

## Importing libraries

The data set includes information about:

- Customers who left within the last month – the column is called Churn
- Services that each customer has signed up for – phone, multiple lines, internet, online security, online backup, device protection, tech support, and streaming TV and movies
- Customer account information – how long they've been a customer, contract, payment method, paperless billing, monthly charges, and total charges
- Demographic info about customers – gender, age range, and if they have partners and dependents

```
import pandas as pd
from pandas.api.types import is_string_dtype
from pandas.api.types import is_numeric_dtype
import numpy as np
import matplotlib.pyplot as plt
import seaborn as sns
from scipy.stats import probplot
import sklearn
!pip install ppscore
import ppscore as pps
```

```
Requirement already satisfied: ppscore in /usr/local/lib/python3.7/dist-packa
Requirement already satisfied: scikit-learn<1.0.0,>=0.20.2 in /usr/local/lib/
Requirement already satisfied: pandas<2.0.0,>=1.0.0 in /usr/local/lib/python3
Requirement already satisfied: python-dateutil>=2.7.3 in /usr/local/lib/pytho
Requirement already satisfied: numpy>=1.15.4 in /usr/local/lib/python3.7/dist
Requirement already satisfied: pytz>=2017.2 in /usr/local/lib/python3.7/dist-
Requirement already satisfied: six>=1.5 in /usr/local/lib/python3.7/dist-pack
Requirement already satisfied: joblib>=0.11 in /usr/local/lib/python3.7/dist-
Requirement already satisfied: scipy>=0.17.0 in /usr/local/lib/python3.7/dist
```

```
df = pd.read_csv("Telco-Customer-Churn.csv")    # Reading the given CSV file
```

## Data Analysis

```
df.head(5)
```

	customerID	gender	SeniorCitizen	Partner	Dependents	tenure	PhoneService
0	7590-VHVEG	Female	0	Yes	No	1	N
1	5575-GNVDE	Male	0	No	No	34	Ye

Description of a few features:

- gender - Whether the customer is a male or a female
- SeniorCitizen - Whether the customer is a senior citizen or not (1, 0)
- Partner - Whether the customer has a partner or not (Yes, No)
- Dependents - Whether the customer has dependents or not (Yes, No)
- tenure - Number of months the customer has stayed with the company
- PhoneService - Whether the customer has a phone service or not (Yes, No)
- MultipleLines - Whether the customer has multiple lines or not, that is capable of holding some calls (Yes, No, No phone service)
- InternetService - Customer's internet service provider (DSL, Fiber optic, No)
- OnlineSecurity - Whether the customer has online security or not (Yes, No, No internet service)

```
df.shape
```

```
(7043, 21)
```

```
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
```

```
df.TotalCharges.unique()
```

```
array(['29.85', '1889.5', '108.15', ..., '346.45', '306.6', '6844.5'],
      dtype=object)
```

```
df['TotalCharges'] = df['TotalCharges'].replace(" ", 0).astype('float32')
```

Most of the data is categorical (might be ordinal or nominal).

