

Customer Churn Prediction: Summary Report

1. Introduction

This project aims to predict customer churn for a telecom company using customer demographic and service-related data. By analyzing customer behavior, the objective is to help the telecom company reduce churn through targeted interventions and strategies for customer retention.

2. Dataset Overview

- Source:** The dataset provided by IBM for telecom customers is publicly available on Kaggle.
- Features:**
 - Demographic (Gender, Partner, Dependents, Senior Citizen)
 - Services (Phone, Internet, Online Security, Device Protection, etc.)
 - Account (Tenure, Contract, Payment method, Monthly Charges)
 - Target: Churn (indicates whether the customer left the company).

3. Data Preprocessing

- Categorical features were converted into numerical form.
- The dataset was cleaned and normalized for numeric columns like Tenure and MonthlyCharges.
- Missing values were handled, and features were scaled using StandardScaler.

4. Exploratory Data Analysis (EDA)

- Churn Distribution:** A substantial portion of customers did not churn, while a smaller percentage (around 26%) churned.

<https://colab.research.google.com/drive/1BvBILUAVZCRVAdptJGL7iarIDTHT#scrollTo=WpuCuCbyhZLuk&printMode=true>

9/28/24, 9:08 PM

- Telecom_cusL_churn_major_project_Academor.ipynb - Colab
- Service Patterns:** Customers with Fiber Optic service showed a higher churn rate, and those on month-to-month contracts were more likely to churn.

5. Model Selection

Three models were selected for prediction:

- Decision Tree Classifier
- Random Forest Classifier
- Logistic Regression

6. Model Performance

- Decision Tree:
 - Accuracy: 78%
 - Strengths: Quick to train, interpretable.
 - Weaknesses: Prone to overfitting.
- Random Forest:
 - Accuracy: 82%
 - Strengths: Reduced overfitting, robust.
 - Weaknesses: Requires more computational power.
- Logistic Regression:
 - Accuracy: 80%
 - Strengths: Simple, good for linear relationships.
 - Weaknesses: Less powerful with complex data interactions.

7. Conclusion

- Best Performing Model:** The Random Forest classifier showed the highest accuracy and was less prone to overfitting compared to the Decision Tree.

<https://colab.research.google.com/drive/1BvBILUAVZCRVAdptJGL7iarIDTHT#scrollTo=WpuCuCbyhZLuk&printMode=true>

9/28/24, 9:08 PM <https://colab.research.google.com/drive/1BvBjUAVZCRVAdptauGL7iarR1DTHf#scrollTo-WpuCvbyhZLkI&printMode=true>

df_raw.isnull().any()

```
Index(['customerID', 'gender', 'SeniorCitizen', 'Partner', 'Dependents', 'tenure', 'Phone', 'Dependents', 'tenure', 'Phone']
      [0]    7590-  Female        0     Yes      No       1
      [1]    5575-  Male         0     No      No      34
      [2]    3668-  Male         0     No      No       2
      [3]    7795-  Male         0     No      No      45
      [4]    9237-  Female        0     No      No       2
   [5 rows × 21 columns]
```

```
df_raw.shape
(7043, 21)
```

```
df_raw.columns
```

```
Index(['customerID', 'gender', 'SeniorCitizen', 'Partner', 'Dependents', 'tenure', 'PhoneService', 'MultipleLines', 'InternetService', 'OnlineSecurity', 'OnlineBackup', 'DeviceProtection', 'TechSupport', 'StreamingTV', 'StreamingMovies', 'Contract', 'PaperlessBilling', 'PaymentMethod', 'MonthlyCharges', 'TotalCharges', 'Churn']
      [0]    7590-  Female        0     Yes      No       1
      [1]    5575-  Male         0     No      No      34
      [2]    3668-  Male         0     No      No       2
      [3]    7795-  Male         0     No      No      45
      [4]    9237-  Female        0     No      No       2
   [5 rows × 21 columns]
```

```
df_raw = pd.read_csv("./content/WA_Fn-UseC-Customer-Churn_ (1).csv")
df_raw.head()
```

```
import pandas as pd
import numpy as np
import matplotlib.pyplot as plt
import seaborn as sns

sns.set()
%matplotlib inline
```

*Predict customer behaviour to reduce customer churn rate and improve customer retention. And churn rate formula is, Assumption: Let's assume that all the records given in this dataset are for particular interval.

Understanding and predicting customer behaviour

Telecom customer churn

Source code - Notebook

Link - [https://github.com/ahmeddas/TELECOM_cust_churn_project%20\(1\).ipynb](https://github.com/ahmeddas/TELECOM_cust_churn_project%20(1).ipynb)

Dataset : Telco customer churn Source : <https://www.kaggle.com/blastchar/telco-customer-churn> (IBM Sample dataset) Here, IBM provided customer data for Telco industry to predict behaviour of the customers. Main objective is that to analyze customer behaviour and develop strategies for customer retention.

