	RangeIndex: 418 entries, 0 to 417 Data columns (total 11 columns): # Column Non-Null Count Dtype
ſ	<pre>8 Fare 417 non-null float64 9 Cabin 91 non-null object 10 Embarked 418 non-null object dtypes: float64(2), int64(4), object(5) memory usage: 36.0+ KB # train_df.describe() # Review survived rate using `percentiles=[.61, .62]` knowing our problem description mentio ns 38% survival rate. # Review Parch distribution using `percentiles=[.75, .8]` # SibSp distribution `[.68, .69]`</pre>
In [8]: Out[8]:	# Age and Fare `[.1, .2, .3, .4, .5, .6, .7, .8, .9, .99]` train_df.describe(include=['0']) Name Sex Ticket Cabin Embarked
<pre>In [9]: Out[9]:</pre>	<pre>top Sage, Mr. Frederick male 347082 G6 S freq 1 577 7 4 644 train_df[['Pclass', 'Survived']].groupby(['Pclass'], as_index=False).mean().sort_values(by= 'Survived', ascending=False)</pre> Pclass Survived 1 0.629630
In [10]: Out[10]:	<pre>1 2 0.472826 2 3 0.242363 train_df[["Sex", "Survived"]].groupby(['Sex'], as_index=False).mean().sort_values(by='Survived', ascending=False)</pre> Sex Survived
In [11]: Out[11]:	<pre>0 female 0.742038 1 male 0.188908 train_df[["SibSp", "Survived"]].groupby(['SibSp'], as_index=False).mean().sort_values(by='Survived', ascending=False) SibSp Survived 1</pre>
In [12]: Out[12]:	<pre>5 5 0.000000 6 8 0.000000 train_df[["Parch", "Survived"]].groupby(['Parch'], as_index=False).mean().sort_values(by='Su rvived', ascending=False)</pre>
In [13]:	<pre>g = sns.FacetGrid(train_df, col='Survived') g.map(plt.hist, 'Age', bins=20) <seaborn.axisgrid.facetgrid 0x1f339c85048="" at=""> Survived = 0 Survived = 1</seaborn.axisgrid.facetgrid></pre>
In []:	
	<pre>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(); # grid = sns.FacetGrid(train_df, col='Embarked') grid = sns.FacetGrid(train_df, row='Embarked', size=2.2, aspect=1.6) grid.map(sns.pointplot, 'Pclass', 'Survived', 'Sex', palette='deep') grid.add_legend() # grid = sns.FacetGrid(train_df, col='Embarked', hue='Survived', palette={0: 'k', 1: 'w'})</pre>
	grid = sns.FacetGrid(train_df, row='Embarked', col='Survived', size=2.2, aspect=1.6) grid.map(sns.barplot, 'Sex', 'Fare', alpha=.5, ci=None) grid.add_legend() C:\ProgramData\Anaconda3\lib\site-packages\seaborn\axisgrid.py:243: UserWarning: The `size` p arameter has been renamed to `height`; please update your code. warnings.warn(msg, UserWarning) C:\ProgramData\Anaconda3\lib\site-packages\seaborn\axisgrid.py:728: UserWarning: Using the ba rplot function without specifying `order` is likely to produce an incorrect plot. warnings.warn(warning) <seaborn.axisgrid.facetgrid 0x1f33a9ffc88="" at=""> Embarked = S Survived = 0 Embarked = S Survived = 1 Embarked = C Survived = 1 Embarked = C Survived = 1</seaborn.axisgrid.facetgrid>
	80 - 60 - 20 - Embarked = Q Survived = 0 Embarked = Q Survived = 1 Surviv
In [17]:	<pre>print("Before", train_df.shape, test_df.shape, combine[0].shape, combine[1].shape) train_df = train_df.drop(['Ticket', 'Cabin'], axis=1) test_df = test_df.drop(['Ticket', 'Cabin'], axis=1) combine = [train_df, test_df] "After", train_df.shape, test_df.shape, combine[0].shape, combine[1].shape</pre>
(<pre>Before (891, 12) (418, 11) (891, 12) (418, 11) ('After', (891, 10), (418, 9), (891, 10), (418, 9)) for dataset in combine: dataset['Title'] = dataset.Name.str.extract(' ([A-Za-z]+)\.', expand=False) pd.crosstab(train_df['Title'], train_df['Sex'])</pre>
