

```
In [1]: #Load data from history
import os
os.chdir('C:/Users/TOTGAOUSER4/Documents/Totago Technologies/Data Science/Projects/BCG')
```

Datasets

```
In [2]: import pandas as pd
cust_data = pd.read_csv('ml_case_training_data.csv')
price_data = pd.read_csv('ml_case_training_hist_data.csv')
churn_data = pd.read_csv('ml_case_training_output.csv')

In [3]: #Customer Dataset
cust_data.head(10)

Out [3]:
```

	id	activity_new	campaign_disc_ele	channel_sales	cons_12m
0	48ada52261e7cf58715202705a0451c9	esouifdkbckeluumfuaadccomkmixw	NaN	lnkebamcaacubkfadiuuecoxiemlma	309275
1	24011ae4ebbc303511f056fa7c15bc57	NaN	NaN	foosdfpfusacimwkcsoibcdckiciaua	0
2	d29c2c54acc38f3c0614da0a653813dd	NaN	NaN	NaN	4680
3	764cf596b1154dac3a6c254a082ea7d	NaN	NaN	foosdfpfusacimwkcsoibcdckiciaua	544
4	bbab0339a292a1ef668029a4c16191cb	NaN	NaN	lnkebamcaacubkfadiuuecoxiemlma	1584
5	568bb38a1af7cd0c49c77b3789b59a3	sfafmfcdpdmckuekkuekuuexkidaouu	NaN	foosdfpfusacimwkcsoibcdckiciaua	121335
6	149d57f9262c41c1f94415803a877cb4b	NaN	NaN	NaN	4425
7	1aa49825382410b098937d5c54ec2cd	NaN	NaN	uslukupasemubtlpokaafesimbrdff	8302
8	7ab4bf4878d8f7661dc20e9b8e18011	scsfopkxpfckuekkuekuuexkmwauub	NaN	foosdfpfusacimwkcsoibcdckiciaua	45097
9	01495c950b7ec5e7f3203406785aae0	NaN	NaN	foosdfpfusacimwkcsoibcdckiciaua	29552

10 rows x 32 columns

```
In [4]: #Pricing Data
price_data.head(10)

Out [4]:
```

	id	price_date	price_p1_var	price_p2_var	price_p3_var	price_p1_fix	price_p2_fix	price_p3_fix
0	038faf19179925da21a25619c5a24b745	2015-01-01	0.151367	0.0	0.0	44.266931	0.0	0.0
1	038faf19179925da21a25619c5a24b745	2015-02-01	0.151367	0.0	0.0	44.266931	0.0	0.0
2	038faf19179925da21a25619c5a24b745	2015-03-01	0.151367	0.0	0.0	44.266931	0.0	0.0
3	038faf19179925da21a25619c5a24b745	2015-04-01	0.149626	0.0	0.0	44.266931	0.0	0.0
4	038faf19179925da21a25619c5a24b745	2015-05-01	0.149626	0.0	0.0	44.266931	0.0	0.0
5	038faf19179925da21a25619c5a24b745	2015-06-01	0.149626	0.0	0.0	44.266930	0.0	0.0
6	038faf19179925da21a25619c5a24b745	2015-07-01	0.150321	0.0	0.0	44.444710	0.0	0.0
7	038faf19179925da21a25619c5a24b745	2015-08-01	0.145859	0.0	0.0	44.444710	0.0	0.0
8	038faf19179925da21a25619c5a24b745	2015-09-01	0.145859	0.0	0.0	44.444710	0.0	0.0
9	038faf19179925da21a25619c5a24b745	2015-10-01	0.145859	0.0	0.0	44.444710	0.0	0.0

In [5]: #Churn data
churn_data.head(10)

