

Task 4 - Modeling & Evaluation

Summary

1. Poor F1 and recall scores due to imbalance classes. Target of interest is churn = 1.
2. Created a model after balancing classes gives better results in these metrics.

Import Libraries

```
In [1]: import numpy as np
import pandas as pd
import matplotlib.pyplot as plt
import seaborn as sns
import sklearn
import xgboost
import sklearn.metrics
from sklearn.metrics import classification_report, DiscriminationThreshold
from sklearn.experimental import enable_hist_gradient_boosting
from sklearn.linear_model import LogisticRegression, LogisticRegressionCV, Ridge
from sklearn.ensemble import RandomForestClassifier, RandomForestRegressor, ExtraTreesClassifier, ExtraTreesRegressor
from sklearn.ensemble import GradientBoostingClassifier, GradientBoostingRegressor, HistGradientBoostingClassifier, HistGradientBoostingRegressor
import plotlylib
sns.set(style='dark')
sns.set(font_scale=1.2)
from sklearn.pipeline import Pipeline
from sklearn.model_selection import RepeatedStratifiedKFold
from sklearn.feature_selection import RFE, RFECV, SelectKBest, f_classif, f_regression, chi2
from sklearn.inspection import permutation_importance
from sklearn.model_selection import cross_val_score, train_test_split, GridSearchCV, RandomizedSearchCV
from sklearn.preprocessing import LabelEncoder, StandardScaler, MinMaxScaler, OneHotEncoder
from sklearn.pipeline import Pipeline
from sklearn.metrics import confusion_matrix, classification_report, mean_absolute_error, mean_squared_error
from sklearn.metrics import plot_confusion_matrix, plot_precision_recall_curve, plot_roc_curve, accuracy_score
from sklearn.metrics import auc, f1_score, precision_score, recall_score, roc_auc_score
from imblearn.under_sampling import RandomUnderSampler
from imblearn.over_sampling import RandomOverSampler
import warnings
warnings.filterwarnings('ignore')
import pickle
from pickle import dump, load
np.random.seed(0)
from pycoret.classification import *
# from pycoret.clustering import *
# from pycoret.regression import *
pd.set_option('display.max_columns',100)
pd.set_option('display.max_rows',100)
pd.set_option('display.width', 1000)
pd.set_printoptions(suppress=True)
```

