# Exp:6 Handling Missing and Inappropriate Data in a Dataset

Aim: Demonstrate an experiment   to handle missing data and inappropriate data in a Data set using Python Pandas Library for Data Preprocessing.

Dataset Given:

**Hotel.csv**

| CustomerID | Age_Group | Rating(1-5) | Hotel | FoodPreference | Bill | NoOfPax | EstimatedSalary | Age_Group |
|---|---|---|---|---|---|---|---|---|
| 1 | 20-25 | 4 | Ibis | veg | 1300 | 2 | 40000 | 20-25 |
| 2 | 30-35 | 5 | LemonTree | Non-Veg | 2000 | 3 | 59000 | 30-35 |
| 3 | 25-30 | 6 | RedFox | Veg | 1322 | 2 | 30000 | 25-30 |
| 4 | 20-25 | -1 | LemonTree | Veg | 1234 | 2 | 120000 | 20-25 |
| 5 | 35+ | 3 | Ibis | Vegetarian | 989 | 2 | 45000 | 35+ |
| 6 | 35+ | 3 | Ibys | Non-Veg | 1909 | 2 | 122220 | 35+ |
| 7 | 35+ | 4 | RedFox | Vegetarian | 1000 | -1 | 21122 | 35+ |
| 8 | 20-25 | 7 | LemonTree | Veg | 2999 | -10 | 345673 | 20-25 |
| 9 | 25-30 | 2 | Ibis | Non-Veg | 3456 | 3 | -99999 | 25-30 |
| 9 | 25-30 | 2 | Ibis | Non-Veg | 3456 | 3 | -99999 | 25-30 |
| 10 | 30-35 | 5 | RedFox | non-Veg | -6755 | 4 | 87777 | 30-35 |

**About Dataset:**

No.of Columns =9 (called as series – CustomerID, Age_Group, Rating(1-5),Hotel, FoodPreference, Bill, NoOfPax, EstimatedSalary)

CutomerID: Numerical Continuous data

Age: Categorical Data

Rating (1-5): Numerical Discrete Data

Hotel: Categorical Data

Food: Categorical Data

Bill: Numerical Continuous data

NoOfPax: Numerical Discrete

EstimatedSalary: Numerical Continuous data

# Python Code:

# Upload Hotel.csv and convert it into dataFrame

**import numpy as np**

**import pandas as pd**

**df=pd.read_csv("Hotel_Dataset.csv")**

**df**

| | CustomerID | Age_Group | Rating(1-5) | Hotel | FoodPreference | Bill | NoOfPax | Estimated Salary | Age_Group.1 |
|---|---|---|---|---|---|---|---|---|---|
| 0 | 1 | 20-25 | 4 | Ibis | veg | 1300 | 2 | 40000 | 20-25 |
| 1 | 2 | 30-35 | 5 | LemonTree | Non-Veg | 2000 | 3 | 59000 | 30-35 |
| 2 | 3 | 25-30 | 6 | RedFox | Veg | 1322 | 2 | 30000 | 25-30 |
| 3 | 4 | 20-25 | -1 | LemonTree | Veg | 1234 | 2 | 120000 | 20-25 |
| 4 | 5 | 35+ | 3 | Ibis | Vegetarian | 989 | 2 | 45000 | 35+ |
| 5 | 6 | 35+ | 3 | Ibys | Non-Veg | 1909 | 2 | 122220 | 35+ |
| 6 | 7 | 35+ | 4 | RedFox | Vegetarian | 1000 | -1 | 21122 | 35+ |
| 7 | 8 | 20-25 | 7 | LemonTree | Veg | 2999 | -10 | 345673 | 20-25 |
| 8 | 9 | 25-30 | 2 | Ibis | Non-Veg | 3456 | 3 | -99999 | 25-30 |
| 9 | 9 | 25-30 | 2 | Ibis | Non-Veg | 3456 | 3 | -99999 | 25-30 |
| 10 | 10 | 30-35 | 5 | RedFox | non-Veg | -6755 | 4 | 87777 | 30-35 |

#From the dataframe identify the duplicate row(i.e row 9)

# The duplicated() method returns a Series with True and False values that describe which rows in the DataFrame are duplicated and not.

**df.duplicated()**

```
0      False
1      False
2      False
3      False
4      False
5      False
6      False
7      False
8      False
9       True
10     False
dtype: bool
```

# The info() method prints information about the DataFrame. The information contains the number of columns, column labels, column data types, memory usage, range index, and the number of cells in each column (non-null values).

