

# Real Estate

```
In [1]: import pandas as pd
import numpy as np
import time
import random
from math import *
import operator
import matplotlib
import matplotlib.pyplot as plt
from pandas.plotting import scatter_matrix
%matplotlib inline
import seaborn as sns
sns.set(style="white", color_codes=True)
sns.set(font_scale=1.5)
import os
os.chdir(r"E:\Simplilearn\Final Project")
```

## 1. Import data

```
In [2]: df_train=pd.read_csv("train.csv")
```

```
In [3]: df_test=pd.read_csv("test.csv")
```

```
In [4]: df_train.columns
```

```
Out[4]: Index(['UID', 'BLOCKID', 'SUMLEVEL', 'COUNTYID', 'STATEID', 'state',
            'state_ab', 'city', 'place', 'type', 'primary', 'zip_code', 'area_code',
            'lat', 'lng', 'ALand', 'AWater', 'pop', 'male_pop', 'female_pop',
            'rent_mean', 'rent_median', 'rent_stdev', 'rent_sample_weight',
            'rent_samples', 'rent_gt_10', 'rent_gt_15', 'rent_gt_20', 'rent_gt_25',
            'rent_gt_30', 'rent_gt_35', 'rent_gt_40', 'rent_gt_50',
            'universe_samples', 'used_samples', 'hi_mean', 'hi_median', 'hi_stdev',
            'hi_sample_weight', 'hi_samples', 'family_mean', 'family_median',
            'family_stdev', 'family_sample_weight', 'family_samples',
            'hc_mortgage_mean', 'hc_mortgage_median', 'hc_mortgage_stdev',
            'hc_mortgage_sample_weight', 'hc_mortgage_samples', 'hc_mean',
            'hc_median', 'hc_stdev', 'hc_samples', 'hc_sample_weight',
            'home_equity_second_mortgage', 'second_mortgage', 'home_equity', 'debt',
            'second_mortgage_cdf', 'home_equity_cdf', 'debt_cdf', 'hs_degree',
            'hs_degree_male', 'hs_degree_female', 'male_age_mean',
            'male_age_median', 'male_age_stdev', 'male_age_sample_weight',
            'male_age_samples', 'female_age_mean', 'female_age_median',
            'female_age_stdev', 'female_age_sample_weight', 'female_age_samples',
            'pct_own', 'married', 'married_snp', 'separated', 'divorced'],
            dtype='object')
```

```
In [5]: df_test.columns
```

```
Out[5]: Index(['UID', 'BLOCKID', 'SUMLEVEL', 'COUNTYID', 'STATEID', 'state',
            'state_ab', 'city', 'place', 'type', 'primary', 'zip_code', 'area_code',
            'lat', 'lng', 'ALand', 'AWater', 'pop', 'male_pop', 'female_pop',
            'rent_mean', 'rent_median', 'rent_stdev', 'rent_sample_weight',
            'rent_samples', 'rent_gt_10', 'rent_gt_15', 'rent_gt_20', 'rent_gt_25',
            'rent_gt_30', 'rent_gt_35', 'rent_gt_40', 'rent_gt_50',
            'universe_samples', 'used_samples', 'hi_mean', 'hi_median', 'hi_stdev',
            'hi_sample_weight', 'hi_samples', 'family_mean', 'family_median',
            'family_stdev', 'family_sample_weight', 'family_samples',
            'hc_mortgage_mean', 'hc_mortgage_median', 'hc_mortgage_stdev',
            'hc_mortgage_sample_weight', 'hc_mortgage_samples', 'hc_mean',
            'hc_median', 'hc_stdev', 'hc_samples', 'hc_sample_weight',
            'home_equity_second_mortgage', 'second_mortgage', 'home_equity', 'debt',
            'second_mortgage_cdf', 'home_equity_cdf', 'debt_cdf', 'hs_degree',
            'hs_degree_male', 'hs_degree_female', 'male_age_mean',
            'male_age_median', 'male_age_stdev', 'male_age_sample_weight',
            'male_age_samples', 'female_age_mean', 'female_age_median',
            'female_age_stdev', 'female_age_sample_weight', 'female_age_samples',
            'pct_own', 'married', 'married_snp', 'separated', 'divorced'],
            dtype='object')
```

```
In [6]: len(df_train)
```

```
Out[6]: 27321
```

```
In [7]: len(df_test)
```

```
Out[7]: 11709
```

```
In [8]: df_train.head()
```

Out[8]:	UID	BLOCKID	SUMLEVEL	COUNTYID	STATEID	state	state_ab	city	place	type	...	female_age_mean	female_age_median	female_age_stdev	femi
	0	267822	NaN	140	53	36	New York	NY	Hamilton	Hamilton	City	...	44.48629	45.33333	22.51276

	UID	BLOCKID	SUMLEVEL	COUNTYID	STATEID	state	state_ab	city	place	type	...	female_age_mean	female_age_median	female_age_stdev	femi
1	246444	NaN	140	141	18	Indiana	IN	South Bend	Roseland	City	...	36.48391	37.58333	23.43353	
2	245683	NaN	140	63	18	Indiana	IN	Danville	Danville	City	...	42.15810	42.83333	23.94119	
3	279653	NaN	140	127	72	Puerto Rico	PR	San Juan	Guaynabo	Urban	...	47.77526	50.58333	24.32015	
4	247218	NaN	140	161	20	Kansas	KS	Manhattan	Manhattan City	City	...	24.17693	21.58333	11.10484	

5 rows × 80 columns

In [9]:

df\_test.head()

Out[9]:

	UID	BLOCKID	SUMLEVEL	COUNTYID	STATEID	state	state_ab	city	place	type	...	female_age_mean	female_age_median	female_age_stdev	femi
0	255504	NaN	140	163	26	Michigan	MI	Detroit	Dearborn Heights City	CDP	...	34.78682	33.75000	21.58531	
1	252676	NaN	140	1	23	Maine	ME	Auburn	Auburn City	City	...	44.23451	46.66667	22.37036	
2	276314	NaN	140	15	42	Pennsylvania	PA	Pine City	Millerton	Borough	...	41.62426	44.50000	22.86213	
3	248614	NaN	140	231	21	Kentucky	KY	Monticello	Monticello City	City	...	44.81200	48.00000	21.03155	
4	286865	NaN	140	355	48	Texas	TX	Corpus Christi	Edroy	Town	...	40.66618	42.66667	21.30900	

5 rows × 80 columns

In [10]:

df\_train.describe()

Out[10]:

	UID	BLOCKID	SUMLEVEL	COUNTYID	STATEID	zip_code	area_code	lat	lng	ALand	...	female_age_mean	female_age_median	female_age_stdev	femi
count	27321.000000	NaN	0.0	27321.0	27321.000000	27321.000000	27321.000000	27321.000000	27321.000000	27321.000000	2.732100e+04	...	27115.000000	27115.000000	
mean	257331.996303	NaN	140.0	85.646426	28.271806	50081.999524	596.507668	37.508813	-91.288394	1.295106e+08	...	40.319803	40.319803	40.319803	
std	21343.859725	NaN	0.0	98.333097	16.392846	29558.115660	232.497482	5.588268	16.343816	1.275531e+09	...	5.886317	5.886317	5.886317	
min	220342.000000	NaN	140.0	1.000000	1.000000	602.000000	201.000000	17.929085	-165.453872	4.113400e+04	...	16.008330	16.008330	16.008330	
25%	238816.000000	NaN	140.0	29.000000	13.000000	26554.000000	405.000000	33.899064	-97.816067	1.799408e+06	...	36.892050	36.892050	36.892050	
50%	257220.000000	NaN	140.0	63.000000	28.000000	47715.000000	614.000000	38.755183	-86.554374	4.866940e+06	...	40.373320	40.373320	40.373320	
75%	275818.000000	NaN	140.0	109.000000	42.000000	77093.000000	801.000000	41.380606	-79.782503	3.359820e+07	...	43.567120	43.567120	43.567120	
max	294334.000000	NaN	140.0	840.000000	72.000000	99925.000000	989.000000	67.074017	-65.379332	1.039510e+11	...	79.837390	79.837390	79.837390	

8 rows × 74 columns

In [11]:

df\_test.describe()

Out[11]:

	UID	BLOCKID	SUMLEVEL	COUNTYID	STATEID	zip_code	area_code	lat	lng	ALand	...	female_age_mean	female_age_median	female_age_stdev	femi
count	11709.000000	NaN	0.0	11709.0	11709.000000	11709.000000	11709.000000	11709.000000	11709.000000	1.170900e+04	...	11613.000000	11613.000000	11613.000000	
mean	257525.004783	NaN	140.0	85.710650	28.489196	50123.418396	593.598514	37.405491	-91.340229	1.095500e+08	...	40.111999	40.111999	40.111999	
std	21466.372658	NaN	0.0	99.304334	16.607262	29775.134038	232.074263	5.625904	16.407818	7.624940e+08	...	5.851192	5.851192	5.851192	
min	220336.000000	NaN	140.0	1.000000	1.000000	601.000000	201.000000	17.965835	-166.770979	8.299000e+03	...	15.360240	15.360240	15.360240	
25%	238819.000000	NaN	140.0	29.000000	13.000000	25570.000000	404.000000	33.919813	-97.816561	1.718660e+06	...	36.729210	36.729210	36.729210	
50%	257651.000000	NaN	140.0	61.000000	28.000000	47362.000000	612.000000	38.618093	-86.643344	4.835000e+06	...	40.196960	40.196960	40.196960	
75%	276300.000000	NaN	140.0	109.000000	42.000000	77406.000000	787.000000	41.232973	-79.697311	3.204540e+07	...	43.496490	43.496490	43.496490	
max	294333.000000	NaN	140.0	810.000000	72.000000	99929.000000	989.000000	64.804269	-65.695344	5.520166e+10	...	90.107940	90.107940	90.107940	

