

# Exploratory Data Analysis Using PYSARK and Machine Learning models

```
In [1]: !python3 --version

Python 3.10.6

In [2]: !jupyter --version

Selected Jupyter core packages...
IPython          : 7.31.1
ipykernel        : 6.7.0
ipywidgets       : 6.0.0
jupyter_client   : 7.1.2
jupyter_core     : 4.0.1
jupyter_server   : not installed
jupyterlab       : not installed
nbclient         : 0.6.0
nbconvert        : 6.4.0
nbformat         : 5.1.3
notebook         : 6.4.8
qtconsole        : not installed
traitlets        : 5.1.1

In [3]: pip install -U notebook-as-pdf

Defaulting to user installation because normal site-packages is not writable
Collecting notebook-as-pdf
  Downloading notebook-as-pdf-0.5.0-py3-none-any.whl (6.5 kB)
Collecting PyPDF2
  Downloading pypdf2-2.11.2-py3-none-any.whl (220 kB)
Requirement already satisfied: nbconvert<0.10.0,>=1.4.3 in /usr/lib/python3/dist-packages (from notebook-as-pdf) (6.4.0)
Collecting pypetter
  Downloading pypetter-1.0.2-py3-none-any.whl (83 kB)
Requirement already satisfied: appdirs<2.0.0,>=1.4.3 in /usr/lib/python3/dist-packages (from pypetter->notebook-as-pdf) (1.2.6.5)
Requirement already satisfied: urllib3<2.0.0,>=1.25.8 in /usr/lib/python3/dist-packages (from pypetter->notebook-as-pdf) (1.4.4)
Requirement already satisfied: pysocks, tqdm, PyPDF2, certifi, pypetter, notebook-as-pdf
Successfully installed PyPDF2-2.11.2 certifi-2022.12.7 notebook-as-pdf-0.5.0 pyee-0.2.2 pypetter-1.0.2 tqdm-4.64.1 websockets-10.4
Note: you may need to restart the kernel to use updated packages.

In [3]: pip install pypetter

Defaulting to user installation because normal site-packages is not writable
Requirement already satisfied: importlib-metadata>=1.4 in /usr/lib/python3.10/site-packages (from pypetter) (4.6.4)
Requirement already satisfied: pypetter-lib<2020.1.0,>=1.0 in /usr/lib/python3/dist-packages (from pypetter) (1.21.5)
Requirement already satisfied: appdirs<2.0.0,>=1.4.3 in /usr/lib/python3/dist-packages (from pypetter) (1.4.4)
Requirement already satisfied: certifi<2021 in /usr/lib/python3.10/site-packages (from pypetter) (2022.12.7)
Requirement already satisfied: pyee<0.0.0,>=0.1.0 in /usr/lib/python3.10/site-packages (from pypetter) (0.2.2)
Requirement already satisfied: nbconvert<0.10.0,>=1.4.3 in /usr/lib/python3/dist-packages (from pypetter->notebook-as-pdf) (1.26.5)
Requirement already satisfied: pysocks, urllib3, PyPDF2, certifi, pypetter, notebook-as-pdf
Successfully installed PyPDF2-2.11.2 certifi-2022.12.7 notebook-as-pdf-0.5.0 pyee-0.2.2 pypetter-1.0.2 tqdm-4.64.1 websockets-10.4
Note: you may need to restart the kernel to use updated packages.

In [2]: pip install pyspark

Defaulting to user installation because normal site-packages is not writable
Collecting pyspark
  Downloading pyspark-3.3.1.tar.gz (281.4 MB)
    Preparing metadata (setup.py) ... done
    Collecting py4j==0.10.9.5
      Downloading py4j-0.10.9.5-py3-none-any.whl (199 kB)
Building wheels for collected packages: pyspark
  Building wheel for pyspark (setup.py) ... done
  Successfully installed py4j-0.10.9.5 pyspark-3.3.1
  Stored in directory: /home/b3dgata/.cache/pip/wheels/0f/f0/3d/517308b8c604868f4f89f214e0a022554e4e4e46969f46
Successfully built pyspark
Installing collected packages: py4j, pyspark
Successfully installed py4j-0.10.9.5 pyspark-3.3.1
Note: you may need to restart the kernel to use updated packages.

In [6]: pip install findspark

Defaulting to user installation because normal site-packages is not writable
Requirement already satisfied: findspark in /usr/lib/python3.10/site-packages (2.0.1)
Note: you may need to restart the kernel to use updated packages.

In [14]: pip install pyspark-dist-explore

Defaulting to user installation because normal site-packages is not writable
Collecting pyspark-dist-explore
  Downloading pyspark-dist-explore-0.1.8-py3-none-any.whl (7.2 kB)
Requirement already satisfied: matplotlib in /usr/lib/python3/dist-packages (from pyspark-dist-explore) (3.5.3)
Requirement already satisfied: pandas in /usr/lib/python3.10/site-packages (from pyspark-dist-explore) (1.5.1)
Requirement already satisfied: numpy in /usr/lib/python3/dist-packages (from pyspark-dist-explore) (1.8.0)
Requirement already satisfied: scipy in /usr/lib/python3/dist-packages (from pyspark-dist-explore) (1.21.5)
Requirement already satisfied: python-dateutil<2.8.1 in /usr/lib/python3/dist-packages (from pandas->pyspark-dist-explore) (2.8.1)
Installing collected packages: pyspark-dist-explore
Successfully installed pyspark-dist-explore-0.1.8
Note: you may need to restart the kernel to use updated packages.
```

```
In [ ]:
```

```
In [36]: import findspark
import warnings
warnings.filterwarnings('ignore')
import matplotlib.pyplot as plt
%matplotlib inline
import pandas as pd
import numpy as np
from pyspark.sql.types import *
from pyspark.sql import Row
import seaborn as sns
from pyspark import SparkContext
from pyspark import SQLContext
findspark.init()

from pyspark.sql import SparkSession
sc = SparkSession.builder \
    .appName("data analysis") \
    .master("local[*]") \
    .config("spark.jars","/spark-tree-plotting/") \
    .config("spark.jars.packages", "--jars /path/to/the/spark-tree-plotting_0.2.jar") \
    .config("spark.driver.memory", "--driver-class-path /path/to/the/spark-tree-plotting_0.2.jar") \
    .getOrCreate()

spark_con = SparkContext(sc=sc)
spark = SQLContext(sc)

print("Explore the data and analysis the assignment")

22/12/05 12:47:11 WARN SparkSession: Using an existing Spark session; only runtime SQL configurations will take effect.
Explore the data and analysis the assignment
```

```
In [2]: df = spark.read.csv('/home/b3dgata/TelecomCustomerAnalysis/Churn.csv', header=True, inferSchema=True)
df.printSchema()
```

[Stage 1: >] (0 + 1) / 1														
-- AccountLength: integer (nullable = true)														
-- VMailMessage: integer (nullable = true)														
-- DayMins: double (nullable = true)														
-- EveMins: double (nullable = true)														
-- NightMins: double (nullable = true)														
-- IntMins: double (nullable = true)														
-- CustServCalls: integer (nullable = true)														
-- Churn: integer (nullable = true)														
-- IntPlan: integer (nullable = true)														
-- VMailPlan: integer (nullable = true)														
-- DayCalls: integer (nullable = true)														
-- DayCharge: double (nullable = true)														
-- EveCalls: integer (nullable = true)														
-- EveCharge: double (nullable = true)														
-- NightCalls: integer (nullable = true)														
-- NightCharge: double (nullable = true)														
-- IntCalls: integer (nullable = true)														
-- IntCharge: double (nullable = true)														
-- State: string (nullable = true)														
-- AreaCode: integer (nullable = true)														
-- Phone: string (nullable = true)														

```
In [3]: df.show(5)
```