Out[18]:	Sex female male Title Capt 0 1 Col 0 2 Countess 1 0 Don 0 1 Dr 1 6 Jonkheer 0 1 Lady 1 0 Major 0 2 Master 0 40 Miss 182 0 Mme 1 0 Mr 0 517 Mrs 125 0 Ms 1 0 Rev 0 6 Sir 0 1
In [19]:	<pre>for dataset in combine: dataset['Title'] = dataset['Title'].replace(['Lady', 'Countess','Capt', 'Col',\</pre>
Out[19]:	<pre>train_df[['Title', 'Survived']].groupby(['Title'], as_index=False).mean() Title Survived 0 Master 0.575000 1 Miss 0.702703 2 Mr 0.156673 3 Mrs 0.793651</pre>
In [20]:	<pre>4 Rare 0.347826 title_mapping = {"Mr": 1, "Miss": 2, "Mrs": 3, "Master": 4, "Rare": 5} for dataset in combine: dataset['Title'] = dataset['Title'].map(title_mapping) dataset['Title'] = dataset['Title'].fillna(0) train_df.head()</pre>
Juc[20]:	PassengerId Survived Pclass Name Sex Age SibSp Parch Fare Embarked Title 0 1 0 3 Braund, Mr. Owen Harris male 22.0 1 0 7.2500 S 1 1 2 1 1 Cumings, Mrs. John Bradley (Florence Briggs Th female 38.0 1 0 71.2833 C 3 2 3 1 3 Heikkinen, Miss. Laina female 26.0 0 0 7.9250 S 2 3 4 1 1 Futrelle, Mrs. Jacques Heath (Lily May Peel) female 35.0 1 0 53.1000 S 3
(
In [22]: Out[22]:	<pre>for dataset in combine: dataset['Sex'] = dataset['Sex'].map({'female': 1, 'male': 0}).astype(int) train_df.head() Survived Pclass Sex Age SibSp Parch Fare Embarked Title 0</pre>
In [23]:	<pre>2 1 3 1 26.0 0 0 7.9250 S 2 3 1 1 1 35.0 1 0 53.1000 S 3 4 0 3 0 35.0 0 0 8.0500 S 1 # grid = sns.FacetGrid(train_df, col='Pclass', hue='Gender') grid = sns.FacetGrid(train_df, row='Pclass', col='Sex', size=2.2, aspect=1.6) grid.map(plt.hist, 'Age', alpha=.5, bins=20) grid.add_legend()</pre>
Out[23]:	C:\ProgramData\Anaconda3\lib\site-packages\seaborn\axisgrid.py:243: UserWarning: The `size` p arameter has been renamed to `height`; please update your code. warnings.warn(msg, UserWarning) <seaborn.axisgrid.facetgrid 0x1f33ac8e188="" at=""> Pclass = 1 Sex = 0 Pclass = 2 Sex = 1 40 Pclass = 2 Sex = 0 Pclass = 2 Sex = 1</seaborn.axisgrid.facetgrid>
	Pclass = 3 Sex = 0
(
	<pre>for i in range(0, 2): for j in range(0, 3): guess_df = dataset[(dataset['Sex'] == i) & \</pre>
	<pre># Convert random age float to nearest .5 age guess_ages[i,j] = int(age_guess/0.5 + 0.5) * 0.5 for i in range(0, 2): for j in range(0, 3): dataset.loc[(dataset.Age.isnull()) & (dataset.Sex == i) & (dataset.Pclass == j+ 1),\</pre>
Out[25]:	<pre>dataset['Age'] = dataset['Age'].astype(int) train_df.head() Survived Pclass Sex Age SibSp Parch Fare Embarked Title 0 0 3 0 22 1 0 7.2500</pre>
In [26]: Out[26]:	<pre>3 1 1 1 35 1 0 53.1000 S 3 4 0 3 0 35 0 0 8.0500 S 1 train_df['AgeBand'] = pd.cut(train_df['Age'], 5) train_df[['AgeBand', 'Survived']].groupby(['AgeBand'], as_index=False).mean().sort_values(by = 'AgeBand', ascending=True) AgeBand Survived 0 (-0.08, 16.0] 0.550000 1 (16.0, 32.0] 0.337374 2 (32.0, 48.0] 0.412037 3 (48.0, 64.0] 0.434783</pre>
In [27]:	<pre>for dataset in combine: dataset.loc[dataset['Age'] <= 16, 'Age'] = 0 dataset.loc[(dataset['Age'] > 16) & (dataset['Age'] <= 32), 'Age'] = 1 dataset.loc[(dataset['Age'] > 32) & (dataset['Age'] <= 48), 'Age'] = 2 dataset.loc[(dataset['Age'] > 48) & (dataset['Age'] <= 64), 'Age'] = 3 dataset.loc[dataset['Age'] > 64, 'Age'] train_df.head()</pre>
Out[27]:	Survived Pclass Sex Age SibSp Parch Fare Embarked Title AgeBand 0 0 3 0 1 1 0 7.2500 S 1 (16.0, 32.0] 1 1 1 2 1 0 71.2833 C 3 (32.0, 48.0] 2 1 3 1 1 0 7.9250 S 2 (16.0, 32.0] 3 1 1 2 1 0 53.1000 S 3 (32.0, 48.0]