Problem / Objective of this analysis:

*Predict customer behaviour to reduce customer churn rate and improve customer retention. And churn rate formula is, Assumption: Let's assume that all the records given in this dataset are for particular interval.

```
import pandas as pd
import numpy as np
import matplotlib.pyplot as plt
import seaborn as sns

sns.set()
%matplotlib inline

df_raw = pd.read_csv("./content/WA_Fn-UseC-Customer-Churn_ (1).csv")
df_raw.head()
```

- **Next Steps:** Based on these models, telecom companies can identify customers at high risk of churning and implement strategies such as personalized offers or improved customer support.

```

customerID False
gender False
SeniorCitizen False
Partner False
Dependents False
tenure False
PhoneService False
MultipleLines False
InternetService False
OnlineSecurity False
OnlineBackup False
DeviceProtection False
TechSupport False
StreamingTV False
StreamingMovies False
Contract False
PaperlessBilling False
PaymentMethod False
MonthlyCharges False
TotalCharges False
Churn False

```

dtype: bool

<https://colab.research.google.com/drive/1BvBILUAVZCRVAdptauGL7arR1DTHT#scrollTo-WpuCvbyhZLuk&printMode=true>

9/28/24, 9:08 PM

5/32

`df_raw.info()`

```

<class 'pandas.core.frame.DataFrame'>
RangeIndex: 7043 entries, 0 to 7042
Data columns (total 21 columns):
 #   Column          Non-Null Count  Dtype  
--- 
 0   customerID      7043 non-null   object 
 1   gender          7043 non-null   object 
 2   SeniorCitizen   7043 non-null   int64  
 3   Partner         7043 non-null   object 
 4   Dependents     7043 non-null   object 
 5   tenure          7043 non-null   int64  
 6   PhoneService    7043 non-null   object 
 7   MultipleLines   7043 non-null   object 
 8   InternetService 7043 non-null   object 
 9   OnlineSecurity  7043 non-null   object 
 10  OnlineBackup    7043 non-null   object 
 11  DeviceProtection 7043 non-null   object 
 12  TechSupport    7043 non-null   object 
 13  StreamingTV     7043 non-null   object 
 14  StreamingMovies 7043 non-null   object 
 15  Contract        7043 non-null   object 
 16  PaperlessBilling 7043 non-null   object 
 17  PaymentMethod   7043 non-null   float64 
 18  MonthlyCharges 7043 non-null   object 
 19  TotalCharges    7043 non-null   object 
 20  Churn           7043 non-null   object 
dtypes: float64(1), int64(2), object(18)
memory usage: 1.1+ MB

```

Data Cleaning

Here, we are going to rename some columns as well as modifying records with proper format for further analysis.

```
df_cal = df_raw.copy()
```

```
df_cal.rename(columns={'gender': 'Gender'}
```

<https://colab.research.google.com/drive/1BvBILUAVZCRVAdptauGL7arR1DTHT#scrollTo-WpuCvbyhZLuk&printMode=true>

6/32

```

9/28/24, 9:08 PM Telecom_cust_churn_major_project_Academor.ipynb - Colab
, 'CustomerID' : 'CustomerID',
, 'Contract' : 'ContractType',
, 'InternetService' : 'InternetServiceType',
, 'tenure' : 'Tenure'
}

, inplace=True)

df_cal[['Partner']] = df_cal.Partner.map({'Yes':1, 'No':0})
df_cal[['Dependents']] = df_cal.Dependents.map({'Yes':1, 'No':0})
df_cal[['PhoneService']] = df_cal.PhoneService.map({'Yes':1, 'No':0})
df_cal[['MultipleLines']] = df_cal.MultipleLines.map({'Yes':1, 'No':0, 'No phone

df_cal['InternetService'] = df_cal.InternetServiceType.map({'DSL':1, 'Fiber optic':2, 'Satellite':3, 'Cable':4, 'No':0})
df_cal['OnlineSecurity'] = df_cal.OnlineSecurity.map({'Yes':1, 'No':0, 'No internet':0})
df_cal[['OnlineBackup']] = df_cal.OnlineBackup.map({'Yes':1, 'No':0, 'No internet':0})
df_cal[['DeviceProtection']] = df_cal.DeviceProtection.map({'Yes':1, 'No':0, 'No protection':0})
df_cal[['TechSupport']] = df_cal.TechSupport.map({'Yes':1, 'No':0, 'No technical support':0})
df_cal[['StreamingTV']] = df_cal.StreamingTV.map({'Yes':1, 'No':0, 'No internet service':0})
df_cal[['StreamingMovies']] = df_cal.StreamingMovies.map({'Yes':1, 'No':0, 'No in-house movies':0})
df_cal[['IsContracted']] = df_cal.ContractType.map({'One year':1, 'Two year':1, 'Three or more years':2, 'Not specified':0})
df_cal[['PaperlessBilling']] = df_cal.PaperlessBilling.map({'Yes':1, 'No':0})
df_cal[['Churn']] = df_cal.Churn.map({'Yes':1, 'No':0})

df_cal.head()

```

<https://colab.research.google.com/drive/1BvBilUAVZCRVAdptauGL7iarIDTHT#scrollTo-WpuCuByhZLkI&printMode=true>

9/28/24, 9:08 PM

Telecom_cust_churn_major_project_Academor.ipynb - Colab

	CustomerID	Gender	SeniorCitizen	Partner	Dependents	Tenure	Phone
0	7590-VHVEG	Female	0	1	0	1	
1	5575-GNVDE	Male	0	0	0	34	
2	3668-QPYBK	Male	0	0	0	2	
3	7795-CFOCW	Male	0	0	0	45	
4	9237-HQITU	Female	0	0	0	2	

5 rows × 23 columns

df_cal.dtypes

10/32

<http://colab.research.google.com/drive/1BvBiiIav7ZP9KvAdoTaIgC7iaB1DTHTh#scrollTo=WouCyuhbzIk&rootMode=true>

```
11 rows x 23 columns
```

	count	7043	7043	7043.000000	7043.000000	7043.000000	7043.000000
	unique	7043	2	NaN	NaN	NaN	NaN
top	VHVEG	750-	Male	NaN	NaN	NaN	NaN
freq	1	3555	NaN	NaN	NaN	NaN	NaN
mean	NaN	NaN	0.162147	0.453033	0.299568	3.	
std	NaN	NaN	0.366612	0.499748	0.458110	2.	
min	NaN	NaN	0.000000	0.000000	0.000000	()
25%	NaN	NaN	0.000000	0.000000	0.000000	:	:
50%	NaN	NaN	0.000000	0.000000	0.000000	2:	2:
75%	NaN	NaN	0.000000	1.000000	1.000000	5:	5:
max	NaN	NaN	1.000000	1.000000	1.000000	7:	7:

```
df['ca'].describe(include='all')
```