[Stage 1: >] (0 + 1) / 1														
[AccountLength][VMailMessage][DayMins][EveMins][NightMins][IntMins][CustServCalls][Churn][IntPlan][VMailPlan][DayCalls][DayCharge][EveCalls][EveCharge][NightCalls][NightCharge][IntCalls][IntCharge][State][AreaCode] Phone]														
1	128	261	265.1	197.4	244.7	10.0	1	0	0	1	120			
2	45.0	98	16.1	26.1	161.6	195.5	254.4	13.7	1	0	0	1	415	382-4657
3	27.7	183	16.62	183	11.45	3	3.7	0	0	415	371-7191			
4	41.38	110	18.3	243.4	122.2	162.6	12.2	0	1	415	358-1921			114
5	50.9	81	0	299.4	61.9	196.9	6.6	1.78	2	0	1	0	1	71
6	75	286	0	166.7	148.3	386.9	10.1	3	0	415	338-9999			113
7	28.34	122	12.61	121	8.41	3	2.73	0	0	415	338-6626			
only showing top 5 rows														

```
In [4]: df.columns
```

```
Out[4]: ['AccountLength',
'VMailMessage',
'DayMins',
'EveMins',
'NightMins',
'IntMins',
'CustServCalls',
'Churn',
'IntPlan',
'VMailPlan',
'DayCalls',
'DayCharge',
'EveCalls',
'EveCharge',
'NightCalls',
'NightCharge',
'IntCalls',
'IntCharge',
'State',
'AreaCode',
'Phone']
```

## Account.Length: how long account has been active.

```
In [5]: df.select('AccountLength').show()

##df['Churn']= df['Churn'].astype('category')
##df['Intl Plan']= df['Intl Plan'].astype('category')
##df['VMail Plan']= df['VMail Plan'].astype('category')

+-----+
| AccountLength |
+-----+
| 128 |
| 45.0 |
| 27.7 |
| 41.38 |
| 50.9 |
| 75 |
| 28.34 |
+-----+
only showing top 5 rows
```

## The sum of all the accountlengt registered

```
In [58]: df.agg({'AccountLength': 'max'}).show()

df.agg({'AccountLength': 'sum'}).show(truncate=False)
df.select(max('AccountLength')).show()
```

```
Out[58]: [306040.9]
```

```
In [59]: df.select(['Churn', 'IntlPlan', 'VMailPlan']).show()
```

```
Out[59]: DataFrame[Churn: string, IntlPlan: string, VMailPlan: string]
```

```
In [6]: from pyspark.sql.functions import col

df_schema = df.withColumn('Churn', col('Churn').cast(IntegerType())) \
    .withColumn('IntlPlan', col('IntlPlan').cast(IntegerType())) \
    .withColumn('VMailPlan', col('VMailPlan').cast(IntegerType())) \
    .withColumn('VMailMessage', col('VMailMessage').cast(IntegerType()))
df_schema.printSchema()
```

```
root
 |-- AccountLength: integer (nullable = true)
 |-- VMailMessage: integer (nullable = true)
 |-- DayMins: double (nullable = true)
 |-- EveMins: double (nullable = true)
 |-- NightMins: double (nullable = true)
 |-- IntMins: double (nullable = true)
 |-- CustServCalls: integer (nullable = true)
 |-- Churn: integer (nullable = true)
 |-- IntPlan: integer (nullable = true)
 |-- VMailPlan: integer (nullable = true)
 |-- DayCalls: integer (nullable = true)
 |-- DayCharge: double (nullable = true)
 |-- EveCalls: integer (nullable = true)
 |-- EveCharge: double (nullable = true)
 |-- NightCalls: integer (nullable = true)
 |-- NightCharge: double (nullable = true)
 |-- IntCalls: integer (nullable = true)
 |-- IntCharge: double (nullable = true)
 |-- State: string (nullable = true)
 |-- AreaCode: integer (nullable = true)
 |-- Phone: string (nullable = true)
```

```
In [7]: df.count()
```

```
Out[7]: 3333
```

```
In [8]: df_schema.describe()
```

```
Out[8]: DataFrame[Summary: string, AccountLength: string, VMailMessage: string, DayMins: string, EveMins: string, NightMins: string, IntMins: string, CustServCalls: string, Churn: string, IntPlan: string, VMailPlan: string, DayCalls: string, DayCharge: string, EveCalls: string, EveCharge: string, NightCalls: string, NightCharge: string, IntCalls: string, IntCharge: string, State: string, AreaCode: string, Phone: string]
```

```
In [9]: #VMail: Number of voicemail messages sent by the customer.
df.select('VMailMessage').show(10)
```

```
Out[9]: [25]
[26]
[0]
[0]
[24]
[0]
[37]
only showing top 10 rows
```

```
In [10]: #convert the column to int and calculate the totaal VMailMessage call in a row
df.agg({'VMailMessage': 'sum'}).show()
df.agg({'VMailMessage': 'max'}).show()
```

```
Out[10]: [26994]
[51]
only showing top 20 rows
```

```
In [65]: #Day.Mins: Time spent on day calls.
df.select(['DayMins', 'DayCalls']).show()
```

```
Out[65]: [265.1]
[161.6]
[299.4]
[166.7]
[223.4]
[218.2]
[184.5]
[258.6]
[123.2]
[187.7]
[128.8]
[156.6]
[128.7]
[332.9]
[196.4]
[198.7]
[189.7]
[224.4]
only showing top 20 rows
```

