

Assignment No. 2

Title - Data Wrangling 2

Problem statement

Create an 'Academic performance' dataset of students and perform following operations using python.

- 1) Scan all variables for missing values inconsistencies. If there are missing values and/or inconsistencies use any of the suitable techniques to deal with it.
- 2) Scan all numeric variables for outliers, If there are outliers use any of the suitable technique to deal with it.
- 3) Apply the transformations on at least one of the variables. The purpose of this transformation should be one of the following reason. to change the scale for better understanding of the variable, to convert a nonlinear relation into linear one, or to decrease the skewness and convert the distribution into a normal distribution.

Learning Objectives

- 1) To understand and implement various data processing and outliers handling techniques.
- 2) Understand and implement various data transformation techniques to change scalar or convert distribution to a normal distribution.

Learning outcomes

By performing this assignment, we will be able to

- 1) Various data preprocessing and outliers handling techniques.
- 3) various data transformation techniques to change scale or convert the distribution to normal distribution.

Software Requirements

- 1) python 3.10 (with pandas, numpy, matplotlib, seaborn, scipy).
- 2) jupyter Notebook

Hardware Requirement

- 1) Intel i5 10th generation
- 2) 8GB RAM

Theory

Data Wrangling

- Data wrangling is the process of cleaning and unifying messy and complex data sets for easy access and analysis.
- This process typically includes manually converting and mapping data from one raw form into another format to allow for more convenient consumption and organisation of the data

package used

- 1) pandas
 - It is a software library in python for data manipulation and analysis.



PICT, PUNE

2) Numpy

It is a library in python used to add support for large collection multidimensional arrays and matrices along with large collection of high level mathematical function to operate on the arrays.

3) scipy-stats

It is a library used to perform advanced statistical operations on dataset.

4) Matplotlib

- It is a library used for plotting simple graphs and also for data visualization to detect outliers, data transformation visualization etc.

5) seaborn

- It is a library used to plot detailed and visually good looking graphs for more advanced data visualization to detect outliers and check if data transformation happened correctly.

Data transformations

Skewness is defined as the measure of asymmetry of the probability distribution of the real valued random variable (feature/ attributes)

the skewness outside the range of -1 to 1 are highly skewed and those in the range of -0.5 to 0.5 are symmetrical or normal, remaining are moderately skewed.

- to convert the distribution to normal distribution means we need to change the skewness and take it in the range -0.5 to 0.5 .

- Some of the data transformation techniques are.

1) logarithmic transformation

2) Reciprocal transformation

3) Square root transformation

4) Exponential transformation

5) Box-cox transformation

- The data transformation techniques can only be applied on positive values.

Detect and eliminate outliers.

- outliers are the data objects that deviate significantly from rest of the data objects and behaves in different manner.

- They are detected on the basis of the kind of dataset attributes we are dealing with.

- Some of the ways to detect outliers are.

1) Using Boxplot or scatter plot (Data Visualization)

2) Using Inter-quartile range

3) Using z score

4) Using standard deviation

Functions used

1) `df.mean()`

- returns the mean of each numeric attribute

2) `df.max()`

- for calculating the maximum numeric value



PICT, PUNE

3) `df.quantile`

- returns the values at the given quantile over requested axis.

4) `df.skew`

- returns the skewness of the values of numeric attributes dataframe `df`.

5) `seaborn.kdeplot(df['001'])`

- plots a KDE plot for column '001'.

6) `nplog(df['001'])`

returns the dataframe array object of logarithmic values.

7) `scipy.stats.boxcox(df['001'])`

- returns the transformed values of each numeric record of column '001' of dataframe 'df' so that the transformed values form a normal distribution after boxcox transformation is applied.

Observations

1) out of the two plots to observe the outliers, boxplot is preferred for univariate while scatterplot is preferred for bivariate / multivariate outliers.

2) out of all the data transformation techniques, logarithmic and reciprocal transformation gives most accurate results.

3) the data transformations techniques can only be applied to positive data points.

Conclusion

Hence, after performing this assignment, we learnt various processes and techniques to handle outliers and apply data transformations to convert distribution to normal distribution, we successfully implemented the techniques on provided dataset.