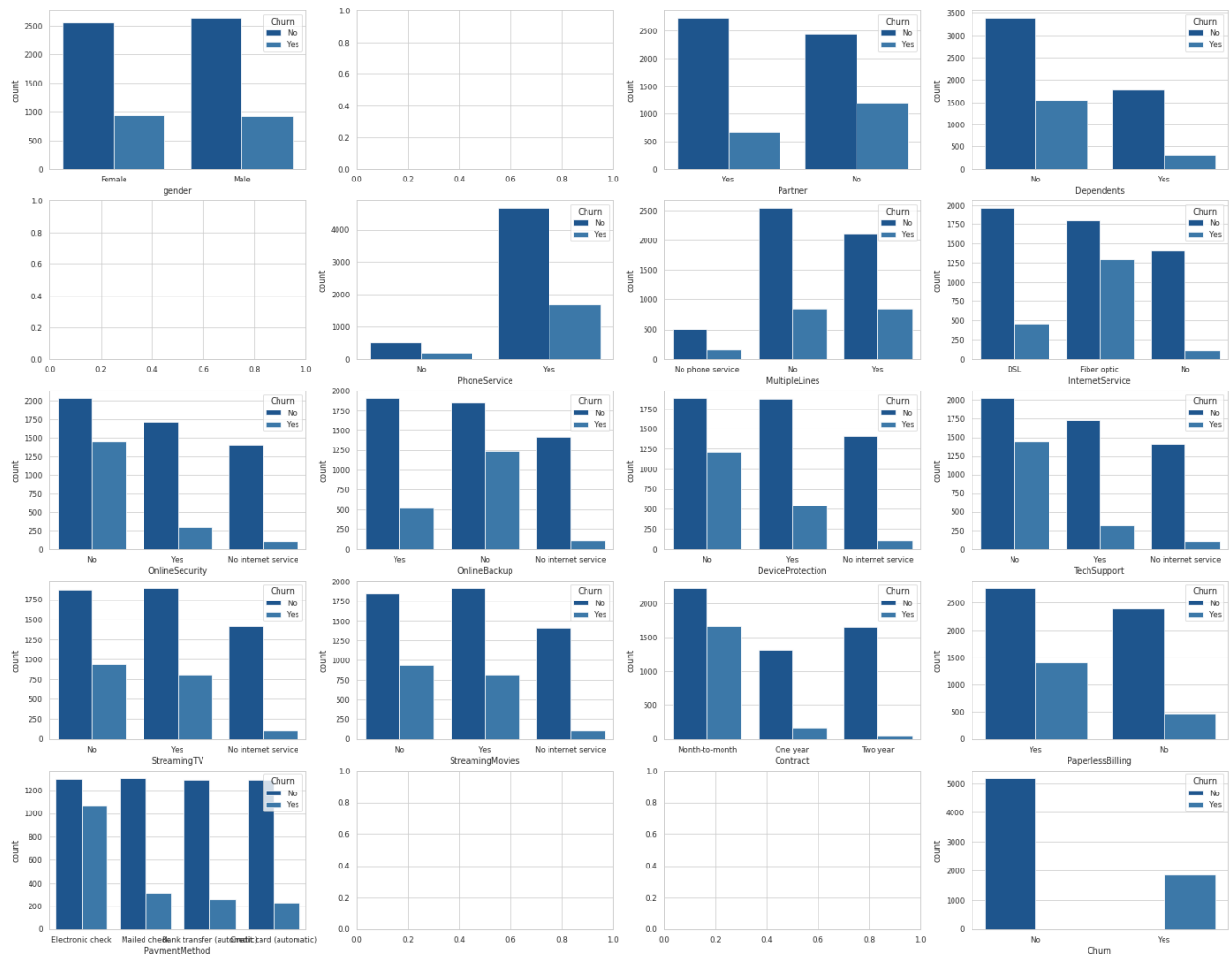
Number of features = 21 (20, if we ignore Customer ID)

There are no missing or null values

```
sns.set(palette=sns.color_palette("Blues_r"),style="whitegrid", context = "paper")
sns.despine(top=True, left=True, bottom=False, right=False)
```

<Figure size 432x288 with 0 Axes>

```
# Plotting barplot for categorical variables only, except customer ID
fig,axes = plt.subplots(5,4, figsize=(25, 20))
for i,ax in zip(df.columns[1:],axes.flatten()):
    if(is_string_dtype(df[i])):
        sns.countplot(x=df[i],ax=ax,hue=df["Churn"])
```



Observation from the above plot:

- In the given data, around 71% people **didn't churn**
- Churn distribution is almost similar for both genders
- Electronic payment method has more churn rate as compared to other payment methods
- Month-to-Month contract has higher churn rate than other contracts. It is logical too, cause once you have bought the service for an year or two, people generally do not prefer changing the service due to the hassle it might cause.
- Customer with no Internet service (wherever this value is present) have a very low churn rate. (Maybe because they are not used to the modern technology, and do not prefer changing services)
- People with device protection and Tech-Support (Yes) have less churn rate than people with no device protection or Tech-Support (No), indicating that people are satisfied with

