```python
In [3]:   import pandas as pd
          import numpy as np
          import seaborn as sns
          import matplotlib.pyplot as plt
          from datetime import datetime
```

```python
In [4]:   # Set plot style
          sns.set(color_codes=True)
```

Load Data

```python
In [5]:   df = pd.read_csv("clean_data_after_eda.csv")
          df["date_activ"]= pd.to_datetime(df["date_activ"], format="%Y-%m-%d")
          df["date_end"]= pd.to_datetime(df["date_end"], format="%Y-%m-%d")
          df["date_modif_prod"] = pd.to_datetime(df["date_modif_prod"], format='%Y-%m-%d')
          df["date_renewal"] = pd.to_datetime(df["date_renewal"], format='%Y-%m-%d')
```

```python
In [6]:   df.head()
```

Out[6]:

| | id | channel_sales | cons_12m | cons_gas_12m | cons_last_month | date_activ | date_end | date_modif_prod | date_renewal | forecast_cons_12m | ... | var_6m_price_off_peak_var | var_... |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| 0 | 24011ae4ebbe3035111d65fa7c15bc57 | foosdfpfkusacimwkcsosbicdxkicaua | 0 | 54946 | 0 | 2013-06-15 | 2016-06-15 | 2015-11-01 | 2015-06-23 | 0.00 | ... | 0.000131 | |
| 1 | d29c2c54acc38ff3c0614d0a653813dd | MISSING | 4660 | 0 | 0 | 2009-08-21 | 2016-08-30 | 2009-08-21 | 2015-08-31 | 189.95 | ... | 0.000003 | |
| 2 | 764c75f661154dac3a6c254cd082ea7d | foosdfpfkusacimwkcsosbicdxkicaua | 544 | 0 | 0 | 2010-04-16 | 2016-04-16 | 2010-04-16 | 2015-04-17 | 47.96 | ... | 0.000004 | |
| 3 | bba03439a292a1e166f80264c16191cb | lmkebamcaaclubfxadlmueccxoimlema | 1584 | 0 | 0 | 2010-03-30 | 2016-03-30 | 2010-03-30 | 2015-03-31 | 240.04 | ... | 0.000003 | |
| 4 | 149d57cf92fc41cf94415803a877cb4b | MISSING | 4425 | 0 | 526 | 2010-01-13 | 2016-03-07 | 2010-01-13 | 2015-03-09 | 445.75 | ... | 0.000011 | |

5 rows × 44 columns

Feature Engineering

```python
In [7]:   price_df = pd.read_csv("price_data.csv")
          price_df["price_date"] = pd.to_datetime(price_df["price_date"], format='%Y-%m-%d')
          price_df.head()
```

Out[7]:

| | id | price_date | price_off_peak_var | price_peak_var | price_mid_peak_var | price_off_peak_fix | price_peak_fix | price_mid_peak_fix |
|---|---|---|---|---|---|---|---|---|
| 0 | 038af19179925da21a25619c5a24b745 | 2015-01-01 | 0.151367 | 0.0 | 0.0 | 44.266931 | 0.0 | 0.0 |
| 1 | 038af19179925da21a25619c5a24b745 | 2015-02-01 | 0.151367 | 0.0 | 0.0 | 44.266931 | 0.0 | 0.0 |
| 2 | 038af19179925da21a25619c5a24b745 | 2015-03-01 | 0.151367 | 0.0 | 0.0 | 44.266931 | 0.0 | 0.0 |
| 3 | 038af19179925da21a25619c5a24b745 | 2015-04-01 | 0.149626 | 0.0 | 0.0 | 44.266931 | 0.0 | 0.0 |
| 4 | 038af19179925da21a25619c5a24b745 | 2015-05-01 | 0.149626 | 0.0 | 0.0 | 44.266931 | 0.0 | 0.0 |

```python
In [9]:   # Group off-peak prices by companies and month
          monthly_price_by_id = price_df.groupby(['id', 'price_date']).agg({'price_off_peak_var': 'mean', 'price_off_peak_fix': 'mean'}).reset_index()

          # Get january and december prices
          jan_prices = monthly_price_by_id.groupby('id').first().reset_index()
          dec_prices = monthly_price_by_id.groupby('id').last().reset_index()

          # Calculate the difference
          diff = pd.merge(dec_prices.rename(columns={'price_off_peak_var': 'dec_1', 'price_off_peak_fix': 'dec_2'}), jan_prices.drop(columns='price_date'), on='id')
          diff['offpeak_diff_dec_january_energy'] = diff['dec_1'] - diff['price_off_peak_var']
          diff['offpeak_diff_dec_january_power'] = diff['dec_2'] - diff['price_off_peak_fix']
          diff = diff[['id', 'offpeak_diff_dec_january_energy','offpeak_diff_dec_january_power']]
          diff.head()
```