Out [5]:

	id	churn
0	48ada52261e7cf58715202705a0451c9	1
1	24011ae4ebbc303511f056fa7c15bc57	0
2	d29c2c54acc38f3c0614da0a653813dd	0
3	764cf596b1154dac3a6c254a082ea7d	0
4	bbab0339a292a1ef668029a4c16191cb	0
5	568bb38a1af7cd0c49c77b3789b59a3	0
6	149d57f9262c41c1f94415803a877cb4b	0
7	1aa49825382410b098937d5c54ec2cd	1
8	7ab4bf4878d8f7661dc20e9b8e18011	1
9	01495c950b7ec5e7f3203406785aae0	0

Data Cleaning

```
Dealing with missing data

In [6]: #Checking for null values in Customer Details Dataset
cust_data.isnull().sum()

Out [6]:
```

id	0
activity_new	9545
campaign_disc_ele	16096
channel_sales	4218
cons_12m	0
cons_gas_12m	0
cons_last_month	0
date_active	0
date_end	2
date_first_activ	1258
date_modif_prod	157
date_renewal	40
forecast_base_bill_ele	12588
forecast_base_bill_year	12588
forecast_bill_12m	12588
forecast_cons	12588
forecast_cons_12m	0
forecast_cons_year	0
forecast_discount_energy	126
forecast_meter_rent_12m	0
forecast_price_energy_p1	126
forecast_price_energy_p2	126
forecast_price_pow_p1	126
has_gas	0
imp_cons	0
margin_gross_pow_ele	13
margin_net_pow_ele	13
nb_prod_act	0
net_margin	15
num_years_antig	0
origin_up	876
pow_max	0
dtype:	int64

Columns with null values - activity_new, campaign_disc_ele, channel_sales, date_end, date_first_activ, date_modif_prod, date_renewal, forecast_base_bill_ele, forecast_base_bill_year, forecast_bill_12m, forecast_cons, forecast_discount_energy, forecast_price_energy_p1, forecast_price_p1_var, forecast_price_p2_var, margin_gross_pow_ele, margin_net_pow_ele, net_margin, origin_up, pow_max

```
In [7]: #Checking for null values in Price Dataset
cust_data.isnull().sum()

Out [7]:
```

id	0
price_p1_var	1359
price_p2_var	1359
price_p3_var	1359
price_p1_fix	1359
price_p2_fix	1359
price_p3_fix	1359
dtype:	int64

Columns with missing data - price_p1_var, price_p2_var, price_p3_var, price_p1_fix, price_p2_fix, price_p3_fix

```
In [8]: #Checking for null values in Churn Dataset
churn_data.isnull().sum()

Out [8]:
```

id	0
churn	0
dtype:	int64

No nulls here

```
In [9]: #Converting the dates to datetime
import datetime
cust_data['date_active', 'date_end', 'date_first_activ', 'date_modif_prod', 'date_renewal'] = \
cust_data[a] = pd.to_datetime(cust_data[a], format='%Y %m %d')
cust_data.head(10)

Out [9]:
```

	id	activity_new	campaign_disc_ele	channel_sales	cons_12m
0	48ada52261e7cf58715202705a0451c9	esouifdkbckeluumfuaadccomkmixw	NaN	lnkebamcaacubkfadiuuecoxiemlma	309275
1	24011ae4ebbc303511f056fa7c15bc57	NaN	NaN	foosdfpfusacimwkcsoibcdckiciaua	0
2	d29c2c54acc38f3c0614da0a653813dd	NaN	NaN	NaN	4680
3	764cf596b1154dac3a6c254a082ea7d	NaN	NaN	foosdfpfusacimwkcsoibcdckiciaua	544
4	bbab0339a292a1ef668029a4c16191cb	NaN	NaN	lnkebamcaacubkfadiuuecoxiemlma	1584
5	568bb38a1af7cd0c49c77b3789b59a3	sfafmfcdpdmckuekkuekuuexkidaouu	NaN	foosdfpfusacimwkcsoibcdckiciaua	121335
6	149d57f9262c41c1f94415803a877cb4b	NaN	NaN	NaN	4425
7	1aa49825382410b098937d5c54ec2cd	NaN	NaN	uslukupasemubtlpokaafesimbrdff	8302
8	7ab4bf4878d8f7661dc20e9b8e18011	scsfopkxpfckuekkuekuuexkmwauub	NaN	foosdfpfusacimwkcsoibcdckiciaua	45097
9	01495c950b7ec5e7f3203406785aae0	NaN	NaN	foosdfpfusacimwkcsoibcdckiciaua	29552

