

Exploratory Data Analysis & Data Cleaning

Author: Patrick Duo

The Datasets

`ml_case_training_output.csv` named as `pco_output` contains:

- id: contact id
- churned: has the client churned over the next 3 months

`ml_case_training_hist_data.csv` named as `pco_hist` contains the history of energy and power consumption per client:

- id: contact id
- activity_new: reference date
- price_p1_var: price of energy for the 1st period
- price_p2_var: price of energy for the 2nd
- periodprice_p3_var: price of energy for the 3rd period
- price_p1_fix: price of power for the 1st period
- price_p2_fix: price of power for the 2nd period
- price_p3_fix: price of power for the 3rd period

`ml_case_training_data.csv` contains:

- id: contact id
- activity_new: category of the company's activity.
- campaign_disc_elec: code of the electricity campaign the customer last subscribed to.
- channel_sales: code of the sales channel
- cons_12m: electricity consumption of the past 12 months
- cons_gas_12m: gas consumption of the past 12 months
- cons_last_month: electricity consumption of the last month
- date_active: date of activation of the contract
- date_end: registered date of the end of the contract
- date_first_activ: date of first contract of the client
- date_modif_prod: date of last modification of the product
- date_renewal: date of the next contract renewal
- forecast_base_bill_ele: forecasted electricity bill baseline for next month
- forecast_base_bill_year: forecasted electricity bill baseline for calendar year
- forecast_bill_12m: forecasted electricity bill baseline for 12 months
- forecast_cons: forecasted electricity consumption for next month
- forecast_cons_12m: forecasted electricity consumption for next 12 months
- forecast_cons_year: forecasted electricity consumption for next calendar year
- forecast_discount_energy: forecasted value of current discount
- forecast_meter_rent_12m: forecasted bill of meter rental for the next 12 months
- forecast_price_energy_p1: forecasted energy price for 1st period
- forecast_price_energy_p2: forecasted energy price for 2nd period
- forecast_price_pow_p1: forecasted power price for 1st period
- has_gas: indicated if client is also a gas client
- imp_cons: current paid consumption
- margin_gross_pow_ele: gross margin on power subscription
- margin_net_pow_ele: net margin on power subscription
- nb_prod_act: number of active products and services
- net_margin: net net margin
- num_years_antig: antiquity of the client (in number of years)
- origin_up: code of the electricity campaign the customer first subscribed to
- pow_max: subscribed power

Import Libraries

```
In [1]: import pandas as pd
import numpy as np
import seaborn as sns
import matplotlib.pyplot as plt
import missingno as msno
from scipy.stats import zscore as zscore
```

Load Data

```
In [2]: # list = ('date_active', 'date_end', 'date_first_activ', 'date_modif_prod', 'date_renewal')
dt_list = ('date_active', 'date_end', 'date_first_activ', 'date_modif_prod', 'date_renewal')

pco_main = pd.read_csv('ml_case_training_data.csv', parse_dates=dt_list)
pco_hist = pd.read_csv('ml_case_training_hist_data.csv', parse_dates=['price_date'])
pco_output = pd.read_csv('ml_case_training_output.csv')
pd.set_option('display.max_columns', None)
```

