	ataset link: https://www.kaggle.com/c/titanic/data  bataset Introduction:  his data set consists of information about the passengers of the RMS Titanic ship And also have info about is that particular passenger has survived that disaster or not. For every individual person, we have information  passenger Id- Id number of passengers in the data set  Name- Name of the passenger  Survival- Person survived or not ( 0 for No & 1 for Yes)  P-class- With what class of ticket that passenger was traveling.  sex- Male Or Female  Age- Age of the person in Years  Sibsp- Number of siblings/spouses on the Titanic  parch- Number of parents/children on the Titanic  Ticket- Ticket number  Fare- Amount of money that person paid to travel
	<ul> <li>Fare- Amount of money that person paid to travel</li> <li>Cabin- Cabin Number of Passenger</li> <li>Embarked- Port of Embarkation ( C = Cherbourg, Q = Queenstown, S = Southampton)</li> <li>Steps we follow:</li> <li>1&gt; Data Collection &amp; Processing</li> <li>2&gt; Data Analysis</li> <li>3&gt; Data Visualization</li> <li>4&gt; Encoding the Categorical Columns</li> <li>5&gt; Model Training</li> <li>6&gt; Model Evaluation</li> </ul> mporting Libraries:
impo impo impo from from from # Suj impo warn	rt numpy as np rt pandas as pd rt matplotlib.pyplot as plt rt seaborn as sns sklearn.model_selection import train_test_split sklearn.linear_model import LogisticRegression sklearn.metrics import accuracy_score  oress warnings rt warnings ings.filterwarnings("ignore")  ca Collection & Processing
tital  • F	rt os
3 4 • C • A • H • V	4 1 1 Futrelle, Mrs. Jacques Heath (Lily May Peel) female 35.0 1 0 113803 53.1000 C123 S 5 0 3 Allen, Mr. William Henry male 35.0 0 0 373450 8.0500 NaN S  Observations:  Is we can see above the data set has 12 columns or features. we already discussed what these features are all about earlier.  Idere the Survived column is the target variable or class label. The target variable is the feature that needs to be predicted by our models.  We have numerical, categorical Types of features  umber of rows and columns  Inic_data.shape
(891,  • g  : tit: <class data<="" range="" td=""><td></td></class>	
11 dtype memor ]: tital <class Range</class 	Name 891 non-null object Sex 891 non-null object Age 714 non-null int64 SibSp 891 non-null int64 Parch 891 non-null int64 Ticket 891 non-null object Fare 891 non-null float64 Cabin 204 non-null object Embarked 889 non-null object Embarked 889 non-null object es: float64(2), int64(5), object(5) ry usage: 83.7+ KB  mic_data.info()  ss 'pandas.core frame.DataFrame'> endex: 891 entries, 0 to 890 columns (total 12 columns):
#  0 1 2 3 4 5 6 7 8 9 10 11 dtype	Column Non-Null Count Dtype  PassengerId 891 non-null int64 Survived 891 non-null int64 Pclass 891 non-null object Sex 891 non-null object Age 714 non-null int64 Sibsp 891 non-null int64 Ticket 891 non-null object Ticket 891 non-null object Fare 891 non-null object Embarked 899 non-null object Ses: float64(2), int64(5), object(5) Ty usage: 83.7+ KB
• T • II • II • S • S	Observation: The shape of data is (891,12) means in our data set we have 891 rows and 12 columns. Each row has info about a passenger so totally we have data of 891 passengers. In the above output, column consists of the name of the column, Non-null Count means How many non-null values we have in that column, Dtype means What type of value that column consists of (int64 means into oat64 means float value, object means string value) In the age column we can see, Out of 891 values, we have 714 non-null values. It implies that we have 177 Null values. (891–714 = 177) Is ame in the Cabin feature Out of 891 values we have only 204 non-null values. It implies that we have 687 Null values. But this is Huge. we have only 23% of values present in the data set and 77% of values are mide on we can drop this feature while making our model. Except for age and cabin, any features do not have any null values.  The cabin feature of missing values in each column  Thic_data.isnull().sum()
Passe Survi Pclas Name Sex Age SibSp Parch Ticke Fare Cabin Emban dtype	engerId 0 ived 0 ss 0 0 177 0 0 0 0 1 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0