EDA - Exploratory data analysis

Telecom_cust_churn_major_project_Academtor.ipynb - Colab
9/28/24, 9:08 PM

dtype: object

InternetService int64
https://colab.research.google.com/drive/1BIVBilAUZQrrnAopTAUGL7tA1R1DTH#scrollTo=WbuCvbyhZUKY&printMode=true

	count
Churn	
0	0.73463
1	0.26537

dtype: float64

```
churn_summary = df_cal.groupby(['Churn'])
print(churn_summary[['Tenure', 'MonthlyCharges']].mean())
```

	Tenure	MonthlyCharges
Churn		
0	37.569965	61.265124
1	17.379133	74.441332

Correlation matrix

Distribution of Tenure, Monthly Charges

```
import matplotlib.pyplot as plt
import seaborn as sns

f, axes = plt.subplots(nrows=2, figsize=(15, 6))

# Customer Tenure Distribution
sns.histplot(df_cal.Tenure,
             color='g', # Color of bars
             kde=False, # Disable Gaussian Kernel density estimate
             ax=axes[0], # Plot on given axis
             bins=30) # Specify number of bins
axes[0].set_title("Customer Tenure Distribution")
axes[0].set_ylabel("Customer counts") # Set ylabel of graph

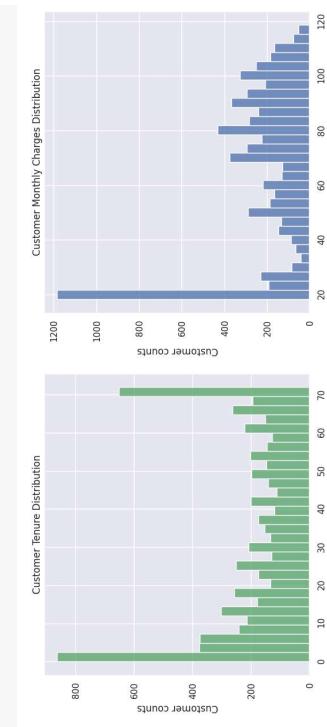
# Customer monthly charges distribution
sns.histplot(df_cal.MonthlyCharges,
             color='b',
             kde=False,
             ax=axes[1],
             bins=30) # Specify number of bins
axes[1].set_title("Customer Monthly Charges Distribution")
axes[1].set_ylabel("Customer counts")
```

<https://colab.research.google.com/drive/1BvJBUAVZCRvAdptuGL7iarIDTHT#scrollTo-WpuCvbyhZlk&printMode=true>

11/32

```
9/28/24, 9:08 PM # Customer monthly charges distribution
sns.histplot(df_cal.MonthlyCharges,
             color='b',
             kde=False,
             ax=axes[1],
             bins=30) # Specify number of bins
axes[1].set_title("Customer Monthly Charges Distribution")
axes[1].set_ylabel("Customer counts")

plt.show()
```



Customer tenure

Distribution shows that customers who has tenure around less than a year left the brand more.

```
import matplotlib.pyplot as plt
import seaborn as sns

fig = plt.figure(figsize=(15, 5))

```

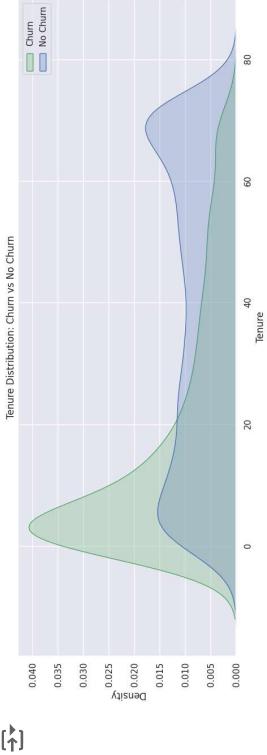
<https://colab.research.google.com/drive/1BvJBUAVZCRvAdptuGL7iarIDTHT#scrollTo-WpuCvbyhZlk&printMode=true>

12/32

```
# KDE plot for customers who churned
ax = sns.kdeplot(df_cal.loc[(df_cal['Churn'] == 1), 'Tenure'],
                  color='g',
                  fill=True, # Use fill instead of shade
                  label='Churn')

# KDE plot for customers who did not churn
ax = sns.kdeplot(df_cal.loc[(df_cal['Churn'] == 0), 'Tenure'],
                  color='b',
                  fill=True, # Use fill instead of shade
                  label='No Churn')

plt.title("Tenure Distribution: Churn vs No Churn")
plt.legend()
plt.show()
```



Customer monthly charges

Distribution shows that customers who has around more than \$65 left brands more

<https://colab.research.google.com/drive/1BvBilUAVZCRVAdptauGL7iarIDTHT#scrollTo-WpuCvbyhZlk&printMode=true>

