

```
In [1]: import pandas as pd
import numpy as np

Customer churn prediction is to measure why customers are leaving a business.

In [2]: df = pd.read_csv("customer_churn.csv")

In [4]: df.sample(5)

Out[4]:
   customerID  gender  SeniorCitizen  Partner  Dependents  tenure  PhoneService  MultipleLines  InternetService  OnlineSecurity  ...  DeviceProtection  TechSupport  StreamingTV  StreamingMovies
3252    2153  FEMALE      0             0      Yes      64         Yes              No              No              No Internet service  ...  No Internet service  No Internet service  No Internet service
6379    1240  FEMALE      0             0      Yes      34         Yes              No              No              No Internet service  ...  No Internet service  No Internet service  No Internet service
659      205  MALE      0             0      Yes      37         Yes              No              No              No Internet service  ...  No Internet service  No Internet service  No Internet service
2124    7802  FEMALE      0             0      Yes      5         No              No phone service  DSL              No ...              No              No              No
642      6970  FEMALE      0             0      No      1         Yes              No              Fiber optic  No ...              No              No              Yes
5 rows x 21 columns

In [5]: df.shape
Out[5]:
(7843, 21)

First of all, drop customerID column as it is of no use

In [6]: df.drop('customerID', axis='columns', inplace=True)

In [7]: df

Out[7]:
   gender  SeniorCitizen  Partner  Dependents  tenure  PhoneService  MultipleLines  InternetService  OnlineSecurity  OnlineBackup  DeviceProtection  TechSupport  StreamingTV  StreamingMovies
5551    Male            0      No          No      20         Yes              No              No              No Internet service  No Internet service  No Internet service  No Internet service
3023    Male            0      Yes        Yes      72         Yes              Yes              DSL              Yes              Yes              Yes              Yes              Yes
2227    Male            0      Yes        Yes      58         Yes              No              No              No Internet service  No Internet service  No Internet service  No Internet service  No Internet service
1061    Male            1      Yes         No      46         Yes              Yes              Fiber optic  No              No              Yes              No              Yes
2225    Male            0      No          No      33         No              No phone service  DSL              No              Yes              No              No              Yes
5 rows x 21 columns

In [9]: df.dtypes

Out[9]:
gender                object
SeniorCitizen         int64
Partner               object
Dependents            object
tenure                int64
PhoneService          object
MultipleLines         object
InternetService       object
OnlineSecurity        object
OnlineBackup          object
DeviceProtection      object
TechSupport           object
StreamingTV           object
StreamingMovies       object
PaperlessBilling      object
PaymentMethod         object
MonthlyCharges        float64
TotalCharges          object
Churn                 object
dtype: object

Quick glance at above makes us realize that TotalCharges should be float but it is an object. Let's check what's going on with this column

In [10]: df.TotalCharges.values
Out[10]:
array(['129.85', '1889.5', '188.15', ..., '346.45', '386.6', '6844.5'],
      dtype=object)

it is string. Lets convert it to numbers

In [11]: pd.to_numeric(df.TotalCharges)

-----
Traceback (most recent call last)
~\Anaconda3\lib\site-packages\pandas\libs\lib.pyx in pandas._libs.lib.maybe_convert_numeric()
ValueError: unable to parse string " "
During handling of the above exception, another exception occurred:
Traceback (most recent call last)
~\AppData\Local\Temp\ipykernel_18012\2112264836.py in <module>
----> 1 pd.to_numeric(df.TotalCharges)
~\Anaconda3\lib\site-packages\pandas\core\tools\numeric.py in to_numeric(arg, errors, downcast)
181     if isinstance(numeric = errors not in ("ignore", "raise"))
182     try:
-> 183         values, _ = lib.maybe_convert_numeric(
184             values, set(), coerce_numeric=coerce_numeric
185         )
~\Anaconda3\lib\site-packages\pandas\libs\lib.pyx in pandas._libs.lib.maybe_convert_numeric()
