

onhrbunbj

April 29, 2024

```
[ ]: import numpy as np
import pandas as pd
import seaborn as sns
import matplotlib.pyplot as plt
from matplotlib import style
```

```
[ ]: test_df = pd.read_csv("test.csv")
train_df = pd.read_csv("train.csv")
```

```
[ ]: test_df.info
```

```
[ ]: <bound method DataFrame.info of      PassengerId  Pclass
Name \
0          892      3              Kelly, Mr. James
1          893      3      Wilkes, Mrs. James (Ellen Needs)
2          894      2      Myles, Mr. Thomas Francis
3          895      3      Wirz, Mr. Albert
4          896      3  Hirvonen, Mrs. Alexander (Helga E Lindqvist)
..      ...      ...
413        1305      3      Spector, Mr. Woolf
414        1306      1      Oliva y Ocana, Dona. Fermina
415        1307      3      Saether, Mr. Simon Sivertsen
416        1308      3      Ware, Mr. Frederick
417        1309      3      Peter, Master. Michael J

      Sex  Age  SibSp  Parch      Ticket     Fare Cabin Embarked
0   male  34.5    0     0      330911    7.8292   NaN      Q
1  female  47.0    1     0      363272    7.0000   NaN      S
2   male  62.0    0     0      240276    9.6875   NaN      Q
3   male  27.0    0     0      315154    8.6625   NaN      S
4  female  22.0    1     1      3101298   12.2875   NaN      S
..      ...  ...  ...  ...
413  male   NaN    0     0      A.5. 3236    8.0500   NaN      S
414  female  39.0    0     0      PC 17758  108.9000  C105      C
415  male  38.5    0     0  SOTON/O.Q. 3101262    7.2500   NaN      S
416  male   NaN    0     0      359309    8.0500   NaN      S
417  male   NaN    1     1      2668    22.3583   NaN      C
```

```
[418 rows x 11 columns]>
```

```
[ ]: train_df.info()
```

```
<class 'pandas.core.frame.DataFrame'>
RangeIndex: 891 entries, 0 to 890
Data columns (total 12 columns):
 #   Column          Non-Null Count  Dtype
---  -
 0   PassengerId     891 non-null   int64
 1   Survived        891 non-null   int64
 2   Pclass          891 non-null   int64
 3   Name            891 non-null   object
 4   Sex             891 non-null   object
 5   Age            714 non-null   float64
 6   SibSp           891 non-null   int64
 7   Parch          891 non-null   int64
 8   Ticket         891 non-null   object
 9   Fare           891 non-null   float64
10   Cabin          204 non-null   object
11   Embarked       889 non-null   object
dtypes: float64(2), int64(5), object(5)
memory usage: 83.7+ KB
```

```
[ ]: total = train_df.isnull().sum()
total
```

```
[ ]: 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
```

```
[ ]: total = train_df.isnull().sum().sort_values(ascending=False)
percent_1 = train_df.isnull().sum()/train_df.isnull().count()*100
percent_2 = (round(percent_1, 1)).sort_values(ascending=False) #rounds the
↳ values in percent_1 to one decimal place and sorts them in descending order.
missing_data = pd.concat([total, percent_2], axis=1, keys=['Total', '%'])
```

```
missing_data
```

```
[ ]:      Total      %
Cabin      687  77.1
Age        177  19.9
Embarked     2   0.2
PassengerId 0   0.0
Survived     0   0.0
Pclass       0   0.0
Name         0   0.0
Sex          0   0.0
SibSp        0   0.0
Parch        0   0.0
Ticket       0   0.0
Fare         0   0.0
```

```
[ ]: train_df.duplicated()
```

```
[ ]: 0      False
1      False
2      False
3      False
4      False
...
886     False
887     False
888     False
889     False
890     False
Length: 891, dtype: bool
```

```
[ ]: train_df.drop_duplicates(inplace=True)
train_df
```

```
[ ]:      PassengerId  Survived  Pclass  \
0                1         0        3
1                2         1        1
2                3         1        3
3                4         1        1
4                5         0        3
..            ...      ...      ...
886            887         0        2
887            888         1        1
888            889         0        3
889            890         1        1
890            891         0        3
```

