points are explained in a simple

[Reply](#)

[analyticsvidhya.com/blog/2016/01/guide-data-exploration/#comment-103511](https://www.analyticsvidhya.com/blog/2016/01/guide-data-exploration/#comment-103511))

([https://trainings.analyticsvidhya.com/courses/course-v1:AnalyticsVidhya+CV101+CV101_T1/about?](https://trainings.analyticsvidhya.com/courses/course-v1:AnalyticsVidhya+CV101+CV101_T1/about?utm_source=CV101-AVBlogBanner&utm_medium=StickyBanner&utm_campaign=CV101Banner)

SATISH

[Reply](#)

January 11, 2016 at 2:15 pm (<https://www.analyticsvidhya.com/blog/2016/01/guide-data-exploration/#comment-103525>).

I haven't come across any other article as detailed as this one. Anyone who is keen about data exploration and Predictive Analytics in general has to go through this. Wondering if you have any data set where in I can work on it.

Bookmarked!



KHALID RIAZ

[Reply](#)

January 11, 2016 at 3:09 pm (<https://www.analyticsvidhya.com/blog/2016/01/guide-data-exploration/#comment-103528>).

[Subscribe!](#)

Hi Ray,

This is a great post. You have treated a fairly vast topic with just the right amount of detail. This makes it very useful, and also very interesting. Thank you for the good work. Keep it up.



[Reply](#)

<https://www.analyticsvidhya.com/blog/2016/01/guide-data-exploration/#comment-103529>

t. I like your blogs, Please continue your good work !

[Reply](#)

<https://www.analyticsvidhya.com/blog/2016/01/guide-data-exploration/#comment-103539>

comprehensive information. Also, I would request some to write a
er than other BI tools like Tableau, Qlikview....gaining more

<https://trainings.analyticsvidhya.com/courses/course->

<https://www.analyticsvidhya.com/blog/2016/01/guide-data-exploration/#comment-103548>



[Reply](#)

<https://www.analyticsvidhya.com/blog/2016/01/guide-data-exploration/#comment-103548>

l. I appreciated if you continue this wonderful work and post an
Python.



JOHN PAUL INERINEDA

<https://trainings.analyticsvidhya.com/courses/course->

<https://www.analyticsvidhya.com/blog/2016/01/guide-data-exploration/#comment-103553>

[Reply](#)

<https://www.analyticsvidhya.com/blog/2016/01/guide-data-exploration/#comment-103553>

utm_source=CV101AVBlogBanner&utm_medium=Stickybanner2utm_campaign=CV101banner)
Thank you Mr. Ray for the very comprehensive discussion on data exploration. I specially liked how you emphasized on the importance of EDA with this statement “quality and efforts invested in data exploration differentiates a good model from a bad model”. Great work Sir! I wish you can tackle dimensionality reduction techniques, principal components analysis, discriminant analysis and the likes in the future. Thanks again Mr. Ray.



SANDRA (HTTP://VWFXPAOAXM.COM)

[Reply](#)

February 9, 2016 at 3:08 pm (<https://www.analyticsvidhya.com/blog/2016/01/guide-data-exploration/#comment-105436>)

I found myself nodding my noggin all the way thruorh.

Subscribe!

**DEBASHIS ROUT**[Reply](#)January 12, 2016 at 3:56 am (<https://www.analyticsvidhya.com/blog/2016/01/guide-data-exploration/#comment-103557>)

Its really worth to read. Very comprehensive and easy to understand . I will be happy to read your article using R on data exploration & Data preparation.

[Reply](#)[analyticsvidhya.com/blog/2016/01/guide-data-exploration/#comment-103559](https://www.analyticsvidhya.com/blog/2016/01/guide-data-exploration/#comment-103559)

I'm glad it helped.

**DARIO ROMERO**[Reply](#)([https://training.analyticsvidhya.com/courses/course-](https://training.analyticsvidhya.com/courses/course-v1:AnalyticsVidhya+DS101+2018T2/about?utm_source=AnalyticsVidhya&utm_medium=BlogBanner&utm_campaign=CV101banner)[v1:AnalyticsVidhya+DS101+2018T2/about?](https://training.analyticsvidhya.com/courses/course-v1:AnalyticsVidhya+DS101+2018T2/about?utm_source=AnalyticsVidhya&utm_medium=BlogBanner&utm_campaign=CV101banner)[analyticsvidhya.com/blog/2016/01/guide-data-exploration/#comment-103560](https://www.analyticsvidhya.com/blog/2016/01/guide-data-exploration/#comment-103560))

important topic. BTW, there is a missing graph on the paragraph analysis. Could you please edit it and add the missing graph. I will be glad to see the ping file. Thanks.

[ANALYTICSVIDHYA.COM/BLOG/2016/01/GUIDE-DATA-](https://www.analyticsvidhya.com/blog/2016/01/guide-data-exploration/#comment-103561)[Reply](#)[analyticsvidhya.com/blog/2016/01/guide-data-exploration/#comment-103561](https://www.analyticsvidhya.com/blog/2016/01/guide-data-exploration/#comment-103561))

here:

http://www.analyticsvidhya.com/wp-content/uploads/2015/02/Data_exploration_4.pnghttp://www.analyticsvidhya.com/wp-content/uploads/2015/02/Data_exploration_4.png)

This picture is the missing one below the paragraph:

utm_source=CV101&utm_medium=BlogBanner&utm_campaign=CV101banner)

“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.”

**AKSHAY KHER (HTTPS://AKSHAYKHER.WORDPRESS.COM/)**[Reply](#)January 12, 2016 at 9:44 am (<https://www.analyticsvidhya.com/blog/2016/01/guide-data-exploration/#comment-103573>)

Hi Sunil,

An intriguing article, I can see the amount of hard work you must have put into it. Its a must read.

Subscribe!

Thanks,
Akshay Kher



[Reply](#)

analyticsvidhya.com/blog/2016/01/guide-data-exploration/#comment-103588)

Sunil. I had All these points scattered across but you got all of
bookmarked this page and this would now be my first page to

[Reply](#)

analyticsvidhya.com/blog/2016/01/guide-data-exploration/#comment-103633),

([https://trainings.analyticsvidhya.com/courses/course-v1:AnalyticsVidhya+DS101+2018T2/about?](https://trainings.analyticsvidhya.com/courses/course-v1:AnalyticsVidhya+DS101+2018T2/about?utm_source=AZIM)



[Reply](#)

analyticsvidhya.com/blog/2016/01/guide-data-exploration/#comment-103871)

a analytics project. Good work keep up.