Out[9]:

| | id | offpeak_diff_dec_january_energy | offpeak_diff_dec_january_power |
|---|---|---|---|
| 0 | 0002203ffbb812588b632b9e628cc38d | -0.006192 | 0.162916 |
| 1 | 0004351ebdd665e6ee664792efc4fd13 | -0.004104 | 0.177779 |
| 2 | 0010bcc39e42b3c2131ed2ce55246e3c | 0.050443 | 1.500000 |
| 3 | 0010ee3855fdea87602a5b7aba8e42de | -0.010018 | 0.162916 |
| 4 | 00114d74e963e47177db89bc70108537 | -0.003994 | -0.000001 |

```python
In [10]:  df = pd.merge(df, diff, on='id')
          df.head()
```

Out[10]:

| | id | channel_sales | cons_12m | cons_gas_12m | cons_last_month | date_activ | date_end | date_modif_prod | date_renewal | forecast_cons_12m | ... | var_6m_price_mid_peak_var | var_... |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| 0 | 24011ae4ebbe3035111d65fa7c15bc57 | foosdfpfkusacimwkcsosbicdxkicaua | 0 | 54946 | 0 | 2013-06-15 | 2016-06-15 | 2015-11-01 | 2015-06-23 | 0.00 | ... | 9.084737e-04 | |
| 1 | d29c2c54acc38ff3c0614d0a653813dd | MISSING | 4660 | 0 | 0 | 2009-08-21 | 2016-08-30 | 2009-08-21 | 2015-08-31 | 189.95 | ... | 0.000000e+00 | |
| 2 | 764c75f661154dac3a6c254cd082ea7d | foosdfpfkusacimwkcsosbicdxkicaua | 544 | 0 | 0 | 2010-04-16 | 2016-04-16 | 2010-04-16 | 2015-04-17 | 47.96 | ... | 0.000000e+00 | |
| 3 | bba03439a292a1e166f80264c16191cb | lmkebamcaaclubfxadlmueccxoimlema | 1584 | 0 | 0 | 2010-03-30 | 2016-03-30 | 2010-03-30 | 2015-03-31 | 240.04 | ... | 0.000000e+00 | |
| 4 | 149d57cf92fc41cf94415803a877cb4b | MISSING | 4425 | 0 | 526 | 2010-01-13 | 2016-03-07 | 2010-01-13 | 2015-03-09 | 445.75 | ... | 4.860000e-10 | |