	Name	Sex	Age	SibSp	\
0	Braund, Mr. Owen Harris	male	22.0	1	
1	Cumings, Mrs. John Bradley (Florence Briggs Th...	female	38.0	1	
2	Heikkinen, Miss. Laina	female	26.0	0	
3	Futrelle, Mrs. Jacques Heath (Lily May Peel)	female	35.0	1	
4	Allen, Mr. William Henry	male	35.0	0	
..	
886	Montvila, Rev. Juozas	male	27.0	0	
887	Graham, Miss. Margaret Edith	female	19.0	0	
888	Johnston, Miss. Catherine Helen "Carrie"	female	NaN	1	
889	Behr, Mr. Karl Howell	male	26.0	0	
890	Dooley, Mr. Patrick	male	32.0	0	

	Parch	Ticket	Fare	Cabin	Embarked
0	0	A/5 21171	7.2500	NaN	S
1	0	PC 17599	71.2833	C85	C
2	0	STON/O2. 3101282	7.9250	NaN	S
3	0	113803	53.1000	C123	S
4	0	373450	8.0500	NaN	S
..	
886	0	211536	13.0000	NaN	S
887	0	112053	30.0000	B42	S
888	2	W./C. 6607	23.4500	NaN	S
889	0	111369	30.0000	C148	C
890	0	370376	7.7500	NaN	Q

[891 rows x 12 columns]

```
[ ]: train_df.isnull()
```

```
[ ]:
   PassengerId  Survived  Pclass   Name    Sex    Age  SibSp  Parch  Ticket  \
0          False     False   False  False  False  False  False  False  False
1          False     False   False  False  False  False  False  False  False
2          False     False   False  False  False  False  False  False  False
3          False     False   False  False  False  False  False  False  False
4          False     False   False  False  False  False  False  False  False
..          ...      ...    ...    ...    ...    ...    ...    ...
886         False     False   False  False  False  False  False  False  False
887         False     False   False  False  False  False  False  False  False
888         False     False   False  False  False  True   False  False  False
889         False     False   False  False  False  False  False  False  False
890         False     False   False  False  False  False  False  False  False

   Fare  Cabin  Embarked
0   False   True    False
1   False  False    False
2   False   True    False
```

```

3    False False    False
4    False  True    False
..    ...    ...    ...
886  False  True    False
887  False False    False
888  False  True    False
889  False False    False
890  False  True    False

```

[891 rows x 12 columns]

```
[ ]: train_df.dropna()
```

```
[ ]:
   PassengerId  Survived  Pclass  \
1             2         1       1
3             4         1       1
6             7         0       1
10            11         1       3
11            12         1       1
..          ...     ...     ...
871           872         1       1
872           873         0       1
879           880         1       1
887           888         1       1
889           890         1       1

```

```

                                     Name    Sex  Age  SibSp  \
1  Cumings, Mrs. John Bradley (Florence Briggs Th...  female  38.0    1
3      Futrelle, Mrs. Jacques Heath (Lily May Peel)  female  35.0    1
6                      McCarthy, Mr. Timothy J    male  54.0    0
10                     Sandstrom, Miss. Marguerite Rut  female   4.0    1
11                     Bonnell, Miss. Elizabeth  female  58.0    0
..          ...     ...     ...
871  Beckwith, Mrs. Richard Leonard (Sallie Monypeny)  female  47.0    1
872                      Carlsson, Mr. Frans Olof    male  33.0    0
879  Potter, Mrs. Thomas Jr (Lily Alexenia Wilson)  female  56.0    0
887                      Graham, Miss. Margaret Edith  female  19.0    0
889                      Behr, Mr. Karl Howell    male  26.0    0