**df.info()**

```
<class 'pandas.core.frame.DataFrame'>
RangeIndex: 11 entries, 0 to 10
Data columns (total 9 columns):
 #   Column          Non-Null Count  Dtype
---  ------          --------------  -----
 0   CustomerID      11 non-null     int64
 1   Age_Group       11 non-null     object
 2   Rating(1-5)     11 non-null     int64
 3   Hotel           11 non-null     object
 4   FoodPreference  11 non-null     object
 5   Bill            11 non-null     int64
 6   NoOfPax         11 non-null     int64
 7   EstimatedSalary 11 non-null     int64
 8   Age_Group.1     11 non-null     object
dtypes: int64(5), object(4)
memory usage: 924.0+ bytes
```

# The drop_duplicates() method removes duplicate rows.

**df.drop_duplicates(inplace=True)**

**df**

|    | CustomerID | Age_Group | Rating(1-5) | Hotel | FoodPreference | Bill | NoOfPax | Estimated Salary | Age_Group.1 |
|----|-----------|-----------|-------------|-----------|----------------|------|---------|------------------|-------------|
| 0  | 1 | 20-25 | 4 | Ibis | veg | 1300 | 2 | 40000 | 20-25 |
| 1  | 2 | 30-35 | 5 | LemonTree | Non-Veg | 2000 | 3 | 59000 | 30-35 |
| 2  | 3 | 25-30 | 6 | RedFox | Veg | 1322 | 2 | 30000 | 25-30 |
| 3  | 4 | 20-25 | -1 | LemonTree | Veg | 1234 | 2 | 120000 | 20-25 |
| 4  | 5 | 35+ | 3 | Ibis | Vegetarian | 989 | 2 | 45000 | 35+ |
| 5  | 6 | 35+ | 3 | Ibys | Non-Veg | 1909 | 2 | 122220 | 35+ |
| 6  | 7 | 35+ | 4 | RedFox | Vegetarian | 1000 | -1 | 21122 | 35+ |
| 7  | 8 | 20-25 | 7 | LemonTree | Veg | 2999 | -10 | 345673 | 20-25 |
| 8  | 9 | 25-30 | 2 | Ibis | Non-Veg | 3456 | 3 | -99999 | 25-30 |
| 10 | 10 | 30-35 | 5 | RedFox | non-Veg | -6755 | 4 | 87777 | 30-35 |

#While removing duplicate record row index also removed

#  The len() function to return the length of an object. With a dataframe, the function returns the number of rows.

**len(df)**

10

#Reset the index

**index=np.array(list(range(0,len(df))))**

**df.set_index(index,inplace=True)**

**index**

```
array([0, 1, 2, 3, 4, 5, 6, 7, 8, 9])
```
**df**

| | CustomerID | Age_Group | Rating(1-5) | Hotel | FoodPreference | Bill | NoOfPax | Estimated Salary | Age_Group.1 |
|---|---|---|---|---|---|---|---|---|---|
| 0 | 1 | 20-25 | 4 | Ibis | veg | 1300 | 2 | 40000 | 20-25 |
| 1 | 2 | 30-35 | 5 | LemonTree | Non-Veg | 2000 | 3 | 59000 | 30-35 |
| 2 | 3 | 25-30 | 6 | RedFox | Veg | 1322 | 2 | 30000 | 25-30 |
| 3 | 4 | 20-25 | -1 | LemonTree | Veg | 1234 | 2 | 120000 | 20-25 |
| 4 | 5 | 35+ | 3 | Ibis | Vegetarian | 989 | 2 | 45000 | 35+ |
| 5 | 6 | 35+ | 3 | Ibys | Non-Veg | 1909 | 2 | 122220 | 35+ |
| 6 | 7 | 35+ | 4 | RedFox | Vegetarian | 1000 | -1 | 21122 | 35+ |
| 7 | 8 | 20-25 | 7 | LemonTree | Veg | 2999 | -10 | 345673 | 20-25 |
| 8 | 9 | 25-30 | 2 | Ibis | Non-Veg | 3456 | 3 | -99999 | 25-30 |
| 9 | 10 | 30-35 | 5 | RedFox | non-Veg | -6755 | 4 | 87777 | 30-35 |

# Axis refers to the dimensions of a DataFrame (index and columns) or Series (index only) Use axis=0 to apply functions row-wise along the index. Use axis=1 to apply functions column-wise across columns.

