

In [279]: #check the customer with LOC_ID of 1000035946 in df2 and df3.
temp_df = pd.DataFrame(np.concatenate([df1[df1.LOC_ID == 1000035946], df2[df2.LOC_ID == 1000035946],
df3[df3.LOC_ID == 1000035946]],axis=0).T)

temp_df.columns = ['df1', 'df2','df3']
temp_df[['label']] = df2.columns
temp_df = temp_df[['label','df1', 'df2','df3']]
temp_df

	label	df1	df2	df3
0	LOC_ID	1000035946	1000035946	1000035946
1	RATE	RSFD	RSFD	RSFD
2	DWEL_UNIT	NaN	NaN	NaN
3	USE1	365	NaN	NaN
4	USE2	25	NaN	NaN
5	USE3	26	0	0
6	USE4	20	1	0
7	USE5	16	2	0
8	USE6	20	0	2
9	USE7	25	0	1
10	USE8	21	0	6
11	USE9	20	0	10
12	USE10	19	1	4
13	USE11	24	3	4
14	USE12	23	4	1
15	READ_DT1	NaT	NaT	NaT
16	READ_DT2	NaT	2019-07-10 00:00:00	2020-07-12 00:00:00
17	READ_DT3	NaT	2019-06-11 00:00:00	2020-06-10 00:00:00
18	READ_DT4	NaT	2019-05-10 00:00:00	2020-05-11 00:00:00
19	READ_DT5	NaT	2019-04-11 00:00:00	2020-04-12 00:00:00
20	READ_DT6	NaT	2019-03-11 00:00:00	2020-03-12 00:00:00
21	READ_DT7	NaT	2019-02-11 00:00:00	2020-02-10 00:00:00
22	READ_DT8	NaT	2019-01-10 00:00:00	2020-01-10 00:00:00
23	READ_DT9	NaT	2018-12-10 00:00:00	2019-12-10 00:00:00
24	READ_DT10	NaT	2018-11-06 00:00:00	2019-11-06 00:00:00
25	READ_DT11	NaT	2018-10-09 00:00:00	2019-10-09 00:00:00
26	READ_DT12	NaT	2018-09-10 00:00:00	2019-09-11 00:00:00
27	READ_DAYS1	NaN	1	1
28	READ_DAYS2	NaN	29	32
29	READ_DAYS3	NaN	32	30
30	READ_DAYS4	NaN	29	29
31	READ_DAYS5	NaN	29	31
32	READ_DAYS6	NaN	30	31
33	READ_DAYS7	NaN	32	31
34	READ_DAYS8	NaN	31	31
35	READ_DAYS9	NaN	34	34
36	READ_DAYS10	NaN	28	29
37	READ_DAYS11	NaN	29	27
38	READ_DAYS12	NaN	33	33

In [280]: #usage data for the customer with the LOC_ID of 1000035946 were recorded but without any READ Date records in year1.
#let's check the residential accounts before and after the customer with LOC_ID of 1000035946.
df1[(df1.LOC_ID == 1000035946) & (df1.DWEL_UNIT!=1)].T

	35946	35946	35947
LOC_ID	1000035946	1000035946	1000035947
RATE	RSFD	RSFD	RSFD
DWEL_UNIT	1	NaN	1
USE1	10	365	14
USE2	8	25	13
USE3	7	26	10
USE4	22	20	11
USE5	7	16	11
USE6	10	20	11
USE7	11	25	13
USE8	11	21	12
USE9	10	20	13
USE10	10	19	13
USE11	10	24	16
USE12	11	23	13
READ_DT1	2018-07-25 00:00:00	NaT	2018-07-30 00:00:00
READ_DT2	2018-06-24 00:00:00	NaT	2018-06-30 00:00:00
READ_DT3	2018-04-25 00:00:00	NaT	2018-04-30 00:00:00
READ_DT4	2018-03-27 00:00:00	NaT	2018-03-31 00:00:00
READ_DT5	2018-02-26 00:00:00	NaT	2018-03-01 00:00:00
READ_DT6	2018-01-25 00:00:00	NaT	2018-01-30 00:00:00
READ_DT7	2017-12-22 00:00:00	NaT	2017-12-28 00:00:00
READ_DT8	2017-11-21 00:00:00	NaT	2017-11-26 00:00:00
READ_DT9	2017-10-23 00:00:00	NaT	2017-10-28 00:00:00
READ_DT10	2017-09-22 00:00:00	NaT	2017-09-27 00:00:00
READ_DT11	2017-08-23 00:00:00	NaT	2017-08-28 00:00:00
READ_DAYS1	30	NaN	33
READ_DAYS2	32	NaN	28
READ_DAYS3	29	NaN	30
READ_DAYS4	29	NaN	31
READ_DAYS5	29	NaN	29
READ_DAYS6	32	NaN	30
READ_DAYS7	34	NaN	33
READ_DAYS8	31	NaN	30
READ_DAYS9	29	NaN	33
READ_DAYS10	31	NaN	29
READ_DAYS11	30	NaN	30
READ_DAYS12	29	NaN	31

In [281]: #use forward fillna approach, since it looks more look like its next account in terms of usage volume.
df1[(df1.LOC_ID == 1000035946, 'DWEL_UNIT') & 1

#the value of USE1 seems to be wrong (data entry error)* because it is too different from its previous usage.
df1[(df1.LOC_ID == 1000035946, 'USE1') = math.ceil(np.mean(df1.loc[df1.LOC_ID == 1000035946, 'USE2':'USE12'].values))
df1[df1.LOC_ID == 1000035946].T

