[39]: df.s (891 (891 (891 (891 (891 (891 (891 (891	Passengerid Survived Polas Manual English Manua
[39]: df.s [39]: we come to see the composition of	The control
[15]: df.r [15]: Pass S S	Transport 1 1 1 1 1 1 1 1 1
[13]: df.i <cla #="" 0="" 1="" 10="" 11="" 2="" 3="" 4="" 44]:="" 46]:="" 5="" 6="" 7="" 8="" 886="" 887="" 888="" 889="" 890="" 891="" 9="" [41]:="" [42]:="" [43]:="" [44]:="" [45]:="" [46]:="" age="" as="" cabi="" data="" df.i="" df1="" df1[="" dtyp="" emba="" funct="" h="" memo="" name="" parc="" pcla="" prir="" prir<="" rang="" rd="" sex="" sibs="" surv="" td="" w="" we=""><td>150 1 1 1 1 1 1 1 1 1 </td></cla>	150 1 1 1 1 1 1 1 1 1
[41]: Pass Survers Perlaments Sex Age Sibs Parce Cabi Embaddtyp We had sex Age	Passenger
[43]: df1 [43]: 0 1 2 3 4 886 887 888 889 890 891 rd We h [44]: df1 [45]: df1 [45]: 0 1 2 3 4 886 887 888 889 890 891 rd We h	Passengerid Survived Polass Name Sex Age SitSsp Parch Ticket Fare Embarked
We h [44]: df1[[45]: df1 [45]: 0 1 2 3 4 886 887 888 889 890 891 rd We h [46]: prin 0 Name As w funct	rave dropped the Cabin column to handle the missing value as we have 687 missing value out of 891 'Age'].fillna(df1['Age'].mean(), inplace=True) #replacing the missing value in "Age" column with mean value PassengerId Survived Pclass Name Sex Age SibSp Parch Ticket Fare Embarked 1 0 3 Braund, Mr. Owen Harris male 22.000000 1 0 A/5 21171 7.2500 S 2 1 1 Cumings, Mrs. John Bradley (Florence Briggs Th female 38.000000 1 0 PC 17599 71.2833 C 3 1 3 Heikkinen, Miss. Laina female 26.000000 0 0 STON/O2. 3101282 7.9250 S 4 1 1 Futrelle, Mrs. Jacques Heath (Lily May Peel) female 35.000000 1 0 113803 53.1000 S
889 890 891 rd We h [46]: prir 0 Name As w funct	5 0 3 Allen, Mr. William Henry male 35.000000 0 0 373450 8.0500 S
0 Name As w funct	889 0 3 Johnston, Miss. Catherine Helen "Carrie" female 29.699118 1 2 W./C. 6607 23.4500 S 890 1 1 1 Behr, Mr. Karl Howell male 26.000000 0 0 111369 30.0000 C 891 0 3 Dooley, Mr. Patrick male 32.000000 0 0 370376 7.7500 Q ows × 11 columns ave replaced the missing data of Age with mean value to fill the missing values
	ave replaced the missing data of Age with mean value to fill the missing values It (df1['Embarked'].mode()) #finding the mode value of "Embarked" column S : Embarked, dtype: object e can observe that S is the most repititive embarked value so we will replace the 2 missing value with S as we can not find the mean value of categorical data, so we will be using the mode ion to handle the missing value
S [48]: df1[We h [50]: df1. [50]: Pass Surv Pcla Name Sex Age SibS Parc Tick Fare Emba	t(df1['Embarked'].mode()[0]) 'Embarked'].fillna(df1['Embarked'].mode()[0], inplace=True) #replacing the missing value in embarked column with mode value ave successfully fill the "Embarked" column with the mode value. Now we will check for other missing values isnull().sum() engerId 0 ived 0 ss 0 0 0 0 p 0 p 0 h 0 et 0 0
1. 1. 1. 1. 1. 1. 1. 1. 1. 1. 1. 1. 1. 1	1.000000 0.000000 1.000000 0.420000 0.000000 0.000000 223.500000 0.000000 22.000000 0.000000 0.000000 7.910400 446.000000 0.000000 3.000000 29.699118 0.000000 0.000000 14.454200 6 668.500000 1.000000 35.000000 1.000000 0.000000 31.000000
1.	By summarizing the data we can observe there is a huge difference 75% and maximum in fair, passengerid and age and little difference in parce amd sibsp hence outliers are present in the data We also get the standard deviation, min, 25th quartile and 75th quartile values from the describe method.