Main Dataset

```
In [3]: pco_main.head()

Out[3]:
```

	id	activity_new	campaign_disc_elec	channel_sales	cons_1
0	48ada5226167cf58715202705a0451c9	essoiHdJbKcsuxmfuacbdckomxiw		NaN	309
1	24011ae4eb3e031165f5a7c15bc57		NaN	fosofdfpfkusacimwksosbdcidckiaua	
2	429c254acc38ff3c0614d0a0653813dd		NaN		41
3	764c7f66119d43ac3ae254cd082ea7d		NaN	foosdfpfkusacimwksosbdcidckiaua	
4	bb0a3439a292a1e160180264c16191cb		NaN	lmkebancaacubf7kadmruccuocniema	1

```
In [4]: pco_main.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_elec 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 16096 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 15939 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 15970 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 float64
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 float64
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 net_margin 16081 non-null float64
29 num_years_antig 16096 non-null float64
30 origin_up 16009 non-null object
31 pow_max 16093 non-null float64
dtypes: datetime64[ns](5), float64(16), int64(6), object(5)
memory usage: 3.1+ MB
```

```
In [5]: # Percentage of nullity by column
missing_perc = pco_main.isnull().mean() * 100
print('Percentage of Missing Values:\n', missing_perc)
```

```
Percentage of Missing Values:
id 0.000000
activity_new 59.300447
campaign_disc_elec 126.000000
channel_sales 26.205268
cons_12m 0.000000
cons_gas_12m 0.000000
cons_last_month 0.000000
date_active 0.000000
date_end 0.012425
date_first_activ 78.205765
date_modif_prod 0.975398
date_renewal 0.246509
forecast_base_bill_ele 78.205765
forecast_base_bill_year 78.205765
forecast_bill_12m 78.205765
forecast_cons 78.205765
forecast_cons_12m 0.000000
forecast_cons_year 0.000000
forecast_discount_energy 0.782803
forecast_meter_rent_12m 0.000000
forecast_price_energy_p1 0.782803
forecast_price_energy_p2 0.782803
forecast_price_pow_p1 0.782803
has_gas 0.000000
imp_cons 0.000000
margin_gross_pow_ele 0.080765
margin_net_pow_ele 0.080765
nb_prod_act 0.000000
net_margin 0.093191
num_years_antig 0.000000
origin_up 0.545007
pow_max 0.018638
dtype: float64
```

```
In [6]: # Descriptive statistics
pco_main.describe()
```

```
Out[6]:
```

	campaign_disc_elec	cons_12m	cons_gas_12m	cons_last_month	forecast_base_bill_ele	forecast_base_bill_year	forecast_bill_12m
count	0.0	1.609600e+04	1.609600e+04	1.609600e+04	3508.000000	3508.000000	3508.000000
mean	NaN	1.948044e+05	3.191864e+05	1.946154e+04	335.843857	335.843857	383.7
std	NaN	6.795151e+05	1.775885e+05	8.238576e+04	649.406000	649.406000	5425.
min	NaN	-1.252760e+05	-3.037008e+03	-9.138607e+04	-364.940000	-364.940000	-2503.
25%	NaN	5.962506e+03	0.000000e+00	0.000000e+00	0.000000	0.000000	1158.
50%	NaN	1.533250e+04	0.000000e+00	9.000000e+02	162.955000	162.955000	2167.
75%	NaN	5.022150e+04	0.000000e+00	4.127000e+03	396.185000	396.185000	4246.
max	NaN	1.609711e+07	4.188440e+06	4.538720e+06	12566.080000	12566.080000	81122.

Observations

- 14 columns have negative minimum values.
- `campaign_disc_elec` column is missing completely.
- `activity_new` column is missing 59.3%.
- `date_first_activ`, `forecast_base_bill_ele`, `forecast_base_bill_year`, `forecast_bill_12m`, and `forecast_cons` columns are each missing 78.2%.

The History Dataset

```
In [7]: pco_hist.head()

Out[7]:
```

	id	price_date	price_p1_var	price_p2_var	price_p3_var	price_p1_fix	price_p2_fix	price_p3_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

```
In [8]: pco_hist.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 datetime64[ns]
2 price_p1_var 191643 non-null float64
3 price_p2_var 191643 non-null float64
4 price_p3_var 191643 non-null float64
5 price_p1_fix 191643 non-null float64
6 price_p2_fix 191643 non-null float64
7 price_p3_fix 191643 non-null float64
dtypes: datetime64[ns](1), float64(6), object(1)
memory usage: 11.8+ MB
```

```
In [9]: # Percentage of nullity by column
missing_perc = pco_hist.isnull().mean() * 100
print('Percentage of Missing Values:\n', missing_perc)
```

```
Percentage of Missing Values:
id 0.000000
price_date 0.000000
price_p1_var 0.704138
price_p2_var 0.704138
price_p3_var 0.704138
price_p1_fix 0.704138
price_p2_fix 0.704138
price_p3_fix 0.704138
dtype: float64
```

```
In [10]: # Descriptive statistics
pco_hist.describe()
```

```
Out[10]:
```