• re	rop the "Cabin" column from the dataframe  nic_data = titanic_data.drop(columns='Cabin', axis=1)  eplacing the missing values in "Age" column with mean value  nic_data['Age'].fillna(titanic_data['Age'].mean(), inplace=True)  nding the mode value of "Embarked" column  t(titanic_data['Embarked'].mode())
Name  Print  S  • re  Tital	Embarked, dtype: object  t(titanic_data['Embarked'].mode()[0])  eplacing the missing values in "Embarked" column with mode value  nic_data['Embarked'].fillna(titanic_data['Embarked'].mode()[0], inplace=True)  heck the number of missing values in each column
Passe Survi Pclas Name Sex Age SibSp Parch Ticke Fare Emban dtype	9
• g	etting some statistical measures about the data
75% max  • C • II is	668.500000 1.000000 35.000000 1.000000 0.000000 31.000000
]: tital  1	nding the number of people survived and not survived  nic_data['Survived'].value_counts()  549 342 Survived, dtype: int64  ca Visualization  set()  naking a count plot for "Survived" column
• C • A	Observation: us we can see 'no' has the value > 500 means these people can't survive in the disaster & approx 350 people survived. It is telling the same story as we have seen above.  Disc_data['Sex'].value_counts()  577
]: sns.	Sex male female
• A	
	countplot('Sex', hue='Survived', data=titanic_data) sSubplot:xlabel='Sex', ylabel='count'>  Survived  0  1
sns.	male female Sex  naking a count plot for "Pclass" column  countplot('Pclass', data=titanic_data) ssubplot:xlabel='Pclass', ylabel='count'>
	1 2 3 Pclass 3  countplot('Pclass', hue='Survived', data=titanic_data) sSubplot:xlabel='Pclass', ylabel='count'>
350 300 250 150 100 50	
• I  n  • F  • F  • S  • A  f  f  g  g  g  g  g  g  g  g  g  g  g	Observations: have heard or read somewhere there is not any value in poor people's lives. The same concept is applying here. In the above output plot, we can see Persons who were traveling in 3rd class, most of them died or of ot survived than other class passengers.  Persons who were traveling in 2nd class, out of them almost equal number of people died and survived.  Persons who were traveling in 1st class, out of a large number of people survived and also a fair number of persons died.  Persons who were traveling in 1st class, out of a large number of people survived and also a fair number of persons who were traveling in higher class like 1st has higher chance to survive.  Persons who were traveling in higher class like 1st has higher chance to survive.  Persons who were traveling in higher class like 1st has higher chance to survive.  Persons who were traveling in higher class like 1st has higher chance to survive.  Persons who were traveling in higher class like 1st has higher chance to survive.
- (	Survived  0 1
• F • F	Observations:  Observ
sns.( - S - S - S - S - S - S - S - S - S - S	
• F	
• F	thers.  Parch Feature  countplot(titanic_data["Parch"], hue = titanic_data["Survived"], data = titanic_data)  ssubplot:xlabel='Parch', ylabel='count'>
100 • C • li	Survived  0 1
sns.: sns.: plt.:	
1.5 1.0 500 400 300 200 100	Survived 0 1
End  tital  male femal Name	ric_data['Sex'].value_counts()  577
• C	### Tichear   Fare   Fare   Embarked   Ticket   Fare   Embarked      1
]: X =	Pclass Sex Age SibSp Parch Fare Embarked 3 0 22.000000 1 0 7.2500 0 1 1 38.000000 1 0 71.2833 1
2 3 4  886 887 888 889 890 [891 ]: prin	3 1 26.000000 0 0 7.9250 0 1 1 35.000000 1 0 53.1000 0 35.00000 0 0 8.0500 0 0
4 886 887 888 889 890 Name • S	1 0 1 0 1 0 1 0 Survived, Length: 891, dtype: int64  pplitting the data into training data & Test data  ain, X_test, Y_train, Y_test = train_test_split(X,Y, test_size=0.2, random_state=2)  t(X.shape, X_train.shape, X_test.shape) 7) (712, 7) (179, 7)
Mo Log : mode.  • tr	del Training  istic Regression  1 = LogisticRegression()  aning the Logistic Regression model with training data  1.fit(X_train, Y_train)  sticRegression()
<ul><li>Mo</li><li>A</li><li>a</li><li>a: X_tra</li><li>prin</li></ul>	del Evaluation  accuracy Score  ccuracy on training data  ain_prediction = model.predict(X_train)  t(X_train_prediction)  0 0 0 0 1 0 0 0 1 0 0 1 0 1 0 1 0 1 0
0 0	0 0 0 0 1 0 0 0 1 0 1 0 1 0 0 0 0 1 0 1