

Name : Tejas V. Redkar

Roll No : PC-44

PRN : 1032210937

Batch : C2, Panel - C

## DFC Assignment - 2

- \* Title : Data Pre-processing using Python.
- \* Aim : Preprocess data using Python.
- \* Problem Statement : Data Pre-processing using Python
- \* Objective :
  - To clean Data & make it noise free.
  - Prepare Data for Analysis.
- \* Conclusion : Understood the steps for data pre-processing.
- \* FAQ's

(Q1) Why data cleaning is important?

Ans Data cleaning, also known as data cleansing or data scrubbing, is a crucial step in the data preparation process for various reasons:

i) Accurate analysis : Clean data ensures that the

information you analyze & use for decision-making is accurate & reliable. Dirty data, containing errors & inconsistencies can lead to incorrect conclusions, affecting the quality of insights & decisions.

- 2) Data Quality: Data quality is a key factor in the success of data-driven projects. High-quality data is more likely to yield accurate & meaningful results, while poor-quality data can lead to costly errors & misinterpretation.
- 3) Data Integration: In many cases, data comes from various sources & systems. Cleaning the data ensures that it is consistent and compatible for integration, making it easier to combine & analyze data from different sources.
- 4) Reducing Errors: Dirty data can contain errors such as duplicate rows or records, missing values, incorrect formatting, & outliers. Cleaning the data helps identify & correct these errors, reducing the chances of making decisions based on flawed information.
- 5) Improved Efficiency: Clean data is easy to work with. Data Analysts & scientists spend less time dealing with errors & inconsistencies, allowing them to focus more on analysis & deriving ~~rights~~ insights.
- 6) Enhancing data Visualization: Clean data is essential for data visualization. Visualizations

are often used to communicate insights, & they rely on accurate & consistent data to be effective.

7) Compliance & Regulation : In some industries, compliance with data protection & privacy regulations is mandatory. Clean data helps ensure that sensitive information is handled correctly & that organizations remain compliant with legal requirements.

8) Customer Satisfaction : If your data pertains to customers, ensuring data accuracy is vital for maintaining trust & delivering a better customer experience. Incorrect or inconsistent customer data can lead to communication errors & dissatisfaction.

9) Cost Savings : Correcting errors & inconsistencies in data can be costly, especially if these errors lead to poor decisions or operational inefficiencies. Investing in data cleaning upfront can lead to significant cost savings in the long run.

10) Predictive Modelling : For machine learning & predictive modelling, clean data is essential. Models trained on dirty data are likely to perform poorly & may not provide reliable predictions.

11) Strategic Decision-Making : Businesses & organizations rely on data to make strategic decisions. Clean data ensures that these decisions

are based on accurate & trust worthy information, which can have a substantial impact on an organization's success.

Q2) Why we need splitting

Ans Splitting data is fundamental & crucial step in data analysis for several reasons:

1) Training & Testing: In machine learning & predictive modelling, you typically split your data into two sets: a training set & a testing set. The training set is used to build & train your model, while the testing set is used to evaluate its performance. This separation helps you assess how your model generalizes to unseen data, preventing overfitting & providing a more accurate estimate of its performance.

2) Validation: In addition to training & testing sets, you may also use a validation set during model development. This helps you fine-tune hyperparameters & make decisions about model selection & feature engineering without leaking information from the test set into the training process.

3) Cross-validation: In situations where you have limited data, cross-validation techniques, such as K-fold cross-validation, are used to split data into multiple subsets (folds). This

allows you to train & test the model on different subsets of the data multiple times, providing a more robust assessment of model performance.

- 4) **Avoiding Data Leakage**: When performing feature engineering or data preprocessing, it's essential to split the data before applying any transformations. If you perform these operations on the entire dataset without splitting, you risk introducing data leakage, where information from the test set inadvertently influences the training process.
- 5) **Exploratory Data Analysis (EDA)**: During exploratory data analysis, you may want to split the data to examine patterns, relationships, & distributions in subsets of the data. This can help you gain insights into specific segments of your data & make informed decisions about subsequent analyses.
- 6) **Comparative analysis**: Splitting data can also be valuable when you want to compare different subsets of your data. For example, you might want to compare performance of different customer segments, time periods, or geographical regions. Splitting the data into these subsets allows you to analyze them separately & draw meaningful comparisons.
- 7) **Anomaly Detection**: when identifying outliers or anomalies in data, splitting it into normal values

and abnormal subsets can make it easier to detect unusual patterns or values in the abnormal subset.

⑧ Sampling: In large datasets, splitting can also be used to create smaller, manageable samples for initial exploration or analysis. Random sampling methods ensure that the sample is representative of the entire dataset.

