# Exploratory Data Analysis

Team ID : PNT2022TMID13523
Date : 07/11/2022
Project Name : Analytics for Hospital's Health Care Data

Required libraries:

```python
import pandas as pd
import numpy as np
import matplotlib.pyplot as plt
import seaborn as sns
%matplotlib inline
```

```python
df = pd.read_csv("~/Downloads/Desktop/Hospital Health care data/case_data.csv")
```

```python
df
```

| | case_id | Hospital_code | Hospital_type_code | City_Code_Hospital | Hospital_region_code | Available Extra Rooms in Hospital | Department | Ward_Type | Ward_Facility_Code | Bed Grade | patientid | City_Code_Patient | Type of Admission | Severity of Illness | Visitors with Patient | Age | A |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| 0 | 1 | 8 | c | 3 | Z | 3 | radiotherapy | R | F | 2.0 | 31397 | 7.0 | Emergency | Extreme | 2 | 51-60 | |
| 1 | 2 | 2 | c | 5 | Z | 2 | radiotherapy | S | F | 2.0 | 31397 | 7.0 | Trauma | Extreme | 2 | 51-60 | |
| 2 | 3 | 10 | e | 1 | X | 2 | anesthesia | S | E | 2.0 | 31397 | 7.0 | Trauma | Extreme | 2 | 51-60 | |
| 3 | 4 | 26 | b | 2 | Y | 2 | radiotherapy | R | D | 2.0 | 31397 | 7.0 | Trauma | Extreme | 2 | 51-60 | |
| 4 | 5 | 26 | b | 2 | Y | 2 | radiotherapy | S | D | 2.0 | 31397 | 7.0 | Trauma | Extreme | 2 | 51-60 | |
| ... | ... | ... | ... | ... | ... | ... | ... | ... | ... | ... | ... | ... | ... | ... | ... | ... | |
| 318433 | 318434 | 6 | a | 6 | X | 3 | radiotherapy | Q | F | 4.0 | 86499 | 23.0 | Emergency | Moderate | 3 | 41-50 | |
| 318434 | 318435 | 24 | a | 1 | X | 3 | anesthesia | Q | E | 4.0 | 325 | 8.0 | Urgent | Moderate | 3 | 41-50 | |
| 318435 | 318436 | 7 | a | 4 | X | 3 | gynecology | R | F | 4.0 | 125235 | 10.0 | Emergency | Minor | 1 | 71-80 | |
| 318436 | 318437 | 11 | b | 2 | Y | 3 | anesthesia | Q | D | 3.0 | 91081 | 8.0 | Trauma | Minor | 1 | 11-20 | |
| 318437 | 318438 | 19 | a | 7 | Y | 5 | gynecology | Q | C | 2.0 | 21641 | 8.0 | Emergency | Minor | 2 | 11-20 | |

318438 rows × 18 columns

```python
df.head()
```

| | case_id | Hospital_code | Hospital_type_code | City_Code_Hospital | Hospital_region_code | Available Extra Rooms in Hospital | Department | Ward_Type | Ward_Facility_Code | Bed Grade | patientid | City_Code_Patient | Type of Admission | Severity of Illness | Visitors with Patient | Age | A |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| 0 | 1 | 8 | c | 3 | Z | 3 | radiotherapy | R | F | 2.0 | 31397 | 7.0 | Emergency | Extreme | 2 | 51-60 | |
| 1 | 2 | 2 | c | 5 | Z | 2 | radiotherapy | S | F | 2.0 | 31397 | 7.0 | Trauma | Extreme | 2 | 51-60 | |
| 2 | 3 | 10 | e | 1 | X | 2 | anesthesia | S | E | 2.0 | 31397 | 7.0 | Trauma | Extreme | 2 | 51-60 | |
| 3 | 4 | 26 | b | 2 | Y | 2 | radiotherapy | R | D | 2.0 | 31397 | 7.0 | Trauma | Extreme | 2 | 51-60 | |
| 4 | 5 | 26 | b | 2 | Y | 2 | radiotherapy | S | D | 2.0 | 31397 | 7.0 | Trauma | Extreme | 2 | 51-60 | |

```python
df.tail()
```

