+ Code ——— + Text

# Data Science @Bharat Intern

## ▾ Done By Harsha Vardhan

```
import numpy as np

import pandas as pd

import seaborn as sns
import matplotlib.pyplot as plt
%matplotlib inline
from matplotlib import style
from sklearn import linear_model
from sklearn.linear_model import LogisticRegression
from sklearn.ensemble import RandomForestClassifier
from sklearn.linear_model import Perceptron
from sklearn.linear_model import SGDClassifier
from sklearn.tree import DecisionTreeClassifier
from sklearn.neighbors import KNeighborsClassifier
from sklearn.naive_bayes import GaussianNB
```

```
train_df = pd.read_csv('train.csv')
test_df = pd.read_csv('test.csv')

train_df['train_test'] = 1
test_df['train_test'] = 0
# test_df['Survived'] = np.NaN
all_data = pd.concat([train_df,test_df])

%matplotlib inline
all_data.columns
```

```
     Index(['PassengerId', 'Survived', 'Pclass', 'Name', 'Sex', 'Age', 'SibSp',
            'Parch', 'Ticket', 'Fare', 'Cabin', 'Embarked', 'train_test'],
           dtype='object')
```

```
train_df.head(10)
```

|   | PassengerId | Survived | Pclass | Name | Sex | Age | SibSp | Parch | Ticket | |
|---|---|---|---|---|---|---|---|---|---|---|
| **0** | 1 | 0 | 3 | Braund, Mr. Owen Harris | male | 22.0 | 1 | 0 | A/5 21171 | 7.: |
| **1** | 2 | 1 | 1 | Cumings, Mrs. John Bradley (Florence Briggs Th... | female | 38.0 | 1 | 0 | PC 17599 | 71.: |
| **2** | 3 | 1 | 3 | Heikkinen, Miss. Laina | female | 26.0 | 0 | 0 | STON/O2. 3101282 | 7.! |
| **3** | 4 | 1 | 1 | Futrelle, Mrs. Jacques Heath (Lily May Peel) | female | 35.0 | 1 | 0 | 113803 | 53. |
| **4** | 5 | 0 | 3 | Allen, Mr. William H.. | male | 35.0 | 0 | 0 | 373450 | 8.( |

```
test_df.head(10)
```

| | PassengerId | Pclass | Name | Sex | Age | SibSp | Parch | Ticket | Fare | Cabin |
|---|---|---|---|---|---|---|---|---|---|---|
| **0** | 892 | 3 | Kelly, Mr. James | male | 34.5 | 0 | 0 | 330911 | 7.8292 | NaN |
| **1** | 893 | 3 | Wilkes, Mrs. James (Ellen Needs) | female | 47.0 | 1 | 0 | 363272 | 7.0000 | NaN |
| **2** | 894 | 2 | Myles, Mr. Thomas | male | 62.0 | 0 | 0 | 240276 | 9.6875 | NaN |

```
train_df.info()
```

```
<class 'pandas.core.frame.DataFrame'>
RangeIndex: 891 entries, 0 to 890
Data columns (total 13 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
 12  train_test   891 non-null    int64
dtypes: float64(2), int64(6), object(5)
memory usage: 90.6+ KB
```
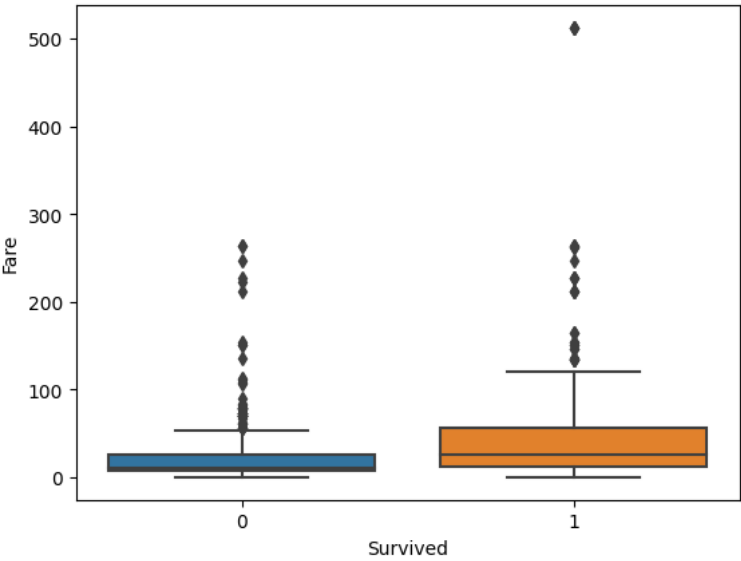
```
train_df.describe()
```