In [28]: Out[28]:	<pre>4 0 3 0 2 0 0 8.0500 S 1 (32.0, 48.0] train_df = train_df.drop(['AgeBand'], axis=1) combine = [train_df, test_df] train_df.head() Survived Pclass Sex Age SibSp Parch</pre>
In [32]:	<pre>1 1 1 1 2 1 0 71.2833</pre>
Out[32]:	train_df[['FamilySize', 'Survived']].groupby(['FamilySize'], as_index=False).mean().sort_values(by='Survived', ascending=False) FamilySize Survived 3
In [33]: Out[33]:	<pre>for dataset in combine: dataset['IsAlone'] = 0 dataset.loc[dataset['FamilySize'] == 1, 'IsAlone'] = 1 train_df[['IsAlone', 'Survived']].groupby(['IsAlone'], as_index=False).mean()</pre>
In [34]:	<pre>train_df = train_df.drop(['Parch', 'SibSp', 'FamilySize'], axis=1) test_df = test_df.drop(['Parch', 'SibSp', 'FamilySize'], axis=1) combine = [train_df, test_df] train_df.head()</pre>
Out[34]:	Survived Pclass Sex Age Fare Embarked Title IsAlone 0 0 3 0 1 7.2500 S 1 0 1 1 1 2 71.2833 C 3 0 2 1 3 1 1 7.9250 S 2 1 3 1 1 2 53.1000 S 3 0 4 0 3 0 2 8.0500 S 1 1
In [35]: Out[35]:	<pre>for dataset in combine: dataset['Age*Class'] = dataset.Age * dataset.Pclass train_df.loc[:, ['Age*Class', 'Age', 'Pclass']].head(10) Age*Class Age Pclass 0</pre>
	1 2 2 1 2 3 1 3 3 2 2 1 4 6 2 3 5 3 1 3 6 3 3 1 7 0 0 3 8 3 1 3
Out[36]:	<pre>freq_port = train_df.Embarked.dropna().mode()[0] freq_port 'S' for dataset in combine:</pre>
In [37]: Out[37]:	<pre>for dataset in combine: dataset['Embarked'] = dataset['Embarked'].fillna(freq_port) train_df[['Embarked', 'Survived']].groupby(['Embarked'], as_index=False).mean().sort_values(by='Survived', ascending=False) Embarked Survived 0</pre>
In [38]: Out[38]:	1 Q 0.389610 2 S 0.339009 for dataset in combine: dataset['Embarked'] = dataset['Embarked'].map({'S': 0, 'C': 1, 'Q': 2}).astype(int) train_df.head() Survived Pclass Sex Age Fare Embarked Title IsAlone Age*Class 0 0 3 0 1 7.2500
In [39]: Out[39]:	4 0 3 0 2 8.0500 0 1 1 6 test_df['Fare'].fillna(test_df['Fare'].dropna().median(), inplace=True) test_df.head() PassengerId Pclass Sex Age Fare Embarked Title IsAlone Age*Class 0 892 3 0 2 7.8292 2 1 1 6 1 893 3 1 2 7.0000 0 3 0 6
In [40]: Out[40]:	<pre>1 893 3 1 2 7.0000 0 3 0 6 2 894 2 0 3 9.6875 2 1 1 6 3 895 3 0 1 8.6625 0 1 1 3 4 896 3 1 1 12.2875 0 3 0 3 train_df['FareBand'] = pd.qcut(train_df['Fare'], 4) train_df[['FareBand', 'Survived']].groupby(['FareBand'], as_index=False).mean().sort_values(by='FareBand', ascending=True)</pre> FareBand Survived
	<pre>0 (-0.001, 7.91] 0.197309 1 (7.91, 14.454] 0.303571 2 (14.454, 31.0] 0.454955 3 (31.0, 512.329] 0.581081 for dataset in combine: dataset.loc[dataset['Fare'] <= 7.91, 'Fare'] = 0</pre>
Out[41]:	Survived Pclass Sex Age Fare Embarked Title IsAlone Age*Class 0 0 3 0 1 0 0 1 0 3 1 1 1 1 2 3 1 3 0 2 2 1 3 1 1 1 0 2 1 3 3 1 1 1 2 3 0 2 4 0 3 0 2 1 1 6 5 0 3 0 1 1 3 1 1 3 6 0 1 0 3 0 1 1 3 7 0 3 0 2 0 4 0 0 8 1 3 1 1 1 3 0 3 0 3 9 </td
In [42]: Out[42]:	9 1 2 1 0 2 1 3 0 0 test_df.head(10) PassengerId Pclass Sex Age Fare Embarked Title IsAlone Age*Class 0 892 3 0 2 0 2 1 1 6 1 893 3 1 2 0 0 3 1 2 6 2 894 2 0 3 1 2 1 1 6 3 895 3 0 1 1 0 0 1 1 0 3 4 896 3 1 1 1 1 0 3 0 3 5 897 3 0 0 1 1 0 0 1 1 0
	<pre>6 898 3 1 1 0</pre>
In [43]:	
Out[43]:	((891, 8), (891,), (418, 8)) 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) acc_log
Out[43]:	Logistic Regression logreg = LogisticRegression() logreg.fit(X_train, Y_train) Y_pred = logreg.predict(X_test) acc_log =