10 rows x 32 columns

```
In [10]: #Checking the data type
cust_data.info()

<class 'pandas.core.frame.DataFrame'>
RangeIndex: 16096 entries, 0 to 16095
Data columns (total 32 columns):
# Column Non-Null Count Dtype
---
0 id 16096 non-null object
1 activity_new 6551 non-null object
2 campaign_disc_ele 0 non-null float64
3 channel_sales 11878 non-null object
4 cons_12m 16096 non-null int64
5 cons_gas_12m 16096 non-null int64
6 cons_last_month 16096 non-null int64
7 date_active 15935 non-null datetime64[ns]
8 date_end 16094 non-null datetime64[ns]
9 date_first_activ 3508 non-null datetime64[ns]
10 date_modif_prod 16096 non-null datetime64[ns]
11 date_renewal 16056 non-null datetime64[ns]
12 forecast_base_bill_ele 3508 non-null float64
13 forecast_base_bill_year 3508 non-null float64
14 forecast_bill_12m 3508 non-null float64
15 forecast_cons 3508 non-null float64
16 forecast_cons_12m 16096 non-null float64
17 forecast_cons_year 16096 non-null int64
18 forecast_discount_energy 15970 non-null float64
19 forecast_meter_rent_12m 16096 non-null float64
20 forecast_price_energy_p1 15970 non-null float64
21 forecast_price_energy_p2 15970 non-null float64
22 forecast_price_pow_p1 15970 non-null float64
23 has_gas 16096 non-null object
24 imp_cons 16096 non-null object
25 margin_gross_pow_ele 16083 non-null float64
26 margin_net_pow_ele 16083 non-null float64
27 nb_prod_act 16096 non-null int64
28 num_years_antig 16083 non-null float64
29 num_years_antig 16096 non-null int64
30 origin_up 16009 non-null object
31 pow_max 16093 non-null object
dtypes: datetime64[ns](5), float64(13), int64(6), object(3)
memory usage: 3.9+ MB
```

```
In [11]: # Filling up nulls in Customer Details Dataset

# Filling with zeros
for col in ['forecast_base_bill_ele', 'forecast_base_bill_year', 'forecast_bill_12m', 'forecast_cons',
'forecast_discount_energy', 'forecast_price_energy_p1', 'forecast_price_energy_p2', 'forecast_price_pow_p1',
'margin_gross_pow_ele', 'margin_net_pow_ele', 'net_margin', 'pow_max']:
    cust_data[col].fillna(0, inplace = True)

# Filling Dates with activation date
for col in ['date_end', 'date_first_activ', 'date_modif_prod', 'date_renewal']:
    cust_data[col].fillna(cust_data['date_active'], inplace = True)

# Dropping Some columns
cust_data.drop(columns=['activity_new', 'campaign_disc_ele', ], inplace=True)
cust_data.head(10)

Out [11]:
```