8 rows × 74 columns

In [12]:

df\_train.info()

<class 'pandas.core.frame.DataFrame'>  
RangeIndex: 27321 entries, 0 to 27320  
Data columns (total 80 columns):  
# Column Non-Null Count Dtype   
--- -  
0 UID 27321 non-null int64   
1 BLOCKID 0 non-null float64  
2 SUMLEVEL 27321 non-null int64   
3 COUNTYID 27321 non-null int64

```

4 STATEID 27321 non-null int64
5 state 27321 non-null object
6 state_ab 27321 non-null object
7 city 27321 non-null object
8 place 27321 non-null object
9 type 27321 non-null object
10 primary 27321 non-null object
11 zip_code 27321 non-null int64
12 area_code 27321 non-null int64
13 lat 27321 non-null float64
14 lng 27321 non-null float64
15 ALand 27321 non-null float64
16 AWater 27321 non-null int64
17 pop 27321 non-null int64
18 male_pop 27321 non-null int64
19 female_pop 27321 non-null int64
20 rent_mean 27007 non-null float64
21 rent_median 27007 non-null float64
22 rent_stdev 27007 non-null float64
23 rent_sample_weight 27007 non-null float64
24 rent_samples 27007 non-null float64
25 rent_gt_10 27007 non-null float64
26 rent_gt_15 27007 non-null float64
27 rent_gt_20 27007 non-null float64
28 rent_gt_25 27007 non-null float64
29 rent_gt_30 27007 non-null float64
30 rent_gt_35 27007 non-null float64
31 rent_gt_40 27007 non-null float64
32 rent_gt_50 27007 non-null float64
33 universe_samples 27321 non-null int64
34 used_samples 27321 non-null int64
35 hi_mean 27053 non-null float64
36 hi_median 27053 non-null float64
37 hi_stdev 27053 non-null float64
38 hi_sample_weight 27053 non-null float64
39 hi_samples 27053 non-null float64
40 family_mean 27023 non-null float64
41 family_median 27023 non-null float64
42 family_stdev 27023 non-null float64
43 family_sample_weight 27023 non-null float64
44 family_samples 27023 non-null float64
45 hc_mortgage_mean 26748 non-null float64
46 hc_mortgage_median 26748 non-null float64
47 hc_mortgage_stdev 26748 non-null float64
48 hc_mortgage_sample_weight 26748 non-null float64
49 hc_mortgage_samples 26748 non-null float64
50 hc_mean 26721 non-null float64
51 hc_median 26721 non-null float64
52 hc_stdev 26721 non-null float64
53 hc_samples 26721 non-null float64
54 hc_sample_weight 26721 non-null float64
55 home_equity_second_mortgage 26864 non-null float64
56 second_mortgage 26864 non-null float64
57 home_equity 26864 non-null float64
58 debt 26864 non-null float64
59 second_mortgage_cdf 26864 non-null float64
60 home_equity_cdf 26864 non-null float64
61 debt_cdf 26864 non-null float64
62 hs_degree 27131 non-null float64
63 hs_degree_male 27121 non-null float64
64 hs_degree_female 27098 non-null float64
65 male_age_mean 27132 non-null float64
66 male_age_median 27132 non-null float64
67 male_age_stdev 27132 non-null float64
68 male_age_sample_weight 27132 non-null float64
69 male_age_samples 27132 non-null float64
70 female_age_mean 27115 non-null float64
71 female_age_median 27115 non-null float64
72 female_age_stdev 27115 non-null float64
73 female_age_sample_weight 27115 non-null float64
74 female_age_samples 27115 non-null float64
75 pct_own 27053 non-null float64
76 married 27130 non-null float64
77 married_snp 27130 non-null float64
78 separated 27130 non-null float64
79 divorced 27130 non-null float64
dtypes: float64(62), int64(12), object(6)
memory usage: 16.7+ MB

```

In [13]:

```
df_test.info()
```

```

<class 'pandas.core.frame.DataFrame'>
RangeIndex: 11709 entries, 0 to 11708
Data columns (total 80 columns):
#   Column                Non-Null Count  Dtype
---  -
0   UID                    11709 non-null  int64
1   BLOCKID                0 non-null      float64
2   SUMLEVEL               11709 non-null  int64
3   COUNTYID               11709 non-null  int64
4   STATEID                11709 non-null  int64
5   state                  11709 non-null  object
6   state_ab               11709 non-null  object
7   city                   11709 non-null  object
8   place                  11709 non-null  object
9   type                   11709 non-null  object
10  primary                11709 non-null  object

```

```

11 zip_code 11709 non-null int64
12 area_code 11709 non-null int64
13 lat 11709 non-null float64
14 lng 11709 non-null float64
15 ALand 11709 non-null int64
16 AWater 11709 non-null int64
17 pop 11709 non-null int64
18 male_pop 11709 non-null int64
19 female_pop 11709 non-null int64
20 rent_mean 11561 non-null float64
21 rent_median 11561 non-null float64
22 rent_stdev 11561 non-null float64
23 rent_sample_weight 11561 non-null float64
24 rent_samples 11561 non-null float64
25 rent_gt_10 11560 non-null float64
26 rent_gt_15 11560 non-null float64
27 rent_gt_20 11560 non-null float64
28 rent_gt_25 11560 non-null float64
29 rent_gt_30 11560 non-null float64
30 rent_gt_35 11560 non-null float64
31 rent_gt_40 11560 non-null float64
32 rent_gt_50 11560 non-null float64
33 universe_samples 11709 non-null int64
34 used_samples 11709 non-null int64
35 hi_mean 11587 non-null float64
36 hi_median 11587 non-null float64
37 hi_stdev 11587 non-null float64
38 hi_sample_weight 11587 non-null float64
39 hi_samples 11587 non-null float64
40 family_mean 11573 non-null float64
41 family_median 11573 non-null float64
42 family_stdev 11573 non-null float64
43 family_sample_weight 11573 non-null float64
44 family_samples 11573 non-null float64
45 hc_mortgage_mean 11441 non-null float64
46 hc_mortgage_median 11441 non-null float64
47 hc_mortgage_stdev 11441 non-null float64
48 hc_mortgage_sample_weight 11441 non-null float64
49 hc_mortgage_samples 11441 non-null float64
50 hc_mean 11419 non-null float64
51 hc_median 11419 non-null float64
52 hc_stdev 11419 non-null float64
53 hc_samples 11419 non-null float64
54 hc_sample_weight 11419 non-null float64
55 home_equity_second_mortgage 11489 non-null float64
56 second_mortgage 11489 non-null float64
57 home_equity 11489 non-null float64
58 debt 11489 non-null float64
59 second_mortgage_cdf 11489 non-null float64
60 home_equity_cdf 11489 non-null float64
61 debt_cdf 11489 non-null float64
62 hs_degree 11624 non-null float64
63 hs_degree_male 11620 non-null float64
64 hs_degree_female 11604 non-null float64
65 male_age_mean 11625 non-null float64
66 male_age_median 11625 non-null float64
67 male_age_stdev 11625 non-null float64
68 male_age_sample_weight 11625 non-null float64
69 male_age_samples 11625 non-null float64
70 female_age_mean 11613 non-null float64
71 female_age_median 11613 non-null float64
72 female_age_stdev 11613 non-null float64
73 female_age_sample_weight 11613 non-null float64
74 female_age_samples 11613 non-null float64
75 pct_own 11587 non-null float64
76 married 11625 non-null float64
77 married_snp 11625 non-null float64
78 separated 11625 non-null float64
79 divorced 11625 non-null float64
dtypes: float64(61), int64(13), object(6)
memory usage: 7.1+ MB

```

## 2. Figure out the primary key and look for the requirement of indexing

```

In [14]: #UID is unique userID value in the train and test dataset. So an index can be created from the UID feature
df_train.set_index(keys=['UID'],inplace=True)#Set the DataFrame index using existing columns.
df_test.set_index(keys=['UID'],inplace=True)

```

```

In [15]: df_train.head(2)

```

```

Out[15]:

```

	BLOCKID	SUMLEVEL	COUNTYID	STATEID	state	state_ab	city	place	type	primary	...	female_age_mean	female_age_median	female_age_stdev	fe
UID															
267822	NaN	140	53	36	New York	NY	Hamilton	Hamilton	City	tract	...	44.48629	45.33333	22.51276	
246444	NaN	140	141	18	Indiana	IN	South Bend	Roseland	City	tract	...	36.48391	37.58333	23.43353	

2 rows × 79 columns

```

In [16]: df_test.head(2)

```

Out[16]:

	BLOCKID	SUMLEVEL	COUNTYID	STATEID	state	state_ab	city	place	type	primary	...	female_age_mean	female_age_median	female_age_stdev	fi
	UID														
	255504	NaN	140	163	26	Michigan	MI	Detroit	Dearborn Heights City	CDP	tract	...	34.78682	33.75000	21.58531
	252676	NaN	140	1	23	Maine	ME	Auburn	Auburn City	City	tract	...	44.23451	46.66667	22.37036

2 rows × 79 columns

3. Gauge the fill rate of the variables and devise plans for missing value treatment. Please explain explicitly the reason for the treatment chosen for each variable.

