

```
In [1]: try:
        from Python import get_ipython
        get_ipython().magic('clear')
        get_ipython().magic('reset -f')
    except
        pass

In [2]: #Importing all the required libraries
import pandas as pd
import numpy as np
import matplotlib.pyplot as plt
%matplotlib inline
import seaborn as sns
import missingno as msno
from scipy.stats import zscore
%ignore the warnings
import warnings as wg
wg.filterwarnings("ignore")

In [3]: # Set plot style
sns.set(color_codes=True)

In [4]: # Set maximum number of columns to be displayed
pd.set_option('display.max_columns', 100)

In [5]: #Read the dataset file of train data
train_data=pd.read_csv("ml_case_training_data.csv")
#Read the dataset file of history data
history_data=pd.read_csv("ml_case_training_hist_data.csv")
#Read the dataset file of churn data
churn_data=pd.read_csv("ml_case_training_output.csv")

In [6]: train = pd.merge(train_data, churn_data, on='id')

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

	id	activity_new	campaign_disc_ele	channel_sales	cons_12m	cons_gas_12m	cons_last_month	date_activ	date_end	date_n
0	48ada5226167058715202705a0451c9	esoiHdtkBtksclumfwjcdkdkommiwv	NaN	lmketbancaacubhadiuueccomiema	309275	0	10025	2012-11-07	2016-11-06	
1	24011a4e4bbe335111059a7150c57	NaN	NaN	foosdfkfusacimwkcsoabckdskcaua	0	54946	0	2013-06-15	2016-06-15	
2	029c2c54acc38f3c0614d0a0653813d0	NaN	NaN	NaN	4660	0	0	2009-08-21	2016-08-21	
3	764c75661154dc3a6c254c0826a7d	NaN	NaN	foosdfkfusacimwkcsoabckdskcaua	544	0	0	2010-04-16	2016-04-16	
4	bba03439a292a1e16680264c16191cb	NaN	NaN	lmketbancaacubhadiuueccomiema	1584	0	0	2010-03-30	2016-03-30	

```
In [8]: pd.DataFrame({"Data type":train.dtypes})

Out[8]:
```

	Data type
	id object
	activity_new object
	campaign_disc_ele float64
	channel_sales object
	cons_12m int64
	cons_gas_12m int64
	cons_last_month int64
	date_activ object
	date_end object
	date_first_activ object
	date_modif_prod object
	date_renewal object
	forecast_base_bill_ele float64
	forecast_base_bill_year float64
	forecast_bill_12m float64
	forecast_cons float64
	forecast_cons_12m float64
	forecast_cons_year int64
	forecast_discount_energy float64
	forecast_meter_rent_12m float64
	forecast_price_energy_p1 float64
	forecast_price_energy_p2 float64
	forecast_price_pow_p1 float64
	has_gas object
	imp_cons float64
	margin_gross_pow_ele float64
	margin_net_pow_ele float64
	nb_prod_act int64
	net_margin float64
	num_years_antig int64
	origin_up object
	pow_max float64
	churn int64

```
In [9]: pd.DataFrame({"Data type":history_data.dtypes})

Out[9]:
```

	Data type
	id object
	price_date object
	price_p1_var float64
	price_p2_var float64
	price_p3_var float64
	price_p1_fix float64
	price_p2_fix float64
	price_p3_fix float64

```
In [10]: train.describe()