9/28/24, 9:08 PM

13/32

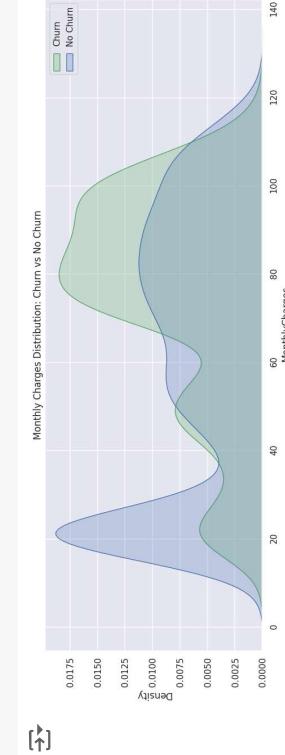
```
import matplotlib.pyplot as plt
import seaborn as sns

fig = plt.figure(figsize=(15, 5))

# KDE plot for customers who churned
ax = sns.kdeplot(df_cal.loc[(df_cal['Churn'] == 1), 'MonthlyCharges'],
                  color='g',
                  fill=True, # Use fill instead of shade
                  label='Churn')

# KDE plot for customers who did not churn
ax = sns.kdeplot(df_cal.loc[(df_cal['Churn'] == 0), 'MonthlyCharges'],
                  color='b',
                  fill=True, # Use fill instead of shade
                  label='No Churn')

plt.title("Monthly Charges Distribution: Churn vs No Churn")
plt.legend()
plt.show()
```



<https://colab.research.google.com/drive/1BvBilUAVZCRVAdptauGL7iarIDTHT#scrollTo-WpuCvbyhZlk&printMode=true>

Tenure vs Monthly charges

```
df_cal['Tenure_norm'] = (df_cal['Tenure'] - df_cal['Tenure'].min()) / (df_cal['Tenure'].max() - df_cal['Tenure'].min())
df_cal['MonthlyCharges_norm'] = (df_cal['MonthlyCharges'] - df_cal['MonthlyCharges'].min()) / (df_cal['MonthlyCharges'].max() - df_cal['MonthlyCharges'].min())
df_cal.head()
```

	CustomerID	Gender	SeniorCitizen	Partner	Dependents	Tenure	Phone
0	7590-VHVEG	Female		0	1	0	1
1	5575-GNVDE	Male		0	0	0	34
2	3668-QPYBK	Male		0	0	0	2
3	7795-CFOCW	Male		0	0	0	45
4	9237-HQIUI	Female		0	0	0	2

5 rows × 25 columns

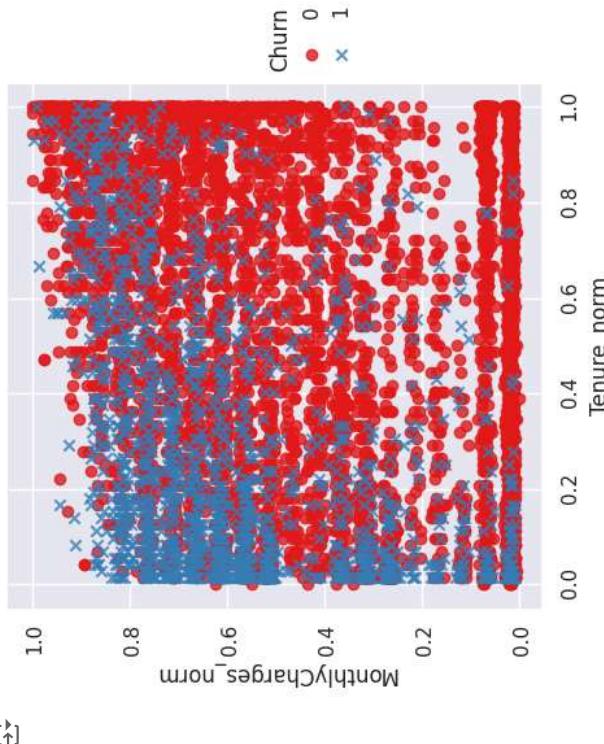
```
sns.lmplot(x='Tenure_norm',y='MonthlyCharges_norm',data=df_cal
            ,hue='Churn'
            ,fit_reg=False
            ,markers=['o','x']
            ,palette="Set1"
            )
plt.show()
```

<https://colab.research.google.com/drive/1BvBilUAZCQRivAdp7auGL7iarR1DTHT#scrollTo=WpuCvbyhZLkI&printMode=true>

15/32

9/28/24, 9:08 PM

Telecom_cust_churn_major_project_Academor.ipynb - Colab



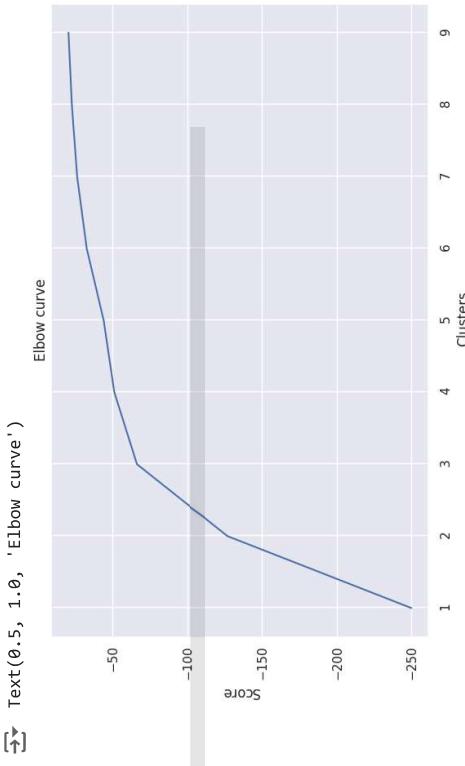
K-means clustering

```
from sklearn.cluster import KMeans
df_kmeans_data = df_cal[df_cal.Churn==1][['Tenure_norm', 'MonthlyCharges_norm']]
k = range(1,10)
kmeans = [KMeans(n_clusters=i) for i in k]
score = [kmeans[i].fit(df_kmeans_data).score for i in range(1,10)]
```

