# INTRODUCTION TO MACHINE LEARNING PROJECT PROGRESS REPORT

ONUR HAKTAN                    30.11.2022                    ESKISEHIR TECHNICAL
GÖKSU TURAC                                                  UNIVERSITY

## ANALYSIS OF THE COMPANY'S CUSTOMER LOSS

We have reviewed the data set and finished the preprocessing process.

- We took the data set for the lost customer analysis from Kaggle. Then we read the csv file(dataset).

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

```python
df=pd.read_csv("WA_Fn-UseC_-Telco-Customer-Churn.csv",sep=',',decimal='.')
```

- A customer data set consisting of 21 characteristics was provided for use in the analysis, and the target variable of these attributes was defined as "Churn".

```python
from numpy import column_stack

df.info(column_stack)
```

```
<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)
```

- We examined the data in the data set and extracted the data that was not useful to us.(Customer ID)



- We have made changes that make it easier for us to do analysis.

- We decided that we need to change the type of the TotalCharges feature.

```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   object
 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   object
 19  Churn             7043 non-null   object
dtypes: float64(1), int64(1), object(18)
memory usage: 1.1+ MB
```

- We have drawn a graph of customer loss.

- We checked for lost data.

```python
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
```
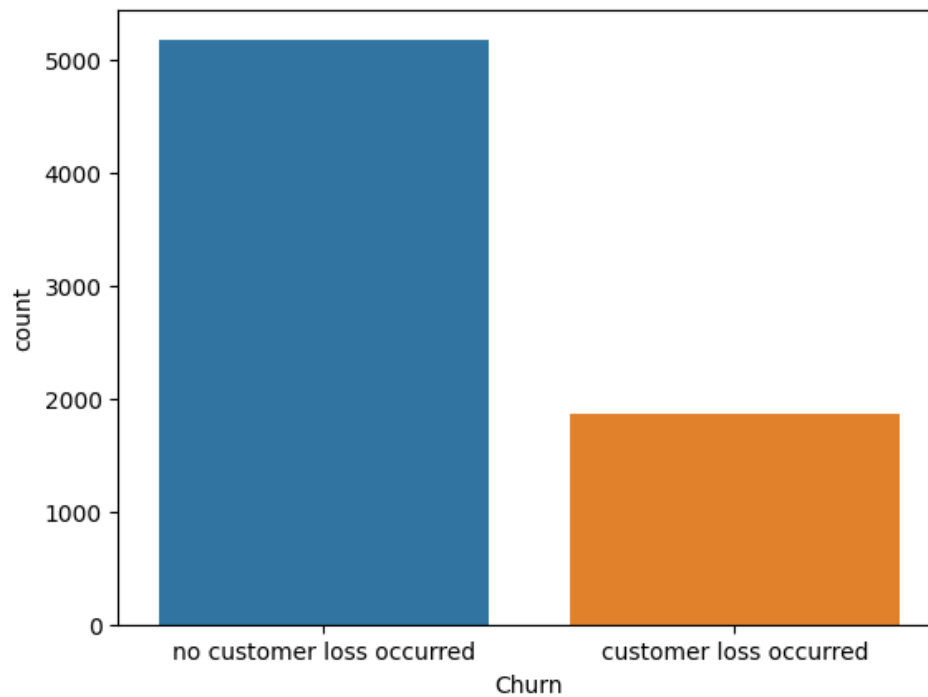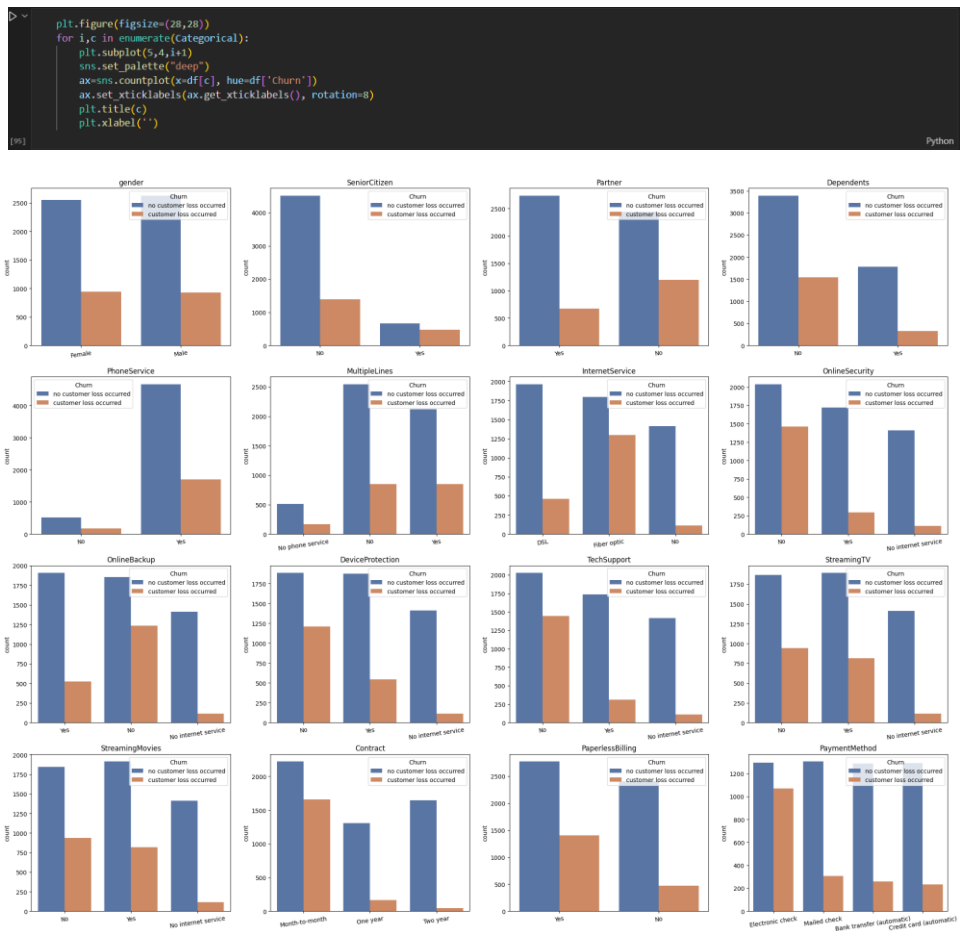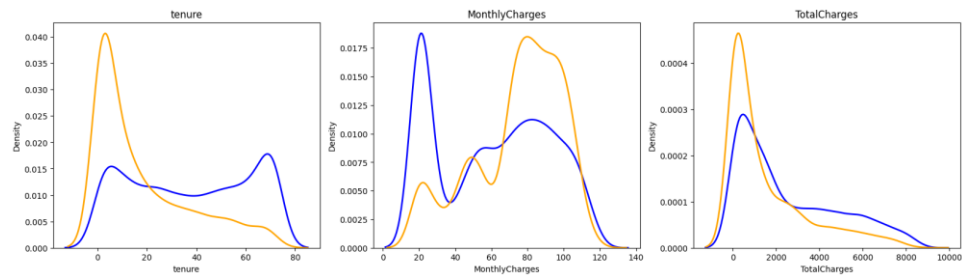
- In order to understand how the variables affect the target value, we had them plotted.

```python
plt.figure(figsize=(28,28))
for i,c in enumerate(Categorical):
    plt.subplot(5,4,i+1)
    sns.set_palette("deep")
    ax=sns.countplot(x=df[c], hue=df['Churn'])
    ax.set_xticklabels(ax.get_xticklabels(), rotation=8)
    plt.title(c)
    plt.xlabel('')
```
