

In [42]: #Filter the null value of dwelling unit for residential duplex group in df1
df1[df1.RATE == "RDUPLEX"] & (df1.DWEL_UNIT.isna()) : T

Out [42]:

	LOC_ID	1200012743	1000065066
	RATE	RDUPLEX	RDUPLEX
DWEL_UNIT		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	NaN	NaN	NaN
READ_DT2	NaN	NaN	NaN
READ_DT3	NaN	NaN	NaN
READ_DT4	NaN	NaN	NaN
READ_DT5	NaN	NaN	NaN
READ_DT6	NaN	NaN	NaN
READ_DT7	NaN	NaN	NaN
READ_DT8	NaN	NaN	NaN
READ_DT9	NaN	NaN	NaN
READ_DT10	NaN	NaN	NaN
READ_DT11	NaN	NaN	NaN
READ_DT12	NaN	NaN	NaN
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 [43]: #check the customer with LOC_ID of 1000012743 and 1000065066 in df2 and df3.
temp = pd.DataFrame(np.concatenate([df1[df1.LOC_ID == 1000012743], df2[df2.LOC_ID == 1000065066],
df3[df3.LOC_ID == 1000012743], df3[df3.LOC_ID == 1000065066],df1[df1.LOC_ID == 1000012743],
df2[df2.LOC_ID == 1000012743], df3[df3.LOC_ID == 1000065066],df3[df3.LOC_ID == 1000012743]],axis=0).T)
temp.columns = ['df1', 'df2', 'df3', 'df1', 'df2', 'df3']
temp['label'] = df2.columns
temp = temp[['label','df1', 'df2', 'df3']]
temp

Out [43]:

	label	df1	df1	df2	df2	df3	df3
0	LOC_ID	1000012743	1000065066	1000012743	1000065066	1000012743	1000012743
1	RATE	RDUPLEX	RDUPLEX	RDUPLEX	RDUPLEX	RDUPLEX	RDUPLEX
2	DWEL_UNIT	NaN	NaN	NaN	NaN	NaN	NaN
3	USE1	NaN	NaN	NaN	NaN	NaN	NaN
4	USE2	NaN	NaN	NaN	NaN	NaN	NaN
5	USE3	NaN	NaN	NaN	NaN	NaN	NaN
6	USE4	NaN	NaN	NaN	NaN	NaN	NaN
7	USE5	NaN	NaN	NaN	NaN	NaN	NaN
8	USE6	NaN	NaN	NaN	NaN	NaN	NaN
9	USE7	NaN	NaN	NaN	NaN	NaN	NaN
10	USE8	NaN	NaN	NaN	NaN	NaN	NaN
11	USE9	NaN	NaN	NaN	NaN	NaN	NaN
12	USE10	NaN	NaN	NaN	NaN	NaN	NaN
13	USE11	NaN	NaN	NaN	NaN	NaN	NaN
14	USE12	NaN	NaN	NaN	NaN	NaN	NaN
15	READ_DT1	NaN	NaN	NaN	NaN	NaN	NaN
16	READ_DT2	NaN	NaN	NaN	NaN	NaN	NaN
17	READ_DT3	NaN	NaN	NaN	NaN	NaN	NaN
18	READ_DT4	NaN	NaN	NaN	NaN	NaN	NaN
19	READ_DT5	NaN	NaN	NaN	NaN	NaN	NaN
20	READ_DT6	NaN	NaN	NaN	NaN	NaN	NaN
21	READ_DT7	NaN	NaN	NaN	NaN	NaN	NaN
22	READ_DT8	NaN	NaN	NaN	NaN	NaN	NaN
23	READ_DT9	NaN	NaN	NaN	NaN	NaN	NaN
24	READ_DT10	NaN	NaN	NaN	NaN	NaN	NaN
25	READ_DT11	NaN	NaN	NaN	NaN	NaN	NaN
26	READ_DT12	NaN	NaN	NaN	NaN	NaN	NaN
27	READ_DAYS1	NaN	NaN	NaN	NaN	NaN	NaN
28	READ_DAYS2	NaN	NaN	NaN	NaN	NaN	NaN
29	READ_DAYS3	NaN	NaN	NaN	NaN	NaN	NaN
30	READ_DAYS4	NaN	NaN	NaN	NaN	NaN	NaN
31	READ_DAYS5	NaN	NaN	NaN	NaN	NaN	NaN
32	READ_DAYS6	NaN	NaN	NaN	NaN	NaN	NaN
33	READ_DAYS7	NaN	NaN	NaN	NaN	NaN	NaN
34	READ_DAYS8	NaN	NaN	NaN	NaN	NaN	NaN
35	READ_DAYS9	NaN	NaN	NaN	NaN	NaN	NaN
36	READ_DAYS10	NaN	NaN	NaN	NaN	NaN	NaN
37	READ_DAYS11	NaN	NaN	NaN	NaN	NaN	NaN
38	READ_DAYS12	NaN	NaN	NaN	NaN	NaN	NaN

These two account of duplex group don't have any data in any datframes, so let's drop them.

In [44]: #Filter the null value of dwelling unit for residential duplex group in df2
df2[(df2.RATE == "RDUPLEX") & (df2.DWEL_UNIT.isna())]

Out [44]:

	LOC_ID	RATE	DWEL_UNIT	USE1	USE2	USE3	USE4	USE5	USE6	USE7	...	READ_DAYS3	READ_DAYS4	READ_DAYS5
12743	1000012743	RDUPLEX	NaN	NaN	NaN	NaN	NaN	NaN	NaN	NaN	...	NaN	NaN	NaN
65066	1000065066	RDUPLEX	NaN	NaN	NaN	NaN	NaN	NaN	NaN	NaN	...	NaN	NaN	NaN

2 rows × 39 columns

In [45]: #check the customer with LOC_ID of 1000012743 and 1000065066 in df1 and df3.
temp = pd.DataFrame(np.concatenate([df1[df1.LOC_ID == 1000012743], df2[df2.LOC_ID == 1000012743],
df3[df3.LOC_ID == 1000012743], df1[df1.LOC_ID == 1000065066], df2[df2.LOC_ID == 1000065066],
df3[df3.LOC_ID == 1000065066],df1[df1.LOC_ID == 1000065066],df2[df2.LOC_ID == 1000065066],axis=0).T)
temp.columns = ['df1', 'df2', 'df3', 'df1', 'df2', 'df3']
temp['label'] = df2.columns
temp = temp[['label','df1', 'df2', 'df3']]
temp