| | PassengerId | Survived | Pclass | Age | SibSp | Parch | F |
|---|---|---|---|---|---|---|---|
| **count** | 891.000000 | 891.000000 | 891.000000 | 714.000000 | 891.000000 | 891.000000 | 891.000 |
| **mean** | 446.000000 | 0.383838 | 2.308642 | 29.699118 | 0.523008 | 0.381594 | 32.204 |
| **std** | 257.353842 | 0.486592 | 0.836071 | 14.526497 | 1.102743 | 0.806057 | 49.693 |
| **min** | 1.000000 | 0.000000 | 1.000000 | 0.420000 | 0.000000 | 0.000000 | 0.000 |
| **25%** | 223.500000 | 0.000000 | 2.000000 | 20.125000 | 0.000000 | 0.000000 | 7.910 |
| **50%** | 446.000000 | 0.000000 | 3.000000 | 28.000000 | 0.000000 | 0.000000 | 14.454 |
| **75%** | 668.500000 | 1.000000 | 3.000000 | 38.000000 | 1.000000 | 0.000000 | 31.000 |
| **max** | 891.000000 | 1.000000 | 3.000000 | 80.000000 | 8.000000 | 6.000000 | 512.329 |

```
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.head(13)
```

```
train_df.columns.values
```

```
array(['PassengerId', 'Survived', 'Pclass', 'Name', 'Sex', 'Age', 'SibSp',
       'Parch', 'Ticket', 'Fare', 'Cabin', 'Embarked', 'train_test'],
      dtype=object)
```

```
sns.boxplot(x='Survived',y='Fare',data=train_df);
```



```
train_df[train_df['Fare']>300]
```

| | PassengerId | Survived | Pclass | Name | Sex | Age | SibSp | Parch | Ticket | F |
|---|---|---|---|---|---|---|---|---|---|---|
| **258** | 259 | 1 | 1 | Ward, Miss. Anna | female | 35.0 | 0 | 0 | PC 17755 | 512.3 |
| **679** | 680 | 1 | 1 | Cardeza, Mr. Thomas Drake Martinez | male | 36.0 | 0 | 1 | PC 17755 | 512.3 |

```
train_df[train_df['Name'].str.contains("Capt")]
```

| | PassengerId | Survived | Pclass | Name | Sex | Age | SibSp | Parch | Ticket | Fare | C |
|---|---|---|---|---|---|---|---|---|---|---|---|
| | | | | Crosby, | | | | | | | |

```
FacetGrid = sns.FacetGrid(train_df, col='Embarked', height=4, aspect=1.2)
FacetGrid.map(sns.pointplot, 'Pclass', 'Survived', 'Sex', ci=95.0, palette='deep', order=None, hue_order=None)
FacetGrid.add_legend();
```

```
/usr/local/lib/python3.10/dist-packages/seaborn/axisgrid.py:848: FutureWarning:

The `ci` parameter is deprecated. Use `errorbar=('ci', 95.0)` for the same effect.

  func(*plot_args, **plot_kwargs)
/usr/local/lib/python3.10/dist-packages/seaborn/axisgrid.py:848: FutureWarning:

The `ci` parameter is deprecated. Use `errorbar=('ci', 95.0)` for the same effect.
```
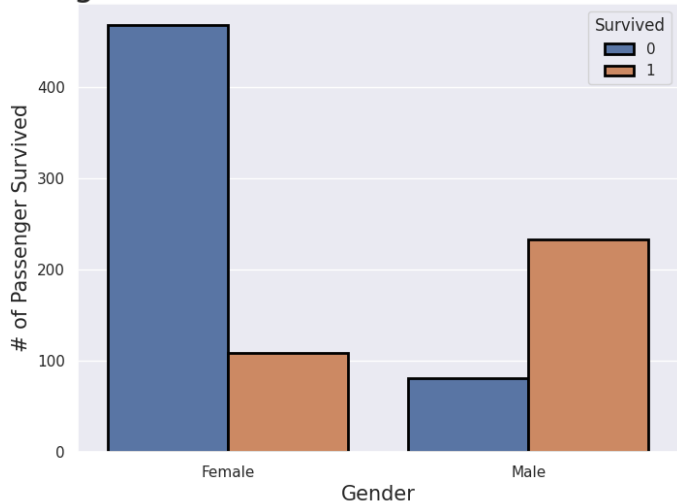
```
sns.set(style='darkgrid')
plt.subplots(figsize = (8,6))
ax=sns.countplot(x='Sex', data = train_df, hue='Survived', edgecolor=(0,0,0), linewidth=2)

plt.title('Passenger distribution of survived vs not-survived', fontsize=25)
plt.xlabel('Gender', fontsize=15)
plt.ylabel("# of Passenger Survived", fontsize = 15)
labels = ['Female', 'Male']

plt.xticks(sorted(train_df.Survived.unique()),labels);
```



```
train_df.groupby(['Sex']).mean()
```

```
<ipython-input-16-0a4a460c27a0>:1: FutureWarning: The default value of numeric_only
  train_df.groupby(['Sex']).mean()
```