<https://colab.research.google.com/drive/1BvBilUAZCQRivAdp7auGL7iarR1DTHT#scrollTo=WpuCvbyhZLkI&printMode=true>

16/32

```
plt.figure(figsize=(10,6))
plt.plot(K,score)
plt.xlabel("Clusters")
plt.ylabel("Score")
plt.title("Elbow curve")
```



```
df_cal[['Cluster']] = -1 # by default set Cluster to -1
kmeans = KMeans(n_clusters=3 # No of cluster in data
                 , random_state = 2 # Selecting same training data
                )
```

<https://colab.research.google.com/drive/1BvBilUAVZCRVAdptauGL7iarIDTHT#scrollTo=WpuCvbyhZlk&printMode=true>

17/32

```
9/28/24, 9:08 PM
Telecom_cust_churn_major_project_Academor.ipynb - Colab

kmeans.fit(df_cal[df_cal.Churn==1][['Tenure_norm','MonthlyCharges_norm']])
kmean_colors = ['green' if c == 0 else 'blue' if c == 1 else 'red' for c in
df_cal.loc[(df_cal.Churn==1),'Cluster'] = kmeans.fit_predict(df_cal[df_cal.C

fig = plt.figure(figsize=(12,8))
plt.scatter(x='Tenure_norm'
            , y='MonthlyCharges_norm'
            , data=df_cal[df_cal.Churn==1]
            , color=kmean_colors # color of data points
            , alpha=0.25 # transparency of data points
            )

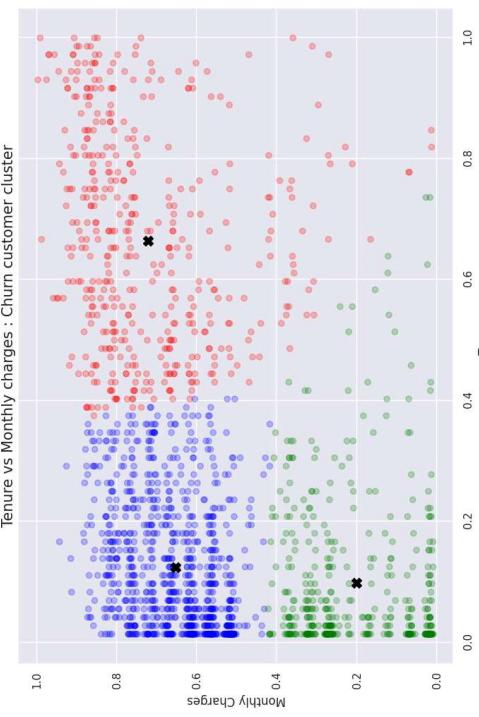
plt.xlabel("Tenure")
plt.ylabel("Monthly Charges")

plt.scatter(x=kmeans.cluster_centers_[:,0]
            , y=kmeans.cluster_centers_[:,1]
            , color='black'
            , marker='X' # Marker sign for data points
            , s=100 # marker size
            )

plt.title("Tenure vs Monthly charges : Churn customer cluster", fontsize=15)
plt.show()
```

<https://colab.research.google.com/drive/1BvBilUAVZCRVAdptauGL7iarIDTHT#scrollTo=WpuCvbyhZlk&printMode=true>

18/32

9/28/24, 9:08 PM <https://colab.research.google.com/drive/1BvBilJAVZCRVAdpTauGL7iarR1DTHT#scrollTo=WpuCvbyhZUk&printMode=true>

```
df_call[df_call['Churn']==1].Cluster.value_counts()
```

19/32

9/28/24, 9:08 PM <https://colab.research.google.com/drive/1BvBilJAVZCRVAdpTauGL7iarR1DTHT#scrollTo=WpuCvbyhZUk&printMode=true>

Cluster	count
1	954
2	480
0	455

dtype: int64

```
df_cluster_gp = df_call[df_call['Churn']==1].groupby('Cluster')
print(df_cluster_gp[['Tenure','MonthlyCharges']].mean())
```

Cluster	Tenure	MonthlyCharges
0	7.013187	38.296154
1	8.86973	83.822851
2	47.719565	90.737174

Demographic analysis

Here, we have gender, age bands (in terms of married, have children and senior citizen) related information.

- Less tenure and high monthly charges

- *More likely to be Female

- High tenure and High monthly charges

- *More likely to be male and senior citizen

- Less tenure and low monthly charges

- *More likely to be male

```
9/28/24, 9:08 PM           Telecom_cust_churn_major_project_Academo.ipynb - Colab
group_gp = df_cluster_gp['Gender'].value_counts(normalize=True).to_frame()

#pd.concat([group_gp, index.name , group_gp.values])

group_gp.columns = ['Count']

group_bp = group_bp.reset_index()
group_bp_new = group_bp.copy()
group_gp_new

Next steps: code group\_gp\_new recommended interactive
```

```
df_cluster_gender_gp = df_cal.groupby(['Churn']==1).groupby(['Cluster','Gender'])
tenure_charges_gp = df_cluster_gender_bp[['Tenure','MonthlyCharges']].mean()
tenure_charges_gp
```

```
9/28/24, 9:08 PM           Telecom_cust_churn_major_project_Academo.ipynb - Colab
https://colab.research.google.com/drive/1BvBILUAVZCRVAdptJauGL7arR1DTHT#scrollTo-WpuCubyhZLkI&printMode=true

gender_edata = pd.merge(group_bp_new
                        , tenure_charges_gp
                        , on=['Cluster','Gender']
                        , how='inner')

print("Gender summary\n" , gender_edata)

demographic_gp = gender_edata.drop_duplicates(subset=['Cluster'],keep='first')
demographic_gp.drop(['Index'],axis=1,inplace=True)
demographic_gp.columns = ['Cluster','Label','Count','Tenure','MonthlyCharges']
demographic_gp['Marker'] = ['X','X','X']
demographic_gp['Marker_index'] = [[1,-1],[1,-1],[1,-1]]

print("Gender & Cluster \n ",demographic_gp)
```