Out [45]:

	label	df1	df1	df2	df2	df3	df3
0	LOC_ID	1000012743	1000065066	1000012743	1000065066	1000012743	1000065066
1	RATE	RDUPLEX	RDUPLEX	RDUPLEX	RDUPLEX	RDUPLEX	RDUPLEX
2	DWEL_UNIT	NaN	NaN	NaN	NaN	NaN	NaN
3	USE1	NaN	NaN	NaN	NaN	NaN	NaN
4	USE2	NaN	NaN	NaN	NaN	NaN	NaN
5	USE3	NaN	NaN	NaN	NaN	NaN	NaN
6	USE4	NaN	NaN	NaN	NaN	NaN	NaN
7	USE5	NaN	NaN	NaN	NaN	NaN	NaN
8	USE6	NaN	NaN	NaN	NaN	NaN	NaN
9	USE7	NaN	NaN	NaN	NaN	NaN	NaN
10	USE8	NaN	NaN	NaN	NaN	NaN	NaN
11	USE9	NaN	NaN	NaN	NaN	NaN	NaN
12	USE10	NaN	NaN	NaN	NaN	NaN	NaN
13	USE11	NaN	NaN	NaN	NaN	NaN	NaN
14	USE12	NaN	NaN	NaN	NaN	NaN	NaN
15	READ_DT1	NaN	NaN	NaN	NaN	NaN	NaN
16	READ_DT2	NaN	NaN	NaN	NaN	NaN	NaN
17	READ_DT3	NaN	NaN	NaN	NaN	NaN	NaN
18	READ_DT4	NaN	NaN	NaN	NaN	NaN	NaN
19	READ_DT5	NaN	NaN	NaN	NaN	NaN	NaN
20	READ_DT6	NaN	NaN	NaN	NaN	NaN	NaN
21	READ_DT7	NaN	NaN	NaN	NaN	NaN	NaN
22	READ_DT8	NaN	NaN	NaN	NaN	NaN	NaN
23	READ_DT9	NaN	NaN	NaN	NaN	NaN	NaN
24	READ_DT10	NaN	NaN	NaN	NaN	NaN	NaN
25	READ_DT11	NaN	NaN	NaN	NaN	NaN	NaN
26	READ_DT12	NaN	NaN	NaN	NaN	NaN	NaN
27	READ_DAYS1	NaN	NaN	NaN	NaN	NaN	NaN
28	READ_DAYS2	NaN	NaN	NaN	NaN	NaN	NaN
29	READ_DAYS3	NaN	NaN	NaN	NaN	NaN	NaN
30	READ_DAYS4	NaN	NaN	NaN	NaN	NaN	NaN
31	READ_DAYS5	NaN	NaN	NaN	NaN	NaN	NaN
32	READ_DAYS6	NaN	NaN	NaN	NaN	NaN	NaN
33	READ_DAYS7	NaN	NaN	NaN	NaN	NaN	NaN
34	READ_DAYS8	NaN	NaN	NaN	NaN	NaN	NaN
35	READ_DAYS9	NaN	NaN	NaN	NaN	NaN	NaN
36	READ_DAYS10	NaN	NaN	NaN	NaN	NaN	NaN
37	READ_DAYS11	NaN	NaN	NaN	NaN	NaN	NaN
38	READ_DAYS12	NaN	NaN	NaN	NaN	NaN	NaN

These two account of duplex group don't have any data in any datframes, so let's drop them.

In [46]: df3[(df3.RATE == "RDUPLEX") & (df3.DWEL_UNIT.isna())]

Out [46]:

	LOC_ID	RATE	DWEL_UNIT	USE1	USE2	USE3	USE4	USE5	USE6	USE7	...	READ_DAYS3	READ_DAYS4	READ_DAYS5
12743	1000012743	RDUPLEX	NaN	NaN	NaN	NaN	NaN	NaN	NaN	NaN	...	NaN	NaN	NaN
63905	1000063905	RDUPLEX	NaN	NaN	NaN	NaN	NaN	NaN	NaN	NaN	...	NaN	NaN	NaN
65066	1000065066	RDUPLEX	NaN	NaN	NaN	NaN	NaN	NaN	NaN	NaN	...	NaN	NaN	NaN

3 rows × 39 columns

In [47]: #check the customer with LOC_ID of 1000012743,1000063905 and 1000065066 in df1 and df2.
temp = pd.DataFrame(np.concatenate([df1[df1.LOC_ID == 1000012743], df2[df2.LOC_ID == 1000012743],
df3[df3.LOC_ID == 1000012743], df1[df1.LOC_ID == 1000063905], df2[df2.LOC_ID == 1000063905],
df3[df3.LOC_ID == 1000063905],df1[df1.LOC_ID == 1000063905],df2[df2.LOC_ID == 1000063905],axis=0).T)
temp.columns = ['df1', 'df2', 'df3', 'df1', 'df2', 'df3']
temp['label'] = df2.columns
temp = temp[['label','df1', 'df2', 'df3']]
temp

Out [47]:

	label	df1	df1	df2	df2	df2	df3	df3	df3
0	LOC_ID	1000012743	1000065066	1000012743	1000063905	1000065066	1000012743	1000063905	1000065066
1	RATE	RDUPLEX	RDUPLEX	RDUPLEX	RDUPLEX	RDUPLEX	RDUPLEX	RDUPLEX	RDUPLEX
2	DWEL_UNIT	NaN	NaN	NaN	NaN	NaN	NaN	NaN	NaN
3	USE1	NaN	11	NaN	NaN	NaN	NaN	NaN	NaN
4	USE2	NaN	12	NaN	NaN	NaN	NaN	NaN	NaN
5	USE3	NaN	10	NaN	NaN	NaN	NaN	NaN	NaN
6	USE4	NaN	9	NaN	NaN	NaN	NaN	NaN	NaN
7	USE5	NaN	9	NaN	NaN	NaN	NaN	NaN	NaN
8	USE6	NaN	8	NaN	NaN	NaN	NaN	NaN	NaN
9	USE7	NaN	10	NaN	NaN	NaN	NaN	NaN	NaN
10	USE8	NaN	8	NaN	NaN	NaN	NaN	NaN	NaN
11	USE9	NaN	10	NaN	NaN	NaN	NaN	NaN	NaN
12	USE10	NaN	8	NaN	NaN	NaN	NaN	NaN	NaN
13	USE11	NaN	9	NaN	NaN	NaN	NaN	NaN	NaN
14	USE12	NaN	8	NaN	NaN	NaN	NaN	NaN	NaN
15	READ_DT1	NaN	2018-07-30 00:00:00	NaN	NaN	NaN	NaN	NaN	NaN
16	READ_DT2	NaN	2018-06-27 00:00:00	NaN	NaN	NaN	NaN	NaN	NaN
17	READ_DT3	NaN	2018-05-30 00:00:00	NaN	NaN	NaN	NaN	NaN	NaN
18	READ_DT4	NaN	2018-04-30 00:00:00	NaN	NaN	NaN	NaN	NaN	NaN
19	READ_DT5	NaN	2018-03-30 00:00:00	NaN	NaN	NaN	NaN	NaN	NaN
20	READ_DT6	NaN	2018-03-01 00:00:00	NaN	NaN	NaN	NaN	NaN	NaN
21	READ_DT7	NaN	2018-01-30 00:00:00	NaN	NaN	NaN	NaN	NaN	NaN
22	READ_DT8	NaN	2017-12-28 00:00:00	NaN	NaN	NaN	NaN	NaN	NaN
23	READ_DT9	NaN	2017-11-28 00:00:00	NaN	NaN	NaN	NaN	NaN	NaN
24	READ_DT10	NaN	2017-10-26 00:00:00	NaN	NaN	NaN	NaN	NaN	NaN
25	READ_DT11	NaN	2017-09-27 00:00:00	NaN	NaN	NaN	NaN	NaN	NaN
26	READ_DT12	NaN	2017-08-28 00:00:00	NaN	NaN	NaN	NaN	NaN	NaN
27	READ_DAYS1	NaN	33	NaN	NaN	NaN	NaN	NaN	NaN
28	READ_DAYS2	NaN	28	NaN	NaN	NaN	NaN	NaN	NaN
29	READ_DAYS3	NaN	30	NaN	NaN	NaN	NaN	NaN	NaN
30	READ_DAYS4	NaN	31	NaN	NaN	NaN	NaN	NaN	NaN
31	READ_DAYS5	NaN	29	NaN	NaN	NaN	NaN	NaN	NaN
32	READ_DAYS6	NaN	30	NaN	NaN	NaN	NaN	NaN	NaN
33	READ_DAYS7	NaN	33	NaN	NaN	NaN	NaN	NaN	NaN
34	READ_DAYS8	NaN	30	NaN	NaN	NaN	NaN	NaN	NaN
35	READ_DAYS9	NaN	33	NaN	NaN	NaN	NaN	NaN	NaN
36	READ_DAYS10	NaN	29	NaN	NaN	NaN	NaN	NaN	NaN
37	READ_DAYS11	NaN	30	NaN	NaN	NaN	NaN	NaN	NaN
38	READ_DAYS12	NaN	31	NaN	NaN	NaN	NaN	NaN	NaN

Let's drop all these accounts since they are not informative

In [48]: df1 = df1[(df1.LOC_ID != 1000065066) & (df1.LOC_ID != 1000012743) & (df1.LOC_ID != 1000024304) & (df1.LOC_ID != 1000063905) & (df1.LOC_ID != 1000054741) & (df1.LOC_ID != 1000035946) & (df1.LOC_ID != 1000062574)]
df2 = df2[(df2.LOC_ID != 1000065066) & (df2.LOC_ID != 1000012743) & (df2.LOC_ID != 1000024304) & (df2.LOC_ID != 1000063905) & (df2.LOC_ID != 1000054741) & (df2.LOC_ID != 1000035946) & (df2.LOC_ID != 1000062574)]
df3 = df3[(df3.LOC_ID != 1000065066) & (df3.LOC_ID != 1000012743) & (df3.LOC_ID != 1000024304) & (df3.LOC_ID != 1000063905) & (df3.LOC_ID != 1000054741) & (df3.LOC_ID != 1000035946) & (df3.LOC_ID != 1000062574)]

In [49]: #check the value of Dwelling units for Duplex houses
temp = df1[(df1.RATE == "RDUPLEX")]
temp.DWEL_UNIT.value_counts(dropna=False, normalize=True)

Out [49]:

2.0 1.0
Name: DWEL_UNIT, dtype: float64

In [50]: #check the value of Dwelling units for Duplex houses
temp = df2[(df2.RATE == "RDUPLEX")]
temp.DWEL_UNIT.value_counts(dropna=False, normalize=True)

Out [50]:

2.0 0.854485
1.0 0.145377
2.0 0.000138
Name: DWEL_UNIT, dtype: float64

In [51]: #For duplex group with number of dwelling unit of 3, set the number of dwelling unit to 2
df3.loc[(df3.RATE == "RDUPLEX") & (df3.DWEL_UNIT == 3), 'DWEL_UNIT'] = 2

In [52]: #For duplex group with number of dwelling unit of 1, set the number of dwelling unit to 2 and double the amount of usage.
df2.loc[(df2.RATE == "RDUPLEX") & (df2.DWEL_UNIT == 1), 'USE1'] = 2*
df2.loc[(df2.RATE == "RDUPLEX") & (df2.DWEL_UNIT == 1), 'USE1'] = 2*
df2.loc[(df2.RATE == "RDUPLEX") & (df2.DWEL_UNIT == 1), 'USE2'] = 2*
df2.loc[(df2.RATE == "RDUPLEX") & (df2.DWEL_UNIT == 1), 'USE2'] = 2*
df2.loc[(df2.RATE == "RDUPLEX") & (df2.DWEL_UNIT == 1), 'USE3'] = 2*
df2.loc[(df2.RATE == "RDUPLEX") & (df2.DWEL_UNIT == 1), 'USE3'] = 2*
df2.loc[(df2.RATE == "RDUPLEX") & (df2.DWEL_UNIT == 1), 'USE4'] = 2*
df2.loc[(df2.RATE == "RDUPLEX") & (df2.DWEL_UNIT == 1), 'USE4'] = 2*
df2.loc[(df2.RATE == "RDUPLEX") & (df2.DWEL_UNIT == 1), 'USE5'] = 2*
df2.loc[(df2.RATE == "RDUPLEX") & (df2.DWEL_UNIT == 1), 'USE5'] = 2*
df2.loc[(df2.RATE == "RDUPLEX") & (df2.DWEL_UNIT == 1

In [73]: temp=df1[(df1.RATE == 'RSFD')]
temp.DWEL_UNIT.value_counts(dropna=False, normalize=True)

Out [73]:

1.0	1.0
Name: DWEL_UNIT, dtype: float64	

In [74]: temp=df2[(df2.RATE == 'RSFD')]
temp.DWEL_UNIT.value_counts(dropna=False, normalize=True)

Out [74]:

1.0	0.999703
2.0	0.000273
4.0	0.000019
Name: DWEL_UNIT, dtype: float64	

In [75]: df2[(df2.RATE == 'RSFD') & (df2.DWEL_UNIT == 2)]