[52]: 0 1 Name Ae w [53]: sns.	'Survived'].value_counts() 549 342 : Survived, dtype: int64 e can observe out of 891 passeneger only 342 survived and 549 did not survive. countplot('Survived', data=df1) sSubplot:vlabel='Survived', vlabel='count'>
50 40 11 20 100	
[55]: df1[[55]: male fema Name [54]: sns. [54]: <axe 100<="" 40="" 50="" 60="" td=""><td>'Sex'].value_counts() 577 1e 314 : Sex, dtype: int64 countplot('Sex',data=df1) sSubplot:xlabel='Sex', ylabel='count'></td></axe>	'Sex'].value_counts() 577 1e 314 : Sex, dtype: int64 countplot('Sex',data=df1) sSubplot:xlabel='Sex', ylabel='count'>
[56]: sns. [56]: <axe 100<="" 40="" td=""><td></td></axe>	
[60]: sns. [60]: <axe 100<="" 400="" 500="" td=""><td>Do - Do -</td></axe>	Do -
[59]: sns. [59]: <axe 35: 30: 25: 49: 20: 15: 10: 5:</axe 	
[61]: sns. [61]: <axe 100<="" 20="" 40="" 50="" 60="" td=""><td>countplot('Embarked', data=df1) sSubplot:xlabel='Embarked', ylabel='count'></td></axe>	countplot('Embarked', data=df1) sSubplot:xlabel='Embarked', ylabel='count'>
[65]: df1 [65]: 0 1 2 3	PassengerId Survived Pclass Name Sex Age SibSp Parch Ticket Fare Embarked 1 0 3 Braund, Mr. Owen Harris 0 22.000000 1 0 A/5 21171 7.2500 0 2 1 1 Cumings, Mrs. John Bradley (Florence Briggs Th 1 38.000000 1 0 PC 17599 71.2833 1 3 1 3 Heikkinen, Miss. Laina 1 26.000000 0 0 STON/O2. 3101282 7.9250 0 4 1 1 Futrelle, Mrs. Jacques Heath (Lily May Peel) 1 35.000000 1 0 13803 53.1000 0
3 4 886 887 888 889	4 1 1 Futrelle, Mrs. Jacques Heath (Lily May Peel) 1 35.000000 1 0 113803 53.1000 0 5 0 3 Allen, Mr. William Henry 0 35.000000 0 0 373450 8.0500 0
890 891 rd As w [68]: X=df	891 0 3 Dooley, Mr. Patrick 0 32.000000 0 0 370376 7.7500 2 bws × 11 columns e can observe that we have encoded the data andd replaced the categorical data into numerical data. 1.drop(columns=['Name', 'Ticket', 'PassengerId', 'Survived'], axis=1)
Y=df [69]: X [69]: 0 1 2 3 4 886 887 888 889	Pclass Sex Age SibSp Parch Fare Embarked 3 0 22.00000
890 891 rd [70]: Y [70]: 0 1 2 3 4 886 887 888 889 890 Name	3 0 32.00000 0 0 7.7500 2 DWS × 7 columns 0 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1
[]: SPLI [82]: X_tr	TTING THE DATA INTO TRAINING AND TESTING DATA ain, X_test, Y_train, Y_test=train_test_split(X,Y,test_size=0.2,random_state=2)
(891 MOD [84]: mode [85]: mode	t(X.shape, X_test.shape, X_train.shape) , 7) (179, 7) (712, 7) PEL TRAINING *!=LogisticRegression() *!.fit(X_train, Y_train) sticRegression()
MOE [86]: X_tr [87]: prir	sticRegression() EL EVALUATION ain_prediction=model.predict(X_train) t(X_train_prediction) 0 0 0 0 1 0 0 0 1 0 1 0 1 0 0 0 0 1 0 0 1 0 1 1 0 0 1 0 1
0 0 0 0 0 1 0 0 0 0 0 0 0 0 0 0 0 0 0 0	0 0 0 0 1 1 0 0 1 0 1 0 1 0 1 0 1 0 0 0 0 0 0 1 0 1 0 1 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 1 1 0 0 1 1 0 0 1 1 0 0 1 1 0 0 1 1 0 0 0 1 0
This [89]: prir ACCU [90]: X_te	ning_data_accuracy=accuracy_score(Y_train, X_train_prediction) is the comparing step of actual data adn predicted data of training it("ACCURACY SCORE OF TRAINING DATA:", training_data_accuracy) RACY SCORE OF TRAINING DATA: 0.8075842696629213 ist_prediction=model.predict(X_test)
[0 0 0 0 1 0 0 0 1 1 0 0 1 1 1 1 1 1 1 1	100000011001011101010110000011
	RACY SCORE OF TEST DATA: 0.7821229050279329 accuracy score of training data and test data is very close which means the model is neither over fitted nor under fitted