	Problem Statement:  ou are the Data Scientist at a telecom company "Neo" whose customers are churning out to its competitors. You have to analyse the data of your company and find insights and stop your customers from churning out to other telecom companies.
	sks to be done: A) Data Manipulation: a. Extract the 5th column & store it in 'customer_5' b. Extract the 15th column & store it in 'customer_15' c. Extract all the male senior citizens whose Payment Method is Electronic check & store the result in 'senior_male_electronic' d. Extract all lose customers whose tenure is greater than 70 months or their Monthly charges is more than 100\$ & store the result in 'customer_total_tenure' e. Extract all the customers whose Contract is of two years, payment method is Mailed check & the value of Churn is 'Yes' & store the result in 'wo_mail_yes' f. Extract 333 random records from the customer_churndataframe& store the result in 'customer_333' g. Get the count of different levels from the 'Churn' column
	mport numpy as np mport pandas as pd mport matplotlib.pyplot as plt mport seaborn as sns  rom sklearn.linear_model import LogisticRegression
	rom sklearn.model_selection import train_test_split  ata = pd.read_csv(r'C:\Users\LENOVO\Downloads\datasets for python\Customer_churn.csv')  ata.head()  customerID gender SeniorCitizen Partner Dependents tenure PhoneService MultipleLines InternetService OnlineSecurity DeviceProtection TechSupport StreamingTV StreamingMovies Contract PaperlessBilling PaymentMethod MonthlyCharges Churn  No phone
	7590-VHVEG         Female         0         Yes         No         1         No         No         No         No         No         No         No         No         Pemale         No         No         No         No         No         No         Pemale         No         No         No         Pemale         No         No         No         No         Pemale         No         Pemale         Pemale         No
	7795- Male 0 No No 45 No No phone service DSL Yes Yes Yes No
In [108	ata1 = data.loc[:,'tenure'] ata1  1 34
	2 45 2 2 2 2 2 2 2 2 2 2 2 2 2 2 2 2 2 2
In [109	042 66 ame: tenure, Length: 7043, dtype: int64  ata2 = data.iloc[:,14] ata2.head()  No No
	No No No ame: StreamingMovies, dtype: object enior_male_electronic=data[(data['gender']=="Male") & (data['SeniorCitizen']== 1) & (data['PaymentMethod'] == "Electronic check")]
Out[110]:	customerID gender SeniorCitizen Partner Dependents tenure PhoneService MultipleLines InternetService OnlineSecurity DeviceProtection TechSupport StreamingTV StreamingTV StreamingTV StreamingTV StreamingTV PaperlessBilling PaymentMethod MonthlyCharges Churn  8779- QRDMV Male 1 1 No No No 1 No Phone service DSL No Yes No No No Yes Month-to- month Month-to- M
	To 5067-XJQFU Male 1 No No 18 Yes Fiber optic No No No No Yes Fiber optic No Yes Yes Yes One year Yes Electronic check 108.45 7076.35 No
In [111	1 2424-WVHPL Male 1 No No 1 Yes No Fiber optic No No Yes No No Month No Electronic check 74.70 74.7 No rows × 21 columns  ustomer_total_tenure= data[(data['tenure']>70)   (data['MonthlyCharges']>100)]  ustomer_total_tenure.head()
Out[111]:	customerID gender GeniorCitizen Partner Dependers tenure Phone Processing Partner Partner Partner Phone Processing Partne
	3 0280-XJGEX Male 0 No No 49 Yes Yes Fiber optic No Yes No Yes Yes Month-to-month Yes Month-to-month Yes Month-to-month Yes Electronic check 105.50 2686.05 No 3655- SNQYZ Female 0 Yes
In [112	rows × 21 columns  wo_mail_yes= data[(data['Contract']=="Month-to-month") & (data['PaymentMethod']=="Mailed check") & (data['Churn']=="Yes")]  wo_mail_yes.head()  customerID gender SeniorCitizen Partner Dependents tenure PhoneService MultipleLines InternetService OnlineSecurity DeviceProtection TechSupport StreamingTV Stream
