Course: Laboratory Practice III

Course Code: 410246

Group Members :

Pravin Jain

Onkar Kulkarni

Class: BE

Mini-Project Title: Build a machine learning model that predicts the type of people who survived the Titanic shipwreck using passenger data (i.e. name, age, gender, socio-economic class, etc.).

Dataset Link: https://www.kaggle.com/competitions/titanic/data

```python
import numpy as np

# data processing
import pandas as pd

# data visualization
import seaborn as sns
%matplotlib inline
from matplotlib import pyplot as plt
from matplotlib import style

# Algorithms
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.svm import SVC, LinearSVC
from sklearn.naive_bayes import GaussianNB
```

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

```python
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
```

```
train_df.head(8)
```

| | PassengerId | Survived | Pclass | Name | Sex | Age | SibSp | Parch | Ticket | Far |
|---|---|---|---|---|---|---|---|---|---|---|
| **0** | 1 | 0 | 3 | Braund, Mr. Owen Harris | male | 22.0 | 1 | 0 | A/5 21171 | 7.250 |
| **1** | 2 | 1 | 1 | Cumings, Mrs. John Bradley (Florence Briggs Th... | female | 38.0 | 1 | 0 | PC 17599 | 71.283 |
| **2** | 3 | 1 | 3 | Heikkinen, Miss. Laina | female | 26.0 | 0 | 0 | STON/O2. 3101282 | 7.925 |

```
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(5)
```
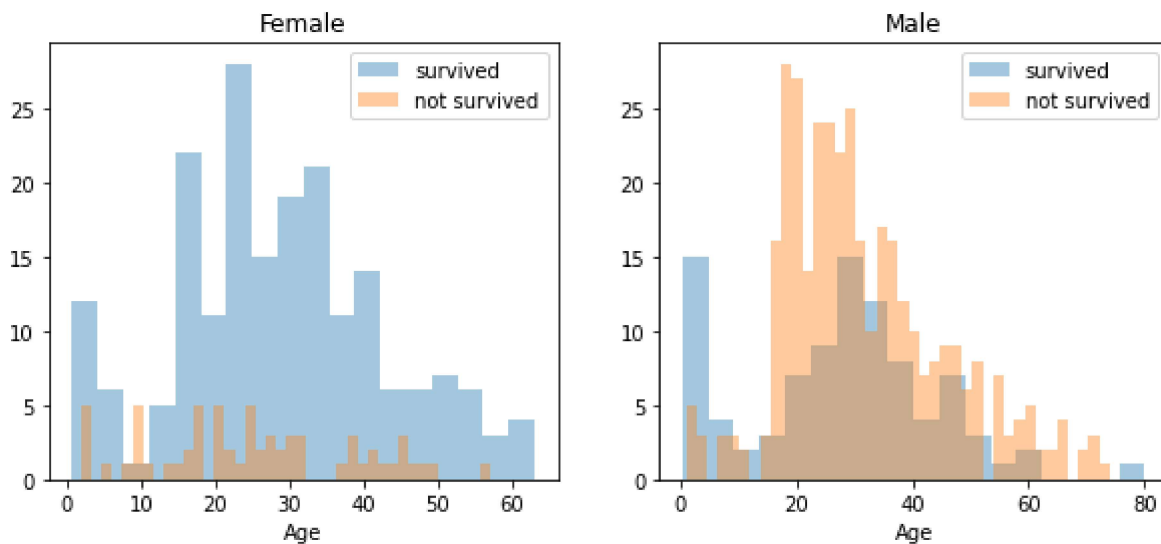
| | Total | % |
|---|---|---|
| **Cabin** | 687 | 77.1 |
| **Age** | 177 | 19.9 |
| **Embarked** | 2 | 0.2 |
| **PassengerId** | 0 | 0.0 |
| **Survived** | 0 | 0.0 |

```
train_df.columns.values
```

```
    array(['PassengerId', 'Survived', 'Pclass', 'Name', 'Sex', 'Age', 'SibSp',
           'Parch', 'Ticket', 'Fare', 'Cabin', 'Embarked'], dtype=object)
```
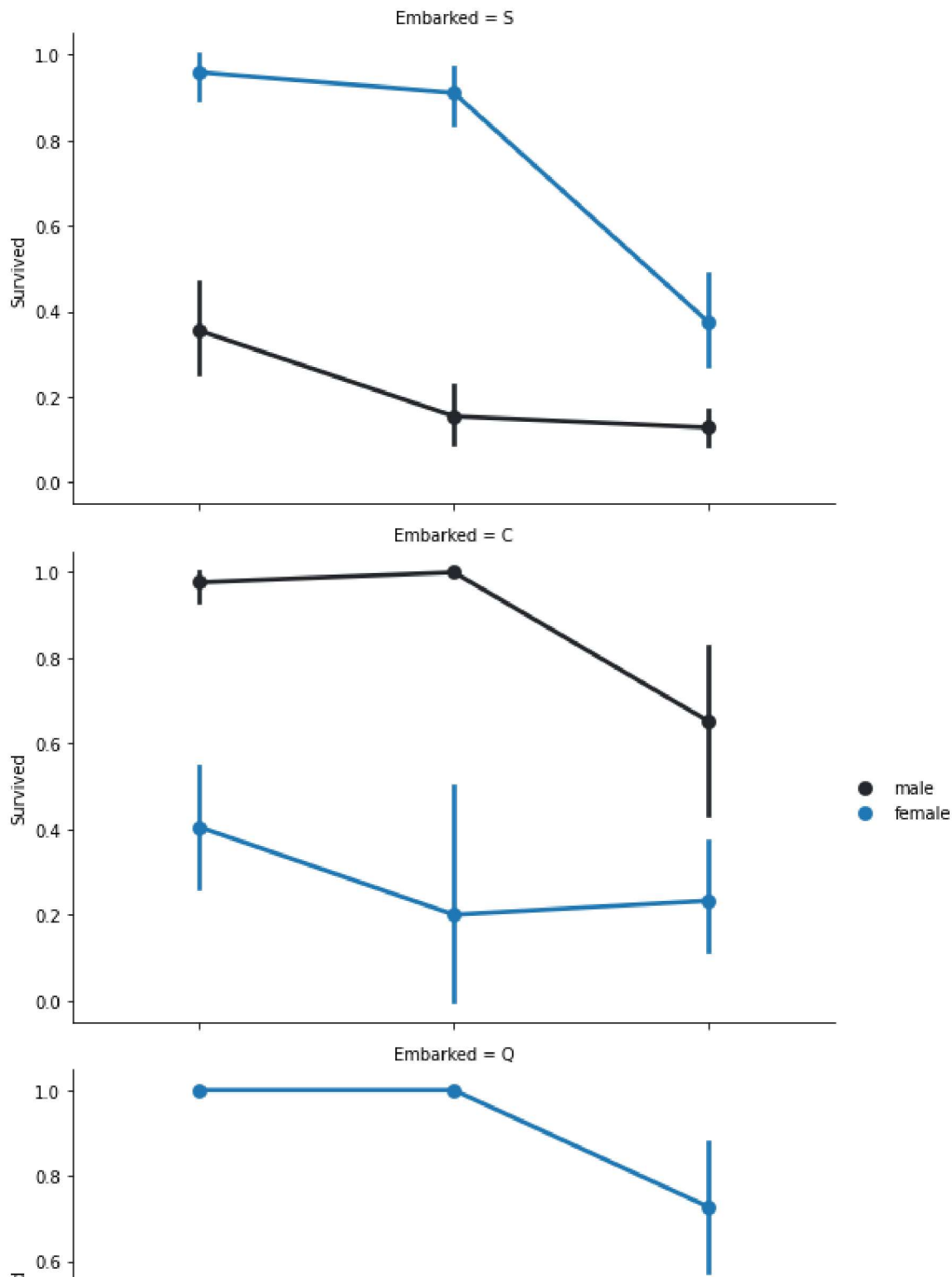
```python
survived = 'survived'
not_survived = 'not survived'