In [17]:

```
#percentage of missing values in train set
missing_list_train=df_train.isnull().sum()*100/len(df_train)
missing_values_df_train=pd.DataFrame(missing_list_train,columns=['Percentage of missing values'])
missing_values_df_train.sort_values(by=['Percentage of missing values'],inplace=True,ascending=False)
missing_values_df_train[missing_values_df_train['Percentage of missing values'] >0][:10]
#BLOCKID can be dropped, since it is 100%missing values
```

Out[17]:

	Percentage of missing values
BLOCKID	100.000000
hc_samples	2.196113
hc_mean	2.196113
hc_median	2.196113
hc_stdev	2.196113
hc_sample_weight	2.196113
hc_mortgage_mean	2.097288
hc_mortgage_stdev	2.097288
hc_mortgage_sample_weight	2.097288
hc_mortgage_samples	2.097288

In [18]:

```
#percentage of missing values in test set
missing_list_test=df_test.isnull().sum()*100/len(df_train)
missing_values_df_test=pd.DataFrame(missing_list_test,columns=['Percentage of missing values'])
missing_values_df_test.sort_values(by=['Percentage of missing values'],inplace=True,ascending=False)
missing_values_df_test[missing_values_df_test['Percentage of missing values'] >0][:10]
#BLOCKID can be dropped, since it is 43%missing values
```

Out[18]:

	Percentage of missing values
BLOCKID	42.857143
hc_samples	1.061455
hc_mean	1.061455
hc_median	1.061455
hc_stdev	1.061455
hc_sample_weight	1.061455
hc_mortgage_mean	0.980930
hc_mortgage_stdev	0.980930
hc_mortgage_sample_weight	0.980930
hc_mortgage_samples	0.980930

In [19]:

```
df_train .drop(columns=['BLOCKID','SUMLEVEL'],inplace=True) #SUMLEVEL and BLOCKID does not seem to contribute
```

In [20]:

```
df_test .drop(columns=['BLOCKID','SUMLEVEL'],inplace=True)
```

In [21]:

```
# Impute with mean
missing_train_cols=[]
for col in df_train.columns:
    if df_train[col].isna().sum() !=0:
        missing_train_cols.append(col)
print(missing_train_cols)
```

['rent\_mean', 'rent\_median', 'rent\_stdev', 'rent\_sample\_weight', 'rent\_samples', 'rent\_gt\_10', 'rent\_gt\_15', 'rent\_gt\_20', 'rent\_gt\_25', 'rent\_gt\_30', 'rent\_gt\_35', 'rent\_gt\_40', 'rent\_gt\_50', 'hi\_mean', 'hi\_median', 'hi\_stdev', 'hi\_sample\_weight', 'hi\_samples', 'family\_mean', 'family\_median', 'family\_stdev', 'family\_sample\_weight', 'family\_samples', 'hc\_mortgage\_mean', 'hc\_mortgage\_median', 'hc\_mortgage\_stdev', 'hc\_mortgage\_sample\_weight', 'hc\_mortgage\_samples', 'hc\_mean', 'hc\_median', 'hc\_stdev', 'hc\_samples', 'hc\_sample\_weight', 'home\_equity\_second\_mortgage', 'second\_mortgage', 'home\_equity', 'debt', 'second\_mortgage\_cdf', 'home\_equity\_cdf', 'debt\_cdf', 'hs\_degree', 'hs\_degree\_male', 'hs\_degree\_female', 'male\_age\_mean', 'male\_age\_median', 'male\_age\_stdev', 'male\_age\_sample\_weight', 'male\_age\_samples', 'female\_age\_mean', 'female\_age\_median', 'female\_age\_stdev', 'female\_age\_sample\_weight', 'female\_age\_samples', 'pct\_own', 'married', 'married\_snp', 'separated', 'divorced']

```
In [22]: # Impute with mean
missing_test_cols=[]
for col in df_test.columns:
    if df_test[col].isna().sum() !=0:
        missing_test_cols.append(col)
print(missing_test_cols)

['rent_mean', 'rent_median', 'rent_stddev', 'rent_sample_weight', 'rent_samples', 'rent_gt_10', 'rent_gt_15', 'rent_gt_20', 'rent_gt_25', 'rent_gt_30', 'rent_gt_35', 'rent_gt_40', 'rent_gt_50', 'hi_mean', 'hi_median', 'hi_stddev', 'hi_sample_weight', 'hi_samples', 'family_mean', 'family_median', 'family_stddev', 'family_sample_weight', 'family_samples', 'hc_mortgage_mean', 'hc_mortgage_median', 'hc_mortgage_stddev', 'hc_mortgage_sample_weight', 'hc_mortgage_samples', 'hc_mean', 'hc_median', 'hc_stddev', 'hc_samples', 'hc_sample_weight', 'home_equity_second_mortgage', 'second_mortgage', 'home_equity', 'debt', 'second_mortgage_cdf', 'home_equity_cdf', 'debt_cdf', 'hs_degree', 'hs_degree_male', 'hs_degree_female', 'male_age_mean', 'male_age_median', 'male_age_stddev', 'male_age_sample_weight', 'male_age_samples', 'female_age_mean', 'female_age_median', 'female_age_stddev', 'female_age_sample_weight', 'female_age_samples', 'pct_own', 'married', 'married_snp', 'separated', 'divorced']
```

```
In [23]: for col in df_train.columns:
        if col in (missing_train_cols):
            df_train[col].replace(np.nan, df_train[col].mean(),inplace=True)
```

```
In [24]: for col in df_test.columns:
        if col in (missing_test_cols):
            df_test[col].replace(np.nan, df_test[col].mean(),inplace=True)
```

```
In [25]: df_train.isna().sum().sum()
```

Out[25]: 0

```
In [26]: df_test.isna().sum().sum()
```

Out[26]: 0

4.Understanding homeowner costs are incredibly valuable because it is positively correlated to consumer spending which drives the economy through disposable income. Perform debt analysis:

a) Explore the top 2,500 locations where the percentage of households with a second mortgage is the highest and percent ownership is above 10 percent. Visualize using geo-map. You may keep the upper limit for the percent of households with a second mortgage to 50 percent

```
In [27]: from pandasql import sqldf
q1 = "select place,pct_own,second_mortgage,lat,lng from df_train where pct_own >0.10 and second_mortgage <0.5 order by second_mortgage DESC LI
pysqldf = lambda q: sqldf(q, globals())
df_train_location_mort_pct=pysqldf(q1)
```

```
In [28]: df_train_location_mort_pct.head()
```

Out[28]:

	place	pct_own	second_mortgage	lat	lng
0	Worcester City	0.20247	0.43363	42.254262	-71.800347
1	Harbor Hills	0.15618	0.31818	40.751809	-73.853582
2	Glen Burnie	0.22380	0.30212	39.127273	-76.635265
3	Egypt Lake-Ieto	0.11618	0.28972	28.029063	-82.495395
4	Lincolnwood	0.14228	0.28899	41.967289	-87.652434

```
In [29]: import plotly.express as px
import plotly.graph_objects as go
```

```
In [30]: fig = go.Figure(data=go.Scattergeo(
    lat = df_train_location_mort_pct['lat'],
    lon = df_train_location_mort_pct['lng'],
))
fig.update_layout(
    geo=dict(
        scope = 'north america',
        showland = True,
        landcolor = "rgb(212, 212, 212)",
        subunitcolor = "rgb(255, 255, 255)",
        countrycolor = "rgb(255, 255, 255)",
        showlakes = True,
        lakecolor = "rgb(255, 255, 255)",
        showsubunits = True,
        showcountries = True,
        resolution = 50,
        projection = dict(
            type = 'conic conformal',
            rotation_lon = -100
        ),
        lonaxis = dict(
            showgrid = True,
            gridwidth = 0.5,
            range= [ -140.0, -55.0 ],
            dtick = 5
        ),
        lataxis = dict (
```

```

        showgrid = True,
        gridwidth = 0.5,
        range= [ 20.0, 60.0 ],
        dtick = 5
    )
),
title='Top 2,500 locations with second mortgage is the highest and percent ownership is above 10 percent')
fig.show()

```

```

In [31]: df_train['bad_debt']=df_train['second_mortgage']+df_train['home_equity']-df_train['home_equity_second_mortgage']

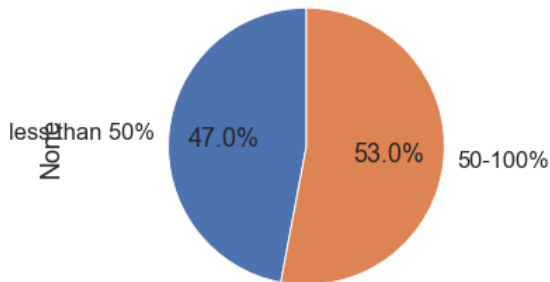
```

```

In [32]: df_train['bins'] = pd.cut(df_train['bad_debt'],bins=[0,0.10,1], labels=["less than 50%","50-100%"])
df_train.groupby(['bins']).size().plot(kind='pie',subplots=True,startangle=90, autopct='%1.1f%%')
plt.axis('equal')

plt.show()
#df.plot.pie(subplots=True,figsize=(8, 3))

