•	Collecting Data  deading in the titanic_dataset.csv file into a pandas dataframe.  A = pd.read_csv("C:/Users/aksha/OneDrive/Desktop/titanic_dataset.csv")  A.head()  Passengerld Survived Pclass Name Gender Age SibSp Parch Ticket Fare Cabin Embarked
3	1 0 3 Braund, Mr. Owen Harris male 22.0 1 0 A/5 21171 7.2500 NaN S 1 2 1 1 Cumings, Mrs. John Bradley (Florence Briggs Th female 38.0 1 0 PC 17599 71.2833 C85 C 2 3 1 3 Heikkinen, Miss. Laina female 26.0 0 5 STON/O2. 3101282 7.9250 NaN S 3 4 1 1 Futrelle, Mrs. Jacques Heath (Lily May Peel) female 35.0 1 0 113803 53.1000 C123 S 4 5 0 3 Allen, Mr. William Henry male 35.0 0 0 373450 8.0500 NaN S
<b>E</b>	Exploratory Data Analysis  Analyzing Data Creating different plots to check the relationship between variables. Simply explore the data by using various columns  sb. set_style("whitegrid") sb. countplot (x = "Survived", data=A)
	<pre><axessubplot:xlabel='survived', ylabel="count">  500 400 400 200</axessubplot:xlabel='survived',></pre>
	#To check how many were men and how many were women sb.set_style("whitegrid") sb.countplot(x = "Survived", hue = "Gender", data=A)
	<pre> AxesSubplot:xlabel='Survived', ylabel='count'&gt;  400  300  100  400  400  400  400  400</pre>
	Survived  *To check the class of passengers who did not survive sb.set_style("whitegrid") sb.countplot(x = "Survived", hue="Pclass", data=A)  *AxesSubplot:xlabel='Survived', ylabel='count'>
	Felas   Felas
	analysis: Number of people who did not survive is more from Class 3.  #To check SibSp sb.set_style("whitegrid") sb.countplot(x = "SibSp", data=A)  *AxesSubplot:xlabel='SibSp', ylabel='count'>
	500
	<pre>#to check the age of passengers travelling in titanic sb.set_style("whitegrid") A["Age"].plot.hist()  </pre> <pre></pre> <pre> </pre> <pre> <pre></pre></pre>
Α	Analysis: People of age ranging 20-30 is more  #To check the average fare of the people who have bought the ticket
<	sb.set_style("whitegrid") A["Fare"].plot.hist(bins=20 , figsize = (10,5)) <pre></pre>
	sb.set_style("whitegrid") sb.boxplot(x = "Pclass", y = "Age",data=A)
	<pre> AxesSubplot:xlabel='Pclass', ylabel='Age'&gt;  80  70  80  40  30  20 </pre>
< F	nalysis: Passengers in class 1 & 2 are older than in class 3.  A. info() <class 'pandas.core.frame.dataframe'=""> RangeIndex: 891 entries, 0 to 890 Data columns (total 12 columns):  # Column Non-Null Count Dtype</class>
	PassengerId 891 non-null int64  1 Survived 891 non-null int64  2 Pclass 891 non-null object  4 Gender 891 non-null object  5 Age 714 non-null float64  6 SibSp 891 non-null int64  7 Parch 891 non-null object  8 Ticket 891 non-null object  8 Ticket 891 non-null object  1 Embarked 889 non-null object  11 Embarked 889 non-null object  dtypes: float64(2), int64(5), object(5)  memory usage: 83.7+ KB
<b>N</b>	Missing Data Treatment  Cleaning the data by removing NaN values and unnecessary columns from the dataset.  A.isnull() #True If it is null #False If it is not null  Passengerid Survived Pclass Name Gender Age SibSp Parch Ticket Fare Cabin Embarked  The state of the state
8	1FalseFalseFalseFalseFalseFalseFalseFalseFalseFalseFalseFalse2False
8: 8:	False
F S H C C H C C	Gender 0 Age 177 SibSp 0 Parch 0 Ticket 0 Fare 0 Cabin 687 Embarked 2 dtype: int64  #To analyze them visually, use seaborn to create a simple heatmap to see where is the data missing. sb.heatmap(A.isnull(), yticklabels = False , cmap = "viridis") <axessubplot:></axessubplot:>
	-10 -08 -06 -04 -02
R In	toughly 20 percent of the Age data is missing. Looking at the cabin column,it looks like too much data is missing.  In order to remove these columns- Either drop them or replace them and put dummy values to it.  A. head ()