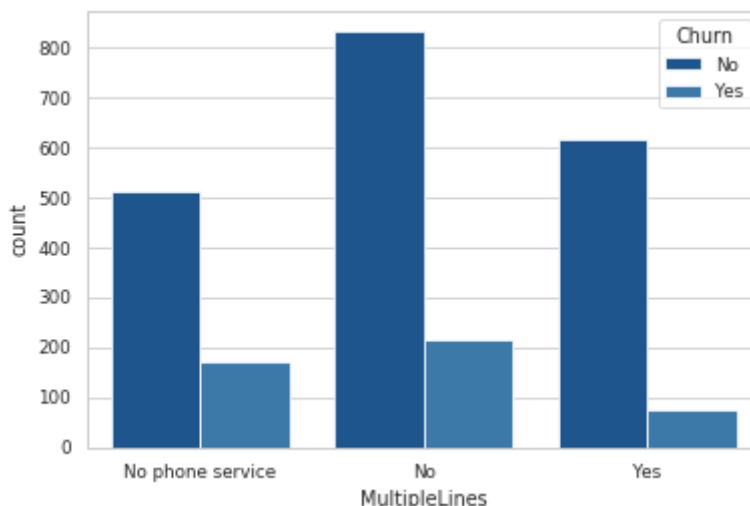
the services provided

- Around 650 people didn't take phone service, number is relatively small compared to the total size of the dataset and won't be analyzed.
- People with dependencies or partners have comparatively less churn rate (Around 12.5% for dependencies and 69.5% for partners), than people with no partners or dependencies.
- Customers with no internet services are generally from rural area or want to use their service for just calling or other purpose, thus leaving a very small margin for dissatisfaction and changing the service.
- Another interesting pattern in internet-services is customer with DSL have relatively very less churn rate as compared to Fiber-optics, indicating dissatisfaction with the latter.
