

Lab Report

Course Code: CSE 316

Course Name: Artificial Intelligence Lab

Lab Report No: 03

Date of submission: 15-08-2023

Submitted To:

Md. Sadekur Rahman

Assistant Professor

Department of CSE

Daffodil International University

Submitted By:

Name: Rayhan Rafin

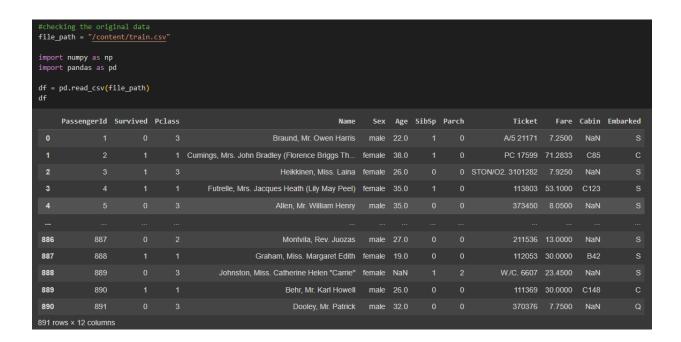
ID: 213-15-4278

Section: 60_B

No	Content	Page No
1.	Getting started	3
2.	Checking Null Values	3
3.	Null Value Handling	4
4.	Label Encoder	6
5.	OneHotEncoder	7

Getting started

Checking original data:



Checking Null Values

Check null:

```
# Here i checked if there is any null value present in the data
df.isnull().sum()
PassengerId
                 0
Survived
                 0
Pclass
                 0
Name
                 0
Sex
                 0
Age
               177
SibSp
                 0
Parch
                 0
Ticket
                 0
Fare
                 0
Cabin
               687
Embarked
                 2
dtype: int64
```

Check null percentage:

```
'''By dividing the total null values
of the data by the total length
we can get the percentage of null'''
(df.isnull().sum())/len(df)*100
PassengerId 0.000000
Survived
              0.000000
Pclass
              0.000000
Name
              0.000000
Sex
              0.000000
Age
             19.865320
SibSp
              0.000000
Parch
              0.000000
Ticket
              0.000000
Fare
              0.000000
Cabin
             77.104377
Embarked
              0.224467
dtype: float64
```

Null Value Handling

Dropna Function:

df_1	= df.copy()		ves all	the null values								
df_1 df_1	= df_1.dropna	1()										
	PassengerId	Survived	Pclass	Name	Sex	Age	SibSp	Parch	Ticket	Fare	Cabin	Embarked
				Cumings, Mrs. John Bradley (Florence Briggs Th	female	38.0			PC 17599	71.2833	C85	С
3				Futrelle, Mrs. Jacques Heath (Lily May Peel)	female	35.0			113803	53.1000	C123	
6				McCarthy, Mr. Timothy J	male	54.0			17463	51.8625	E46	
10	11			Sandstrom, Miss. Marguerite Rut	female	4.0			PP 9549	16.7000	G6	
11	12			Bonnell, Miss. Elizabeth	female	58.0			113783	26.5500	C103	
871	872			Beckwith, Mrs. Richard Leonard (Sallie Monypeny)	female	47.0			11751	52.5542	D35	
872	873			Carlsson, Mr. Frans Olof	male	33.0			695	5.0000	B51 B53 B55	
879	880			Potter, Mrs. Thomas Jr (Lily Alexenia Wilson)	female	56.0			11767	83.1583	C50	С
887	888			Graham, Miss. Margaret Edith	female	19.0			112053	30.0000	B42	
889	890			Behr, Mr. Karl Howell	male	26.0			111369	30.0000	C148	С
183 rd	ows × 12 column	S										

```
#null percentage after dropna function
(df_1.isnull().sum())/len(df_1)*100
PassengerId
               0.0
Survived
               0.0
Pclass
              0.0
Name
              0.0
Sex
              0.0
              0.0
Age
SibSp
              0.0
Parch
               0.0
Ticket
              0.0
Fare
               0.0
Cabin
              0.0
Embarked
               0.0
dtype: float64
```

Fillna Function:

```
#Replacing the null values of age column with 1001
df_2 = df.copy()
df_2['Age'].fillna(1001,inplace = True)
df_2['Age']
         22.0
         38.0
         26.0
         35.0
4
         35.0
886
         27.0
         19.0
887
       1001.0
888
889
         26.0
         32.0
890
Name: Age, Length: 891, dtype: float64
```

```
# Replced the null values with the mean of age column
df_x = df.copy()
df_x['Age'].fillna(df_x["Age"].mean(),inplace = True)
df_x['Age']
0
       22.000000
1
       38.000000
2
       26.000000
       35.000000
       35.000000
       27.000000
886
887
       19.000000
       29.699118
889
       26.000000
       32.000000
Name: Age, Length: 891, dtype: float64
```