	2 3668-QPYBK Male 0 No No 2 Yes No DSL Yes No No No No No No No Month-to-month Yes Mailed check 53.85 108.15 Yes  22 1066-JKSGK Male 0 No No No 1 Yes No No No No No internet service No internet service service service service when the service with the service with the service when the service with the service with the service when the service with the service with the service with the service when the service with the service with the service when the service when the service with the service when the service with the service when the service with the service with the service with the service when the service with the s
	97 VANOG Male 0 No No 5 Yes No No Fiber optic No
In [113	customer_333 = data.sample(n=333) ustomer_333.head()  customerID gender SeniorCitizen Partner Dependents tenure PhoneService MultipleLines InternetService OnlineSecurity DeviceProtection TechSupport StreamingTV StreamingTV StreamingMovies Contract PaperlessBilling PaymentMethod MonthlyCharges Churn
	4589-IUAJB Male 0 Yes No 70 Yes Yes No No No internet service No i
	O75 SO81- NWSUP Female 0 No No 10 Yes No DSL No
In [114 Out[114]:	rows × 21 columns  ata['Churn'].value_counts()  c 5174 es 1869 ame: Churn, dtype: int64
	Build a bar-plot for the 'InternetService' column: i. Set x-axis label to 'Categories of Internet Service' ii. Set y-axis label to 'Count of Categories' iii. Set the title of plot to be 'Distribution of Internet Service' iv. Set the color of the bars to be 'orange' b. Build a histogram for the 'tenure' column: i. Set the number of bins to be 30 ii. Set the color of the bins to be 'green' iii. Assign the title 'Distribution of tenure' c. Build a scatter-plot between 'MonthlyCharges' & 'tenure'. Map 'MonthlyCharges' to the y-axis & 'tenure' to the 'x-axis': i. Assign the points a color of 'brown' ii.
In [115	the x-axis label to 'Tenure of customer' iii. Set the y-axis label to 'Monthly Charges of customer' iv. Set the title to 'Tenure vs Monthly Charges' d. Build a box-plot between 'tenure' & 'Contract'. Map 'tenure' on the y-axis & 'Contract' on the x-axis.    t. bar (data['InternetService'].value_counts().keys().tolist(), data['InternetService'].value_counts().tolist(), color="orange")   t. xlabel("Categories of Internet Service")   t. ylabel("Count of Categories")   t. title("Distribution of Internet Service")   ext(0.5, 1.0, 'Distribution of Internet Service')
Out[115]:	Distribution of Internet Service  3000 - 2500 -
	1500 - 1000 - 500 -
	Fiber optic DSL No Categories of Internet Service  1t.hist(data['tenure'], bins=30, color="green") 1t.title("Distribution of tenure")
Out[116]:	Distribution of tenure  Distribution of tenure
	lt.scatter(x=data['tenure'], y=data['MonthlyCharges'],color="brown") lt.xlabel('Tenure of customer') lt.ylabel('Monthly Charges of customer')
	ext(0.5, 1.0, 'Tenure vs Monthly Charges')  Tenure vs Monthly Charges  Tenure vs Monthly Charges  120
	100 - 80 - 60 -
In 「118	20
	lt.xlabel('Contract') lt.ylabel('tenure') ext(0, 0.5, 'tenure')  Boxplot grouped by Contract
	70
	20 10 0 Month-to-month One year Two year Contract
	E) Linear Regression:  Build a simple linear model where dependent variable is 'MonthlyCharges' and independent variable is 'tenure' i. Divide the dataset into train and test sets in 70:30 ratio. ii. Build the model on train set and predict the values on test set iii. After predicting the values, find the root mean quare error iv. Find out the error in prediction & store the result in 'error' v. Find the root mean square error
	rom sklearn.linear_model import LogisticRegression rom sklearn import linear_model rom sklearn.linear_model import LinearRegression rom sklearn.model_selection import train_test_split  = data[['MonthlyCharges']]
	<ol> <li>34</li> <li>2</li> <li>45</li> <li>2</li> </ol>
	<ul> <li></li> <li>038</li></ul>
	43 rows × 1 columns  _train, x_test, y_train, y_test =train_test_split(x, y, test_size=0.30, random_state=0)
Out[122]: In [123	_train.shape,y_train.shape,x_test.shape,y_test.shape (4930, 1), (4930, 1), (2113, 1), (2113, 1))  egression= LinearRegression()