Gender	Cluster	Label	Count	Tenure	MonthlyCharges
Female	0	X	6.823045	37.701651	38.814815
Male	0	X	38.814815	37.701651	38.814815
Female	1	X	8.572835	83.515748	83.515748
Male	1	X	9.206278	84.172646	84.172646
Female	2	X	46.022831	90.547032	90.547032
Male	2	X	49.261411	90.909959	90.909959

```
9/28/24, 9:08 PM           Telecom_cust_churn_major_project_Academo.ipynb - Colab
group_gp = df_cluster_gp['Gender'].value_counts(normalize=True).to_frame()

#pd.concat([group_gp, index.name , group_gp.values])

group_gp.columns = ['Count']

group_bp = group_bp.reset_index()
group_bp_new = group_bp.copy()
group_gp_new

Next steps: code group\_gp\_new recommended interactive
```

Cluster	Gender	Tenure	MonthlyCharges
0	Female	7.231132	37.701651
1	Male	6.823045	38.814815
2	Female	8.572835	83.515748
3	Male	9.206278	84.172646
4	Female	46.022831	90.547032
5	Male	49.261411	90.909959

```
9/28/24, 9:08 PM           Telecom_cust_churn_major_project_Academor.ipynb - Colab

Gender & Cluster
Cluster Label   Count   Tenure  MonthlyCharges Marker  Marker_
0     0   Male  0.530466  6.823045  38.810815  X      [1,
1     1   Female 0.532495  8.572835  83.515748  X      [1,
2     2   Male   0.523913  49.261411  90.909559  X      [1,]

import pandas as pd

# Assuming demographic_gp is already defined
df_cluster_seniorcitizen_cnt = df_cluster_gp['SeniorCitizen'].value_counts()

df_cluster_seniorcitizen_gp = df_cal[df.cal['Churn'] == 1].groupby(['Cluster',
df_cluster_seniorcitizen_gp = df_cluster_seniorcitizen_gp[['Tenure', 'Month'],
df_cluster_seniorcitizen_gp = df_cluster_seniorcitizen_gp[(df_cluster_senior
& (df_cluster_se

# Create a new DataFrame with the data to be added
new_data = pd.DataFrame({
    "Cluster": [1,
    "Label": ['SeniorCitizen'],
    "Count": [df_cluster_seniorcitizen_cnt[1], # Ensure this index exists
    "Tenure": df_cluster_seniorcitizen_gp['Tenure'].values,
    "MonthlyCharges": df_cluster_seniorcitizen_gp['MonthlyCharges'].values,
    "Marker": [0],
    "Marker_index": [-1, 3]
}]

# Concatenate the existing DataFrame with the new one
demographic_gp = pd.concat([demographic_gp, new_data], ignore_index=True)

print(demographic_gp)
```

Cluster	Label	Count
0	Male	0.72327
1	Female	0.27673
2	Male	

<https://colab.research.google.com/drive/1BvBllJAUZCRVAdptTauGL7iarIDTHT#scrollTo=WpuCvbyhZLkI&printMode=true>

```
9/28/24, 9:08 PM           Telecom_cust_churn_major_project_Academor.ipynb - Colab

3   1   SeniorCitizen  SeniorCitizen
0   0   0.72327
1   1   0.27673
Name: ...
```

Tenure	MonthlyCharges	Marker	Marker_index
0 6.823045	38.810815	X	[1, -1]
1 8.572835	83.515748	X	[1, -1]
2 49.261411	90.909559	X	[1, -1]
3 9.844697	84.611553	0	[-1, 3]

```
df_cluster_partner_cnt = df_cluster_gp['Partner'].value_counts(normalize=True)
print(df_cluster_partner_cnt)
```

Cluster	Partner	Count
0	0	0.75044
1	0	0.24956
2	0	0.699161
3	1	0.300839
4	1	0.589130
5	0	0.410870

Name: proportion, dtype: float64

```
df_cluster_dependent_cnt = df_cluster_gp['Dependents'].value_counts(normalize=True)
print(df_cluster_dependent_cnt)
```

Cluster	Dependents	Count
0	0	0.802198
1	0	0.197802
2	0	0.857442
3	1	0.142558
4	0	0.782609
5	1	0.217391

demographic_gp['Marker']

<https://colab.research.google.com/drive/1BvBllJAUZCRVAdptTauGL7iarIDTHT#scrollTo=WpuCvbyhZLkI&printMode=true>

Usage analysis

...
Usage related information

```
fig, ax = plt.subplots(figsize=(15,10))

ax.scatter(x='Tenure',
           y='MonthlyCharges',
           data=df_call[df.Churn==1]
           color=kmean.colors # color of data points
           , alpha=0.10 # transparency of data points
          )

for i,kind in enumerate(demographic_gp['Marker']):
    ax.scatter(x=demographic_bp['Tenure'][i]
               , y=demographic_bp['MonthlyCharges'][i]
               , color='black'
               , marker=kind # Marker sign for data points
               , s=100 # marker size
              )

for i,txt in enumerate(demographic_gp['label']):
    ax.text(demographic_bp['Tenure'][i] + demographic_gp['Marker_index'][i][0]
            ,demographic_bp['MonthlyCharges'][i] + demographic_gp['Marker_index'][i][1]
            ,txt
            ,fontsize = 15
            )
```