Out[10]:
```

	campaign_disc_ele	cons_12m	cons_gas_12m	cons_last_month	forecast_base_bill_ele	forecast_base_bill_year	forecast_bill_12m	forecast_cons	forecast_cons_12m	forecast_cons_year	for
count	0.0	1.609600e+04	1.609600e+04	1.609600e+04	3508.000000	3508.000000	3508.000000	3508.000000	16096.000000	16096.000000	
mean	NaN	1.948044e+05	3.191164e+04	1.946154e+04	335.843857	3837.441866	206.845165	2370.555949	1907.347229		
std	NaN	6.795191e+05	1.775885e+05	8.236670e+04	649.400000	649.400000	5425.744327	455.834288	4305.085664	5257.364759	
min	NaN	1.252700e+05	-3.037000e+03	-9.138600e+00	-0.000000e+00	-0.000000e+00	-364.840000	-2503.480000	-0.000000	-16089.260000	-86527.000000
25%	NaN	5.906250e+03	0.000000e+00	0.000000e+00	0.000000	0.000000	1158.175000	0.000000	513.230000	0.000000	
50%	NaN	1.533250e+04	0.000000e+00	9.010000e+02	162.955000	162.955000	2187.230000	42.215000	1179.180000	378.000000	
75%	NaN	5.022150e+04	0.000000e+00	4.127000e+03	396.185000	396.185000	4246.555000	228.117000	2692.977500	1994.250000	
max	NaN	1.609711e+07	4.188440e+06	4.538720e+06	12566.080000	12566.080000	61122.630000	9682.890000	103801.930000	175375.000000	

```
In [11]: train["campaign_disc_ele"].isnull().values.all()

Out[11]: True

In [12]: history_data.describe()

Out[12]:
```

	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	43.325546	10.698201	6.455436
std	0.025117	0.050033	0.036335	5.437952	12.856046	7.782279
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	0.085483	0.000000	44.266930	0.000000	0.000000
75%	0.151635	0.101780	0.072558	44.444710	24.335681	16.226389
max	0.280700	0.229788	0.114102	59.444710	36.490692	17.458221

```
In [13]: pd.DataFrame({"Missing values (%)": train.isnull().sum()/len(train.index)*100})

Out[13]:
```

	Missing values (%)
id	0.000000
activity_new	59.300447
campaign_disc_ele	100.000000
channel_sales	26.202688
cons_12m	0.000000
cons_gas_12m	0.000000
cons_last_month	0.000000
date_activ	0.000000
date_end	0.012425
date_first_activ	78.205765
date_modif_prod	0.975398
date_renewal	0.248509
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	0.782083
forecast_meter_rent_12m	0.000000
forecast_price_energy_p1	0.782083
forecast_price_energy_p2	0.782083
forecast_price_pow_p1	0.782083
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	0.003191
num_years_antig	0.000000
origin_up	0.540507
pow_max	0.018638
churn	0.000000

```
In [14]: pd.DataFrame({"Missing values (%)": history_data.isnull().sum()/len(history_data.index)*100})

Out[14]:
```

	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

```
In [15]: train[train['date_end']==].isnull()

Out[15]:
```

	id	activity_new	campaign_disc_ele	channel_sales	cons_12m	cons_gas_12m	cons_last_month	date_activ	date_end	date
4452	f7ae022079494de7607c74fc6a353	fkhsaopvcpkpswhbkdzfpamfwelpwi	NaN	foosdfkfusacimwkcsoabckdskcaua	107905	0	0	2013-06-19	NaN	
12880	edc0b46866d67e9a233cd9f834835	mioxdtbtoxfuakuewewdtdtfs	NaN	NaN	18954	0	1596	2010-09-06	NaN	

```
In [16]: train.drop(['campaign_disc_ele','date_first_activ','forecast_base_bill_ele','forecast_cons','forecast_bill_12m','forecast_base_bill_year'],inplace=True,axis=1)

In [17]: train.isnull().sum()

Out[17]:
```

id	activity_new	channel_sales	cons_12m	cons_gas_12m	cons_last_month	date_activ	date_end	date_modif_prod	date_renewal	forecast_cons_12m	forecast_cons_year	forecast_discount_energy	forecast_meter_rent_12m	forecast_price_energy_p1	forecast_price_energy_p2	forecast_price_pow_p1	has_gas	imp_cons	margin_gross_pow_ele	margin_net_pow_ele	nb_prod_act	net_margin	num_years_antig	origin_up	pow_max	churn	dtype
0	9545	4218	0	0	0	0	0	157	49	0	0	126	0	126	126	126	0	0	13	13	0	0	15	87	3	0	int64

```
In [18]: df=train.head()

In [19]: df['input_nan(df.variable,median):
df(variable)=df(variable).fillna(median)