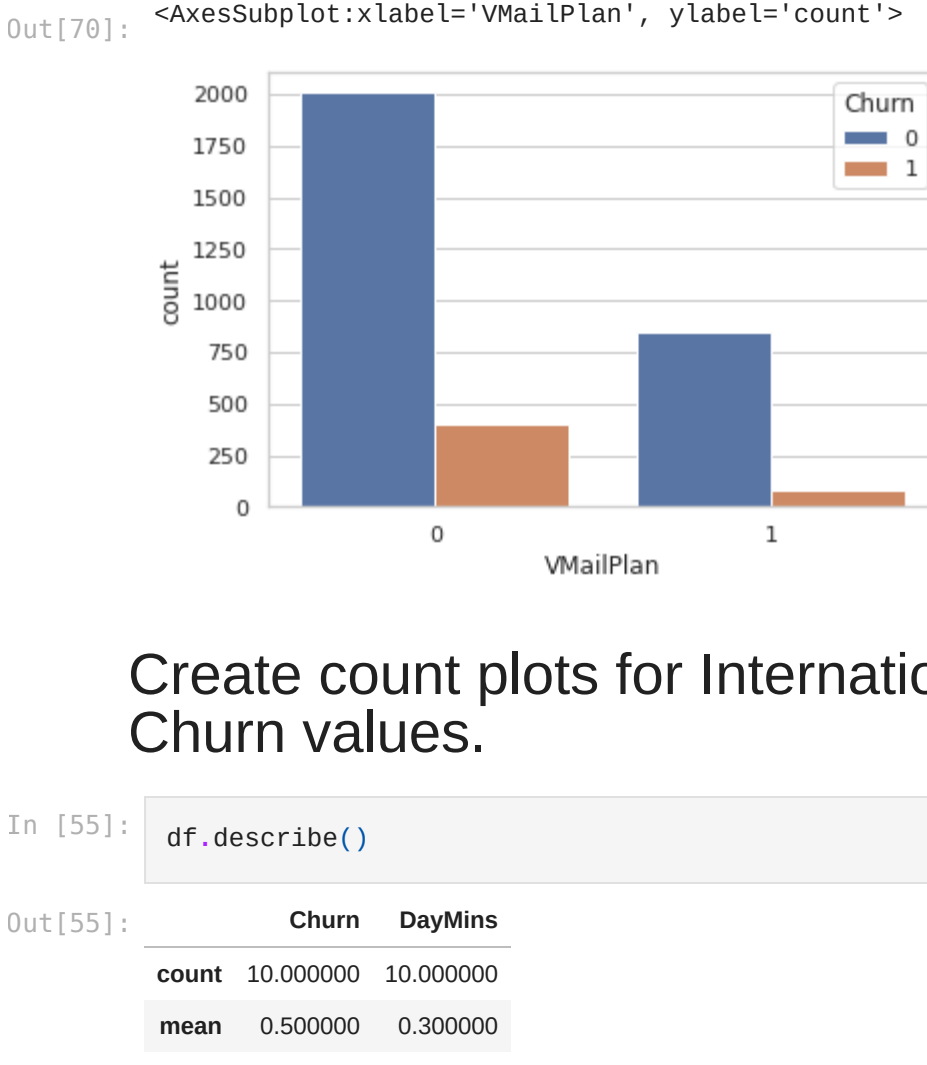
```
In [66]: df.groupby('DayMins').agg(['DayCalls': 'min']).show()
```

```
Out[66]: [0]
[180]
[100.1]
[106.8]
[101.1]
[101.2]
[101.4]
[101.9]
[102.1]
[102.3]
[102.7]
[102.8]
[103.1]
[103.2]
[103.3]
[103.4]
[103.6]
[103.8]
[103.9]
[104.1]
[104.2]
[104.3]
[104.4]
[104.5]
[104.6]
[104.7]
[104.8]
[104.9]
[105.0]
[105.1]
[105.2]
[105.3]
[105.4]
[105.5]
[105.6]
[105.7]
[105.8]
[105.9]
[106.0]
[106.1]
[106.2]
[106.3]
[106.4]
[106.5]
[106.6]
[106.7]
[106.8]
[106.9]
[107.0]
[107.1]
[107.2]
[107.3]
[107.4]
[107.5]
[107.6]
[107.7]
[107.8]
[107.9]
[108.0]
[108.1]
[108.2]
[108.3]
[108.4]
[108.5]
[108.6]
[108.7]
[108.8]
[108.9]
[109.0]
[109.1]
[109.2]
[109.3]
[109.4]
[109.5]
[109.6]
[109.7]
[109.8]
[109.9]
[110.0]
[110.1]
[110.2]
[110.3]
[110.4]
[110.5]
[110.6]
[110.7]
[110.8]
[110.9]
[111.0]
[111.1]
[111.2]
[111.3]
[111.4]
[111.5]
[111.6]
[111.7]
[111.8]
[111.9]
[112.0]
[112.1]
[112.2]
[112.3]
[112.4]
[112.5]
[112.6]
[112.7]
[112.8]
[112.9]
[113.0]
[113.1]
[113.2]
[113.3]
[113.4]
[113.5]
[113.6]
[113.7]
[113.8]
[113.9]
[114.0]
[114.1]
[114.2]
[114.3]
[114.4]
[114.5]
[114.6]
[114.7]
[114.8]
[114.9]
[115.0]
[115.1]
[115.2]
[115.3]
[115.4]
[115.5]
[115.6]
[115.7]
[115.8]
[115.9]
[116.0]
[116.1]
[116.2]
[116.3]
[116.4]
[116.5]
[116.6]
[116.7]
[116.8]
[116.9]
[117.0]
[117.1]
[117.2]
[117.3]
[117.4]
[117.5]
[117.6]
[117.7]
[117.8]
[117.9]
[118.0]
[118.1]
[118.2]
[118.3]
[118.4]
[118.5]
[118.6]
[118.7]
[118.8]
[118.9]
[119.0]
[119.1]
[119.2]
[119.3]
[119.4]
[119.5]
[119.6]
[119.7]
[119.8]
[119.9]
[120.0]
[120.1]
[120.2]
[120.3]
[120.4]
[120.5]
[120.6]
[120.7]
[120.8]
[120.9]
[121.0]
[121.1]
[121.2]
[121.3]
[121.4]
[121.5]
[121.6]
[121.7]
[121.8]
[121.9]
[122.0]
[122.1]
[122.2]
[122.3]
[122.4]
[122.5]
[122.6]
[122.7]
[122.8]
[122.9]
[123.0]
[123.1]
[123.2]
[123.3]
[123.4]
[123.5]
[123.6]
[123.7]
[123.8]
[123.9]
[124.0]
[124.1]
[124.2]
[124.3]
[124.4]
[124.5]
[124.6]
[124.7]
[124.8]
[124.9]
[125.0]
[125.1]
[125.2]
[125.3]
[125.4]
[125.5]
[125.6]
[125.7]
[125.8]
[125.9]
[126.0]
[126.1]
[126.2]
[126.3]
[126.4]
[126.5]
[126.6]
[126.7]
[126.8]
[126.9]
[127.0]
[127.1]
[127.2]
[127.3]
[127.4]
[127.5]
[127.6]
[127.7]
[127.8]
[127.9]
[128.0]
[128.1]
[128.2]
[128.3]
[128.4]
[128.5]
[128.6]
[128.7]
[128.8]
[128.9]
[129.0]
[129.1]
[129.2]
[129.3]
[129.4]
[129.5]
[129.6]
[129.7]
[129.8]
[129.9]
[130.0]
[130.1]
[130.2]
[130.3]
[130.4]
[130.5]
[130.6]
[130.7]
[130.8]
[130.9]
[131.0]
[131.1]
[131.2]
[131.3]
[131.4]
[131.5]
[131.6]
[131.7]
[131.8]
[131.9]
[132.0]
[132.1]
[132.2]
[132.3]
[132.4]
[132.5]
[132.6]
[132.7]
[132.8]
[132.9]
[133.0]
[133.1]
[133.2]
[133.3]
[133.4]
[133.5]
[133.6]
[133.7]
[133.8]
[133.9]
[134.0]
[134.1]
[134.2]
[134.3]
[134.4]
[134.5]
[134.6]
[134.7]
[134.8]
[134.9]
[135.0]
[135.1]
[135.2]
[135.3]
[135.4]
[135.5]
[135.6]
[135.7]
[135.8]
[135.9]
[136.0]
[136.1]
[136.2]
[136.3]
[136.4]
[136.5]
[136.6]
[136.7]
[136.8]
[136.9]
[137.0]
[137.1]
[137.2]
[137.3]
[137.4]
[137.5]
[137.6]
[137.7]
[137.8]
[137.9]
[138.0]
[138.1]
[138.2]
[138.3]
[138.4]
[138.5]
[138.6]
[138.7]
[138.8]
[138.9]
[139.0]
[139.1]
[139.2]
[139.3]
[139.4]
[139.5]
[139.6]
[139.7]
[139.8]
[139.9]
[140.0]
[140.1]
[140.2]
[140.3]
[140.4]
[140.5]
[140.6]
[140.7]
[140.8]
[140.9]
[141.0]
[141.1]
[141.2]
[141.3]
[141.4]
[141.5]
[141.6]
[141.7]
[141.8]
[141.9]
[142.0]
[142.1]
[142.2]
[142.3]
[142.4]
[142.5]
[142.6]
[142.7]
[142.8]
[142.9]
[143.0]
[143.1]
[143.2]
[143.3]
[143.4]
[143.5]
[143.6]
[143.7]
[143.8]
[143.9]
[144.0]
[144.1]
[144.2]
[144.3]
[144.4]
[144.5]
[144.6]
[144.7]
[144.8]
[144.9]
[145.0]
[145.1]
[145.2]
[145.3]
[145.4]
[145.5]
[145.6]
[145.7]
[145.8]
[145.9]
[146.0]
[146.1]
[146.2]
[146.3]
[146.4]
[146.5]
[146.6]
[146.7]
[146.8]
[146.9]
[147.0]
[147.1]
[147.2]
[147.3]
[147.4]
[147.5]
[147.6]
[147.7]
[147.8]
[147.9]
[148.0]
[148.1]
[148.2]
[148.3]
[148.4]
[148.5]
[148.6]
[148.7]
[148.8]
[148.9]
[149.0]
[149.1]
[149.2]
[149.3]
[149.4]
[149.5]
[149.6]
[149.7]
[149.8]
[149.9]
[150.0]
[150.1]
[150.2]
[150.3]
[150.4
```

