In [2]:
```python
# QUESTION
#Problem Statement:
#You 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.
# .............................................................
# Tasks 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 those 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 'two_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
```

```
  File "<ipython-input-2-f41482e7624f>", line 8
    A) Data Manipulation:
     ^
SyntaxError: unmatched ')'
```

In [3]:
```python
# ANSWER BELLOW
```

In [11]:
```python
import pandas as pd
import numpy as np
import seaborn as sns
import matplotlib.pyplot as plt
import warnings
warnings.filterwarnings('ignore')
sns.set(style="darkgrid")
%matplotlib inline
```

In [12]:
```python
customer = pd.read_csv("customer_churn.csv")
```

In [13]:
```python
customer
```

Out[13]:

| | customerID | gender | SeniorCitizen | Partner | Dependents | tenure | PhoneService | MultipleLines | InternetService | OnlineSecurity | ... | DeviceP |
|---|---|---|---|---|---|---|---|---|---|---|---|---|
| 0 | 7590-VHVEG | Female | 0 | Yes | No | 1 | No | No phone service | DSL | No | ... | |
| 1 | 5575-GNVDE | Male | 0 | No | No | 34 | Yes | No | DSL | Yes | ... | |
| 2 | 3668-QPYBK | Male | 0 | No | No | 2 | Yes | No | DSL | Yes | ... | |
| 3 | 7795-CFOCW | Male | 0 | No | No | 45 | No | No phone service | DSL | Yes | ... | |
| 4 | 9237-HQITU | Female | 0 | No | No | 2 | Yes | No | Fiber optic | No | ... | |
| ... | ... | ... | ... | ... | ... | ... | ... | ... | ... | ... | ... | |
| 7038 | 6840-RESVB | Male | 0 | Yes | Yes | 24 | Yes | Yes | DSL | Yes | ... | |
| 7039 | 2234-XADUH | Female | 0 | Yes | Yes | 72 | Yes | Yes | Fiber optic | No | ... | |
| 7040 | 4801-JZAZL | Female | 0 | Yes | Yes | 11 | No | No phone service | DSL | Yes | ... | |
| 7041 | 8361-LTMKD | Male | 1 | Yes | No | 4 | Yes | Yes | Fiber optic | No | ... | |
| 7042 | 3186-AJIEK | Male | 0 | No | No | 66 | Yes | No | Fiber optic | Yes | ... | |

7043 rows × 21 columns

In [7]:
```python
customer.head(5)
```

Out[7]:

| | customerID | gender | SeniorCitizen | Partner | Dependents | tenure | PhoneService | MultipleLines | InternetService | OnlineSecurity | ... | DeviceProte |
|---|---|---|---|---|---|---|---|---|---|---|---|---|

| | | | | | | | | | | | |
|---|---|---|---|---|---|---|---|---|---|---|---|
| **0** | 7590-VHVEG | Female | 0 | Yes | No | 1 | No | No phone service | DSL | No | ... |
| **1** | 5575-GNVDE | Male | 0 | No | No | 34 | Yes | No | DSL | Yes | ... |
| **2** | 3668-QPYBK | Male | 0 | No | No | 2 | Yes | No | DSL | Yes | ... |
| **3** | 7795-CFOCW | Male | 0 | No | No | 45 | No | No phone service | DSL | Yes | ... |
| **4** | 9237-HQITU | Female | 0 | No | No | 2 | Yes | No | Fiber optic | No | ... |

5 rows × 21 columns

In [14]:
```python
customer.iloc[:,5]
```

Out[14]:
```
0        1
1       34
2        2
3       45
4        2
        ..
7038    24
7039    72
7040    11
7041     4
7042    66
Name: tenure, Length: 7043, dtype: int64
```

In [15]:
```python
customer.iloc[:,15]
```

Out[15]:
```
0       Month-to-month
1             One year
2       Month-to-month
3             One year
4       Month-to-month
             ...
7038          One year
7039          One year
7040    Month-to-month
7041    Month-to-month
7042          Two year
Name: Contract, Length: 7043, dtype: object
```