	35946
LOC_ID	1000035946
RATE	RSFD
DWEL_UNIT	1
USE1	25
USE2	22
USE3	26
USE4	20
USE5	16
USE6	20
USE7	25
USE8	21
USE9	20
USE10	19
USE11	24
USE12	23
READ_DT1	NaT
READ_DT2	NaT
READ_DT3	NaT
READ_DT4	NaT
READ_DT5	NaT
READ_DT6	NaT
READ_DT7	NaT
READ_DT8	NaT
READ_DT9	NaT
READ_DT10	NaT
READ_DT11	NaT
READ_DT12	NaT
READ_DAYS1	NaN
READ_DAYS2	NaN
READ_DAYS3	NaN
READ_DAYS4	NaN
READ_DAYS5	NaN
READ_DAYS6	NaN
READ_DAYS7	NaN
READ_DAYS8	NaN
READ_DAYS9	NaN
READ_DAYS10	NaN
READ_DAYS11	NaN
READ_DAYS12	NaN

In [282]: #check the two missing values of dwelling unit among Duplex customers in df1.
df1[(df1.RATE == 'RDUPLEX') & (df1.DWEL_UNIT.isna())].T

	12743	65066
LOC_ID	1000012743	1000065066
RATE	RDUPLEX	RDUPLEX
DWEL_UNIT	NaN	NaN
USE1	NaN	NaN
USE2	NaN	NaN
USE3	NaN	NaN
USE4	NaN	NaN
USE5	NaN	NaN
USE6	NaN	NaN
USE7	NaN	NaN
USE8	NaN	NaN
USE9	NaN	NaN
USE10	NaN	NaN
USE12	NaN	NaN
READ_DT1	NaT	NaT
READ_DT2	NaT	NaT
READ_DT3	NaT	NaT
READ_DT4	NaT	NaT
READ_DT5	NaT	NaT
READ_DT6	NaT	NaT
READ_DT7	NaT	NaT
READ_DT8	NaT	NaT
READ_DT9	NaT	NaT
READ_DT10	NaT	NaT
READ_DT11	NaT	NaT
READ_DT12	NaT	NaT
READ_DAYS1	NaN	NaN
READ_DAYS2	NaN	NaN
READ_DAYS3	NaN	NaN
READ_DAYS4	NaN	NaN
READ_DAYS5	NaN	NaN
READ_DAYS6	NaN	NaN
READ_DAYS7	NaN	NaN
READ_DAYS8	NaN	NaN
READ_DAYS9	NaN	NaN
READ_DAYS10	NaN	NaN
READ_DAYS11	NaN	NaN
READ_DAYS12	NaN	NaN

In [283]: #check the customer with LOC_ID of 1000012743 in df2 and df3.
temp_df = pd.DataFrame(np.concatenate([df1[df1.LOC_ID == 1000012743], df2[df2.LOC_ID == 1000012743],
df3[df3.LOC_ID == 1000012743]],axis=0).T)

temp_df.columns = ['df1', 'df2','df3']
temp_df[['label']] = df2.columns
temp_df = temp_df[['label','df1', 'df2','df3']]
temp_df

	label	df1	df2	df3
0	LOC_ID	1000012743	1000012743	1000012743
1	RATE	RDUPLEX	RDUPLEX	RDUPLEX
2	DWEL_UNIT	NaN	NaN	NaN
3	USE1	NaN	NaN	NaN
4	USE2	NaN	NaN	NaN
5	USE3	NaN	NaN	NaN
6	USE4	NaN	NaN	NaN
7	USE5	NaN	NaN	NaN
8	USE6	NaN	NaN	NaN
9	USE7	NaN	NaN	NaN
10	USE8	NaN	NaN	NaN
11	USE9	NaN	NaN	NaN
12	USE10	NaN	NaN	NaN
13	USE11	NaN	NaN	NaN
14	USE12	NaN	NaN	NaN
15	READ_DT1	NaT	NaT	NaT
16	READ_DT2	NaT	NaT	NaT
17	READ_DT3	NaT	NaT	NaT
18	READ_DT4	NaT	NaT	NaT
19	READ_DT5	NaT	NaT	NaT
20	READ_DT6	NaT	NaT	NaT
21	READ_DT7	NaT	NaT	NaT
22	READ_DT8	NaT	NaT	NaT
23	READ_DT9	NaT	NaT	NaT
24	READ_DT10	NaT	NaT	NaT
25	READ_DT11	NaT	NaT	NaT
26	READ_DT12	NaT	NaT	NaT
27	READ_DAYS1	NaN	NaN	NaN
28	READ_DAYS2	NaN	NaN	NaN
29	READ_DAYS3	NaN	NaN	NaN
30	READ_DAYS4	NaN	NaN	NaN
31	READ_DAYS5	NaN	NaN	NaN
32	READ_DAYS6	NaN	NaN	NaN
33	READ_DAYS7	NaN	NaN	NaN
34	READ_DAYS8	NaN	NaN	NaN
35	READ_DAYS9	NaN	NaN	NaN
36	READ_DAYS10	NaN	NaN	NaN
37	READ_DAYS11	NaN	NaN	NaN
38	READ_DAYS12	NaN	NaN	NaN