fig, axes = plt.subplots(nrows=1, ncols=2,figsize=(10, 4))
women = train_df[train_df['Sex']=='female']
men = train_df[train_df['Sex']=='male']
ax = sns.distplot(women[women['Survived']==1].Age.dropna(), bins=18, label = survived, ax = a
ax = sns.distplot(women[women['Survived']==0].Age.dropna(), bins=40, label = not_survived, ax
ax.legend()
ax.set_title('Female')
ax = sns.distplot(men[men['Survived']==1].Age.dropna(), bins=18, label = survived, ax = axes[
ax = sns.distplot(men[men['Survived']==0].Age.dropna(), bins=40, label = not_survived, ax = a
ax.legend()
_ = ax.set_title('Male')
```

```
/usr/local/lib/python3.7/dist-packages/seaborn/distributions.py:2619: FutureWarning: `di
  warnings.warn(msg, FutureWarning)
```
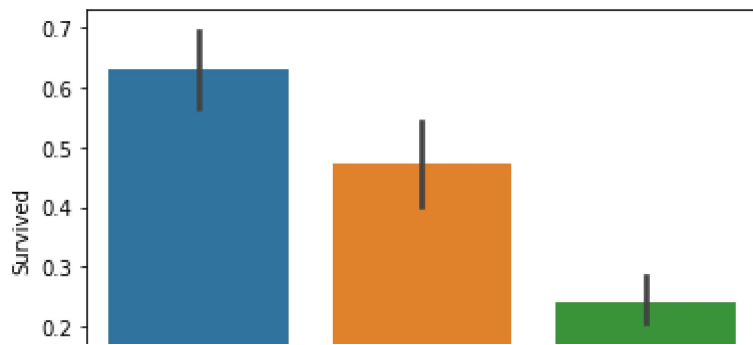


```python
FacetGrid = sns.FacetGrid(train_df, row='Embarked', size=4.5, aspect=1.6)
FacetGrid.map(sns.pointplot, 'Pclass', 'Survived', 'Sex', palette=None,  order=None, hue_orde
FacetGrid.add_legend()
```

