```python
import numpy as np
import pandas as pd
import matplotlib.pyplot as plt
import seaborn as sns
from sklearn.preprocessing import LabelEncoder
from imblearn.over_sampling import SMOTE
from sklearn.model_selection import train_test_split, cross_val_score
from sklearn.tree import DecisionTreeClassifier
from sklearn.ensemble import RandomForestClassifier
from xgboost import XGBClassifier
from sklearn.metrics import accuracy_score, confusion_matrix, classification_report
import pickle
```

```python
# load teh csv data to a pandas dataframe
df = pd.read_csv(r'/content/WA_Fn-UseC_-Telco-Customer-Churn.csv')
```

```python
df.shape
```

    (7043, 21)

```python
df.head()
```

|   | customerID | gender | SeniorCitizen | Partner | Dependents | tenure | PhoneService | MultipleLines | InternetService | On |
|---|---|---|---|---|---|---|---|---|---|---|
| 0 | 7590-VHVEG | Female | 0 | Yes | No | 1 | No | No phone service | DSL | |
| 1 | 5575-GNVDE | Male | 0 | No | No | 34 | Yes | No | DSL | |
| 2 | 3668-QPYBK | Male | 0 | No | No | 2 | Yes | No | DSL | |
| 3 | 7795-CFOCW | Male | 0 | No | No | 45 | No | No phone service | DSL | |
| 4 | 9237-HQITU | Female | 0 | No | No | 2 | Yes | No | Fiber optic | |

5 rows × 21 columns

```python
pd.set_option("display.max_columns", None)
```

```python
df.head(2)
```

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

```python
df.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    object
 18  MonthlyCharges     7043 non-null    float64
 19  TotalCharges       7043 non-null    object
 20  Churn              7043 non-null    object
dtypes: float64(1), int64(2), object(18)
memory usage: 1.1+ MB
```

```python
df = df.drop(columns=["customerID"])
```

```python
df.head(2)
```

| | gender | SeniorCitizen | Partner | Dependents | tenure | PhoneService | MultipleLines | InternetService | OnlineSecurity |
|---|---|---|---|---|---|---|---|---|---|
| 0 | Female | 0 | Yes | No | 1 | No | No phone service | DSL | No |
| 1 | Male | 0 | No | No | 34 | Yes | No | DSL | Yes |

Next steps:   Generate code with `df`   ⬤ View recommended plots   New interactive sheet

```python
df.columns
```

```
Index(['gender', 'SeniorCitizen', 'Partner', 'Dependents', 'tenure',
       'PhoneService', 'MultipleLines', 'InternetService', 'OnlineSecurity',
       'OnlineBackup', 'DeviceProtection', 'TechSupport', 'StreamingTV',
       'StreamingMovies', 'Contract', 'PaperlessBilling', 'PaymentMethod',
       'MonthlyCharges', 'TotalCharges', 'Churn'],
      dtype='object')
```

```python
print(df["gender"].unique())
```

```
['Female' 'Male']
```

```python
print(df["SeniorCitizen"].unique())
```

```
[0 1]
```

```python
# printing the unique values in all the columns

numerical_features_list = ["tenure", "MonthlyCharges", "TotalCharges"]

for col in df.columns:
  if col not in numerical_features_list:
    print(col, df[col].unique())
    print("-"*50)
```

```
gender ['Female' 'Male']
--------------------------------------------------
SeniorCitizen [0 1]
--------------------------------------------------
Partner ['Yes' 'No']
--------------------------------------------------
Dependents ['No' 'Yes']
--------------------------------------------------
PhoneService ['No' 'Yes']
--------------------------------------------------
MultipleLines ['No phone service' 'No' 'Yes']
--------------------------------------------------
InternetService ['DSL' 'Fiber optic' 'No']
--------------------------------------------------
OnlineSecurity ['No' 'Yes' 'No internet service']
```

```
--------------------------------------------------
OnlineBackup ['Yes' 'No' 'No internet service']
--------------------------------------------------
DeviceProtection ['No' 'Yes' 'No internet service']
--------------------------------------------------
TechSupport ['No' 'Yes' 'No internet service']
--------------------------------------------------
StreamingTV ['No' 'Yes' 'No internet service']
--------------------------------------------------
StreamingMovies ['No' 'Yes' 'No internet service']
--------------------------------------------------
Contract ['Month-to-month' 'One year' 'Two year']
--------------------------------------------------
PaperlessBilling ['Yes' 'No']
--------------------------------------------------
PaymentMethod ['Electronic check' 'Mailed check' 'Bank transfer (automatic)'
 'Credit card (automatic)']
--------------------------------------------------
Churn ['No' 'Yes']
--------------------------------------------------
```

```python
print(df.isnull().sum())
```

```
gender              0
SeniorCitizen       0
Partner             0
Dependents          0
tenure              0
PhoneService        0
MultipleLines       0
InternetService     0
OnlineSecurity      0
OnlineBackup        0
DeviceProtection    0
TechSupport         0
StreamingTV         0
StreamingMovies     0
Contract            0
PaperlessBilling    0
PaymentMethod       0
MonthlyCharges      0
TotalCharges        0
Churn               0
dtype: int64
```

```python
df[df["TotalCharges"]==" "]
```