[Reply](#)

analyticsvidhya.com/blog/2016/01/guide-data-exploration/#comment-104271)

helped me a lot....Thanks a lot

([https://trainings.analyticsvidhya.com/courses/course-v1:AnalyticsVidhya+CVDL101+CVDL101_T1/about?](https://trainings.analyticsvidhya.com/courses/course-v1:AnalyticsVidhya+CVDL101+CVDL101_T1/about?utm_source=AZIM)

[utm_source=AZIM](https://trainings.analyticsvidhya.com/courses/course-v1:AnalyticsVidhya+CVDL101+CVDL101_T1/about?utm_source=AZIM)101AVBlogBanner&utm_medium=Stickybanner2utm_campaign=CV101banner)

[Reply](#)

January 26, 2016 at 12:06 pm (<https://www.analyticsvidhya.com/blog/2016/01/guide-data-exploration/#comment-104638>).

when we create new variable like var_male and var_female we assign 0,1 to them? how is this 0,1 is used in our model? can we assign 200 instead of 0 and 2000 instead of 1?

Please help .



BRAJENDRA GOUDA

[Reply](#)

February 4, 2016 at 6:43 am (<https://www.analyticsvidhya.com/blog/2016/01/guide-data-exploration/#comment-105186>).

clear, Concise and Very well explained. !!

[Subscribe!](#)

**SUHEL**[Reply](#)

February 13, 2016 at 6:57 pm (<https://www.analyticsvidhya.com/blog/2016/01/guide-data-exploration/#comment-105638>)



for zero or negative values.

one to all values (if data has lots of zeros), take log, then finally negative.

([https://trainings.analyticsvidhya.com/courses/course-v1:AnalyticsVidhya+DS101+2018T2/about?](https://trainings.analyticsvidhya.com/courses/course-v1:AnalyticsVidhya+DS101+2018T2/about?utm_source=CV101AVBlogBanner&utm_medium=StickyBanner2utm_campaign=CV101banner)

FRANK SAUVAGE[Reply](#)

[w.analyticsvidhya.com/blog/2016/01/guide-data-exploration/#comment-105714](https://www.analyticsvidhya.com/blog/2016/01/guide-data-exploration/#comment-105714))



pedagogic and comprehensive. Two thumbs up!
g a new data project...

[Reply](#)

[analyticsvidhya.com/blog/2016/01/guide-data-exploration/#comment-105802](https://www.analyticsvidhya.com/blog/2016/01/guide-data-exploration/#comment-105802))

basic math /statistics understanding can also understand subject

([https://trainings.analyticsvidhya.com/courses/course-v1:AnalyticsVidhya+CVDL101+CVDL101_T1/about?](https://trainings.analyticsvidhya.com/courses/course-v1:AnalyticsVidhya+CVDL101+CVDL101_T1/about?utm_source=CV101AVBlogBanner&utm_medium=StickyBanner2utm_campaign=CV101banner)

BIDHAN[Reply](#)

February 21, 2016 at 11:07 am (<https://www.analyticsvidhya.com/blog/2016/01/guide-data-exploration/#comment-106000>)

Concise and comprehensive. Great article.

**WHY STATISTICS**[Reply](#)

February 25, 2016 at 6:55 am (<https://www.analyticsvidhya.com/blog/2016/01/guide-data-exploration/#comment-106218>)

Very well written.

**MATHU**[Reply](#)

Subscribe!

March 6, 2016 at 8:48 pm (<https://www.analyticsvidhya.com/blog/2016/01/guide-data-exploration/#comment-106727>)

One of the best blogs I have ever read till date!



[Reply](#)

[analyticsvidhya.com/blog/2016/01/guide-data-exploration/#comment-107505](https://www.analyticsvidhya.com/blog/2016/01/guide-data-exploration/#comment-107505)),

[Reply](#)

[analyticsvidhya.com/blog/2016/01/guide-data-exploration/#comment-108501](https://www.analyticsvidhya.com/blog/2016/01/guide-data-exploration/#comment-108501)),

Outliers and Missing Values??

([https://trainings.analyticsvidhya.com/courses/course-v1:AnalyticsVidhya+DL101+CV101+CV101T2/about?](https://trainings.analyticsvidhya.com/courses/course-v1:AnalyticsVidhya+DL101+CV101+CV101T2/about?utm_source=CV101AVBlogBanner&utm_medium=Stickybanner2utm_campaign=CV101banner)

[Reply](#)



[analyticsvidhya.com/blog/2016/01/guide-data-exploration/#comment-108686](https://www.analyticsvidhya.com/blog/2016/01/guide-data-exploration/#comment-108686)),

[Reply](#)

[analyticsvidhya.com/blog/2016/01/guide-data-exploration/#comment-108689](https://www.analyticsvidhya.com/blog/2016/01/guide-data-exploration/#comment-108689)),

([https://trainings.analyticsvidhya.com/courses/course-v1:AnalyticsVidhya+DL101+CV101+CV101T1/about?](https://trainings.analyticsvidhya.com/courses/course-v1:AnalyticsVidhya+DL101+CV101+CV101T1/about?utm_source=CV101AVBlogBanner&utm_medium=Stickybanner2utm_campaign=CV101banner)

[Reply](#)

April 12, 2016 at 6:22 pm (<https://www.analyticsvidhya.com/blog/2016/01/guide-data-exploration/#comment-109365>),
<https://www.analyticsvidhya.com/blog/2016/01/guide-data-exploration/#comment-109365>
 utm_source=CV101AVBlogBanner&utm_medium=Stickybanner2utm_campaign=CV101banner)
 Amazing guide.. very structured and simplistic. enjoyed and learnt a lot reading this article.