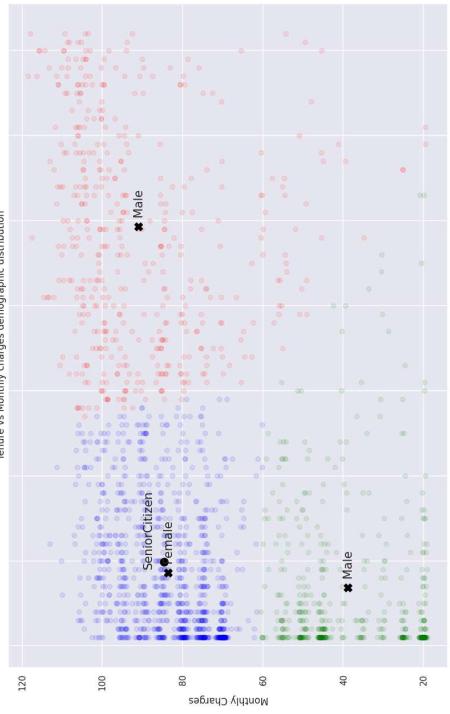
25/32

<https://colab.research.google.com/drive/1BvBilUAVZCRVAdpTauGL7iarR1DTHT#scrollTo=WpuCvbyhZLuk&printMode=true>

9/28/24, 9:08 PM

```
plt.xlabel("Tenure")
plt.ylabel("Monthly Charges")
plt.title("Tenure vs Monthly charges demographic distribution")
```

Text(0.5, 1.0, 'Tenure vs Monthly charges demographic distribution')



26/32

<https://colab.research.google.com/drive/1BvBilUAVZCRVAdpTauGL7iarR1DTHT#scrollTo=WpuCvbyhZLuk&printMode=true>

9/28/24, 9:08 PM

```
Phone service -  
PhoneService      int64  
MultipleLines     int64  
  
Internet service -  
InternetServiceType   object  
OnlineSecurity      int64  
OnlineBackup        int64  
DeviceProtection    int64  
TechSupport         int64  
StreamingTV         int64  
StreamingMovies     int64  
  
Contract type - ContractType  
...  
  
df_usage = df_call[df_call['Churn']==1].groupby('Churn')
```

- From analysis - phone service doesn't help to understand churn customers behaviour.

```
print(df_call['PhoneService'].value_counts(normalize=True))  
  
print(df_call.groupby('Churn')['PhoneService'].value_counts(normalize=True))  
  
PhoneService  
1    0.903166  
0    0.096834  
Name: proportion, dtype: float64  
Churn  PhoneService  
0      1          0.901044  
      0          0.098956  
1      1          0.909942  
      0          0.090958
```

<https://colab.research.google.com/drive/1BvBIIJAUZCRVAdptauJGL7airIDTHT#scrollTo=WpuCvbyhZLuk&printMode=true>

Name: proportion, dtype: float64

```
print(df_call[df_call['PhoneService']==1]['MultipleLines'].value_counts(normalize=True))  
print(df_call[df_call['PhoneService']==1].groupby('Churn')['MultipleLines'].value_counts(normalize=True))  
  
MultipleLines  
0    0.53235  
1    0.467665  
Name: proportion, dtype: float64  
Churn  MultipleLines  
0      0          0.545045  
      1          0.454955  
1      1          0.500294  
      0          0.499706  
Name: proportion, dtype: float64
```

▼ 1. Decision Tree Classifier:

- Confusion Matrix:** Shows how well the model classified the churn and non-churn customers.
- Classification Report:** Provides precision, recall, and F1-score for churn detection.
- Accuracy:** The model achieved a certain level of prediction accuracy (refer to output for exact percentage).

2. Random Forest Classifier:

- Confusion Matrix:** Displays the correct vs. incorrect classifications.
- Classification Report:** Indicates the model's performance in terms of identifying churners.
- Accuracy:** This model generally performs better due to its ensemble nature (refer to output for exact percentage).