| | gender | SeniorCitizen | Partner | Dependents | tenure | PhoneService | MultipleLines | InternetService | OnlineSecur |
|------|--------|---------------|---------|------------|--------|--------------|---------------|-----------------|-------------|
| 488 | Female | 0 | Yes | Yes | 0 | No | No phone service | DSL | |
| 753 | Male | 0 | No | Yes | 0 | Yes | No | No | No inte ser |
| 936 | Female | 0 | Yes | Yes | 0 | Yes | No | DSL | |
| 1082 | Male | 0 | Yes | Yes | 0 | Yes | Yes | No | No inte ser |
| 1340 | Female | 0 | Yes | Yes | 0 | No | No phone service | DSL | |
| 3331 | Male | 0 | Yes | Yes | 0 | Yes | No | No | No inte ser |
| 3826 | Male | 0 | Yes | Yes | 0 | Yes | Yes | No | No inte ser |
| 4380 | Female | 0 | Yes | Yes | 0 | Yes | No | No | No inte ser |
| 5218 | Male | 0 | Yes | Yes | 0 | Yes | No | No | No inte ser |
| 6670 | Female | 0 | Yes | Yes | 0 | Yes | Yes | DSL | |
| 6754 | Male | 0 | No | Yes | 0 | Yes | Yes | DSL | |

```python
len(df[df["TotalCharges"]==" "])
```

```
11
```

```python
df["TotalCharges"] = df["TotalCharges"].replace({" ": "0.0"})
```

```python
df["TotalCharges"] = df["TotalCharges"].astype(float)
```

```python
df.info()
```

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

```python
# checking the class distribution of target column
print(df["Churn"].value_counts())
```

```
Churn
No      5174
Yes     1869
Name: count, dtype: int64
```

```
df.shape
```

```
(7043, 20)
```

```
df.columns
```

```
Index(['gender', 'SeniorCitizen', 'Partner', 'Dependents', 'tenure',
       'PhoneService', 'MultipleLines', 'InternetService', 'OnlineSecurity',
       'OnlineBackup', 'DeviceProtection', 'TechSupport', 'StreamingTV',
       'StreamingMovies', 'Contract', 'PaperlessBilling', 'PaymentMethod',
       'MonthlyCharges', 'TotalCharges', 'Churn'],
      dtype='object')
```

```
df.head(2)
```

| | gender | SeniorCitizen | Partner | Dependents | tenure | PhoneService | MultipleLines | InternetService | OnlineSecurity |
|---|---|---|---|---|---|---|---|---|---|
| 0 | Female | 0 | Yes | No | 1 | No | No phone service | DSL | No |
| 1 | Male | 0 | No | No | 34 | Yes | No | DSL | Yes |