In [18]:
```python
seniorcitizen=customer[(customer["gender"]=="Male")&(customer["SeniorCitizen"]==1)&(customer["PaymentMethod"]=="E
```

In [19]:
```python
seniorcitizen
```

Out[19]:

| | customerID | gender | SeniorCitizen | Partner | Dependents | tenure | PhoneService | MultipleLines | InternetService | OnlineSecurity | ... | DeviceP |
|---|---|---|---|---|---|---|---|---|---|---|---|---|
| **20** | 8779-QRDMV | Male | 1 | No | No | 1 | No | No phone service | DSL | No | ... | |
| **55** | 1658-BYGOY | Male | 1 | No | No | 18 | Yes | Yes | Fiber optic | No | ... | |
| **57** | 5067-XJQFU | Male | 1 | Yes | Yes | 66 | Yes | Yes | Fiber optic | No | ... | |
| **78** | 0191-ZHSKZ | Male | 1 | No | No | 30 | Yes | No | DSL | Yes | ... | |
| **91** | 2424-WVHPL | Male | 1 | No | No | 1 | Yes | No | Fiber optic | No | ... | |
| **...** | ... | ... | ... | ... | ... | ... | ... | ... | ... | ... | ... | |
| **6837** | 6229-LSCKB | Male | 1 | No | No | 6 | Yes | No | Fiber optic | No | ... | |
| **6894** | 1400-MMYXY | Male | 1 | Yes | No | 3 | Yes | Yes | Fiber optic | No | ... | |
| **6914** | 7142-HVGBG | Male | 1 | Yes | No | 43 | Yes | Yes | Fiber optic | No | ... | |
| **6967** | 8739-WWKDU | Male | 1 | No | No | 25 | Yes | Yes | Fiber optic | No | ... | |

| | customerID | gender | SeniorCitizen | Partner | Dependents | tenure | PhoneService | MultipleLines | InternetService | OnlineSecurity | ... | DeviceP |
|---|---|---|---|---|---|---|---|---|---|---|---|---|
| 7032 | 6894-LFHLY | Male | 1 | No | No | 1 | Yes | Yes | Fiber optic | No | ... | |

298 rows × 21 columns

In [16]:
```python
seniorcitizen=customer[(customer["tenure"]>70 )|(customer["MonthlyCharges"]<100)]
```

In [39]:
```python
seniorcitizen
```

Out[39]:

| | customerID | gender | SeniorCitizen | Partner | Dependents | tenure | PhoneService | MultipleLines | InternetService | OnlineSecurity | ... | DeviceP |
|---|---|---|---|---|---|---|---|---|---|---|---|---|
| 20 | 8779-QRDMV | Male | 1 | No | No | 1 | No | No phone service | DSL | No | ... | |
| 55 | 1658-BYGOY | Male | 1 | No | No | 18 | Yes | Yes | Fiber optic | No | ... | |
| 57 | 5067-XJQFU | Male | 1 | Yes | Yes | 66 | Yes | Yes | Fiber optic | No | ... | |
| 78 | 0191-ZHSKZ | Male | 1 | No | No | 30 | Yes | No | DSL | Yes | ... | |
| 91 | 2424-WVHPL | Male | 1 | No | No | 1 | Yes | No | Fiber optic | No | ... | |
| ... | ... | ... | ... | ... | ... | ... | ... | ... | ... | ... | ... | |
| 6837 | 6229-LSCKB | Male | 1 | No | No | 6 | Yes | No | Fiber optic | No | ... | |
| 6894 | 1400-MMYXY | Male | 1 | Yes | No | 3 | Yes | Yes | Fiber optic | No | ... | |
| 6914 | 7142-HVGBG | Male | 1 | Yes | No | 43 | Yes | Yes | Fiber optic | No | ... | |
| 6967 | 8739-WWKDU | Male | 1 | No | No | 25 | Yes | Yes | Fiber optic | No | ... | |
| 7032 | 6894-LFHLY | Male | 1 | No | No | 1 | Yes | Yes | Fiber optic | No | ... | |