	id	channel_sales	cons_12m	cons_gas_12m	cons_last_month	date_active	date_end
0	48ada52261e7cf58715202705a0451c9	lnkebamcaacubkfadiuuecoxiemlma	309275	0	10025	2012-11-07	2016-10-16
1	24011ae4ebbc303511f056fa7c15bc57	foosdfpfusacimwkcsoibcdckiciaua	0	54946	0	2013-06-15	2016-10-16
2	d29c2c54acc38f3c0614da0a653813dd	NaN	4680	0	0	2009-08-21	2016-10-16
3	764cf596b1154dac3a6c254a082ea7d	foosdfpfusacimwkcsoibcdckiciaua	544	0	0	2010-04-16	2016-10-16
4	bbab0339a292a1ef668029a4c16191cb	lnkebamcaacubkfadiuuecoxiemlma	1584	0	0	2010-03-30	2016-10-16
5	568bb38a1af7cd0c49c77b3789b59a3	foosdfpfusacimwkcsoibcdckiciaua	121335	0	12400	2010-04-08	2016-10-16
6	149d57f9262c41c1f94415803a877cb4b	NaN	4425	0	526	2010-01-13	2016-10-16
7	1aa49825382410b098937d5c54ec2cd	uslukupasemubtlpokaafesimbrdff	8302	0	1998	2011-12-09	2016-10-16
8	7ab4bf4878d8f7661dc20e9b8e18011	foosdfpfusacimwkcsoibcdckiciaua	45097	0	0	2011-12-02	2016-10-16
9	01495c950b7ec5e7f3203406785aae0	foosdfpfusacimwkcsoibcdckiciaua	29552	0	1260	2010-04-21	2016-10-16

10 rows x 30 columns

```
In [12]: # Filling up nulls in Price Dataset

for col in ['price_p1_var', 'price_p2_var', 'price_p3_var', 'price_p1_fix', 'price_p2_fix', 'price_p3_fix']:
    price_data[col].fillna(0, inplace = True)
price_data.head(10)

Out [12]:
```

	id	price_date	price_p1_var	price_p2_var	price_p3_var	price_p1_fix	price_p2_fix	price_p3_fix
0	038faf19179925da21a25619c5a24b745	2015-01-01	0.151367	0.0	0.0	44.266931	0.0	0.0
1	038faf19179925da21a25619c5a24b745	2015-02-01	0.151367	0.0	0.0	44.266931	0.0	0.0
2	038faf19179925da21a25619c5a24b745	2015-03-01	0.151367	0.0	0.0	44.266931	0.0	0.0
3	038faf19179925da21a25619c5a24b745	2015-04-01	0.149626	0.0	0.0	44.266931	0.0	0.0
4	038faf19179925da21a25619c5a24b745	2015-05-01	0.149626	0.0	0.0	44.266931	0.0	0.0
5	038faf19179925da21a25619c5a24b745	2015-06-01	0.149626	0.0	0.0	44.266930	0.0	0.0
6	038faf19179925da21a25619c5a24b745	2015-07-01	0.150321	0.0	0.0	44.444710	0.0	0.0
7	038faf19179925da21a25619c5a24b745	2015-08-01	0.145859	0.0	0.0	44.444710	0.0	0.0
8	038faf19179925da21a25619c5a24b745	2015-09-01	0.145859	0.0	0.0	44.444710	0.0	0.0
9	038faf19179925da21a25619c5a24b745	2015-10-01	0.145859	0.0	0.0	44.444710	0.0	0.0

```
In [13]: #Cross checking nulls and checking variable types for Customer Details Dataset

<class 'pandas.core.frame.DataFrame'>
RangeIndex: 16096 entries, 0 to 16095
Data columns (total 30 columns):
# Column Non-Null Count Dtype
---
0 id 16096 non-null object
1 channel_sales 11878 non-null object
2 cons_12m 16096 non-null int64
3 cons_gas_12m 16096 non-null int64
4 cons_last_month 16096 non-null int64
5 date_active 16096 non-null datetime64[ns]
6 date_end 16094 non-null datetime64[ns]
7 date_first_activ 3508 non-null datetime64[ns]
8 date_modif_prod 16096 non-null datetime64[ns]
9 date_renewal 16056 non-null datetime64[ns]
10 forecast_base_bill_ele 3508 non-null float64
11 forecast_base_bill_year 3508 non-null float64
12 forecast_bill_12m 3508 non-null float64
13 forecast_cons 3508 non-null float64
14 forecast_cons_12m 16096 non-null float64
15 forecast_cons_year 16096 non-null int64
16 forecast_discount_energy 15970 non-null float64
17 forecast_meter_rent_12m 16096 non-null float64
18 forecast_price_energy_p1 15970 non-null float64
19 forecast_price_energy_p2 15970 non-null float64
20 forecast_price_pow_p1 15970 non-null float64
21 has_gas 16096 non-null object
22 imp_cons 16096 non-null object
23 margin_gross_pow_ele 16083 non-null float64
24 margin_net_pow_ele 16083 non-null float64
25 nb_prod_act 16096 non-null int64
26 net_margin 16096 non-null float64
27 num_years_antig 16096 non-null int64
28 origin_up 16009 non-null object
29 pow_max 16093 non-null object
dtypes: datetime64[ns](5), float64(13), int64(6), object(4)
memory usage: 3.7+ MB
```

```
In [14]: #Cross checking nulls and checking variable types for Price Dataset
price_data.info()

<class 'pandas.core.frame.DataFrame'>
RangeIndex: 193002 entries, 0 to 193001
Data columns (total 8 columns):
# Column Non-Null Count Dtype
---
0 id 193002 non-null object
1 price_date 193002 non-null object
2 price_p1_var 193002 non-null float64
3 price_p2_var 193002 non-null float64
4 price_p3_var 193002 non-null float64
5 price_p1_fix 193002 non-null float64
6 price_p2_fix 193002 non-null float64
7 price_p3_fix 193002 non-null float64
dtypes: float64(6), object(2)
memory usage: 11.8+ MB
```