In [284]: #check the customer with LOC_ID of 1000012743 in df2 and df3.
temp_df = pd.DataFrame(np.concatenate([df1[df1.LOC_ID == 1000065066], df2[df2.LOC_ID == 1000065066],
df3[df3.LOC_ID == 1000065066]],axis=0).T)

temp_df.columns = ['df1', 'df2','df3']
temp_df[['label']] = df2.columns
temp_df = temp_df[['label','df1', 'df2','df3']]
temp_df

	label	df1	df2	df3
0	LOC_ID	1000065066	1000065066	1000065066
1	RATE	RDUPLEX	RDUPLEX	RDUPLEX
2	DWEL_UNIT	NaN	NaN	NaN
3	USE1	NaN	NaN	NaN
4	USE2	NaN	NaN	NaN
5	USE3	NaN	NaN	NaN
6	USE4	NaN	NaN	NaN
7	USE5	NaN	NaN	NaN
8	USE6	NaN	NaN	NaN
9	USE7	NaN	NaN	NaN
10	USE8	NaN	NaN	NaN
11	USE9	NaN	NaN	NaN
12	USE10	NaN	NaN	NaN
13	USE11	NaN	NaN	NaN
14	USE12	NaN	NaN	NaN
15	READ_DT1	NaT	NaT	NaT
16	READ_DT2	NaT	NaT	NaT
17	READ_DT3	NaT	NaT	NaT
18	READ_DT4	NaT	NaT	NaT
19	READ_DT5	NaT	NaT	NaT
20	READ_DT6	NaT	NaT	NaT
21	READ_DT7	NaT	NaT	NaT
22	READ_DT8	NaT	NaT	NaT
23	READ_DT9	NaT	NaT	NaT
24	READ_DT10	NaT	NaT	NaT
25	READ_DT11	NaT	NaT	NaT
26	READ_DT12	NaT	NaT	NaT
27	READ_DAYS1	NaN	NaN	NaN
28	READ_DAYS2	NaN	NaN	NaN
29	READ_DAYS3	NaN	NaN	NaN
30	READ_DAYS4	NaN	NaN	NaN
31	READ_DAYS5	NaN	NaN	NaN
32	READ_DAYS6	NaN	NaN	NaN
33	READ_DAYS7	NaN	NaN	NaN
34	READ_DAYS8	NaN	NaN	NaN
35	READ_DAYS9	NaN	NaN	NaN
36	READ_DAYS10	NaN	NaN	NaN
37	READ_DAYS11	NaN	NaN	NaN
38	READ_DAYS12	NaN	NaN	NaN

There are no values for customer with LOC_ID of 1000012743 and 1000065066 over three years . So those rows are not informative and let's drop from all DataFrames.

In [285]: df1 = df1[(df1.LOC_ID != 1000012743) & (df1.LOC_ID != 1000065066)]
df2 = df2[(df2.LOC_ID != 1000035946, 'USE1') = df1.loc[df1.LOC_ID == 1000035946, 'DWEL_UNIT']].values
df3 = df3[(df3.LOC_ID != 1000012743) & (df3.LOC_ID != 1000065066)]
df1.shape

Out[285]: (78891, 39)

In [287]: #Calculate the market share of each Rate(customer type), the number of Rate based on the number of Dwelling unit
#and % of related missing values.
dict2 = (rate:count(df2.rate)) for rate in df2.RATE.unique()
customer = pd.DataFrame(dict2)
customer.columns = ['Single_family', 'Multi_family','Duplex','Commercial','Irrigation','Industrial','Commercial and Residual']
customer = customer.T
customer.columns = ['count', '% missing_value', '% missing_value']
print(customer)
customer

Out[287]:

	count	# missing_value	% missing_value
Single_family	53839.0	2.0	0.0
Multi_family	88047.0	1.0	0.0
Duplex	13460.0	0.0	0.0
Commercial	5269.0	0.0	0.0
Irrigation	1124.0	0.0	0.0
Industrial	194.0	0.0	0.0
Commercial and Residual	1580.0	0.0	0.0

In [288]: #check the two missing values of dwelling unit among single family in df2.
df2[(df2.RATE == 'RSFD') & (df2.DWEL_UNIT.isna())].T

	35946	63905
LOC_ID	1000035946	1000063905
RATE	RSFD	RSFD
DWEL_UNIT	NaN	NaN
USE1	NaN	NaN
USE2	NaN	NaN
USE3	0	NaN
USE4	1	NaN
USE5	2	NaN
USE6	0	NaN
USE7	0	NaN
USE8	0	NaN
USE9	0	NaN
USE10	1	NaN
USE11	3	NaN
USE12	4	NaN
READ_DT1	NaT	NaT
READ_DT2	2019-07-10 00:00:00	NaT
READ_DT3	2019-06-11 00:00:00	NaT
READ_DT4	2019-05-10 00:00:00	NaT
READ_DT5	2019-04-11 00:00:00	NaT
READ_DT6	2018-03-11 00:00:00	NaT
READ_DT7	2018-02-11 00:00:00	NaT
READ_DT8	2018-01-10 00:00:00	NaT
READ_DT9	2018-12-10 00:00:00	NaT
READ_DT10	2018-11-06 00:00:00	NaT
READ_DT11	2018-10-09 00:00:00	NaT
READ_DT12	2018-09-10 00:00:00	NaT
READ_DAYS1	1	NaN
READ_DAYS2	29	NaN
READ_DAYS3	32	NaN
READ_DAYS4	29	NaN
READ_DAYS5	29	NaN
READ_DAYS6	30	NaN
READ_DAYS7	32	NaN
READ_DAYS8	31	NaN
READ_DAYS9	34	NaN
READ_DAYS10	28	NaN
READ_DAYS11	29	NaN
READ_DAYS12	33	NaN