ANDRII

[Reply](#)

May 31, 2016 at 6:54 pm (<https://www.analyticsvidhya.com/blog/2016/01/guide-data-exploration/#comment-111650>),

Many thanks for the guide, very useful. Would you advise R packages that help with data exploration?
 Thanks



GUSTAVO

[Reply](#)

June 1, 2016 at 2:26 pm (<https://www.analyticsvidhya.com/blog/2016/01/guide-data-exploration/#comment-111695>)

[Subscribe!](#)

THANK YOU FOR SHARING THIS CONCEPTS AND METHOD.



ARIJIT

[Reply](#)



<https://www.analyticsvidhya.com/blog/2016/01/guide-data-exploration/#comment-112328>

How should we treat them in a logistic regression framework?

[Reply](#)

<https://www.analyticsvidhya.com/blog/2016/01/guide-data-exploration/#comment-112364>

([https://trainings.analyticsvidhya.com/courses/course-v1:AnalyticsVidhya+CVDL101+CVDL101_T2/about?](https://trainings.analyticsvidhya.com/courses/course-v1:AnalyticsVidhya+CVDL101+CVDL101_T2/about?utm_source=CVDL101+CVDL101_T2)

[Reply](#)



<https://www.analyticsvidhya.com/blog/2016/01/guide-data-exploration/#comment-112395>

[Reply](#)

<https://www.analyticsvidhya.com/blog/2016/01/guide-data-exploration/#comment-112716>

([https://trainings.analyticsvidhya.com/courses/course-v1:AnalyticsVidhya+CVDL101+CVDL101_T1/about?](https://trainings.analyticsvidhya.com/courses/course-v1:AnalyticsVidhya+CVDL101+CVDL101_T1/about?utm_source=CVDL101+CVDL101_T1)

[Reply](#)

<https://www.analyticsvidhya.com/blog/2016/01/guide-data-exploration/#comment-112794>

I open a file in google drive to keep this page alone as a cheatsheet...Thank you so much..



MARKETING ANALYST (HTTP://WWW.DATANANALYTICS.COM)

[Reply](#)

July 6, 2016 at 11:55 am (<https://www.analyticsvidhya.com/blog/2016/01/guide-data-exploration/#comment-113148>)

This is very useful summary, thank you for that!

I particularly liked the before-after comparisons to demonstrate the importance of the process steps.

Thanks,

Chill

[Subscribe!](#)

**NIRAV**[Reply](#)

July 16, 2016 at 6:52 pm (<https://www.analyticsvidhya.com/blog/2016/01/guide-data-exploration/#comment-113543>)

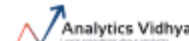
Great article! Few questions:

1) Do you run your data exploration on sample or full data set? If sample then what percentage and any



([https://trainings.analyticsvidhya.com/courses/course-](https://trainings.analyticsvidhya.com/courses/course-v1:AnalyticsVidhya+DS101+2019T2/about?utm_source=QV1012016BlogBanner&utm_medium=StickyBanner2016_campaign=QV101banner)

v1:AnalyticsVidhya+DS101+2019T2/about?

[Reply](#)

([https://trainings.analyticsvidhya.com/courses/course-](https://trainings.analyticsvidhya.com/courses/course-v1:AnalyticsVidhya+DS101+2019T2/about?utm_source=QV1012016BlogBanner&utm_medium=StickyBanner2016_campaign=QV101banner)

v1:AnalyticsVidhya+DS101+2019T2/about?

[Reply](#)

utm_source=QV1012016BlogBanner&utm_medium=StickyBanner2016_campaign=QV101banner) ([comment-113920](https://www.analyticsvidhya.com/blog/2016/01/guide-data-exploration/#comment-113920))

utm_source=QV1012016BlogBanner&utm_medium=StickyBanner2016_campaign=QV101banner) ([comment-114585](https://www.analyticsvidhya.com/blog/2016/01/guide-data-exploration/#comment-114585))

[Reply](#)

([https://trainings.analyticsvidhya.com/courses/course-](https://trainings.analyticsvidhya.com/courses/course-v1:AnalyticsVidhya+DS101+2019T2/about?utm_source=QV1012016BlogBanner&utm_medium=StickyBanner2016_campaign=QV101banner)

v1:AnalyticsVidhya+DS101+2019T2/about?utm_source=QV1012016BlogBanner&utm_medium=StickyBanner2016_campaign=QV101banner) ([comment-114806](https://www.analyticsvidhya.com/blog/2016/01/guide-data-exploration/#comment-114806))

reading your blog and learned a lot!!!! Thanks a lot for investing



([https://trainings.analyticsvidhya.com/courses/course-](https://trainings.analyticsvidhya.com/courses/course-v1:AnalyticsVidhya+DS101+2019T2/about?utm_source=QV1012016BlogBanner&utm_medium=StickyBanner2016_campaign=QV101banner)

v1:AnalyticsVidhya+DS101+2019T2/about?

[Reply](#)

utm_source=QV1012016BlogBanner&utm_medium=StickyBanner2016_campaign=QV101banner) ([comment-114806](https://www.analyticsvidhya.com/blog/2016/01/guide-data-exploration/#comment-114806))

utm_source=QV1012016BlogBanner&utm_medium=StickyBanner2016_campaign=QV101banner) ([comment-114806](https://www.analyticsvidhya.com/blog/2016/01/guide-data-exploration/#comment-114806))

Well Written. it really shows how to tackle the data

**RAJESH SRINIVASAN**[Reply](#)

August 24, 2016 at 8:33 am (<https://www.analyticsvidhya.com/blog/2016/01/guide-data-exploration/#comment-115083>)

Excellent read on EDA simple and to the point. Great Help to newbie like me.

**MANGESH PANCHWAGH**[Reply](#)

August 29, 2016 at 10:35 am (<https://www.analyticsvidhya.com/blog/2016/01/guide-data-exploration/#comment-115281>)

Subscribe

Thank you for sharing knowledge. It helps a lot.



AARON

[Reply](#)

September 12, 2016 at 1:33 am (<https://www.analyticsvidhya.com/blog/2016/01/guide-data-exploration/#comment-115291>)



[Reply](#)

[v.analyticsvidhya.com/blog/2016/01/guide-data-exploration/#comment-115690](https://www.analyticsvidhya.com/blog/2016/01/guide-data-exploration/#comment-115690)).