- A general pattern can be observed: Customers who didn't take services like Online-Security, Online-Backup, Device-Support, Tech-Support, have higher churn rate.

```
# sns.countplot(df[df["InternetService"]=="Fiber optic"]["OnlineSecurity"])
# sns.countplot(df[df["InternetService"]=="No"]["OnlineSecurity"])
# sns.countplot(df[df["InternetService"]=="No"]["OnlineBackup"])
sns.countplot(df[df["InternetService"]=="DSL"]["MultipleLines"],hue = df["Churn"])
plt.show()
# sns.countplot(df[df["InternetService"]=="Fiber optic"]["OnlineSecurity"])

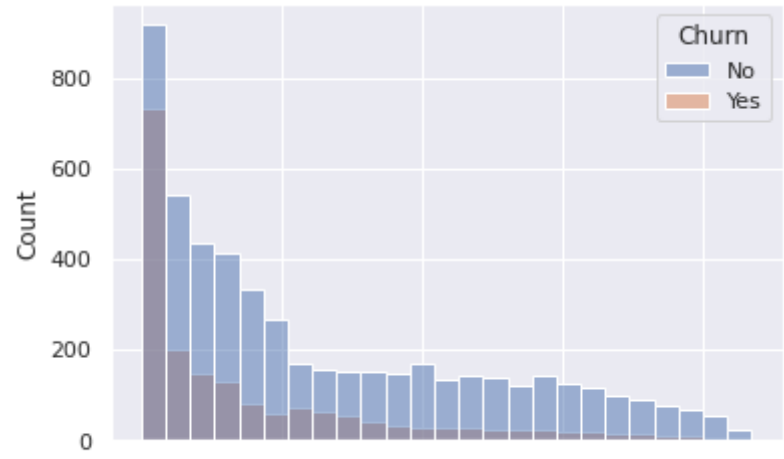
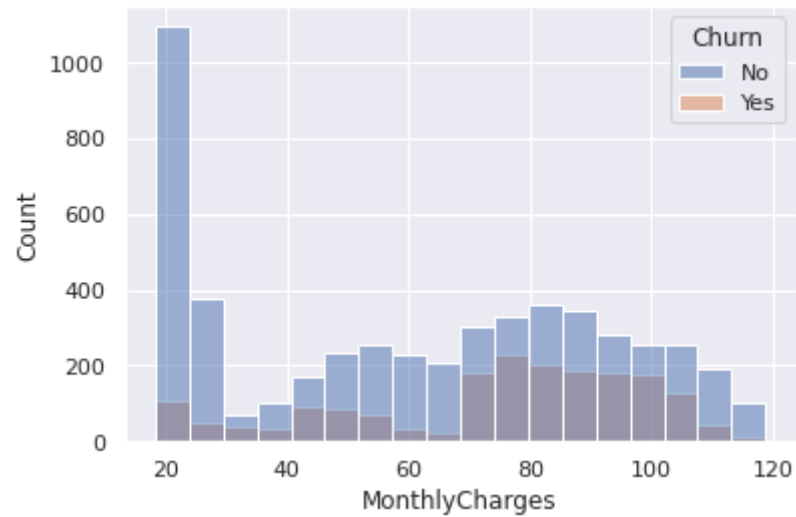
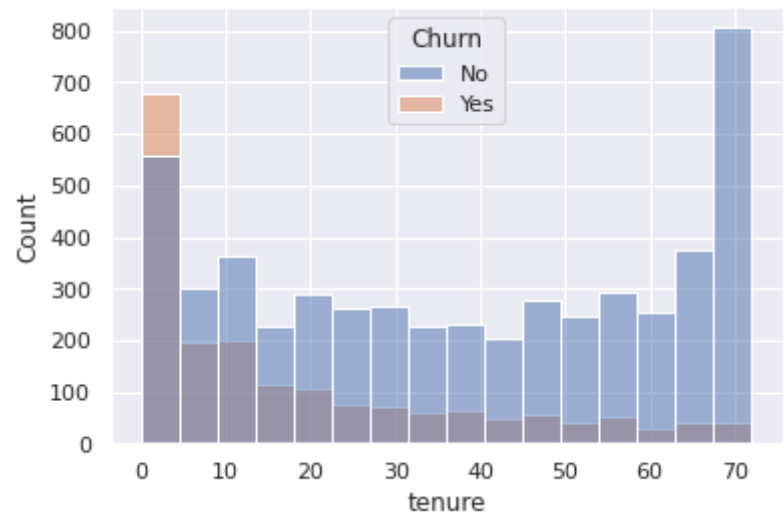
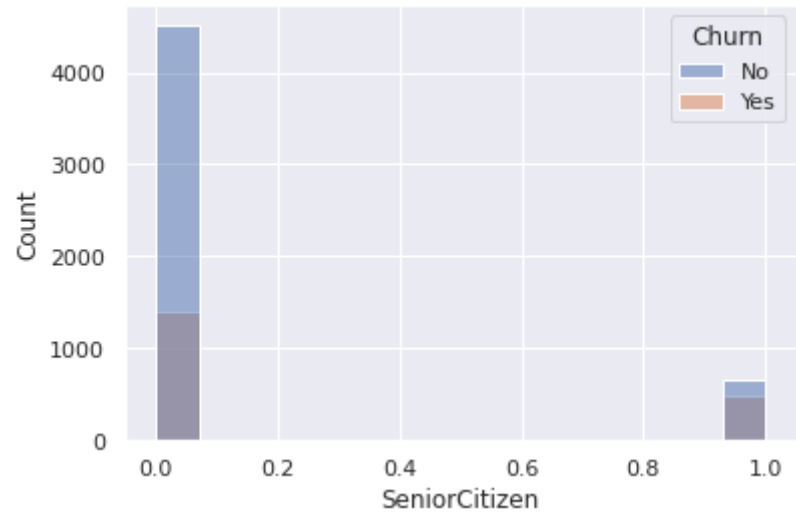
# Customers with no internet service, do not have onlineSecurity, Online Backup, S
```