⑨ Privacy & Security: In some cases, data may be need to be split to maintain privacy & security. Sensitive information may be partitioned to restrict access to certain subsets of data, protecting individuals privacy & complying with data protection regulations.

Q3) what is dummy variables?

Ans: Dummy variables, also known as indicator variables or binary variables, are a way to represent categorical data in statistical & regression analysis. Categorical data are variables that can take on a limited, fixed number of distinct values or categories. Since many statistical & machine learning algorithms require numerical inputs, dummy variables are used to convert categorical data into a format that can be easily incorporated into models.

Q4) What are the criteria to identify an outlier

Ans Identifying outliers in a dataset is an important step in data analysis & can help in understanding data quality, detecting errors, & potentially revealing interesting patterns or anomalies. There are several criteria & techniques to identify outliers, & the choice of method depends on the nature of your data & your specific goals. Here are some common criteria & techniques used to identify outliers:

- 1) Visual inspection
- 2) Z-score
- ✓ 3) Modified Z-score
- 4) Interquartile range (IQR)
- 5) Tukey's Fences
- 6) Mahalanobis Distance
- 7) Visualizing Residuals
- 8) Domain Knowledge
- 9) Machine learning algorithms
- 10) Statistical Tests.

Q5) What is the reason for an outlier to exist in a dataset?

Ans Outliers can exist in a dataset for various reasons & understanding the underlying causes of outliers is important for effective data analysis.

Here are some common reasons why outliers may

occur:

- 1) Measurement Errors
- 2) Natural Variation
- 3) Data Entry Errors
- 4) Data Transformation
- 5) Sampling Errors
- 6) Genuine Anomalies
- 7) Data Integrity Issues
- 8) Changes in data generation process
- 9) Data collection Bias
- 10) Extreme events

(Q6) What are the impacts of having outliers in a dataset?

Ans The presence of outliers in a dataset can have various impacts which can both be positive & negative, depending on the context & the goals of your analysis. Here are some key impacts of having outliers in a dataset:

Negative Impacts:

1) Skewed descriptive analysis:

Outliers can significantly affect summary statistics like the mean & standard deviation, causing them to be biased or misleading. The mean, in particular, is sensitive to extreme values.

2) Distorted Data Distribution:

- Outliers can distort the shape & interpretation of data distributions. For example, in a positively skewed distribution with outliers on the right side, the distribution may appear more skewed than it actually is.

### 3) Inaccurate models:

- When outliers are present, statistical models that assume normality or homoscedasticity (constant variance) may produce biased or inaccurate results. This can impact regression analysis & other modelling techniques.

### 4) Reduced Predictive Accuracy:

- In machine learning models, outliers can lead to overfitting. The model may learn to overly accommodate the extreme values, resulting in reduced generalization to new data.

### 5) Misleading Visualizations:

- Outliers can make visualizations less informative or misleading. For example, a scatter plot with outliers might obscure patterns in the main cluster of data points.

### 6) Inflated Error Metrics:

- Outliers can inflate error metrics such as the Mean Absolute Error (MAE) or Root Mean Square Error (RMSE), making model performance seem worse than it actually is.

Positive Impacts:

### 1) Detection of Anomalies:

- Outliers can be indicative of anomalies or rare events in the data. Detecting these anomalies is essential in fields like fraud detection, network security, & quality control.

### 2) Identification of interesting insights:

- Outliers can reveal interesting & valuable insights about the data. They might represent unique phenomena or exceptional cases that warrant further investigation.

### 3) Improved Robustness:

- In some cases, robust statistical methods that are less sensitive to outliers can be employed to provide more reliable estimates & models.

### 4) Increased Data Quality Awareness:

- The presence of outliers can draw attention to data quality issues & the need for data cleaning, leading to improved data integrity.

### 5) Enhanced Model Robustness:

- In machine learning, some algorithms are designed to be robust to outliers, which can help in cases where outliers are informative rather than erroneous.

(Q7) Discuss measuring techniques of the dispersion of data.

**Ans** Measuring the dispersion of data is crucial in statistics & data analysis because it helps us understand how spread out our variable a dataset is. Dispersion measures provide valuable insights into the variability & distribution of data points. There are several common techniques for measuring dispersion:

1) Range:

- The range is the simplest value measure of dispersion & is calculated by subtracting the minimum value from the maximum value in the dataset.
- Range = Max value - Min value
- It provides a rough idea of the spread but is highly influenced by outliers & may not give a complete picture of the distribution.

2) Inter Quartile Range (IQR)

3) Variance

4) Standard Deviation

5) Coefficient of variation

6) Mean Absolute Deviation

7) Range - Based Measures

8) Coefficient of Quartile Deviation

9) Z-score.

