


```
#fill na
df_clean[isna(df_clean)].fillna('No', inplace=True)
#built in unit tests
print(df_clean[col].Value_counts(), len(df_clean[col]))
```

```
No      8743
Yes     1257
Name: Techie, dtype: int64 10000
Yes     8128
No       1872
Name: Phone, dtype: int64 10000
No        6626
Yes       3374
Name: TechSupport, dtype: int64 10000
```

Issue 4

Empty unnamed first column that corresponds to index should be dropped and Case Order should be set to index

Using drop and set_index, a proper index will be made. This de-dupes data and structures it better.

```
In [47]: # View Unnamed column
df_clean.iloc[:, 0][1:].head()
```

```
Out[47]:
Unnamed: 0
0          1
1          2
2          3
3          4
4          5
```

```
In [48]: # Drop column inplace
df_clean.drop(df_clean.columns[0],axis=1,inplace=True)
```

```
In [49]: # drop and reset index
df_clean = df_clean.reset_index(drop=True)
df_clean.set_index("CaseOrder",inplace=True)
```

```
In [50]: #unit test verification
df_clean.head()
```

```
[In] [50]: Out[50]:
   Customer_id  Interaction    City State County Zip Lat Lng Population Area Timezone
CaseOrder
1 K409198 ae905026- Point Baker Prince of Wales-Hyder 99927 56.251 -133.37571 38 Urban America/Sitka Alaska 1
2 S120509 bf76459f- West Branch MI Ogemaw 48661 44.32893 -84.2408 10446 Urban America/Detroit Michigan 1
3 K191035 344d114c- Yamhill OR Yamhill 97148 45.35589 -123.24657 3735 Urban America/Los_Angeles Oregon 1
4 D908050 adfa2b40- Del Mar CA San Diego 92014 32.96687 -117.24798 13863 Suburban America/Los_Angeles California 1
5 K662701 b6a8e16f- Needville TX Fort Bend 77461 29.38012 -95.80673 11352 Suburban America/Chicago Texas 1
```

Issue 5

Education levels should be standardized. Using a dictionary created here, I am standardizing the customer's education levels based on US Census standards. The current list is unwieldy and too granular to make general insights.

```
In [51]: df_clean["Education"].value_counts()
Regular High School Diploma         2421
Bachelor's Degree                   1703
Some College, 1 or More Years, No Degree    1562
9th Grade to 12th Grade, No Diploma    970
Master's Degree                       764
Associate's Degree                    760
Some College, Less than 1 Year         652
Nursery School to 8th Grade           449
GED or Alternative Credential          387
Professional School Degree            198
No Diploma                           118
Doctorate Degree                      116
Name: Education, dtype: int64
```

```
In [52]: education_levels = {'Master's Degree': 'Graduate Degree',
                          'Regular High School Diploma': 'Highschool',
                          'Doctorate Degree': 'Graduate Degree',
                          'No Schooling Completed': 'No Diploma',
                          "Associate's Degree": 'Associate's/Some College',
                          'Bachelor's Degree': 'Bachelor's Degree',
                          'Some College, Less than 1 Year': 'Associate's/Some College',
                          "GED or Alternative Credential": 'Graduate Degree',
                          '9th Grade to 12th Grade, No Diploma': 'No Diploma',
                          'Nursery School to 8th Grade': 'No Diploma',
                          'Professional School Degree':'Graduate Degree'}

#use mapping to the above dictionary
df_clean['Education'] = df_clean['Education'].map(education_levels)
```

```
In [53]: df_clean["Education"].value_counts()
Associate's/Some College             2974
HighSchool                         2421
Bachelor's Degree                   1703
Graduate Degree                     1465
No Diploma                        1437
Name: Education, dtype: int64
```

Issue 6

Renaming Survey Item columns to be more descriptive. This allows easier parsing of data meaning

```
In [54]: df_clean.rename(columns={'item1':"Timeliness",
                              "item2":"Fishes",
                              "item3":"Replacements",
                              "item4":"Reliability",
                              "item5":"Options",
                              "item6":"'Respectfulness'",
                              "item7":"Courteous",
                              "item8":"'Listening'"},
                               inplace=True)
```

```
In [55]: #Unit Test
def clean_customers():
    """Clean Customers"""
    ["Age", "Population", "Area", "Timezone", "Job", "Children", "Age",
     "Education", "Employment", "Income", "Marital", "Gender", "Churn",
     "Outage_sec_perweek", "Email", "Contacts", "Yearly equip failure",
     "Techie", "Contract", "Fort_modern", "Tablet", "InternetService",
     "Multiplex", "OnlineSecurity", "OnlineBackup",
     "FraudProtection", "Techsupport", "StreamingTV", "StreamingMovies",
     "PaperlessBilling", "PaymentMethod", "Tenure", "MonthlyCharge",
     "Bandwidth_GB_Year", "Timeliness", "Fishes", "Replacements",
     "Reliability", "vgc1n8", "Respectfulness", "Courteous", "Listening"],
    dtype=object))
```

D5 Provide a copy of the cleaned data set.

The data has been cleaned and structured. A CSV with the cleaned structure will be saved.

```
In [56]: df_clean.to_csv('churn_clean.csv')
```

D6 Summarize the limitations of the data-cleaning process.

"The limitations of this dataset can be traced back to the source material. The data was voluntarily reported and therefore quality and accuracy can only be verified from non-strenuous income matching via credit checks or employment verification. My imputation techniques also may have less effective than K-Nearest Neighbors or decision tree classifiers such as random forest or a decision tree. Those methods however are more computationally expensive and may be impacted further by the quality of self-reported data. This seriously impacts the ability of our meaningful hypothesis testing or ANOVA with the current cleaned data."