<seaborn.axisgrid.FacetGrid at 0x7f9c5b3e21d0>



```
sns.barplot(x='Pclass', y='Survived', data=train_df)
```

<matplotlib.axes._subplots.AxesSubplot at 0x7f9c5b3db810>



```
grid = sns.FacetGrid(train_df, col='Survived', row='Pclass', size=2.2, aspect=1.6)
grid.map(plt.hist, 'Age', alpha=.5, bins=20)
grid.add_legend();
```

```
/usr/local/lib/python3.7/dist-packages/seaborn/axisgrid.py:337: UserWarning: The `size`
  warnings.warn(msg, UserWarning)
```
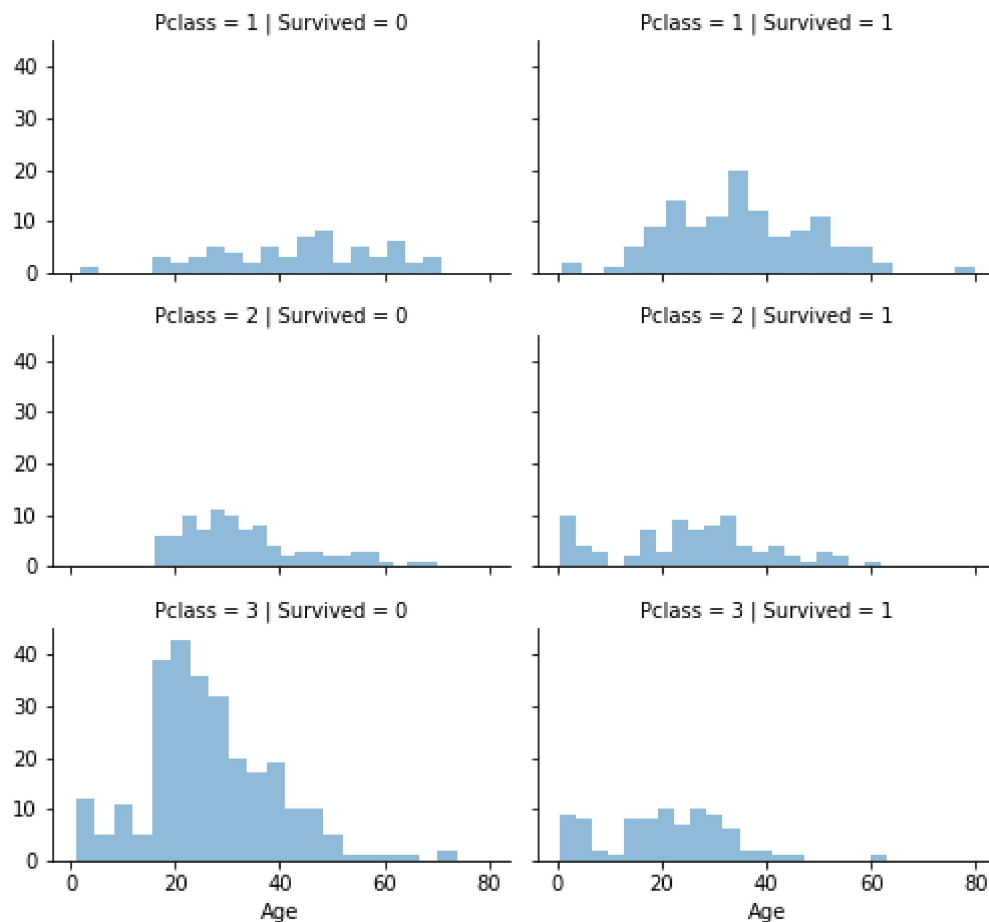


```
data = [train_df, test_df]
for dataset in data:
    dataset['relatives'] = dataset['SibSp'] + dataset['Parch']
    dataset.loc[dataset['relatives'] > 0, 'not_alone'] = 0
    dataset.loc[dataset['relatives'] == 0, 'not_alone'] = 1
    dataset['not_alone'] = dataset['not_alone'].astype(int)
train_df['not_alone'].value_counts()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

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

## DATA PRE-PROCESSING

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


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).grou
    dataset['Deck'] = dataset['Deck'].map(deck)
    dataset['Deck'] = dataset['Deck'].fillna(0)
    dataset['Deck'] = dataset['Deck'].astype(int)
# we can now drop the cabin feature
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()
    # compute random numbers between the mean, std and is_null
    rand_age = np.random.randint(mean - std, mean + std, size = is_null)
    # fill NaN values in Age column with random values generated
    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


common_value = 'S'
```

```
data = [train_df, test_df]

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

## CONVERTING FEATURES

```
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   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  relatives  891 non-null    int64
 11  not_alone  891 non-null    int64
 12  Deck       891 non-null    int64
dtypes: float64(1), int64(8), object(4)
memory usage: 90.6+ KB
```

```
data = [train_df, test_df]

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

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

for dataset in data:
    # extract titles
    dataset['Title'] = dataset.Name.str.extract(' ([A-Za-z]+)\.', expand=False)
    # replace titles with a more common title or as Rare
    dataset['Title'] = dataset['Title'].replace(['Lady', 'Countess','Capt', 'Col','Don', 'Dr'
                                                'Major', 'Rev', 'Sir', 'Jonkheer', 'Dona'], 'Rare
    dataset['Title'] = dataset['Title'].replace('Mlle', 'Miss')
    dataset['Title'] = dataset['Title'].replace('Ms', 'Miss')
    dataset['Title'] = dataset['Title'].replace('Mme', 'Mrs')
    # convert titles into numbers
```

```python
    dataset['Title'] = dataset['Title'].map(titles)
    # filling NaN with 0, to get safe
    dataset['Title'] = dataset['Title'].fillna(0)
train_df = train_df.drop(['Name'], axis=1)
test_df = test_df.drop(['Name'], axis=1)