	price_p1_var	price_p2_var	price_p3_var	price_p1_fix	price_p2_fix	price_p3_fix
count	191643.000000	191643.000000	191643.000000	191643.000000	191643.000000	191643.000000
mean	0.140991	0.054412	0.030712	43.325563	10.698201	6.455436
std	0.025117	0.050033	0.036335	5.437952	12.856039	7.882273
min	0.000000	0.000000	0.000000	-0.177779	-0.097752	-0.065172
25%	0.125976	0.000000	0.000000	40.728885	0.000000	0.000000
50%	0.146033	0.088483	0.000000	44.266930	0.000000	0.000000
75%	0.151635	0.101780	0.072558	44.444710	24.339581	16.226389
max	0.280700	0.229788	0.114102	59.444710	36.490692	17.458221

Observations

- `price_p1_var`, `price_p2_var`, `price_p3_var`, `price_p1_fix`, `price_p2_fix`, `price_p3_fix` are missing 70.4% values.
- `price_p1_fix`, `price_p2_fix`, `price_p3_fix` contain negative values, which doesn't make sense for price of power.

The Output Dataset

```
In [11]: pco_output.head()

Out[11]:
```

	id	churn
0	48ada5226167cf58715202705a0451c9	0
1	24011ae4eb3e031165f5a7c15bc57	1
2	429c254acc38ff3c0614d0a0653813dd	0
3	764c7f66119d43ac3ae254cd082ea7d	0
4	bb0a3439a292a1e160180264c16191cb	0

```
In [12]: pco_output.info()

<class 'pandas.core.frame.DataFrame'>
RangeIndex: 16096 entries, 0 to 16095
Data columns (total 2 columns):
# Column Non-Null Count Dtype
---
0 id 16096 non-null object
1 churn 16096 non-null int64
dtypes: int64(1), object(1)
memory usage: 251.6+ KB
```

```
In [13]: # Percentage of nullity by column
missing_perc = pco_output.isnull().mean() * 100
print('Percentage of Missing Values:\n', missing_perc)
```

```
Percentage of Missing Values:
id 0.0
churn 0.0
dtype: float64
```

```
In [14]: # Descriptive statistics
pco_output.describe()
```

```
Out[14]:
```

	churn
count	16096.000000
mean	0.099093
std	0.298796
min	0.000000
25%	0.000000
50%	0.000000
75%	0.000000
max	1.000000

Observations

- Complete dataset.

Data Cleaning and Imputation

Missing Data

Types of missingness

Missing Completely at Random (MCAR)

Missingness has no relationship between any values, observed or missing

Missing at Random (MAR)

There is a systematic relationship between missingness and other observed data, but not the missing data

Missing Not at Random (MNAR)

There is a relationship between missingness and its values, missing or non-missing