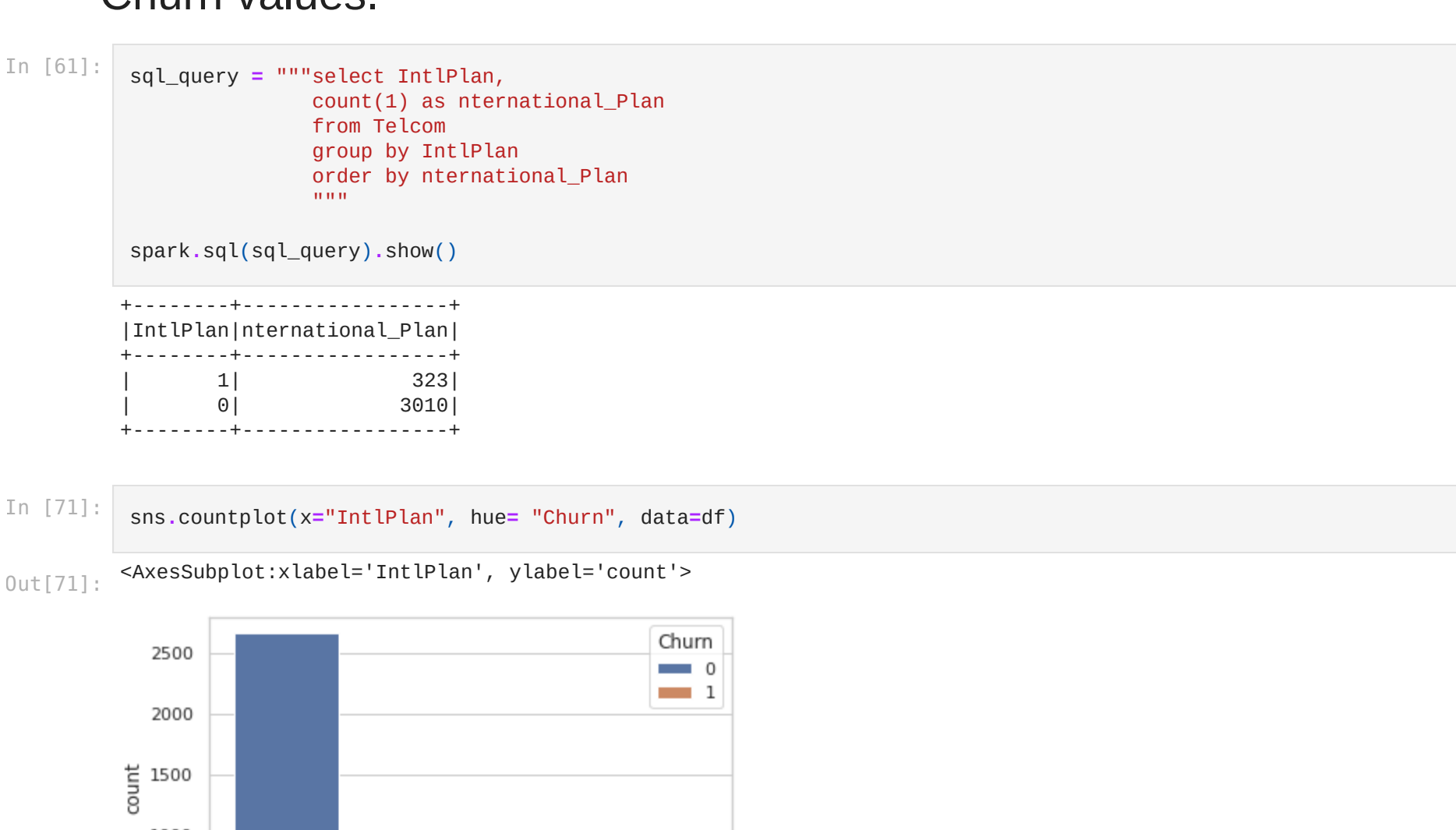




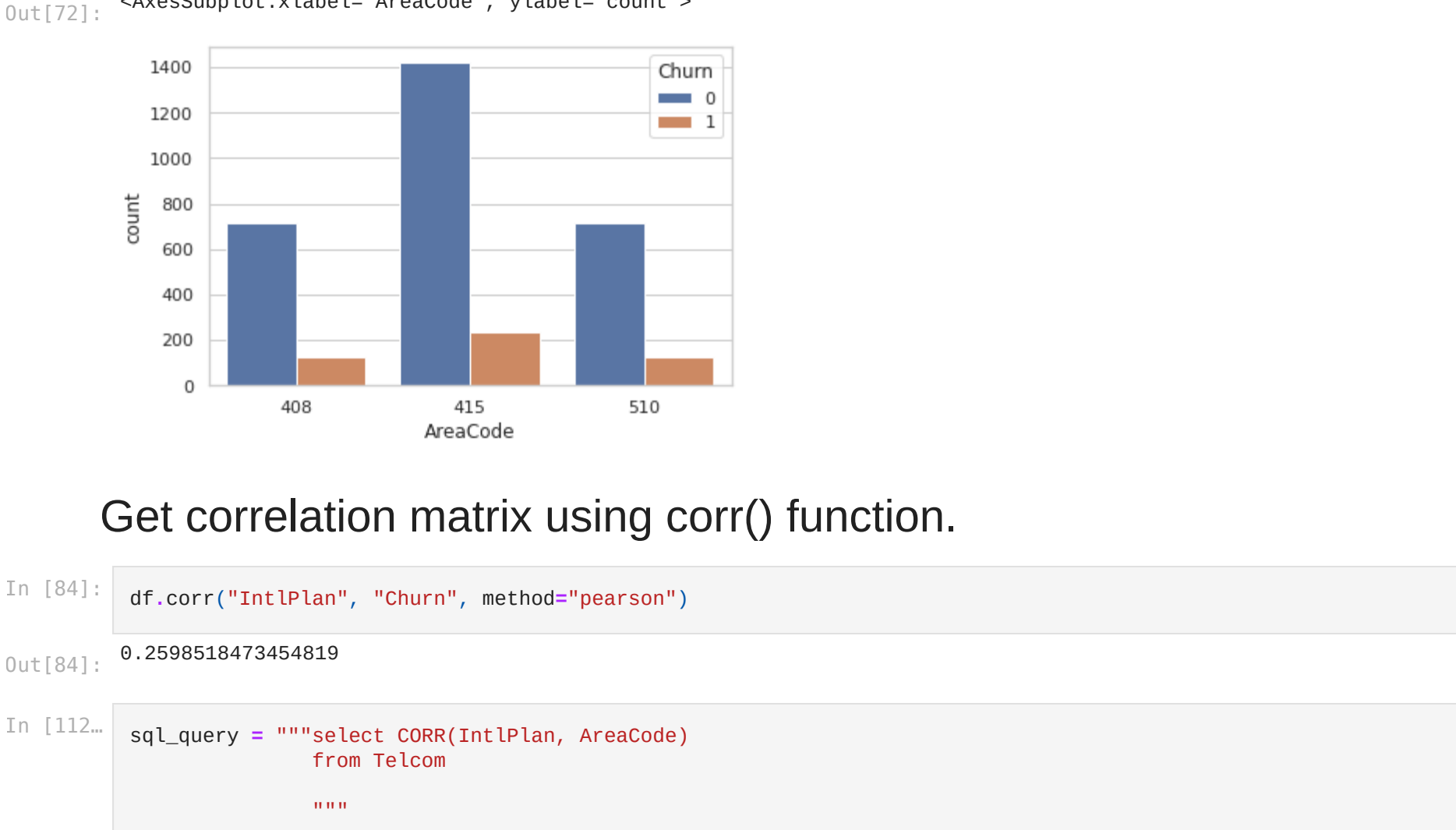




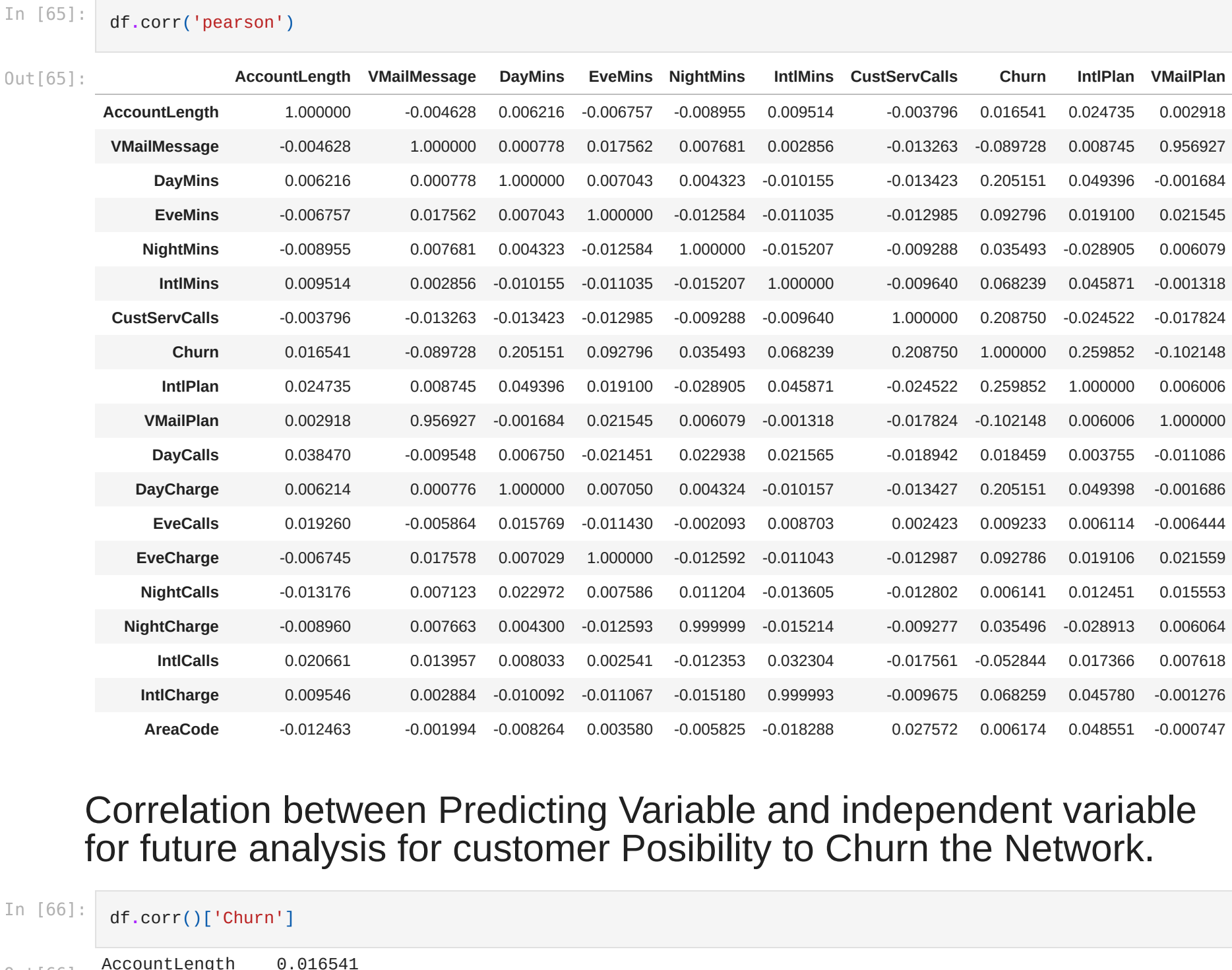
Create count plots for International