298 rows × 21 columns

In [ ]:
```python
# E)Extract all the customers whose contract is of two years,payment method is mailed check & the value of churn
```

In [41]:
```python
seniorcitizen_mailcheck=seniorcitizen[(seniorcitizen['Contract']=='Two year') & (seniorcitizen['PaymentMethod']==
```

In [42]:
```python
seniorcitizen
```

Out[42]:

| | customerID | gender | SeniorCitizen | Partner | Dependents | tenure | PhoneService | MultipleLines | InternetService | OnlineSecurity | ... | DeviceP |
|---|---|---|---|---|---|---|---|---|---|---|---|---|
| 20 | 8779-QRDMV | Male | 1 | No | No | 1 | No | No phone service | DSL | No | ... | |
| 55 | 1658-BYGOY | Male | 1 | No | No | 18 | Yes | Yes | Fiber optic | No | ... | |
| 57 | 5067-XJQFU | Male | 1 | Yes | Yes | 66 | Yes | Yes | Fiber optic | No | ... | |
| 78 | 0191-ZHSKZ | Male | 1 | No | No | 30 | Yes | No | DSL | Yes | ... | |
| 91 | 2424-WVHPL | Male | 1 | No | No | 1 | Yes | No | Fiber optic | No | ... | |
| ... | ... | ... | ... | ... | ... | ... | ... | ... | ... | ... | ... | |
| 6837 | 6229-LSCKB | Male | 1 | No | No | 6 | Yes | No | Fiber optic | No | ... | |
| 6894 | 1400-MMYXY | Male | 1 | Yes | No | 3 | Yes | Yes | Fiber optic | No | ... | |
| 6914 | 7142-HVGBG | Male | 1 | Yes | No | 43 | Yes | Yes | Fiber optic | No | ... | |
| 6967 | 8739-WWKDU | Male | 1 | No | No | 25 | Yes | Yes | Fiber optic | No | ... | |
| 7032 | 6894-LFHLY | Male | 1 | No | No | 1 | Yes | Yes | Fiber optic | No | ... | |

298 rows × 21 columns

```
In [ ]:   ## F)Extract 333 random records from the customer_churn dataframe & store the result in 'customer_333'

In [50]:  customer_333=customer.sample(n=333)

In [82]:  customer_333
```

Out[82]:

| | customerID | gender | SeniorCitizen | Partner | Dependents | tenure | PhoneService | MultipleLines | InternetService | OnlineSecurity | ... | DeviceP |
|---|---|---|---|---|---|---|---|---|---|---|---|---|
| 4808 | 8144-DGHXP | Female | 0 | No | No | 54 | Yes | No | DSL | Yes | ... | |
| 1039 | 9948-YPTDG | Male | 0 | Yes | No | 38 | Yes | No | Fiber optic | Yes | ... | |
| 2783 | 4760-THGOT | Female | 0 | Yes | No | 43 | Yes | Yes | Fiber optic | Yes | ... | |
| 6208 | 0909-SELIE | Male | 0 | Yes | No | 61 | Yes | Yes | DSL | Yes | ... | |
| 2658 | 7473-ZBDSN | Female | 0 | Yes | Yes | 14 | Yes | No | No | No internet service | ... | N |
| ... | ... | ... | ... | ... | ... | ... | ... | ... | ... | ... | ... | |
| 218 | 2040-LDIWQ | Male | 0 | Yes | Yes | 65 | Yes | Yes | DSL | No | ... | |
| 3032 | 5696-EXCYS | Male | 0 | No | No | 17 | Yes | No | No | No internet service | ... | N |
| 5678 | 4797-AXPXK | Female | 0 | No | Yes | 1 | Yes | No | DSL | Yes | ... | |
| 5690 | 0336-KXKFK | Male | 0 | No | No | 72 | No | No phone service | DSL | Yes | ... | |
| 3479 | 8229-TNIQA | Female | 0 | No | No | 4 | Yes | No | DSL | No | ... | |