**df.drop(['Age_Group.1'],axis=1,inplace=True)**

**df**

| | CustomerID | Age_Group | Rating(1-5) | Hotel | FoodPreference | Bill | NoOfPax | Estimated Salary |
|---|---|---|---|---|---|---|---|---|
| 0 | 1 | 20-25 | 4 | Ibis | veg | 1300 | 2 | 40000 |
| 1 | 2 | 30-35 | 5 | LemonTree | Non-Veg | 2000 | 3 | 59000 |
| 2 | 3 | 25-30 | 6 | RedFox | Veg | 1322 | 2 | 30000 |
| 3 | 4 | 20-25 | -1 | LemonTree | Veg | 1234 | 2 | 120000 |
| 4 | 5 | 35+ | 3 | Ibis | Vegetarian | 989 | 2 | 45000 |
| 5 | 6 | 35+ | 3 | Ibys | Non-Veg | 1909 | 2 | 122220 |
| 6 | 7 | 35+ | 4 | RedFox | Vegetarian | 1000 | -1 | 21122 |
| 7 | 8 | 20-25 | 7 | LemonTree | Veg | 2999 | -10 | 345673 |
| 8 | 9 | 25-30 | 2 | Ibis | Non-Veg | 3456 | 3 | -99999 |
| 9 | 10 | 30-35 | 5 | RedFox | non-Veg | -6755 | 4 | 87777 |

# The function . loc is typically used for label indexing and can access multiple columns.

**df.CustomerID.loc[df.CustomerID<0]=np.nan**

**df.Bill.loc[df.Bill<0]=np.nan**

**df.EstimatedSalary.loc[df.EstimatedSalary<0]=np.nan**

**df**

```
C:\Users\Ayyadurai\AppData\Local\Temp\ipykernel_5300\2580639570.py:1: S
ettingWithCopyWarning:
A value is trying to be set on a copy of a slice from a DataFrame
See the caveats in the documentation: https://pandas.pydata.org/pandas-
docs/stable/user_guide/indexing.html#returning-a-view-versus-a-copy
  df.CustomerID.loc[df.CustomerID<0]=np.nan
df.Bill.loc[df.Bill<0]=np.nan
C:\Users\Ayyadurai\AppData\Local\Temp\ipykernel_5300\2580639570.py:2: S
ettingWithCopyWarning:
A value is trying to be set on a copy of a slice from a DataFrame
df.EstimatedSalary.loc[df.EstimatedSalary<0]=np.nan
```

| | CustomerID | Age_Group | Rating(1-5) | Hotel | FoodPreference | Bill | NoOfPax | Estimated Salary |
|---|---|---|---|---|---|---|---|---|
| 0 | 1.0 | 20-25 | 4.0 | Ibis | veg | 1300.0 | 2 | 40000.0 |
| 1 | 2.0 | 30-35 | 5.0 | LemonTree | Non-Veg | 2000.0 | 3 | 59000.0 |
| 2 | 3.0 | 25-30 | NaN | RedFox | Veg | 1322.0 | 2 | 30000.0 |
| 3 | 4.0 | 20-25 | NaN | LemonTree | Veg | 1234.0 | 2 | 120000.0 |
| 4 | 5.0 | 35+ | 3.0 | Ibis | Vegetarian | 989.0 | 2 | 45000.0 |
| 5 | 6.0 | 35+ | 3.0 | Ibys | Non-Veg | 1909.0 | 2 | 122220.0 |
| 6 | 7.0 | 35+ | 4.0 | RedFox | Vegetarian | 1000.0 | -1 | 21122.0 |
| 7 | 8.0 | 20-25 | NaN | LemonTree | Veg | 2999.0 | -10 | 345673.0 |
| 8 | 9.0 | 25-30 | 2.0 | Ibis | Non-Veg | 3456.0 | 3 | NaN |
| 9 | 10.0 | 30-35 | 5.0 | RedFox | non-Veg | NaN | 4 | 87777.0 |

**df['NoOfPax'].loc[(df['NoOfPax']<1) | (df['NoOfPax']>20)]=np.nan**

**df**

```
C:\Users\Ayyadurai\AppData\Local\Temp\ipykernel_5300\2129877948.py:1: S
ettingWithCopyWarning:
A value is trying to be set on a copy of a slice from a DataFrame

See the caveats in the documentation: https://pandas.pydata.org/pandas-
docs/stable/user_guide/indexing.html#returning-a-view-versus-a-copy
  df['NoOfPax'].loc[(df['NoOfPax']<1) | (df['NoOfPax']>20)]=np.nan
```