[Reply](#)

([https://trainings.analyticsvidhya.com/courses/course-v1:AnalyticsVidhya+DS101+2018T2/about?](https://trainings.analyticsvidhya.com/courses/course-v1:AnalyticsVidhya+DS101+2018T2/about?utm_source=CV101AVBlogBanner&utm_medium=Stickybanner2utm_campaign=CV101banner)



tion of missing values. I once had a dataset with missing values replacing missing values with the most frequent value of that values and found that they were all uniformly distributed. With ue by randomly choosing a value among the set of unique to hear if this was statistically the right thing to do?

[Reply](#)

[w.analyticsvidhya.com/blog/2016/01/guide-data-exploration/#comment-115929](https://www.analyticsvidhya.com/blog/2016/01/guide-data-exploration/#comment-115929)).

The Best Period
([https://trainings.analyticsvidhya.com/courses/course-v1:AnalyticsVidhya+CVDL101+CVDL101_T1/about?](https://trainings.analyticsvidhya.com/courses/course-v1:AnalyticsVidhya+CVDL101+CVDL101_T1/about?utm_source=CV101AVBlogBanner&utm_medium=Stickybanner2utm_campaign=CV101banner)

utm_source=CV101AVBlogBanner&utm_medium=Stickybanner2utm_campaign=CV101banner)



NEERAJA

[Reply](#)

September 12, 2016 at 1:33 am (<https://www.analyticsvidhya.com/blog/2016/01/guide-data-exploration/#comment-115953>)

Hi Sunil

Thank you very much for really useful and clear structure.



GAURAV

[Reply](#)

September 15, 2016 at 7:48 am (<https://www.analyticsvidhya.com/blog/2016/01/guide-data-exploration/#comment-116078>)

Great explanation, would be better. If you could give us some sample data and then explain step by step on that.

[Subscribe!](#)

**ANUJ JAIN**[Reply](#)

September 22, 2016 at 4:58 pm (<https://www.analyticsvidhya.com/blog/2016/01/guide-data-exploration/#comment-116356>)



manner. 😊

[Reply](#)

[analyticsvidhya.com/blog/2016/01/guide-data-exploration/#comment-116883](https://www.analyticsvidhya.com/blog/2016/01/guide-data-exploration/#comment-116883))

topic of data exploration with enough details to understand.

([https://trainings.analyticsvidhya.com/courses/course-v1:AnalyticsVidhya+DS101+2018T2/about?](https://trainings.analyticsvidhya.com/courses/course-v1:AnalyticsVidhya+DS101+2018T2/about?utm_source=AnalyticsVidhyaBanner&utm_medium=Stickybanner2&utm_campaign=CV101banner)

[Reply](#)

[analyticsvidhya.com/blog/2016/01/guide-data-exploration/#comment-117051](https://www.analyticsvidhya.com/blog/2016/01/guide-data-exploration/#comment-117051))

very much for sharing

[ANALYTICSVIDHYA.GITHUB.IO](https://www.analyticsvidhya.com/blog/2016/01/guide-data-exploration/#comment-117927))

[Reply](#)

[analyticsvidhya.com/blog/2016/01/guide-data-exploration/#comment-117927](https://www.analyticsvidhya.com/blog/2016/01/guide-data-exploration/#comment-117927))

step by step explanation of EDA process.

time as an organized flow.

([https://trainings.analyticsvidhya.com/courses/course-v1:AnalyticsVidhya+CVDL101+CVDL101_T1/about?](https://trainings.analyticsvidhya.com/courses/course-v1:AnalyticsVidhya+CVDL101+CVDL101_T1/about?utm_source=CV101AVBlogBanner&utm_medium=Stickybanner2&utm_campaign=CV101banner)

[v1:AnalyticsVidhya+CVDL101+CVDL101_T1/about?](https://trainings.analyticsvidhya.com/courses/course-v1:AnalyticsVidhya+CVDL101+CVDL101_T1/about?utm_source=CV101AVBlogBanner&utm_medium=Stickybanner2&utm_campaign=CV101banner)

[utm_source=CV101AVBlogBanner&utm_medium=Stickybanner2utm_campaign=CV101banner\)](https://trainings.analyticsvidhya.com/courses/course-v1:AnalyticsVidhya+CVDL101+CVDL101_T1/about?utm_source=CV101AVBlogBanner&utm_medium=Stickybanner2&utm_campaign=CV101banner)

**CAUI**[Reply](#)

January 11, 2017 at 12:13 am (<https://www.analyticsvidhya.com/blog/2016/01/guide-data-exploration/#comment-120843>)

I've started to study Data Science few months ago, this tutorial was one of the most clarifying for me, the step by step guide introduced the theory that can easily be used at practice. Thanks for the advices.

**POONAM LATA**[Reply](#)

January 25, 2017 at 8:03 am (<https://www.analyticsvidhya.com/blog/2016/01/guide-data-exploration/#comment-121584>)

Great! Very crisp, yet comprehensive.

[Subscribe!](#)

**BILL**[Reply](#)

January 30, 2017 at 1:17 am (<https://www.analyticsvidhya.com/blog/2016/01/guide-data-exploration/#comment-121813>)

“Though, It can’t be applied to zero or negative values as well”. Did you mean “can” and not “can’t”

[Reply](#)

[analyticsvidhya.com/blog/2016/01/guide-data-exploration/#comment-124381](https://www.analyticsvidhya.com/blog/2016/01/guide-data-exploration/#comment-124381))

[Reply](#)

[yticsvidhya.com/blog/2016/01/guide-data-exploration/#comment-126073](https://www.analyticsvidhya.com/blog/2016/01/guide-data-exploration/#comment-126073))

([https://trainings.analyticsvidhya.com/courses/course-](https://trainings.analyticsvidhya.com/courses/course-v1:AnalyticsVidhya+DS401+2018T2/about?utm_source=CV101AVBlogBanner&utm_medium=Stickybanner2utm_campaign=CV101banner)

[v1:AnalyticsVidhya+DS401+2018T2/about?](https://trainings.analyticsvidhya.com/courses/course-v1:AnalyticsVidhya+DS401+2018T2/about?utm_source=CV101AVBlogBanner&utm_medium=Stickybanner2utm_campaign=CV101banner)