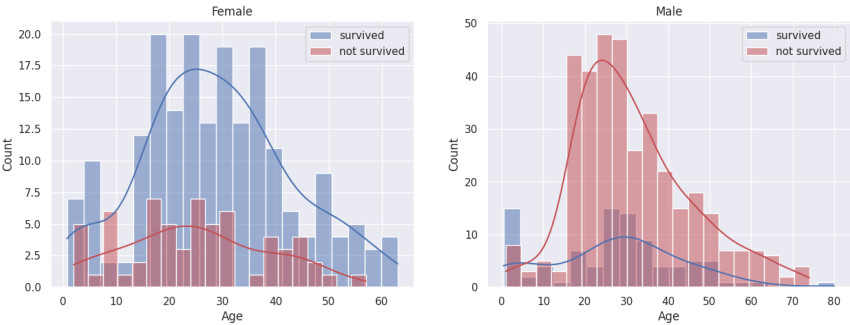|  | PassengerId | Survived | Pclass | Age | SibSp | Parch | Fare | trai |
|---|---|---|---|---|---|---|---|---|
| **Sex** | | | | | | | | |
| **female** | 431.028662 | 0.742038 | 2.159236 | 27.915709 | 0.694268 | 0.649682 | 44.479818 | |
| **male** | 454.147314 | 0.188908 | 2.389948 | 30.726645 | 0.429809 | 0.235702 | 25.523893 | |

```
train_df.groupby(['Sex','Pclass']).mean()
```

```
survived = 'survived'
not_survived = 'not survived'

fig, axes = plt.subplots(nrows=1, ncols=2, figsize=(15, 5))

women = train_df[train_df['Sex']=='female']
men = train_df[train_df['Sex']=='male']

ax = sns.histplot(women[women['Survived']==1].Age.dropna(), bins=20, label = survived, ax = axes[0],color='b', kde=True)
ax = sns.histplot(women[women['Survived']==0].Age.dropna(), bins=20, label = not_survived, ax = axes[0],color='r', kde=True)
ax.legend()
ax.set_title('Female')

ax = sns.histplot(men[men['Survived']==1].Age.dropna(), bins=20, label = survived, ax = axes[1],color='b', kde=True)
ax = sns.histplot(men[men['Survived']==0].Age.dropna(), bins=20, label = not_survived, ax = axes[1],color='r', kde=True)
ax.legend()
ax.set_title('Male');
```



```
train_df[train_df['Age']<18].groupby(['Sex','Pclass']).mean()
```
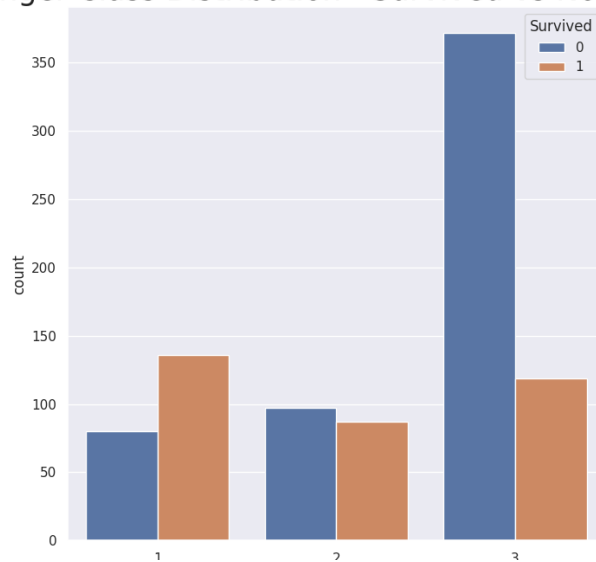
```
<ipython-input-20-828df17eba2f>:1: FutureWarning: The default value of numeric_only
  train_df[train_df['Age']<18].groupby(['Sex','Pclass']).mean()
```

| Sex | Pclass | PassengerId | Survived | Age | SibSp | Parch | Fare | trai |
|-----|--------|-------------|----------|-----|-------|-------|------|------|
| female | 1 | 525.375000 | 0.875000 | 14.125000 | 0.500000 | 0.875000 | 104.083337 | |
| | 2 | 369.250000 | 1.000000 | 8.333333 | 0.583333 | 1.083333 | 26.241667 | |
| | 3 | 374.942857 | 0.542857 | 8.428571 | 1.571429 | 1.057143 | 18.727977 | |
| male | 1 | 526.500000 | 1.000000 | 8.230000 | 0.500000 | 2.000000 | 116.072900 | |
| | 2 | 527.818182 | 0.818182 | 4.757273 | 0.727273 | 1.000000 | 25.659473 | |
| | 3 | 437.953488 | 0.232558 | 9.963256 | 2.069767 | 1.000000 | 22.752523 | |

```
plt.subplots(figsize = (8,8))
ax=sns.countplot(x='Pclass',hue='Survived',data=train_df)
plt.title("Passenger Class Distribution - Survived vs Non-Survived", fontsize = 25);
```

