

# Predictive Analysis

Import libraries and packages

```
In [24]: import pandas as pd
import numpy as np
from sklearn.model_selection import train_test_split
from sklearn.ensemble import RandomForestRegressor
from sklearn.metrics import mean_squared_error, r2_score

In [25]: print('pandas version ' + pd.__version__)
print('numpy version ' + np.__version__)
```

pandas version 1.1.3  
numpy version 1.19.2

Load and check data set

```
In [26]: df = pd.read_csv('/Users/ebeth/Desktop/Churn Data/churn_clean.csv')

In [27]: df.head()
```

	CaseOrder	Customer_id	Interaction	UID	City	State	County	Zip	Lat	Longitude
0	1	K409198	aa90260b-4141-4a24-8e36-b04ce1f4f77b	e885b299883d4f9fb18e39c75155d990	Point Baker	AK	Prince of Wales-Hyder	99927	56.25100	-133.375
1	2	S120509	fb76459f-c047-4a9d-8af9-e0f7d4ac2524	f2de8bef964785f41a2959829830fb8a	West Branch	MI	Ogemaw	48661	44.32893	-84.240
2	3	K191035	344d114c-3736-4be5-98f7-c72c281e2d35	f1784cfa9f6d92ae816197eb175d3c71	Yamhill	OR	Yamhill	97148	45.35589	-123.246
3	4	D90850	abfa2b40-2d43-4994-b15a-989b8c79e311	dc8a365077241bb5cd5ccd305136b05e	Del Mar	CA	San Diego	92014	32.96687	-117.247
4	5	K662701	68a861fd-0d20-4e51-a587-8a90407ee574	aabb64a116e83fdc4bfc1fbab1663f9	Needville	TX	Fort Bend	77461	29.38012	-95.806

5 rows x 50 columns

```
In [28]: df.describe()
```

	CaseOrder	Zip	Lat	Lng	Population	Children	Age	Income	Outage_sec
count	10000.00000	10000.000000	10000.000000	10000.000000	10000.000000	10000.0000	10000.000000	10000.000000	10000.000000
mean	5000.50000	49153.319600	38.757567	-90.782536	9756.562400	2.0877	53.078400	39806.926771	10000.000000
std	2886.89568	27532.196108	5.437389	15.156142	14432.698671	2.1472	20.698882	28199.916702	10000.000000
min	1.00000	601.000000	17.966120	-171.688150	0.000000	0.0000	18.000000	348.670000	10000.000000
25%	2500.75000	26292.500000	35.341828	-97.082813	738.000000	0.0000	35.000000	19224.717500	10000.000000
50%	5000.50000	48869.500000	39.395800	-87.918800	2910.500000	1.0000	53.000000	33170.605000	10000.000000
75%	7500.25000	71866.500000	42.106908	-80.088745	13168.000000	3.0000	71.000000	53246.170000	10000.000000
max	10000.00000	99929.000000	70.640660	-65.667850	111850.000000	10.0000	89.000000	258900.700000	10000.000000

8 rows x 23 columns

## Data Preparation

Check for null values

```
In [29]: df.isnull().values.any()

Out[29]: False
```

Check for duplicates

```
In [30]: df.duplicated().values.any()

Out[30]: False
```

Drop unused columns

```
In [31]: df2 = df.drop(['CaseOrder', 'Customer_id', 'Interaction', 'UID', 'Lat', 'Lng', 'TimeZone', 'Job', 'City', 'County', 'State'])

In [32]: df2.head()
```