333 rows × 21 columns

```
In [ ]:   #g) Get the count of different levels from the 'churn' column

In [83]:  customer['Churn'].value_counts()

Out[83]:  No     5174
          Yes    1869
          Name: Churn, dtype: int64

In [4]:   # QUESTION
          # 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. Set 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

            File "<ipython-input-4-4bf9dc5dc80a>", line 2
              i. Set x-axis label to 'Categories of Internet Service'
                 ^
          SyntaxError: invalid syntax


In [53]:  #i)set x-axis label to 'Categories of internet Service
```

In [84]:
```python
import matplotlib.pyplot as plt
```

In [85]:
```python
customer['InternetService'].value_counts()
```

Out[85]:
```
Fiber optic     3096
DSL             2421
No              1526
Name: InternetService, dtype: int64
```

In [86]:
```python
x=customer['InternetService'].value_counts()
```
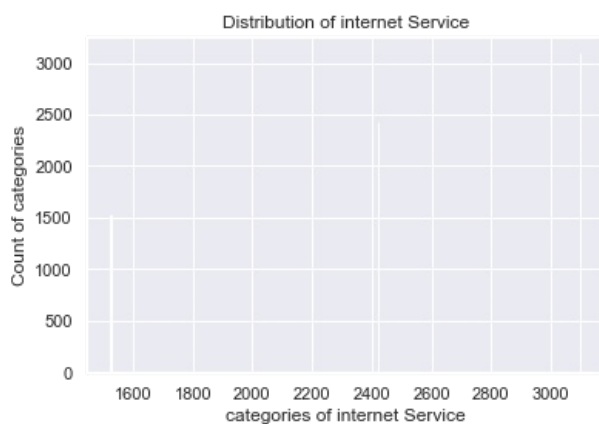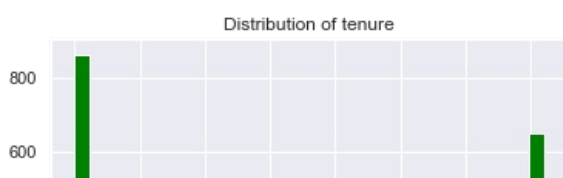
In [87]:
```python
x
```

Out[87]:
```
Fiber optic     3096
DSL             2421
No              1526
Name: InternetService, dtype: int64
```

In [77]:
```python
y=customer['InternetService'].value_counts().tolist()
```

In [78]:
```python
y
```

Out[78]:
```
[3096, 2421, 1526]
```

In [98]:
```python
import numpy as np
import matplotlib.pyplot as plt

plt.bar(x,y,color='red')
plt.xlabel("categories of internet Service")
plt.ylabel("Count of categories")
plt.title("Distribution of internet Service")
plt.show()
```



In [ ]:
```python
# b)Build a histogram for the 'tenure' column:
## i)set num of bins to be 30
## ii)set the color of bins to be 'green'
## iii)Assign the title 'Distribution of tenure'
```

In [94]:
```python
plt.hist(customer['tenure'],bins=30,color='green')
plt.title("Distribution of tenure")
plt.show()
```

```
### Built a scatter-plot between 'MonthlyCharges' & 'tenure'. Map 'MonthlyCharges' to the y-axis & 'tenure' to th
#### i)Assign the points a color of 'brown'
#### ii)Set the x-axis label to tenure of customer
#### iii)Set the y-axis label to 'Monthly Charges of customer'
#### iV)Set the title to the 'Tenure vs Monthly Charges'
```

In [ ]:

```python
import pandas as pd
```

In [99]:

```python
data=pd.read_csv("customer_churn.csv")
```

In [100...

```python
data
```

Out[100...