Plan opt by the customer with Churn values.



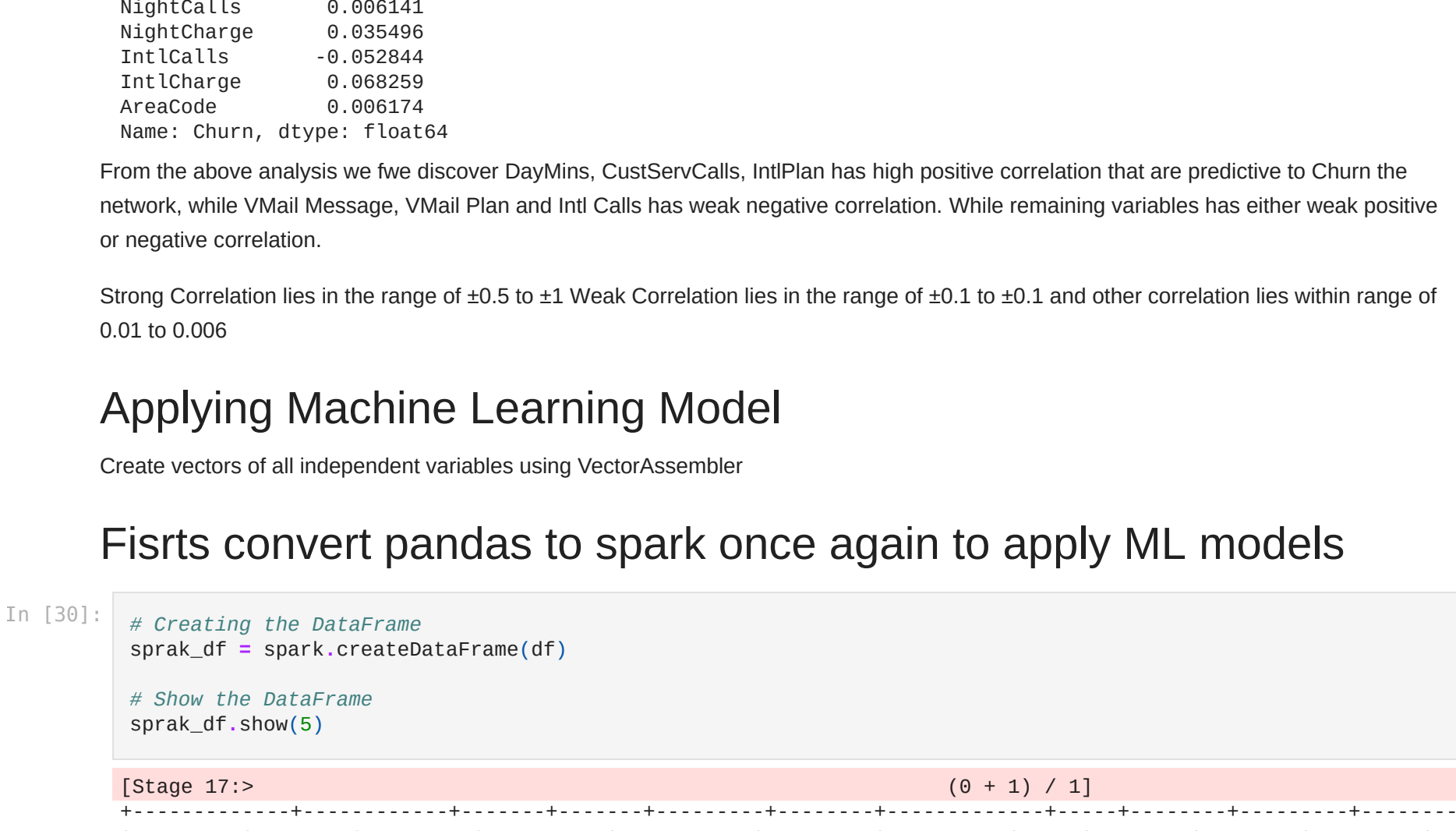
Create count plots for International Plan opt by the customer with Churn values.