	CaseOrder	Customer_id	Interaction	UID	City	State	County	Zip	Lat	Longitude
0	1	K409198	aa90260b-4141-4a24-8e36-b04ce1f4f77b	e885b299883d4f9fb18e39c75155d990	Point Baker	AK	Prince of Wales-Hyder	99927	56.25100	-133.375
1	2	S120509	fb76459f-c047-4a9d-8af9-e0f7d4ac2524	f2de8bef964785f41a2959829830fb8a	West Branch	MI	Ogemaw	48661	44.32893	-84.240
2	3	K191035	344d114c-3736-4be5-98f7-c72c281e2d35	f1784cfa9f6d92ae816197eb175d3c71	Yamhill	OR	Yamhill	97148	45.35589	-123.246
3	4	D90850	abfa2b40-2d43-4994-b15a-989b8c79e311	dc8a365077241bb5cd5ccd305136b05e	Del Mar	CA	San Diego	92014	32.96687	-117.247
4	5	K662701	68a861fd-0d20-4e51-a587-8a90407ee574	aabb64a116e83fdc4bfc1fbab1663f9	Needville	TX	Fort Bend	77461	29.38012	-95.806

5 rows x 50 columns

```
In [33]: df2.info()

<class 'pandas.core.frame.DataFrame'>
RangeIndex: 10000 entries, 0 to 9999
Data columns (total 30 columns):
 #   Column                                Non-Null Count  Dtype
---  -
 0   Population                            10000 non-null  int64
 1   Area                                  10000 non-null  object
 2   Children                             10000 non-null  int64
 3   Age                                   10000 non-null  int64
 4   Income                               10000 non-null  float64
 5   Marital                             10000 non-null  object
 6   Gender                               10000 non-null  object
 7   Churn                                10000 non-null  object
 8   Outage_sec_perweek                   10000 non-null  float64
 9   Email                                10000 non-null  int64
10   Contacts                             10000 non-null  int64
11   Yearly equip failure                  10000 non-null  int64
12   Techie                               10000 non-null  object
13   Contract                             10000 non-null  object
14   Port_modem                           10000 non-null  object
15   Tablet                                10000 non-null  object
16   InternetService                      10000 non-null  object
17   Phone                                10000 non-null  object
18   Multiple                             10000 non-null  object
19   OnlineSecurity                       10000 non-null  object
20   OnlineBackup                         10000 non-null  object
21   DeviceProtection                     10000 non-null  object
22   TechSupport                          10000 non-null  object
23   StreamingTV                          10000 non-null  object
24   StreamingMovies                      10000 non-null  object
25   PaperlessBilling                     10000 non-null  object
26   PaymentMethod                        10000 non-null  object
27   Tenure                               10000 non-null  float64
28   MonthlyCharge                        10000 non-null  float64
29   Bandwidth_GB_Year                    10000 non-null  float64
dtypes: float64(5), int64(6), object(19)
memory usage: 2.3+ MB
```