ValueError: unable to parse string " " at position 488

some values seems to be not numbers but blank string. Let's find out such rows

In [13]: df[pd.to_numeric(df.TotalCharges, errors='coerce').isnull()]

Out[13]:
   gender  SeniorCitizen  Partner  Dependents  tenure  PhoneService  MultipleLines  InternetService  OnlineSecurity  OnlineBackup  DeviceProtection  TechSupport  StreamingTV  StreamingMovies
488  Female            0      Yes         No      0         No              No phone service  DSL              Yes              No              Yes              Yes              Yes
793   Male            0      No          Yes      0         Yes              No              No              No Internet service  No Internet service  No Internet service  No Internet service  No Internet service
936  Female            0      Yes         Yes      0         Yes              No              DSL              Yes              Yes              Yes              Yes              Yes
1062  Male            0      Yes        Yes      0         Yes              Yes              No              No Internet service  No Internet service  No Internet service  No Internet service  No Internet service
1340  Female            0      Yes        Yes      0         No              No phone service  DSL              Yes              Yes              Yes              Yes              Yes
3321  Male            0      Yes         Yes      0         Yes              No              No              No Internet service  No Internet service  No Internet service  No Internet service  No Internet service
3626  Male            0      Yes        Yes      0         Yes              Yes              No              No Internet service  No Internet service  No Internet service  No Internet service  No Internet service
4390  Female            0      Yes        Yes      0         Yes              No              No              No Internet service  No Internet service  No Internet service  No Internet service  No Internet service
5218  Male            0      Yes         Yes      0         Yes              No              No              No Internet service  No Internet service  No Internet service  No Internet service  No Internet service
6670  Female            0      Yes         Yes      0         Yes              DSL              No              Yes              Yes              Yes              Yes              Yes
6754  Male            0      No          Yes      0         Yes              Yes              DSL              Yes              Yes              No              Yes              No
5 rows x 21 columns

In [14]: df.iloc[488].TotalCharges
Out[14]:
''

Remove rows with space in TotalCharges

In [15]: df1 = df[df.TotalCharges != '']

In [16]: df1.TotalCharges = pd.to_numeric(df1.TotalCharges)

C:\Users\user\Anaconda3\lib\site-packages\pandas\core\generic.py:5518: SettingWithCopyWarning:
A value is being stored to a variable that is a copy of a slice from a DataFrame.
Try using .loc[row_indexer,col_indexer] = value instead
See the caveats in the documentation: https://pandas.pydata.org/pandas-docs/stable/user_guide/indexing.html#returning-a-view-versus-a-copy
self[name] = value

In [17]: df1[df1.Churn == 'No']

Out[17]:
   gender  SeniorCitizen  Partner  Dependents  tenure  PhoneService  MultipleLines  InternetService  OnlineSecurity  OnlineBackup  DeviceProtection  TechSupport  StreamingTV  StreamingMovies
0  Female            0      Yes         No      1         No              No phone service  DSL              No              Yes              No              No              No
1  Male            0      No          Yes      0         Yes              No              No              No Internet service  No Internet service  No Internet service  No Internet service  No Internet service
3  Male            0      No          No      45         No              No phone service  DSL              Yes              No              Yes              Yes              No
6  Male            0      No          Yes      22         Yes              Yes              Fiber optic  No              Yes              No              No              Yes
7  Female            0      No          No      10         No              No phone service  DSL              Yes              No              No              No              No