```



```

In [33]: cols=[]
df_train.columns

```

```

Out[33]: Index(['COUNTYID', 'STATEID', 'state', 'state_ab', 'city', 'place', 'type',
        'primary', 'zip_code', 'area_code', 'lat', 'lng', 'ALand', 'AWater',
        'pop', 'male_pop', 'female_pop', 'rent_mean', 'rent_median',
        'rent_stdev', 'rent_sample_weight', 'rent_samples', 'rent_gt_10',
        'rent_gt_15', 'rent_gt_20', 'rent_gt_25', 'rent_gt_30', 'rent_gt_35',
        'rent_gt_40', 'rent_gt_50', 'universe_samples', 'used_samples',
        'hi_mean', 'hi_median', 'hi_stdev', 'hi_sample_weight', 'hi_samples',
        'family_mean', 'family_median', 'family_stdev', 'family_sample_weight',
        'family_samples', 'hc_mortgage_mean', 'hc_mortgage_median',
        'hc_mortgage_stdev', 'hc_mortgage_sample_weight', 'hc_mortgage_samples',
        'hc_mean', 'hc_median', 'hc_stdev', 'hc_samples', 'hc_sample_weight',
        'home_equity_second_mortgage', 'second_mortgage', 'home_equity', 'debt',
        'second_mortgage_cdf', 'home_equity_cdf', 'debt_cdf', 'hs_degree',
        'hs_degree_male', 'hs_degree_female', 'male_age_mean',
        'male_age_median', 'male_age_stdev', 'male_age_sample_weight',
        'male_age_samples', 'female_age_mean', 'female_age_median',
        'female_age_stdev', 'female_age_sample_weight', 'female_age_samples',
        'pct_own', 'married', 'married_snp', 'separated', 'divorced',
        'bad_debt', 'bins'],
        dtype='object')

```

```

In [34]: #Taking Hamilton and Manhattan cities data
cols=['second_mortgage','home_equity','debt','bad_debt']
df_box_hamilton=df_train.loc[df_train['city'] == 'Hamilton']
df_box_manhattan=df_train.loc[df_train['city'] == 'Manhattan']

```

```
df_box_city=pd.concat([df_box_hamilton,df_box_manhattan])
df_box_city.head(4)
```

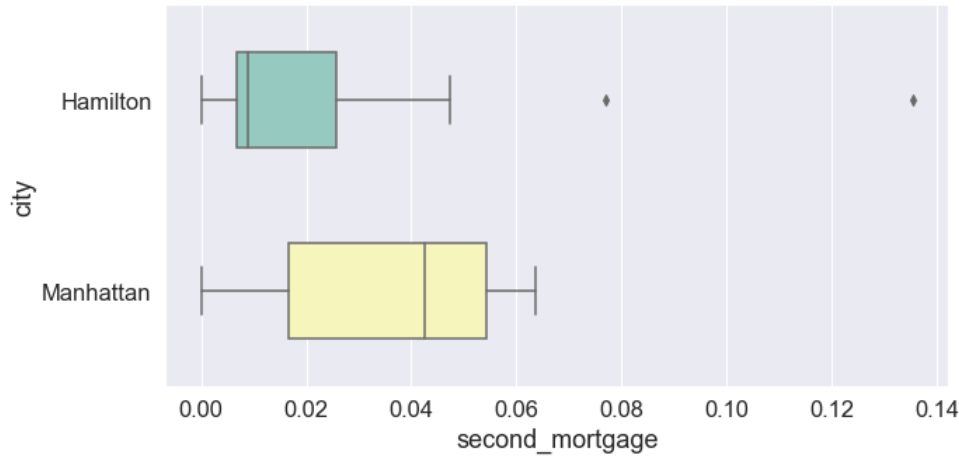
Out[34]:

	COUNTYID	STATEID	state	state_ab	city	place	type	primary	zip_code	area_code	...	female_age_stdev	female_age_sample_weight	female_ag
	UID													
	267822	53	36	New York	NY	Hamilton	Hamilton	City	tract	13346	315	...	22.51276	685.33845
	263797	21	34	New Jersey	NJ	Hamilton	Yardville	City	tract	8610	609	...	24.05831	732.58443
	270979	17	39	Ohio	OH	Hamilton	Hamilton City	Village	tract	45015	513	...	22.66500	565.32725
	259028	95	28	Mississippi	MS	Hamilton	Hamilton	CDP	tract	39746	662	...	22.79602	483.01311

4 rows × 79 columns

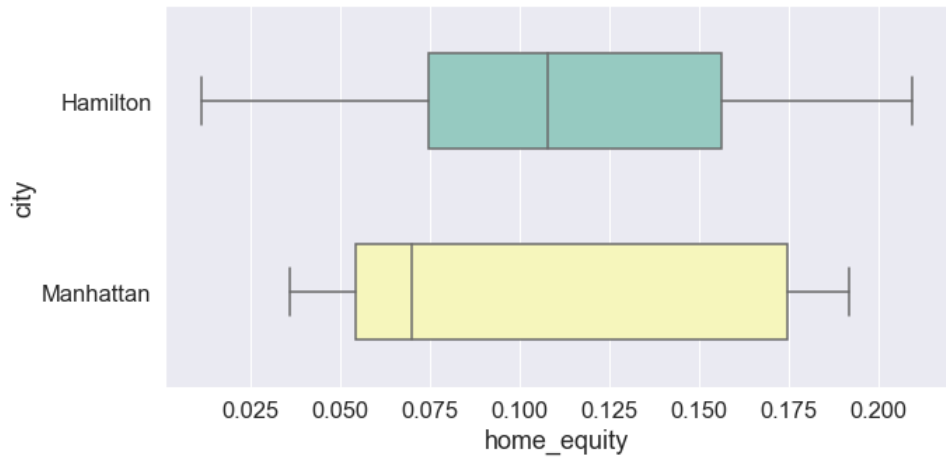
In [35]:

```
plt.figure(figsize=(10,5))
sns.boxplot(data=df_box_city,x='second_mortgage', y='city',width=0.5,palette="Set3")
plt.show()
```



In [36]:

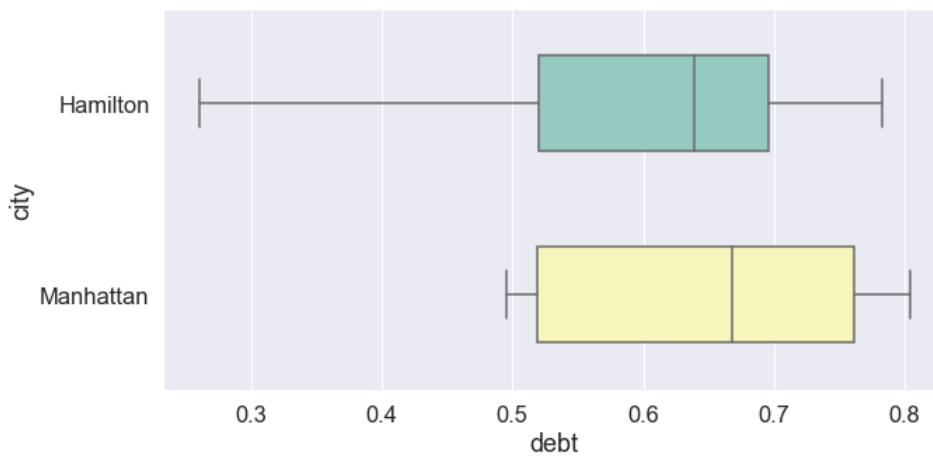
```
plt.figure(figsize=(10,5))
sns.boxplot(data=df_box_city,x='home_equity', y='city',width=0.5,palette="Set3")
plt.show()
```



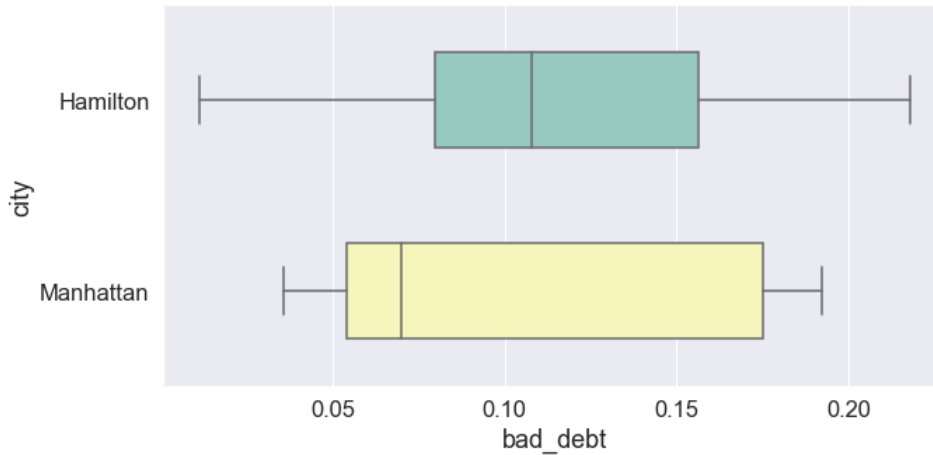
In [37]:

```
plt.figure(figsize=(10,5))
sns.boxplot(data=df_box_city,x='debt', y='city',width=0.5,palette="Set3")
plt.show()
```





```
In [38]: plt.figure(figsize=(10,5))
sns.boxplot(data=df_box_city,x='bad_debt', y='city',width=0.5,palette="Set3")
plt.show()
```

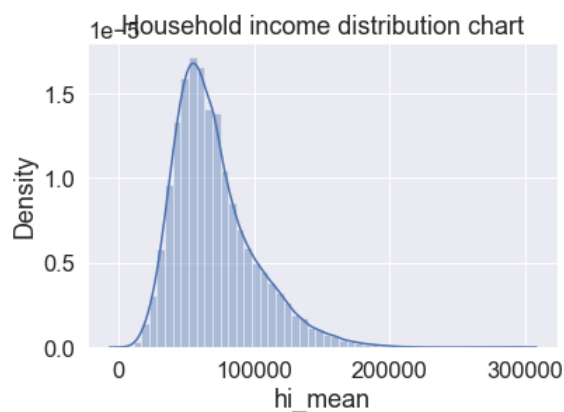


Manhattan has higher metrics compared to Hamilton

```
In [39]: sns.distplot(df_train['hi_mean'])
plt.title('Household income distribution chart')
plt.show()
```