```

```

   Parch  Ticket   Fare  Cabin Embarked
1      0  PC 17599  71.2833    C85        C
3      0  113803  53.1000   C123        S
6      0   17463  51.8625   E46        S
10     1  PP 9549  16.7000    G6        S
11     0  113783  26.5500   C103        S
..    ...     ...     ...
871    1   11751  52.5542   D35        S

```

872	0	695	5.0000	B51	B53	B55	S
879	1	11767	83.1583			C50	C
887	0	112053	30.0000			B42	S
889	0	111369	30.0000			C148	C

[183 rows x 12 columns]

```
[ ]: train_df.fillna(0)
```

```
[ ]:
      PassengerId  Survived  Pclass  \
0                1         0       3
1                2         1       1
2                3         1       3
3                4         1       1
4                5         0       3
..            ...         ...     ...
886            887         0       2
887            888         1       1
888            889         0       3
889            890         1       1
890            891         0       3
```

	Name	Sex	Age	SibSp	\
0	Braund, Mr. Owen Harris	male	22.0	1	
1	Cumings, Mrs. John Bradley (Florence Briggs Th...	female	38.0	1	
2	Heikkinen, Miss. Laina	female	26.0	0	
3	Futrelle, Mrs. Jacques Heath (Lily May Peel)	female	35.0	1	
4	Allen, Mr. William Henry	male	35.0	0	
..	
886	Montvila, Rev. Juozas	male	27.0	0	
887	Graham, Miss. Margaret Edith	female	19.0	0	
888	Johnston, Miss. Catherine Helen "Carrie"	female	0.0	1	
889	Behr, Mr. Karl Howell	male	26.0	0	
890	Dooley, Mr. Patrick	male	32.0	0	

	Parch	Ticket	Fare	Cabin	Embarked
0	0	A/5 21171	7.2500	0	S
1	0	PC 17599	71.2833	C85	C
2	0	STON/O2. 3101282	7.9250	0	S
3	0	113803	53.1000	C123	S
4	0	373450	8.0500	0	S
..	
886	0	211536	13.0000	0	S
887	0	112053	30.0000	B42	S
888	2	W./C. 6607	23.4500	0	S
889	0	111369	30.0000	C148	C
890	0	370376	7.7500	0	Q

[891 rows x 12 columns]

```
[ ]: train_df.dropna()
```

```
[ ]:
   PassengerId  Survived  Pclass \
1             2         1       1
3             4         1       1
6             7         0       1
10            11         1       3
11            12         1       1
..          ...         ...     ...
871           872         1       1
872           873         0       1
879           880         1       1
887           888         1       1
889           890         1       1
```

```

              Name      Sex  Age  SibSp \
1  Cumings, Mrs. John Bradley (Florence Briggs Th...  female  38.0      1
3      Futrelle, Mrs. Jacques Heath (Lily May Peel)  female  35.0      1
6                      McCarthy, Mr. Timothy J      male  54.0      0
10                     Sandstrom, Miss. Marguerite Rut  female   4.0      1
11                     Bonnell, Miss. Elizabeth      female  58.0      0
..          ...         ...     ...
871  Beckwith, Mrs. Richard Leonard (Sallie Monypeny)  female  47.0      1
872                      Carlsson, Mr. Frans Olof      male  33.0      0
879  Potter, Mrs. Thomas Jr (Lily Alexenia Wilson)  female  56.0      0
887                      Graham, Miss. Margaret Edith  female  19.0      0
889                      Behr, Mr. Karl Howell      male  26.0      0
```

```

   Parch  Ticket   Fare      Cabin Embarked
1      0  PC 17599  71.2833      C85        C
3      0  113803  53.1000     C123        S
6      0   17463  51.8625     E46        S
10     1  PP 9549  16.7000      G6        S
11     0  113783  26.5500     C103        S
..     ...     ...     ...     ...     ...
871     1   11751  52.5542     D35        S
872     0     695   5.0000  B51 B53 B55        S
879     1   11767  83.1583     C50        C
887     0  112053  30.0000     B42        S
889     0  111369  30.0000     C148        C
```