7037  Female            0      No          No      72         Yes              No              No              No Internet service  No Internet service  No Internet service  No Internet service  No Internet service
7038  Male            0      Yes        Yes      24         Yes              Yes              DSL              Yes              No              Yes              Yes              Yes
7039  Female            0      Yes        Yes      72         Yes              Yes              Fiber optic  No              Yes              Yes              No              Yes
7040  Female            0      Yes         Yes      11         No              No phone service  DSL              Yes              No              No              No              No
7042  Male            0      No          No      66         Yes              No              Fiber optic  Yes              No              Yes              Yes              Yes
5163 rows x 20 columns

Many of the columns are yes, no etc. Let's print unique values in object columns to see data values

In [94]: def print_unique(df1):
for column in df1:
    if df1[column].dtypes == 'object':
        print(f'{column}: {df1[column].unique()}')

In [91]: print_unique(df1)

gender: ['Female' 'Male']
PhoneService: ['No' 'Yes']
MultipleLines: ['No' 'Yes']
InternetService: ['DSL' 'Fiber optic' 'No']
OnlineSecurity: ['No' 'Yes']
OnlineBackup: ['Yes' 'No']
DeviceProtection: ['No' 'Yes']
TechSupport: ['No' 'Yes']
StreamingTV: ['No' 'Yes']
StreamingMovies: ['No' 'Yes']
Contract: ['Month-to-month' 'One year' 'Two year']
PaymentMethod: ['Electronic check' 'Mailed check' 'Bank transfer (automatic)']
Credit card (automatic): ['No' 'Yes']
Churn: ['No' 'Yes']

Some of the columns have no internet service or no phone service, that can be replaced with a simple No

In [94]: df1.replace('No Internet service', 'No', inplace=True)
df1.replace('No phone service', 'No', inplace=True)

In [96]: print_unique(df1)

gender: ['Female' 'Male']
PhoneService: ['No' 'Yes']
MultipleLines: ['No' 'Yes']
InternetService: ['DSL' 'Fiber optic' 'No']
OnlineSecurity: ['No' 'Yes']
OnlineBackup: ['Yes' 'No']
DeviceProtection: ['No' 'Yes']
TechSupport: ['No' 'Yes']
StreamingTV: ['No' 'Yes']
StreamingMovies: ['No' 'Yes']
Contract: ['Month-to-month' 'One year' 'Two year']
PaymentMethod: ['Electronic check' 'Mailed check' 'Bank transfer (automatic)']
Credit card (automatic): ['No' 'Yes']
Churn: ['No' 'Yes']

Convert Yes and No to 1 or 0

In [96]: yes_no_column = ['Partner', 'Dependents', 'PaperlessBilling']
for col in yes_no_column:
    df1[col].replace({'Yes': 1, 'No': 0}, inplace=True)

In [38]: df1.dtypes

Out[38]:
gender                object
SeniorCitizen         int64
Partner               int64
Dependents            int64
tenure                int64
PhoneService          object
MultipleLines         object
InternetService       object
OnlineSecurity        object
OnlineBackup          object
DeviceProtection      object
TechSupport           object
StreamingTV           object
StreamingMovies       object
Contract              object
PaperlessBilling      object
PaymentMethod         object
MonthlyCharges        float64
TotalCharges          float64
Churn                 object
dtype: object

Applying get dummies on categorical Variables

In [39]: df2 = pd.get_dummies(data=df1, columns=['InternetService', 'Contract', 'PaymentMethod'])

In [40]: df2

Out[40]:
   gender  SeniorCitizen  Partner  Dependents  tenure  PhoneService  MultipleLines  OnlineSecurity  OnlineBackup  DeviceProtection  ...  InternetService_DSL  InternetService_Fiber optic  InternetService_N
0  Female            0      1          0      1         No              No              No              Yes              No ...              1              0
1  Male            0      0      0      34         Yes              No              Yes              No              Yes ...              1              0
2  Male            0      0      0      2         Yes              No              Yes              Yes              No ...              1              0
3  Male            0      0      0      45         No              No              No              Yes              No ...              1              0
4  Female            0      0      0      2         Yes              No              No              No              No ...              0              1
...  ...              ...      ...      ...      ...      ...      ...      ...      ...      ...  ...  ...
