

# pre-processing-techniques

February 11, 2025

```
[1]: import pandas as pd
import numpy as np
```

```
[2]: dataset = pd.read_csv('tested.csv')
dataset.head()
```

```
[2]:
```

	PassengerId	Survived	Pclass	\
0	892	0	3	
1	893	1	3	
2	894	0	2	
3	895	0	3	
4	896	1	3	

  

	Name	Sex	Age	SibSp	Parch	\
0	Kelly, Mr. James	male	34.5	0	0	
1	Wilkes, Mrs. James (Ellen Needs)	female	47.0	1	0	
2	Myles, Mr. Thomas Francis	male	62.0	0	0	
3	Wirz, Mr. Albert	male	27.0	0	0	
4	Hirvonen, Mrs. Alexander (Helga E Lindqvist)	female	22.0	1	1	

  

	Ticket	Fare	Cabin	Embarked
0	330911	7.8292	NaN	Q
1	363272	7.0000	NaN	S
2	240276	9.6875	NaN	Q
3	315154	8.6625	NaN	S
4	3101298	12.2875	NaN	S

```
[3]: dataset.info()
```

```
<class 'pandas.core.frame.DataFrame'>
RangeIndex: 418 entries, 0 to 417
Data columns (total 12 columns):
#   Column          Non-Null Count  Dtype
---  -
0   PassengerId     418 non-null   int64
1   Survived        418 non-null   int64
2   Pclass          418 non-null   int64
3   Name            418 non-null   object
```

```

4   Sex          418 non-null   object
5   Age          332 non-null   float64
6   SibSp        418 non-null   int64
7   Parch        418 non-null   int64
8   Ticket       418 non-null   object
9   Fare         417 non-null   float64
10  Cabin        91 non-null    object
11  Embarked     418 non-null   object
dtypes: float64(2), int64(5), object(5)
memory usage: 39.3+ KB

```

```
[4]: dataset.describe()
```

```

[4]:      PassengerId  Survived  Pclass    Age  SibSp  \
count    418.000000    418.000000    418.000000    332.000000    418.000000
mean     1100.500000     0.363636     2.265550     30.272590     0.447368
std       120.810458     0.481622     0.841838     14.181209     0.896760
min        892.000000     0.000000     1.000000      0.170000     0.000000
25%       996.250000     0.000000     1.000000     21.000000     0.000000
50%      1100.500000     0.000000     3.000000     27.000000     0.000000
75%      1204.750000     1.000000     3.000000     39.000000     1.000000
max      1309.000000     1.000000     3.000000     76.000000     8.000000

      Parch    Fare
count    418.000000    417.000000
mean         0.392344    35.627188
std         0.981429    55.907576
min         0.000000     0.000000
25%         0.000000     7.895800
50%         0.000000    14.454200
75%         0.000000    31.500000
max         9.000000   512.329200

```

```
[5]: dataset.isnull().sum()
```

```

[5]: PassengerId    0
     Survived      0
     Pclass        0
     Name          0
     Sex           0
     Age          86
     SibSp         0
     Parch         0
     Ticket        0
     Fare          1
     Cabin        327
     Embarked      0

```

dtype: int64

```
[6]: d1 = dataset['Age'].mean()  
d1
```

```
[6]: np.float64(30.272590361445783)
```

```
[7]: d1 =round(d1)  
d1
```

```
[7]: 30
```

```
[8]: dataset.isnull().sum()
```

```
[8]: PassengerId      0  
Survived           0  
Pclass            0  
Name              0  
Sex               0  
Age              86  
SibSp            0  
Parch            0  
Ticket           0  
Fare             1  
Cabin           327  
Embarked         0  
dtype: int64
```

```
[10]: d2 =round(dataset['Fare'].mean())  
dataset['Fare'] =dataset['Fare'].fillna(d2)  
dataset.isnull().sum()
```

```
[10]: PassengerId      0  
Survived           0  
Pclass            0  
Name              0  
Sex               0  
Age              86  
SibSp            0  
Parch            0  
Ticket           0  
Fare             0  
Cabin           327  
Embarked         0  
dtype: int64
```

```
[11]: dataset['Cabin'].unique()
```

```
[11]: array([nan, 'B45', 'E31', 'B57 B59 B63 B66', 'B36', 'A21', 'C78', 'D34',
        'D19', 'A9', 'D15', 'C31', 'C23 C25 C27', 'F G63', 'B61', 'C53',
        'D43', 'C130', 'C132', 'C101', 'C55 C57', 'B71', 'C46', 'C116',
        'F', 'A29', 'G6', 'C6', 'C28', 'C51', 'E46', 'C54', 'C97', 'D22',
        'B10', 'F4', 'E45', 'E52', 'D30', 'B58 B60', 'E34', 'C62 C64',
        'A11', 'B11', 'C80', 'F33', 'C85', 'D37', 'C86', 'D21', 'C89',
        'F E46', 'A34', 'D', 'B26', 'C22 C26', 'B69', 'C32', 'B78',
        'F E57', 'F2', 'A18', 'C106', 'B51 B53 B55', 'D10 D12', 'E60',
        'E50', 'E39 E41', 'B52 B54 B56', 'C39', 'B24', 'D28', 'B41', 'C7',
        'D40', 'D38', 'C105'], dtype=object)
```

```
[12]: dataset['Cabin'].value_counts()
```

```
[12]: Cabin
B57 B59 B63 B66      3
B45                  2
C23 C25 C27          2
C78                  2
C31                  2
..
B41                  1
C7                   1
D40                  1
D38                  1
C105                 1
Name: count, Length: 76, dtype: int64
```

```
[13]: dataset['Cabin'] =dataset['Cabin'].ffill()
dataset.isnull().sum()
```

```
[13]: PassengerId      0
Survived              0
Pclass               0
Name                 0
Sex                  0
Age                 86
SibSp               0
Parch               0
Ticket              0
Fare                0
Cabin               12
Embarked            0
dtype: int64
```

```
[14]: dataset['Cabin'] =dataset['Cabin'].bfill()
dataset.isnull().sum()
```

```
[14]: PassengerId    0
      Survived      0
      Pclass        0
      Name          0
      Sex           0
      Age          86
      SibSp         0
      Parch         0
      Ticket        0
      Fare          0
      Cabin         0
      Embarked      0
      dtype: int64
```

```
[15]: dataset.head()
```

```
[15]:   PassengerId  Survived  Pclass  \
0          892         0        3
1          893         1        3
2          894         0        2
3          895         0        3
4          896         1        3
```

```

                                Name    Sex  Age  SibSp  Parch  \
0                        Kelly, Mr. James  male  34.5    0    0
1      Wilkes, Mrs. James (Ellen Needs)  female  47.0    1    0
2                Myles, Mr. Thomas Francis  male  62.0    0    0
3                        Wirz, Mr. Albert  male  27.0    0    0
4  Hirvonen, Mrs. Alexander (Helga E Lindqvist)  female  22.0    1    1
```

```

      Ticket    Fare  Cabin  Embarked
0   330911    7.8292   B45         Q
1   363272    7.0000   B45         S
2   240276    9.6875   B45         Q
3   315154    8.6625   B45         S
4   3101298  12.2875   B45         S
```

```
[ ]:
```