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

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


train_df['Ticket'].describe()

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


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


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

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


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    164
    6    158
    5    153
    3    138
    2    114
    1     96
```

```
0    68
Name: Age, dtype: int64
```

train_df.head(10)

| | Survived | Pclass | Sex | Age | SibSp | Parch | Fare | Embarked | relatives | not_alone | Deck |
|---|---|---|---|---|---|---|---|---|---|---|---|
| **0** | 0 | 3 | 0 | 2 | 1 | 0 | 7 | 0 | 1 | 0 | 8 |
| **1** | 1 | 1 | 1 | 5 | 1 | 0 | 71 | 1 | 1 | 0 | 3 |
| **2** | 1 | 3 | 1 | 3 | 0 | 0 | 7 | 0 | 0 | 1 | 8 |
| **3** | 1 | 1 | 1 | 5 | 1 | 0 | 53 | 0 | 1 | 0 | 3 |
| **4** | 0 | 3 | 0 | 5 | 0 | 0 | 8 | 0 | 0 | 1 | 8 |
| **5** | 0 | 3 | 0 | 1 | 0 | 0 | 8 | 2 | 0 | 1 | 8 |
| **6** | 0 | 1 | 0 | 6 | 0 | 0 | 51 | 0 | 0 | 1 | 5 |
| **7** | 0 | 3 | 0 | 0 | 3 | 1 | 21 | 0 | 4 | 0 | 8 |
| **8** | 1 | 3 | 1 | 3 | 0 | 2 | 11 | 0 | 2 | 0 | 8 |
| **9** | 1 | 2 | 1 | 1 | 1 | 0 | 30 | 1 | 1 | 0 | 8 |

```python
data = [train_df, test_df]

for dataset in data:
    dataset.loc[ dataset['Fare'] <= 7.91, 'Fare'] = 0
    dataset.loc[(dataset['Fare'] > 7.91) & (dataset['Fare'] <= 14.454), 'Fare'] = 1
    dataset.loc[(dataset['Fare'] > 14.454) & (dataset['Fare'] <= 31), 'Fare']   = 2
    dataset.loc[(dataset['Fare'] > 31) & (dataset['Fare'] <= 99), 'Fare']   = 3
    dataset.loc[(dataset['Fare'] > 99) & (dataset['Fare'] <= 250), 'Fare']   = 4
    dataset.loc[ dataset['Fare'] > 250, 'Fare'] = 5
    dataset['Fare'] = dataset['Fare'].astype(int)


data = [train_df, test_df]
for dataset in data:
    dataset['Age_Class']= dataset['Age']* dataset['Pclass']


for dataset in data:
    dataset['Fare_Per_Person'] = dataset['Fare']/(dataset['relatives']+1)
    dataset['Fare_Per_Person'] = dataset['Fare_Per_Person'].astype(int)
# Let's take a last look at the training set, before we start training the models.
train_df.head(10)
```

| | Survived | Pclass | Sex | Age | SibSp | Parch | Fare | Embarked | relatives | not_alone | Deck |
|---|---|---|---|---|---|---|---|---|---|---|---|
| **0** | 0 | 3 | 0 | 2 | 1 | 0 | 0 | 0 | 1 | 0 | 8 |
| **1** | 1 | 1 | 1 | 5 | 1 | 0 | 3 | 1 | 1 | 0 | 3 |
| **2** | 1 | 3 | 1 | 3 | 0 | 0 | 0 | 0 | 0 | 1 | 8 |
| **3** | 1 | 1 | 1 | 5 | 1 | 0 | 3 | 0 | 1 | 0 | 3 |
| **4** | 0 | 3 | 0 | 5 | 0 | 0 | 1 | 0 | 0 | 1 | 8 |
| **5** | 0 | 3 | 0 | 1 | 0 | 0 | 1 | 2 | 0 | 1 | 8 |
| **6** | 0 | 1 | 0 | 6 | 0 | 0 | 3 | 0 | 0 | 1 | 5 |
| **7** | 0 | 3 | 0 | 0 | 3 | 1 | 2 | 0 | 4 | 0 | 8 |

## BUILDING MACHINE LEARNING MODELS

| **9** | 1 | 2 | 1 | 1 | 1 | 0 | 2 | 1 | 1 | 0 | 8 |

## Random Forest

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


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)
```