7038  Male            0      1      1      24         Yes              Yes              Yes              No              Yes ...              1              0
7039  Female            0      1      1      72         Yes              Yes              No              Yes              Yes ...              0              1
7040  Female            0      1      1      11         No              No              No              No              No ...              1              0
7041  Male            1      1      0      4         Yes              Yes              No              No              No ...              0              1
7042  Male            0      0      0      66         Yes              No              Yes              No              Yes ...              0              1
7032 rows x 27 columns

In [41]: df2.dtypes

Out[41]:
gender                object
SeniorCitizen         int64
Partner               int64
Dependents            int64
tenure                int64
PhoneService          object
MultipleLines         object
InternetService       object
OnlineSecurity        object
OnlineBackup          object
DeviceProtection      object
TechSupport           object
StreamingTV           object
StreamingMovies       object
Contract              object
PaperlessBilling      object
PaymentMethod         object
MonthlyCharges        float64
TotalCharges          float64
Churn                 object
InternetService_DSL  uint8
InternetService_Fiber optic  uint8
InternetService_N      uint8
Contract_One year      uint8
Contract_Two year      uint8
PaymentMethod_Bank transfer (automatic)  uint8
PaymentMethod_Credit card (automatic)  uint8
PaymentMethod_Electronic check  uint8
PaymentMethod_Mailed check  uint8
dtype: object

In [ ]:

Matplotlib is a comprehensive library for creating static, animated, and interactive visualizations in Python.

In [43]: from matplotlib import pyplot as plt
import matplotlib inline

In [ ]:

In [58]: fig = plt.figure()
fig.set_figheight(7)
fig.set_figwidth(11)
x = [1, 51, 120, 138]
y = [24, 37, 55, 99]
x1 = [25, 52, 88, 150]
y1 = [2, 59, 48, 98]
plt.plot(x, y, color='red', linewidth=3, label='Income in $ per day')
plt.plot(x1, y1, color='blue', linewidth=3, label='Income in $ per day')
plt.title("My Data", fontdict = {'fontname': 'Comic Sans MS', 'fontsize': 28, 'fontweight': 'bold'})
plt.xlabel("X", fontdict = {'fontname': 'Comic Sans MS', 'fontsize': 14, 'fontweight': 'bold'})
plt.ylabel("Y", fontdict = {'fontname': 'Comic Sans MS', 'fontsize': 14, 'fontweight': 'bold'})
plt.legend(loc='lower right', ncol=2)

Out[58]:
Text(0.5, 0, 'X')


In [60]: fig = plt.figure()
fig.set_figheight(7)
fig.set_figwidth(11)
x = [1, 51, 120, 138]
y = [24, 37, 55, 99]
x1 = [25, 52, 88, 150]
y1 = [2, 59, 48, 98]
plt.plot(x, y, color='red', linewidth=3, label='Income in $ per day')
plt.plot(x1, y1, color='blue', linewidth=3, label='Income in $ per day')
plt.title("My Data", fontdict = {'fontname': 'Comic Sans MS', 'fontsize': 28, 'fontweight': 'bold'})
plt.xlabel("X", fontdict = {'fontname': 'Comic Sans MS', 'fontsize': 14, 'fontweight': 'bold'})
plt.ylabel("Y", fontdict = {'fontname': 'Comic Sans MS', 'fontsize': 14, 'fontweight': 'bold'})
plt.legend(loc='lower right', ncol=2)

Out[60]:
Text(0.5, 0, 'X')


In [64]: fig = plt.figure()
fig.set_figheight(7)
fig.set_figwidth(11)
x = [1, 51, 120, 138]
y = [24, 37, 55, 99]
x1 = [25, 52, 88, 150]
y1 = [2, 59, 48, 98]