Out[123]: In [124…	egression.fit(x_train,y_train) inearRegression()  pridict= regression.predict(x_test)  rom sklearn.metrics import mean_squared_error
Out[125]:	rom sklearn.metrics import mean_squared_error p.sqrt (mean_squared_error(y_test,y_pridict)) 3.816464447907975  D) Logistic Regression:
In [126	Build a simple logistic regression modelwhere dependent variable is 'Churn' & independent variable is 'MonthlyCharges' i. Divide the dataset in 65:35 ratio ii. Build the model on train set and predict the values on test set iii. Build the confusion matrix and get the accuracy score b. Build multiple logistic regression model where dependent variable is 'Churn' & independent variables are 'tenure' & 'MonthlyCharges' i. Divide the dataset in 80:20 ratio ii. Build the model on train set and predict the values on test set iii. Build the confusion matrix and get the accuracy score    Churn'   Ch
Out[126]:	Churn  0 No  1 No
	1 No 2 Yes 3 No 4 Yes
	038 No 039 No 040 No 041 Yes
In [127	No  43 rows × 1 columns  _train, x_test, y_train, y_test = train_test_split(x, y, test_size=0.35, random_state=0)  rint(x_train, shape , x_test_shape , y_train, shape, y_train, shape , y_train,
In [129	rint(x_train.shape , x_test.shape ,y_train.shape, y_test.shape) 4577, 1) (2466, 1) (4577, 1) (2466, 1)  rom sklearn.linear_model import LogisticRegression  og = LogisticRegression()
	og.fit(x_train, y_train) :\Users\LENOVO\anaconda3\lib\site-packages\sklearn\utils\validation.py:993: DataConversionWarning: A column-vector y was passed when a 1d array was expected. Please change the shape of y to (n_samples, ), for example using rave: 0. y = column_or_ld(y, warn=True) orgisticRegression()
In [132 Out[132]:	<pre>pred = log.predict(x_test) pred  rray(['No', 'No', 'No', 'No', 'No', 'No'], dtype=object)  rom sklearn.metrics import confusion_matrix,accuracy_score</pre>
In [134 Out[134]:	onfusion_matrix(y_test, y_pred), accuracy_score(y_test, y_pred)  array([[1815, 0],         [651, 0]], dtype=int64), 0.7360097323600974)
Out[135]: In [136…	1815)/(1815+651)  .7360097323600974  =data[['MonthlyCharges','tenure']] =data[['Churn']]
Out[136]:	Churn  0 No  1 No  2 Yes
	2 Yes 3 No 4 Yes 038 No
	039 No 040 No 041 Yes 042 No
In [137 In [138	43 rows × 1 columns  _train, x_test, y_train, y_test = train_test_split(x, y, test_size=0.20, random_state=0)  rint(x_train.shape , x_test.shape , y_train.shape, y_test.shape)  5634 2) (1409 2) (5634 1) (1409 1)
In [139 In [140	rom sklearn.linear_model import LogisticRegression  og = LogisticRegression()  og.fit(x_train, y_train)
Out[140]: In [141	:\Users\LENOVO\anaconda3\lib\site-packages\sklearn\utils\validation.py:993: DataConversionWarning: A column-vector y was passed when a 1d array was expected. Please change the shape of y to (n_samples, ), for example using rave.  y = column_or_ld(y, warn=True)  orgisticRegression()  pred = log.predict(x_test)
Out[141]: In [142	pred  rray(['No', 'No', 'No', 'No', 'No', 'No'], dtype=object)  rom sklearn.metrics import confusion_matrix,accuracy_score  onfusion_matrix(y_test, y_pred), accuracy_score(y_test, y_pred)
Out[143]: In [144	array([[934, 107],
	) Decision Tree:  Build a decision tree model where dependent variable is 'Churn' & independent variable is 'tenure' i. Divide the dataset in 80:20 ratio ii. Build the model on train set and predict the values on test set iii. Build the confusion matrix and calculate the accuracy
	=data[['tenure']] =data[['Churn']]  Churn  O No
	<ul> <li>No</li> <li>No</li> <li>Yes</li> <li>No</li> <li>Yes</li> </ul>
	4 Yes 038 No 040 No
	041       Yes         042       No         43 rows × 1 columns       1 columns
In [147	<pre>rom sklearn.tree import DecisionTreeClassifier _train, x_test, y_train, y_test =train_test_split(x,y,test_size=0.20,random_state=0)  y_tree = DecisionTreeClassifier() y_tree.fit(x_train,y_train)</pre>
Out[147]:	_pred= my_tree.predict(x_test)
	rom sklearn.metrics import confusion_matrix,accuracy_score