LabelEncoder

Encoding District:

```
# we encoded the District column with label encoder to New District
# this will assign numeric values by dictionary order in New District
from sklearn.preprocessing import LabelEncoder
lb = LabelEncoder()
df["New District"] = lb.fit_transform(df["District"])
df
       District Size Population
                                      Speciality New District
 0
          Dhaka 1432
                           2250000 Administrative
         Gazipur
                   879
                            567984
                                          Industry
                                                               9
     Narayanganj
                   576
                              53426
                                          Industry
 3
          Rajbari
                   897
                              65899
                                        Agriculture
                                                              12
                  1234
      Chittagong
                           1345566
                                          Industry
     Cox's Bazer
                   456
                              46567
                                            Tourist
      Bandarban
                   345
                             67579
                                           Tourist
 7
         Khustia
                   432
                              57798
                                         Business
                                                              10
 8
            Feni
                   543
                              67890
                                        Agriculture
 9
                              77898
         Comilla
                   564
                                            Tourist
 10
          Barisal
                   577
                              89750
                                        Agriculture
        Faridpur
                              77650
                   567
                                        Agriculture
 12
        Rangpur
                   575
                              78966
                                        Agriculture
                                                              13
 13
        Chadpur
                   876
                              67789
                                            Tourist
```

Encoding Speciality:

				lumn the same nsform(df["Spe		
	District	Size	Population	Speciality	New District	New Speciality
0	Dhaka	1432	2250000	Administrative	6	0
1	Gazipur	879	567984	Industry	9	3
2	Narayanganj	576	53426	Industry	11	3
3	Rajbari	897	65899	Agriculture	12	1
4	Chittagong	1234	1345566	Industry	3	3
5	Cox's Bazer	456	46567	Tourist	5	4
6	Bandarban	345	67579	Tourist	0	4
7	Khustia	432	57798	Business	10	2
8	Feni	543	67890	Agriculture	8	1
9	Comilla	564	77898	Tourist	4	4
10	Barisal	577	89750	Agriculture	1	1
11	Faridpur	567	77650	Agriculture	7	1
12	Rangpur	575	78966	Agriculture	13	1
13	Chadpur	876	67789	Tourist	2	4

OneHotEncoder

Setting up:

Encoding District:

```
# Encodes the column to array with binary values
dfn = oh.fit transform(df[['District']])
dfn
/usr/local/lib/python3.10/dist-packages/sklearn/preprocessing/_encod
 warnings.warn(
array([[0., 0., 0., 0., 0., 1., 0., 0., 0., 0., 0., 0., 0.],
    [0., 0., 0., 0., 0., 0., 0., 1., 0., 0., 0., 0.],
    [0., 0., 0., 0., 0., 0., 0., 0., 0., 1., 0., 0.],
    [0., 0., 0., 0., 1., 0., 0., 0., 0., 0., 0., 0., 0.],
    [0., 0., 0., 0., 0., 0., 0., 0., 1., 0., 0., 0.],
    [0., 0., 0., 0., 0., 0., 1., 0., 0., 0., 0., 0.],
    [0., 0., 0., 1., 0., 0., 0., 0., 0., 0., 0., 0., 0.],
    [0., 0., 0., 0., 0., 0., 1., 0., 0., 0., 0., 0., 0.],
```

Creating new DataFrame:

```
#Creating a new dataframe with the OneHotEncoded district data
ohendf = pd.DataFrame(dfn)
ohendf
     10
      11
```

Joining DataFrame:

	#Here we join both df and ohendf dataframe and put it in df dataframe df=df.join(ohendf) df																		
	District	Size	Population	Speciality	New District	New Speciality	0	1	2	3	4	5	6	7	8	9	10	11	12
0	Dhaka	1432	2250000	Administrative	6		0.0	0.0	0.0	0.0	0.0	1.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0
1	Gazipur	879	567984	Industry	9	3	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	1.0	0.0	0.0	0.0	0.0
2	Narayanganj	576	53426	Industry	11	3	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	1.0	0.0	0.0
3	Rajbari	897	65899	Agriculture	12	1	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	1.0	0.0
4	Chittagong	1234	1345566	Industry	3	3	0.0	0.0	1.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0
5	Cox's Bazer	456	46567	Tourist	5	4	0.0	0.0	0.0	0.0	1.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0
6	Bandarban	345	67579	Tourist		4	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0
7	Khustia	432	57798	Business	10	2	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	1.0	0.0	0.0	0.0
8	Feni	543	67890	Agriculture	8	1	0.0	0.0	0.0	0.0	0.0	0.0	0.0	1.0	0.0	0.0	0.0	0.0	0.0
9	Comilla	564	77898	Tourist	4	4	0.0	0.0	0.0	1.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0
10	Barisal	577	89750	Agriculture	1	1	1.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0
11	Faridpur	567	77650	Agriculture	7	1	0.0	0.0	0.0	0.0	0.0	0.0	1.0	0.0	0.0	0.0	0.0	0.0	0.0
12	Rangpur	575	78966	Agriculture	13	1	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	1.0
13	Chadpur	876	67789	Tourist	2	4	0.0	1.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0