E:\Programs\Anaconda\lib\site-packages\seaborn\distributions.py:2619: FutureWarning:

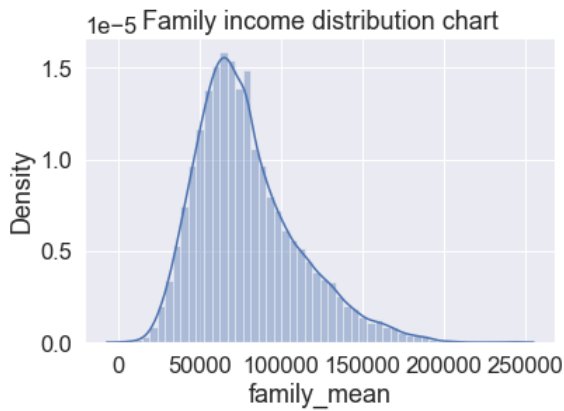
`distplot` is a deprecated function and will be removed in a future version. Please adapt your code to use either `displot` (a figure-level function with similar flexibility) or `histplot` (an axes-level function for histograms).



```
In [40]: sns.distplot(df_train['family_mean'])
plt.title('Family income distribution chart')
plt.show()
```

E:\Programs\Anaconda\lib\site-packages\seaborn\distributions.py:2619: FutureWarning:

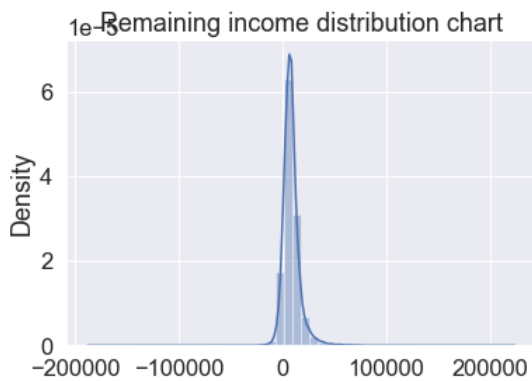
`distplot` is a deprecated function and will be removed in a future version. Please adapt your code to use either `displot` (a figure-level function with similar flexibility) or `histplot` (an axes-level function for histograms).



```
In [41]: sns.distplot(df_train['family_mean']-df_train['hi_mean'])
plt.title('Remaining income distribution chart')
plt.show()
```

E:\Programs\Anaconda\lib\site-packages\seaborn\distributions.py:2619: FutureWarning:

`distplot` is a deprecated function and will be removed in a future version. Please adapt your code to use either `displot` (a figure-level function with similar flexibility) or `histplot` (an axes-level function for histograms).



Income distribution almost has normality in its distribution

**Perform EDA and come out with insights into population density and age. You may have to derive new fields (make sure to weight averages for accurate measurements):**

```
In [42]: #plt.figure(figsize=(25,10))
fig,(ax1,ax2,ax3)=plt.subplots(3,1)
sns.distplot(df_train['pop'],ax=ax1)
sns.distplot(df_train['male_pop'],ax=ax2)
sns.distplot(df_train['female_pop'],ax=ax3)
plt.subplots_adjust(wspace=0.8,hspace=0.8)
plt.tight_layout()
plt.show()
```

E:\Programs\Anaconda\lib\site-packages\seaborn\distributions.py:2619: FutureWarning:

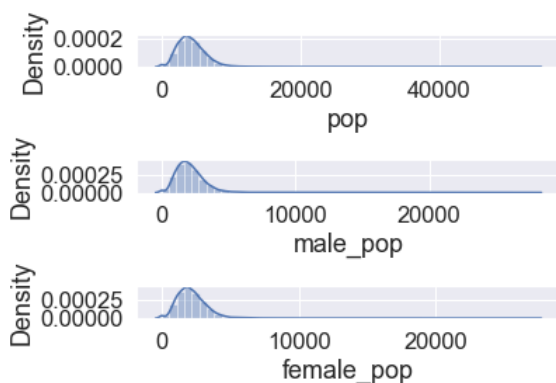
`distplot` is a deprecated function and will be removed in a future version. Please adapt your code to use either `displot` (a figure-level function with similar flexibility) or `histplot` (an axes-level function for histograms).

E:\Programs\Anaconda\lib\site-packages\seaborn\distributions.py:2619: FutureWarning:

`distplot` is a deprecated function and will be removed in a future version. Please adapt your code to use either `displot` (a figure-level function with similar flexibility) or `histplot` (an axes-level function for histograms).

E:\Programs\Anaconda\lib\site-packages\seaborn\distributions.py:2619: FutureWarning:

`distplot` is a deprecated function and will be removed in a future version. Please adapt your code to use either `displot` (a figure-level function with similar flexibility) or `histplot` (an axes-level function for histograms).



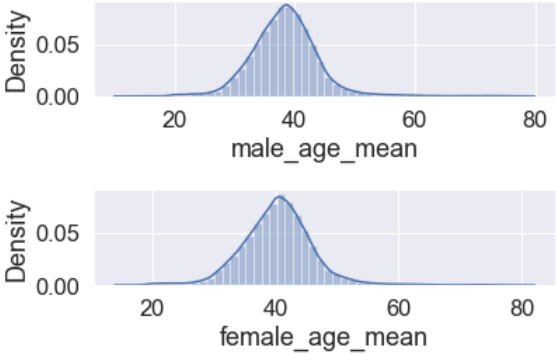
```
In [43]: #plt.figure(figsize=(25,10))
fig,(ax1,ax2)=plt.subplots(2,1)
sns.distplot(df_train['male_age_mean'],ax=ax1)
sns.distplot(df_train['female_age_mean'],ax=ax2)
plt.subplots_adjust(wspace=0.8,hspace=0.8)
plt.tight_layout()
plt.show()
```

E:\Programs\Anaconda\lib\site-packages\seaborn\distributions.py:2619: FutureWarning:

`distplot` is a deprecated function and will be removed in a future version. Please adapt your code to use either `displot` (a figure-level function with similar flexibility) or `histplot` (an axes-level function for histograms).

E:\Programs\Anaconda\lib\site-packages\seaborn\distributions.py:2619: FutureWarning:

`distplot` is a deprecated function and will be removed in a future version. Please adapt your code to use either `displot` (a figure-level function with similar flexibility) or `histplot` (an axes-level function for histograms).



a) Use pop and ALand variables to create a new field called population density

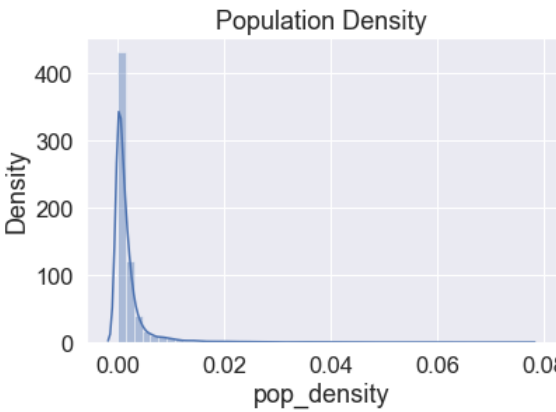
```
In [44]: df_train['pop_density']=df_train['pop']/df_train['ALand']

In [45]: df_test['pop_density']=df_test['pop']/df_test['ALand']

In [46]: sns.distplot(df_train['pop_density'])
plt.title('Population Density')
plt.show() # Very Less density is noticed
```

E:\Programs\Anaconda\lib\site-packages\seaborn\distributions.py:2619: FutureWarning:

`distplot` is a deprecated function and will be removed in a future version. Please adapt your code to use either `displot` (a figure-level function with similar flexibility) or `histplot` (an axes-level function for histograms).



Use male\_age\_median, female\_age\_median, male\_pop, and female\_pop to create a new field called median age c) Visualize the findings using appropriate chart type

```
In [47]: df_train['age_median']=(df_train['male_age_median']+df_train['female_age_median'])/2
df_test['age_median']=(df_test['male_age_median']+df_test['female_age_median'])/2

In [48]: df_train[['male_age_median', 'female_age_median', 'male_pop', 'female_pop', 'age_median']].head()
```

Out[48]:

	male_age_median	female_age_median	male_pop	female_pop	age_median
UID					
267822	44.00000	45.33333	2612	2618	44.666665
246444	32.00000	37.58333	1349	1284	34.791665
245683	40.83333	42.83333	3643	3238	41.833330
279653	48.91667	50.58333	1141	1559	49.750000

	male_age_median	female_age_median	male_pop	female_pop	age_median
--	-----------------	-------------------	----------	------------	------------

UID					
247218	22.41667	21.58333	2586	3051	22.000000

```
In [49]: sns.distplot(df_train['age_median'])
plt.title('Median Age')
plt.show()
# Age of population is mostly between 20 and 60
# Majority are of age around 40
# Median age distribution has a gaussian distribution
# Some right skewness is noticed
```

E:\Programs\Anaconda\lib\site-packages\seaborn\distributions.py:2619: FutureWarning:

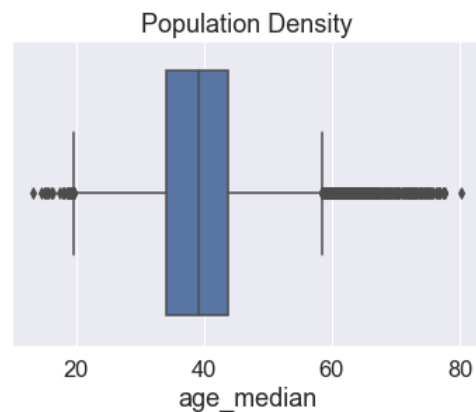
`distplot` is a deprecated function and will be removed in a future version. Please adapt your code to use either `displot` (a figure-level function with similar flexibility) or `histplot` (an axes-level function for histograms).



```
In [50]: sns.boxplot(df_train['age_median'])
plt.title('Population Density')
plt.show()
```

E:\Programs\Anaconda\lib\site-packages\seaborn\\_decorators.py:36: FutureWarning:

Pass the following variable as a keyword arg: x. From version 0.12, the only valid positional argument will be `data`, and passing other arguments without an explicit keyword will result in an error or misinterpretation.