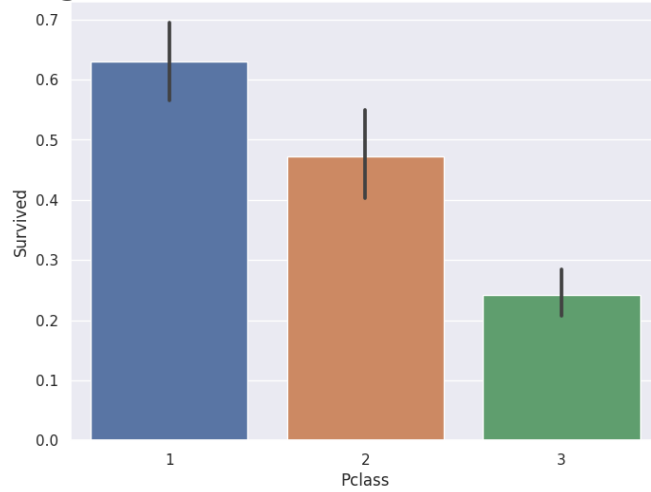## Passenger Class Distribution - Survived vs Non-Survived



```
plt.subplots(figsize=(10,8))
ax=sns.kdeplot(train_df.loc[(train_df['Survived'] == 0),'Pclass'],shade=True,color='r',label='Not Survived')
ax.legend()
ax=sns.kdeplot(train_df.loc[(train_df['Survived'] == 1),'Pclass'],shade=True,color='b',label='Survived')
ax.legend()

plt.title("Passenger Class Distribution - Survived vs Non-Survived", fontsize = 25)
labels = ['First', 'Second', 'Third']
plt.xticks(sorted(train_df.Pclass.unique()),labels);
```

```
<ipython-input-22-89cb45b403e7>:2: FutureWarning:
```
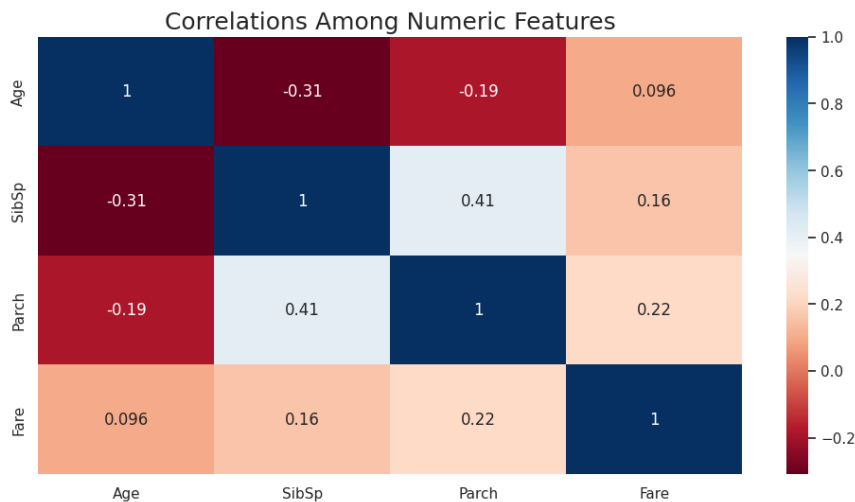
```
plt.subplots(figsize = (8,6))
sns.barplot(x='Pclass', y='Survived', data=train_df);
plt.title("Passenger Class Distribution - Survived Passengers", fontsize = 25);
```

## Passenger Class Distribution - Survived Passengers



```
df_num = train_df[['Age','SibSp','Parch','Fare']]
df_cat = train_df[['Survived','Pclass','Sex','Ticket','Cabin','Embarked']]
```

```
plt.subplots(figsize = (12,6))
sns.heatmap(df_num.corr(), annot=True,cmap="RdBu")
plt.title("Correlations Among Numeric Features", fontsize = 18);
```
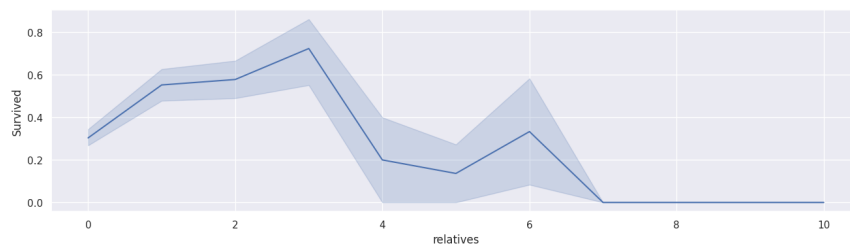
## Correlations Among Numeric Features



```
train_df = train_df.drop(['PassengerId'], axis=1)
train_df.head()
```