| | customerID | gender | SeniorCitizen | Partner | Dependents | tenure | PhoneService | MultipleLines | InternetService | OnlineSecurity | ... | DeviceP |
|---|---|---|---|---|---|---|---|---|---|---|---|---|
| 0 | 7590-VHVEG | Female | 0 | Yes | No | 1 | No | No phone service | DSL | No | ... | |
| 1 | 5575-GNVDE | Male | 0 | No | No | 34 | Yes | No | DSL | Yes | ... | |
| 2 | 3668-QPYBK | Male | 0 | No | No | 2 | Yes | No | DSL | Yes | ... | |
| 3 | 7795-CFOCW | Male | 0 | No | No | 45 | No | No phone service | DSL | Yes | ... | |
| 4 | 9237-HQITU | Female | 0 | No | No | 2 | Yes | No | Fiber optic | No | ... | |
| ... | ... | ... | ... | ... | ... | ... | ... | ... | ... | ... | ... | |
| 7038 | 6840-RESVB | Male | 0 | Yes | Yes | 24 | Yes | Yes | DSL | Yes | ... | |
| 7039 | 2234-XADUH | Female | 0 | Yes | Yes | 72 | Yes | Yes | Fiber optic | No | ... | |
| 7040 | 4801-JZAZL | Female | 0 | Yes | Yes | 11 | No | No phone service | DSL | Yes | ... | |
| 7041 | 8361-LTMKD | Male | 1 | Yes | No | 4 | Yes | Yes | Fiber optic | No | ... | |
| 7042 | 3186-AJIEK | Male | 0 | No | No | 66 | Yes | No | Fiber optic | Yes | ... | |

7043 rows × 21 columns

In [101...

```python
data.head(2)
```

Out[101...

| | customerID | gender | SeniorCitizen | Partner | Dependents | tenure | PhoneService | MultipleLines | InternetService | OnlineSecurity | ... | DeviceProte |
|---|---|---|---|---|---|---|---|---|---|---|---|---|
| 0 | 7590-VHVEG | Female | 0 | Yes | No | 1 | No | No phone service | DSL | No | ... | |
| 1 | 5575-GNVDE | Male | 0 | No | No | 34 | Yes | No | DSL | Yes | ... | |

2 rows × 21 columns

In [102...

```python
import matplotlib.pyplot as plt
```

In [105...

```python
plt.scatter(x=data['tenure'].head(20),y=data['MonthlyCharges'].head(20),color='brown')
plt.xlabel('tenure of customer')
plt.ylabel('Monthly Charges of customer')
plt.title('Tenure vs Monthly Charges')
plt.show()
```
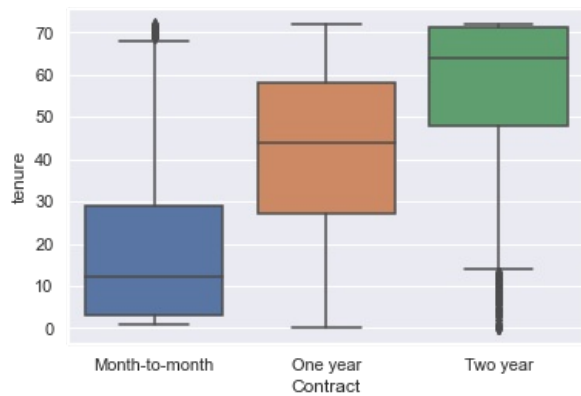
Tenure vs Monthly Charges

Tenure vs Monthly Charges

In [106... `## d)built a boxplot between tenure & contract. Map 'tenure' on y-axis and 'Contract' on the x-axis`

In [108... 
```python
import seaborn as sns
```

In [109... 
```python
sns.boxplot(x=data['Contract'],y=data['tenure'])
plt.show()
```



In [110... 
```python
import seaborn as sns
import matplotlib.pyplot as plt
```

In [111... 
```python
data.boxplot(by=['Contract'],column="tenure",figsize=(7,5),color='olive')
plt.xlabel("Contract")
plt.ylabel("Tenure")
plt.title("Tenure vs Contract")
plt.show()
```



In [ ]: 
```python
# QUESETION
# C) Linear Regression:
# a. 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 square error
#iv. Find out the error in prediction & store the result in 'error'
#v. Find the root mean square error
```