Data Exploration and Analysis

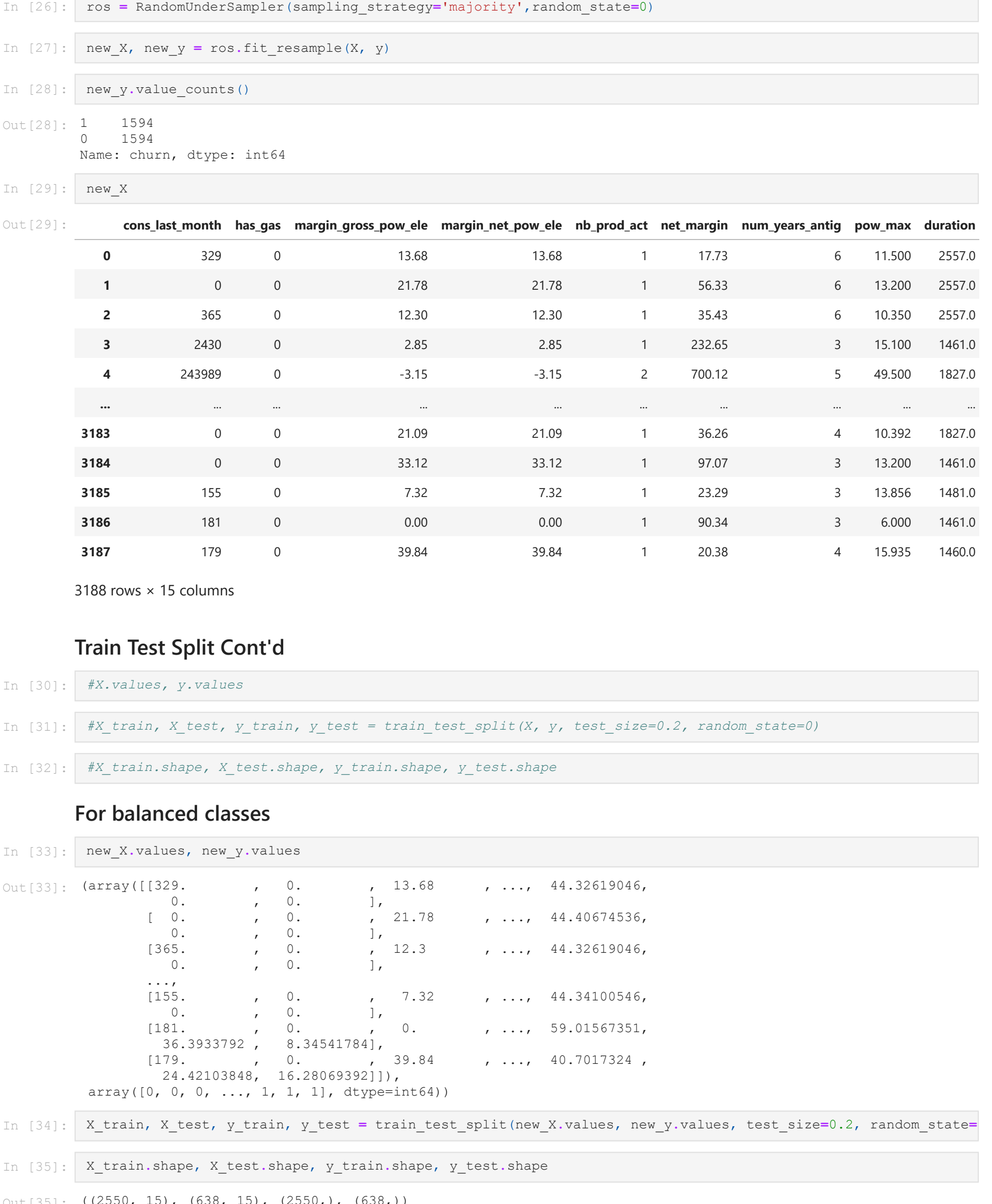
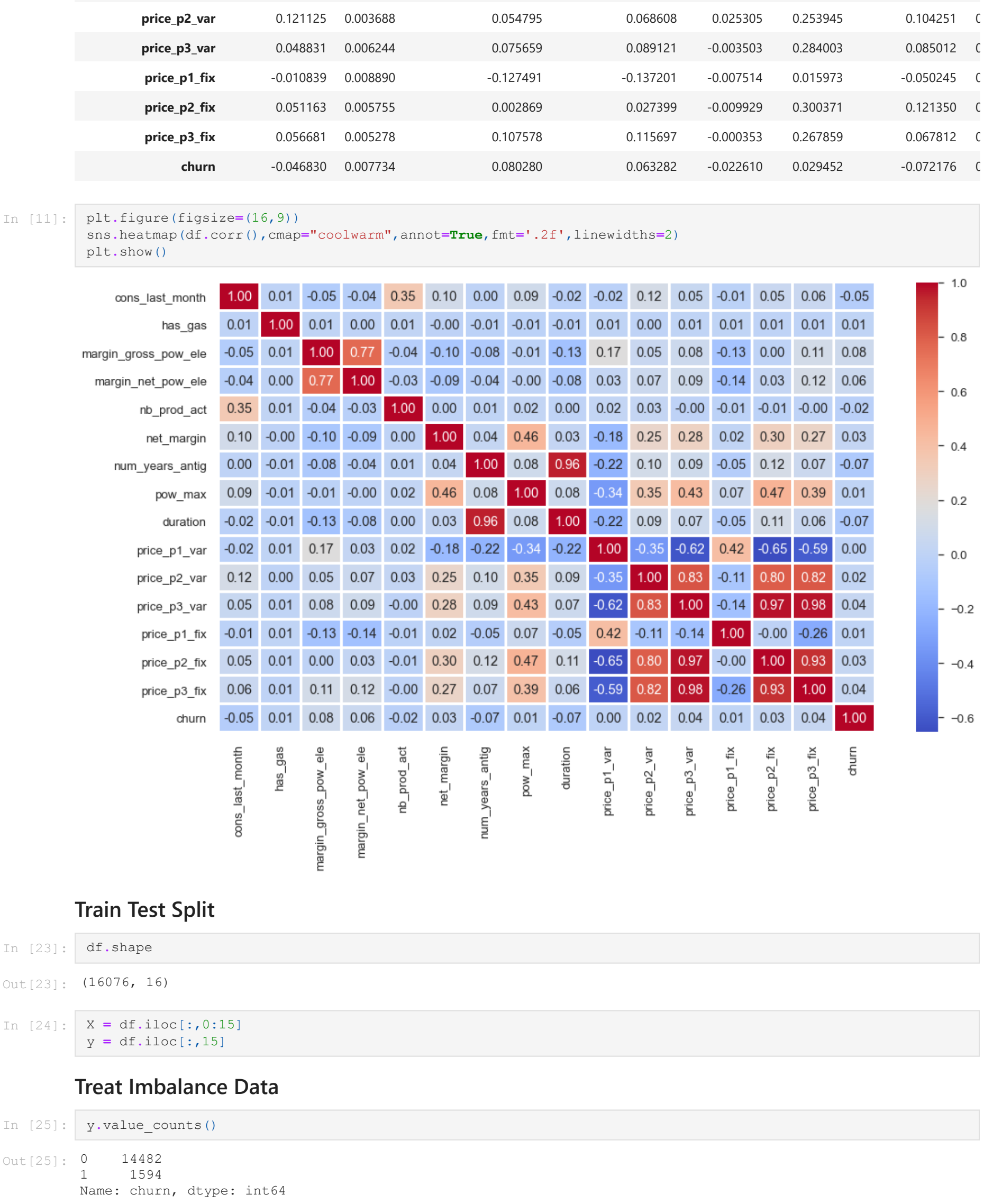
```
In [2]: df = pd.read_csv('final2.csv')
In [3]: df
Out[1]:
```

