

## **EXPERIMENT NO: 2**

**AIM:** Apply data cleaning techniques (e.g. Data Imputation).

### **THEORY:**

#### 1. Dealing with missing data:

Identifying missing data: Use methods like `isnull()` or `info()` to identify missing values in your dataset.

Data Imputation: Fill missing values using techniques like mean, median, mode, or advanced imputation methods such as KNN imputation or interpolation.

#### 2. Check for Duplicates:

Use `duplicated()` method to identify duplicate rows in your dataset.

Remove duplicates using `drop_duplicates()` method.

#### 3. Detect Outliers:

Outliers can be detected using statistical methods like z-score or IQR (Interquartile Range).

Visualizations such as box plots or scatter plots can also help in identifying outliers.

#### 4. Normalize Casing:

Normalize string data to a consistent format, such as converting all text to lowercase or uppercase, using `str.lower()` or `str.upper()`.

#### 5. Check for Trailing Whitespaces:

Use `strip()` method to remove leading and trailing whitespaces from string columns.

#### 6. Extracting Additional Variables:

Extract additional variables from existing columns using string manipulation functions or by splitting columns.

For example, extracting day, month, or year from a date column.

#### 7. Joining Cleaned Datasets:

Use `merge()` function in pandas to join cleaned datasets based on a common key.

Ensure that both datasets have a common key to merge on.

**CODE:**

```
import pandas as pd

# Read the hotel dataset
df_hotels = pd.read_csv('hotels.csv')

# Rename the 'Name' column to 'Hotel' for consistency
df_hotels.rename(columns={'Name': 'Hotel'}, inplace=True)

# Display the first few rows of the hotel dataset
print("Hotel Dataset:")
print(df_hotels)

# Check for duplicates in the hotel dataset
duplicates_hotels = df_hotels[df_hotels.duplicated()]
print("\nDuplicate rows in hotel dataset:")
print(duplicates_hotels)

# Remove duplicates from the hotel dataset
df_hotels = df_hotels.drop_duplicates()

# Read the reviews dataset
df_reviews = pd.read_csv('reviews.csv')

# Display the first few rows of the reviews dataset
print("\nReviews Dataset:")
```

```
print(df_reviews)
```

```
# Check for duplicates in the reviews dataset
```

```
duplicates_reviews = df_reviews[df_reviews.duplicated()]
```

```
print("\nDuplicate rows in reviews dataset:")
```

```
print(duplicates_reviews)
```

```
# Remove duplicates from the reviews dataset
```

```
df_reviews = df_reviews.drop_duplicates()
```

```
# Dealing with missing data
```

```
print("\nMissing data in hotel dataset:")
```

```
print(df_hotels.isnull().sum())
```

```
print("\nMissing data in reviews dataset:")
```

```
print(df_reviews.isnull().sum())
```

```
# Detect outliers
```

```
outliers_hotels = df_hotels[df_hotels['Price'] > 1000]
```

```
print("\nOutliers in hotel dataset:")
```

```
print(outliers_hotels)
```

```
# Merge cleaned datasets on a common key, which is now 'Hotel' in both
```

```
df_hotels and df_reviews
```

```
merged_df = pd.merge(df_hotels, df_reviews, on='Hotel', how='inner')
```

```
# Display the merged dataset
```

```
print("\nMerged Dataset:")
```

```
print(merged_df)
```