WHICH TEACHES DATA EXPLORATION IN DETAIL |
[WWW.SHUJIANLIU.COM/BLOGS/A-COMPLETE-TUTORIAL-](https://www.shujianliu.com/blogs/a-complete-tutorial-data-exploration-in-detail-2/)
[DATA-EXPLORATION-IN-DETAIL-2/](https://www.shujianliu.com/blogs/a-complete-tutorial-data-exploration-in-detail-2/)

[ticsvidhya.com/blog/2016/01/guide-data-exploration/#comment-126084](https://www.analyticsvidhya.com/blog/2016/01/guide-data-exploration/#comment-126084))

[blog/2016/01/guide-data-exploration/?utm_content=buffer087f0](https://www.analyticsvidhya.com/blog/2016/01/guide-data-exploration/?utm_content=buffer087f0)
[/guide-data-exploration/?utm_content=buffer087f0](https://www.analyticsvidhya.com/blog/2016/01/guide-data-exploration/?utm_content=buffer087f0)) [...]

[Reply](#)

[ticsvidhya.com/blog/2016/01/guide-data-exploration/#comment-126101](https://www.analyticsvidhya.com/blog/2016/01/guide-data-exploration/#comment-126101))

([https://trainings.analyticsvidhya.com/courses/course-](https://trainings.analyticsvidhya.com/courses/course-v1:AnalyticsVidhya+DS401+2018T2/about?utm_source=CV101AVBlogBanner&utm_medium=Stickybanner2utm_campaign=CV101banner)

[v1:AnalyticsVidhya+DS401+2018T2/about?](https://trainings.analyticsvidhya.com/courses/course-v1:AnalyticsVidhya+DS401+2018T2/about?utm_source=CV101AVBlogBanner&utm_medium=Stickybanner2utm_campaign=CV101banner)

[utm_source=CV101AVBlogBanner&utm_medium=Stickybanner2utm_campaign=CV101banner](https://trainings.analyticsvidhya.com/courses/course-v1:AnalyticsVidhya+DS401+2018T2/about?utm_source=CV101AVBlogBanner&utm_medium=Stickybanner2utm_campaign=CV101banner))

**JOSEPH MACHADO**[Reply](#)

August 19, 2017 at 10:55 pm (<https://www.analyticsvidhya.com/blog/2016/01/guide-data-exploration/#comment-134701>)

Hi Jack,

I am working on a prediction problem for which I am using this post as a guide for EDA. If you want some code examples please check out https://github.com/JosephKevin/sales_prediction

(https://github.com/JosephKevin/sales_prediction)

Regards,

Joseph

[Subscribe!](#)



HIRENDRASINGH CHAUHAN

[Reply](#)

April 18, 2017 at 8:07 am (<https://www.analyticsvidhya.com/blog/2016/01/guide-data-exploration/#comment-127075>)

Very well written article. One suggestion for next Enhanced version of the Article



along with example from same data set is provided.

[V.THECRAZYANALYST.COM](#))

[Reply](#)

[analyticsvidhya.com/blog/2016/01/guide-data-exploration/#comment-127082](https://www.analyticsvidhya.com/blog/2016/01/guide-data-exploration/#comment-127082))

Exploration process very lucidly. Kudos !

[Reply](#)

[analyticsvidhya.com/blog/2016/01/guide-data-exploration/#comment-130364](https://www.analyticsvidhya.com/blog/2016/01/guide-data-exploration/#comment-130364))

(<https://trainings.analyticsvidhya.com/courses/course->

[H1:Smaltravelsva+DS101+CV101+CVDL101+T1+about?](#)



[Reply](#)

[analyticsvidhya.com/blog/2016/01/guide-data-exploration/#comment-130519](https://www.analyticsvidhya.com/blog/2016/01/guide-data-exploration/#comment-130519))

Very helpful for base understanding.

[Reply](#)

[analyticsvidhya.com/blog/2016/01/guide-data-exploration/#comment-130719](https://www.analyticsvidhya.com/blog/2016/01/guide-data-exploration/#comment-130719))

Thank you for the article. It is super helpful!
(<https://trainings.analyticsvidhya.com/courses/course->

[Do you mind providing the download of the dataset as well? Thanks! As a beginner, I'd like to follow your tutorial step by step!](#)
[V1:AnalyticsVidhya+CVDE101+CVDL101+T1+about?](#)
[utm_source=CV101AVBlogBanner&utm_medium=Stickybanner2utm_campaign=CV101banner](#))



KISHORE

[Reply](#)

June 22, 2017 at 11:48 am (<https://www.analyticsvidhya.com/blog/2016/01/guide-data-exploration/#comment-130915>)

Hello Sunil,

Really an amazing stuff . Appreciate you for sharing your hard work..



AKASH GOYAL

[Reply](#)

June 26, 2017 at 3:19 pm (<https://www.analyticsvidhya.com/blog/2016/01/guide-data-exploration/#comment-131134>)

[Subscribe!](#)

please tell me ,which course are better for statistical and exploratory analysis in sense of industry.



LAUTARO

[Reply](#)



[https://trainings.analyticsvidhya.com/courses/course-v1:AnalyticsVidhya+DS101+2018T2/about?](https://trainings.analyticsvidhya.com/courses/course-v1:AnalyticsVidhya+DS101+2018T2/about?utm_source=CV101AVBlogBanner&utm_medium=Stickybanner2utm_campaign=CV101banner)



Really great help for beginners in data exploration and feature engineering!
([https://trainings.analyticsvidhya.com/courses/course-v1:AnalyticsVidhya+CVDL101+CVDL101_T1/about?](https://trainings.analyticsvidhya.com/courses/course-v1:AnalyticsVidhya+CVDL101+CVDL101_T1/about?utm_source=CV101AVBlogBanner&utm_medium=Stickybanner2utm_campaign=CV101banner)

[utm_source=CV101AVBlogBanner&utm_medium=Stickybanner2utm_campaign=CV101banner\)](https://trainings.analyticsvidhya.com/courses/course-v1:AnalyticsVidhya+CVDL101+CVDL101_T1/about?utm_source=CV101AVBlogBanner&utm_medium=Stickybanner2utm_campaign=CV101banner)



BHUVANA NARAYANAN

[Reply](#)

July 21, 2017 at 12:27 pm (<https://www.analyticsvidhya.com/blog/2016/01/guide-data-exploration/#comment-132558>)

Very clear and concise as well as informative . Well done.