Next steps:    Generate code with `df`      ◉ View recommended plots      New interactive sheet

```
df.describe()
```

| | SeniorCitizen | tenure | MonthlyCharges | TotalCharges |
|---|---|---|---|---|
| count | 7043.000000 | 7043.000000 | 7043.000000 | 7043.000000 |
| mean | 0.162147 | 32.371149 | 64.761692 | 2279.734304 |
| std | 0.368612 | 24.559481 | 30.090047 | 2266.794470 |
| min | 0.000000 | 0.000000 | 18.250000 | 0.000000 |
| 25% | 0.000000 | 9.000000 | 35.500000 | 398.550000 |
| 50% | 0.000000 | 29.000000 | 70.350000 | 1394.550000 |
| 75% | 0.000000 | 55.000000 | 89.850000 | 3786.600000 |
| max | 1.000000 | 72.000000 | 118.750000 | 8684.800000 |

```python
def plot_histogram(df, column_name):

  plt.figure(figsize=(5, 3))
  sns.histplot(df[column_name], kde=True)
  plt.title(f"Distribution of {column_name}")

  # calculate the mean and median values for the columns
  col_mean = df[column_name].mean()
  col_median = df[column_name].median()

  # add vertical lines for mean and median
  plt.axvline(col_mean, color="red", linestyle="--", label="Mean")
  plt.axvline(col_median, color="green", linestyle="-", label="Median")

  plt.legend()

  plt.show()


plot_histogram(df, "tenure")
```

Distribution of tenure

```python
plot_histogram(df, "MonthlyCharges")
```



Distribution of MonthlyCharges
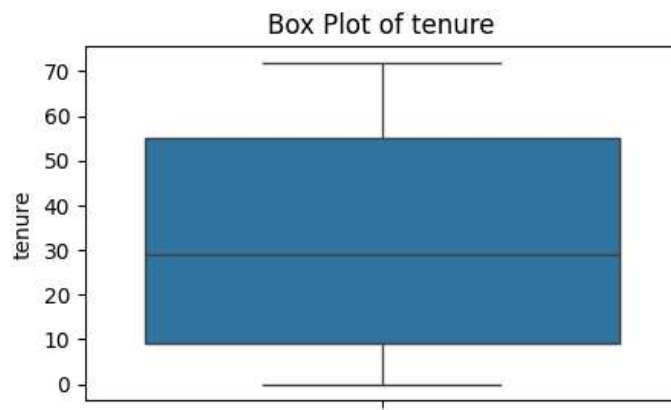
```python
plot_histogram(df, "TotalCharges")
```
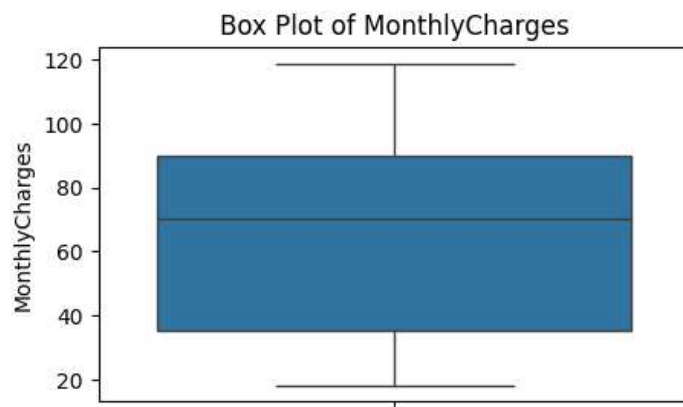