Create bins for population into a new variable by selecting appropriate class interval so that the number of categories don't exceed 5 for the ease of analysis.

```
In [51]: df_train['pop'].describe()
```

```
Out[51]: count    27321.000000
mean      4316.032685
std       2169.226173
min        0.000000
25%      2885.000000
50%      4042.000000
75%      5430.000000
max      53812.000000
Name: pop, dtype: float64
```

```
In [52]: df_train['pop_bins']=pd.cut(df_train['pop'],bins=5,labels=['very low','low','medium','high','very high'])
```

```
In [53]: df_train[['pop', 'pop_bins']]
```

```
Out[53]:
```

	pop	pop_bins
--	-----	----------

UID		
267822	5230	very low

	pop	pop_bins
UID		
246444	2633	very low
245683	6881	very low
279653	2700	very low
247218	5637	very low
...	...	...
279212	1847	very low
277856	4155	very low
233000	2829	very low
287425	11542	low
265371	3726	very low

27321 rows × 2 columns

```
In [54]: df_train['pop_bins'].value_counts()
```

```
Out[54]: very low    27058
low         246
medium      9
high        7
very high   1
Name: pop_bins, dtype: int64
```

Analyze the married, separated, and divorced population for these population brackets

```
In [55]: df_train.groupby(by='pop_bins')[['married', 'separated', 'divorced']].count()
```

```
Out[55]:
```

	married	separated	divorced
pop_bins			
very low	27058	27058	27058
low	246	246	246
medium	9	9	9
high	7	7	7
very high	1	1	1

```
In [56]: df_train.groupby(by='pop_bins')[['married', 'separated', 'divorced']].agg(["mean", "median"])
```

```
Out[56]:
```

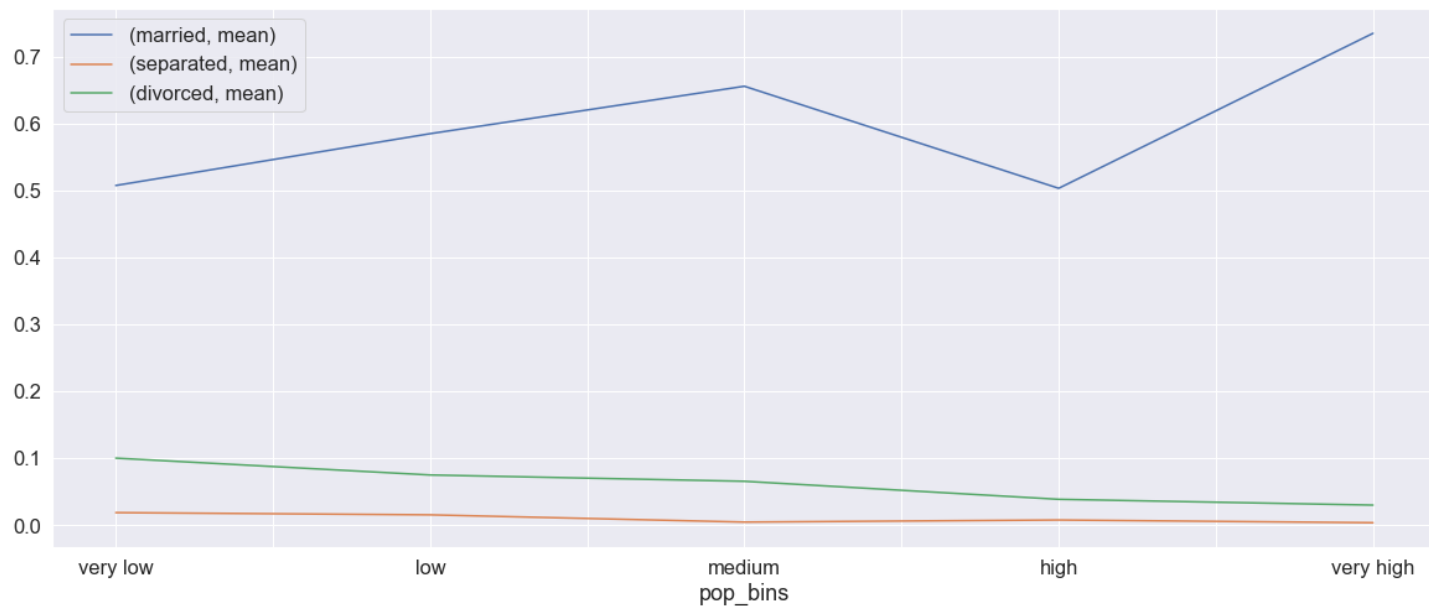
	married		separated		divorced	
	mean	median	mean	median	mean	median
pop_bins						
very low	0.507548	0.524680	0.019126	0.013650	0.100504	0.096020
low	0.584894	0.593135	0.015833	0.011195	0.075348	0.070045
medium	0.655737	0.618710	0.005003	0.004120	0.065927	0.064890
high	0.503359	0.335660	0.008141	0.002500	0.039030	0.010320
very high	0.734740	0.734740	0.004050	0.004050	0.030360	0.030360

- 1. Very high population group has more married people and less percentage of separated and divorced couples
- 2. In very low population groups, there are more divorced people

Visualize using appropriate chart type

```
In [57]: plt.figure(figsize=(10,5))
pop_bin_married=df_train.groupby(by='pop_bins')[['married', 'separated', 'divorced']].agg(["mean"])
pop_bin_married.plot(figsize=(20,8))
plt.legend(loc='best')
plt.show()
```

<Figure size 720x360 with 0 Axes>



Please detail your observations for rent as a percentage of income at an overall level, and for different states.

```
In [58]: rent_state_mean=df_train.groupby(by='state')['rent_mean'].agg(["mean"])
rent_state_mean.head()
```

```
Out[58]:
```

	mean
state	
Alabama	774.004927
Alaska	1185.763570
Arizona	1097.753511
Arkansas	720.918575
California	1471.133857

```
In [59]: income_state_mean=df_train.groupby(by='state')['family_mean'].agg(["mean"])
income_state_mean.head()
```

```
Out[59]:
```

	mean
state	
Alabama	67030.064213
Alaska	92136.545109
Arizona	73328.238798
Arkansas	64765.377850
California	87655.470820

```
In [60]: rent_perc_of_income=rent_state_mean['mean']/income_state_mean['mean']
rent_perc_of_income.head(10)
```

```
Out[60]:
```

state	
Alabama	0.011547
Alaska	0.012870
Arizona	0.014970
Arkansas	0.011131
California	0.016783
Colorado	0.013529
Connecticut	0.012637
Delaware	0.012929
District of Columbia	0.013198
Florida	0.015772

Name: mean, dtype: float64

```
In [61]: #overall level rent as a percentage of income
sum(df_train['rent_mean']/sum(df_train['family_mean']))
```

```
Out[61]: 0.013358170721473864
```

Perform correlation analysis for all the relevant variables by creating a heatmap. Describe your findings.

```
In [62]: df_train.columns
```