In [20]: median=df.net_margin.median()

In [21]: median

Out[21]: 25.46

In [22]: [input_nan(df,'forecast_discount_energy',median)
input_nan(df,'forecast_price_energy_p1',median)
input_nan(df,'forecast_price_energy_p2',median)
input_nan(df,'forecast_price_pow_p1',median)
input_nan(df,'margin_gross_pow_ele',median)
input_nan(df,'margin_net_pow_ele',median)
input_nan(df,'net_margin',median)
input_nan(df,'pow_max',median)]

In [23]: def input_nan(df,variable):
    most_frequent_category=df[variable].mode()[0]
    df[variable].fillna(most_frequent_category,inplace=True)

In [24]: for feature in ['activity_new','channel_sales','date_modif_prod']:
    input_nan(df,feature)

In [25]: df.isnull().sum()

Out[25]:
```

id	activity_new	channel_sales	cons_12m	cons_gas_12m	cons_last_month	date_activ	date_end	date_modif_prod	date_renewal	forecast_cons_12m	forecast_cons_year	forecast_discount_energy	forecast_meter_rent_12m	forecast_price_energy_p1	forecast_price_energy_p2	forecast_price_pow_p1	has_gas	imp_cons	margin_gross_pow_ele	margin_net_pow_ele	nb_prod_act	net_margin	num_years_antig	origin_up	pow_max	churn	dtype
0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	int64

```
In [26]: import matplotlib.pyplot as plt
%matplotlib inline

In [27]: df.forecast_discount_energy.hist(bins=50)

Out[27]: <AxesSubplot:~>
```

```
In [28]: import seaborn as sns
sns.boxplot('forecast_price_energy_p1',data=df)

Out[28]: <AxesSubplot:~label='forecast_price_energy_p1'~>
```

```
In [29]: import seaborn as sns
sns.boxplot('forecast_price_energy_p2',data=df)

Out[29]: <AxesSubplot:~label='forecast_price_energy_p2'~>
```

```
In [30]: df.forecast_price_energy_p2.hist(bins=50)

Out[30]: <AxesSubplot:~>
```

```
In [31]: # Fill missing values with mean(history_data.mean(), inplace=True)
# Count the number of NaN values in each column
print(history_data.isnull().sum())

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

In [32]: history_data['price_p1_var'].hist(bins=50)

Out[32]: <AxesSubplot:~>
```

```
In [33]: import seaborn as sns
sns.boxplot('price_p2_var',data=history_data)

Out[33]: <AxesSubplot:~label='price_p2_var'~>
```

```
In [34]: import seaborn as sns
sns.boxplot('price_p3_var',data=history_data)

Out[34]: <AxesSubplot:~label='price_p3_var'~>
```

```
In [35]: import seaborn as sns
sns.boxplot('price_p1_fix',data=history_data)

Out[35]: <AxesSubplot:~label='price_p1_fix'~>
```

```
In [36]: import seaborn as sns
sns.boxplot('price_p2_fix',data=history_data)

Out[36]: <AxesSubplot:~label='price_p2_fix'~>
```

```
In [37]: import seaborn as sns
sns.boxplot('price_p3_fix',data=history_data)

Out[37]: <AxesSubplot:~label='price_p3_fix'~>
```

```
In [40]: min_threshod = history_data['price_p1_fix'].quantile(0.95)
max_threshod

Out[40]: 46.44471036

In [42]: history_data[history_data['price_p1_fix']>max_threshod]

