

## **SPRINT 2**

Date	07-11-2022
Team ID	PNT2022TMID04531
Project Title	Analytics for Hospitals' Health-Care Data
Team Members	Haris S, Gokula Kannan G, Abishek A, hari vignesh K G

### **Data Cleaning and Preparation**

In this data set, variables “City\_code\_patient” and “Bed Grade” have missing values. These missing values must be treated before feeding to the algorithm as they distort the model performance.

So, the missing values are replaced using the “mode” of the column.  
Since most of the variables in the dataset have ordinal data, we transformed them into levels by using a label encoder to perform further analysis on the data.

#### Distinct Observations of Ordinal Data

Variables	Number of distinct observations
Hospital_type_code	7
Hospital_region_code	3
Department	5
Ward_Type	6
Ward_Facility_Code	6
Type of Admission	3
Severity of Illness	3
Age	10
Stay	11

SPRINT2.ipynb

train.info()

```
<class 'pandas.core.frame.DataFrame'>
Int64Index: 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  int64
 3   City_Code_Hospital                      318438 non-null  int64
 4   Hospital_region_code                   318438 non-null  int64
 5   Available Extra Rooms in Hospital      318438 non-null  int64
 6   Department                             318438 non-null  int64
 7   Ward_Type                               318438 non-null  int64
 8   Ward_Facility_Code                     318438 non-null  int64
 9   Bed_Grade                               318438 non-null  float64
10   patientid                               318438 non-null  int64
11   City_Code_Patient                      318438 non-null  float64
12   Type of Admission                       318438 non-null  int64
13   Severity of illness                     318438 non-null  int64
14   Visitors with Patient                   318438 non-null  int64
15   Age                                     318438 non-null  int64
16   Admission_Deposit                       318438 non-null  float64
17   Stay                                   318438 non-null  int64
dtypes: float64(3), int64(15)
memory usage: 46.2 MB
```

test.info()

```
<class 'pandas.core.frame.DataFrame'>
Int64Index: 137857 entries, 0 to 137856
Data columns (total 18 columns):
 #   Column                                     Non-Null Count  Dtype
---  -
 0   case_id                                   137857 non-null  int64
 1   Hospital_code                             137857 non-null  int64
 2   Hospital_type_code                       137857 non-null  int64
 3   City_Code_Hospital                      137857 non-null  int64
 4   Hospital_region_code                   137857 non-null  int64
 5   Available Extra Rooms in Hospital      137857 non-null  int64
 6   Department                             137857 non-null  int64
 7   Ward_Type                               137857 non-null  int64
 8   Ward_Facility_Code                     137857 non-null  int64
 9   Bed_Grade                               137857 non-null  float64
```

test.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
0	318439	21	2	3	2	3	2	3	0	2.0	17006	2.0	0	
1	318440	29	0	4	0	2	2	3	5	2.0	17006	2.0	1	
2	318441	26	1	2	1	3	2	1	3	4.0	17006	2.0	0	
3	318442	6	0	6	0	3	2	1	5	2.0	17006	2.0	1	
4	318443	28	1	11	0	2	2	2	5	2.0	17006	2.0	1	

train.shape

```
(318438, 18)
```

test.shape

```
(137857, 18)
```

train.info()

```
<class 'pandas.core.frame.DataFrame'>
Int64Index: 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  int64
 3   City_Code_Hospital                      318438 non-null  int64
 4   Hospital_region_code                   318438 non-null  int64
 5   Available Extra Rooms in Hospital      318438 non-null  int64
 6   Department                             318438 non-null  int64
 7   Ward_Type                               318438 non-null  int64
```

SPRINT2.ipynb

File Edit View Insert Runtime Tools Help All changes saved

Files

- sample\_data
- sample\_sub.csv
- test\_data.csv
- train\_data.csv
- train\_data\_dictionary.csv

```
[36] for i in ['Hospital_type_code', 'Hospital_region_code', 'Department',
            'Ward_Type', 'Ward_Facility_Code', 'Type of Admission', 'Severity of Illness', 'Age']:
    le = LabelEncoder()
    df[i] = le.fit_transform(df[i].astype(str))

[37] #Separating train and test Datasets
train = df[df['Stay']!=1]
test = df[df['Stay']==1]

[38] train.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	Sever Illn
0	1	8	2	3	2	3	3	2	5	2.0	31397	7.0	0	
1	2	2	2	5	2	2	3	3	5	2.0	31397	7.0	1	
2	3	10	4	1	0	2	1	3	4	2.0	31397	7.0	1	
3	4	26	1	2	1	2	3	2	3	2.0	31397	7.0	1	
4	5	26	1	2	1	2	3	3	3	2.0	31397	7.0	1	

