In [ ]:	<pre>import numpy as np import matplotlib.pyplot as plt import statsmodels.api as sm</pre>
	<pre>import sklearn as skl import random from scipy.stats import iqr import scipy.stats as stats import seaborn as sns</pre> <pre>Importing the Dataset</pre>
In []:	<pre>print(data.head()) data.dtypes  Survived Pclass \ PassengerId 1</pre>
	Name Sex Age \ PassengerId  Braund, Mr. Owen Harris male 22.0  Cumings, Mrs. John Bradley (Florence Briggs Th female 38.0  Heikkinen, Miss. Laina female 26.0  Futrelle, Mrs. Jacques Heath (Lily May Peel) female 35.0  Allen, Mr. William Henry male 35.0  SibSp Parch Ticket Fare Cabin Embarked
Out[]:	PassengerId  1
	Age float64 SibSp int64 Parch int64 Ticket object Fare float64 Cabin object Embarked object dtype: object
In [ ]:	Creating Filters  # Cutting down the dataset
	<pre>mod_data = data[['Survived', 'Pclass', 'Age', 'Sex', 'SibSp', 'Parch', 'Fare' ]]  # Classifiers male_df = mod_data[mod_data.Sex == 'male'].drop(['Sex'], axis = 1) female_df = mod_data[mod_data.Sex == 'female'].drop(['Sex'], axis = 1)  # Explanatory Variables in_df_m = male_df.iloc[:,1:7] in_df_f = female_df.iloc[:,1:7]</pre>
In [ ]: Out[ ]:	<pre>#Exploratory Variable dep_m = male_df.Survived dep_f = female_df.Survived  male_df</pre>
	1       0       3       22.0       1       0       7.2500         5       0       3       35.0       0       0       8.0500         6       0       3       0.0       0       0       8.4583         7       0       1       54.0       0       0       51.8625         8       0       3       2.0       3       1       21.0750
	884 0 2 28.0 0 0 10.5000  885 0 3 25.0 0 0 7.0500  887 0 2 27.0 0 0 13.0000  890 1 1 26.0 0 0 30.0000  891 0 3 32.0 0 0 7.7500
	Exploratory Data Analysis  Male Population
In [ ]:	PClass  As seen from the data, majority of the male passengers onboard the Titanic are shown to be in Passenger Class 3 having over 60%. Passenger Class 1 coming in second with just over 20% and Passenger Class 2 having the least with just shy of 20%.  print(male_df.Pclass.value_counts()) print('========') index_pclass = list(map(str,np.array(male_df.Pclass.value_counts().reset_index()).T[0])) count_pclass = np.array(male_df.Pclass.value_counts().reset_index()).T[1]
	<pre>norm_pclass = count_pclass/count_pclass.sum() # Histogram of the PClasses fig = plt.figure() ax = fig.add_axes([0,0,1,1]) ax.bar(index_pclass, norm_pclass) plt.show  3</pre>
Out[ ]:	======================================
	0.4 - 0.3 - 0.2 - 0.1 -
In [ ]:	Applying the Freedman-Diacones Rule  def bin_range(df):     vals = df.values     bin_width = 2 * (iqr(vals) / (len(vals) ** (1./3)))     num_bins = int((np.amax(vals) + bin_width) / bin_width)
	<pre>return bin_width, num_bins  def histo_viz(df):     vals = df.values     bin_width, num_bins = bin_range(df)     hist = df.hist(bins=num_bins)     return hist</pre> Ages
In [ ]:	<pre># Bar graph of the ages  width_age, numbin_age = bin_range(male_df.Age) print('Bin width: {}'.format(round(width_age,4))) print('Number of bins: {}'.format(numbin_age)) print('====================================</pre>
	<pre>Bin width: 7.4473 Number of bins: 11 ===================================</pre>
	100 - 80 - 60 - 40 - 20 - 30 - 40 - 50 - 60 - 70 - 80
In [ ]:	<pre>SibSp print(male_df.SibSp.value_counts()) print('=========') index_sibsp = list(map(str,np.array(male_df.SibSp.value_counts().reset_index()).T[0])) count_sibsp =np.array(male_df.SibSp.value_counts().reset_index()).T[1]  # Histogram of the SibSp</pre>
	<pre>fig = plt.figure() ax = fig.add_axes([0,0,1,1]) ax.bar(index_sibsp, count_sibsp.sum()) plt.show  0     434 1     103 2     15 4     12 3     5 5     4</pre>
Out[ ]:	Name: SibSp, dtype: int64 ====================================
	0.5 - 0.4 - 0.3 - 0.2 - 0.1 - 0.1 - 0.1 - 0.1 - 0.1 - 0.1 - 0.2 - 0.1 - 0
In [ ]:	0.0 0 1 2 4 3 5 8  Parch