	LOC_ID	RATE	DWEL_UNIT	USE1	USE2	USE3	USE4	USE5	USE6	USE7	...	READ_DAYS3	READ_DAYS4	READ_DAYS5
1959	100001959	RSFD	2.0	26.0	30.0	27.0	22.0	9.0	11.0	24.0	...	28.0	30.0	31.0
5070	100005070	RSFD	2.0	7.0	7.0	6.0	6.0	6.0	6.0	5.0	...	28.0	30.0	31.0
5946	1000005946	RSFD	2.0	41.0	29.0	14.0	16.0	13.0	19.0	19.0	...	32.0	29.0	29.0
7286	1000007286	RSFD	2.0	23.0	18.0	17.0	18.0	7.0	8.0	12.0	...	30.0	31.0	29.0
7835	1000007835	RSFD	2.0	5.0	5.0	5.0	3.0	3.0	6.0	5.0	...	31.0	28.0	29.0
29827	1000029827	RSFD	2.0	12.0	11.0	11.0	9.0	7.0	8.0	9.0	...	32.0	29.0	29.0
29443	1000029443	RSFD	2.0	17.0	14.0	13.0	13.0	10.0	8.0	12.0	...	32.0	29.0	29.0
31848	1000031848	RSFD	2.0	29.0	38.0	33.0	47.0	36.0	106.0	100.0	...	30.0	29.0	26.0
35104	1000035104	RSFD	2.0	8.0	9.0	10.0	9.0	6.0	7.0	10.0	...	30.0	31.0	30.0
46745	1000046745	RSFD	2.0	11.0	10.0	13.0	11.0	10.0	12.0	10.0	...	32.0	29.0	29.0
50331	1000050331	RSFD	2.0	7.0	6.0	6.0	6.0	5.0	5.0	8.0	...	30.0	29.0	31.0
51160	1000051160	RSFD	2.0	29.0	28.0	24.0	25.0	9.0	21.0	20.0	...	31.0	29.0	29.0
66500	1000066500	RSFD	2.0	11.0	10.0	8.0	9.0	8.0	10.0	11.0	...	29.0	31.0	29.0
68845	1000068845	RSFD	2.0	8.0	10.0	8.0	8.0	8.0	24.0	30.0	...	29.0	31.0	29.0
70919	1000070919	RSFD	2.0	6.0	8.0	5.0	4.0	5.0	5.0	5.0	...	29.0	31.0	30.0

15 rows × 39 columns

It seems the number of dwelling units are incorrectly recorded and they have 1 dwelling, since the amount of usage is similar to the one dwelling's usage.

```

In [76]: df2.loc[(df2.RATE == 'RSFD') & (df2.DWEL_UNIT == 2), 'DWEL_UNIT'] = 1

In [77]: df2[(df2.RATE == 'RSFD') & (df2.DWEL_UNIT == 4)]

Out[77]:

```

	LOC_ID	RATE	DWEL_UNIT	USE1	USE2	USE3	USE4	USE5	USE6	USE7	...	READ_DAYS3	READ_DAYS4	READ_DAYS5
9349	1000009349	RSFD	4.0	19.0	31.0	29.0	31.0	31.0	37.0	37.0	...	31.0	30.0	29.0

1 rows × 39 columns

```

In [78]: df2.loc[(df2.RATE == 'RSFD') & (df2.DWEL_UNIT == 4), 'DWEL_UNIT'] = 1

In [79]: temp = df2[(df2.RATE == 'RSFD')]
temp.DWEL_UNIT.value_counts(dropna=False, normalize=True)

Out[79]: 1.0    1.0
Name: DWEL_UNIT, dtype: float64

In [80]: temp = df3[(df3.RATE == 'RSFD')]
temp.DWEL_UNIT.value_counts(dropna=False, normalize=True)

Out[80]: 1.0    0.999702

```

15 rows x 39 columns

It seems the number of dwelling units are incorrectly recorded and they have 1 dwelling, since the amount of usage is similar to the one dwelling's usage.

In [76]: df2.loc[(df2.RATE == 'RSFD') & (df2.DWEL_UNIT == 2), 'DWEL_UNIT'] = 1

In [77]: df2[(df2.RATE == 'RSFD') & (df2.DWEL_UNIT == 4)]

Out [77]:

	LOC_ID	RATE	DWEL_UNIT	USE1	USE2	USE3	USE4	USE5	USE6	USE7	...	READ_DAYS3	READ_DAYS4	READ_DAYS5
9349	100009349	RSFD	4.0	30.0	29.0	29.0	28.0	21.0	21.0	26.0	...	29.0	31.0	30.0

1 rows x 39 columns

In [78]: df2.loc[(df2.RATE == 'RSFD') & (df2.DWEL_UNIT == 4), 'DWEL_UNIT'] = 1

In [79]: temp=df2[(df2.RATE == 'RSFD')]
temp.DWEL_UNIT.value_counts(dropna=False, normalize=True)

Out [79]:

1.0	1.0
Name: DWEL_UNIT, dtype: float64	

In [80]: temp=df3[(df3.RATE == 'RSFD')]
temp.DWEL_UNIT.value_counts(dropna=False, normalize=True)

Out [80]:

1.0	0.999702
2.0	0.000273
4.0	0.000019
Name: DWEL_UNIT, dtype: float64	

In [81]: df3[(df3.RATE == 'RSFD') & (df3.DWEL_UNIT == 4)]

Out [81]:

	LOC_ID	RATE	DWEL_UNIT	USE1	USE2	USE3	USE4	USE5	USE6	USE7	...	READ_DAYS3	READ_DAYS4	READ_DAYS5
9349	100009349	RSFD	4.0	30.0	29.0	29.0	28.0	21.0	21.0	26.0	...	29.0	31.0	30.0

1 rows x 39 columns

In [82]: df3.loc[(df3.RATE == 'RSFD') & (df3.DWEL_UNIT == 4), 'DWEL_UNIT'] = 1

In [83]: df3[(df3.RATE == 'RSFD') & (df3.DWEL_UNIT == 2)]

29443	1000029443	RSFD	2.0	21.0	21.0	19.0	15.0	11.0	11.0	12.0	...	30.0	31.0	29.0
31848	1000031848	RSFD	2.0	14.0	19.0	19.0	17.0	15.0	14.0	17.0	...	32.0	29.0	29.0
35104	1000035104	RSFD	2.0	10.0	9.0	9.0	7.0	7.0	6.0	9.0	...	30.0	28.0	31.0

15 rows x 39 columns

It seems the number of dwelling units are incorrectly recorded and they have 1 dwelling, since the amount of usage is similar to the one dwelling's usage.

In [84]: df3.loc[(df3.RATE == 'RSFD') & (df3.DWEL_UNIT == 2), 'DWEL_UNIT'] = 1

In [85]: temp=df3[(df3.RATE == 'RSFD')]
temp.DWEL_UNIT.value_counts(dropna=False, normalize=True)

Out [85]:

1.0	1.0
Name: DWEL_UNIT, dtype: float64	

Now all dwelling units in single family group have value of 1.