The History Dataset

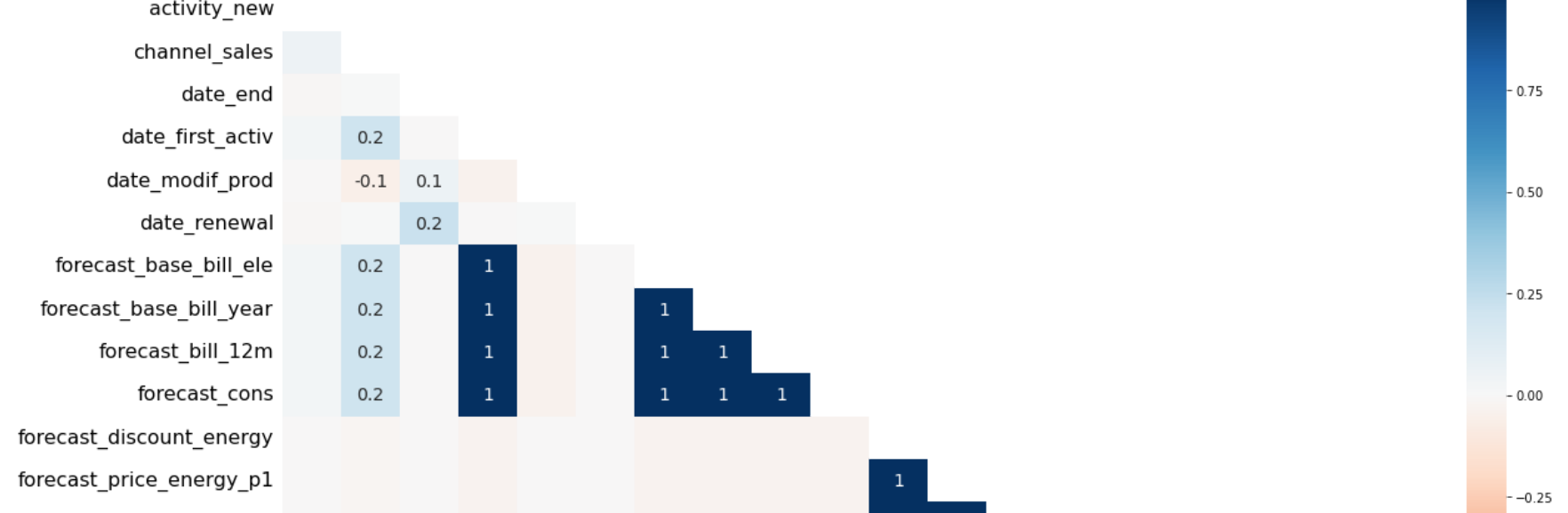
```
In [15]: # Identify negative columns
negative_cols = ['price_p1_fix', 'price_p2_fix', 'price_p3_fix']
# Convert to positive the negative columns in pco_hist
pco_hist[negative_cols] = pco_hist[negative_cols].apply(abs)

pco_hist.describe()
```

```
Out[15]:
```

	price_p1_var	price_p2_var	price_p3_var	price_p1_fix	price_p2_fix	price_p3_fix
count	191643.000000	191643.000000	191643.000000	191643.000000	191643.000000	191643.000000
mean	0.140991	0.054412	0.030712	43.325563	10.698201	6.455436
std	0.025117	0.050033	0.036335	5.437952	12.856039	7.882273
min	0.000000	0.000000	0.000000	0.000000	0.000000	0.000000
25%	0.125976	0.000000	0.000000	40.728885	0.000000	0.000000
50%	0.146033	0.088483	0.000000	44.266930	0.000000	0.000000
75%	0.151635	0.101780	0.072558	44.444710	24.339581	16.226389
max	0.280700	0.229788	0.114102	59.444710	36.490692	17.458221

```
In [16]: # Visualize the completeness of the dataframe
msno.bar(pco_hist)
plt.show()
```



```
In [17]: # Visualize the locations of the missing values of the dataset
sorted = pco_hist.sort_values(by = ['id', 'price_date'])
msno.matrix(sorted)
plt.show()
```



```
In [18]: # Visualize the correlation between the numeric variables of the dataframe
plt.heatmap(pco_hist)
plt.show()
```



```
In [19]: # Identify the index of the IDs containing missing values.
hist_NAN_index = pco_hist[pco_hist.isnull().any(axis=1)].index.values.tolist()

# Obtain a dataframe with the missing values
pco_hist_missing = pco_hist.loc[hist_NAN_index:]

# Glimpse at the NaN cases of the pco_hist dataset
pco_hist_missing.head(10)
```

```
Out[19]:
```