	cons_last_month	has_gas	margin_gross_pow_ele	margin_net_pow_ele	nb_prod_act	net_margin	num_years_antig	pow_max	duration
0	10025	0	21.76	-0.37766	1	1732.36	3	180.000	14601
1	0	1	45.44	25.44	2	678.99	6	13.800	10961
2	0	0	16.38	16.38	1	18.89	6	13.800	205486
3	0	0	28.60	28.60	1	6.60	6	13.856	21921
4	0	0	30.22	30.22	1	25.46	6	13.200	21921
...
16071	0	0	27.88	27.88	2	381.77	4	15.000	14451
16072	181	0	0.00	0.00	1	90.34	3	6.000	14611
16073	179	0	39.84	39.84	1	20.38	4	15.935	14601
16074	0	0	13.08	13.08	1	96.34	4	11.000	14611
16075	0	1	11.84	11.84	1	0.96	6	10.392	25561

16076 rows x 16 columns

```
In [4]: df.info()
<class 'pandas.core.frame.DataFrame'>
Int64Index: 16076 entries, 0 to 16075
Data columns (total 16 columns):
#   Column                Non-Null Count  Dtype
---  -
0   cons_last_month        16076 non-null   int64
1   has_gas                16076 non-null   int64
2   margin_gross_pow_ele   16076 non-null   float64
3   margin_net_pow_ele     16076 non-null   float64
4   nb_prod_act            16076 non-null   int64
5   net_margin             16076 non-null   float64
6   num_years_antig        16076 non-null   int64
7   pow_max                16076 non-null   float64
8   duration              16076 non-null   float64
9   price_p1_var           16076 non-null   float64
10  price_p2_var           16076 non-null   float64
11  price_p3_var           16076 non-null   float64
12  price_p1_fix           16076 non-null   float64
13  price_p2_fix           16076 non-null   float64
14  price_p3_fix           16076 non-null   float64
15  churn                  16076 non-null   int64
dtypes: float64(11), int64(5)
memory usage: 2.0 MB
In [5]: df.describe(include='all')
```