4.2.2) Filter the accounts in multifamily group with dwelling unit with null value.

In [86]: temp=df1[(df1.RATE == 'RMF')]
temp.DWEL_UNIT.value_counts(dropna=False, normalize=True)

Out [86]:

4.0	0.240540
3.0	0.169509
1.0	0.103219
8.0	0.082885
6.0	0.067572
92.0	0.000091
120.0	0.000091
201.0	0.000091
91.0	0.000091
183.0	0.000091
Name: DWEL_UNIT, Length: 126, dtype: float64	

In [87]: df2[(df2.RATE == 'RMF') & (df2.DWEL_UNIT.isna())]

Out [87]:

	LOC_ID	RATE	DWEL_UNIT	USE1	USE2	USE3	USE4	USE5	USE6	USE7	...	READ_DAYS3	READ_DAYS4	READ_DAYS5
52668	1000052668	RMF	NaN	NaN	NaN	NaN	NaN	NaN	NaN	NaN	...	NaN	NaN	NaN

1 rows x 39 columns

In [88]: df3[(df3.RATE == 'RMF') & (df3.DWEL_UNIT.isna())]

Out [88]:

	LOC_ID	RATE	DWEL_UNIT	USE1	USE2	USE3	USE4	USE5	USE6	USE7	...	READ_DAYS3	READ_DAYS4	READ_DAYS5
52668	1000052668	RMF	NaN	NaN	NaN	NaN	NaN	NaN	NaN	NaN	...	NaN	NaN	NaN

1 rows x 39 columns

In [89]: df1 = df1[(df1.LOC_ID != 1000052668)]
df2 = df2[(df2.LOC_ID != 1000052668)]
df3 = df3[(df3.LOC_ID != 1000052668)]

In [90]: temp=df1[(df1.RATE == 'RMF')]
temp.DWEL_UNIT.value_counts(dropna=False, normalize=True)

Out [90]:

4.0	0.240542
3.0	0.169433
1.0	0.103228
8.0	0.082893
6.0	0.067572
92.0	0.000091
120.0	0.000091
201.0	0.000091
91.0	0.000091
183.0	0.000091
Name: DWEL_UNIT, Length: 126, dtype: float64	

In [91]: temp=df2[(df1.RATE == 'RMF')]
temp.DWEL_UNIT.value_counts(dropna=False, normalize=True)

Out [91]:

4.0	0.240744
3.0	0.169524
1.0	0.103411
8.0	0.082801
6.0	0.067572
92.0	0.000091
120.0	0.000091
201.0	0.000091
91.0	0.000091
183.0	0.000091
Name: DWEL_UNIT, Length: 126, dtype: float64	

In [92]: temp=df2[(df1.RATE == 'RMF')]
temp.DWEL_UNIT.value_counts(dropna=False, normalize=True)

Out [92]:

4.0	0.240744
3.0	0.169524
1.0	0.103411
8.0	0.082801
6.0	0.067572
92.0	0.000091
120.0	0.000091
201.0	0.000091
91.0	0.000091
183.0	0.000091
Name: DWEL_UNIT, Length: 126, dtype: float64	

There is no Null value in dwelling unit column of multifamily group in all 3 dataframes.

4.2.2) Filter the dwelling unit with zero value

In [93]: #Filter the RATEs that their dwelling units are zero.
print(len(df1.DWEL_UNIT == 0)['RATE'].unique())
print(len(df2.DWEL_UNIT == 0)['RATE'].unique())
print(len(df3.DWEL_UNIT == 0)['RATE'].unique())

[IND]
Categories (1, object): [IND]
[IND]
Categories (1, object): [IND]
[IND]
Categories (1, object): [IND]

Industrial group has some accounts with dwelling unit of zero. Let's look at them.

4.2.2.1) Accounts in industrial group with zero dwelling

In [94]: #Check the number of null in DWEL_UNIT column with the RATE of Industrial in df1.
print(df1.RATE == 'IND')['DWEL_UNIT'].isna().sum())
#Calculate sum of DWEL_UNIT with the RATE of Industrial.
print(df1.RATE == 'IND')['DWEL_UNIT'].sum())

Out [94]:

0	0.0
---	-----

In [95]: #Check the number of null in DWEL_UNIT column with the RATE of Industrial in df2.
print(df2.RATE == 'IND')['DWEL_UNIT'].isna().sum())
#Calculate sum of DWEL_UNIT with the RATE of Industrial.
print(df2.RATE == 'IND')['DWEL_UNIT'].sum())

Out [95]:

0	0.0
---	-----

In [96]: #Check the number of null in DWEL_UNIT column with the RATE of Industrial in df3.
print(df3.RATE == 'IND')['DWEL_UNIT'].isna().sum())
#Calculate sum of DWEL_UNIT with the RATE of Industrial.
print(df3.RATE == 'IND')['DWEL_UNIT'].sum())

Out [96]:

0	0.0
---	-----

Number of dwelling units for industrial customers in all three dataframes is zero. Let's set the value one for dwelling units of industrial customers.

In [97]: #Set the value one for dwelling units of industrial customers.
df1.loc[df1.RATE == 'IND', 'DWEL_UNIT'] = 1
df2.loc[df2.RATE == 'IND', 'DWEL_UNIT'] = 1
df3.loc[df3.RATE == 'IND', 'DWEL_UNIT'] = 1

In [98]: #Filter the RATEs that their dwelling units are zero.
print(len(df1.DWEL_UNIT == 0))
print(len(df2.DWEL_UNIT == 0))
print(len(df3.DWEL_UNIT == 0))

Out [98]:

0	0	0
---	---	---

There is no zero value in dwelling unit column in all 3 dataframes.

5) What is the market share of different RATE?