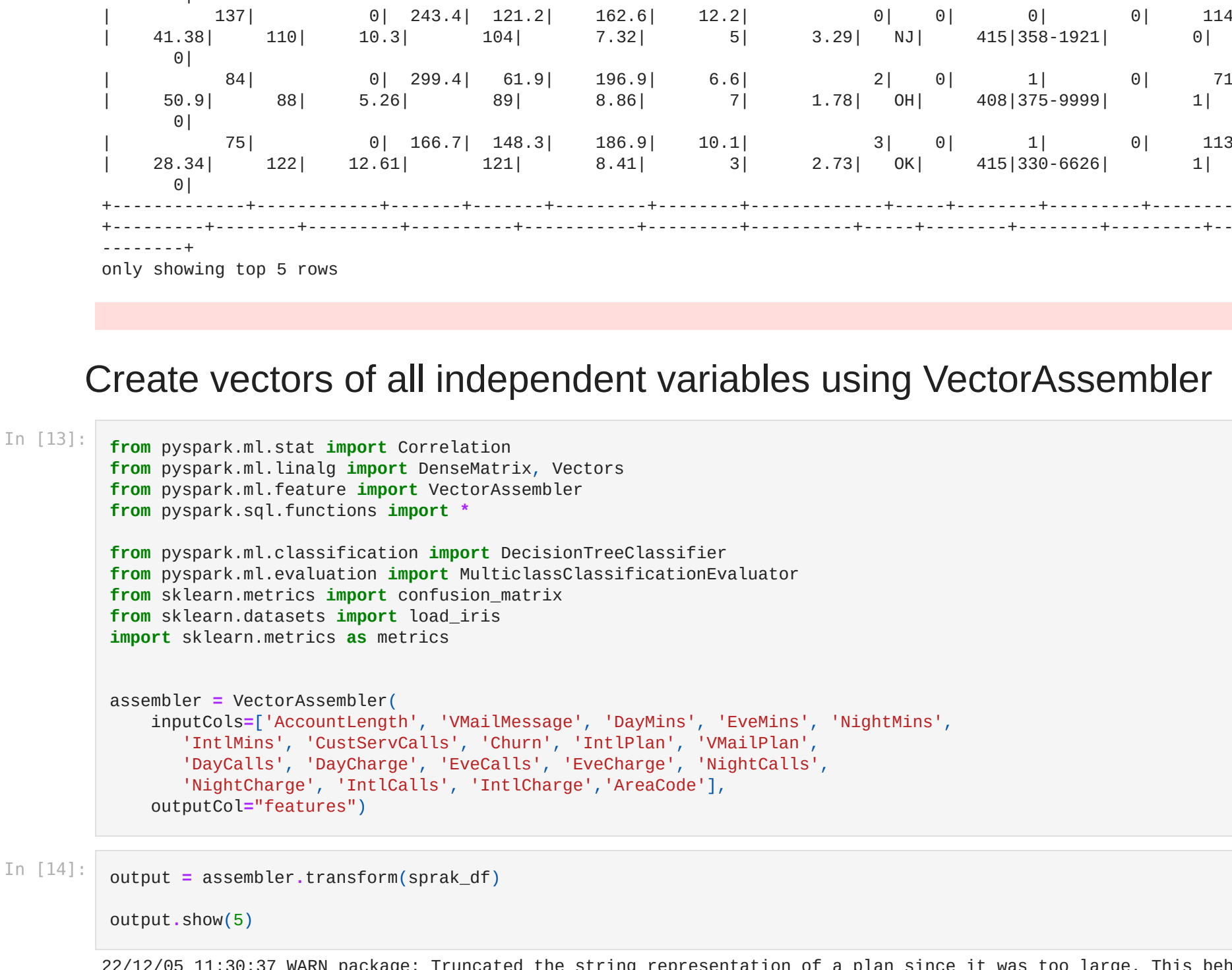
Plot Area Wise churner and non-churner.



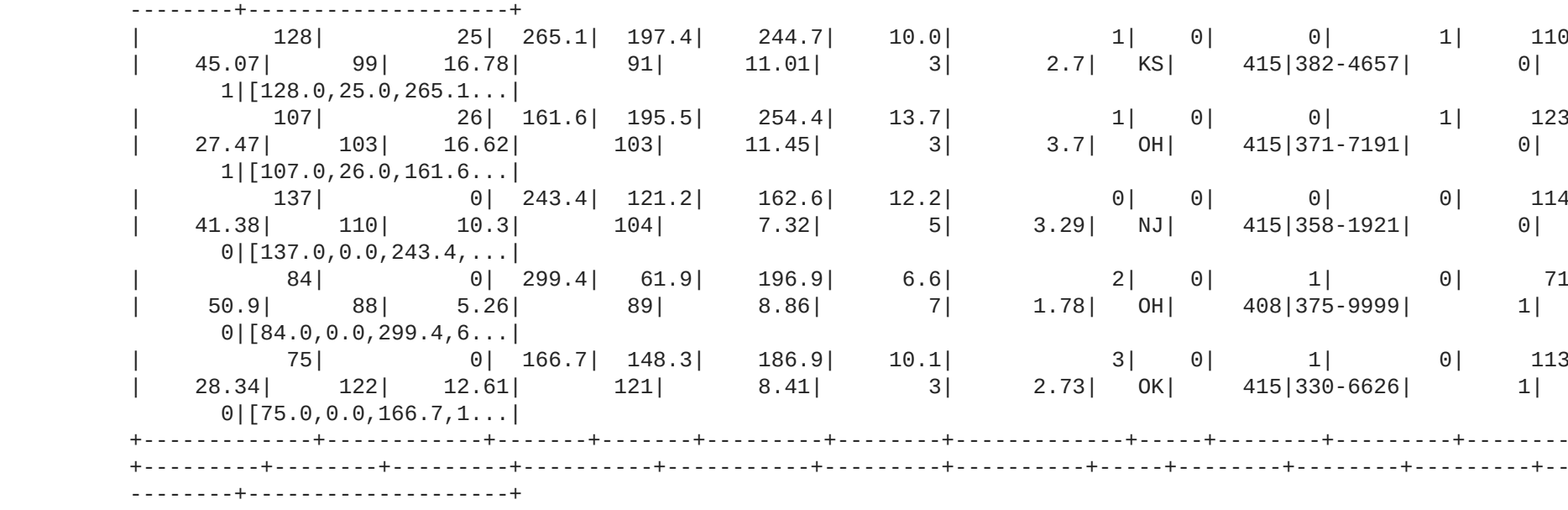
Get correlation matrix using corr() function.