```
[39] test.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	Sever Illn
0	318439	21	2	3	2	3	2	3	0	2.0	17006	2.0	0	
1	318440	29	0	4	0	2	2	3	5	2.0	17006	2.0	1	

0s completed at 11:39 PM

SPRINT2.ipynb

File Edit View Insert Runtime Tools Help All changes saved

Files

- sample\_data
- sample\_sub.csv
- test\_data.csv
- train\_data.csv
- train\_data\_dictionary.csv

```
[32] for i in test.columns:
    print(i, " - ", test[i].nunique())

case_id - 137857
Hospital_code - 32
Hospital_type_code - 7
City_Code_Hospital - 11
Hospital_region_code - 3
Available Extra Rooms in Hospital - 15
Department - 5
Ward_Type - 6
Ward_Facility_Code - 6
Bed Grade - 4
patientid - 39687
City_Code_Patient - 37
Type of Admission - 3
Severity of Illness - 3
Visitors with Patient - 27
Age - 10
Admission_Deposit - 6689

[33] train['Bed Grade'].fillna(train['Bed Grade'].mode()[0], inplace = True)
test['Bed Grade'].fillna(test['Bed Grade'].mode()[0], inplace = True)
train['City_Code_Patient'].fillna(train['City_Code_Patient'].mode()[0], inplace = True)
test['City_Code_Patient'].fillna(test['City_Code_Patient'].mode()[0], inplace = True)

[34] from sklearn.preprocessing import LabelEncoder
le = LabelEncoder()
train['Stay'] = le.fit_transform(train['Stay'].astype('str'))

[35] test['Stay'] = -1
df = pd.concat([train, test])
df.shape

(455495, 18)

[36] for i in ['Hospital_type_code', 'Hospital_region_code', 'Department',
            'Ward_Type', 'Ward_Facility_Code', 'Type of Admission', 'Severity of Illness', 'Age']:
    le = LabelEncoder()
    df[i] = le.fit_transform(df[i].astype(str))
```

0s completed at 11:39 PM

SPRINT2.ipynb

File Edit View Insert Runtime Tools Help All changes saved

Files

- sample\_data
- sample\_sub.csv
- test\_data.csv
- train\_data.csv
- train\_data\_dictionary.csv

```
[30] train.shape
test.shape

(137057, 17)

[31] for i in train.columns:
    print(i, '-', train[i].nunique())

case_id - 310438
Hospital_code - 32
Hospital_type_code - 7
City_Code_Hospital - 11
Hospital_region_code - 3
Available Extra Rooms in Hospital - 18
Department - 5
Ward_Type - 6
Ward_Facility_Code - 6
Bed_Grade - 4
patientid - 92017
City_Code_Patient - 37
Type of Admission - 3
Severity of Illness - 3
Visitors with Patient - 28
Age - 10
Admission_Deposit - 7300
Stay - 11

[32] for i in test.columns:
    print(i, '-', test[i].nunique())

case_id - 137057
Hospital_code - 32
Hospital_type_code - 7
City_Code_Hospital - 11
Hospital_region_code - 3
Available Extra Rooms in Hospital - 15
Department - 5
Ward_Type - 6
Ward_Facility_Code - 6
Bed_Grade - 4
patientid - 39607
City_Code_Patient - 37
Type of Admission - 3
```