Data Analysis

```
In [15]: #For better visualization import matplotlib and seaborn
import matplotlib.pyplot as plt
plt.style.use('classic')
import matplotlib
import numpy as np
import seaborn as sns
sns.set()
```

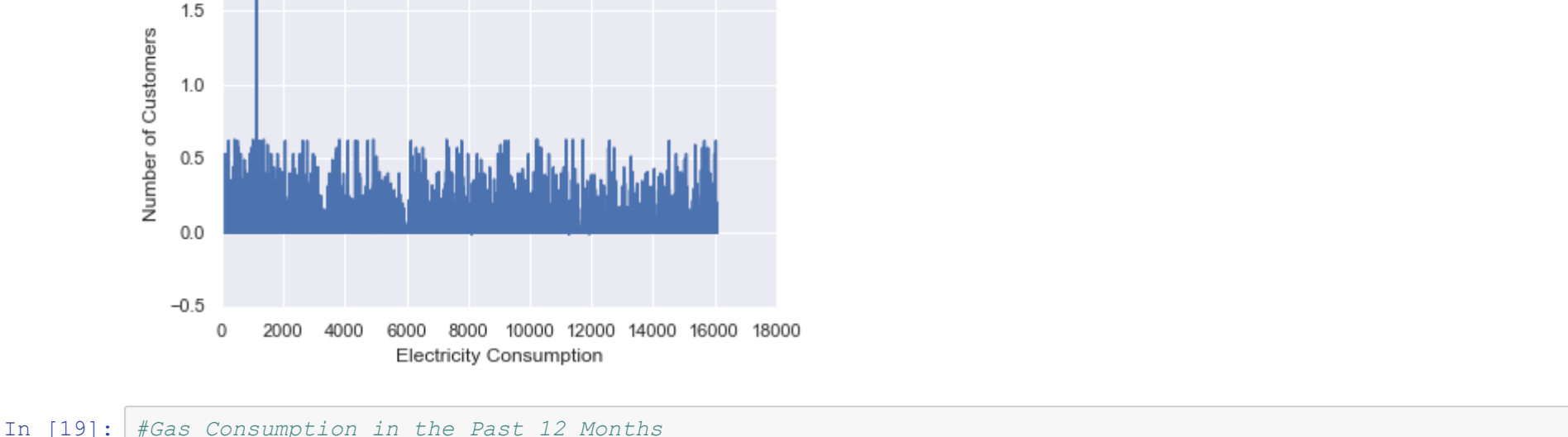
Looking at the rate of those who are also gas clients.

```
In [16]: #Histogram showing number of customers that have gas
plt.hist(cust_data['has_gas'], alpha=1.0)
plt.title('Customers that are Gas Clients')
plt.xlabel('Number of Customers', fontdict = None, labelpad = None)
plt.ylabel('True/False', fontdict = None, labelpad = None)
plt.plot()
```



From the histogram above, we see that a few of the customers are also gas clients. Now lets compare that to those that churn.