In []: # Gaussian Naive Bayes

acc_gaussian

acc_perceptron

In []: # Stochastic Gradient Descent

acc_decision_tree

acc_random_forest

In [54]: models = pd.DataFrame({

sgd = SGDClassifier()
sgd.fit(X_train, Y_train)

Y_pred = sgd.predict(X_test)

decision_tree = DecisionTreeClassifier() decision_tree.fit(X_train, Y_train)
Y_pred = decision_tree.predict(X_test)

In []: # Perceptron

In []: # Linear SVC

acc_sgd

In []: # Decision Tree

In []: # Random Forest

Out[54]:

Out[61]:

3

8

1

2

5

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)

perceptron = Perceptron()
perceptron.fit(X_train, Y_train)
Y_pred = perceptron.predict(X_test)
acc_perceptron = round(perceptron.score(X_train, Y_train) * 100, 2)

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)
acc_linear_svc

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

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

'Score': [acc_svc, acc_knn, acc_log, acc_random_forest, acc_gaussian, acc_perceptron, acc_sgd, acc_linear_svc, acc_decision_tree]})

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

random_forest = RandomForestClassifier(n_estimators=100)
random_forest.fit(X_train, Y_train)
Y_pred = random_forest.predict(X_test)

'Decision Tree'],

models.sort_values(by='Score', ascending=False)

Model Score

KNN 84.74

Random Forest 86.76

Decision Tree 86.76

Linear SVC 79.01

Perceptron 78.34

Naive Bayes 72.28

"Survived": Y_pred

0

0 0

0

1

0 1

"PassengerId": test_df["PassengerId"],

In [61]: pd.set_option("display.max_columns", 100, "display.width", 1000)

Logistic Regression 80.36

Support Vector Machines 78.23

6 Stochastic Gradient Decent 54.10

Passengerld Survived

892 893

894

895 896

1305 1306

1307

1308

1309

In []: submission = pd.DataFrame({

}) # submission

submission

3

413

415

416

417

In [1]: # ------##<<<<-----#

import pandas as pd import numpy as np import random as rnd

import seaborn as sns

%matplotlib inline

import matplotlib.pyplot as plt

In []: train_df = pd.read_csv('train.csv')
 test_df = pd.read_csv('test.csv')

In [65]: print(train_df.columns.values)

preview the Data

train_df.head()

In [63]: train_df.tail()

886

887

888 889

890

In [6]: train_df.info()

print('%'*40)
test_df.info()

Column

Survived

Pclass

Name

Sex

Age

0

2

3

4

Out[63]:

combine = [train_df, test_df]

'Ticket' 'Fare' 'Cabin' 'Embarked']

from sklearn.svm import SVC, LinearSVC

from sklearn.naive_bayes import GaussianNB

from sklearn.linear_model import Perceptron

from sklearn.linear_model import SGDClassifier

from sklearn.tree import DecisionTreeClassifier

['PassengerId' 'Survived' 'Pclass' 'Name' 'Sex' 'Age' 'SibSp' 'Parch'

Survived Pclass Sex Age Fare Embarked Title IsAlone Age*Class

0 2

1 1

1

2 0 1 1

1 1 1 2

0

<class 'pandas.core.frame.DataFrame'> RangeIndex: 891 entries, 0 to 890 Data columns (total 12 columns):

PassengerId 891 non-null

0 1 2

1

Non-Null Count Dtype

891 non-null

891 non-null

891 non-null

891 non-null

714 non-null

int64

int64

int64

object

object

float64

2

1

1

from sklearn.linear_model import LogisticRegression

from sklearn.ensemble import RandomForestClassifier

from sklearn.neighbors import KNeighborsClassifier

DESCRIPTION : THE PROGRAM PREDICTES THAT IF A PASSENGER WILL SURVIVE ON THE TITANIC OR NO