| | case_id | Hospital_code | Hospital_type_code | City_Code_Hospital | Hospital_region_code | Available Extra Rooms in Hospital | Department | Ward_Type | Ward_Facility_Code | Bed Grade | patientid | City_Code_Patient | Type of Admission | Severity of Illness | Visitors with Patient | Age | A |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| 318433 | 318434 | 6 | a | 6 | X | 3 | radiotherapy | Q | F | 4.0 | 86499 | 23.0 | Emergency | Moderate | 3 | 41-50 | |
| 318434 | 318435 | 24 | a | 1 | X | 3 | anesthesia | Q | E | 4.0 | 325 | 8.0 | Urgent | Moderate | 3 | 41-50 | |
| 318435 | 318436 | 7 | a | 4 | X | 3 | gynecology | R | F | 4.0 | 125235 | 10.0 | Emergency | Minor | 1 | 71-80 | |
| 318436 | 318437 | 11 | b | 2 | Y | 3 | anesthesia | Q | D | 3.0 | 91081 | 8.0 | Trauma | Minor | 1 | 11-20 | |
| 318437 | 318438 | 19 | a | 7 | Y | 5 | gynecology | Q | C | 2.0 | 21641 | 8.0 | Emergency | Minor | 2 | 11-20 | |

```python
df.info()
```

```
<class 'pandas.core.frame.DataFrame'>
RangeIndex: 318438 entries, 0 to 318437
Data columns (total 18 columns):
 #   Column                              Non-Null Count   Dtype
---  ------                              --------------   -----
 0   case_id                             318438 non-null  int64
 1   Hospital_code                       318438 non-null  int64
 2   Hospital_type_code                  318438 non-null  object
 3   City_Code_Hospital                  318438 non-null  int64
 4   Hospital_region_code                318438 non-null  object
 5   Available Extra Rooms in Hospital   318438 non-null  int64
 6   Department                          318438 non-null  object
 7   Ward_Type                           318438 non-null  object
 8   Ward_Facility_Code                  318438 non-null  object
 9   Bed Grade                           318325 non-null  float64
 10  patientid                           318438 non-null  int64
 11  City_Code_Patient                   313906 non-null  float64
 12  Type of Admission                   318438 non-null  object
 13  Severity of Illness                 318438 non-null  object
 14  Visitors with Patient               318438 non-null  int64
 15  Age                                 318438 non-null  object
 16  Admission_Deposit                   318438 non-null  float64
 17  Stay                                318438 non-null  object
dtypes: float64(3), int64(6), object(9)
memory usage: 43.7+ MB
```

```python
df.dtypes
```

```
case_id                               int64
Hospital_code                         int64
Hospital_type_code                    object
City_Code_Hospital                    int64
Hospital_region_code                  object
Available Extra Rooms in Hospital     int64
Department                            object
Ward_Type                             object
Ward_Facility_Code                    object
Bed Grade                             float64
patientid                             int64
City_Code_Patient                     float64
Type of Admission                     object
Severity of Illness                   object
Visitors with Patient                 int64
Age                                   object
Admission_Deposit                     float64
Stay                                  object
dtype: object
```

```python
df.shape
```

```
(318438, 18)
```

## Before Null Values checking :

```python
df.isnull().sum()
```

```
```

```python
df.isnull()
```

| | case_id | Hospital_code | Hospital_type_code | City_Code_Hospital | Hospital_region_code | Available Extra Rooms in Hospital | Department | Ward_Type | Ward_Facility_Code | Bed Grade | patientid | City_Code_Patient | Type of Admission | Severity of Illness | Visitors with Patient | Age | A |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| 0 | False | False | False | False | False | False | False | False | False | False | False | False | False | False | False | False | |
| 1 | False | False | False | False | False | False | False | False | False | False | False | False | False | False | False | False | |
| 2 | False | False | False | False | False | False | False | False | False | False | False | False | False | False | False | False | |
| 3 | False | False | False | False | False | False | False | False | False | False | False | False | False | False | False | False | |
| 4 | False | False | False | False | False | False | False | False | False | False | False | False | False | False | False | False | |
| ... | ... | ... | ... | ... | ... | ... | ... | ... | ... | ... | ... | ... | ... | ... | ... | ... | |
| 318433 | False | False | False | False | False | False | False | False | False | False | False | False | False | False | False | False | |
| 318434 | False | False | False | False | False | False | False | False | False | False | False | False | False | False | False | False | |
| 318435 | False | False | False | False | False | False | False | False | False | False | False | False | False | False | False | False | |
| 318436 | False | False | False | False | False | False | False | False | False | False | False | False | False | False | False | False | |
| 318437 | False | False | False | False | False | False | False | False | False | False | False | False | False | False | False | False | |