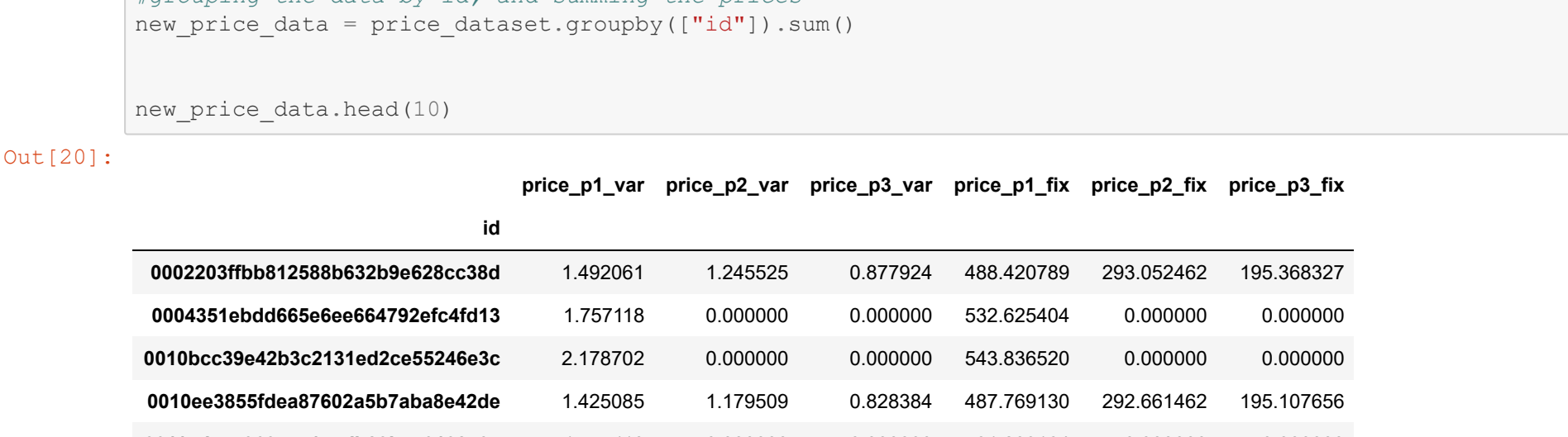
```
In [17]: plt.hist(churn_data['churn'], alpha=1.0)
plt.title('Customer Churn Rate')
plt.xlabel('Number of Customers', fontdict = None, labelpad = None)
plt.ylabel('Churn (No-0, Yes-1)', fontdict = None, labelpad = None)
plt.plot()
```