In [ ]:
```
# i)Divide the dataset into train and test sets is 70:30 ratio
```

In [112...
```python
from sklearn.model_selection import train_test_split
from sklearn.linear_model import LinearRegression
```

In [113...
```python
data.head(2)
```

Out[113...

| | customerID | gender | SeniorCitizen | Partner | Dependents | tenure | PhoneService | MultipleLines | InternetService | OnlineSecurity | ... | DeviceProte |
|---|---|---|---|---|---|---|---|---|---|---|---|---|
| 0 | 7590-VHVEG | Female | 0 | Yes | No | 1 | No | No phone service | DSL | No | ... | |
| 1 | 5575-GNVDE | Male | 0 | No | No | 34 | Yes | No | DSL | Yes | ... | |

2 rows × 21 columns

In [114...
```python
x=pd.DataFrame(data['tenure']) #independent variable
```

In [116...
```python
y=data['MonthlyCharges']        # dependent variable
```

In [117...
```python
x_train,x_test,y_train,y_test=train_test_split(x,y,test_size=0.3,random_state=0)
print(data.shape)
print(x_train.shape)
print(x_test.shape)
print(y_train.shape)
print(y_test.shape)
```

```
(7043, 21)
(4930, 1)
(2113, 1)
(4930,)
(2113,)
```

In [118...
```python
import numpy as np
import pandas as pd
```

In [119...
```python
data=pd.read_csv("customer_churn.csv")
```

In [120...
```python
data
```

Out[120...

| | customerID | gender | SeniorCitizen | Partner | Dependents | tenure | PhoneService | MultipleLines | InternetService | OnlineSecurity | ... | DeviceP |
|---|---|---|---|---|---|---|---|---|---|---|---|---|
| 0 | 7590-VHVEG | Female | 0 | Yes | No | 1 | No | No phone service | DSL | No | ... | |
| 1 | 5575-GNVDE | Male | 0 | No | No | 34 | Yes | No | DSL | Yes | ... | |
| 2 | 3668-QPYBK | Male | 0 | No | No | 2 | Yes | No | DSL | Yes | ... | |
| 3 | 7795-CFOCW | Male | 0 | No | No | 45 | No | No phone service | DSL | Yes | ... | |
| 4 | 9237-HQITU | Female | 0 | No | No | 2 | Yes | No | Fiber optic | No | ... | |
| ... | ... | ... | ... | ... | ... | ... | ... | ... | ... | ... | ... | ... |
| 7038 | 6840-RESVB | Male | 0 | Yes | Yes | 24 | Yes | Yes | DSL | Yes | ... | |
| 7039 | 2234-XADUH | Female | 0 | Yes | Yes | 72 | Yes | Yes | Fiber optic | No | ... | |
| 7040 | 4801-JZAZL | Female | 0 | Yes | Yes | 11 | No | No phone service | DSL | Yes | ... | |
| 7041 | 8361-LTMKD | Male | 1 | Yes | No | 4 | Yes | Yes | Fiber optic | No | ... | |

| | customerID | gender | SeniorCitizen | Partner | Dependents | tenure | PhoneService | MultipleLines | InternetService | OnlineSecurity | ... | DeviceProte |
|---|---|---|---|---|---|---|---|---|---|---|---|---|
| 7042 | 3186-AJIEK | Male | 0 | No | No | 66 | Yes | No | Fiber optic | Yes | ... | |

7043 rows × 21 columns

```
In [122... data.head(3)
```

| | customerID | gender | SeniorCitizen | Partner | Dependents | tenure | PhoneService | MultipleLines | InternetService | OnlineSecurity | ... | DeviceProte |
|---|---|---|---|---|---|---|---|---|---|---|---|---|
| 0 | 7590-VHVEG | Female | 0 | Yes | No | 1 | No | No phone service | DSL | No | ... | |
| 1 | 5575-GNVDE | Male | 0 | No | No | 34 | Yes | No | DSL | Yes | ... | |
| 2 | 3668-QPYBK | Male | 0 | No | No | 2 | Yes | No | DSL | Yes | ... | |

