# Experiment # 8: Implement various Data preprocessing techniques on a given data set

**Aim/Objective:**

This experiment aims to implement data pre-processing techniques to clean, transform, and prepare raw data for further analysis or machine learning tasks

**Description:**

In this experiment, students will learn the importance of data pre-processing in the data science workflow. They will understand the various steps involved in cleaning and transforming raw data to make it suitable for analysis or model building. Students will implement a data pre-processing pipeline using Python and relevant libraries, gaining hands-on experience in handling missing values, outliers, categorical variables, feature scaling, and more.

**Pre-Requisites:**

Basic understanding of data types, including numerical and categorical variables.

Familiarity with Python programming and data manipulation libraries such as pandas

**Pre-Lab:**

1. Why data are dirty?
2. What is data preprocessing? Why is it important in machine learning?
3. What are some common problems that occur during data processing? How can they be fixed?
4. How do you handle the missing data?
5. What is the difference between missing value treatment and outliers treatment?

① Data is considered "dirty" when it contain errors, inconsistencies, coy is incomplete leading to ear inaccurate dy unrealible result in analysis!-

reasons data become dirty:- !) Missing values
ii) Duplicates
iii) Incorrect data

② Data Preprocessing is a Process of cleaning & transforming raw data before feeding it into a machine learning model. It involves handling missing values, normalizing data, encoding categorical variables.

③ Common Problems:- i) Missing Data
   ii) Outliers
   iii) Inconsistent Data
   iv) Noise
   v) Duplicate data.

④ Handling missing data can be done in several ways:- i) Removal
   ii) Imputation
   iii) Indicator for Missing values.

(5.)

Missing outlier treatment:- It focus on filling in (or) dealing with missing data points.

Outlier treatment:- Deals with extreme (or) unusual data points that may skew the model results.

**In-Lab:**

Implement a Python program to find and impute the missing data in the following dataset.

Dataset Link:
https://www.kaggle.com/datasets/bharatnatrayn/movies-dataset-for-feature-extracion-prediction/data?select=movies.csv

**Procedure/Program:**

```
import Pandas as id

from    Sklearn.impute import simpleImputer

    df = Pd. read-civ ('movies. csv)
    Print ("missing values before imputation:):
    Print (df.isnull(). sum())

numerical-colowns = df. spect.dtypes (include= (flout64, int64]

    imputer-num= simple Impu ten (strategy = 'mean')

df [numerical-columns] = inputer-num. fit-transform

imputer-cat = simple Inputor (strategy = most-
                                        freqvent).
```

```python
Print ("missing values after imputation"):
 Print (df. isnull().sum())
df.to-csv('movies.imputed. csv', index=false)
```

- Data and Results:

the dataset is loaded from csv files using Pandas. The Program Points the number of missing values for each column.

- Analysis and Inferences:

For numerical values, missing values are imputed using mean.

The imputed dataset is saved as movies_imputed.csv

## VIVA-VOCE Questions (In-Lab):

1. What is the difference between normalization and standardization?
2. What are the different encoding techniques for categorical data?
3. What are some common techniques for data reduction?
4. How do you preprocess time-series data?
5. What is data integration and what challenges are associated with it?

① Normalization :- Rescales data to a range of [0, 1], useful when data needs to be on same scale for comparison.

standardization :- Rescal data to have a mean of 0 and a standard deviation of 1.

② encoding techniques :- i) label encoding
ii) one-Hot Encoding
iii) ordinal Encoding
iv) target Encoding

③ Techniques:- i) Principle component Analysis
　　　　　ii) feature selection
　　　　　iii) Sampling
　　　　　iv) Aggregation

④ Preprocessing time-series data-:
　　　　　i) Handling Missing data
　　　　　ii) Smoothing
　　　　　iii) Resampling
　　　　　iv) Differencing
　　　　　v) Normalization

⑤ Data integration:- combining data from
different sources into a unfied dataset

challenges:- i) schema Mismatch
　　　　　ii) Data Quality
　　　　　iii) Duplicate data
　　　　　iv) Scaling

**Post-Lab:**

Implement a Python program to apply various data preprocessing techniques on the following dataset.

Dataset Link:

https://catalog.data.gov/dataset/electric-vehicle-population-data/resource/fa51be35-691f-45d2-9f3e-535877965e69

**Procedure/Program:**

```
import Pandas as pid.

from sklearn.impute impart SimpleImputer

from sklearn.Preprocening impor labelEncoder

from sklearn.model-selection import train.ksl-split.

url = 'https://data.wa.gov/api/views/frwt-92d2/
                        rows.csv? acesstype = Download

df = df.read_csv(url)
Print (Intial dataset information:)
Print (df.info())

numerical_columns = df.select_dtypes(include = int y)
    inputer-num = SimpleImputer(strategy = mean)

categorical_columns = df.select_types(include = ['object'])
```

```
imputer_cat = Simple Imputer(strategy= most_frequent)
Print("\n missing value after imputation:")
Print(df. isnull(). sum())

label_encoder= Label Encoder()

for column in categorical_columns:
    df[column]=label_encoder. fit_transform(df[column])

Print("\n first 5 rows after encoding:")
    Print(df. head()).
```

| Experiment # | | | Student ID | |
|---|---|---|---|---|
| Date | | | Student Name | |

**Data and Results:**

Imputed missing values using mean for numerical data and the most frequent value for categorical data.

**Analysis and Inferences:**

The preprocessing steps ensured the dataset was free of missing values and duplicates.