completed at 11:39 PM

SPRINT2.ipynb

File Edit View Insert Runtime Tools Help All changes saved

Files

- sample\_data
- sample\_sub.csv
- test\_data.csv
- train\_data.csv
- train\_data\_dictionary.csv

```
train.info()
train.Stay.unique()

<class 'pandas.core.frame.DataFrame'>
RangeIndex: 310438 entries, 0 to 310437
Data columns (total 18 columns):
 #   Column                                Non-Null Count  Dtype
---  --
 0   case_id                             310438 non-null  int64
 1   Hospital_code                       310438 non-null  int64
 2   Hospital_type_code                  310438 non-null  object
 3   City_Code_Hospital                  310438 non-null  int64
 4   Hospital_region_code                310438 non-null  object
 5   Available Extra Rooms in Hospital    310438 non-null  int64
 6   Department                          310438 non-null  object
 7   Ward_Type                          310438 non-null  object
 8   Ward_Facility_Code                  310438 non-null  object
 9   Bed_Grade                           310438 non-null  float64
10   patientid                           310438 non-null  int64
11   City_Code_Patient                   310438 non-null  float64
12   Type of Admission                    310438 non-null  object
13   Severity of Illness                  310438 non-null  object
14   Visitors with Patient                310438 non-null  int64
15   Age                                 310438 non-null  object
16   Admission_Deposit                   310438 non-null  float64
17   Stay                                310438 non-null  object
dtypes: float64(3), int64(6), object(9)
memory usage: 43.7+ MB
array(['0-10', '11-20', '21-30', '31-40', '41-50', '51-60', '61-70', '71-80',
       'More than 80 Days', '81-90', '91-100'], dtype=object)

[28] train.isnull().sum().sort_values(ascending = False)

City_Code_Patient    4532
Bed_Grade            112
Hospital_code         0
Admission_Deposit    0
Age                  0
Visitors with Patient 0
Severity of Illness  0
Type of Admission    0
patientid            0
case_id              0
Ward_Facility_Code   0
Ward_Type            0
```

completed at 11:39 PM

SPRINT2.ipynb

File Edit View Insert Runtime Tools Help All changes saved

Files

- sample\_data
- sample\_sub.csv
- test\_data.csv
- train\_data.csv
- train\_data\_dictionary.csv

```
import numpy as np
import pandas as pd
import matplotlib.pyplot as plt
np.set_printoptions(suppress=True)
import warnings
warnings.filterwarnings('ignore')

# Importing datasets
train = pd.read_csv('train_data.csv')
test = pd.read_csv('test_data.csv')

train.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
0	1	8	c	3	Z	3	radiotherapy	R	F	2.0	31997	7.0	Emergency	Extm
1	2	2	c	5	Z	2	radiotherapy	S	F	2.0	31997	7.0	Trauma	Extm
2	3	10	e	1	X	2	anesthesia	S	E	2.0	31997	7.0	Trauma	Extm
3	4	26	b	2	Y	2	radiotherapy	R	D	2.0	31997	7.0	Trauma	Extm
4	5	26	b	2	Y	2	radiotherapy	S	D	2.0	31997	7.0	Trauma	Extm

```
train.info()
train.Stay.unique()

<class 'pandas.core.frame.DataFrame'>
RangeIndex: 318438 entries, 0 to 318437
```

0s completed at 11:39 PM

---

SPRINT2.ipynb

File Edit View Insert Runtime Tools Help All changes saved

Files

- sample\_data
- sample\_sub.csv
- test\_data.csv
- train\_data.csv
- train\_data\_dictionary.csv

```
16 Admission_Deposit 318438 non-null float64
17 Stay 318438 non-null int64
dtypes: float64(3), int64(15)
memory usage: 48.2 MB

test.info()

<class 'pandas.core.frame.DataFrame'>
Int64Index: 137857 entries, 0 to 137856
Data columns (total 18 columns):
# Column Non-Null Count Dtype
---
0 case_id 137857 non-null int64
1 Hospital_code 137857 non-null int64
2 Hospital_type_code 137857 non-null int64
3 City_Code_Hospital 137857 non-null int64
4 Hospital_region_code 137857 non-null int64
5 Available Extra Rooms in Hospital 137857 non-null int64
6 Department 137857 non-null int64
7 Ward_Type 137857 non-null int64
8 Ward_Facility_Code 137857 non-null int64
9 Bed Grade 137857 non-null float64
10 patientid 137857 non-null int64
11 City_Code_Patient 137857 non-null float64
12 Type of Admission 137857 non-null int64
13 Severity of Illness 137857 non-null int64
14 Visitors with Patient 137857 non-null int64
15 Age 137857 non-null int64
16 Admission_Deposit 137857 non-null float64
17 Stay 137857 non-null int64
dtypes: float64(3), int64(15)
memory usage: 15.9 MB
```

0s completed at 11:39 PM