plt.plot(x, y, color='red', linewidth=3, label='Age')
plt.plot(x1, y1, color='blue', linewidth=3, label='Income in $ per day')
plt.title("My Data", fontdict = {'fontname': 'Comic Sans MS', 'fontsize': 28, 'fontweight': 'bold'})
plt.xlabel("X", fontdict = {'fontname': 'Comic Sans MS', 'fontsize': 14, 'fontweight': 'bold'})
plt.ylabel("Y", fontdict = {'fontname': 'Comic Sans MS', 'fontsize': 14, 'fontweight': 'bold'})
plt.legend(loc='lower right', ncol=2)

Out[64]:
<matplotlib.legend.Legend at 0x19b48ac97f9>


In [69]: fig = plt.figure()
fig.set_figheight(7)
fig.set_figwidth(11)
x = [1, 51, 120, 138]
y = [24, 37, 55, 99]
x1 = [25, 52, 88, 150]
y1 = [2, 59, 48, 98]
plt.plot(x, y, color='red', linewidth=3, label='Age')
plt.plot(x1, y1, color='blue', linewidth=3, label='Income in $ per day')
plt.title("My Data", fontdict = {'fontname': 'Comic Sans MS', 'fontsize': 28, 'fontweight': 'bold'})
plt.xlabel("X", fontdict = {'fontname': 'Comic Sans MS', 'fontsize': 14, 'fontweight': 'bold'})
plt.ylabel("Y", fontdict = {'fontname': 'Comic Sans MS', 'fontsize': 14, 'fontweight': 'bold'})
plt.legend(bbox_to_anchor=(0.75, 1.15), ncol=2)

Out[69]:
<matplotlib.legend.Legend at 0x19b422e6d9d>


In [70]: from matplotlib import style
style.use('ggplot')
fig = plt.figure()
fig.set_figheight(7)
fig.set_figwidth(11)
x = [1, 51, 120, 138]
y = [24, 37, 55, 99]
x1 = [25, 52, 88, 150]
y1 = [2, 59, 48, 98]
plt.plot(x, y, color='red', linewidth=3, label='Age')
plt.plot(x1, y1, color='blue', linewidth=3, label='Income in $ per day')
plt.title("My Data", fontdict = {'fontname': 'Comic Sans MS', 'fontsize': 28, 'fontweight': 'bold'})
plt.xlabel("X", fontdict = {'fontname': 'Comic Sans MS', 'fontsize': 14, 'fontweight': 'bold'})
plt.ylabel("Y", fontdict = {'fontname': 'Comic Sans MS', 'fontsize': 14, 'fontweight': 'bold'})
plt.legend(bbox_to_anchor=(0.75, 1.15), ncol=2)

Out[70]:
<matplotlib.legend.Legend at 0x19b422ea49d>


In [71]: style.use('dark_background')
fig = plt.figure()
fig.set_figheight(7)
fig.set_figwidth(11)
x = [1, 51, 120, 138]
y = [24, 37, 55, 99]
x1 = [25, 52, 88, 150]
y1 = [2, 59, 48, 98]
plt.plot(x, y, color='red', linewidth=3, label='Age')
plt.plot(x1, y1, color='blue', linewidth=3, label='Income in $ per day')
plt.title("My Data", fontdict = {'fontname': 'Comic Sans MS', 'fontsize': 28, 'fontweight': 'bold'})
plt.xlabel("X", fontdict = {'fontname': 'Comic Sans MS', 'fontsize': 14, 'fontweight': 'bold'})
plt.ylabel("Y", fontdict = {'fontname': 'Comic Sans MS', 'fontsize': 14, 'fontweight': 'bold'})
plt.legend(bbox_to_anchor=(0.75, 1.15), ncol=2)

Out[71]:
<matplotlib.legend.Legend at 0x19b4236491d>


In [73]: style.use('dark_background')
fig = plt.figure()
fig.set_figheight(7)
fig.set_figwidth(11)
x = [1, 51, 120, 138]
y = [24, 37, 55, 99]
x1 = [25, 52, 88, 150]
y1 = [2, 59, 48, 98]
plt.scatter(x, y, color='red', linewidth=3, label='Age', linestyle='dotted')
plt.scatter(x1, y1, color='blue', linewidth=3, label='Income in $ per day', linestyle='dotted')
plt.title("My Data", fontdict = {'fontname': 'Comic Sans MS', 'fontsize': 28, 'fontweight': 'bold'})
plt.xlabel("X", fontdict = {'fontname': 'Comic Sans MS', 'fontsize': 14, 'fontweight': 'bold'})
plt.ylabel("Y", fontdict = {'fontname': 'Comic Sans MS', 'fontsize': 14, 'fontweight': 'bold'})
plt.legend(bbox_to_anchor=(0.75, 1.15), ncol=2)

Out[73]:
<matplotlib.legend.Legend at 0x19b3f5da2e0>


In [77]: style.use('dark_background')