/usr/local/lib/python3.7/dist-packages/seaborn/\_decorators.py:43: FutureWarning  
FutureWarning



Double-click (or enter) to edit

```
sns.set()
for i in (df.columns):
    if(is_numeric_dtype(df[i])):
        sns.histplot(x=df[i],hue=df["Churn"])
        plt.show()
```



0 2000 4000 6000 8000  
TotalCharges

Observations from the above plots:

- People with less value of tenure are more likely to churn.
- Customer's whose monthly bill is less ( $\leq 70$ ) or among the greatest ( $\geq 110$ ) are less likely to churn, than customers in the middle range.
- Senior citizens are more likely to churn

Creating a copy of Dataframe, and replacing categorical values with numerical, for ease in further analysis.

```
df1 = df
df1.drop(["customerID"],axis=1, inplace=True)
df1.replace(to_replace = ["Yes", "Two year","Female", "DSL"],value=2, inplace=True)
df1.replace(to_replace = ["No","One year", "Male", "Fiber optic"],value=1, inplace=True)
df1.replace(to_replace = ["No internet service", "Month-to-month", "No phone servi
```

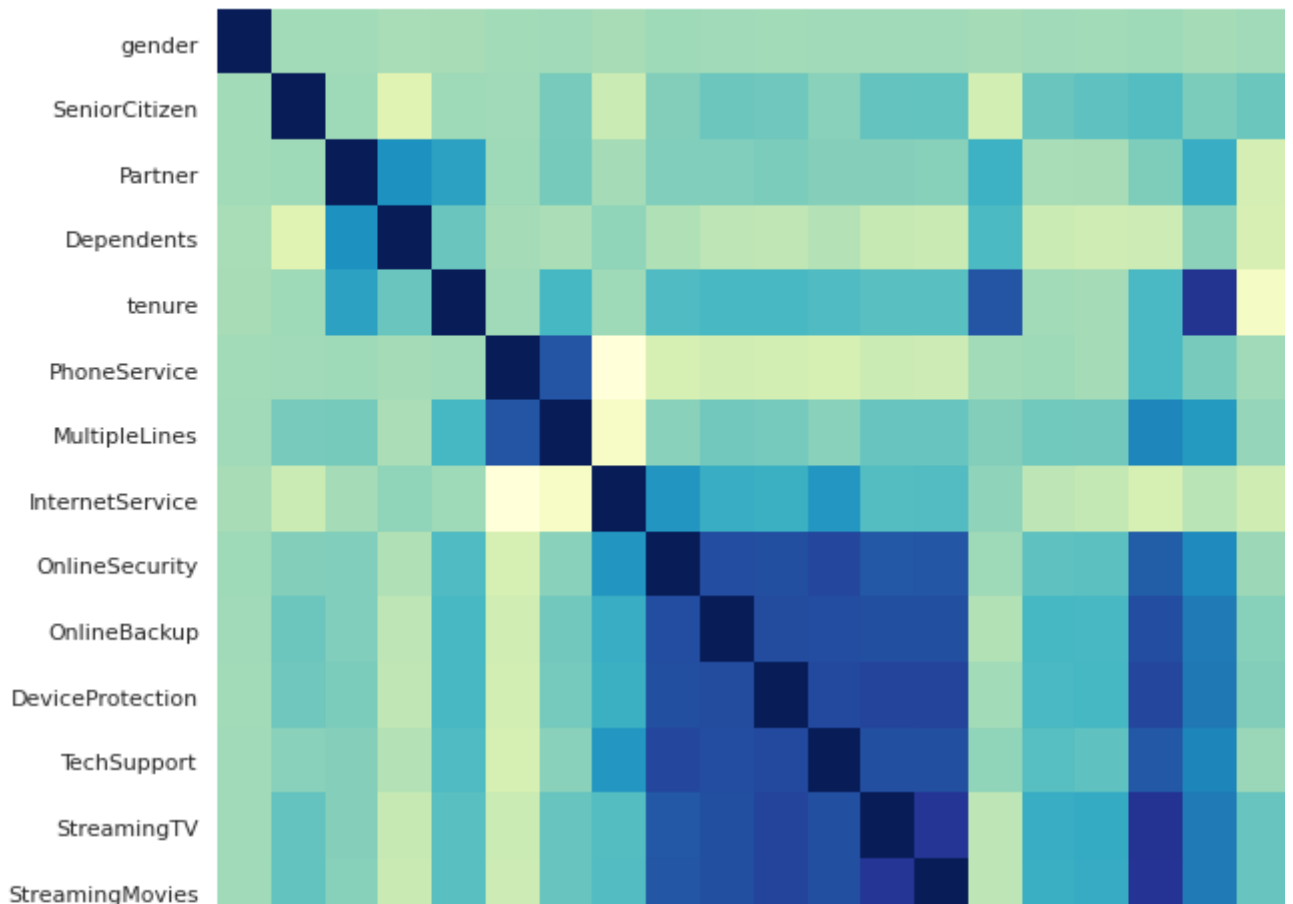
```
df1.PaymentMethod.unique()
```

```
array([3, 0, 1, 2])
```

```
df1.replace(to_replace = "Electronic check",value=3, inplace=True)
df1.replace(to_replace = "Credit card (automatic)",value=2, inplace=True)
df1.replace(to_replace = "Bank transfer (automatic)",value=1, inplace=True)
df1.replace(to_replace = "Mailed check",value=0, inplace=True)
```

```
plt.figure(figsize=(12, 12))
sns.heatmap(df1.corr(),cmap="YlGnBu")
```