RAFAEL

[Reply](#)

July 23, 2017 at 11:54 am (<https://www.analyticsvidhya.com/blog/2016/01/guide-data-exploration/#comment-132696>)

Very good article! Comprehensive and very easy to understand. Do you guys have any ebooks with all of this content?

[Subscribe!](#)



ANU (HTTP://WWW.ANALYTICSVIDHYA.COM)

[Reply](#)

July 31, 2017 at 8:52 am (<https://www.analyticsvidhya.com/blog/2016/01/guide-data-exploration/#comment-133264>),

great article. precisely written. Thanks for the clarity in the explanation given. keep up the good work.



([https://www.analyticsvidhya.com/courses/course-v1:AnalyticsVidhya+DS101+2018T2/about?](https://www.analyticsvidhya.com/courses/course-v1:AnalyticsVidhya+DS101+2018T2/about?utm_source=AnalyticsVidhyaBanner&utm_medium=Stickybanner2utm_campaign=CV101banner)



([https://www.analyticsvidhya.com/courses/course-v1:AnalyticsVidhya+CVDL101+CVDL101_T1/about?](https://www.analyticsvidhya.com/courses/course-v1:AnalyticsVidhya+CVDL101+CVDL101_T1/about?utm_source=CV101AVBlogBanner&utm_medium=Stickybanner2utm_campaign=CV101banner)

[RJB](#) ([https://www.analyticsvidhya.com/courses/course-v1:AnalyticsVidhya+CVDL101+CVDL101_T1/about?](https://www.analyticsvidhya.com/courses/course-v1:AnalyticsVidhya+CVDL101+CVDL101_T1/about?utm_source=CV101AVBlogBanner&utm_medium=Stickybanner2utm_campaign=CV101banner)



VIVEK

[Reply](#)

August 11, 2017 at 6:31 pm (<https://www.analyticsvidhya.com/blog/2016/01/guide-data-exploration/#comment-134102>),

The guide is super. IF you can take a sample dataset and apply all the steps to make dataset more informative then it would be very helpful.



JOSEPH MACHADO

[Reply](#)

August 19, 2017 at 10:43 pm (<https://www.analyticsvidhya.com/blog/2016/01/guide-data-exploration/#comment-134700>)

[Subscribe!](#)

Hi Sunil,

Thank you for the amazing article, very organized and clear. I have a question

In the 'Categorical & Continuous' bivariate analysis part, if ANOVA shows a statistically significant difference between various groups in one variable, how do we incorporate this knowledge into the prediction process ?



[Reply](#)

analyticsvidhya.com/blog/2016/01/guide-data-exploration/#comment-135043)

[Reply](#)

(<https://trainings.analyticsvidhya.com/courses/course-v1:AnalyticsVidhya+DS101+2018T2/about?>
August 20, 2017 at 10:10 pm (<https://www.analyticsvidhya.com/blog/2016/01/guide-data-exploration/#comment-135117>).



[RULEUR.COM](#))

[Reply](#)

analyticsvidhya.com/blog/2016/01/guide-data-exploration/#comment-135526)

ing down each individual concept. Adding some actual code to
actical standpoint.

[Reply](#)

(https://trainings.analyticsvidhya.com/courses/course-v1:AnalyticsVidhya+CVDL101+CVDL101_T1/about?
September 18, 2017 at 10:11 pm (<https://www.analyticsvidhya.com/blog/2016/01/guide-data-exploration/#comment-135996>).

Wonderful and Descriptive but can I get some Working Codes which can highlight the procedure "what if the
utm_source=CV101AVBlogBanner&utm_medium=StickyBanner2utm_campaign=CV101Banner")
data is heterogeneous..?" (I mean to say multi-valued data and mixture of numeric and text form). Does
Python, R or Matlab provide any help in this regard..?

• **FEATURE ENGINEERING 特徴工程中常見的方法 – I FAILED THE TURING TEST** **([HTTPS://VINTA.WS/CODE/FEATURE-ENGINEERING.HTML](https://vinta.ws/code/feature-engineering.html))**

September 18, 2017 at 4:32 pm (<https://www.analyticsvidhya.com/blog/2016/01/guide-data-exploration/#comment-137404>)

[...] ref: <https://www.analyticsvidhya.com/blog/2016/01/guide-data-exploration/>
(<https://www.analyticsvidhya.com/blog/2016/01/guide-data-exploration/>) [...]



SRINI

[Reply](#)

October 25, 2017 at 11:07 am (<https://www.analyticsvidhya.com/blog/2016/01/guide-data-exploration/#comment-140750>).

[Subscribe](#)

Thanks alot. Great article.



JOHAN (HTTP://WWW.MEDISENTIO.COM)

[Reply](#)

May 14, 2018 at 6:44 pm (<https://www.analyticsvidhya.com/blog/2016/01/guide-data-exploration/#comment-151818>),



[Reply](#)

[analyticsvidhya.com/blog/2016/01/guide-data-exploration/#comment-152181](https://www.analyticsvidhya.com/blog/2016/01/guide-data-exploration/#comment-152181)),

([https://trainings.analyticsvidhya.com/courses/course-](https://trainings.analyticsvidhya.com/courses/course-v1:AnalyticsVidhya+DL101_T1/about?utm_source=Cy4121A+Blog+Banner+utm_medium=Sticky+Banner+utm_campaign=Cy4121A+Banner)

[v1:AnalyticsVidhya+DL101_T1/about?](https://trainings.analyticsvidhya.com/courses/course-v1:AnalyticsVidhya+DL101_T1/about?utm_source=Cy4121A+Blog+Banner+utm_medium=Sticky+Banner+utm_campaign=Cy4121A+Banner)

[Reply](#)



[analyticsvidhya.com/blog/2016/01/guide-data-exploration/#comment-152208](https://www.analyticsvidhya.com/blog/2016/01/guide-data-exploration/#comment-152208)),