| | Survived | Pclass | Name | Sex | Age | SibSp | Parch | Ticket | Fare | Cabin | |
|---|---|---|---|---|---|---|---|---|---|---|---|
| **0** | 0 | 3 | Braund, Mr. Owen Harris | male | 22.0 | 1 | 0 | A/5 21171 | 7.2500 | NaN | |
| **1** | 1 | 1 | Cumings, Mrs. John Bradley (Florence | female | 38.0 | 1 | 0 | PC 17599 | 71.2833 | C85 | |

```
data = [train_df, test_df]
for dataset in data:
    dataset['relatives'] = dataset['SibSp'] + dataset['Parch']
    dataset.loc[dataset['relatives'] > 0, 'not_alone'] = 0
    dataset.loc[dataset['relatives'] == 0, 'not_alone'] = 1
    dataset['not_alone'] = dataset['not_alone'].astype(int)
train_df['not_alone'].value_counts()
```

```
    1    537
    0    354
    Name: not_alone, dtype: int64
```

```
plt.subplots(figsize = (16,4))
ax = sns.lineplot(x='relatives',y='Survived', data=train_df)
```



```
import re
deck = {"A": 1, "B": 2, "C": 3, "D": 4, "E": 5, "F": 6, "G": 7, "U": 8}
data = [train_df, test_df]

for dataset in data:
    dataset['Cabin'] = dataset['Cabin'].fillna("U0")
    dataset['Deck'] = dataset['Cabin'].map(lambda x: re.compile("([a-zA-Z]+)").search(x).group())
    dataset['Deck'] = dataset['Deck'].map(deck)
    dataset['Deck'] = dataset['Deck'].fillna(0)
    dataset['Deck'] = dataset['Deck'].astype(int)


train_df = train_df.drop(['Cabin'], axis=1)
test_df = test_df.drop(['Cabin'], axis=1)


data = [train_df, test_df]

for dataset in data:
    mean = train_df["Age"].mean()
    std = test_df["Age"].std()
    is_null = dataset["Age"].isnull().sum()

    rand_age = np.random.randint(mean - std, mean + std, size = is_null)

    age_slice = dataset["Age"].copy()
    age_slice[np.isnan(age_slice)] = rand_age
    dataset["Age"] = age_slice
    dataset["Age"] = train_df["Age"].astype(int)


train_df["Age"].isnull().sum()
```

```
    0
```

```
train_df['Embarked'].describe()
```

```
    count     889
    unique      3
```

```
        top        S
        freq     644
        Name: Embarked, dtype: object
```

```python
common_value = 'S'
data = [train_df, test_df]

for dataset in data:
    dataset['Embarked'] = dataset['Embarked'].fillna(common_value)
train_df['Embarked'].isnull().sum()
```

```
    0
```

```python
train_df.info()
```

```
    <class 'pandas.core.frame.DataFrame'>
    RangeIndex: 891 entries, 0 to 890
    Data columns (total 14 columns):
     #   Column      Non-Null Count  Dtype
    ---  ------      --------------  -----
     0   Survived    891 non-null    int64
     1   Pclass      891 non-null    int64
     2   Name        891 non-null    object
     3   Sex         891 non-null    object
     4   Age         891 non-null    int64
     5   SibSp       891 non-null    int64
     6   Parch       891 non-null    int64
     7   Ticket      891 non-null    object
     8   Fare        891 non-null    float64
     9   Embarked    891 non-null    object
     10  train_test  891 non-null    int64
     11  relatives   891 non-null    int64
     12  not_alone   891 non-null    int64
     13  Deck        891 non-null    int64
    dtypes: float64(1), int64(9), object(4)
    memory usage: 97.6+ KB
```

```python
data = [train_df, test_df]

for dataset in data:
    dataset['Fare'] = dataset['Fare'].fillna(0)
    dataset['Fare'] = dataset['Fare'].astype(int)
train_df.info()
```

```
    <class 'pandas.core.frame.DataFrame'>
    RangeIndex: 891 entries, 0 to 890
    Data columns (total 14 columns):
     #   Column      Non-Null Count  Dtype
    ---  ------      --------------  -----
     0   Survived    891 non-null    int64
     1   Pclass      891 non-null    int64
     2   Name        891 non-null    object
     3   Sex         891 non-null    object
     4   Age         891 non-null    int64
     5   SibSp       891 non-null    int64
     6   Parch       891 non-null    int64
     7   Ticket      891 non-null    object
     8   Fare        891 non-null    int64
     9   Embarked    891 non-null    object
     10  train_test  891 non-null    int64
     11  relatives   891 non-null    int64
     12  not_alone   891 non-null    int64
     13  Deck        891 non-null    int64
    dtypes: int64(10), object(4)
    memory usage: 97.6+ KB
```

```python
data = [train_df, test_df]
titles = {"Mr": 1, "Miss": 2, "Mrs": 3, "Master": 4, "Other": 5}

for dataset in data:

    dataset['Title'] = dataset.Name.str.extract('([A-Za-z]+)\.', expand=False)


    dataset['Title'] = dataset['Title'].replace(['Lady', 'Countess','Capt', 'Col','Don', 'Dr','Major', 'Rev', 'Sir', 'Jonkheer', 'Dona']
    dataset['Title'] = dataset['Title'].replace('Mlle', 'Miss')
    dataset['Title'] = dataset['Title'].replace('Ms', 'Miss')
    dataset['Title'] = dataset['Title'].replace('Mme', 'Mrs')


    dataset['Title'] = dataset['Title'].map(titles)
```

```python
    dataset['Title'] = dataset['Title'].fillna(0)
```

```python
train_df = train_df.drop(['Name'], axis=1)
test_df = test_df.drop(['Name'], axis=1)
```

```python
train_df.head()
```