[183 rows x 12 columns]

```
[ ]: train_df.fillna(0,inplace=True)
train_df
```

```
[ ]:      PassengerId  Survived  Pclass  \
0             1         0         3
1             2         1         1
2             3         1         3
3             4         1         1
4             5         0         3
..          ...         ...         ...
886          887         0         2
887          888         1         1
888          889         0         3
889          890         1         1
890          891         0         3
```

```

                                Name      Sex  Age  SibSp  \
0                Braund, Mr. Owen Harris   male  22.0      1
1  Cumings, Mrs. John Bradley (Florence Briggs Th... female  38.0      1
2                Heikkinen, Miss. Laina   female  26.0      0
3  Futrelle, Mrs. Jacques Heath (Lily May Peel)   female  35.0      1
4                Allen, Mr. William Henry   male  35.0      0
..          ...         ...         ...         ...
886                Montvila, Rev. Juozas   male  27.0      0
887                Graham, Miss. Margaret Edith   female  19.0      0
888    Johnston, Miss. Catherine Helen "Carrie"   female   0.0      1
889                Behr, Mr. Karl Howell   male  26.0      0
890                Dooley, Mr. Patrick   male  32.0      0
```

```

      Parch      Ticket    Fare Cabin Embarked
0         0      A/5 21171    7.2500     0        S
1         0      PC 17599   71.2833    C85        C
2         0  STON/O2. 3101282    7.9250     0        S
3         0      113803   53.1000   C123        S
4         0      373450    8.0500     0        S
..      ...         ...         ...         ...
886        0      211536   13.0000     0        S
887        0      112053   30.0000   B42        S
888        2      W./C. 6607   23.4500     0        S
889        0      111369   30.0000   C148        C
890        0      370376    7.7500     0        Q
```

[891 rows x 12 columns]

```
[ ]: train_df.head(8)
```



```
[ ]: PassengerId  Survived  Pclass  \
0            1         0         3
1            2         1         1
2            3         1         3
3            4         1         1
4            5         0         3
5            6         0         3
6            7         0         1
7            8         0         3
```

```

                                Name      Sex  Age  SibSp  \
0                        Braund, Mr. Owen Harris    male  22.0      1
1  Cumings, Mrs. John Bradley (Florence Briggs Th... female  38.0      1
2                        Heikkinen, Miss. Laina  female  26.0      0
3  Futrelle, Mrs. Jacques Heath (Lily May Peel)  female  35.0      1
4                        Allen, Mr. William Henry    male  35.0      0
5                        Moran, Mr. James          male   NaN      0
6                        McCarthy, Mr. Timothy J    male  54.0      0
7                        Palsson, Master. Gosta Leonard  male   2.0      3
```

```

    Parch      Ticket    Fare Cabin Embarked
0      0   A/5 21171    7.2500   NaN        S
1      0    PC 17599   71.2833   C85        C
2      0 STON/O2. 3101282    7.9250   NaN        S
3      0    113803   53.1000  C123        S
4      0    373450    8.0500   NaN        S
5      0    330877    8.4583   NaN        Q
6      0    17463   51.8625   E46        S
7      1    349909   21.0750   NaN        S
```

Adressing Missing values

```
[ ]: total = train_df.isnull().sum().sort_values(ascending=False)
percent_1 = train_df.isnull().sum()/train_df.isnull().count()*100
percent_2 = (round(percent_1, 1)).sort_values(ascending=False)
missing_data = pd.concat([total, percent_2], axis=1, keys=['Total', '%'])
missing_data
```

```
[ ]:
      Total  %
Cabin     687  77.1
Age       177  19.9
Embarked     2   0.2
PassengerId    0   0.0
Survived       0   0.0
Pclass        0   0.0
Name          0   0.0
Sex           0   0.0
```

SibSp	0	0.0
Parch	0	0.0
Ticket	0	0.0
Fare	0	0.0