In [99]: def count(data,rate:str):
'''this function takes the type of customer and returns the
number of customer, number and percentage of missing values associated to.'''
temp = data[data.RATE == rate]
number = temp['DWEL_UNIT'].sum()
missing_num = temp['DWEL_UNIT'].isna().sum()
missing_per = round(100* temp['DWEL_UNIT'].isna().sum()/len(df1), 2)
return (number, missing_num, missing_per)

df1

In [100]: #Calculate the market share of each Rate in first year, the number of Rate based on the number of Dwel
ling unit
#and % of related missing values.
dict1= {rate:count(df1,rate) for rate in df1.RATE.unique()}
customer1 = pd.DataFrame(dict1)
customer1.columns = ['Singel_family', 'Multi_family','Duplex','Commercial','Irrigation','Industrial',
'Commercial and Residential']
customer1.columns = ['count', '% missing_value', '% missing_value']
customer1['market_share'] = list(round(100* df1.RATE.value_counts()/len(df1),2))

Out [100]:

	market_share	count	# missing_value	% missing_value
Singel_family	68.23	53827.0	0.0	0.0
Multi_family	13.90	88061.0	0.0	0.0
Duplex	9.19	14498.0	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

df2

In [101]: #Calculate the market share of each Ratein second year, the number of Rate based on the number of Dwel
ling 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 = customer2.2
customer2.columns = ['count', '% missing_value', '% missing_value']
customer2['market_share'] = list(round(100* df2.RATE.value_counts()/len(df2),2))

Out [101]:

	market_share	count	# missing_value	% missing_value
Singel_family	68.23	53827.0	0.0	0.0
Multi_family	13.90	88061.0	0.0	0.0
Duplex	9.20	14498.0	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

df3

In [102]: #Calculate the market share of each Ratein second year, the number of Rate based on the number of Dwel
ling 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 = customer3.3
customer3.columns = ['count', '% missing_value', '% missing_value']
customer3['market_share'] = list(round(100* df3.RATE.value_counts()/len(df3),2))

Out [102]:

	market_share	count	# missing_value	% missing_value
Singel_family	68.11	53827.0	0.0	0.0
Multi_family	13.90	88061.0	0.0	0.0
Duplex	9.31	14498.0	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	


```
In [119]: #Create a list of DataFrames, each dataframe with index of month(s) between previous read date and read date,
#and extracting the month and year of index(s)
def yr_month(smali):
    result = []
    for tup in small.iteruples():
        #find the month between previous read date and read date
        data['month_start'] = pd.PeriodIndex([tup.Previous_READ_DT, tup.READ_DT], freq='M')
        temp = pd.DataFrame({(tup.LOC_ID, tup.RATE, tup.Period, tup.DWEL_UNIT, tup.USE, tup.READ_DT,
        temp['month'] = temp.index.month
        temp['year'] = temp.index.year
        result.append(temp)
    result = pd.concat(result, axis=0)
    return result

In [120]: def days_in_period(data):
#Create the month start for each row
index = pd.PeriodIndex([data.year * 10000 + data.month * 100 + 1].apply(str), format='%Ym%d')

#Find number of days in the month
data['days_in_month'] = data[['year', 'month']].apply(lambda x: calendar.monthrange(x[0],x[1])[1], axis=1)

#Count number of days between start of month and read date
data['read_date_to_month_start'] = data[['month_start','READ_DT']].apply(
    lambda x: x[2] - pd.Timedelta(x[1]-x[0], 'D').days, axis=1)

#Count number of days between previous read date and end of the month
data['previous_read_date_to_month_end'] = data[['previous_read_date','month_start','days_in_month'
]].apply(
    lambda x: x[2] - pd.Timedelta(x[0]-x[1], 'D').days, axis=1)

#Count number of days in each period
data['days_to_count_in_period'] = data[['previous_read_date_to_month_end','month_start','
'read_date_to_month_start','days_in_month']].apply(
    lambda x: x[3] if x[0] > x[3] & x[2] >= x[3]
else x[0] if x[0] <= x[3] & else x[2], axis=1)

return data[['RATE','LOC_ID','DWEL_UNIT','month_start','days_in_month','previous_read_date',
'READ_DAYS','READ_DT','month','year','days_to_count_in_period','USE']]
```

df1

```
In [121]: #Chunk the long df1
chunk1 = long_df1.iloc[0:300000,:]
chunk2 = long_df1.iloc[300000,:]
```

```
In [122]: result1 = yr_month(chunk1)
```

```
In [123]: result1.head()
```

```
Out[123]:
```

		0	1	2	3	4	5	6	7	month	year
2018-06	1000000000	RSFD	1	1.0	22.0	2018-07-11	30.0	2018-06-11	6	2018	
2018-07	1000000000	RSFD	1	1.0	22.0	2018-07-11	30.0	2018-06-11	7	2018	
2018-05	1000000000	RSFD	2	1.0	22.0	2018-06-11	32.0	2018-05-10	5	2018	
2018-06	1000000000	RSFD	2	1.0	22.0	2018-06-11	32.0	2018-05-10	6	2018	
2018-04	1000000000	RSFD	3	1.0	15.0	2018-05-10	29.0	2018-04-11	4	2018	

```
In [124]: result2 = yr_month(chunk2)
```

```
In [125]: #Concatenate the results of chunks, and rename the columns
long_df1 = pd.concat([result1,result2], axis=0).reset_index()

#Label columns appropriately
long_df1.rename(columns = {(0:'LOC_ID', 1:'RATE', 2:'period',3:'DWEL_UNIT',4:'USE',
5:'READ_DT',6:'READ_DAYS', 7:'previous_read_date'), inplace=True)

print(long_df1.shape)
long_df1.head()
(1832948, 11)
```

```
Out[125]:
```

	index	LOC_ID	RATE	period	DWEL_UNIT	USE	READ_DT	READ_DAYS	previous_read_date	month	year
0	2018-06	1000000000	RSFD	1	1.0	22.0	2018-07-11	30.0	2018-06-11	6	2018
1	2018-07	1000000000	RSFD	1	1.0	22.0	2018-07-11	30.0	2018-06-11	7	2018
2	2018-05	1000000000	RSFD	2	1.0	22.0	2018-06-11	32.0	2018-05-10	5	2018
3	2018-06	1000000000	RSFD	2	1.0	22.0	2018-06-11	32.0	2018-05-10	6	2018
4	2018-04	1000000000	RSFD	3	1.0	15.0	2018-05-10	29.0	2018-04-11	4	2018

```
In [126]: long_df1.reset_index(inplace=True)
small_df1 = days_in_period(long_df1)
small_df1.head()
```

```
Out[126]:
```

	RATE	LOC_ID	DWEL_UNIT	month_start	days_in_month	previous_read_date	READ_DAYS	READ_DT	month	year	days_to_count
0	RSFD	1000000000		2018-06-01	30	2018-06-11	30.0	2018-07-11	6	2018	
1	RSFD	1000000000		2018-07-01	31	2018-06-11	30.0	2018-07-11	7	2018	
2	RSFD	1000000000		2018-05-01	31	2018-05-10	32.0	2018-06-11	5	2018	
3	RSFD	1000000000		2018-06-01	30	2018-05-10	32.0	2018-06-11	6	2018	
4	RSFD	1000000000		2018-04-01	30	2018-04-11	29.0	2018-05-10	4	2018	

```
In [127]: small_df1['usage_to_count_in_period'] = round((small_df1.days_to_count_in_period / small_df1.READ_DAYS) * small_df1.USE, 1)
small_df1.head()
```

```
<ipython-input-127-13baa3d7d8d5>:1: SettingWithCopyWarning:
A value is trying to be set on a copy of a DataFrame.
Try using .loc[row_indexer,col_indexer] = value instead

See the caveats in the documentation: https://pandas.pydata.org/pandas-docs/stable/user_guide/indexin
g.html#returning-a-view-versus-a-copy
small_df1['usage_to_count_in_period'] = round((small_df1.days_to_count_in_period / small_df1.READ_D
AYS) * small_df1.USE, 1)
```

```
Out[127]:
```