|   | Survived | Pclass | Sex | Age | SibSp | Parch | Ticket | Fare | Embarked | train_test |
|---|---|---|---|---|---|---|---|---|---|---|
| **0** | 0 | 3 | male | 22 | 1 | 0 | A/5 21171 | 7 | S | 1 |
| **1** | 1 | 1 | female | 38 | 1 | 0 | PC 17599 | 71 | C | 1 |
| **2** | 1 | 3 | female | 26 | 0 | 0 | STON/O2. 3101282 | 7 | S | 1 |
| **3** | 1 | 1 | female | 35 | 1 | 0 | 113803 | 53 | S | 1 |

```python
genders = {"male": 0, "female": 1}
data = [train_df, test_df]

for dataset in data:
    dataset['Sex'] = dataset['Sex'].map(genders)
```

```python
train_df.head()
```

|   | Survived | Pclass | Sex | Age | SibSp | Parch | Ticket | Fare | Embarked | train_test | re |
|---|---|---|---|---|---|---|---|---|---|---|---|
| **0** | 0 | 3 | 0 | 22 | 1 | 0 | A/5 21171 | 7 | S | 1 | |
| **1** | 1 | 1 | 1 | 38 | 1 | 0 | PC 17599 | 71 | C | 1 | |
| **2** | 1 | 3 | 1 | 26 | 0 | 0 | STON/O2. 3101282 | 7 | S | 1 | |
| **3** | 1 | 1 | 1 | 35 | 1 | 0 | 113803 | 53 | S | 1 | |

```python
train_df['Ticket'].describe()
```

```
count        891
unique       681
top       347082
freq           7
Name: Ticket, dtype: object
```

```python
train_df = train_df.drop(['Ticket'], axis=1)
test_df = test_df.drop(['Ticket'], axis=1)
```

```python
train_df.head()
```

|   | Survived | Pclass | Sex | Age | SibSp | Parch | Fare | Embarked | train_test | relatives | not_alone | Deck | Tit |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| **0** | 0 | 3 | 0 | 22 | 1 | 0 | 7 | S | 1 | 1 | 0 | 8 | |
| **1** | 1 | 1 | 1 | 38 | 1 | 0 | 71 | C | 1 | 1 | 0 | 3 | |
| **2** | 1 | 3 | 1 | 26 | 0 | 0 | 7 | S | 1 | 0 | 1 | 8 | |
| **3** | 1 | 1 | 1 | 35 | 1 | 0 | 53 | S | 1 | 1 | 0 | 3 | |
| **4** | 0 | 3 | 0 | 35 | 0 | 0 | 8 | S | 1 | 0 | 1 | 8 | |

```python
ports = {"S": 0, "C": 1, "Q": 2}
data = [train_df, test_df]

for dataset in data:
    dataset['Embarked'] = dataset['Embarked'].map(ports)
```

```python
train_df.head()
```

| | Survived | Pclass | Sex | Age | SibSp | Parch | Fare | Embarked | train_test | relatives | not_alone | Deck | Tit |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| 0 | 0 | 3 | 0 | 22 | 1 | 0 | 7 | 0 | 1 | 1 | 0 | 8 | |

```
data = [train_df, test_df]
for dataset in data:
    dataset['Age'] = dataset['Age'].astype(int)
    dataset.loc[ dataset['Age'] <= 11, 'Age'] = 0
    dataset.loc[(dataset['Age'] > 11) & (dataset['Age'] <= 18), 'Age'] = 1
    dataset.loc[(dataset['Age'] > 18) & (dataset['Age'] <= 22), 'Age'] = 2
    dataset.loc[(dataset['Age'] > 22) & (dataset['Age'] <= 27), 'Age'] = 3
    dataset.loc[(dataset['Age'] > 27) & (dataset['Age'] <= 33), 'Age'] = 4
    dataset.loc[(dataset['Age'] > 33) & (dataset['Age'] <= 40), 'Age'] = 5
    dataset.loc[(dataset['Age'] > 40) & (dataset['Age'] <= 66), 'Age'] = 6
    dataset.loc[ dataset['Age'] > 66, 'Age'] = 6


train_df['Age'].value_counts()
```

```
4    163
6    157
5    151
3    136
2    116
1    100
0     68
Name: Age, dtype: int64
```

```
train_df.head()
```