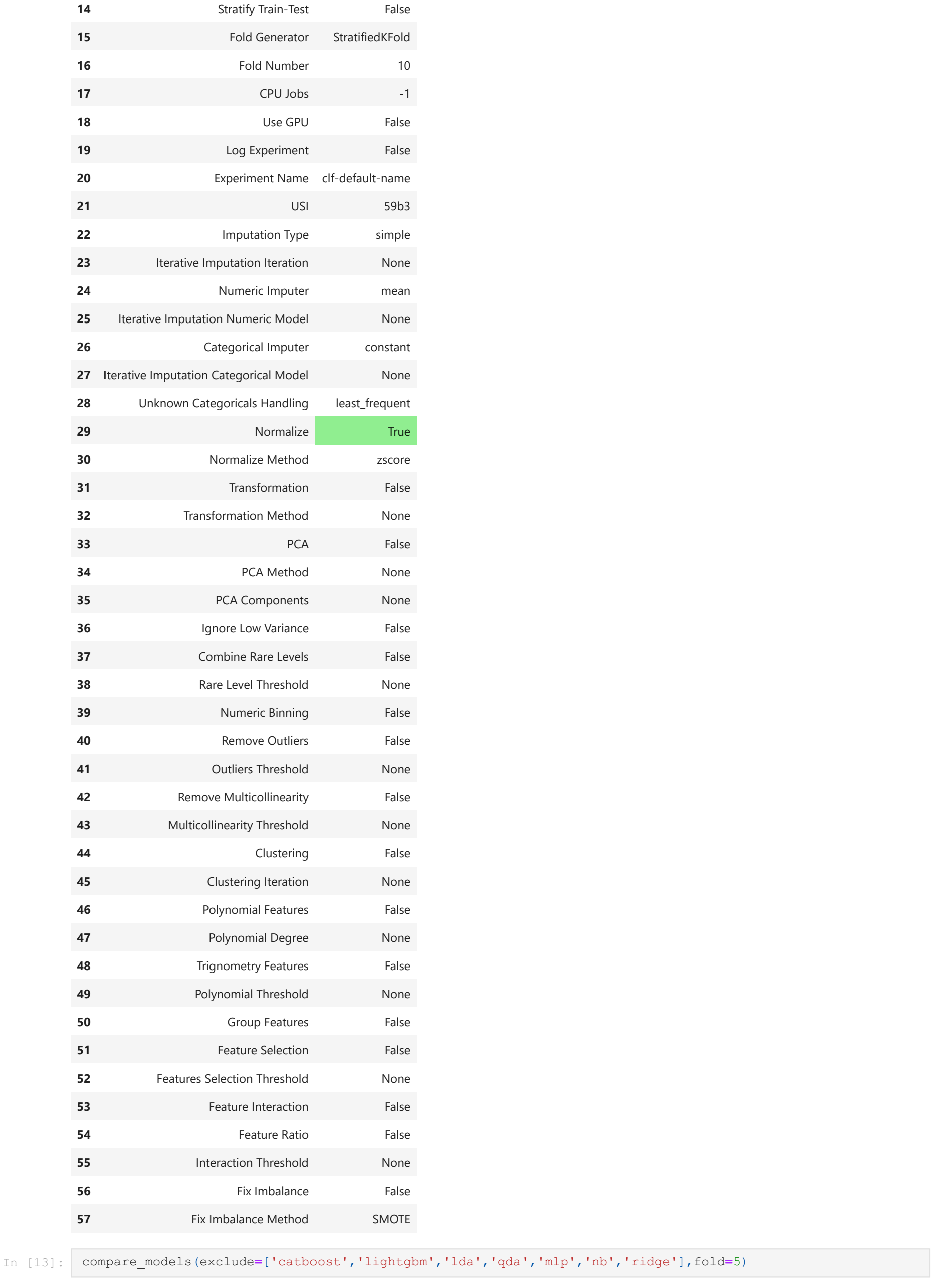
	cons_last_month	has_gas	margin_gross_pow_ele	margin_net_pow_ele	nb_prod_act	net_margin	num_years_antig	pow_max
count	1.607600e+04	16076.000000	16076.000000	16076.000000	16076.000000	16076.000000	16076.000000	16076.000000
mean	1.94347e+04	0.141888	22.470364	21.497969	1.347972	217.938200	5.029547	20.5486
std	8.233817e+04	0.376749	52.554800	27.920127	1.460638	366.673412	1.675753	21.48370
min	-1.138600e+04	0.000000	-23.754023	-61.660000	1.000000	-4148.990000	1.000000	1.000000
25%	0.000000e+00	0.000000	12.050000	11.950000	1.000000	51.985000	4.000000	12.50000
50%	0.000000e+02	0.000000	22.090000	20.970000	1.000000	119.685000	5.000000	13.85000
75%	4.124000e+03	0.000000	29.640000	29.640000	1.000000	275.772500	6.000000	19.80000
max	4.538720e+06	1.000000	374.640000	374.640000	32.000000	24570.650000	16.000000	500.00000