## 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)
```

## 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)
```

## Linear Support Vector Machine

```
linear_svc = LinearSVC()
linear_svc.fit(X_train, Y_train)

Y_pred = linear_svc.predict(X_test)

acc_linear_svc = round(linear_svc.score(X_train, Y_train) * 100, 2)
```

```
/usr/local/lib/python3.7/dist-packages/sklearn/svm/_base.py:1208: ConvergenceWarning: Li
  ConvergenceWarning,
```

```
results = pd.DataFrame({
    'Model': ['Support Vector Machines', 'Logistic Regression',
              'Random Forest', 'Naive Bayes'],
    'Score': [acc_linear_svc,  acc_log,
              acc_random_forest, acc_gaussian]})
result_df = results.sort_values(by='Score', ascending=False)
result_df = result_df.set_index('Score')
result_df.head(9)
```

| Score | Model |
|-------|-------|
| 93.04 | Random Forest |
| 81.37 | Logistic Regression |
| 81.26 | Support Vector Machines |
| 78.68 | Naive Bayes |

```
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: [0.75555556 0.82022472 0.75280899 0.85393258 0.86516854 0.86516854
 0.84269663 0.7752809  0.87640449 0.84269663]
Mean: 0.8249937578027465
Standard Deviation: 0.044566518656266645
```

## Feature Importance

```
importances = pd.DataFrame({'feature':X_train.columns,'importance':np.round(random_forest.fea
importances = importances.sort_values('importance',ascending=False).set_index('feature')
importances.head(15)
```

| feature | importance |
|---|---|
| Title | 0.205 |
| Sex | 0.165 |
| Age_Class | 0.093 |
| Deck | 0.085 |
| Age | 0.075 |
| Pclass | 0.074 |
| Fare | 0.065 |
| relatives | 0.061 |
| Embarked | 0.057 |
| Fare_Per_Person | 0.044 |
| SibSp | 0.039 |
| Parch | 0.025 |
| not_alone | 0.012 |

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

```
<matplotlib.axes._subplots.AxesSubplot at 0x7f9c580b2210>
```

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


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

```
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(round(acc_random_forest,2,), "%")
```

        93.04 %


```
# Random Forest
random_forest = RandomForestClassifier(criterion = "gini",
                                       min_samples_leaf = 1,
                                       min_samples_split = 10,
                                       n_estimators=100,
                                       max_features='auto',
                                       oob_score=True,
                                       random_state=1,
                                       n_jobs=-1)

random_forest.fit(X_train, Y_train)
Y_prediction = random_forest.predict(X_test)

random_forest.score(X_train, Y_train)

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

        oob score: 82.49 %


## Evaluation


```
from sklearn.model_selection import cross_val_predict
from sklearn.metrics import confusion_matrix
predictions = cross_val_predict(random_forest, X_train, Y_train, cv=3)
confusion_matrix(Y_train, predictions)
```

        array([[491,  58],

```
        [ 97, 245]])
```

```python
from sklearn.metrics import precision_score, recall_score

print("Precision:", precision_score(Y_train, predictions))
print("Recall:",recall_score(Y_train, predictions))
```

```
    Precision: 0.8085808580858086
    Recall: 0.716374269005848
```

```python
from sklearn.metrics import f1_score
f1_score(Y_train, predictions)
```