<matplotlib.axes.\_subplots.AxesSubplot at 0x7fe818d17f10>



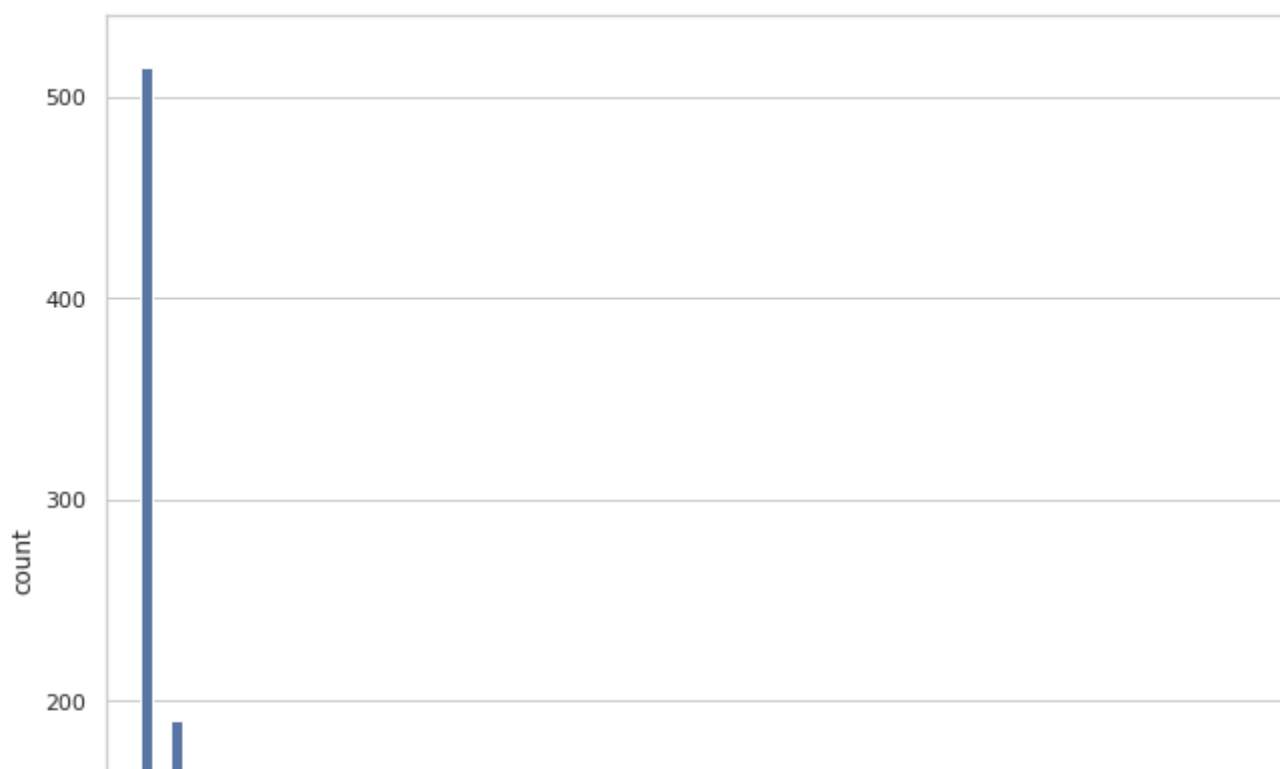
Observations from the above heatmap:

- Contract, Total Charges and Tenure are highly correlated
- Monthly charges increase when customers bought other services (Tech-Support, Online backup, Streaming, etc.)
- Most Senior citizens have no dependents
- All the offered services are correlated, we can just replace them with one attribute Services
- Tenure and Churn are negatively correlated as expected

```
plt.figure(figsize=(20,10))
sns.set(style="whitegrid")
sns.countplot(df['tenure'], hue = df["Partner"])
```



```
/usr/local/lib/python3.7/dist-packages/seaborn/_decorators.py:43: FutureWarning
FutureWarning
<matplotlib.axes._subplots.AxesSubplot at 0x7fe807db8790>
```