Lets check the Correlation Matrix



Correlation between Predicting Variable and independent variable for future analysis for customer Possibility to Churn the Network.



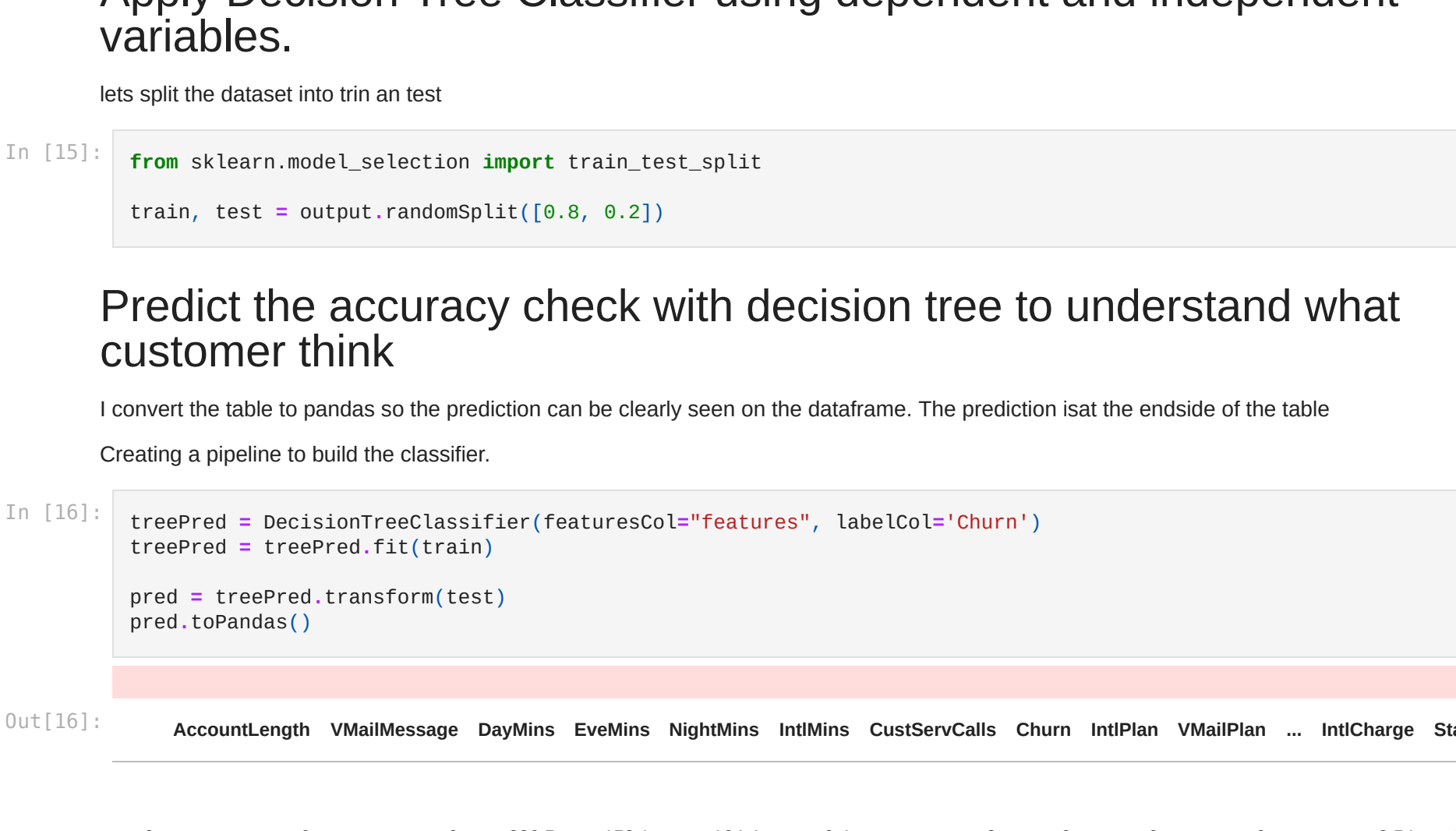
From the above analysis we have discover DayMins, CustServCalls, IntlPlan has high positive correlation that are predictive to Churn the network, while VMail Message, VMail Plan and Intl Calls has weak negative correlation. While remaining variables has either weak positive or negative correlation.

Strong Correlation lies in the range of  $\pm 0.5$  to  $\pm 1$  Weak Correlation lies in the range of  $\pm 0.1$  to  $\pm 0.1$  and other correlation lies within range of 0.01 to 0.006