3 rows × 21 columns

```
In [123... lr=LinearRegression()
         lr.fit(x_train,y_train)
Out[123... LinearRegression()
```

```
In [124... y_pred=lr.predict(x_test)
```

```
In [125... y_pred
Out[125... array([60.95089608, 72.98096699, 59.1903979 , ..., 75.62171426,
                70.63363608, 65.6455579 ])
```

```
In [126... y_test.values
Out[126... array([ 58.2 , 116.6 ,  71.95, ..., 109.95,  24.55,  81.6 ])
```

```
In [ ]:  #D)Logistic Regression
         # a)Built a simple logistic regressin model where dependent variable is 'churn' & independent variable is 'Monthl
         #i)Divide the dataset into 65:35 ratio
         #ii)Build the model on train set and predict the values on test set
         # iV)Build the confusion matrix and get the accuracy score¶
```

```
In [127... from sklearn.metrics import mean_squared_error
         import numpy as np
```

```
In [128... msc=mean_squared_error(y_test,y_pred)
```

```
In [129... error=np.sqrt(msc)
```

```
In [130... error
Out[130... 29.394584027273893
```

```
In [131...
```

```
In [132... from sklearn.linear_model import LogisticRegression
```

```
In [133... x=pd.DataFrame(data['MonthlyCharges'])
```

```python
In [134… y=data['Churn']
```

```python
In [139… y_pred
```

```
Out[139… array([60.95089608, 72.98096699, 59.1903979 , ..., 75.62171426,
               70.63363608, 65.6455579 ])
```

```python
In [140… y_test.values
```

```
Out[140… array([ 58.2 , 116.6 ,  71.95, ..., 109.95,  24.55,  81.6 ])
```

```python
In [145… from sklearn.metrics import confusion_matrix,accuracy_score
```

```python
In [146… (1815+0)/(1815+651+0+0)
```

```
Out[146… 0.7360097323600974
```

```python
In [148… x=pd.DataFrame(data.loc[:,['tenure','MonthlyCharges']])
```

```python
In [151… y=data['Churn']
```

```python
In [153… y
```

```
Out[153… 0        No
         1        No
         2       Yes
         3        No
         4       Yes
                ...
         7038     No
         7039     No
         7040     No
         7041    Yes
         7042     No
         Name: Churn, Length: 7043, dtype: object
```

```python
In [154… mlr=LogisticRegression()
```

```python
In [155… mlr
```

```
Out[155… LogisticRegression()
```

```python
In [159… y_pred
```

```
Out[159… array([60.95089608, 72.98096699, 59.1903979 , ..., 75.62171426,
               70.63363608, 65.6455579 ])
```

```python
In [160… y_test.values
```

```
Out[160… array([ 58.2 , 116.6 ,  71.95, ..., 109.95,  24.55,  81.6 ])
```

```python
In [162… (934+156)/(934+156+212+107)
```

`0.7735982966643009`

```python
 # QUESTION
# D) Logistic Regression:
# a. Build a simple logistic regression modelwhere dependent variable is 'Churn' & independent variable is 'Month
# 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
```

```python
# QUESTION
# b. Build a multiple logistic regression model where dependent variable is 'Churn' & independent variables are
# 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
```

```python
x=pd.DataFrame(data.loc[:,['tenure']])
y=data['Churn']
```

```python
x_train,x_test,y_train,y_test=train_test_split(x,y,train_size=0.8,random_state=0)
```

```python
from sklearn.tree import DecisionTreeClassifier
```

```python
# Create Decision Tree classifer object
clf = DecisionTreeClassifier()

# Train Decision Tree Classifer
clf = clf.fit(x_train,y_train)

#Predict the response for test dataset
y_pred = clf.predict(x_test)
```

```python
y_pred
```