Distribution of TotalCharges

```python
def plot_boxplot(df, column_name):

  plt.figure(figsize=(5, 3))
  sns.boxplot(y=df[column_name])
  plt.title(f"Box Plot of {column_name}")
  plt.ylabel(column_name)
  plt.show


plot_boxplot(df, "tenure")
```
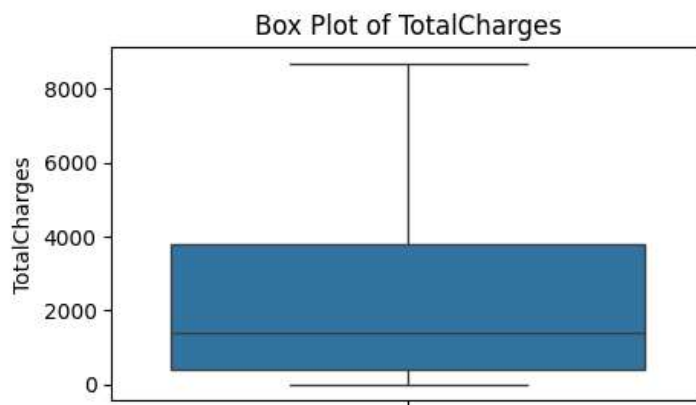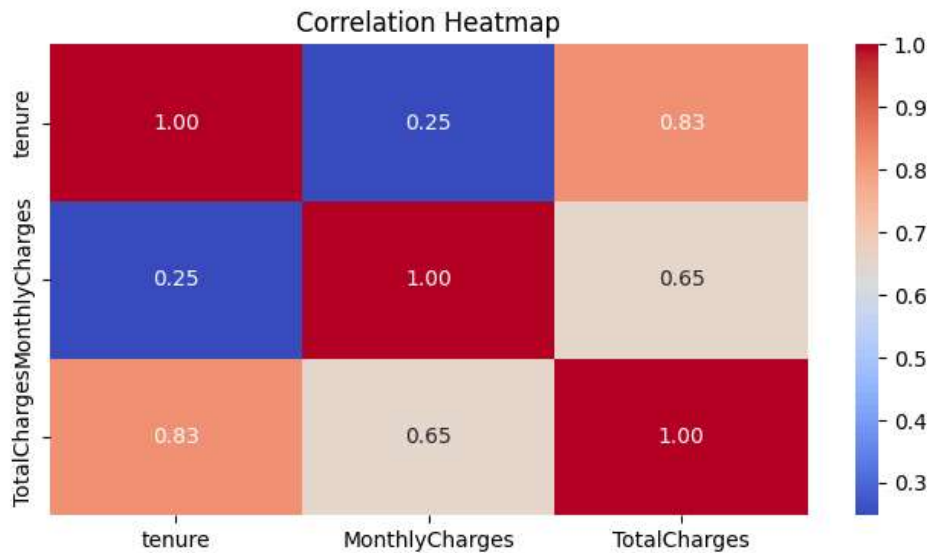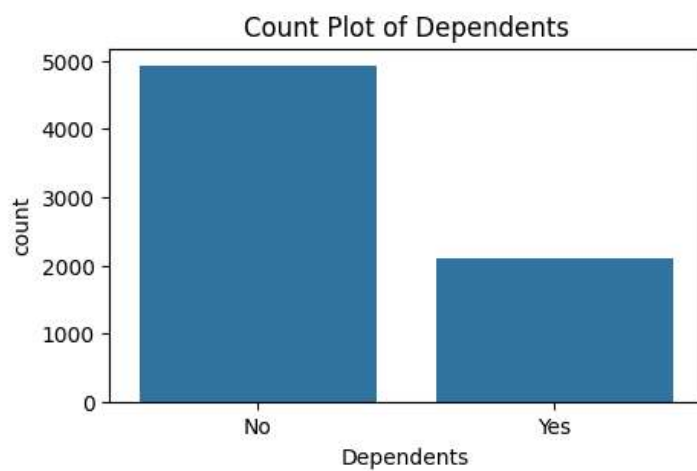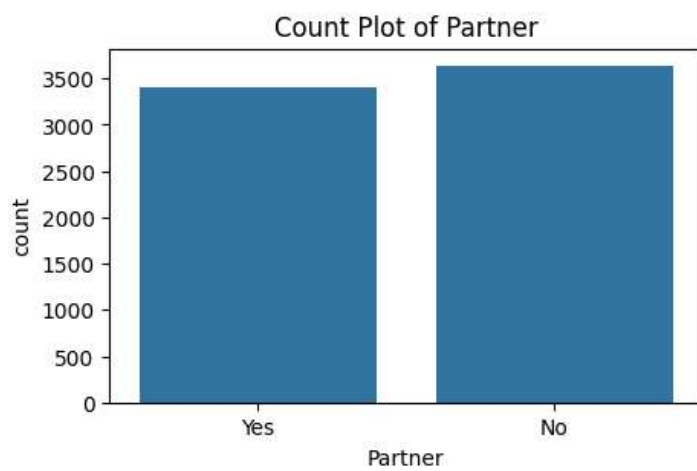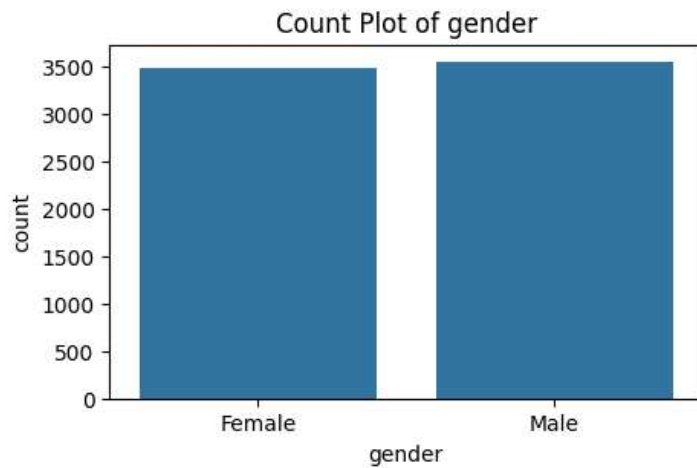
Box Plot of tenure

```
plot_boxplot(df, "MonthlyCharges")
```