318438 rows × 18 columns

```python
df.describe()
```

| | case_id | Hospital_code | City_Code_Hospital | Available Extra Rooms in Hospital | Bed Grade | patientid | City_Code_Patient | Visitors with Patient | Admission_Deposit |
|---|---|---|---|---|---|---|---|---|---|
| count | 318438.000000 | 318438.000000 | 318438.000000 | 318438.000000 | 318325.000000 | 318438.000000 | 313906.000000 | 318438.000000 | 318438.000000 |
| mean | 159219.500000 | 18.318841 | 4.777177 | 3.197627 | 2.625807 | 65747.079472 | 7.251906 | 3.284098 | 4880.746393 |
| std | 91924.076847 | 8.633738 | 3.102638 | 1.168171 | 0.873146 | 37975.936440 | 4.745266 | 1.764061 | 1086.776254 |
| min | 1.000000 | 1.000000 | 1.000000 | 0.000000 | 1.000000 | 1.000000 | 1.000000 | 0.000000 | 1800.000000 |
| 25% | 79610.250000 | 11.000000 | 2.000000 | 2.000000 | 2.000000 | 32667.000000 | 4.000000 | 2.000000 | 4186.000000 |
| 50% | 159219.500000 | 19.000000 | 5.000000 | 3.000000 | 3.000000 | 65724.500000 | 8.000000 | 3.000000 | 4741.000000 |
| 75% | 238828.750000 | 26.000000 | 7.000000 | 4.000000 | 3.000000 | 98472.000000 | 8.000000 | 4.000000 | 5409.000000 |
| max | 318438.000000 | 32.000000 | 13.000000 | 24.000000 | 4.000000 | 131624.000000 | 38.000000 | 32.000000 | 11008.000000 |

```python
df.isnull().sum()
```

```
case_id                               0
Hospital_code                         0
Hospital_type_code                    0
City_Code_Hospital                    0
Hospital_region_code                  0
Available Extra Rooms in Hospital     0
Department                            0
Ward_Type                             0
Ward_Facility_Code                    0
Bed Grade                             113
patientid                             0
City_Code_Patient                     4532
Type of Admission                     0
Severity of Illness                   0
Visitors with Patient                 0
Age                                   0
Admission_Deposit                     0
Stay                                  0
dtype: int64
```

```python
df.corr()
```

| | case_id | Hospital_code | City_Code_Hospital | Available Extra Rooms in Hospital | Bed Grade | patientid | City_Code_Patient | Visitors with Patient | Admission_Deposit |
|---|---|---|---|---|---|---|---|---|---|
| case_id | 1.000000 | -0.043023 | -0.011362 | 0.042990 | 0.013152 | -0.004150 | 0.000196 | 0.001336 | -0.049373 |
| Hospital_code | -0.043023 | 1.000000 | 0.128204 | -0.046711 | -0.043658 | 0.001839 | -0.015320 | -0.003688 | -0.045446 |
| City_Code_Hospital | -0.011362 | 0.128204 | 1.000000 | -0.045711 | -0.040006 | 0.002730 | -0.023368 | 0.018194 | -0.034455 |
| Available Extra Rooms in Hospital | 0.042990 | -0.046711 | -0.045711 | 1.000000 | 0.095681 | 0.001725 | 0.005468 | 0.041846 | 0.149278 |
| Bed Grade | 0.013152 | -0.043658 | -0.040006 | 0.095681 | 1.000000 | -0.003920 | 0.001051 | -0.008105 | 0.089848 |
| patientid | -0.004150 | 0.001839 | 0.002730 | 0.001725 | -0.003920 | 1.000000 | 0.003302 | -0.006849 | -0.000877 |
| City_Code_Patient | 0.000196 | -0.015320 | -0.023368 | 0.005468 | 0.001051 | 0.003302 | 1.000000 | -0.010253 | 0.025687 |
| Visitors with Patient | 0.001336 | -0.003688 | 0.018194 | 0.041846 | -0.008105 | -0.006849 | -0.010253 | 1.000000 | 0.150358 |
| Admission_Deposit | -0.049373 | -0.045446 | -0.034455 | 0.149278 | 0.075633 | -0.000877 | 0.025687 | 0.150358 | 1.000000 |