In [289]: #check the customer with LOC_ID of 1000035946 in df1
temp_df = pd.DataFrame(np.concatenate([df1[df1.LOC_ID == 1000035946], df2[df2.LOC_ID == 1000035946],
df3[df3.LOC_ID == 1000035946]],axis=0).T)

temp_df.columns = ['df1', 'df2','df3']
temp_df[['label']] = df2.columns
temp_df = temp_df[['label','df1', 'df2','df3']]
temp_df

21	READ_DT8	NaT	2019-01-10 00:00:00	2020-01-10 00:00:00
22	READ_DT8	NaT	2018-12-10 00:00:00	2019-12-10 00:00:00
23	READ_DT9	NaT	2018-12-10 00:00:00	2019-12-10 00:00:00
24	READ_DT10	NaT	2018-11-06 00:00:00	2019-11-06 00:00:00
25	READ_DT11	NaT	2018-10-09 00:00:00	2019-10-09 00:00:00
26	READ_DT12	NaT	2018-09-10 00:00:00	2019-09-11 00:00:00
27	READ_DAYS1	NaN	1	1
28	READ_DAYS2	NaN	29	32
29	READ_DAYS3	NaN	32	30
30	READ_DAYS4	NaN	29	29
31	READ_DAYS5	NaN	29	31
32	READ_DAYS6	NaN	30	31
33	READ_DAYS7	NaN	32	31
34	READ_DAYS8	NaN	31	31
35	READ_DAYS9	NaN	34	34
36	READ_DAYS10	NaN	28	29
37	READ_DAYS11	NaN	29	27
38	READ_DAYS12	NaN	33	33

```
In [290]: #Assign the value of dwelling unit in year1 to dwelling unit in year2 and year3
df2.loc[df2.LOC_ID == 1000035946, 'DWEL_UNIT'] = df1.loc[df1.LOC_ID == 1000035946, 'DWEL_UNIT'].values
df3.loc[df3.LOC_ID == 1000035946, 'DWEL_UNIT'] = df1.loc[df1.LOC_ID == 1000035946, 'DWEL_UNIT'].values

#Assign the mean value of usages in period 3 to 12 to use in period 1 and 2 in year2
df2.loc[df2.LOC_ID == 1000035946, 'USE1'] = np.mean(df2.loc[df2.LOC_ID == 1000035946, 'USE3':'USE12']).va
lues()
df2.loc[df2.LOC_ID == 1000035946, 'USE2'] = np.mean(df2.loc[df2.LOC_ID == 1000035946, 'USE3':'USE12']).va
lues()

#Assign the mean value of usages in period 3 to 12 to usge in period 1 and 2 in year3
df3.loc[df3.LOC_ID == 1000035946, 'USE1'] = np.mean(df3.loc[df3.LOC_ID == 1000035946, 'USE3':'USE12']).va
lues()
df3.loc[df3.LOC_ID == 1000035946, 'USE2'] = np.mean(df3.loc[df3.LOC_ID == 1000035946, 'USE3':'USE12']).va
lues()

In [291]: #Select the customer with LOC_ID of 1000035946 in df1
temp_df = pd.DataFrame(np.concatenate([df1[df1.LOC_ID == 1000035946], df2[df2.LOC_ID == 1000035946],
df3[df3.LOC_ID == 1000035946]]),axis=0).T)
temp_df.columns = ['df1', 'df2', 'df3']
temp_df['label'] = df2.columns
temp_df = temp_df[['label', 'df1', 'df2', 'df3']]
temp_df
```

	label	df1	df2	df3
0	LOC_ID	1000035946	1000035946	1000035946
1	RATE	RSFD	RSFD	RSFD
2	DWEL_UNIT	1	NaN	NaN
3	USE1	22	NaN	NaN
4	USE2	25	NaN	NaN
5	USE3	26	0	0
6	USE4	27	0	0
7	USE5	28	0	0
8	USE6	29	0	0
9	USE7	30	0	0
10	USE8	31	0	0
11	USE9	32	0	0
12	USE10	33	0	0
13	USE11	34	0	0
14	USE12	35	0	0


```
In [297]: #check the custom with LOC_ID of 1000052668 in df2,df3
temp_df = pd.DataFrame(np.concatenate([df1[df1.LOC_ID == 1000052668], df2[df2.LOC_ID == 1000052668],
                                         df3[df3.LOC_ID == 1000052668]],axis=0).T)
temp_df.columns = ['df1', 'df2','df3']
temp_df['label'] = df3.columns
temp_df = temp_df[['label','df1', 'df2','df3']]
temp_df
```

	label	df1	df2	df3
0	LOC_ID	1000052668	1000052668	1000052668
1	RATE	RMSF	RMSF	RMSF
2	DWEL_UNIT	3	NaN	NaN
3	USE1	19	NaN	NaN
4	USE2	19	NaN	NaN
5	USE3	20	NaN	NaN
6	USE4	19	NaN	NaN
7	USE5	17	NaN	NaN
8	USE6	18	NaN	NaN
9	USE7	17	NaN	NaN
10	USE8	17	NaN	NaN
11	USE9	18	NaN	NaN
12	USE10	17	NaN	NaN
13	USE11	18	NaN	NaN
14	USE12	18	NaN	NaN
15	READ_DT1	2018-08-02 00:00:00	NaT	NaT
16	READ_DT2	2018-07-03 00:00:00	NaT	NaT
17	READ_DT3	2018-06-04 00:00:00	NaT	NaT
18	READ_DT4	2018-05-03 00:00:00	NaT	NaT
19	READ_DT5	2018-04-04 00:00:00	NaT	NaT
20	READ_DT6	2018-03-06 00:00:00	NaT	NaT
21	READ_DT7	2018-02-02 00:00:00	NaT	NaT
22	READ_DT8	2018-01-03 00:00:00	NaT	NaT
23	READ_DT9	2017-12-04 00:00:00	NaT	NaT
24	READ_DT10	2017-10-31 00:00:00	NaT	NaT
25	READ_DT11	2017-10-02 00:00:00	NaT	NaT
26	READ_DT12	2017-08-31 00:00:00	NaT	NaT
27	READ_DAYS1	30	NaN	NaN
28	READ_DAYS2	29	NaN	NaN
29	READ_DAYS3	32	NaN	NaN
30	READ_DAYS4	29	NaN	NaN
31	READ_DAYS5	29	NaN	NaN
32	READ_DAYS6	32	NaN	NaN
33	READ_DAYS7	30	NaN	NaN
34	READ_DAYS8	30	NaN	NaN
35	READ_DAYS9	34	NaN	NaN
36	READ_DAYS10	29	NaN	NaN
37	READ_DAYS11	32	NaN	NaN
38	READ_DAYS12	29	NaN	NaN