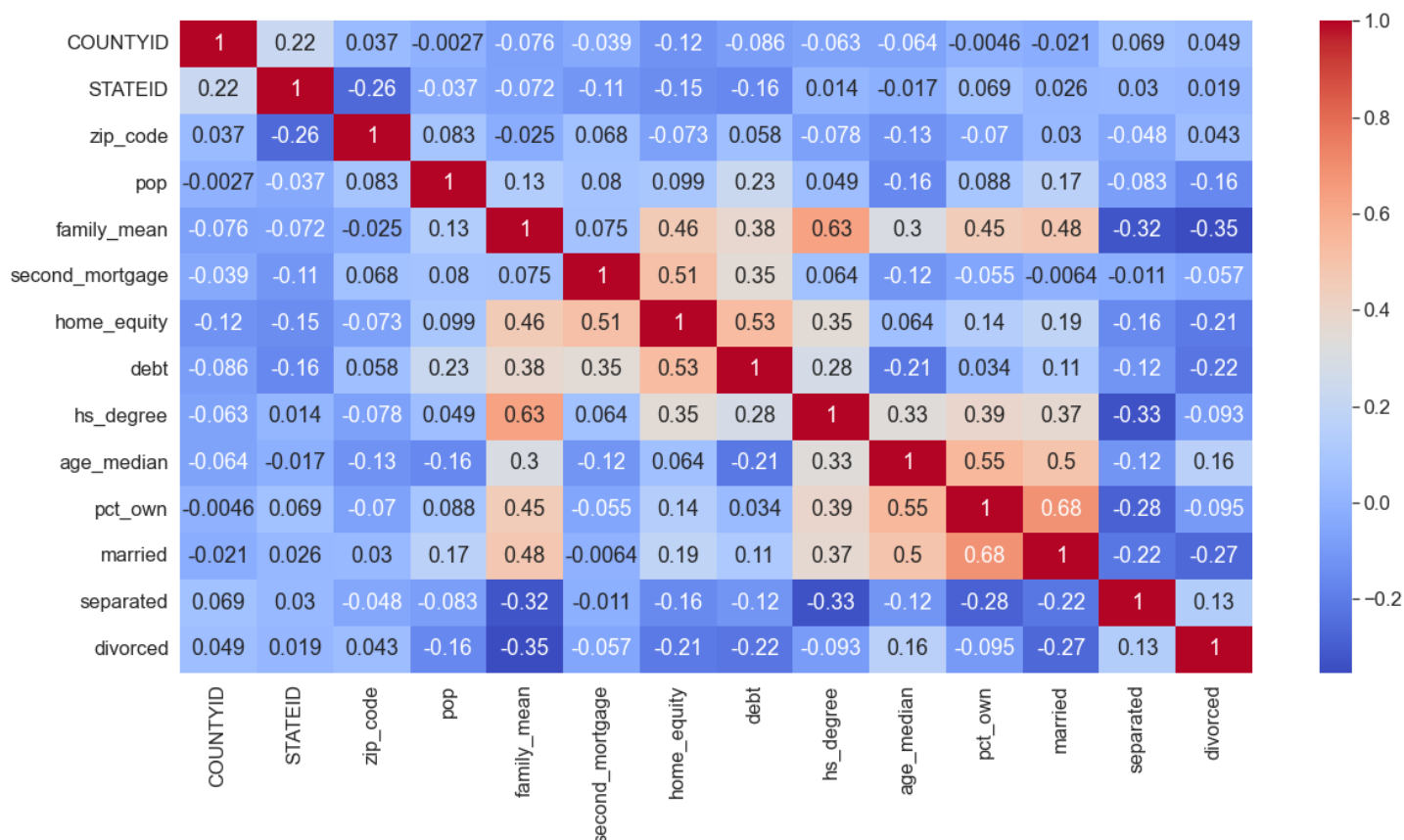
```
Out[62]:
```

Index(['COUNTYID', 'STATEID', 'state', 'state\_ab', 'city', 'place', 'type', 'primary', 'zip\_code', 'area\_code', 'lat', 'lng', 'ALand', 'AWater', 'pop', 'male\_pop', 'female\_pop', 'rent\_mean', 'rent\_median', 'rent\_stdev', 'rent\_sample\_weight', 'rent\_samples', 'rent\_gt\_10',

```
'rent_gt_15', 'rent_gt_20', 'rent_gt_25', 'rent_gt_30', 'rent_gt_35',
'rent_gt_40', 'rent_gt_50', 'universe_samples', 'used_samples',
'hi_mean', 'hi_median', 'hi_stdev', 'hi_sample_weight', 'hi_samples',
'family_mean', 'family_median', 'family_stdev', 'family_sample_weight',
'family_samples', 'hc_mortgage_mean', 'hc_mortgage_median',
'hc_mortgage_stdev', 'hc_mortgage_sample_weight', 'hc_mortgage_samples',
'hc_mean', 'hc_median', 'hc_stdev', 'hc_samples', 'hc_sample_weight',
'home_equity_second_mortgage', 'second_mortgage', 'home_equity', 'debt',
'second_mortgage_cdf', 'home_equity_cdf', 'debt_cdf', 'hs_degree',
'hs_degree_male', 'hs_degree_female', 'male_age_mean',
'male_age_median', 'male_age_stdev', 'male_age_sample_weight',
'male_age_samples', 'female_age_mean', 'female_age_median',
'female_age_stdev', 'female_age_sample_weight', 'female_age_samples',
'pct_own', 'married', 'married_snp', 'separated', 'divorced',
'bad_debt', 'bins', 'pop_density', 'age_median', 'pop_bins'],
dtype='object')
```

```
In [63]: cor=df_train[['COUNTYID','STATEID','zip_code','type','pop', 'family_mean',
'second_mortgage', 'home_equity', 'debt','hs_degree',
'age_median','pct_own', 'married','separated', 'divorced']].corr()
```

```
In [64]: plt.figure(figsize=(20,10))
sns.heatmap(cor,annot=True,cmap='coolwarm')
plt.show()
```



1. High positive correaltion is noticed between pop, male\_pop and female\_pop
2. High positive correaltion is noticed between rent\_mean,hi\_mean, family\_mean,hc\_mean

1. The economic multivariate data has a significant number of measured variables. The goal is to find where the measured variables depend on a number of smaller unobserved common factors or latent variables. 2. Each variable is assumed to be dependent upon a linear combination of the common factors, and the coefficients are known as loadings. Each measured variable also includes a component due to independent random variability, known as “specific variance” because it is specific to one variable. Obtain the common factors and then plot the loadings. Use factor analysis to find latent variables in our dataset and gain insight into the linear relationships in the data. Following are the list of latent variables:

• Highschool graduation rates • Median population age • Second mortgage statistics • Percent own • Bad debt expense

```
In [65]: from sklearn.decomposition import FactorAnalysis
from factor_analyzer import FactorAnalyzer
```

```
In [66]: fa=FactorAnalyzer(n_factors=5)
fa.fit_transform(df_train.select_dtypes(exclude= ('object', 'category')))
fa.loadings_
```

```
Out[66]: array([[ -1.12589166e-01,  1.95646468e-02, -2.39331083e-02,
-6.27632623e-02,  4.23474749e-02],
[ -1.10186763e-01,  1.33506215e-02,  2.79651243e-02,
-1.49825864e-01,  1.10838807e-01],
[ -8.28678646e-02,  5.16372377e-02, -1.36451871e-01,
-4.98918634e-02, -1.04024841e-01],
[ 1.80961149e-02,  1.92013753e-02,  5.81329827e-03,
```