Correlation

```
In [10]: df.corr()
Out[10]:
```

	cons_last_month	has_gas	margin_gross_pow_ele	margin_net_pow_ele	nb_prod_act	net_margin	num_years_antig	pow_max	duration
cons_last_month	1.000000	0.005240	-0.054205	-0.037766	0.351237	0.096491	0.004668	0.007612	-0.001287
has_gas	0.005240	1.000000	0.009682	0.004285	0.014025	-0.002258	-0.008432	-0.001787	-0.000425
margin_gross_pow_ele	-0.054205	0.009682	1.000000	0.766458	-0.043861	-0.098432	-0.081287	-0.007702	-0.001774
margin_net_pow_ele	-0.037766	0.004285	0.766458	1.000000	-0.032268	-0.086216	-0.035457	-0.037702	-0.001774
nb_prod_act	0.351237	0.014025	-0.043861	-0.032268	1.000000	0.004580	0.009473	0.004735	0.000943
net_margin	0.096491	-0.002258	-0.098432	-0.086216	0.004580	1.000000	0.035457	0.004580	0.000943
num_years_antig	0.004668	-0.001287	-0.081287	-0.035457	0.009473	0.035457	1.000000	0.004668	0.000943
pow_max	0.007612	-0.000425	-0.001787	-0.001774	0.000943	0.004580	0.004668	1.000000	0.000943
duration	-0.001287	-0.000425	-0.001787	-0.001774	0.000943	0.004580	0.004668	0.000943	1.000000
price_p1_var	-0.016011	-0.013632	-0.132225	-0.076409	0.002700	0.026245	0.998616	0.998616	0.998616
price_p2_var	-0.017634	0.005187	0.174420	0.031320	0.023888	-0.177704	-0.219561	-0.219561	-0.219561
price_p3_var	0.121125	0.003688	0.057495	0.069608	0.025305	0.253945	0.104251	0.104251	0.104251
price_p1_fix	0.048311	0.006244	0.075659	0.089121	-0.003503	0.284033	0.085012	0.085012	0.085012
price_p2_fix	-0.010839	0.008890	-0.127491	-0.137201	-0.007514	0.015973	-0.050245	-0.050245	-0.050245
price_p3_fix	0.051163	0.005758	0.002869	0.027399	-0.000923	0.300371	0.121350	0.121350	0.121350
price_p1_fix	0.056681	0.005275	0.107578	0.115697	-0.000553	0.267859	0.067812	0.067812	0.067812
churn	-0.046830	0.007734	0.080280	0.063282	-0.022610	0.029452	-0.072176	-0.072176	-0.072176



Train Test Split

```
In [23]: df.shape
Out[23]: (16076, 16)
In [24]: X = df.iloc[:,0:15]
y = df.iloc[:,15]
```

Treat Imbalance Data

```
In [25]: y.value_counts()
Out[25]:
0    14482
1     1594
Name: churn, dtype: int64
In [26]: ros = RandomUnderSampler(sampling_strategy='majority', random_state=0)
In [27]: new_X, new_y = ros.fit_resample(X, y)
In [28]: new_X.value_counts()
Out[28]:
0    1594
1     1594
Name: churn, dtype: int64
In [29]: new_X
```

	cons_last_month	has_gas	margin_gross_pow_ele	margin_net_pow_ele	nb_prod_act	net_margin	num_years_antig	pow_max	duration
0	329	0	13.68	13.68	1	57.33	6	11.500	25570
1	0	0	21.78	21.78	1	16.73	6	13.200	25570
2	365	0	12.30	12.30	1	35.43	6	10.350	25570
3	2430	0	28.85	28.85	1	232.65	3	15.100	14610
4	24369	0	-3.15	-3.15	2	700.12	5	49.500	18270
...
3183	0	0	21.09	21.09	1	36.26	4	10.392	18270
3184	0	0	73.32	73.32	1	97.07	3	13.200	14610
3185	195	0	3.72	3.72	1	23.29	3	13.856	14610
3186	981	0	70.00	70.00	1	90.34	3	6.000	14610
3187	179	0	39.84	39.84	1	20.38	4	15.935	14600