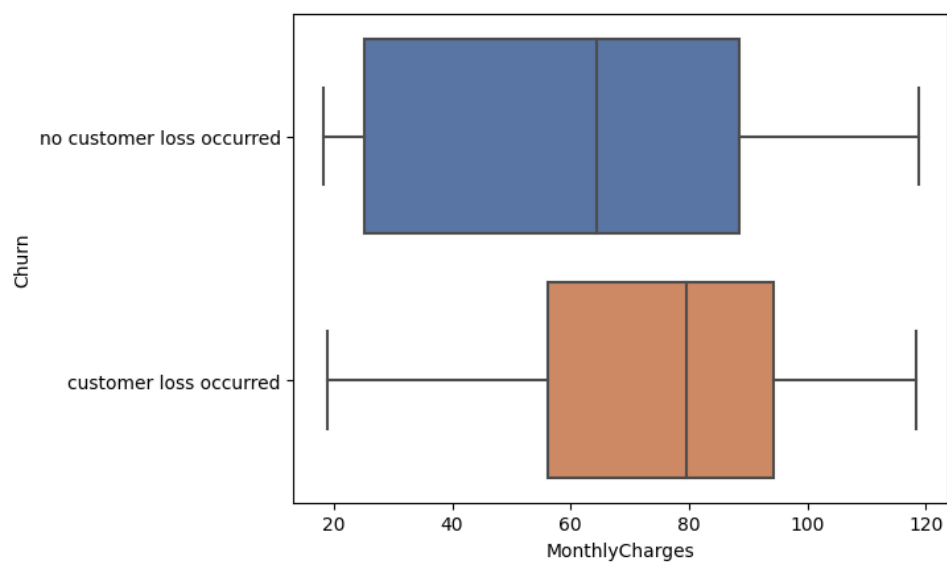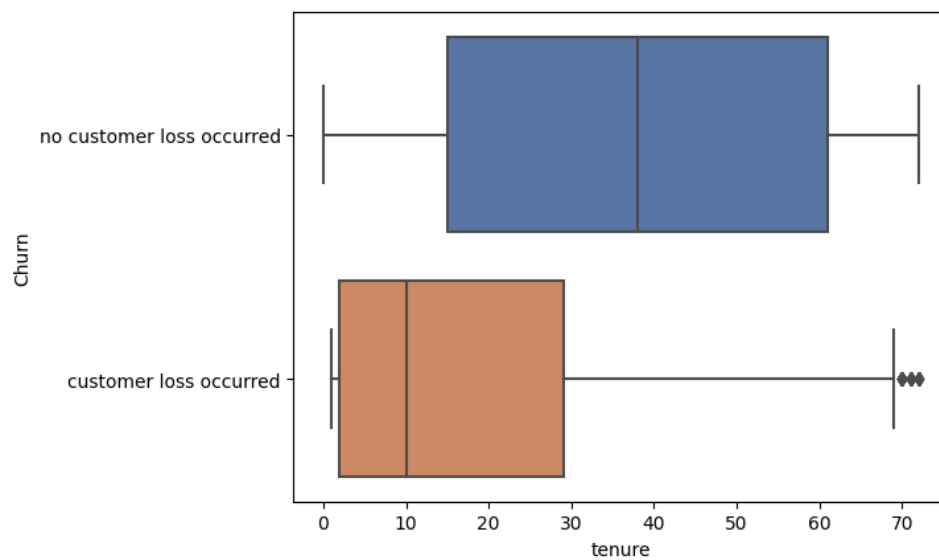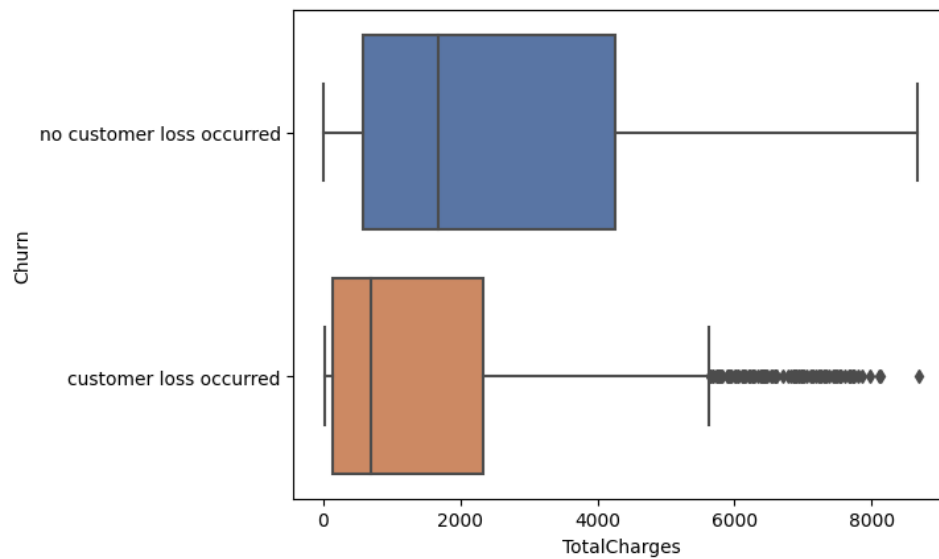
```
plt.figure(figsize=(20,5))
for i,c in enumerate(['tenure', 'MonthlyCharges', 'TotalCharges']):
    plt.subplot(1,3,i+1)
    sns.distplot(df[df['Churn'] == 'no customer loss occurred'][c], kde=True, color='blue', hist=False, kde_kws=dict(linewidth=2), label='no customer loss occurred')
    sns.distplot(df[df['Churn'] == 'customer loss occurred'][c], kde=True, color='Orange', hist=False, kde_kws=dict(linewidth=2), label='customer loss occurred')
    plt.title(c)
```
Python



- We looked to see if there were any outlier data.

In conclusion; In the part of the Total Charges that resulted in customer loss, we saw some outlier data. There are some outlier data also appears in the tenure. We have cleared the outlier data from the dataset with  LabelEncoder. In the next part, we will divide the final dataset into two (test and train) and perform machine learning.