Box Plot of MonthlyCharges

```
plot_boxplot(df, "TotalCharges")
```



Box Plot of TotalCharges

```
# correlation matrix - heatmap
plt.figure(figsize=(8, 4))
sns.heatmap(df[["tenure", "MonthlyCharges", "TotalCharges"]].corr(), annot=True, cmap="coolwarm", fmt=".2f")
plt.title("Correlation Heatmap")
plt.show()
```

## Correlation Heatmap



```
df.columns
```

```
Index(['gender', 'SeniorCitizen', 'Partner', 'Dependents', 'tenure',
       'PhoneService', 'MultipleLines', 'InternetService', 'OnlineSecurity',
       'OnlineBackup', 'DeviceProtection', 'TechSupport', 'StreamingTV',
       'StreamingMovies', 'Contract', 'PaperlessBilling', 'PaymentMethod',
       'MonthlyCharges', 'TotalCharges', 'Churn'],
      dtype='object')
```

```
df.info()
```

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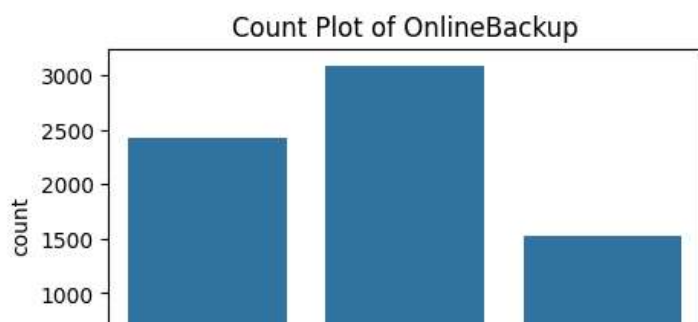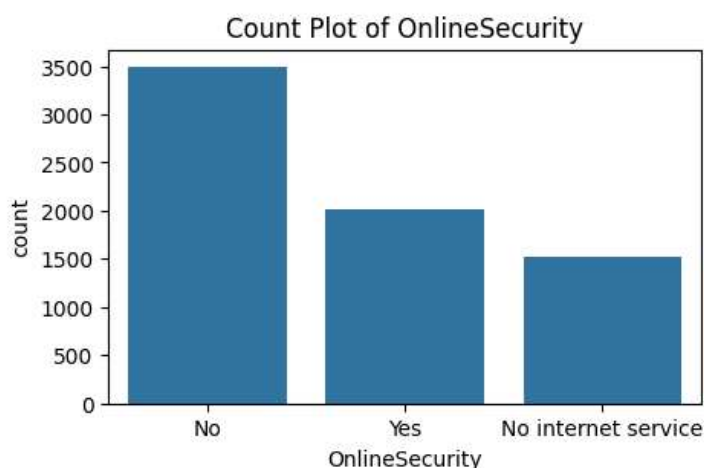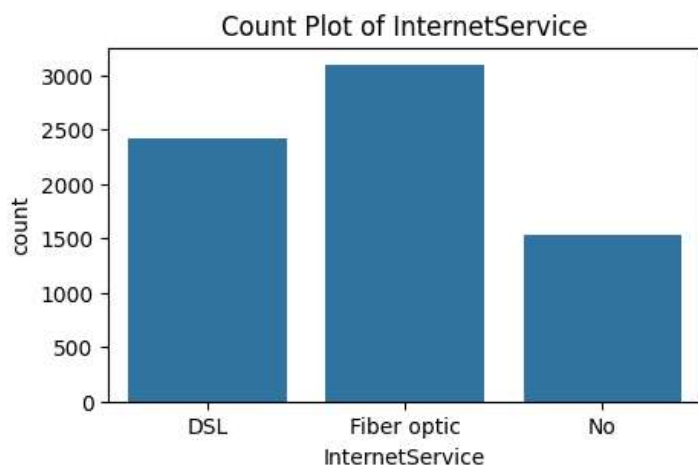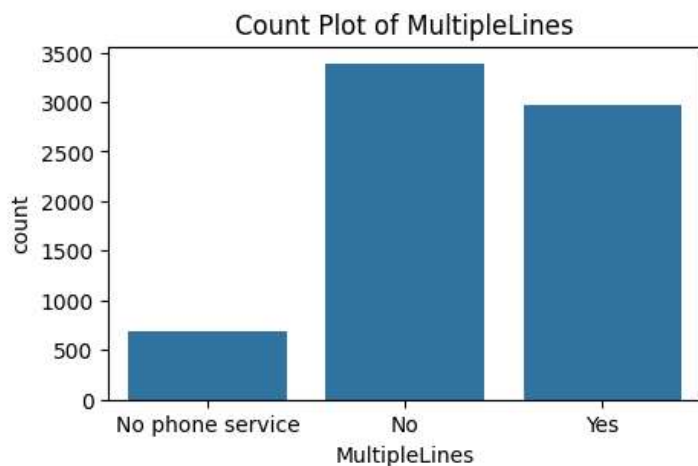
```python