There is no data related to customer with LOC_ID of 1000052668 in year2 and year3. Let's drop from all dataframes.

```
In [298]: df1 = df1[df1.LOC_ID != 1000052668]
df2 = df2[df2.LOC_ID != 1000052668]
df3 = df3[df3.LOC_ID != 1000052668]
```

```
df3

In [299]: #Calculate the market share of each Rate(customer type), the number of Rate based on the number of Dwe
liling unit
#and % of related missing values.
dict3 = (rate:count(df3.rate) for rate in df3.RATE.unique())
customer = pd.DataFrame(dict3)
customer.columns = ['Singel_family', 'Multi_family','Duplex','Commercial','Irrigation','Industrial','C
ommercial and Residential']
customer = customer.T
customer.columns = ['count', '% missing_value', '% missing_value']
print(df3)
customer
```

	count	% missing_value	% missing_value
Singel_family	5374.0	3.0	0.0
Multi_family	88126.0	0.0	0.0
Duplex	13597.0	0.0	0.0
Commercial	5267.0	0.0	0.0
Irrigation	1125.0	0.0	0.0
Industrial	195.0	0.0	0.0
Commercial and Residential	1308.0	0.0	0.0

```
In [300]: #check the two missing values of dwelling unit among single family in df2.
df3[df3.RATE == "RSFD"] & (df3.DWEL_UNIT.isna())>.T
```

	24304	54741	62574
LOC_ID	1000024304	1000054741	1000062574
RATE	RSFD	RSFD	RSFD
DWEL_UNIT	NaN	NaN	NaN
USE1	NaN	NaN	NaN
USE2	NaN	NaN	NaN
USE3	NaN	NaN	NaN
USE4	NaN	NaN	NaN
USE5	NaN	NaN	NaN
USE6	NaN	NaN	NaN
USE7	NaN	NaN	NaN
USE8	NaN	NaN	NaN
USE9	NaN	NaN	NaN
USE10	NaN	NaN	NaN
USE11	NaN	NaN	NaN
USE12	NaN	NaN	NaN
READ_DT1	NaT	NaT	NaT
READ_DT2	NaT	NaT	NaT
READ_DT3	NaT	NaT	NaT
READ_DT4	NaT	NaT	NaT
READ_DT5	NaT	NaT	NaT
READ_DT6	NaT	NaT	NaT
READ_DT7	NaT	NaT	NaT
READ_DT8	NaT	NaT	NaT
READ_DT9	NaT	NaT	NaT
READ_DT10	NaT	NaT	NaT
READ_DT11	NaT	NaT	NaT
READ_DT12	NaT	NaT	NaT
READ_DAYS1	NaN	NaN	NaN
READ_DAYS2	NaN	NaN	NaN
READ_DAYS3	NaN	NaN	NaN
READ_DAYS4	NaN	NaN	NaN
READ_DAYS5	NaN	NaN	NaN
READ_DAYS6	NaN	NaN	NaN
READ_DAYS7	NaN	NaN	NaN
READ_DAYS8	NaN	NaN	NaN
READ_DAYS9	NaN	NaN	NaN
READ_DAYS10	NaN	NaN	NaN
READ_DAYS11	NaN	NaN	NaN
READ_DAYS12	NaN	NaN	NaN

```
In [301]: #check the customer with LOC_ID of 1000024304 in df1,df2
temp_df = pd.DataFrame(np.concatenate([df1[df1.LOC_ID == 1000024304], df2[df2.LOC_ID == 1000024304],
                                         df3[df3.LOC_ID == 1000024304]],axis=0).T)
temp_df.columns = ['df1', 'df2','df3']
temp_df['label'] = df3.columns
temp_df = temp_df[['label','df1', 'df2','df3']]
temp_df
```