[Reply](#)

[analyticsvidhya.com/blog/2016/01/guide-data-exploration/#comment-152555](https://www.analyticsvidhya.com/blog/2016/01/guide-data-exploration/#comment-152555)),

to me

([https://trainings.analyticsvidhya.com/courses/course-](https://trainings.analyticsvidhya.com/courses/course-v1:AnalyticsVidhya+DL101_T1/about?utm_source=Cy4121A+Blog+Banner+utm_medium=Sticky+Banner+utm_campaign=Cy4121A+Banner)

[v1:AnalyticsVidhya+DL101_T1/about?](https://trainings.analyticsvidhya.com/courses/course-v1:AnalyticsVidhya+DL101_T1/about?utm_source=Cy4121A+Blog+Banner+utm_medium=Sticky+Banner+utm_campaign=Cy4121A+Banner)

[Reply](#)

<https://www.analyticsvidhya.com/blog/2016/01/guide-data-exploration/#comment-152915>),

It is very useful. Thank you for your efforts Sunil.



NADA B

[Reply](#)

May 14, 2018 at 7:41 pm (<https://www.analyticsvidhya.com/blog/2016/01/guide-data-exploration/#comment-153252>)

Very complete and useful ! Thank you !



BHAGWAT

[Reply](#)

May 20, 2018 at 11:50 pm (<https://www.analyticsvidhya.com/blog/2016/01/guide-data-exploration/#comment-153413>)

[Subscribe!](#)

Extremely useful article, can someone guide me to a link or any resource where all steps mentioned above are applied on real dataset.



ANSHUWAR SINGH



[Reply](#)

analyticsvidhya.com/blog/2016/01/guide-data-exploration/#comment-153442)

[n/courses/course-](#)

[ut](#)) is a training course on R for big mart sales dataset. A similar

[Reply](#)

analyticsvidhya.com/blog/2016/01/guide-data-exploration/#comment-153471)

(<https://trainings.analyticsvidhya.com/courses/course->

[Create article in Analytics Vidhya DS101+2018T2/about?](#)



(<https://trainings.analyticsvidhya.com/courses/course->

[v1:AnalyticsVidhya+CVDL101+CVDL101_T1/about?](#)

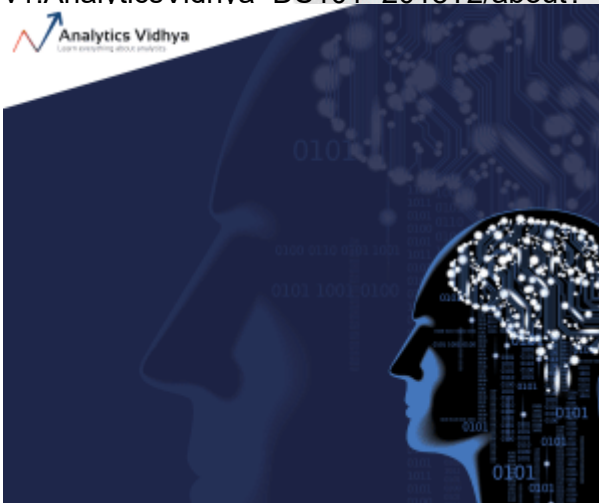
[utm_source=CV101AVBlogBanner&utm_medium=Stickybanner2utm_campaign=CV101banner](#))

Subscribe!



(<https://www.analyticsvidhya.com/datahack-summit-2018/>?)

(<https://trainings.analyticsvidhya.com/courses/course-v1:AnalyticsVidhya+DS101+2018T2/about?>



(https://trainings.analyticsvidhya.com/courses/course-v1:AnalyticsVidhya+CVDL101+CVDL101_T1/about?utm_source=CV101AVBlogBanner&utm_medium=Stickybanner2utm_campaign=CV101banner)

POPULAR POSTS

24 Ultimate Data Science Projects To Boost Your Knowledge and Skills (& can be accessed freely)
(<https://www.analyticsvidhya.com/blog/2018/05/24-ultimate-data-science-projects-to-boost-your-knowledge-and-skills/>)

A Complete Tutorial to Learn Data Science with Python from Scratch
(<https://www.analyticsvidhya.com/blog/2016/01/complete-tutorial-learn-data-science-python-scratch-2/>)

Essentials of Machine Learning Algorithms (with Python and R Codes)
(<https://www.analyticsvidhya.com/blog/2017/09/common-machine-learning-algorithms/>)

Subscribe!

Understanding Support Vector Machine algorithm from examples (along with code)

(<https://www.analyticsvidhya.com/blog/2017/09/understaing-support-vector-machine-example-code/>)

7 Types of Regression Techniques you should know!

(<https://www.analyticsvidhya.com/blog/2015/08/comprehensive-guide-regression/>)



thm (with codes in Python and R)

([7/09/naive-bayes-explained/](https://www.analyticsvidhya.com/blog/2017/09/naive-bayes-explained/))

te a Time Series Forecast (with Codes in Python)

([6/02/time-series-forecasting-codes-python/](https://www.analyticsvidhya.com/blog/2016/02/time-series-forecasting-codes-python/))

lified (with implementation in Python)

([8/03/introduction-k-neighbours-algorithm-clustering/](https://www.analyticsvidhya.com/blog/2018/03/introduction-k-neighbours-algorithm-clustering/))

(<https://trainings.analyticsvidhya.com/courses/course->

W1:AnalyticsVidhya+DS101+2018T2/about?

Analytics Vidhya+DS101+2018T2/about?utm_source=CV101AVBlogBanner&utm_medium=Stickybanner2utm_campaign=CV101banner) **2018/11/data-engineer-comprehensive-list-resources-get-**



-CNN Algorithm for Object Detection (Part 2 – with Python

[blog/2018/11/implementation-faster-r-cnn-python-object-](https://www.analyticsvidhya.com/blog/2018/11/implementation-faster-r-cnn-python-object-)

ies & Reddit Discussions (October 2018)

[2018/11/best-machine-learning-github-repositories-reddit-](https://www.analyticsvidhya.com/blog/2018/11/best-machine-learning-github-repositories-reddit-)

(<https://www.analyticsvidhya.com/courses/course->

W1:AnalyticsVidhya+CVDL101+CVDL101_T1/about?