| | CustomerID | Age_Group | Rating(1-5) | Hotel | FoodPreference | Bill | NoOfPax | EstimatedSalary |
|---|---|---|---|---|---|---|---|---|
| 0 | 1.0 | 20-25 | 4.0 | Ibis | veg | 1300.0 | 2.0 | 40000.0 |
| 1 | 2.0 | 30-35 | 5.0 | LemonTree | Non-Veg | 2000.0 | 3.0 | 59000.0 |
| 2 | 3.0 | 25-30 | NaN | RedFox | Veg | 1322.0 | 2.0 | 30000.0 |
| 3 | 4.0 | 20-25 | NaN | LemonTree | Veg | 1234.0 | 2.0 | 120000.0 |
| 4 | 5.0 | 35+ | 3.0 | Ibis | Vegetarian | 989.0 | 2.0 | 45000.0 |
| 5 | 6.0 | 35+ | 3.0 | Ibys | Non-Veg | 1909.0 | 2.0 | 122220.0 |
| 6 | 7.0 | 35+ | 4.0 | RedFox | Vegetarian | 1000.0 | NaN | 21122.0 |
| 7 | 8.0 | 20-25 | NaN | LemonTree | Veg | 2999.0 | NaN | 345673.0 |
| 8 | 9.0 | 25-30 | 2.0 | Ibis | Non-Veg | 3456.0 | 3.0 | NaN |
| 9 | 10.0 | 30-35 | 5.0 | RedFox | non-Veg | NaN | 4.0 | 87777.0 |

**df.Age_Group.unique()**

```
array(['20-25', '30-35', '25-30', '35+'], dtype=object)
```

**df.Hotel.unique()**

```
array(['Ibis', 'LemonTree', 'RedFox', 'Ibys'], dtype=object)
```

# Using the inplace=True keyword in a pandas method changes the default behaviour such that the operation on the dataframe doesn't return anything, it instead 'modifies the underlying data

**df.Hotel.replace(['Ibys'],'Ibis',inplace=True)**

**df.FoodPreference.unique**

```
<bound method Series.unique of 0          veg
1       Non-Veg
2           Veg
3           Veg
4     Vegetarian
5       Non-Veg
6     Vegetarian
7           Veg
8       Non-Veg
9       non-Veg
Name: FoodPreference, dtype: object>
```
**df.FoodPreference.replace(['Vegetarian','veg'],'Veg',inplace=True)**

**df.FoodPreference.replace(['non-Veg'],'Non-Veg',inplace=True)**

# Fillna is a Pandas function to fill the NA/NaN values with the specified method.

# If column or feature is numerical continuous data then replace the missing(NaN) value by taking mean value.

# If column or feature is numerical discrete data then replace  the missing(NaN) value by taking median value.

# If column or feature is non-numerical i.e Categorical  data then  replace the missing(NaN) value by taking mode value.

**df.EstimatedSalary.fillna(round(df.EstimatedSalary.mean()),inplace=True)**

**df.NoOfPax.fillna(round(df.NoOfPax.median()),inplace=True)**

**df['Rating(1-5)'].fillna(round(df['Rating(1-5)'].median()), inplace=True)**

**df.Bill.fillna(round(df.Bill.mean()),inplace=True)**

**df**

| | CustomerID | Age_Group | Rating(1-5) | Hotel | FoodPreference | Bill | NoOfPax | Estimated Salary |
|---|---|---|---|---|---|---|---|---|
| 0 | 1.0 | 20-25 | 4.0 | Ibis | Veg | 1300.0 | 2.0 | 40000.0 |
| 1 | 2.0 | 30-35 | 5.0 | LemonTree | Non-Veg | 2000.0 | 3.0 | 59000.0 |
| 2 | 3.0 | 25-30 | 4.0 | RedFox | Veg | 1322.0 | 2.0 | 30000.0 |
| 3 | 4.0 | 20-25 | 4.0 | LemonTree | Veg | 1234.0 | 2.0 | 120000.0 |
| 4 | 5.0 | 35+ | 3.0 | Ibis | Veg | 989.0 | 2.0 | 45000.0 |
| 5 | 6.0 | 35+ | 3.0 | Ibis | Non-Veg | 1909.0 | 2.0 | 122220.0 |
| 6 | 7.0 | 35+ | 4.0 | RedFox | Veg | 1000.0 | 2.0 | 21122.0 |
| 7 | 8.0 | 20-25 | 4.0 | LemonTree | Veg | 2999.0 | 2.0 | 345673.0 |
| 8 | 9.0 | 25-30 | 2.0 | Ibis | Non-Veg | 3456.0 | 3.0 | 96755.0 |
| 9 | 10.0 | 30-35 | 5.0 | RedFox | Non-Veg | 1801.0 | 4.0 | 87777.0 |