	label	df1	df2	df3
0	LOC_ID	1000024304	1000024304	1000024304
1	RATE	RSFD	RSFD	RSFD
2	DWEL_UNIT	1	1	NaN
3	USE1	18	8	NaN
4	USE2	25	11	NaN
5	USE3	17	8	NaN
6	USE4	13	9	NaN
7	USE5	16	34	NaN
8	USE6	17	12	NaN
9	USE7	14	16	NaN
10	USE8	12	17	NaN
11	USE9	18	16	NaN
12	USE10	NaN	17	NaN
13	USE11	NaN	15	NaN
14	USE12	NaN	14	NaN
15	READ_DT1	2018-07-17 00:00:00	2019-07-16 00:00:00	NaT
16	READ_DT2	2018-06-15 00:00:00	2019-06-17 00:00:00	NaT
17	READ_DT3	2018-06-16 00:00:00	2019-06-15 00:00:00	NaT
18	READ_DT4	2018-04-17 00:00:00	2019-06-16 00:00:00	NaT
19	READ_DT5	2018-03-19 00:00:00	2019-03-18 00:00:00	NaT
20	READ_DT6	2018-02-15 00:00:00	2019-02-15 00:00:00	NaT
21	READ_DT7	2018-01-17 00:00:00	2019-01-16 00:00:00	NaT
22	READ_DT8	2017-12-14 00:00:00	2018-12-14 00:00:00	NaT
23	READ_DT9	2017-11-13 00:00:00	2018-11-13 00:00:00	NaT
24	READ_DT10	NaT	2018-10-15 00:00:00	NaT
25	READ_DT11	NaT	2018-09-14 00:00:00	NaT
26	READ_DT12	NaT	2018-08-15 00:00:00	NaT
27	READ_DAYS1	32	29	NaN
28	READ_DAYS2	30	33	NaN
29	READ_DAYS3	29	29	NaN
30	READ_DAYS4	29	29	NaN
31	READ_DAYS5	32	31	NaN
32	READ_DAYS6	29	30	NaN
33	READ_DAYS7	34	33	NaN
34	READ_DAYS8	31	31	NaN
35	READ_DAYS9	38	29	NaN
36	READ_DAYS10	NaN	31	NaN
37	READ_DAYS11	NaN	30	NaN
38	READ_DAYS12	NaN	29	NaN

```
In [302]: #check the customer with LOC_ID of 1000062574 in df1,df2
temp_df = pd.DataFrame(np.concatenate([df1[df1.LOC_ID == 1000062574], df2[df2.LOC_ID == 1000062574],
                                         df3[df3.LOC_ID == 1000062574]],axis=0).T)
temp_df.columns = ['df1', 'df2','df3']
temp_df['label'] = df3.columns
temp_df = temp_df[['label','df1', 'df2','df3']]
temp_df
```

	label	df1	df2	df3
0	LOC_ID	1000062574	1000062574	1000062574
1	RATE	RSFD	RSFD	RSFD
2	DWEL_UNIT	1	1	NaN
3	USE1	24	9	NaN
4	USE2	21	7	NaN
5	USE3	8	15	NaN
6	USE4	7	7	NaN
7	USE5	5	5	NaN
8	USE6	14	4	NaN
9	USE7	10	8	NaN
10	USE8	15	6	NaN
11	USE9	8	8	NaN
12	USE10	7	8	NaN
13	USE11	12	7	NaN
14	USE12	8	12	NaN
15	READ_DT1	2018-07-17 00:00:00	2019-07-17 00:00:00	NaT
16	READ_DT2	2018-06-15 00:00:00	2019-06-17 00:00:00	NaT
17	READ_DT3	2018-06-16 00:00:00	2019-06-16 00:00:00	NaT
18	READ_DT4	2018-04-17 00:00:00	2019-04-17 00:00:00	NaT
19	READ_DT5	2018-03-19 00:00:00	2019-03-19 00:00:00	NaT
20	READ_DT6	2018-02-15 00:00:00	2019-02-15 00:00:00	NaT
21	READ_DT7	2018-01-17 00:00:00	2019-01-16 00:00:00	NaT
22	READ_DT8	2017-12-14 00:00:00	2018-12-14 00:00:00	NaT
23	READ_DT9	2017-11-13 00:00:00	2018-11-13 00:00:00	NaT
24	READ_DT10	2017-10-13 00:00:00	2018-10-15 00:00:00	NaT
25	READ_DT11	2017-09-14 00:00:00	2018-09-14 00:00:00	NaT
26	READ_DT12	2017-08-15 00:00:00	2018-08-15 00:00:00	NaT
27	READ_DAYS1	32	30	NaN
28	READ_DAYS2	30	32	NaN
29	READ_DAYS3	29	29	NaN
30	READ_DAYS4	29	29	NaN
31	READ_DAYS5	32	32	NaN
32	READ_DAYS6	29	30	NaN
33	READ_DAYS7	34	33	NaN
34	READ_DAYS8	31	31	NaN
35	READ_DAYS9	31	29	NaN
36	READ_DAYS10	29	31	NaN
37	READ_DAYS11	30	30	NaN
38	READ_DAYS12	29	29	NaN

There is no data related to customers with LOC_ID of 1000024304, 1000054741 and 1000062574 in year3. Let's drop from all dataframes.

```
In [304]: df1 = df1[(df1.LOC_ID != 1000024304) & (df1.LOC_ID != 1000054741) & (df1.LOC_ID != 1000062574)]
df2 = df2[(df2.LOC_ID != 1000024304) & (df2.LOC_ID != 1000054741) & (df2.LOC_ID != 1000062574)]
df3 = df3[(df3.LOC_ID != 1000024304) & (df3.LOC_ID != 1000054741) & (df3.LOC_ID != 1000062574)]
```

```
In [305]: #Calculate the market share of each Rate in first year, the number of Rate based on the number of Dwe
liling unit
#and % of related missing values.
dict3 = (rate:count(df3.rate) for rate in df3.RATE.unique())
customer1 = pd.DataFrame(dict1)
customer1.columns = ['Singel_family', 'Multi_family','Duplex','Commercial','Irrigation','Industrial',
'Commercial and Residential']
customer1 = customer1.T
customer1.columns = ['count', '% missing_value', '% missing_value']
customer1['market_share'] = list(round(100* df3.RATE.value_counts() / len(df1),2))

customer1 = customer1[['market_share','count', '% missing_value', '% missing_value']]
printmd("\nspan style='color:Blue'>>>Information about RATE(customer types) over the Year 2017-2018**
</span>")
customer1
```