```
    0.7596899224806202
```

```python
from sklearn.metrics import precision_recall_curve

# getting the probabilities of our predictions
y_scores = random_forest.predict_proba(X_train)
y_scores = y_scores[:,1]

precision, recall, threshold = precision_recall_curve(Y_train, y_scores)
def plot_precision_and_recall(precision, recall, threshold):
    plt.plot(threshold, precision[:-1], "r-", label="precision", linewidth=5)
    plt.plot(threshold, recall[:-1], "b", label="recall", linewidth=5)
    plt.xlabel("threshold", fontsize=19)
    plt.legend(loc="upper right", fontsize=19)
    plt.ylim([0, 1])

plt.figure(figsize=(14, 7))
plot_precision_and_recall(precision, recall, threshold)
plt.show()
```
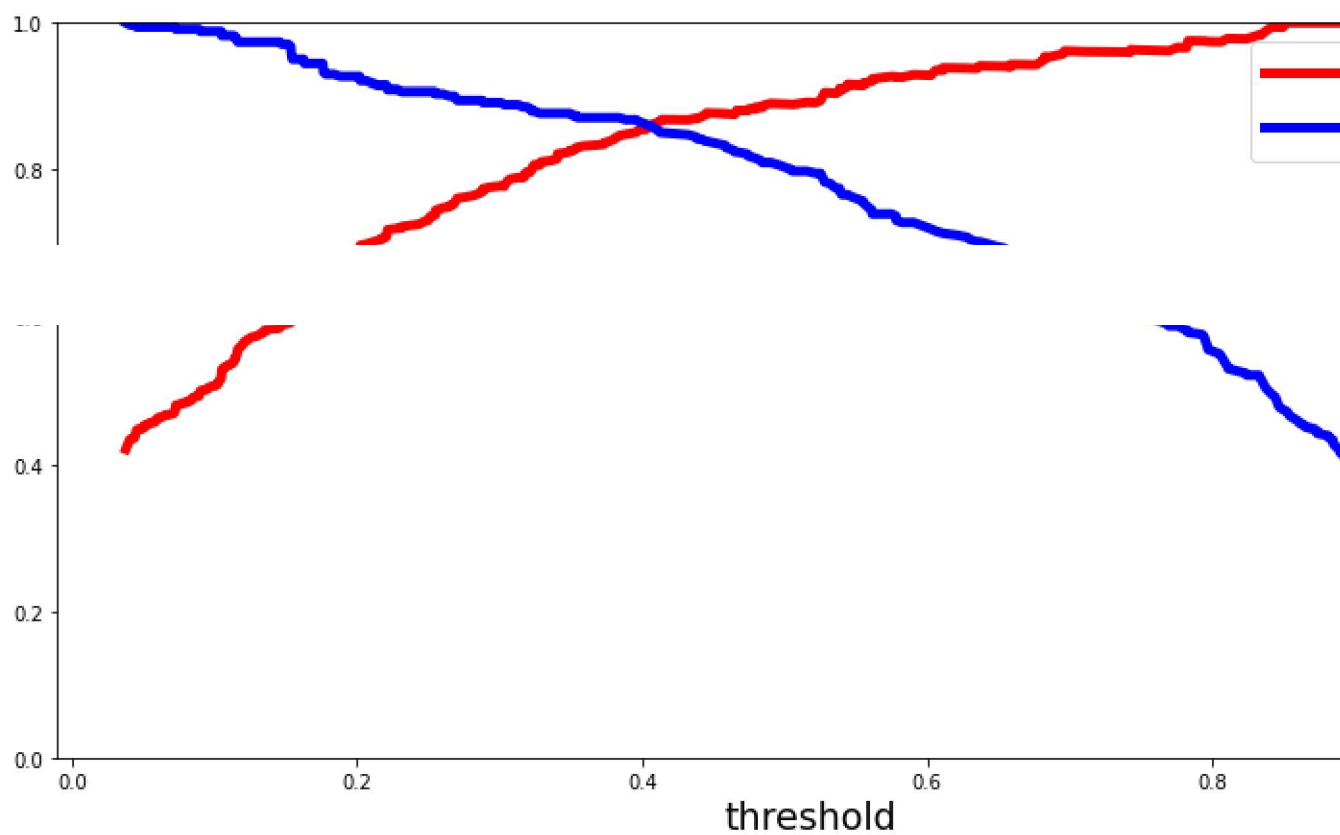
⤷

threshold