object_cols = df.select_dtypes(include="object").columns.to_list()

object_cols = ["SeniorCitizen"] + object_cols

for col in object_cols:
  plt.figure(figsize=(5, 3))
  sns.countplot(x=df[col])
  plt.title(f"Count Plot of {col}")
  plt.show()
```
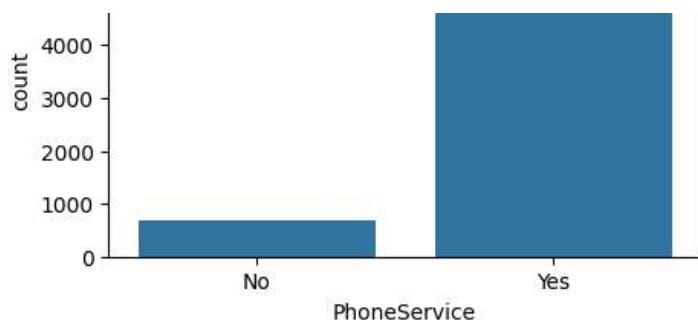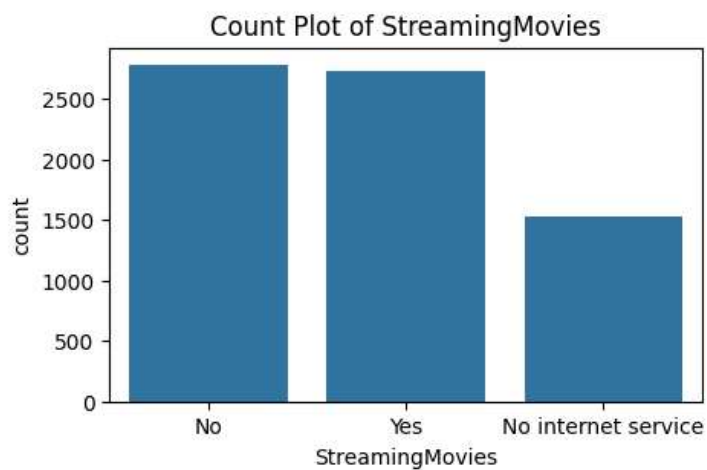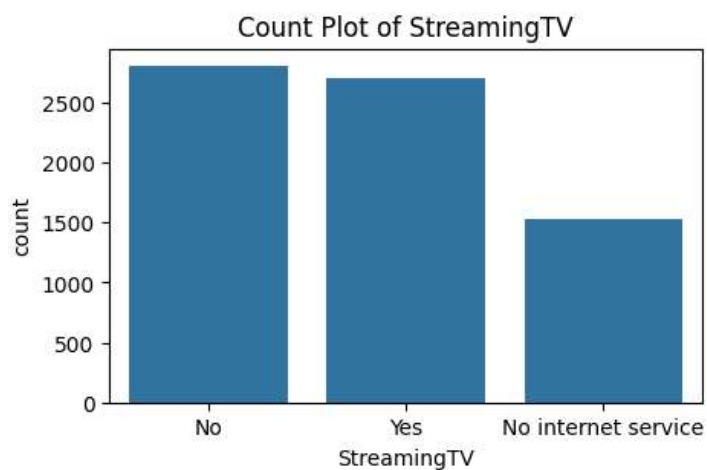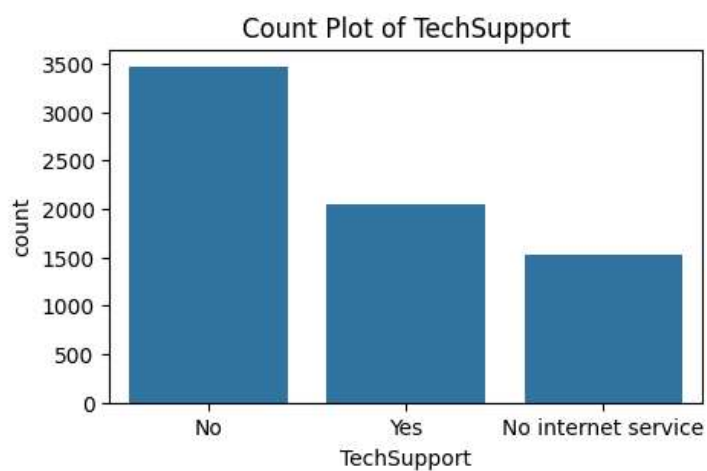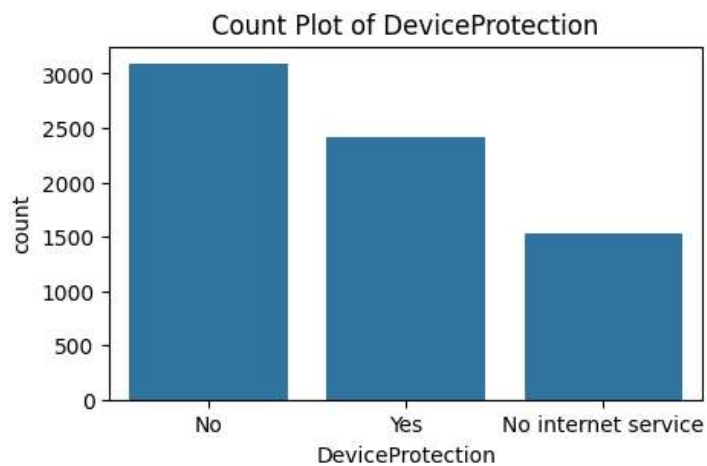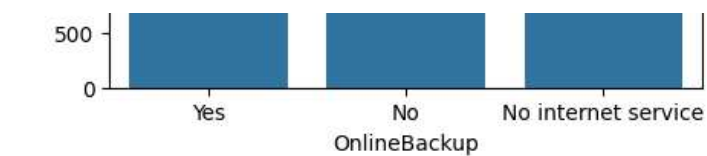
## Count Plot of SeniorCitizen