3188 rows x 15 columns

Train Test Split Cont'd

```
In [30]: #X.values, y.values
In [31]: #X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.2, random_state=0)
In [32]: #X_train.shape, X_test.shape, y_train.shape, y_test.shape
```

For balanced classes

```
In [33]: new_X.values, new_y.values
Out[33]: (array([[329, 0, 13.68, ..., 44.32619046, 0, 1, 57.33, 6, 11.5, 25570],
[0, 0, 21.78, ..., 44.40674536, 0, 1, 16.73, 6, 13.2, 14610],
[365, 0, 12.3, ..., 44.32619046, 0, 1, 35.43, 6, 10.35, 25570],
[2430, 0, 28.85, ..., 44.32619046, 0, 1, 232.65, 3, 15.1, 14610],
[24369, 0, -3.15, ..., 44.32619046, 0, 2, 700.12, 5, 49.5, 18270],
[21.09, 0, 36.26, ..., 44.32619046, 0, 1, 97.07, 3, 13.2, 14610],
[73.32, 0, 97.07, ..., 44.32619046, 0, 1, 23.29, 3, 13.856, 14610],
[981, 0, 70.0, ..., 44.32619046, 0, 1, 90.34, 3, 6.0, 14610],
[39.84, 0, 20.38, ..., 44.32619046, 0, 1, 15.935, 4, 15.935, 14600],
[21.09, 0, 36.26, ..., 44.32619046, 0, 1, 97.07, 3, 13.2, 14610],
[73.32, 0, 97.07, ..., 44.32619046, 0, 1, 23.29, 3, 13.856, 14610],
[981, 0, 70.0, ..., 44.32619046, 0, 1, 90.34, 3, 6.0, 14610],
[39.84, 0, 20.38, ..., 44.32619046, 0, 1, 15.935, 4, 15.935, 14600],
[21.09, 0, 36.26, ..., 44.32619046, 0, 1, 97.07, 3, 13.2, 14610],
[73.32, 0, 97.07, ..., 44.32619046, 0, 1, 23.29, 3, 13.856, 14610],
[981, 0, 70.0, ..., 44.32619046, 0, 1, 90.34, 3, 6.0, 14610],
[39.84, 0, 20.38, ..., 44.32619046, 0, 1, 15.935, 4, 15.935, 14600],
[21.09, 0, 36.26, ..., 44.32619046, 0, 1, 97.07, 3, 13.2, 14610],
[73.32, 0, 97.07, ..., 44.32619046, 0, 1, 23.29, 3, 13.856, 14610],
[981, 0, 70.0, ..., 44.32619046, 0, 1, 90.34, 3, 6.0, 14610],
[39.84, 0, 20.38, ..., 44.32619046, 0, 1, 15.935, 4, 15.935, 14600],
[21.09, 0, 36.26, ..., 44.32619046, 0, 1, 97.07, 3, 13.2, 14610],
[73.32, 0, 97.07, ..., 44.32619046, 0, 1, 23.29, 3, 13.856, 14610],
[981, 0, 70.0, ..., 44.32619046, 0, 1, 90.34, 3, 6.0, 14610],
[39.84, 0, 20.38, ..., 44.32619046, 0, 1, 15.935, 4, 15.935, 14600],
[21.09, 0, 36.26, ..., 44.32619046, 0, 1, 97.07, 3, 13.2, 14610],
[73.32, 0, 97.07, ..., 44.32619046, 0, 1, 23.29, 3, 13.856, 14610],
[981, 0, 70.0, ..., 44.32619046, 0, 1, 90.34, 3, 6.0, 14610],
[39.84, 0, 20.38, ..., 44.32619046, 0, 1, 15.935, 4, 15.935, 14600],
[21.09, 0, 36.26, ..., 44.32619046, 0, 1, 97.07, 3, 13.2, 14610],
[73.32, 0, 97.07, ..., 44.32619046, 0, 1, 23.29, 3, 13.856, 14610],
[981, 0, 70.0, ..., 44.32619046, 0, 1, 90.34, 3, 6.0, 14610],
[39.84, 0, 20.38, ..., 44.32619046, 0, 1, 15.935, 4, 15.935, 14600],
[21.09, 0, 36.26, ..., 44.32619046, 0, 1, 97.07, 3, 13.2, 14610],