	<pre>count_parch =np.array(male_df.Parch.value_counts().reset_index()).T[1]  # Histogram of the SibSp fig = plt.figure() ax = fig.add_axes([0,0,1,1]) ax.bar(index_parch, count_parch/count_parch.sum()) plt.show</pre> 0 484 1 58 2 31
Out[ ]:	4 2 5 1 3 1 Name: Parch, dtype: int64 ====================================
	0.7 - 0.6 - 0.5 - 0.4 - 0.3 -
	Visualization of the Data
In [ ]:	<pre>pd.plotting.parallel_coordinates(male_df, 'Survived', color=('#FFA500', '#0000FF'))</pre>
Out[]:	<pre><axessubplot:></axessubplot:></pre>
In [ ]:	# norm_m_df = (male_df.iloc[:,2:6] -male_df.iloc[:,2:6].mean())/male_df.iloc[:,2:6].std() # norm_m_df = pd.merge(norm_m_df, dep_m, right_index= True, left_index= True) # norm_m_df.head()
Out[ ]:	PassengerId           1         0.2750         0.125         0.0         0.014151         0
	5 0.4375 0.000 0.0 0.0 10.015713 0 6 0.0000 0.00 0.00 0.0 0.016510 0 7 0.6750 0.000 0.0 0.0 1.01229 0 8 0.0250 0.375 0.2 0.041136 0  Normalizing the data now shows that most points can be distinguished from one another throughout the 4 variables used. Now we move on to checking the correlation of the data using a heatmap. Thus, we can conclude that all 4 variables can be applicable in distinguishing the data.
In [ ]: Out[ ]:	
	0.4 0.2 0.0 Age SibSp Parch Fare
In [ ]:	<pre>Heat Map (Correlation Coefficient)  temp = norm_m_df.drop(['Survived'], axis=1)     corr = temp.corr()     print(corr)     adj_norm_m_df = corr[((corr &gt; 0.3) &amp; (corr != 1))   ((corr &lt; -0.3) &amp; (corr != 1))]     print('============')     print(adj_norm_m_df.count())     print('============')     print('Total number of correlations that fit into the desire correlation conditions: {}'. format(sum(adj_norm_m_df.count())))</pre>
	Age SibSp Parch Fare  Age 1.000000 -0.173230 -0.079041 0.106134  SibSp -0.173230 1.000000 0.524849 0.181804  Parch -0.079041 0.524849 1.000000 0.312197  Fare 0.106134 0.181804 0.312197 1.000000  ==============================
In [ ]:	heatmap = sns.heatmap(corr, vmin=-1, vmax=1, annot=True)  -100
	Exp.       - 0.17       1       - 0.079       0.11       - 0.75         - 0.50       - 0.50       - 0.25         - 0.00       - 0.079       0.52       1       0.31       - 0.25
	eg - 0.11 0.18 0.31 10.751.001.00
In [ ]: In [ ]:	<pre>from sklearn.linear_model import LogisticRegression from sklearn import metrics from statsmodels.stats.anova import anova_lm # One-way ANOVA</pre>
Out[ ]:	Age         SibSp         Parch         Fare         Survived           PassengerId         0.2750         0.125         0.0         0.014151         0           5         0.4375         0.000         0.0         0.015713         0           6         0.0000         0.00         0.016510         0           7         0.6750         0.000         0.0         0.101229         0
	8       0.0250       0.375       0.2       0.041136       0                 884       0.3500       0.000       0.0       0.020495       0         885       0.3125       0.000       0.0       0.013761       0         887       0.3375       0.000       0.0       0.025374       0         890       0.3250       0.000       0.0       0.058556       1         891       0.4000       0.00       0.0       0.015127       0
In [ ]:	<pre>577 rows x 5 columns  logit_model=sm.Logit(dep_m,norm_m_df.iloc[:,0:4]) result=logit_model.fit() print(result.summary2())  Optimization terminated successfully.</pre>
	Titerations 6  Results: Logit  Dependent Variable: Survived  AIC: 619.6782  Date: 2022-06-01 18:24 BIC: 637.1096  No. Observations: 577 Log-Likelihood: -305.84  Df Model: 3 LL-Null: -279.63  Df Residuals: 573 LLR p-value: 1.0000  No. Iterations: 6.0000
	Coef. Std.Err. z $P> z $ [0.025 0.975]  Age $-3.5075$ 0.3524 $-9.9520$ 0.0000 $-4.1983$ $-2.8167$ SibSp $-4.0503$ 1.1152 $-3.6318$ 0.0003 $-6.2361$ $-1.8645$ Parch 1.3693 0.9465 1.4467 0.1480 $-0.4858$ 3.2244  Fare 2.0469 1.1371 1.8002 0.0718 $-0.1817$ 4.2756
In [ ]:	result=logit_model.fit() print(result.summary2())  Optimization terminated successfully.