## Work With Null Values :

```python
df["Bed Grade"].fillna(df["Bed Grade"].mean(),inplace=True)
```

```python
df.isnull().sum()
```

```python
df.describe()
```

```
case_id                               0
Hospital_code                         0
Hospital_type_code                    0
City_Code_Hospital                    0
Hospital_region_code                  0
Available Extra Rooms in Hospital     0
Department                            0
Ward_Type                             0
Ward_Facility_Code                    0
Bed Grade                             0
patientid                             0
City_Code_Patient                     4532
Type of Admission                     0
Severity of Illness                   0
Visitors with Patient                 0
Age                                   0
Admission_Deposit                     0
Stay                                  0
dtype: int64
```

```python
df["City_Code_Patient"].fillna(df["City_Code_Patient"].mean(),inplace=True)
```

```python
df.isnull().sum()
```

## After Cleaning Process :

## Total Null Values Checking :

```python
df.isnull().sum()
```

```
case_id                               0
Hospital_code                         0
Hospital_type_code                    0
City_Code_Hospital                    0
Hospital_region_code                  0
Available Extra Rooms in Hospital     0
Department                            0
Ward_Type                             0
Ward_Facility_Code                    0
Bed Grade                             0
patientid                             0
City_Code_Patient                     0
Type of Admission                     0
Severity of Illness                   0
Visitors with Patient                 0
Age                                   0
Admission_Deposit                     0
Stay                                  0
dtype: int64
```

## Total Null Values :

```python
df.isnull().sum().sum()
```

```
0
```

```python
df.corr()
```

| | case_id | Hospital_code | City_Code_Hospital | Available Extra Rooms in Hospital | Bed Grade | patientid | City_Code_Patient | Visitors with Patient | Admission_Deposit |
|---|---|---|---|---|---|---|---|---|---|
| case_id | 1.000000 | 0.430257e+00 | -341.61.259836 | -0241.01187 | 4572.486177 | 0.068462e+00 | 0.468686e+07 | 28036.806476 | -241.258814 | -4.592706e+08 |
| Hospital_code | 0.430257e+00 | 74.541723 | 3.436341 | 0.436941 | -0.601465 | -0.120559 | 7.551946 | -0.420561 | -0.434073 | -2.994139e+02 |
| City_Code_Hospital | -3.23701e+03 | 3.436541 | 9.635738 | -0.165887 | -0.120848 | 8.941996e+01 | -6.345149 | 0.068635 | -1.181705e+02 |
| Available Extra Rooms in Hospital | 4.57248e+03 | -0.601465 | -0.165887 | 1.364624 | -0.116146 | 4.689539e+01 | 0.408208 | 0.182350 | 4.732048e+02 |
| Bed Grade | 1.084464e+03 | -0.120559 | -0.120848 | -0.116146 | 1.000000 | -0.030071 | 0.042021 | -0.033073 | 0.196902 | 7.004526e+01 |
| patientid | -1.468868e+07 | 751.114984 | 86.418979 | 40.809395 | 64.528836 | 1.482479e+09 | 355.708801 | 461.570860 | -5.025211e+04 |
| City_Code_Patient | 2.80359e+04 | -0.627256 | -0.345199 | 0.200208 | 0.053580 | 0.632075 | 1.217595e+02 | 0.069346 | 1.312736e+02 |
| Visitors with Patient | 2.122506e+04 | -0.434073 | 0.095640 | 0.158252 | 0.156362 | 4.575799e+02 | -0.069346 | 3.111913 | 2.882567e+02 |
| Admission_Deposit | -4.592730e+06 | -401.91924 | -116.17091 | -182.488676 | 70.046318 | -5.025837e+04 | 131.273635 | 288.256678 | 1.181583e+06 |

```python
sns.heatmap(df.corr(),annot=True)
plt.title("Correlation Matrix")
plt.show()
```
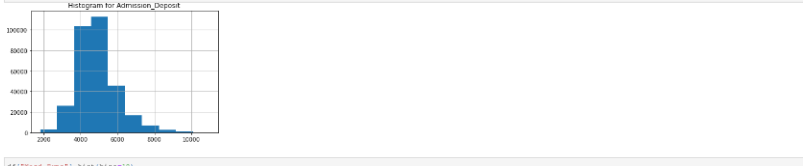


```python
df["Admission_Deposit"].hist(bins=20)
plt.title("Histogram for Admission_Deposit")
plt.show()
```



```python
df["Ward_Type"].hist(bins=20)
plt.title("Histogram for Ward_Type")
plt.show()
```



```python
df["patientid"].hist(bins=20)
plt.title("Histogram for patientid")
plt.show()
```