	id	price_date	price_p1_var	price_p2_var	price_p3_var	price_p1_fix	price_p2_fix	price_p3_fix
75	e7716222b9d7b79a8dfc0b3c31e3575560	2015-04-01	NaN	NaN	NaN	NaN	NaN	NaN
221	0f52310b2ebac8b628da8eaab343	2015-06-01	NaN	NaN	NaN	NaN	NaN	NaN
377	2793639cde582fndf3e386ce0c8d8f35	2015-06-01	NaN	NaN	NaN	NaN	NaN	NaN
413	783c1ab3caf1902b1df4072820243c	2015-06-01	NaN	NaN	NaN	NaN	NaN	NaN
461	307f6c64d060e12a049017009b09d3f3	2015-06-01	NaN	NaN	NaN	NaN	NaN	NaN
471	33bb3af90650ac2e9eac8ff2c975a6b	2015-04-01	NaN	NaN	NaN	NaN	NaN	NaN
472	33bb3af90650ac2e9eac8ff2c975a6b	2015-05-01	NaN	NaN	NaN	NaN	NaN	NaN
475	33bb3af90650ac2e9eac8ff2c975a6b	2015-08-01	NaN	NaN	NaN	NaN	NaN	NaN
476	33bb3af90650ac2e9eac8ff2c975a6b	2015-09-01	NaN	NaN	NaN	NaN	NaN	NaN
874	0e90710b08183c09546e827e4025647	2015-12-01	NaN	NaN	NaN	NaN	NaN	NaN

```
In [20]: # extract the unique dates of missing data
date_list = pco_hist_missing['price_date'].unique()
id_list = pco_hist_missing['id'].unique()

# Create a time dataframe with the unique dates
time_df = pd.DataFrame(data=date_list, columns=['price_date'])

# Glimpse the time dataframe
time_df.sort_values(by=['price_date'])
```

```
Out[20]:
```

	price_date
9	2015-01-01
11	2015-02-01
8	2015-03-01
0	2015-04-01
2	2015-05-01
1	2015-06-01
10	2015-07-01
3	2015-08-01
4	2015-09-01
7	2015-10-01
6	2015-11-01
5	2015-12-01

Observations

- `activity_new` is MCAR with low correlation with other variables. Can drop this column
- `campaign_disc_elec` is MCAR. Can drop this column. Suggests that subscribers are not subscribing through campaign offers.
- `date_first_activ` cannot replace `date_active`. MAR
- `net_margin` has strong correlation between `margin_gross_pow_ele` and `margin`. `net_pow_ele`. Suggests multi-collinearity.
- `origin_up` and `pow_max` is MCAR. Can drop.
- `Forecast_base_bill_ele`, `forecast_base_bill_year`, `forecast_bill_12m` and `forecast_cons` variables are highly-correlated with `date_first_activ`. MNAR

```
In [31]: # Choose the columns without missing values
incomplete_cols = ['channel_sales', 'date_first_activ', 'forecast_base_bill_ele', 'forecast_base_bill_year', 'forecast_cons', 'forecast_discount_energy', 'forecast_meter_rent_12m', 'forecast_price_energy_p1', 'forecast_price_energy_p2', 'forecast_price_pow_p1']
complete_cols = [column_name for column_name in pco_main.drop(columns=incomplete_cols) if column_name not in incomplete_cols]
pco_main_cc = pco_main.drop(complete_cols)

# Fix negative numeric variables
numeric = [column_name for column_name in pco_main_cc.columns if pco_main_cc[column_name].dtype == 'float64']
pco_main_cc[numeric] = pco_main_cc[numeric].abs()

# Overwrite positive values on negative values
pco_main_cc[numeric] = pco_main_cc[numeric].apply(abs)

# Describe
pco_main_cc.describe()
```

```
Out[31]:
```

	date_active	date_first_activ
count	3508	3508
mean	2011-09-03 07:45:05.13128832	2011-06-19 20:20:23.26117440
min	2003-09-23 00:00:00	2001-01-10 00:00:00
25%	2012-01-26 00:00:00	2010-08-04 18:00:00
50%	2012-01-03 00:00:00	2011-10-28 00:00:00
75%	2012-08-08 00:00:00	2012-06-22 00:00:00
max	2014-09-01 00:00:00	2014-09-01 00:00:00

```
In [28]: # Drop the column activity_new and campaign_disc_elec
pco_main_drop = pco_main.drop(labels=['activity_new', 'campaign_disc_elec'], axis=1)

# Remove date_end and date_modif_prod date_renewal origin_up pow_max margin_gross_pow_ele margin_net_pow_ele net_margin
brush = ['date_end', 'date_modif_prod', 'date_renewal', 'origin_up', 'pow_max', 'margin_gross_pow_ele', 'margin_net_pow_ele', 'net_margin', 'forecast_discount_energy', 'forecast_price_energy_p1', 'forecast_price_energy_p2', 'forecast_price_pow_p1
```