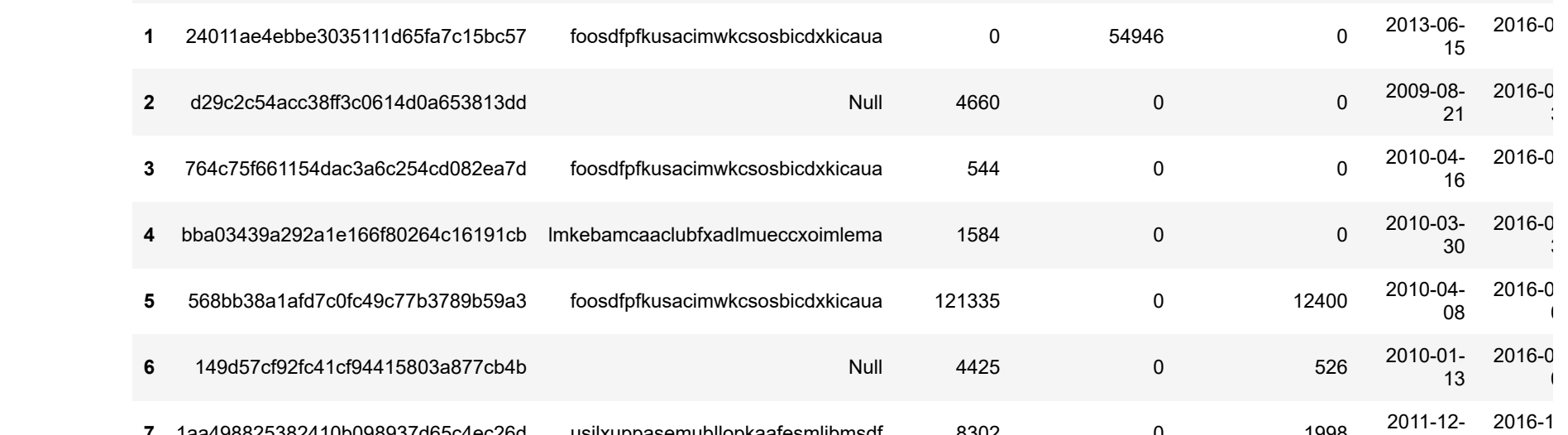
A few Customers have churned. Which means that many Customers haven't churned.

Now Lets compare the Electricity and Gas Consumption for the past 12 months.

```
In [18]: #Electricity Consumption in the past 12 months
plt.plot(cust_data['cons_12m'], alpha=1.0)
plt.title('Electricity Consumption in Past 12 Months')
plt.xlabel('Number of Customers', fontdict = None, labelpad = None)
plt.ylabel('Electricity Consumption', fontdict = None, labelpad = None)
plt.plot()
```



```
In [19]: #Gas Consumption in the Past 12 Months
plt.plot(cust_data['cons_gas_12m'], alpha=1.0)
plt.title('Gas Consumption in Past 12 Months')
plt.xlabel('Number of Customers', fontdict = None, labelpad = None)
plt.ylabel('Gas Consumption', fontdict = None, labelpad = None)
plt.plot()
```



Merge the Dataset

```
In [20]: #dropping the price date, so as to group the data by id
price_dataset = price_data.drop(['price_date'], axis = 1)

#grouping the data by id, and summing the prices
new_price_data = price_dataset.groupby(['id']).sum()
new_price_data.head(10)

Out [20]:
```

	id	price_p1_var	price_p2_var	price_p3_var	price_p1_fix	price_p2_fix	price_p3_fix
0	0002301fbb6b12588b6329a962c8c3bd	1.482081	1.245525	0.877924	488.420789	293.052462	195.368327
1	000435f6b2b3c1313d1e62c54524b23ae3	2.178702	0.000000	0.000000	532.625404	0.000000	0.000000
2	0010cc339a292a1ef668029a4c16191cb	1.437678	0.000000	0.000000	543.936520	0.000000	0.000000
3	0010a385f9a292a1ef668029a4c16191cb	1.425085	1.179699	0.828384	487.769130	292.661462	195.107656
4	0011474983a47f77d78989c70108537	1.775110	0.000000	0.000000	531.203164	0.000000	0.000000
5	001326c87f76427a6d7f1080a6b6b6e	1.437678	1.193008	0.843651	487.932042	292.759214	195.178288
6	001376a29b32f6ad87a1f05899a2d27	1.512911	1.266503	0.899048	488.746620	293.247860	195.498660
7	00194e99727feef73a7b7563f9ab40d6	1.717648	0.000000	0.000000	531.203164	0.000000	0.000000
8	001987ad64ef23a747e7c7233a1f4	1.473007	1.227485	0.8			