[73.32, 0, 97.07, ..., 44.32619046, 0, 1, 23.29, 3, 13.856, 14610],
[981, 0, 70.0, ..., 44.32619046, 0, 1, 90.34, 3, 6.0, 14610],
[39.84, 0, 20.38, ..., 44.32619046, 0, 1, 15.935, 4, 15.935, 14600],
[21.09, 0, 36.26, ..., 44.32619046, 0, 1, 97.07, 3, 13.2, 14610],
[73.32, 0, 97.07, ..., 44.32619046, 0, 1, 23.29, 3, 13.856, 14610],
[981, 0, 70.0, ..., 44.32619046, 0, 1, 90.34, 3, 6.0, 14610],
[39.84, 0, 20.38, ..., 44.32619046, 0, 1, 15.935, 4, 15.935, 14600],
[21.09, 0, 36.26, ..., 44.32619046, 0, 1, 97.07, 3, 13.2, 14610],
[73.32, 0, 97.07, ..., 44.32619046, 0, 1, 23.29, 3, 13.856, 14610],
[981, 0, 70.0, ..., 44.32619046, 0, 1, 90.34, 3, 6.0, 14610],
[39.84, 0, 20.38, ..., 44.32619046, 0, 1, 15.935, 4, 15.935, 14600],
[21.09, 0, 36.26, ..., 44.32619046, 0, 1, 97.07, 3, 13.2, 14610],
[73.32, 0, 97.07, ..., 44.32619046, 0, 1, 23.29, 3, 13.856, 14610],
[981, 0, 70.0, ..., 44.32619046, 0, 1, 90.34, 3, 6.0, 14610],
[39.84, 0, 20.38, ..., 44.32619046, 0, 1, 15.935, 4, 15.935, 14600],
[21.09, 0, 36.26, ..., 44.32619046, 0, 1, 97.07, 3, 13.2, 14610],
[73.32, 0, 97.07, ..., 44.32619046, 0, 1, 23.29, 3, 13.856, 14610],
[981, 0, 70.0, ..., 44.32619046, 0, 1, 90.34, 3, 6.0, 14610],
[39.84, 0, 20.38, ..., 44.32619046, 0, 1, 15.935, 4, 15.935, 14600],
[21.09, 0, 36.26, ..., 44.32619046, 0, 1, 97.07, 3, 13.2, 14610],
[73.32, 0, 97.07, ..., 44.32619046, 0, 1, 23.29, 3, 13.856, 14610],
[981, 0, 70.0, ..., 44.32619046, 0, 1, 90.34, 3, 6.0, 14610],
[39.84, 0, 20.38, ..., 44.32619046, 0, 1, 15.935, 4, 15.935, 14600],
[21.09, 0, 36.26, ..., 44.32619046, 0, 1, 97.07, 3, 13.2, 14610],
[73.32, 0, 97.07, ..., 44.32619046, 0, 1, 23.29, 3, 13.856, 14610],
[981, 0, 70.0, ..., 44.32619046, 0, 1, 90.34, 3, 6.0, 14610],
[39.84, 0, 20.38, ..., 44.32619046, 0, 1, 15.935, 4, 15.935, 14600],
[21.09, 0, 36.26, ..., 44.32619046, 0, 1, 97.07, 3, 13.2, 14610],
[73.32, 0, 97.07, ..., 44.32619046, 0, 1, 23.29, 3, 13.856, 14610],
[981, 0, 70.0, ..., 44.32619046, 0, 1, 90.34, 3, 6.0, 14610],
[39.84, 0, 20.38, ..., 44.32619046, 0, 1, 15.935, 4, 15.935, 14600],
[21.09, 0, 36.26, ..., 44.32619046, 0, 1, 97.07, 3, 13.2, 14610],
[73.32, 0, 97.07, ..., 44.32619046, 0, 1, 23.29, 3, 13.856, 14610],
[981, 0, 70.0, ..., 44.32619046, 0, 1, 90.34, 3, 6.0, 14610],
[39.84, 0, 20.38, ..., 44.32619046, 0, 1, 15.935, 4, 15.935, 14600],
[21.09, 0, 36.26, ..., 44.32619046, 0, 1, 97.07, 3, 13.2, 14610],
[73.32, 0, 97.07, ..., 44.32619046, 0, 1, 23.29, 3, 13.856, 14610],
[981, 0, 70.0, ..., 44.32619046, 0, 1, 90.34, 3, 6.0, 14610],
[39.8
```