Some Quality issues include age pairing with employment and children. The customer 18 year of age labeled as retired with and multiple kids. Some of these younger people are labeled as widowed or divorced as well. There is not a reasonable way to adjust these entries due to being self reported. K-Nearest Neighbors would also cause the dataset to impute incorrectly with these quality issues.

Other anomalies come from the data firm should have. A telecom provider giving a customer dataset should not have null values for its services per customer. A re-evaluation of data collection techniques should be used. It is hard to know if the imputed data is consistent with real world data.

D7 Discuss how the limitations in part D6 affect the analysis of the question or decision from part A.

These limitations point to using this dataset only for tentative testing. For example, Age had 25% of its data missing. With that missing data assuming the distribution was valid, bias can be introduced into the dataset that would not be reflective of the real customer ages. This could also impact machine learning training in ways causing bad predictions.

PCA

Approach

PCA will be applied to the continuous variables in the dataset. Before applying PCA, the data will be scaled via RobustScaler. Many of the items in this dataset were self-reported causing anomalies in data quality and mismatched categorization. For example, a customer in Yamhill Oregon is noted to be a CFO but Income is listed \$33000. Due to these outliers and the nature of this data being a self-reported survey providing treatment corrections to this data without another similar sample is not recommended. Instead, RobustScaler limits the effects of outliers by assuming the interquartile range of data correctly correspond to the data distribution (Hale, 2019).

The scaled data will consist of the following continuous variables: ['Outage_sec_perweek', 'Yearly_equip_failure', 'Age', 'Contacts', 'Tenure', 'MonthlyCharge', 'Bandwidth_GB_Year']

```
In [57]: # Set scaler
from sklearn.preprocessing import RobustScaler, StandardScaler

# instantiate scaler
scaler = RobustScaler()
```

```
In [58]: # Set scaling metrics
cont_cols = ['Outage_sec_perweek', 'Yearly_equip_failure', 'Age', 'Contacts', 'Tenure', 'MonthlyCharge', 'Bandw
```

```
In [59]: # a copy of the clean dataframe will be made to preserve dataframe integrity

df_scaled = df_clean.copy()
```

```
In [60]: # test copy
df_scaled.head()
```

```
[In] [60]: Out[60]:
   Customer_id  Interaction    City State County Zip Lat Lng Population Area Timezone
CaseOrder
1 K409198 ae905026- Point Baker Prince of Wales-Hyder 99927 56.251 -133.37571 38 Urban America/Sitka Alaska 1
2 S120509 bf76459f- West Branch MI Ogemaw 48661 44.32893 -84.2408 10446 Urban America/Detroit Michigan 1
3 K191035 344d114c- Yamhill OR Yamhill 97148 45.35589 -123.24657 3735 Urban America/Los_Angeles Oregon 1
4 D908050 adfa2b40- Del Mar CA San Diego 92014 32.96687 -117.24798 13863 Suburban America/Los_Angeles California 1
5 K662701 b6a8e16f- Needville TX Fort Bend 77461 29.38012 -95.80673 11352 Suburban America/Chicago Texas 1
```

```
In [61]: # preview continuous variables descriptive stats as a baseline to review after scaling
df_scaled.cont_vals
```

```
[In] [61]: Out[61]:
               Outage_sec_perweek  Yearly_equip_failure    Age  Contacts    Tenure  MonthlyCharge  Bandwidth_GB_Year
count  10000.000000          10000.000000  10000.000000  10000.000000  10000.000000  10000.000000  10000.000000
mean      11453592.8                0.398000      53.284800      0.994200      36.966114      174.076305      3173.812307
std       7.0248482              0.635953      20.747402      0.988466      26.329346      43.335473      2173.85637
min       0.000000                0.000000      18.000000      0.000000      1.000259      77.505230      155.506715
25%      8.054382                0.000000      35.000000      0.000000      8.700329      147.017078      1234.000000
50%     10.220896                0.000000      53.000000      1.000000      47.444800      169.915400      2046.732140
75%     12.478920                0.000000      71.000000      2.000000      61.990000      203.777441      5466.284500
max      47.049480              6.000000      89.000000      7.000000      91.999280      315.878600      7158.982000
```

```
In [62]: # apply scaler
df_scaled(cont_vals) = scaler.fit_transform(df_scaled(cont_vals))
```

```
In [63]: #Unit test scaling
df_scaled(cont_vals).head()
```

```
[In] [63]: Out[63]:
               Outage_sec_perweek  Yearly_equip_failure    Age  Contacts    Tenure  MonthlyCharge  Bandwidth_GB_Year
count  1.000000e+04          1.000000e+00  1.000000e+00  1.000000e+00  1.000000e+00  1.000000e+00  1.000000e+04
mean      2.82115e-01                0.398000      0.007894      -0.002900      -0.200359      0.066355      2.672505e-01
std       1.584569e+00              0.635953      0.576317      0.494233      0.503432      0.691086      5.136459e-01
min       -2.301432e+00              0.000000      -0.972222      -0.880407      -0.740817      -0.1473697      -4.468569e-01
25%      -4.846374e-01                0.000000      -0.500000      -0.500000      -0.470817      -0.459990      -1.920315e-01
50%     -2.036365e-16                 0.000000      0.000000      0.000000      0.000000      0.000000      2.688821e-17
75%     5.153656e-16                 0.000000      0.500000      0.000000      0.259183      0.540010      8.796856e-01
max      8.311317e+00                6.000000      1.000000      
```