	Model: Logit Pseudo R-squared: -0.097 Dependent Variable: Survived AIC: 619.6619 Date: 2022-06-01 18:25 BIC: 632.7354 No. Observations: 577 Log-Likelihood: -306.83 Df Model: 2 LL-Null: -279.63 Df Residuals: 574 LLR p-value: 1.0000 Converged: 1.0000 Scale: 1.0000 No. Iterations: 6.0000  Coef. Std.Err. z P> z  [0.025 0.975]
In [ ]:	Age
In [ ]: Out[ ]:	Age         SibSp         Parch         Fare         Survived           PassengerId         1         0.2750         0.125         0.0         0.014151         0           5         0.4375         0.000         0.0         0.015713         0
	6 0.0000 0.000 0.00 0.0 0.016510 0 7 0.6750 0.000 0.0 0.101229 0 8 0.0250 0.375 0.2 0.041136 0 884 0.3500 0.000 0.0 0.020495 0 885 0.3125 0.000 0.0 0.0 0.013761 0 887 0.3375 0.000 0.0 0.0025374 0
In [ ]:	890 0.3250 0.000 0.0 0.058556 1  891 0.4000 0.000 0.0 0.015127 0  577 rows × 5 columns  X_train = norm_m_df.iloc[:,[0,1,3]] y_train = dep_m
	<pre>X_test = pd.read_csv('/Users/raph/Desktop/titanic/test.csv').set_index('PassengerId') X_test.drop(['Pclass','Name','Sex','Cabin','Embarked','Ticket','Parch'], axis = 1, inplace = True)  X_test.Age = X_test.Age.replace(np.nan, 0) X_test.Fare = X_test.Fare.replace(np.nan, 0) X_test.dropna(inplace = True)  y_test = pd.read_csv('/Users/raph/Desktop/titanic/gender_submission.csv').set_index('PassengerId')</pre>
Out[]: In []:	<pre>logreg = LogisticRegression() logreg.fit(X_train, y_train)  v LogisticRegression LogisticRegression()  X_test = X_test/X_test.max() X_test</pre>
Out[]:	Age         SibSp         Fare           PassengerId         892         0.453947         0.000         0.015282           893         0.618421         0.125         0.013663           894         0.815789         0.000         0.018909           895         0.355263         0.000         0.016908
	896       0.289474       0.125       0.023984               1305       0.000000       0.000       0.015713         1306       0.513158       0.000       0.212559         1307       0.506579       0.00       0.014151         1308       0.000000       0.000       0.015713         1309       0.000000       0.125       0.043640
In [ ]: In [ ]:	418 rows × 3 columns  y_pred = logreg.predict(X_test) print('Accuracy of logistic regression classifier on test set: {:.2f}'.format(logreg.score(X_test, y_test)))  Accuracy of logistic regression classifier on test set: 0.64  from sklearn.metrics import confusion_matrix
In [ ]:	<pre>confusion_matrix = confusion_matrix(y_test, y_pred) print(confusion_matrix)  [[266   0]    [151   1]]  from sklearn.metrics import classification_report print(classification_report(y_test, y_pred))</pre>
	0 0.64 1.00 0.78 266 1 1.00 0.01 0.01 152 accuracy 0.64 418 macro avg 0.82 0.50 0.40 418 weighted avg 0.77 0.64 0.50 418