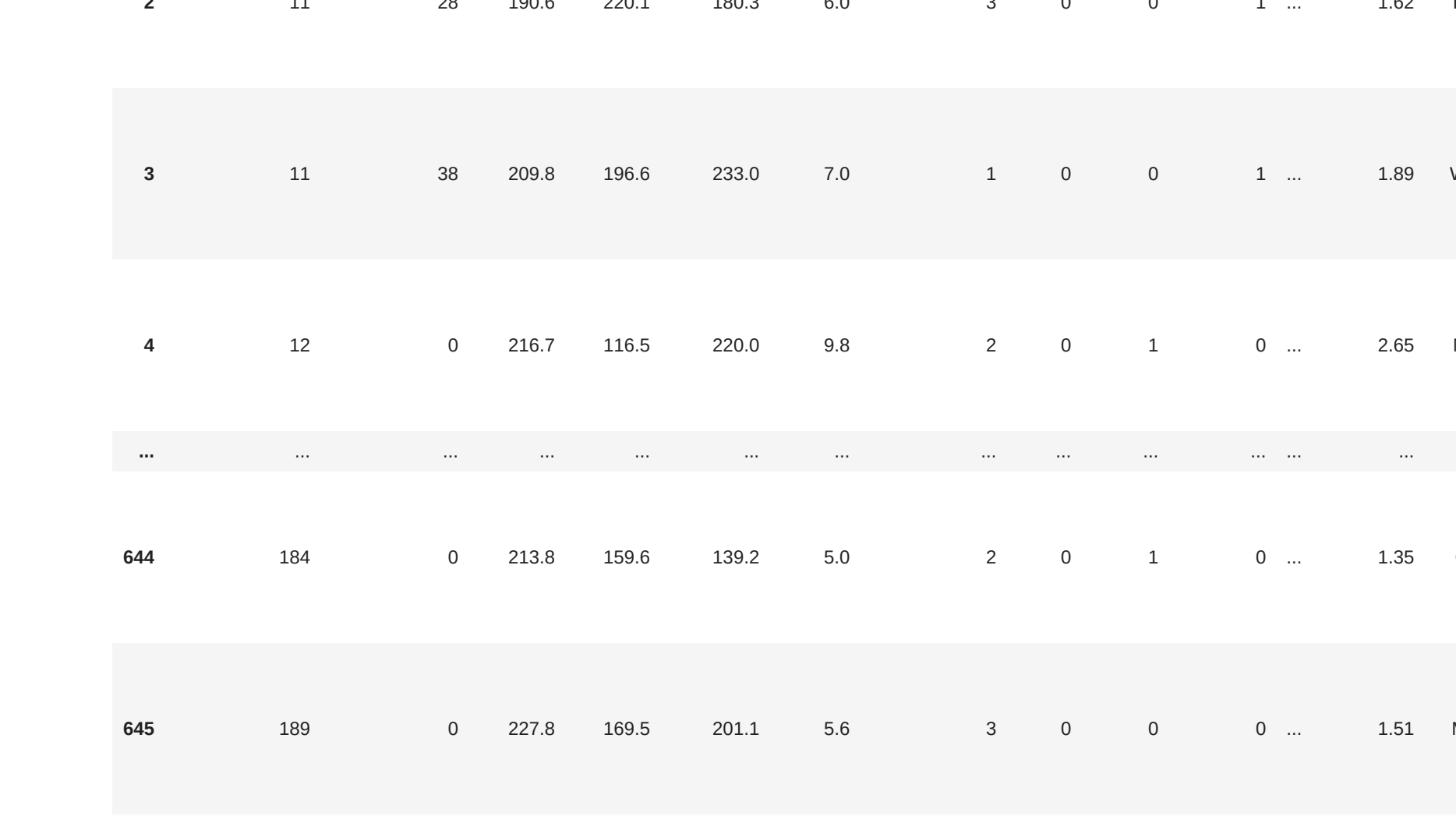
Applying Machine Learning Model

Create vectors of all independent variables using VectorAssembler

Firsrts convert pandas to spark once again to apply ML models



Create vectors of all independent variables using VectorAssembler

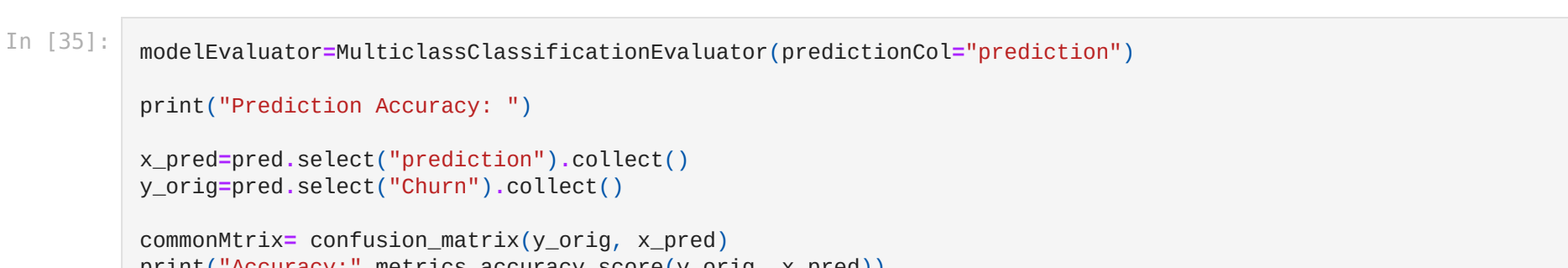


Each independent Variables from the Dataframe printed first 5 features



Apply Decision Tree Classifier using dependent and independent variables.

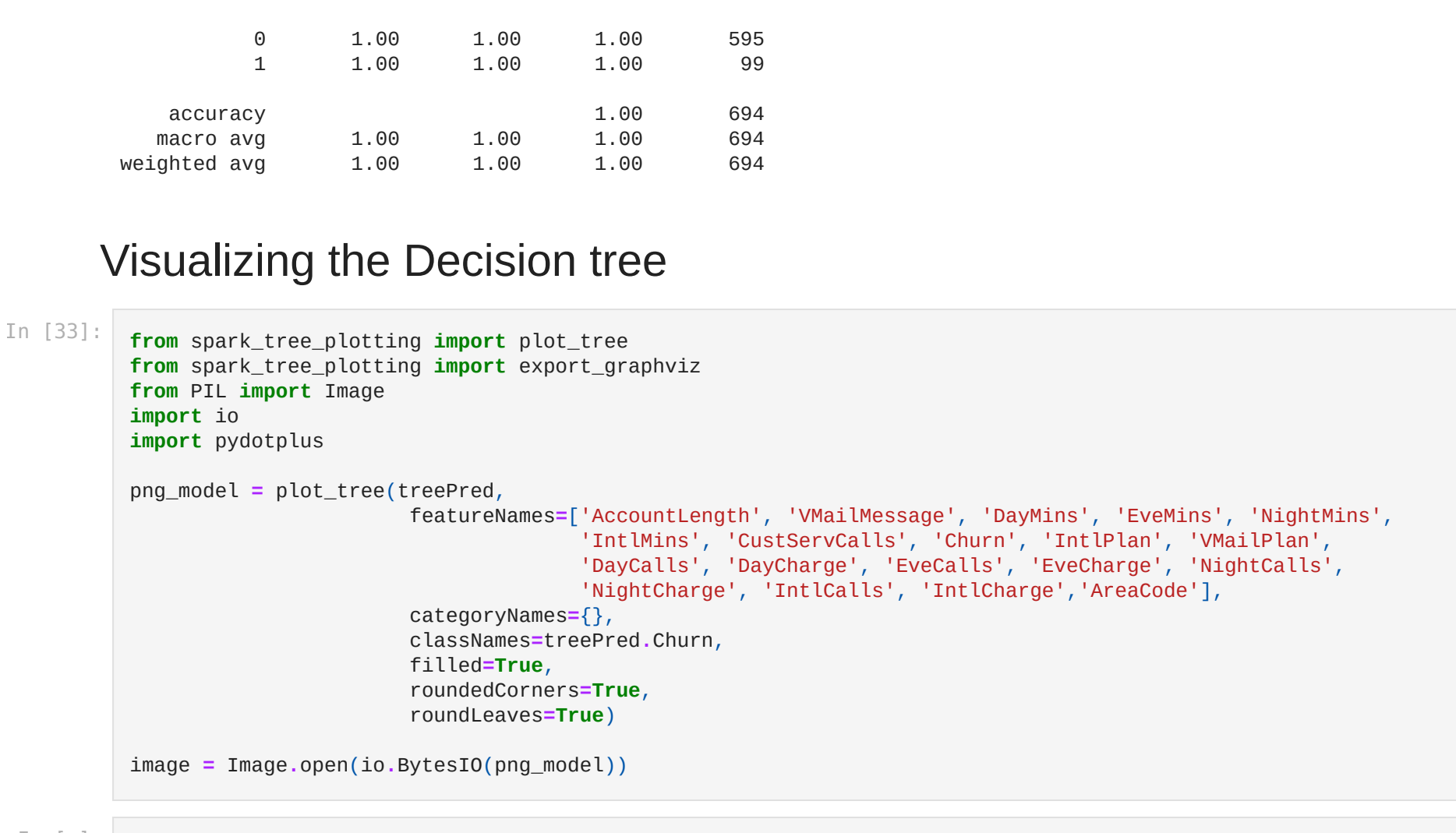
lets split the dataset into train and test



Predict the accuracy check with decision tree to understand what customer think

I convert the table to pandas so the prediction can be clearly seen on the dataframe. The prediction is at the end of the table

Creating a pipeline to build the classifier.



Evaluating the Model with muticlassification matrix

After going through the models, we can see the predicted test data and accuracy metrics. we used the function

MulticlassClassificationEvaluator from pyspark to check the accuracy matrix and the confusion matrix can also be created using confusion\_matrix function from sklearn.metrics module.



Using metric classification method to alculate recall, precision, score and individual support which gives us the total sample size.