Information about RATE(customer types) over the Year 2017-2018**

	market_share	count	# missing_value	% missing_value
Singel_family	68.24	53842.0	0.0	0.0
Multi_family	13.90	88061.0	0.0	0.0
Duplex	9.31	13461.0	0.0	0.0
Commercial	6.68	5268.0	0.0	0.0
Irrigation	1.42	1122.0	0.0	0.0
Industrial	0.32	197.0	0.0	0.0
Commercial and Residential	0.25	1578.0	0.0	0.0

```
In [306]: #Calculate the market share of each Rate in second year, the number of Rate based on the number of Dwe
liling unit
#and % of related missing values.
dict2 = (rate:count(df2.rate) for rate in df2.RATE.unique())
customer2 = pd.DataFrame(dict1)
customer2.columns = ['Singel_family', 'Multi_family','Duplex','Commercial','Irrigation','Industrial',
'Commercial and Residential']
customer2 = customer2.T
customer2.columns = ['count', '% missing_value', '% missing_value']
customer2['market_share'] = list(round(100* df2.RATE.value_counts() / len(df2),2))

customer2 = customer2[['market_share','count', '% missing_value', '% missing_value']]
printmd("\nspan style='color:Blue'>>>Information about RATE(customer types) over the Year 2018-2019**
</span>")
customer2
```

Information about RATE(customer types) over the Year 2018-2019**

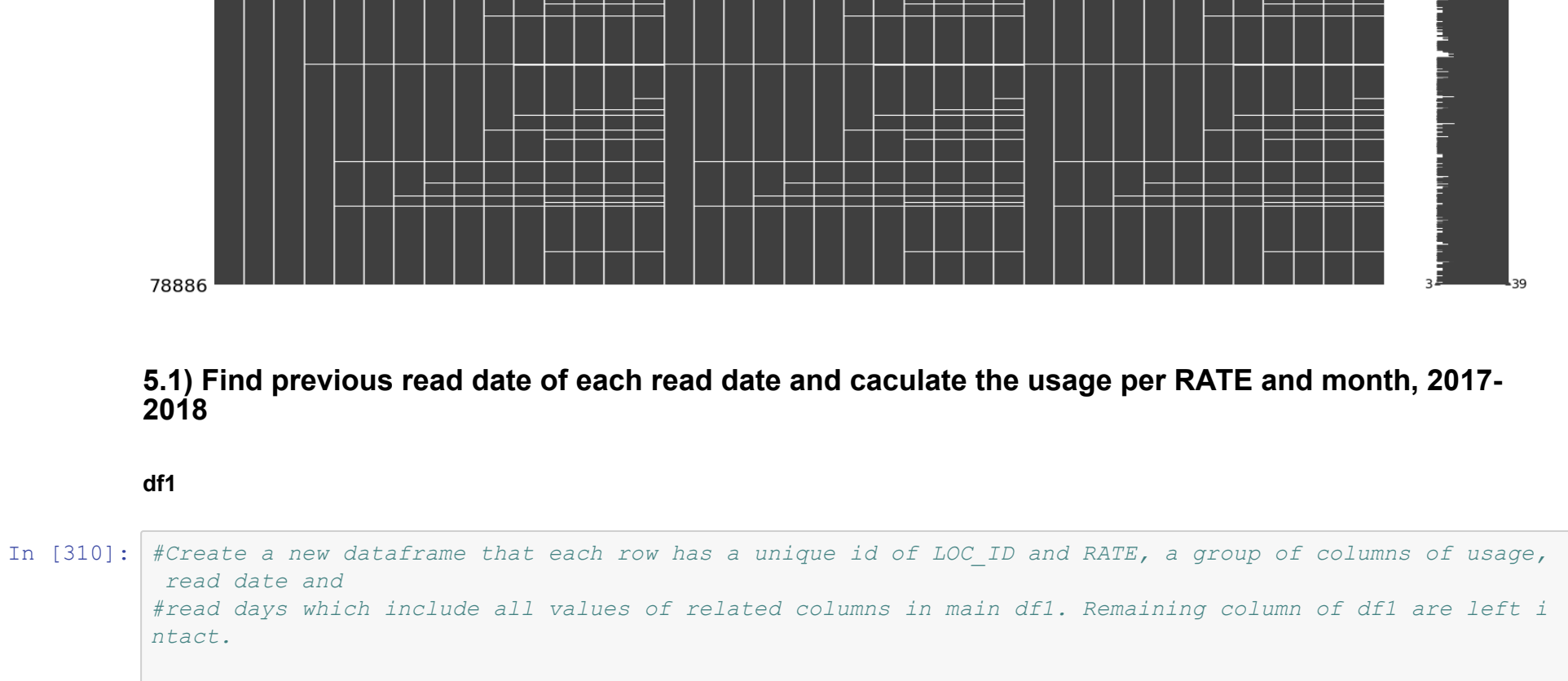
	market_share	count	# missing_value	% missing_value
Singel_family	68.23	53842.0	0.0	0.0
Multi_family	13.90	88061.0	0.0	0.0
Duplex	9.20	13461.0	0.0	0.0
Commercial	6.68	5268.0	0.0	0.0
Irrigation	1.42	1122.0	0.0	0.0
Industrial	0.33	197.0	0.0	0.0
Commercial and Residential	0.25	1578.0	0.0	0.0

```
In [307]: #Calculate the market share of each Rate in second year, the number of Rate based on the number of Dwe
liling unit
#and % of related missing values.
dict3 = (rate:count(df3.rate) for rate in df3.RATE.unique())
customer3 = pd.DataFrame(dict1)
customer3.columns = ['Singel_family', 'Multi_family','Duplex','Commercial','Irrigation','Industrial',
'Commercial and Residential']
customer3 = customer3.T
customer3.columns = ['count', '% missing_value', '% missing_value']
customer3['market_share'] = list(round(100* df3.RATE.value_counts() / len(df3),2))

customer3 = customer3[['market_share','count', '% missing_value', '% missing_value']]
printmd("\nspan style='color:Blue'>>>Information about RATE(customer types) over the Year 2019-2020**
</span>")
customer3
```