Create Dummy variables for all categorical columns and drop unneeded columns. (code used from: <https://towardsdatascience.com/the-dummys-guide-to-creating-dummy-variables-f21faddb1d40>)

```
In [34]: dummy1 = pd.get_dummies(df2.Area, prefix = 'Area', drop_first = True)
dummy2 = pd.get_dummies(df2.Marital, prefix = 'Marital', drop_first = True)
dummy3 = pd.get_dummies(df2.Gender, prefix = 'Gender', drop_first = True)
dummy4 = pd.get_dummies(df2.Churn, prefix = 'Churn', drop_first = True)
dummy5 = pd.get_dummies(df2.Techie, prefix = 'Techie', drop_first = True)
dummy6 = pd.get_dummies(df2.Contract, prefix = 'Contract', drop_first = True)
dummy7 = pd.get_dummies(df2.Port_modem, prefix = 'Port_modem', drop_first = True)
dummy8 = pd.get_dummies(df2.Tablet, prefix = 'Tablet', drop_first = True)
dummy9 = pd.get_dummies(df2.InternetService, prefix = 'InternetService', drop_first = True)
dummy10 = pd.get_dummies(df2.Phone, prefix = 'Phone', drop_first = True)
dummy11 = pd.get_dummies(df2.Multiple, prefix = 'Multiple', drop_first = True)
dummy12 = pd.get_dummies(df2.OnlineSecurity, prefix = 'OnlineSecurity', drop_first = True)
dummy13 = pd.get_dummies(df2.OnlineBackup, prefix = 'OnlineBackup', drop_first = True)
dummy14 = pd.get_dummies(df2.DeviceProtection, prefix = 'DeviceProtection', drop_first = True)
dummy15 = pd.get_dummies(df2.TechSupport, prefix = 'TechSupport', drop_first = True)
dummy16 = pd.get_dummies(df2.StreamingTV, prefix = 'StreamingTV', drop_first = True)
dummy17 = pd.get_dummies(df2.StreamingMovies, prefix = 'StreamingMovies', drop_first = True)
dummy18 = pd.get_dummies(df2.PaperlessBilling, prefix = 'PaperlessBilling', drop_first = True)
dummy19 = pd.get_dummies(df2.PaymentMethod, prefix = 'PaymentMethod', drop_first = True)

df2 = df2.drop(columns = 'Area').merge(dummy1, left_index = True, right_index = True)
df2 = df2.drop(columns = 'Marital').merge(dummy2, left_index = True, right_index = True)
df2 = df2.drop(columns = 'Gender').merge(dummy3, left_index = True, right_index = True)
df2 = df2.drop(columns = 'Churn').merge(dummy4, left_index = True, right_index = True)
df2 = df2.drop(columns = 'Techie').merge(dummy5, left_index = True, right_index = True)
df2 = df2.drop(columns = 'Contract').merge(dummy6, left_index = True, right_index = True)
df2 = df2.drop(columns = 'Port_modem').merge(dummy7, left_index = True, right_index = True)
df2 = df2.drop(columns = 'Tablet').merge(dummy8, left_index = True, right_index = True)
df2 = df2.drop(columns = 'InternetService').merge(dummy9, left_index = True, right_index = True)
df2 = df2.drop(columns = 'Phone').merge(dummy10, left_index = True, right_index = True)
df2 = df2.drop(columns = 'Multiple').merge(dummy11, left_index = True, right_index = True)
df2 = df2.drop(columns = 'OnlineSecurity').merge(dummy12, left_index = True, right_index = True)
df2 = df2.drop(columns = 'OnlineBackup').merge(dummy13, left_index = True, right_index = True)
df2 = df2.drop(columns = 'DeviceProtection').merge(dummy14, left_index = True, right_index = True)
df2 = df2.drop(columns = 'TechSupport').merge(dummy15, left_index = True, right_index = True)
df2 = df2.drop(columns = 'StreamingTV').merge(dummy16, left_index = True, right_index = True)
df2 = df2.drop(columns = 'StreamingMovies').merge(dummy17, left_index = True, right_index = True)
df2 = df2.drop(columns = 'PaperlessBilling').merge(dummy18, left_index = True, right_index = True)
df2 = df2.drop(columns = 'PaymentMethod').merge(dummy19, left_index = True, right_index = True)
```

Create copy of prepared data

```
In [35]: df2.to_csv('predictive_prepared_churn.csv')
```

## Analysis

Split the data set into training and test sets

```
In [36]: y = df2['MonthlyCharge']
X = df2.loc[:, df2.columns != 'MonthlyCharge']

In [37]: X_train, X_test, y_train, y_test = train_test_split(X, y, test_size = 0.3, random_state = 13)
```

Create copy of training and test data

```
In [38]: X_train.to_csv('X_train.csv')
X_test.to_csv('X_test.csv')
y_train.to_csv('y_train.csv')
y_test.to_csv('y_test.csv')
```

Build the model

```
In [39]: rf = RandomForestRegressor(random_state = 13)
rf.fit(X_train, y_train)
y_pred = rf.predict(X_test)
r_squared = r2_score(y_test, y_pred)
mse_test = mean_squared_error(y_test, y_pred)
rmse_test = mse_test**(1/2)
print('Test set R squared score: {:.2f}'.format(r_squared))
print('Test set Mean Squared Error: {:.2f}'.format(mse_test))
print('Test set Root Mean Squared Error: {:.2f}'.format(rmse_test))
```

Test set R squared score: 0.999990  
Test set Mean Squared Error: 0.018317  
Test set Root Mean Squared Error: 0.135339

```
In [ ]:
```

```
In [ ]:
```

```
In [ ]:
```

```
In [ ]:
```

```
In [ ]:
```