

⑦ Some things to Read & Consider.

VD Prep

⇒ Steps to Data Exploration & Preparation.

- Variable Identification
- Univariate Analysis
- Bi-variate Analysis
- Missing Values Treatment
- Outlier Treatment
- Variable Transformation] Feature Engineering
- Variable Creation.

Generalized
EDA

[We need to iterate over steps 4-7 times to get refined model]

① Predictor Variables → Input Variables

Univariate Analysis for -

→ Continuous Variables

We need to understand the Central
tendency & the spread of the variable

② Note - Univariate analysis is also performed to highlight
the missing & outlier value.

→ Categorical Variable

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 & Count% against each category. Bar chart can be used as Visualization.

Bi-variate - relation b/w two variables.

We look for association & dissociation b/w two variables at pre-significance level.

b/w -

① Continuous & Continuous

To check relation b/w two continuous variable we use scatter plot.

② Categorical & Categorical

→ Two-way table of count & count%.

The rows represent category of one variable & the column represents category of the other variable.

→ Stacked-Column chart

Stacked-bar plot, visual of a two way table,

two way table \leftrightarrow cross-tab, sns.countplot(hue=?)

→ Chi-Square test

It is used to derive the statistical significance b/w two variables. It returns probability for the computed chi-square distribution with the degree of freedom -

Probability \rightarrow 0 - are dependent, 1 - are independent

Note

The chi-square test statistic for a test of independence of two categorical variables is found by

$$\chi^2 = \sum \frac{(O-E)^2}{E}$$

where O - observed frequency

E - expected frequency under the null hypothesis -

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

→ (Most freq) (Box plots for smaller levels)

③ Categorical & Continuous

We can draw box-plots for each level of categorical variable.

If levels are smaller in number it'll not show 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.

$$z = \frac{|\bar{x}_1 - \bar{x}_2|}{\sqrt{\frac{s_1^2}{n_1} + \frac{s_2^2}{n_2}}}$$

If the probability of Z is small, then the difference b/w the averages is more significant.

The t -test is very similar to Z -test but it is used when no. of observations for both categories is less than 30.

$$t = \frac{|\bar{x}_1 - \bar{x}_2|}{\sqrt{\frac{s^2}{n_1} + \frac{s^2}{n_2}}}$$

$$s^2 = \frac{(n_1-1)s_1^2 + (n_2-1)s_2^2}{n_1 + n_2 - 2}$$

$\bar{x}_1, \bar{x}_2 \rightarrow$ averages

$s_1^2, s_2^2 \rightarrow$ variances of sub 1 & 2

$n_1, n_2 \rightarrow$ counts

$t \rightarrow$ has t distribution with $n_1 + n_2 - 2$ degree of freedom

ANOVA

It assesses whether the avg. of more than two groups is statistically different.

处理 the missing values

① Deletion

details are in
PC/mobile

② Mean / Mode / Median Imputation

③ Prediction model - predictive model to estimate missing data.

④ KNN Imputation

3
#

Outliers Detection & treatment.

6
(cont.)

Types

- Artificial / Non-natural
- Natural

Examples -

- Data Entry Errors (Human)
- Measurement Errors (Instrumental)
- Experimental Error &
- Intentional Outliers (Self reported measure)
- Data Processing Errors (Data mining)
- ~~Simple Sampling Error~~
- Natural Outliers.

Impact due to Outliers.

- It increases the lower variance and reduces the power of Statistical test.
- 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 statistical model assumptions.

Detection of Outliers.

Visualization using -

- Box plots
- Histograms
- Scatter plots.

Remaining Outliers

- Deleting Observations
 - Transforming and binning values.
 - [eg → Natural log of a value reduces the variation caused by extreme values.]
 - Imputing. (If artificial outliers, then we can impute it, ~~else~~)
 - Treat Separately. just like missing values
- If there are significant number of outliers we need to treat them separately in the statistical model.

final Stage - Feature Engineering. (if)

We perform feature Engineering after completing the 5 steps -

- Variable Identification
- Univariate, Bivariate Analysis
- Missing values imputation
- Outliers treatment
- feature Engineering

Variable transformation. Variable/feature creation.

Variable Transformation

It refers to the replacement of a variable by a function. for instance, replacing a variable by square, cube or logarithm transformation.

When to use these transformations?

- + When we want to change the scale of a variable or standardize the value of a variable for better understanding.
- + When we can transform the complex non-linear relationships into the linear relationships.
- + Log transform is commonly used in this situation.
- + Symmetric distribution is preferred over a skew distribution as it is easier to interpret and generate inferences. It is used to reduce Skewness.

Note - for right skewed, we take square / cube root or logarithm for variable.

- For left skewed, we take square / cube or exponential of variables.

Common Methods of transformation

- Logarithm
- Square / Cube roots
- Binning - Used to categorize Variable. It is performed on original value, percentile or frequency.

Variable / feature Creation

feature variable generation is a process to generate new variables / feature based on existing variables.

for eg, if we have date (dd-mm-yy) as an input variable in a dataset. We can generate new variables like day, month, year, week, weekday that may have better relationship with target variable.

④ Various techniques to create new features.

1) Creating derived Variables.

This refers to creating new variables from existing variable using Set of functions or different methods.

2) Creating dummy Variables.

One of the most common application of dummy variable is to convert categorical to numerical Variable.

Dummy Variables are also called Indicator variables. It is useful to take categorical variable as a predictor in statistical models.