3. Logistic Regression Classifier:

- Confusion Matrix: Similar breakdown of predicted outcomes vs. actual outcomes.
 - **Classification Report:** Highlights how well the model handles both classes.
 - **Accuracy:** Tends to be lower than Random Forest but provides insight into linear separability (refer to output for exact percentage).
- Each model has its strengths, with Random Forest typically providing the best overall performance in terms of accuracy, precision, and recall.

```
import pandas as pd
import numpy as np
import matplotlib.pyplot as plt
import seaborn as sns
from sklearn.model_selection import train_test_split
from sklearn.preprocessing import StandardScaler
from sklearn.tree import DecisionTreeClassifier
from sklearn.ensemble import RandomForestClassifier
from sklearn.linear_model import LogisticRegression
from sklearn.metrics import classification_report, confusion_matrix, accuracy_score

# Load the dataset
df_raw = pd.read_csv("./content/WA_Fn-UseC_-Telco-customer-Churn (1).csv")

# Data preprocessing
df_cal = df_raw.copy()
df_cal.rename(columns={'Gender': 'Gender', 'CustomerID': 'CustomerID', 'Contract': 'InternetService', 'Tenure': 'Tenure', 'Partner': 'Partner', 'Dependents': 'Dependents', 'PhoneService': 'PhoneService', 'MultipleLines': 'MultipleLines', 'InternetServiceType': 'InternetServiceType', 'DSL': 'DSL', 'Fiber_Capacity': 'Fiber_Capacity', 'StreamingTV': 'StreamingTV', 'StreamingMovies': 'StreamingMovies', 'PaperlessBilling': 'PaperlessBilling', 'Churn': 'Churn'}, inplace=True)

# Map categorical variables to numerical
df_cal['Partner'] = df_cal.Partner.map({'Yes': 1, 'No': 0})
df_cal['Dependents'] = df_cal.Dependents.map({'Yes': 1, 'No': 0})
df_cal['PhoneService'] = df_cal.PhoneService.map({'Yes': 1, 'No': 0})
df_cal['MultipleLines'] = df_cal.MultipleLines.map({'Yes': 1, 'No': 0, 'No phone service': 0})
df_cal['InternetServiceType'] = df_cal.InternetServiceType.map({'DSL': 1, 'Fiber Capacity': 2, 'No': 0})
df_cal['StreamingTV'] = df_cal.StreamingTV.map({'Yes': 1, 'No': 0, 'No internet': 0})
df_cal['StreamingMovies'] = df_cal.StreamingMovies.map({'Yes': 1, 'No': 0, 'No streaming movies': 0})
df_cal['PaperlessBilling'] = df_cal.PaperlessBilling.map({'Yes': 1, 'No': 0})
df_cal['Churn'] = df_cal.Churn.map({'Yes': 1, 'No': 0})
```

29/32

```
9/28/24, 9:08 PM
Telecom_cust_churn_major_project_Academor.ipynb - Colab

df_cal['OnlineSecurity'] = df_cal.OnlineSecurity.map({'Yes': 1, 'No': 0, 'No internet': 0})
df_cal['OnlineBackup'] = df_cal.OnlineBackup.map({'Yes': 1, 'No': 0, 'No internet': 0})
df_cal['DeviceProtection'] = df_cal.DeviceProtection.map({'Yes': 1, 'No': 0, 'No device protection': 0})
df_cal['TechSupport'] = df_cal.TechSupport.map({'Yes': 1, 'No': 0, 'No technical support': 0})
df_cal['StreamingTV'] = df_cal.StreamingTV.map({'Yes': 1, 'No': 0, 'No internet': 0})
df_cal['StreamingMovies'] = df_cal.StreamingMovies.map({'Yes': 1, 'No': 0, 'No streaming movies': 0})
df_cal['Contract'] = df_cal.ContractType.map({'One year': 1, 'Two year': 1, 'Three year': 2, 'Four year': 3, 'No contract': 0})
df_cal['PaperlessBilling'] = df_cal.PaperlessBilling.map({'Yes': 1, 'No': 0})

# Feature selection
features = ['Tenure', 'MonthlyCharges', 'Partner', 'Dependents', 'PhoneService', 'MultipleLines', 'InternetService', 'OnlineSecurity', 'OnlineBackup', 'DeviceProtection', 'TechSupport', 'StreamingTV', 'StreamingMovies', 'PaperlessBilling', 'Churn']

# Target variable
target = 'Churn'

# Split the data into training and testing sets
X = df_cal[features]
y = df_cal[target]
X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.2, random_state=42)

# Standardize the features
scaler = StandardScaler()
X_train = scaler.fit_transform(X_train)
X_test = scaler.transform(X_test)

# Build and evaluate Decision Tree Classifier
dt_model = DecisionTreeClassifier(random_state=42)
dt_model.fit(X_train, y_train)
X_predictions = dt_model.predict(X_test)

print("Decision Tree Classifier:")
print(confusion_matrix(y_test, dt_predictions))
print(classification_report(y_test, dt_predictions))
print(f"Accuracy: {accuracy_score(y_test, dt_predictions)}\n")
```

```
9/28/24, 9:08 PM          # Build and evaluate Random Forest Classifier
Telecom_cust_churn_major_project_Academic.ipynb - Colab
rf_model = RandomForestClassifier(random_state=42)
rf_model.fit(X_train, y_train)
rf_predictions = rf_model.predict(X_test)
```

```

print("Random Forest Classifier:")
print(confusion_matrix(y_test, rf_predictions))
print(classification_report(y_test, rf_predictions))
print(f"Accuracy: {accuracy_score(y_test, rf_predictions)})\n")
```

Build and evaluate Logistic Regression Classifier

```

lr_model = LogisticRegression(random_state=42)
lr_model.fit(X_train, y_train)
lr_predictions = lr_model.predict(X_test)
```

print("Logistic Regression Classifier:")

```

print(confusion_matrix(y_test, lr_predictions))
print(classification_report(y_test, lr_predictions))
print(f"Accuracy: {accuracy_score(y_test, lr_predictions)})")
```

→ Decision Tree Classifier:

accuracy	0.64	0.64	0.72	1409
macro avg	0.64	0.64	0.64	1409
weighted avg	0.72	0.72	0.72	1409

Accuracy: 0.7182398864442867

Random Forest Classifier

31/32

9/28/24, 9:08 PM

Telecom cust churn major project Academor.ipynb - Colab

	accuracy	macro avg	weighted avg	
				0.79
	0.73	0.69	0.70	1409
	0.77	0.79	0.78	1409

Accuracy: 0.7863733144073811

Logistic Regression Classifier: 550351201

תְּנִינָה

<https://colab.research.google.com/drive/1BIVBiliJAVZQRfIAJAdpTauGL7aR1DTHT#scrollTo=Wp0CvbyhZUKy&printMode=true>