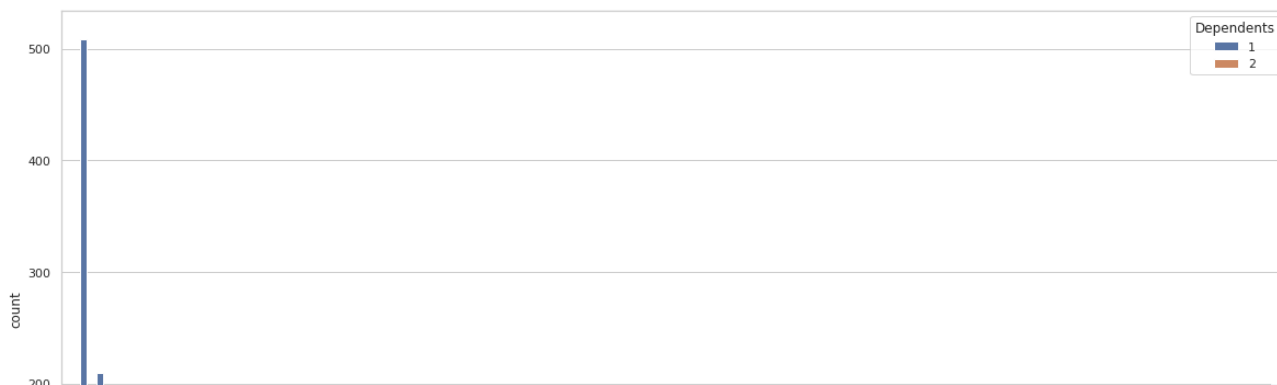


If customers have partner, then their tenure is more likely to be greater as compared to the ones without a partner



```
plt.figure(figsize=(20,10))
sns.set(style="whitegrid")
sns.countplot(df['tenure'], hue = df["Dependents"])
```

```
/usr/local/lib/python3.7/dist-packages/seaborn/_decorators.py:43: FutureWarning
FutureWarning
<matplotlib.axes._subplots.AxesSubplot at 0x7fe7bb78c4d0>
```

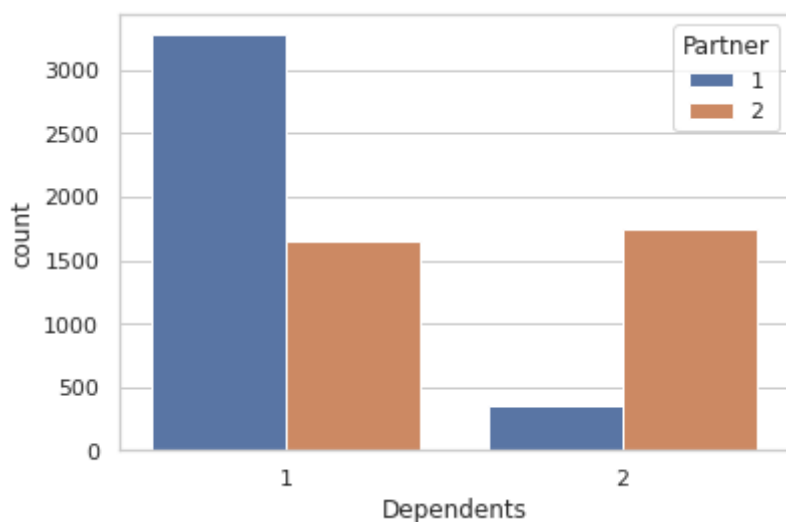


Customers with no dependent are likely to have less tenure.



```
sns.countplot(x = "Dependents", data=df1, hue="Partner")
```

```
<matplotlib.axes._subplots.AxesSubplot at 0x7fe7bae07390>
```



We can club Dependents and Partner into one attribute

Plotting Heatmap of Predictive Power Score

```
matrix_df = pps.matrix(df)[['x', 'y', 'ppscore']].pivot(columns='x', index='y', va
```

```
plt.figure(figsize=(12,12))
sns.heatmap(matrix_df, vmin=0, vmax=1, cmap="YlGnBu", linewidths=0.5)
```

Heatmap showing the correlation matrix for 20 variables. The color scale ranges from 0.0 (yellow) to 1.0 (dark blue). The diagonal is dark blue (1.0). The color scale ranges from 0.0 (yellow) to 1.0 (dark blue).

- Services are good predictor of monthly charges,
- total charges and tenure are good predictors of each other
- Dependents can predict Partners , but the opposite is not True
- most relation observed are among services provided, charges (monthly and total) and tenure.

✓ 0s completed at 13:31

● ×