Information about RATE(customer types) over the Year 2019-2020**

	market_share	count	# missing_value	% missing_value
Singel_family	68.11	53842.0	0.0	0.0
Multi_family	13.90	88061.0	0.0	0.0
Duplex	9.31	13461.0	0.0	0.0
Commercial	6.68	5268.0	0.0	0.0
Irrigation	1.43	1122.0	0.0	0.0
Industrial	0.33	197.0	0.0	0.0
Commercial and Residential	0.25	1578.0	0.0	0.0



Now,there is no missing value in Dwelling unit column of dataframes.

And about market share:

1) There is very small changes in the market share of Singel_family and Duplex.

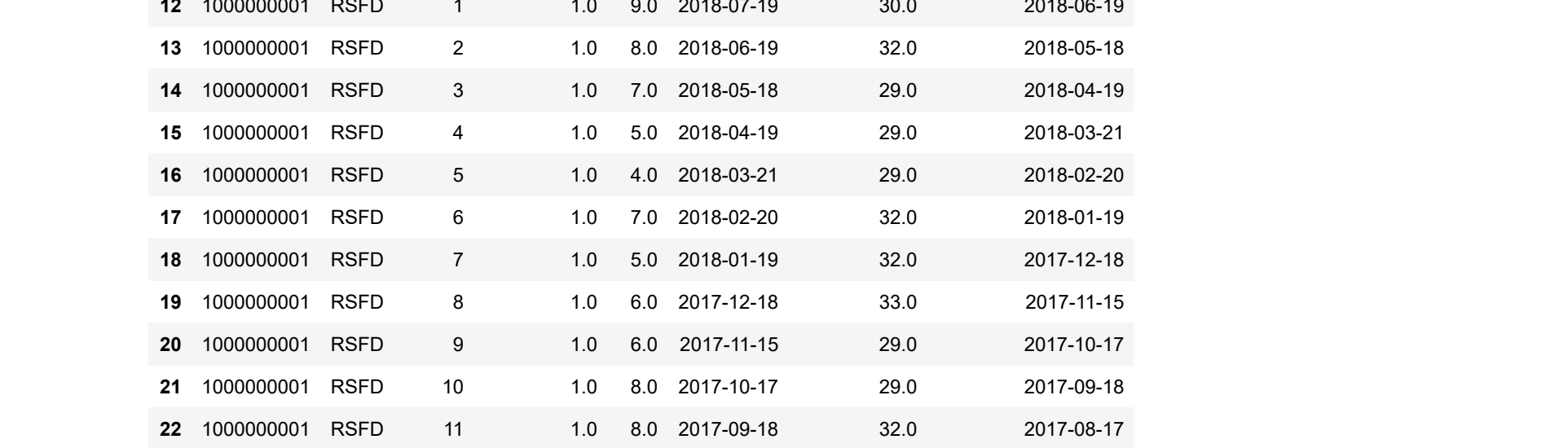
2) The majority of customers (more than 91%) are Residential including single family dwelling, multi-family and duplex houses.

3) There are 6.68% commercial customers. Only 2% of customers used water for commercial and residential, irrigation and industrial usage.

5) Visualizing missing values, checking validity of numeric variables

```
df1

In [309]: #Find matrix of missing values and plot it
mono.matrix(df1)
plt.title('Matrix of Missing Values in df1', fontsize=30, color='red')
plt.show()
```



```
5.1) Find previous read date of each read date and caculate the usage per RATE and month, 2017-
2018

In [310]: #Create a new dataframe that each row has a unique id of LOC_ID and RATE, a group of columns of usage,
read date and
#and days which include all values of related columns in main df1. Remaining column of df1 are left i
n read.

long_df1=pd.concat(long_df1,columns=['USE','READ_DT','READ_DAYS'],
                    i= ['LOC_ID','RATE'],
                    j='period',reset_index())
#long_df1['previous_read_date'] = long_df1['READ_DT'] - pd.to_timedelta(long_df1['READ_DT'], unit='D')
long_df1.head()
```

3						NaT			
34	1000000002	RSFD	11		1.0	NaH	NaT	NaH	NaT

In [317]: `#Dropping the rows that their Read date and previous read date are null
missing_read_date = long_dfl[long_dfl["READ_DT"].isna()]
long_dfl = long_dfl[~(long_dfl["READ_DT"].isna())]
long_dfl.head()`

Out[317]:

	LOC_ID	RATE	period	DWEL_UNIT	USE	READ_DT	READ_DAYS	Previous_READ_DT
0	1000000000	RSFD	1	1.0	2.0	2018-07-11	30.0	2018-06-11
1	1000000000	RSFD	2	1.0	2.0	2018-06-11	32.0	2018-06-10
2	1000000000	RSFD	3	1.0	15.0	2018-05-10	29.0	2018-04-11
3	1000000000	RSFD	4	1.0	9.0	2018-04-11	29.0	2018-03-13