## Count Plot of gender



## Count Plot of Partner



## Count Plot of Dependents



## Count Plot of PhoneService

## Count Plot of MultipleLines



## Count Plot of InternetService



## Count Plot of OnlineSecurity



## Count Plot of OnlineBackup

OnlineBackup

## Count Plot of DeviceProtection



DeviceProtection

## Count Plot of TechSupport



TechSupport

## Count Plot of StreamingTV



StreamingTV

## Count Plot of StreamingMovies



StreamingMovies

### Count Plot of Contract



### Count Plot of PaperlessBilling



### Count Plot of PaymentMethod



### Count Plot of Churn

```
df.head(3)
```

|   | gender | SeniorCitizen | Partner | Dependents | tenure | PhoneService | MultipleLines | InternetService | OnlineSecurity |
|---|--------|---------------|---------|------------|--------|--------------|---------------|-----------------|----------------|
| 0 | Female | 0 | Yes | No | 1 | No | No phone service | DSL | No |
| 1 | Male | 0 | No | No | 34 | Yes | No | DSL | Yes |
| 2 | Male | 0 | No | No | 2 | Yes | No | DSL | Yes |

Next steps:  **Generate code with** `df`    ● **View recommended plots**    **New interactive sheet**

```
df["Churn"] = df["Churn"].replace({"Yes": 1, "No": 0})
```

```
<ipython-input-100-b6eb27bc3ee0>:1: FutureWarning: Downcasting behavior in `replace` is deprecated and will be rem
  df["Churn"] = df["Churn"].replace({"Yes": 1, "No": 0})
```

```
df.head(3)
```

|   | gender | SeniorCitizen | Partner | Dependents | tenure | PhoneService | MultipleLines | InternetService | OnlineSecurity |
|---|--------|---------------|---------|------------|--------|--------------|---------------|-----------------|----------------|
| 0 | Female | 0 | Yes | No | 1 | No | No phone service | DSL | No |
| 1 | Male | 0 | No | No | 34 | Yes | No | DSL | Yes |
| 2 | Male | 0 | No | No | 2 | Yes | No | DSL | Yes |

Next steps:  **Generate code with** `df`    ● **View recommended plots**    **New interactive sheet**

```
print(df["Churn"].value_counts())
```

```
Churn
0    5174
1    1869
Name: count, dtype: int64
```

```
# identifying columns with object data type
object_columns = df.select_dtypes(include="object").columns
```

```
print(object_columns)
```

```
Index(['gender', 'Partner', 'Dependents', 'PhoneService', 'MultipleLines',
       'InternetService', 'OnlineSecurity', 'OnlineBackup', 'DeviceProtection',
       'TechSupport', 'StreamingTV', 'StreamingMovies', 'Contract',
       'PaperlessBilling', 'PaymentMethod'],
      dtype='object')
```

```
# initialize a dictionary to save the encoders
encoders = {}

# apply label encoding and store the encoders
for column in object_columns:
  label_encoder = LabelEncoder()
  df[column] = label_encoder.fit_transform(df[column])
  encoders[column] = label_encoder

# save the encoders to a pickle file
with open("encoders.pkl", "wb") as f:
  pickle.dump(encoders, f)
```