	RATE	LOC_ID	DWEL_UNIT	month_start	days_in_month	previous_read_date	READ_DAYS	READ_DT	month	year	days_to_count
0	RSFD	1000000000		2018-06-01	30	2018-06-11	30.0	2018-07-11	6	2018	
1	RSFD	1000000000		2018-07-01	31	2018-06-11	30.0	2018-07-11	7	2018	
2	RSFD	1000000000		2018-05-01	31	2018-05-10	32.0	2018-06-11	5	2018	
3	RSFD	1000000000		2018-06-01	30	2018-05-10	32.0	2018-06-11	6	2018	
4	RSFD	1000000000		2018-04-01	30	2018-04-11	29.0	2018-05-10	4	2018	

```
In [128]: small_df1.isnull().sum()
```

```
Out[128]:
```

	RATE	LOC_ID	DWEL_UNIT	month_start	days_in_month	previous_read_date	READ_DAYS	READ_DT	month	year	days_to_count
0	RSFD	1000000000		2018-06-01	30	2018-06-11	30.0	2018-07-11	6	2018	
1	RSFD	1000000000		2018-07-01	31	2018-06-11	30.0	2018-07-11	7	2018	
2	RSFD	1000000000		2018-05-01	31	2018-05-10	32.0	2018-06-11	5	2018	
3	RSFD	1000000000		2018-06-01	30	2018-05-10	32.0	2018-06-11	6	2018	
4	RSFD	1000000000		2018-04-01	30	2018-04-11	29.0	2018-05-10	4	2018	

```
In [129]: small_df1.describe()
```

```
Out[129]:
```

	LOC_ID	DWEL_UNIT	days_in_month	READ_DAYS	month	year	days_to_count_in_period	USE_u
count	1.8832948e+06	1.832948e+06	1.832948e+06	1.8832948e+06	1.832948e+06	1.832948e+06	1.832948e+06	1.832948e+06
mean	1.000039e+09	2.091431e+00	3.044337e+01	3.045039e+01	6.466700e+00	2.017569e+03	1.524685e+01	2.191004e+01
std	2.277445e+04	7.366812e+00	8.410156e+01	2.887039e+00	3.413716e+00	4.951835e+01	8.482162e+00	1.772232e+02
min	1.000000e+09	1.000000e+00	2.800000e+01	8.800000e+00	1.000000e+00	2.016000e+03	0.000000e+00	0.000000e+00
25%	1.000020e+09	1.000000e+00	3.000000e+01	2.900000e+01	4.000000e+00	2.017000e+03	8.000000e+00	6.000000e+00
50%	1.000039e+09	1.000000e+00	3.100000e+01	3.000000e+01	6.000000e+00	2.018000e+03	1.500000e+01	1.000000e+01
75%	1.000059e+09	1.000000e+00	3.100000e+01	3.200000e+01	9.000000e+00	2.018000e+03	2.200000e+01	1.700000e+01
max	1.000079e+09	5.700000e+02	3.100000e+01	4.830000e+02	1.200000e+01	2.018000e+03	3.100000e+01	4.817200e+04

Min of days_to_count_in_period is zero.

```
In [130]: small_df1 = small_df1.groupby(['RATE', 'month_start'])\
[ ['DWEL_UNIT', 'days_to_count_in_period', 'usage_to_count_in_period']].sum
().reset_index()
small_df1.rename(columns= ('days_to_count_in_period': 'days_in_period', 'usage_to_count_in_period': 'u
sage_in_period'))
small_df1.head()
```

```
Out[130]:
```

	RATE	month_start	DWEL_UNIT	days_to_count_in_period	usage_to_count_in_period
0	COM	2016-04-01	1.0	19	1.4
1	COM	2016-05-01	1.0	12	0.5
2	COM	2016-11-01	3.0	58	5.7
3	COM	2016-12-01	4.0	62	6.7
4	COM	2017-01-01	5.0	69	8.4

df2

```
In [131]: #Chunk the long df2
chunk3 = long_df2.iloc[0:450000,:]
chunk4 = long_df2.iloc[450000,:]
```

```
In [132]: result3 = yr_month(chunk3)
```

```
In [133]: result4 = yr_month(chunk4)
```

```
In [134]: long_df2 = pd.concat([result3,result4], axis=0).reset_index()
long_df2.rename(columns = {(0:'LOC_ID', 1:'RATE', 2:'period',3:'DWEL_UNIT',4:'USE',
5:'READ_DT',6:'READ_DAYS', 7:'previous_read_date'), inplace=True)
print(long_df2.shape)
long_df2.head()
```

```
Out[134]:
```

	index	LOC_ID	RATE	period	DWEL_UNIT	USE	READ_DT	READ_DAYS	previous_read_date	month	year
0	2019-06	1000000000	RSFD	1	1.0	20.0	2019-07-10	30.0	2019-06-10	6	2019
1	2019-07	1000000000	RSFD	1	1.0	20.0	2019-07-10	30.0	2019-06-10	7	2019
2	2019-05	1000000000	RSFD	2	1.0	13.0	2019-06-10	31.0	2019-05-10	5	2019
3	2019-06	1000000000	RSFD	2	1.0	13.0	2019-06-10	31.0	2019-05-10	6	2019
4	2019-04	1000000000	RSFD	3	1.0	6.0	2019-05-10	29.0	2019-04-11	4	2019

```
In [135]: long_df2 = long_df2.reset_index()
small_df2 = days_in_period(long_df2)
small_df2.head()
```

```
Out[135]:
```

	RATE	LOC_ID	DWEL_UNIT	month_start	days_in_month	previous_read_date	READ_DAYS	READ_DT	month	year	days_to_count
0	RSFD	1000000000		2019-06-01	30	2019-06-10	30.0	2019-07-10	6	2019	
1	RSFD	1000000000		2019-07-01	31	2019-06-10	30.0	2019-07-10	7	2019	
2	RSFD	1000000000		2019-05-01	31	2019-05-10	31.0	2019-06-10	5	2019	
3	RSFD	1000000000		2019-06-01	30	2019-06-10	31.0	2019-06-10	6	2019	
4	RSFD	1000000000		2019-04-01	30	2019-04-11	29.0	2019-05-10	4	2019	

```
In [139]: small_df2['usage_to_count_in_period'] = round((small_df2.days_to_count_in_period / small_df2.READ_DAYS) * small_df2.USE, 1)
small_df2.head()
```

```
<ipython-input-139-9a18942feaf5>:1: SettingWithCopyWarning:
A value is trying to be set on a copy of a slice from a DataFrame.
Try using .loc[row_indexer,col_indexer] = value instead

See the caveats in the documentation: https://pandas.pydata.org/pandas-docs/stable/user_guide/indexin
g.html#returning-a-view-versus-a-copy
small_df2['usage_to_count_in_period'] = round((small_df2.days_to_count_in_period / small_df2.READ_D
AYS) * small_df2.USE, 1)
```

```
Out[139]:
```