| | Survived | Pclass | Sex | Age | SibSp | Parch | Fare | Embarked | train_test | relatives | n |
|---|---|---|---|---|---|---|---|---|---|---|---|
| 0 | 0 | 3 | 0 | 2 | 1 | 0 | 7 | 0 | 1 | 1 | |
| 1 | 1 | 1 | 1 | 5 | 1 | 0 | 71 | 1 | 1 | 1 | |
| 2 | 1 | 3 | 1 | 3 | 0 | 0 | 7 | 0 | 1 | 0 | |
| 3 | 1 | 1 | 1 | 5 | 1 | 0 | 53 | 0 | 1 | 1 | |
| 4 | 0 | 3 | 0 | 5 | 0 | 0 | 8 | 0 | 1 | 0 | |

```
pd.qcut(train_df['Fare'], q=6)
```

```
0        (-0.001, 7.0]
1        (52.0, 512.0]
2        (-0.001, 7.0]
3        (52.0, 512.0]
4          (7.0, 8.0]
             ...
886        (8.0, 14.0]
887       (26.0, 52.0]
888       (14.0, 26.0]
889       (26.0, 52.0]
890      (-0.001, 7.0]
Name: Fare, Length: 891, dtype: category
Categories (6, interval[float64, right]): [(-0.001, 7.0] < (7.0, 8.0] < (8.0, 14.0] < (14.0, 26.0] <
                                            (26.0, 52.0] < (52.0, 512.0]]
```

```
data = [train_df, test_df]

for dataset in data:
    dataset.loc[ dataset['Fare'] <= 7, 'Fare'] = 0
    dataset.loc[(dataset['Fare'] > 7) & (dataset['Fare'] <= 8), 'Fare'] = 1
    dataset.loc[(dataset['Fare'] > 8) & (dataset['Fare'] <= 14), 'Fare']   = 2
    dataset.loc[(dataset['Fare'] > 14) & (dataset['Fare'] <= 26), 'Fare']   = 3
    dataset.loc[(dataset['Fare'] > 26) & (dataset['Fare'] <= 52), 'Fare']   = 4
    dataset.loc[dataset['Fare'] > 52, 'Fare'] = 5
    dataset['Fare'] = dataset['Fare'].astype(int)


train_df.head(10)
```

| | Survived | Pclass | Sex | Age | SibSp | Parch | Fare | Embarked | train_test | relatives | n |
|---|---|---|---|---|---|---|---|---|---|---|---|
| **0** | 0 | 3 | 0 | 2 | 1 | 0 | 0 | 0 | 1 | 1 | |
| **1** | 1 | 1 | 1 | 5 | 1 | 0 | 5 | 1 | 1 | 1 | |
| **2** | 1 | 3 | 1 | 3 | 0 | 0 | 0 | 0 | 1 | 0 | |
| **3** | 1 | 1 | 1 | 5 | 1 | 0 | 5 | 0 | 1 | 1 | |

```
X_train = train_df.drop("Survived", axis=1)
Y_train = train_df["Survived"]
X_test  = test_df.drop("PassengerId", axis=1).copy()
```

## SGD

```
            9       1   2   1   1    1    0   4         1          1          1
sgd = linear_model.SGDClassifier(max_iter=5, tol=None)
sgd.fit(X_train, Y_train)

Y_pred = sgd.predict(X_test)

sgd.score(X_train, Y_train)
acc_sgd = round(sgd.score(X_train, Y_train) * 100, 2)

# Print score
print(round(acc_sgd,2,), "%")
```

```
    75.98 %
```

## Decision Tree

```
decision_tree = DecisionTreeClassifier()
decision_tree.fit(X_train, Y_train)

Y_pred = decision_tree.predict(X_test)

acc_decision_tree = round(decision_tree.score(X_train, Y_train) * 100, 2)

# Print score
print(round(acc_decision_tree,2,), "%")
```

```
    92.82 %
```

## Random Forest

```
random_forest = RandomForestClassifier(n_estimators=100)
random_forest.fit(X_train, Y_train)

Y_prediction = random_forest.predict(X_test)

random_forest.score(X_train, Y_train)
acc_random_forest = round(random_forest.score(X_train, Y_train) * 100, 2)

# Print score
print(round(acc_random_forest,2,), "%")
```

```
    92.82 %
```

## Logistic Regression

```
logreg = LogisticRegression()
logreg.fit(X_train, Y_train)

Y_pred = logreg.predict(X_test)

acc_log = round(logreg.score(X_train, Y_train) * 100, 2)

# Print score
print(round(acc_log,2,), "%")
```

```
    81.14 %
```