fig = plt.figure()
fig.set_figheight(7)
fig.set_figwidth(11)
x = [1, 51, 120, 138]
y = [24, 37, 55, 99]
x1 = [25, 52, 88, 150]
y1 = [2, 59, 48, 98]
plt.scatter(x, y, color='red', linewidth=3, label='Age', linestyle='dotted')
plt.scatter(x1, y1, color='blue', linewidth=3, label='Income in $ per day', linestyle='dotted')
plt.title("My Data", fontdict = {'fontname': 'Comic Sans MS', 'fontsize': 28, 'fontweight': 'bold'})
plt.xlabel("X", fontdict = {'fontname': 'Comic Sans MS', 'fontsize': 14, 'fontweight': 'bold'})
plt.ylabel("Y", fontdict = {'fontname': 'Comic Sans MS', 'fontsize': 14, 'fontweight': 'bold'})
plt.legend(bbox_to_anchor=(0.75, 1.15), ncol=2)

Out[77]:
<matplotlib.legend.Legend at 0x19b3f5da2e0>


In [ ]:

In [79]: df = pd.read_csv("customer_churn.csv")

In [81]: df

Out[81]:
   customerID  gender  SeniorCitizen  Partner  Dependents  tenure  PhoneService  MultipleLines  InternetService  OnlineSecurity  ...  DeviceProtection  TechSupport  StreamingTV  StreamingMovies
0    7590  FEMALE      0             0      Yes      1         No              No phone service  DSL              No ...              No              No              No
1     655  FEMALE      0             0      No      34         Yes              No              DSL              Yes ...              No              No              No
2    3668  MALE      0             0      No      2         Yes              No              No              Yes ...              No              No              No
3     796  FEMALE      0             0      No      45         No              No phone service  DSL              Yes ...              Yes              No              No
4    8237  FEMALE      0             0      No      2         Yes              No              Fiber optic  No ...              No              No              No
...  ...  ...      ...      ...      ...      ...      ...      ...      ...      ...  ...  ...
7038 680-RESVB  MALE      0      Yes        Yes      24         Yes              Yes              DSL              Yes ...              Yes              Yes              Yes
7039 2224-AQJH  FEMALE      0      Yes        Yes      72         Yes              Yes              Fiber optic  No ...              No              No              Yes
7040 4803-IZATL  FEMALE      0      Yes         Yes      11         No              No phone service  DSL              Yes ...              No              No              No
7041 688-LTHMD  MALE      1      Yes         No      4         Yes              Yes              Fiber optic  No ...              No              No              No
7042 3185-A3EKK  MALE      0             0      No      66         Yes              No              Fiber optic  Yes ...              Yes              Yes              Yes
7043 rows x 21 columns

In [82]: tenureChurnNo = df[df.Churn == 'No'].tenure
tenureChurnYes = df[df.Churn == 'Yes'].tenure
plt.hist([tenureChurnNo, tenureChurnYes], color=['red', 'green'], label=['Churn=No', 'Churn=Yes'])
plt.legend()

Out[82]:
<matplotlib.legend.Legend at 0x19b433f8280>


In [85]: tenureChurnNo = df[df.Churn == 'No'].MonthlyCharges
tenureChurnYes = df[df.Churn == 'Yes'].MonthlyCharges
plt.xlabel("Monthly Charges")
plt.ylabel("Number of Customers")
plt.hist([tenureChurnNo, tenureChurnYes], color=['red', 'green'], label=['Churn=No', 'Churn=Yes'])
plt.legend()

Out[85]:
<matplotlib.legend.Legend at 0x19b433f8280>

```