	RATE	LOC_ID	DWEL_UNIT	month_start	days_in_month	previous_read_date	READ_DAYS	READ_DT	month	year	days_to_count
0	RSFD	1000000000		2019-06-01	30	2019-06-10	30.0	2019-07-10	6	2019	
1	RSFD	1000000000		2019-07-01	31	2019-06-10	30.0	2019-07-10	7	2019	
2	RSFD	1000000000		2019-05-01	31	2019-05-10	31.0	2019-06-10	5	2019	
3	RSFD	1000000000		2019-06-01	30	2019-05-10	31.0	2019-06-10	6	2019	
4	RSFD	1000000000		2019-04-01	30	2019-04-11	29.0	2019-05-10	4	2019	

```
In [140]: small_df2.info()
```

```
<class 'pandas.core.frame.DataFrame'>
RangeIndex: 1886324 entries, 0 to 1886323
Data columns (total 13 columns):
# Column            Dtype
---  ---
0 RATE              object
1 LOC_ID            int64
2 DWEL_UNIT         float64
3 month_start       datetime64[ns]
4 days_in_month     int64
5 previous_read_date datetime64[ns]
6 READ_DAYS         float64
7 READ_DT           datetime64[ns]
8 month            int64
9 year             int64
10 days_to_count_in_period int64
11 USE             float64
12 usage_to_count_in_period float64
dtypes: datetime64[ns](3), float64(4), int64(5), object(1)
memory usage: 187.1+ MB
```

```
In [141]: small_df2.describe()
```

```
Out[141]:
```

	LOC_ID	DWEL_UNIT	days_in_month	READ_DAYS	month	year	days_to_count_in_period	USE_u
count	1.886324e+06	1.886350e+06	1.886350e+06	1.886324e+06	1.886350e+06	1.886324e+06	1.886350e+06	1.886324e+06
mean	9.343115e+05	2.089459e+00	3.044947e+01	3.048122e+01	6.497316e+00	2.018563e+03	1.527481e+01	2.070709e+01
std	2.277606e+04	7.360700e+00	8.403949e+01	2.162080e+00	3.415857e+00	4.960451e+01	8.458260e+00	1.638603e+02
min	1.000000e+09	1.000000e+00	2.800000e+01	8.800000e+00	1.000000e+00	2.016000e+03	0.000000e+00	0.000000e+00
25%	1.000020e+09	1.000000e+00	3.000000e+01	2.900000e+01	4.000000e+00	2.016000e+03	8.000000e+00	6.000000e+00
50%	1.000039e+09	1.000000e+00	3.100000e+01	3.000000e+01	6.000000e+00	2.018000e+03	1.500000e+01	1.000000e+01
75%	1.000059e+09	1.000000e+00	3.100000e+01	3.200000e+01	9.000000e+00	2.018000e+03	2.200000e+01	1.600000e+01
max	1.000079e+09	5.700000e+02	3.100000e+01	4.830000e+02	1.200000e+01	2.018000e+03	3.100000e+01	2.609700e+04

```
In [142]: small_df2 = small_df2.groupby(['RATE', 'month_start'])\
[ ['DWEL_UNIT', 'days_to_count_in_period', 'usage_to_count_in_period']].sum
().reset_index()
small_df2.rename(columns= ('days_to_count_in_period': 'days_in_period', 'usage_to_count_in_period': 'u
sage_in_period'))
small_df2.head()
```

```
Out[142]:
```

	RATE	month_start	DWEL_UNIT	days_to_count_in_period	usage_to_count_in_period
0	COM	2016-04-01	1.0	20	4.7
1	COM	2016-05-01	2.0	30	7.3
2	COM	2016-08-01	3.0	23	14.7
3	COM	2016-03-01	7.0	107	38.8
4	COM	2018-04-01	6.0	107	163.8

df3

```
In [143]: #Chunk the long df3
chunks = long_df3.iloc[0:450000,:]
chunk6 = long_df3.iloc[450000,:]
```

```
In [144]: result5 = yr_month(chunk5)
```

```
In [145]: result6 = yr_month(chunk6)
```

```
In [146]: long_df3 = pd.concat([result5, result6], axis=0).reset_index()
long_df3.rename(columns = {(0:'LOC_ID', 1:'RATE', 2:'period',3:'DWEL_UNIT',4:'USE',
5:'READ_DT',6:'READ_DAYS', 7:'previous_read_date'), inplace=True)
long_df3.head()
```

```
Out[146]:
```

	index	LOC_ID	RATE	period	DWEL_UNIT	USE	READ_DT	READ_DAYS	previous_read_date	month	year
0	2020-06	1000000000	RSFD	1	1.0	16.0	2020-07-10	31.0	2020-06-09	6	2020
1	2020-07	1000000000	RSFD	1	1.0	16.0	2020-07-10	31.0	2020-06-09	7	2020
2	2020-05	1000000000	RSFD	2	1.0	13.0	2020-06-09	29.0	2020-05-11	5	2020
3	2020-06	1000000000	RSFD	2	1.0	13.0	2020-06-09	29.0	2020-05-11	6	2020
4	2020-04	1000000000	RSFD	3	1.0	20.0	2020-05-11	31.0	2020-04-10	4	2020

```
In [147]: long_df3 = long_df3.reset_index()
small_df3 = days_in_period(long_df3)
small_df3.head()
```

```
Out[147]:
```

	RATE	LOC_ID	DWEL_UNIT	month_start	days_in_month	previous_read_date	READ_DAYS	READ_DT	month	year	days_to_count
0	RSFD	1000000000		2020-06-01	30	2020-06-09	31.0	2020-07-10	6	2020	
1	RSFD	1000000000		2020-07-01	31	2020-06-09	31.0	2020-07-10	7	2020	
2	RSFD	1000000000		2020-05-01	31	2020-05-11	29.0	2020-06-09	5	2020	
3	RSFD	1000000000		2020-06-01	30	2020-05-11	29.0	2020-06-09	6	2020	
4	RSFD	1000000000		2020-04-01	30</						