5 rows × 46 columns

```python
In [11]:  # Aggregate average prices per period by company
          mean_prices = price_df.groupby(['id']).agg({
              'price_off_peak_var': 'mean',
              'price_peak_var': 'mean',
              'price_mid_peak_var': 'mean',
              'price_off_peak_fix': 'mean',
              'price_peak_fix': 'mean',
              'price_mid_peak_fix': 'mean'
          }).reset_index()
```

```python
In [12]:  # Calculate the mean difference between consecutive periods
          mean_prices['off_peak_peak_var_mean_diff'] = mean_prices['price_off_peak_var'] - mean_prices['price_peak_var']
          mean_prices['peak_mid_peak_var_mean_diff'] = mean_prices['price_peak_var'] - mean_prices['price_mid_peak_var']
          mean_prices['off_peak_mid_peak_var_mean_diff'] = mean_prices['price_off_peak_var'] - mean_prices['price_mid_peak_var']
          mean_prices['off_peak_peak_fix_mean_diff'] = mean_prices['price_off_peak_fix'] - mean_prices['price_peak_fix']
          mean_prices['peak_mid_peak_fix_mean_diff'] = mean_prices['price_peak_fix'] - mean_prices['price_mid_peak_fix']
          mean_prices['off_peak_mid_peak_fix_mean_diff'] = mean_prices['price_off_peak_fix'] - mean_prices['price_mid_peak_fix']
```

```python
In [13]:  columns = [
              'id',
              'off_peak_peak_var_mean_diff',
              'peak_mid_peak_var_mean_diff',
              'off_peak_mid_peak_var_mean_diff',
              'off_peak_peak_fix_mean_diff',
              'peak_mid_peak_fix_mean_diff',
              'off_peak_mid_peak_fix_mean_diff'
          ]
          df = pd.merge(df, mean_prices[columns], on='id')
          df.head()
```

Out[13]:

| | id | channel_sales | cons_12m | cons_gas_12m | cons_last_month | date_activ | date_end | date_modif_prod | date_renewal | forecast_cons_12m | ... | var_6m_price_mid_peak | churn |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| 0 | 24011ae4ebbe3035111d65fa7c15bc57 | foosdfpfkusacimwkcsosbicdxkicaua | 0 | 54946 | 0 | 2013-06-15 | 2016-06-15 | 2015-11-01 | 2015-06-23 | 0.00 | ... | 4.423670e+01 | 1 |
| 1 | d29c2c54acc38ff3c0614d0a653813dd | MISSING | 4660 | 0 | 0 | 2009-08-21 | 2016-08-30 | 2009-08-21 | 2015-08-31 | 189.95 | ... | 0.000000e+00 | 0 |
| 2 | 764c75f661154dac3a6c254cd082ea7d | foosdfpfkusacimwkcsosbicdxkicaua | 544 | 0 | 0 | 2010-04-16 | 2016-04-16 | 2010-04-16 | 2015-04-17 | 47.96 | ... | 0.000000e+00 | 0 |
| 3 | bba03439a292a1e166f80264c16191cb | lmkebamcaaclubfxadlmueccxoimlema | 1584 | 0 | 0 | 2010-03-30 | 2016-03-30 | 2010-03-30 | 2015-03-31 | 240.04 | ... | 0.000000e+00 | 0 |
| 4 | 149d57cf92fc41cf94415803a877cb4b | MISSING | 4425 | 0 | 526 | 2010-01-13 | 2016-03-07 | 2010-01-13 | 2015-03-09 | 445.75 | ... | 4.860000e-10 | 0 |