## OUTPUT:

```
Hotel Dataset:
  Hotel  Place  Price  Rating
0  Hotel A  Amsterdam  100.0    4.5
1  Hotel B  Amsterdam  150.0    4.0
2  Hotel C    Paris  120.0    4.2
3  Hotel D    Paris  200.0    4.7
4  Hotel E  London  180.0    4.3
5  Hotel F  London  140.0    4.6
6  Hotel A  Amsterdam   NaN    4.5
7  Hotel G    Wilan   NaN    4.9
8  Hotel H    Berlin  220.0    5.0
9  Hotel I    Rome  180.0    3.8
10 Hotel J  Amsterdam  110.0    4.1
11 Hotel K    Munich  250.0    5.5
12 Hotel L    Paris  150.0    4.4
13 Hotel M  Barcelone  190.0    4.2
14 Hotel N  Barcelone  200.0    4.6
15 Hotel O    Berlin  160.0    4.0
16 Hotel P    Paris  300.0    6.2
17 Hotel Q    London  170.0    4.7
18 Hotel R  Amsterdam  210.0    3.5
19 Hotel S    Paris  230.0    4.8

Duplicate rows in hotel dataset:
Empty DataFrame
Columns: [Hotel, Place, Price, Rating]
Index: []

Reviews Dataset:
  Hotel  Rating  Review
0  Hotel A    4.5  Great experience!
1  Hotel B    4.0  Decent stay
2  Hotel C    4.2  Nice hotel
3  Hotel D    4.7  Amazing service
4  Hotel E    4.3  Fantastic location
5  Hotel F    4.6  Excellent stay
6  Hotel G    4.9  Highly recommend
7  Hotel H    5.0  Outstanding service and amenities
8  Hotel I    3.8  Disappointing experience
9  Hotel J    4.1  Good value for money
10 Hotel K    5.0  Unbelievable experience
11 Hotel L    4.4  Very comfortable stay
12 Hotel M    4.2  Enjoyable stay overall
13 Hotel N    4.6  Would definitely visit again
14 Hotel O    4.0  Average stay
15 Hotel P    6.2  Exceptional service beyond expectation
16 Hotel Q    4.7  Superb location
17 Hotel R    3.5  Not as expected
18 Hotel S    4.8  Highly satisfied with the stay

Duplicate rows in reviews dataset:
Empty DataFrame
Columns: [Hotel, Rating, Review]
Index: []

Missing data in hotel dataset:
Hotel      0
Place      0
Price      2
Rating     0
dtype: int64
```

```

Missing data in reviews dataset:
Hotel      0
Rating     0
Review     0
dtype: int64

Outliers in hotel dataset:
Empty DataFrame
Columns: [Hotel, Place, Price, Rating]
Index: []

Merged Dataset:
   Hotel  Place  Price  Rating_x  Rating_y  \
0  Hotel A  Amsterdam  100.0     4.5     4.5
1  Hotel A  Amsterdam   NaN     4.5     4.5
2  Hotel B  Amsterdam  150.0     4.0     4.0
3  Hotel C    Paris  120.0     4.2     4.2
4  Hotel D    Paris  200.0     4.7     4.7
5  Hotel E   London  180.0     4.3     4.3
6  Hotel F   London  140.0     4.6     4.6
7  Hotel G    Wilan   NaN     4.9     4.9
8  Hotel H   Berlin  220.0     5.0     5.0
9  Hotel I    Rome  180.0     3.8     3.8
10 Hotel J  Amsterdam  110.0     4.1     4.1
11 Hotel K    Munich  250.0     5.5     5.0
12 Hotel L    Paris  150.0     4.4     4.4
13 Hotel M  Barcelone  190.0     4.2     4.2
14 Hotel N  Barcelone  200.0     4.6     4.6
15 Hotel O    Berlin  160.0     4.0     4.0
16 Hotel P    Paris  300.0     6.2     6.2
17 Hotel Q   London  170.0     4.7     4.7
18 Hotel R  Amsterdam  210.0     3.5     3.5
19 Hotel S    Paris  230.0     4.8     4.8

   Review
0  Great experience!
1  Great experience!
2  Decent stay
3  Nice hotel
4  Amazing service
5  Fantastic location
6  Excellent stay
7  Highly recommend
8  Outstanding service and amenities
9  Disappointing experience
10 Good value for money
11 Unbelievable experience
12 Very comfortable stay
13 Enjoyable stay overall
14 Would definitely visit again
15 Average stay
16 Exceptional service beyond expectation
17 Superb location
18 Not as expected
19 Highly satisfied with the stay

```

## CONCLUSION:

In this experiment, we cleaned two datasets: one containing hotel information and another with customer reviews. We removed duplicates, handled missing data, and checked for outliers. Both datasets were mostly clean, with only a few missing values in the hotel dataset. After cleaning, we merged the datasets using the hotel names.