## ▾ KNN

```
knn = KNeighborsClassifier(n_neighbors = 3)
knn.fit(X_train, Y_train)

Y_pred = knn.predict(X_test)

acc_knn = round(knn.score(X_train, Y_train) * 100, 2)

# Print score
print(round(acc_knn,2,), "%")
```

```
      85.19 %
```

## ▾ Gaussian Naive Bayes

```
gaussian = GaussianNB()
gaussian.fit(X_train, Y_train)

Y_pred = gaussian.predict(X_test)

acc_gaussian = round(gaussian.score(X_train, Y_train) * 100, 2)

# Print score
print(round(acc_gaussian,2,), "%")
```

```
      78.56 %
```

## ▾ Perceptron

```
perceptron = Perceptron(max_iter=1000)
perceptron.fit(X_train, Y_train)

Y_pred = perceptron.predict(X_test)

acc_perceptron = round(perceptron.score(X_train, Y_train) * 100, 2)

# Print score
print(round(acc_perceptron,2,), "%")
```

```
      64.09 %
```

## ▾ Model evaluation

```
results = pd.DataFrame({
    'Model': ['KNN', 'Logistic Regression',
             'Random Forest', 'Naive Bayes', 'Perceptron',
             'Stochastic Gradient Decent',
             'Decision Tree'],
    'Score': [acc_knn, acc_log,
             acc_random_forest, acc_gaussian, acc_perceptron,
             acc_sgd, acc_decision_tree]})

result_df = results.sort_values(by='Score', ascending=False)
result_df = result_df.set_index('Score')
result_df.head(9)
```

Model  🪄  📊

Score

## ▾ K-Fold Cross Validation

85.49                        KNN

```
from sklearn.model_selection import cross_val_score

rf = RandomForestClassifier(n_estimators=100)
scores = cross_val_score(rf, X_train, Y_train, cv=10, scoring = "accuracy")
print("Scores:", scores)
print("Mean:", scores.mean())
print("Standard Deviation:", scores.std())
```

```
Scores: [0.78888889 0.83146067 0.75280899 0.86516854 0.86516854 0.85393258
 0.82022472 0.7752809  0.86516854 0.80898876]
Mean: 0.8227091136079899
Standard Deviation: 0.0386440089058461
```

## ▾ Random Forest

```
importances = pd.DataFrame({'feature':X_train.columns,'importance':np.round(random_forest.feature_importances_,3)})
importances = importances.sort_values('importance',ascending=False).set_index('feature')
importances.head(12)
```

|  | importance 🪄 📊 |
| --- | --- |
| **feature** | |
| **Title** | 0.206 |
| **Sex** | 0.160 |
| **Age** | 0.157 |
| **Fare** | 0.113 |
| **Deck** | 0.089 |
| **Pclass** | 0.087 |
| **relatives** | 0.067 |
| **Embarked** | 0.047 |
| **SibSp** | 0.039 |
| **Parch** | 0.023 |
| **not_alone** | 0.012 |
| **train_test** | 0.000 |

```
importances.plot.bar();
```

```python
train_df  = train_df.drop("not_alone", axis=1)
test_df  = test_df.drop("not_alone", axis=1)
```

```python
train_df  = train_df.drop("Parch", axis=1)
test_df  = test_df.drop("Parch", axis=1)
```

```python
X_train = train_df.drop("Survived", axis=1)
Y_train = train_df["Survived"]
X_test  = test_df.drop("PassengerId", axis=1).copy()
```

```python
random_forest = RandomForestClassifier(n_estimators=100, oob_score = True)
random_forest.fit(X_train, Y_train)

Y_prediction = random_forest.predict(X_test)

random_forest.score(X_train, Y_train)
acc_random_forest = round(random_forest.score(X_train, Y_train) * 100, 2)

# Print scores
print(round(acc_random_forest,2,), "%")
```

```
92.82 %
```

```python
importances = pd.DataFrame({'feature':X_train.columns,'importance':np.round(random_forest.feature_importances_,3)})
importances = importances.sort_values('importance',ascending=False).set_index('feature')
importances.head(12)
```

| feature | importance |
|---|---|
| Title | 0.193 |
| Sex | 0.185 |
| Age | 0.157 |
| Fare | 0.108 |
| Pclass | 0.092 |
| Deck | 0.087 |
| relatives | 0.082 |
| SibSp | 0.050 |
| Embarked | 0.046 |
| train_test | 0.000 |

```python
print("oob score:", round(random_forest.oob_score_, 4)*100, "%")
```

```
oob score: 80.92 %
```

## ▾ Hyperparameter Tuning

```python
def clf_performance(classifier, model_name):
    print(model_name)
    print('Best Score: ' + str(classifier.best_score_))
    print('Best Parameters: ' + str(classifier.best_params_))
```

1m 21s    completed at 8:40 PM

1m 21s    completed at 8:40 PM