3	Passengerid         Survived         Pelas         Name         Gender         Age         SibSp         Parch         Ticket         Fase         Cabin         Embarked           0         1         0         3         Braund, Mr. Owen Harris         male         22.0         1         0         A/S 2171         7.250         NaN         S           1         2         1         1         Cumings, Mrs. John Bradley (Florence Briggs Th         female         8.0         1         0         PC 1759         7.283         C85         C           2         3         1         3         Heikkinen, Miss. Laina         female         8.0         1         0         STON/O2. 3101282         7.9250         NaN         S           3         4         1         1         Futurelle, Mrs. Jacques Heath (Lily May Peel)         female         35.0         1         0         113803         53.1000         NaN         S           4         5         0         3         Allen, Mr. William Henry         male         35.0         0         373450         8.0500         NaN         S
3	Passenger Id         Survive Id         Pclass         Pclass         Name         Gender Id         Ticket         Fare         Embarket           0         1         0         3         Braund, Mr. Owen Harris         male         2.0         1         0         A/5 2117         7.2500         S           1         2         1         1         Cumings, Mrs. John Bradley (Florence Briggs Th
	<pre>#to check if the values have been removed or not. sb.heatmap(A.isnull(), yticklabels=False , cbar = False)  <axessubplot:></axessubplot:></pre>
	Passengerid Survived Pulass Name Gender Age SBSp Parch Ticket Fare Fare Embarked
H S S H S S H S S H S S H S S H S S H S S H S S H S S H S S H S S S H S S S H S S S H S	A.isnull().sum()  PassengerId 0  Survived 0  Pclass 0  Name 0  Gender 0  Age 0  SibSp 0  Parch 0  Ticket 0
	Coverting Categorical Features  A.head (2)  Passengerid Vurvived Pclass Name Gender Age SibSp Parch Ticket Fare Embarked  1 0 3 Braund, Mr. Owen Harris male 2.0 1 0 NAS Dr. Name Gender Name (20 1 0 NAS Dr. Name (20 1 0 NAS Dr. Name Name (20 1 0 NAS Dr. Name Name Name Name Name Name Name Name
ТІ	The data has lot of string values. Have to convert categorical featues to dummy variables using pandas. Otherwise the machine learning algorithm won't be able to directly take in those features as input.  G = pd.get_dummies(A["Gender"], drop_first=True)  male  0 1 1 0
8	2 0 3 0 4 1 885 0 886 1 887 0 889 1
7	### 10
8	2 0 1 3 0 1 4 0 1 5 885 1 0 886 0 1 887 0 1 889 0 0
7	### ### ##############################
8	1 0 0 2 0 1 3 0 0 4 0 1 885 0 1 886 1 0
7	889 0 0 890 0 1 12 rows × 2 columns  #concatenate all the new rows into the dataset A = pd.concat([A,G,E,Pcl],axis=1) A.head()  Passengerld Survived Pclass  Name Gender Age SibSp Parch  Ticket Fare Embarked male Q S 2 3
3	1 0 3 Braund, Mr. Owen Harris male 22.0 1 0 A/5 21171 7.2500 S 1 0 1 0 1 0 1 0 1 0 1 0 1 0 1 0 1 0
1 2 3	Survivet         Age         SibSp         Parch         Fare         male         Q         S         2         S         S         Parch         Fare         male         Q         S         D         S         D         S         D
S <sub> </sub>	Building a Logistic Regression Model  plitting the data into training set and testing set.  Train Test Split   X = A.drop("Survived", axis=1) Y = A[["Survived"]]
	<pre>from sklearn.model_selection import train_test_split xtrain,xtest,ytrain,ytest = train_test_split(X,Y,test_size=0.2,random_state=31)  from warnings import filterwarnings filterwarnings("ignore")  from sklearn.linear_model import LogisticRegression lr = LogisticRegression() lrmodel = lr.fit(xtrain,ytrain)  pred = lrmodel.predict(xtest)</pre>
	pred  array([0, 1, 1, 0, 1, 0, 0, 1, 1, 0, 0, 0, 1, 0, 1, 1, 1, 0, 1, 1, 0, 0, 1, 0, 0, 1, 0, 0, 1, 0, 1, 1, 0, 1, 1, 1, 0, 0, 0, 0, 0, 0, 1, 0, 0, 1, 1, 0, 0, 1, 0, 1, 1, 1, 0, 0, 0, 0, 0, 0, 0, 0, 1, 1, 0, 0, 0, 0, 1, 1, 0, 0, 0, 0, 0, 1, 1, 0, 0, 0, 0, 1, 1, 0, 0, 0, 0, 1, 1, 0, 0, 0, 0, 1, 1, 0, 0, 0, 0, 1, 1, 0, 0, 0, 0, 1, 1, 0, 0, 1, 1, 0, 0, 1, 1, 0, 0, 1, 1, 0, 0, 1, 1, 0, 0, 1, 1, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 1, 1, 0, 1, 1, 0, 0, 0, 0, 0, 0, 0, 0, 0, 1, 1, 0, 1, 1, 0, 0, 0, 0, 0, 0, 0, 0, 1, 1, 0, 1, 1, 0, 0, 0, 0, 0, 0, 0, 0, 1, 1, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0,
	acro avg 0.79 0.78 0.78 143\nweighted avg 0.79 0.79 143\n'
Ca	#to carculate accuracy using confusion matrix
<b>C</b> a	<pre>#to calculate accuracy using confusion matrix from sklearn.metrics import confusion_matrix confusion_matrix(ytest,pred)  array([[71, 13],</pre>