Out[42]:
```

	id	price_date	price_p1_var	price_p2_var	price_p3_var	price_p1_fix	price_p2_fix	price_p3_fix
528	45f7bc48222a962099877796a76114	2015-01-01	0.107491	0.094000	0.068317	59.173468	36.490669	8.367731
529	45f7bc48222a962099877796a76114	2015-02-01	0.107491	0.094000	0.068317	59.173468	36.490669	8.367731
530	45f7bc48222a962099877796a76114	2015-03-01	0.107491	0.094000	0.068317	59.173468	36.490669	8.367731
531	45f7bc48222a962099877796a76114	2015-04-01	0.107491	0.094000	0.068317	59.173468	36.490669	8.367731
532	45f7bc48222a962099877796a76114	2015-05-01	0.107491	0.094000	0.068317	59.173468	36.490669	8.367731
192913	4ce4a65723a68718f950d070e5cd04129a	2015-08-01	0.104531	0.094642	0.071661	58.936780	36.344721	8.334263
192914	4ce4a65723a68718f950d070e5cd04129a	2015-09-01	0.100262	0.090362	0.067368	58.936780	36.344721	8.334263
192915	4ce4a65723a68718f950d070e5cd04129a	2015-10-01	0.100262	0.090362	0.067368	58.936780	36.344721	8.334263
192916	4ce4a65723a68718f950d070e5cd04129a	2015-11-01	0.098414	0.089947	0.066059	58.936780	36.490669	8.367731
192917	4ce4a65723a68718f950d070e5cd04129a	2015-12-01	0.098414	0.089947	0.066059	59.173468	36.490669	8.367731

8971 rows x 8 columns

```
In [44]: min_threshod = history_data['price_p1_fix'].quantile(0.05)
min_threshod

Out[44]: 46.5659694

In [46]: history_data[history_data['price_p1_fix']<min_threshod]

Out[46]:
```

	id	price_date	price_p1_var	price_p2_var	price_p3_var	price_p1_fix	price_p2_fix	price_p3_fix
12	31f0c549924679a3cb02d12f8a9ea43	2015-01-01	0.125976	0.103395	0.071536	40.565969	24.339581	16.226389
13	31f0c549924679a3cb02d12f8a9ea43	2015-02-01	0.125976	0.103395	0.071536	40.565969	24.339581	16.226389
14	31f0c549924679a3cb02d12f8a9ea43	2015-03-01	0.125976	0.103395	0.071536	40.565969	24.339581	16.226389
15	31f0c549924679a3cb02d12f8a9ea43	2015-04-01	0.125976	0.103395	0.071536	40.565969	24.339581	16.226389
16	31f0c549924679a3cb02d12f8a9ea43	2015-05-01	0.124815	0.102234	0.070375	40.565969	24.339581	16.226389
192998	10f51cd02baa19a0c940e1b3b01705	2015-01-01	0.119916	0.100682	0.087111	40.565969	24.339581	16.226389
192999	10f51cd02baa19a0c940e1b3b01705	2015-02-01	0.129444	0.100682	0.075004	40.565969	24.339581	16.226389
192999	10f51cd02baa19a0c940e1b3b01705	2015-03-01	0.129444	0.106863	0.075004	40.565969	24.339581	16.226389
192999	10f51cd02baa19a0c940e1b3b01705	2015-04-01	0.129444	0.106863	0.075004	40.565969	24.339581	16.226389

24125 rows x 8 columns

```
In [46]: history_data[history_data['price_p1_fix']<max_threshod & (history_data['price_p1_fix']>min_threshod)]

Out[46]:
```

	id	price_date	price_p1_var	price_p2_var	price_p3_var	price_p1_fix	price_p2_fix	price_p3_fix
0	038af19179925da21a25619c5a246745	2015-01-01	0.151367	0.000000	0.000000	44.266931	0.000000	0.000000
1	038af19179925da21a25619c5a246745	2015-02-01	0.151367	0.000000	0.000000	44.266931	0.000000	0.000000
2	038af19179925da21a25619c5a246745	2015-03-01	0.151367	0.000000	0.000000	44.266931	0.000000	0.000000
3	038af19179925da21a25619c5a246745	2015-04-01	0.151367	0.000000	0.000000	44.266931	0.000000	0.000000
4	038af19179925da21a25619c5a246745	2015-05-01	0.149626	0.000000	0.000000	44.266931	0.000000	0.000000
...
192997	10f51cd02baa19a0c940e1b3b01705	2015-08-01	0.119916	0.102232	0.076257	40.728885	24.43733	16.291555
192998	10f51cd02baa19a0c940e1b3b01705	2015-09-01	0.119916	0.102232	0.076257	40.72		