`array(['No', 'No', 'No', ..., 'No', 'No', 'Yes'], dtype=object)`

```python
y_test
```

```
2200     No
4627     No
3225     No
2828     No
3768     No
        ...
2631    Yes
5333    Yes
6972    Yes
4598     No
3065     No
Name: Churn, Length: 1409, dtype: object
```

```python
#Import scikit-learn metrics module for accuracy calculation
from sklearn.metrics import confusion_matrix,accuracy_score
from sklearn import metrics
```

```python
# Model Accuracy, how often is the classifier correct?
```
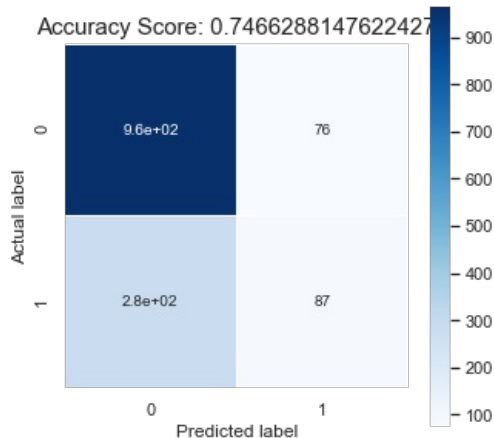
```
# Model Accuracy, how often is the classifier correct?
print("Accuracy:",metrics.accuracy_score(y_test, y_pred))
```

Accuracy: 0.7466288147622427

In [177...
```
cm = confusion_matrix(y_test, y_pred)
plt.figure(figsize=(5,5))
sns.heatmap(data=cm,linewidths=.5, annot=True,square = True,  cmap = 'Blues')
plt.ylabel('Actual label')
plt.xlabel('Predicted label')
all_sample_title = 'Accuracy Score: {0}'.format(clf.score(x_test, y_test))
plt.title(all_sample_title, size = 15)
```

Out[177... Text(0.5, 1.0, 'Accuracy Score: 0.7466288147622427')



In [178...
```
!pip install graphviz
```

```
Collecting graphviz
  Downloading graphviz-0.20-py3-none-any.whl (46 kB)
     |████████████████████████████████| 46 kB 313 kB/s eta 0:00:01
Installing collected packages: graphviz
Successfully installed graphviz-0.20
```

In [179...
```
!pip install pydotplus
```

```
Collecting pydotplus
  Downloading pydotplus-2.0.2.tar.gz (278 kB)
     |████████████████████████████████| 278 kB 503 kB/s eta 0:00:01
Requirement already satisfied: pyparsing>=2.0.1 in /opt/anaconda3/lib/python3.8/site-packages (from pydotplus) (2
.4.7)
Building wheels for collected packages: pydotplus
  Building wheel for pydotplus (setup.py) ... done
  Created wheel for pydotplus: filename=pydotplus-2.0.2-py3-none-any.whl size=24566 sha256=12e67bb213ac49306445a6
75f4026cbb2e3ddd3192bc720bba2ab611a8332b75
  Stored in directory: /Users/surajjamadar/Library/Caches/pip/wheels/fe/cd/78/a7e873cc049759194f8271f780640cf96b3
5e5a48bef0e2f36
Successfully built pydotplus
Installing collected packages: pydotplus
Successfully installed pydotplus-2.0.2
```

In [ ]:

In [ ]:

In [ ]:

In [ ]:

In [ ]:

```
In [ ]:

In [9]:  # QUESTION
         E) Decision Tree:
         a. 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

           File "<ipython-input-9-f3425134eba2>", line 2
             E) Decision Tree:
              ^
         SyntaxError: unmatched ')'


In [ ]:

In [ ]:

In [ ]:

In [ ]:

In [ ]:

In [ ]:

In [ ]:  # QUESTION
         F) Random Forest:
         a. Build a Random Forest model where dependent variable is 'Churn' & independent
         variables are 'tenure' and 'MonthlyCharges'
         i. Divide the dataset in 70:30 ratio
         ii. Build the model on train set and predict the values on test set
         iii. Build the confusion matrix and calculate the accuracy
```

Loading [MathJax]/jax/output/CommonHTML/fonts/TeX/fontdata.js