5 rows × 52 columns

```python
In [14]:  # Aggregate average prices per period by company
          mean_prices_by_month = price_df.groupby(['id', 'price_date']).agg({
              'price_off_peak_var': 'mean',
              'price_peak_var': 'mean',
              'price_mid_peak_var': 'mean',
              'price_off_peak_fix': 'mean',
              'price_peak_fix': 'mean',
              'price_mid_peak_fix': 'mean'
          }).reset_index()
```

```python
In [15]:  # Calculate the mean difference between consecutive periods
          mean_prices_by_month['off_peak_peak_var_mean_diff'] = mean_prices_by_month['price_off_peak_var'] - mean_prices_by_month['price_peak_var']
          mean_prices_by_month['peak_mid_peak_var_mean_diff'] = mean_prices_by_month['price_peak_var'] - mean_prices_by_month['price_mid_peak_var']
          mean_prices_by_month['off_peak_mid_peak_var_mean_diff'] = mean_prices_by_month['price_off_peak_var'] - mean_prices_by_month['price_mid_peak_var']
          mean_prices_by_month['off_peak_peak_fix_mean_diff'] = mean_prices_by_month['price_off_peak_fix'] - mean_prices_by_month['price_peak_fix']
          mean_prices_by_month['peak_mid_peak_fix_mean_diff'] = mean_prices_by_month['price_peak_fix'] - mean_prices_by_month['price_mid_peak_fix']
          mean_prices_by_month['off_peak_mid_peak_fix_mean_diff'] = mean_prices_by_month['price_off_peak_fix'] - mean_prices_by_month['price_mid_peak_fix']
```

```python
In [16]:  # Calculate the maximum monthly difference across time periods
          max_diff_across_periods_months = mean_prices_by_month.groupby(['id']).agg({
              'off_peak_peak_var_mean_diff': 'max',
              'peak_mid_peak_var_mean_diff': 'max',
              'off_peak_mid_peak_var_mean_diff': 'max',
              'off_peak_peak_fix_mean_diff': 'max',
              'peak_mid_peak_fix_mean_diff': 'max',
              'off_peak_mid_peak_fix_mean_diff': 'max'
          }).reset_index().rename(
              columns={
                  'off_peak_peak_var_mean_diff': 'off_peak_peak_var_max_monthly_diff',
                  'peak_mid_peak_var_mean_diff': 'peak_mid_peak_var_max_monthly_diff',
                  'off_peak_mid_peak_var_mean_diff': 'off_peak_mid_peak_var_max_monthly_diff',
                  'off_peak_peak_fix_mean_diff': 'off_peak_peak_fix_max_monthly_diff',
                  'peak_mid_peak_fix_mean_diff': 'peak_mid_peak_fix_max_monthly_diff',
                  'off_peak_mid_peak_fix_mean_diff': 'off_peak_mid_peak_fix_max_monthly_diff'
              }
          )
```

```python
In [17]:  columns = [
              'id',
              'off_peak_peak_var_max_monthly_diff',
              'peak_mid_peak_var_max_monthly_diff',
              'off_peak_mid_peak_var_max_monthly_diff',
              'off_peak_peak_fix_max_monthly_diff',
              'peak_mid_peak_fix_max_monthly_diff',
              'off_peak_mid_peak_fix_max_monthly_diff'
          ]

          df = pd.merge(df, max_diff_across_periods_months[columns], on='id')
          df.head()
```

Out[17]:

| | id | channel_sales | cons_12m | cons_gas_12m | cons_last_month | date_activ | date_end | date_modif_prod | date_renewal | forecast_cons_12m | ... | off_peak_mid_peak_var_max_diff |
|---|---|---|---|---|---|---|---|---|---|---|---|---|
| 0 | 24011ae4ebbe3035111d65fa7c15bc57 | foosdfpfkusacimwkcsosbicdxkicaua | 0 | 54946 | 0 | 2013-06-15 | 2016-06-15 | 2015-11-01 | 2015-06-23 | 0.00 | ... | 0.058253 |
| 1 | d29c2c54acc38ff3c0614d0a653813dd | MISSING | 4660 | 0 | 0 | 2009-08-21 | 2016-08-30 | 2009-08-21 | 2015-08-31 | 189.95 | ... | 0.149609 |
| 2 | 764c75f661154dac3a6c254cd082ea7d | foosdfpfkusacimwkcsosbicdxkicaua | 544 | 0 | 0 | 2010-04-16 | 2016-04-16 | 2010-04-16 | 2015-04-17 | 47.96 | ... | 0.170512 |
| 3 | bba03439a292a1e166f80264c16191cb | lmkebamcaaclubfxadlmueccxoimlema | 1584 | 0 | 0 | 2010-03-30 | 2016-03-30 | 2010-03-30 | 2015-03-31 | 240.04 | ... | 0.151210 |
| 4 | 149d57cf92fc41cf94415803a877cb4b | MISSING | 4425 | 0 | 526 | 2010-01-13 | 2016-03-07 | 2010-01-13 | 2015-03-09 | 445.75 | ... | 0.051309 |

5 rows × 58 columns

```python
In [ ]:
```