NOVEMBER 1, 2018 utm_source=CV101AVBlogBanner&utm_medium=Stickybanner2utm_campaign=CV101banner)

An Introduction to Text Summarization using the TextRank Algorithm (with Python implementation)

(<https://www.analyticsvidhya.com/blog/2018/11/introduction-text-summarization-textrank-python/>)

NOVEMBER 1, 2018

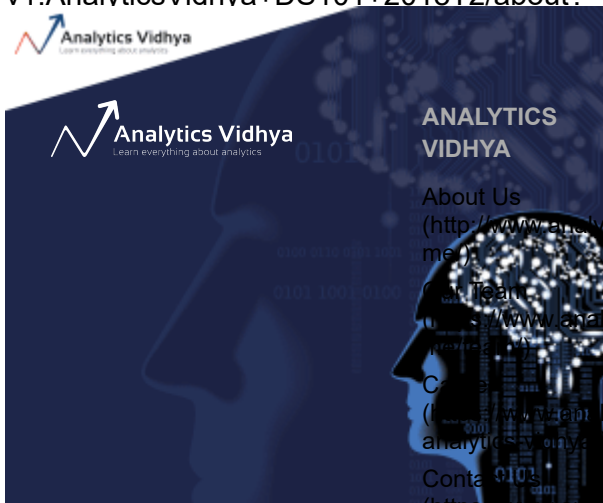
Subscribe!



(<http://www.edvancer.in/certified-data-scientist-with-python->

[s&utm_campaign=AVadsnonfc&utm_content=pythonavad](#))

(<https://trainings.analyticsvidhya.com/courses/course-v1:AnalyticsVidhya+DS101+2018T2/about?>



(https://trainings.analyticsvidhya.com/courses/course-v1:AnalyticsVidhya+CVDL101+2018T1/about?utm_source=CV101AVBlogBanner&utm_medium=Stickybanner2utm_campaign=CV101banner)

DATA SCIENTISTS

Blog (<https://www.analyticsvidhya.com/blog/>)
 Hackathon (<https://trainings.analyticsvidhya.com/>)
 Discussions (<https://datahack.analyticsvidhya.com/>)
 Apply Jobs (<https://www.analyticsvidhya.com/careers/>)
 Leaderboard (<https://datahack.analyticsvidhya.com/leaderboard/>)

COMPANIES

Post Jobs

(<https://www.analyticsvidhya.com/corporate/>)

(<https://www.analyticsvidhya.com/blog/>)

(<https://trainings.analyticsvidhya.com/>)

(<https://datahack.analyticsvidhya.com/>)

(<https://www.analyticsvidhya.com/>)

(<https://discuss.analyticsvidhya.com/>)

(<https://datahack.analyticsvidhya.com/>)

(<https://www.analyticsvidhya.com/careers/>)

(<https://www.analyticsvidhya.com/contact/>)

(<https://www.analyticsvidhya.com/>)

(<https://datahack.analyticsvidhya.com/leaderboard/>)

(<https://www.analyticsvidhya.com/contact/>)

(<https://www.analyticsvidhya.com/>)

JOIN OUR COMMUNITY :

f

(<https://www.analyticsvidhya.com/corporate/>)

(<https://www.analyticsvidhya.com/blog/>)

(<https://trainings.analyticsvidhya.com/>)

(<https://datahack.analyticsvidhya.com/>)

(<https://www.analyticsvidhya.com/>)

(<https://discuss.analyticsvidhya.com/>)

(<https://datahack.analyticsvidhya.com/>)

(<https://www.analyticsvidhya.com/careers/>)

(<https://www.analyticsvidhya.com/contact/>)

(<https://www.analyticsvidhya.com/>)

(<https://datahack.analyticsvidhya.com/leaderboard/>)

(<https://www.analyticsvidhya.com/contact/>)

(<https://www.analyticsvidhya.com/>)

(<https://plus.google.com/+AnalyticsVidhya>)

(<https://plus.google.com/+AnalyticsVidhya>)

(<https://plus.google.com/+AnalyticsVidhya>)

(<https://plus.google.com/+AnalyticsVidhya>)

(<https://plus.google.com/+AnalyticsVidhya>)

(<https://plus.google.com/+AnalyticsVidhya>)

(<https://plus.google.com/+AnalyticsVidhya>)

(<https://plus.google.com/+AnalyticsVidhya>)

(<https://plus.google.com/+AnalyticsVidhya>)

(<https://plus.google.com/+AnalyticsVidhya>)

(<https://plus.google.com/+AnalyticsVidhya>)

(<https://plus.google.com/+AnalyticsVidhya>)

(<https://plus.google.com/+AnalyticsVidhya>)

(<https://plus.google.com/+AnalyticsVidhya>)

(<https://plus.google.com/+AnalyticsVidhya>)

(<https://plus.google.com/+AnalyticsVidhya>)

(<https://plus.google.com/+AnalyticsVidhya>)

(<https://plus.google.com/+AnalyticsVidhya>)

(<https://plus.google.com/+AnalyticsVidhya>)

(<https://plus.google.com/+AnalyticsVidhya>)

(<https://plus.google.com/+AnalyticsVidhya>)

(<https://plus.google.com/+AnalyticsVidhya>)

(<https://plus.google.com/+AnalyticsVidhya>)

(<https://plus.google.com/+AnalyticsVidhya>)

(<https://plus.google.com/+AnalyticsVidhya>)

(<https://plus.google.com/+AnalyticsVidhya>)

(<https://plus.google.com/+AnalyticsVidhya>)

(<https://plus.google.com/+AnalyticsVidhya>)

(<https://plus.google.com/+AnalyticsVidhya>)

(<https://plus.google.com/+AnalyticsVidhya>)

Subscribe to emailer

>

© Copyright 2013-2018 Analytics Vidhya.

Privacy Policy (<https://www.analyticsvidhya.com/privacy-policy/>)

Don't have an account? Sign up ([https://www.analyticsvidhya.com/sign-up/](#))

Terms of Use (<https://www.analyticsvidhya.com/terms/>)

Refund Policy (<https://www.analyticsvidhya.com/refund-policy/>)

Subscribe!