2.64842740e-02, -6.12442486e-03],  
[ 9.02324715e-02, -9.72544297e-02, -6.54601315e-02,  
-1.33145899e-01, -1.48594601e-01],  
[-1.07335697e-02, -4.12376818e-02, 1.45853484e-01,  
8.80433327e-03, 1.08227565e-01],  
[-4.28796971e-02, -2.09780214e-02, 3.66726851e-02,  
-9.45597383e-02, 5.91380520e-02],  
[-2.44243003e-03, -1.53245409e-02, -2.68300902e-03,  
-4.52473044e-02, 2.37240659e-02],  
[ 7.92164339e-02, 9.57453331e-01, -8.71151642e-02,  
-6.59923845e-03, -3.97273184e-02],  
[ 7.39808211e-02, 9.18750524e-01, -1.08834840e-01,  
-2.79371590e-02, -3.93153640e-02],  
[ 8.06598896e-02, 9.47839220e-01, -6.08006509e-02,  
1.53627095e-02, -3.86977277e-02],  
[ 7.70052137e-01, 9.84675329e-03, -3.71249754e-02,  
1.14949046e-01, -1.23784684e-01],  
[ 7.18615881e-01, 6.24980464e-03, -4.59787407e-02,  
1.09109689e-01, -1.35301911e-01],  
[ 7.07647246e-01, 2.46625399e-02, -1.00860846e-02,  
1.04472488e-01, 7.72381251e-02],  
[-1.34545492e-01, 3.36809297e-01, -4.87894959e-01,  
-4.15446166e-02, 3.17608532e-01],  
[ 2.31079697e-01, 4.37729787e-01, -6.40209196e-01,  
-2.52310925e-02, 3.47216216e-01],  
[-4.52068133e-02, 3.51263844e-02, 3.07537041e-02,  
4.44793508e-01, -1.63273411e-01],  
[-2.50717066e-02, 1.70166796e-02, 4.57227280e-02,  
6.76083904e-01, -1.55256767e-01],  
[-3.90694439e-02, -1.67460866e-02, 8.13962673e-02,  
8.36389105e-01, -9.18259792e-02],  
[-5.14161936e-02, -3.57207135e-02, 1.10795184e-01,  
9.25123775e-01, -4.44866508e-02],  
[-6.08589973e-02, -4.41860610e-02, 1.35794031e-01,  
9.53019931e-01, -2.21548653e-02],  
[-4.57771154e-02, -5.25526111e-02, 1.41019868e-01,  
9.32702618e-01, -5.83072683e-07],  
[-4.19486050e-02, -5.90387622e-02, 1.28851766e-01,  
8.87316645e-01, 1.05894326e-02],  
[-2.47894627e-02, -7.29670546e-02, 9.41510444e-02,  
7.79023669e-01, 2.95352834e-02],  
[ 2.12258459e-01, 4.65992346e-01, -6.14495951e-01,  
-2.47660018e-02, 3.66644539e-01],  
[ 2.33057249e-01, 4.47057850e-01, -6.28263427e-01,  
-2.71547710e-02, 3.43419624e-01],  
[ 7.85157098e-01, 4.91249258e-02, 1.44540484e-01,  
-2.05217633e-01, -1.54523363e-01],  
[ 7.10324880e-01, 4.99730440e-02, 1.32239990e-01,  
-2.19171864e-01, -2.10505574e-01],  
[ 8.61780953e-01, 4.35044836e-02, 1.65839099e-01,  
-1.19850816e-01, 3.16733610e-02],  
[-2.23443274e-01, 8.46259549e-01, -4.61177183e-02,  
6.88599272e-02, 2.27742314e-01],  
[ 1.43837557e-01, 9.53197424e-01, 2.27887469e-02,  
-4.57890445e-02, 1.00796449e-01],  
[ 8.30286496e-01, 3.42026003e-02, 1.61106001e-01,  
-2.04570330e-01, -7.48710486e-02],  
[ 7.94476585e-01, 2.83818596e-02, 1.51219548e-01,  
-2.07681498e-01, -9.12497108e-02],  
[ 8.11481669e-01, 4.32314899e-02, 1.43645563e-01,  
-1.07778264e-01, 5.79540241e-02],  
[-3.37741907e-01, 8.64927625e-01, 3.58933705e-02,  
9.07183972e-02, 4.46327264e-02],  
[ 5.03572654e-02, 9.35515351e-01, 1.51475403e-01,  
-2.51501256e-02, -9.34471627e-02],  
[ 9.78242247e-01, -3.31490292e-02, -1.05261173e-01,  
4.50364254e-02, 7.37362061e-02],  
[ 9.59137204e-01, -3.90023014e-02, -1.20630340e-01,  
4.52591439e-02, 6.64877295e-02],  
[ 8.14087167e-01, 2.23057235e-03, 7.66518536e-02,  
2.02747428e-02, 1.27634815e-01],  
[-4.15353984e-01, 7.18339585e-01, 3.40068065e-01,  
-7.18402769e-02, -2.77950513e-01],  
[ 7.64912662e-02, 7.24900629e-01, 2.74193203e-01,  
-4.83952643e-02, -3.52988282e-01],  
[ 9.10390833e-01, -5.36541214e-02, -4.68641819e-02,  
-7.64183441e-04, 1.63870440e-01],  
[ 8.73011872e-01, -5.30302307e-02, -5.89943125e-02,  
-1.58989743e-03, 1.52417545e-01],  
[ 7.55087663e-01, -3.56133824e-03, 5.39542590e-02,  
4.24181436e-03, 2.58043475e-01],  
[-1.23469882e-01, 6.07438109e-01, 6.33039195e-01,  
-2.14798897e-02, 2.47973902e-01],  
[-3.42866889e-01, 5.59526278e-01, 5.88213005e-01,  
-2.51533548e-02, 2.18419885e-01],  
[-1.60867206e-01, -1.53062590e-02, -1.57026584e-01,  
1.09243754e-01, -6.61660805e-01],  
[-1.37306764e-01, -2.17250646e-02, -1.58408933e-01,  
1.25156195e-01, -6.71630806e-01],  
[ 2.45096182e-01, -2.54584590e-02, -2.66691452e-02,  
9.53148496e-02, -6.42510840e-01],  
[ 2.03988656e-01, 7.85172835e-02, -3.01656228e-01,  
2.28379491e-02, -6.29223365e-01],  
[ 1.08926110e-01, -6.34332375e-02, -3.36565241e-02,  
-9.49480582e-02, 6.81473893e-01],  
[-2.63787624e-01, -6.43281163e-03, -3.58792147e-02,  
-9.37962446e-02, 6.47816997e-01],  
[-2.15717044e-01, -7.36588960e-02, 3.50113237e-01,



```
-1.95201626e-02, 6.36783769e-01],
[ 3.94306145e-01, 6.09565687e-02, 2.55337865e-01,
-2.20362099e-01, -1.84248084e-01],
[ 4.07877887e-01, 6.27256518e-02, 2.23926910e-01,
-2.10028737e-01, -1.71989227e-01],
[ 3.53156874e-01, 5.36715654e-02, 2.69603566e-01,
-2.16933217e-01, -1.80072068e-01],
[ 2.33537263e-01, -4.91732963e-02, 8.14450798e-01,
9.36688947e-02, 3.27131934e-01],
[ 2.40298212e-01, -3.38140094e-02, 8.31497001e-01,
7.52417674e-02, 2.46323616e-01],
[-6.71839510e-02, 6.58504550e-02, 5.86207693e-01,
8.72955244e-02, 9.12541350e-02],
[ 5.59835557e-02, 8.17918708e-01, -1.78458352e-01,
-1.55949438e-02, -3.34299731e-02],
[ 7.16426399e-02, 9.23428534e-01, -1.07142695e-01,
-2.78635371e-02, -4.35991120e-02],
[ 1.92496943e-01, -4.75870407e-02, 8.03173194e-01,
1.43492710e-01, 3.33862148e-01],
[ 1.87644429e-01, -3.29941023e-02, 8.58024491e-01,
1.31329954e-01, 2.55679719e-01],
[-1.02263658e-01, 6.03984260e-02, 4.72982256e-01,
7.36848384e-02, 1.12273907e-01],
[ 6.14776655e-02, 8.77962760e-01, -1.50410288e-01,
2.20991044e-02, -4.17158177e-02],
[ 7.83728218e-02, 9.54508791e-01, -5.91095909e-02,
1.64800936e-02, -4.32590999e-02],
[-3.24381907e-02, 1.11167165e-01, 7.84467399e-01,
-4.37718588e-02, -2.80931233e-01],
[ 1.76682389e-01, 1.90494237e-01, 5.61405482e-01,
-1.20746167e-01, -1.32570785e-01],
[-6.37386592e-02, -7.03047926e-02, -2.68934069e-01,
1.28589794e-01, 1.88507865e-01],
[-1.56051271e-01, -7.08033942e-02, -1.45964500e-01,
1.24253735e-01, 1.46293116e-01],
[-3.56716299e-01, -5.29910748e-02, 1.47771610e-01,
2.87196214e-02, 1.13159576e-01],
[ 2.42173821e-01, -2.86199139e-02, -3.25958384e-02,
1.05027822e-01, -6.55406092e-01],
[ 3.50196758e-01, -1.05016411e-02, -3.95274124e-01,
5.92876786e-02, 2.91651801e-01],
[ 2.25671546e-01, -3.42672751e-02, 8.92876642e-01,
1.12426818e-01, 2.67065205e-01]]))
```

## Data Modeling : Linear Regression

1. Build a linear Regression model to predict the total monthly expenditure for home mortgages loan. Please refer 'deplotment\_RE.xlsx'. Column hc\_mortgage\_mean is predicted variable. This is the mean monthly mortgage and owner costs of specified geographical location. Note: Exclude loans from prediction model which have NaN (Not a Number) values for hc\_mortgage\_mean.

```
In [67]: df_train.columns
```

```
Out[67]: Index(['COUNTYID', 'STATEID', 'state', 'state_ab', 'city', 'place', 'type',
        'primary', 'zip_code', 'area_code', 'lat', 'lng', 'ALand', 'AWater',
        'pop', 'male_pop', 'female_pop', 'rent_mean', 'rent_median',
        'rent_stdev', 'rent_sample_weight', 'rent_samples', 'rent_gt_10',
        'rent_gt_15', 'rent_gt_20', 'rent_gt_25', 'rent_gt_30', 'rent_gt_35',
        'rent_gt_40', 'rent_gt_50', 'universe_samples', 'used_samples',
        'hi_mean', 'hi_median', 'hi_stdev', 'hi_sample_weight', 'hi_samples',
        'family_mean', 'family_median', 'family_stdev', 'family_sample_weight',
        'family_samples', 'hc_mortgage_mean', 'hc_mortgage_median',
        'hc_mortgage_stdev', 'hc_mortgage_sample_weight', 'hc_mortgage_samples',
        'hc_mean', 'hc_median', 'hc_stdev', 'hc_samples', 'hc_sample_weight',
        'home_equity_second_mortgage', 'second_mortgage', 'home_equity', 'debt',
        'second_mortgage_cdf', 'home_equity_cdf', 'debt_cdf', 'hs_degree',
        'hs_degree_male', 'hs_degree_female', 'male_age_mean',
        'male_age_median', 'male_age_stdev', 'male_age_sample_weight',
        'male_age_samples', 'female_age_mean', 'female_age_median',
        'female_age_stdev', 'female_age_sample_weight', 'female_age_samples',
        'pct_own', 'married', 'married_snp', 'separated', 'divorced',
        'bad_debt', 'bins', 'pop_density', 'age_median', 'pop_bins'],
        dtype='object')
```

```
In [68]: df_train['type'].unique()
type_dict={'type':{'City':1,
                  'Urban':2,
                  'Town':3,
                  'CDP':4,
                  'Village':5,
                  'Borough':6}
           }
df_train.replace(type_dict,inplace=True)
```

```
In [69]: df_train['type'].unique()
```

```
Out[69]: array([1, 2, 3, 4, 5, 6], dtype=int64)
```

```
In [70]: df_test.replace(type_dict,inplace=True)
```

```
In [71]: df_test['type'].unique()
```

Out[71]: array([4, 1, 6, 3, 5, 2], dtype=int64)

In [72]: feature\_cols=['COUNTYID','STATEID','zip\_code','type','pop', 'family\_mean',  
                  'second\_mortgage', 'home\_equity', 'debt','hs\_degree',  
                  'age\_median','pct\_own', 'married','separated', 'divorced']

In [73]: x\_train=df\_train[feature\_cols]  
         y\_train=df\_train['hc\_mortgage\_mean']

In [74]: x\_test=df\_test[feature\_cols]  
         y\_test=df\_test['hc\_mortgage\_mean']

In [75]: from sklearn.preprocessing import StandardScaler  
         from sklearn.linear\_model import LinearRegression  
         from sklearn.metrics import r2\_score, mean\_absolute\_error,mean\_squared\_error,accuracy\_score

In [76]: x\_train.head()

Out[76]:

	COUNTYID	STATEID	zip_code	type	pop	family_mean	second_mortgage	home_equity	debt	hs_degree	age_median	pct_own	married	separated	divorc
UID															
267822	53	36	13346	1	5230	67994.14790	0.02077	0.08919	0.52963	0.89288	44.666665	0.79046	0.57851	0.01240	0.087
246444	141	18	46616	1	2633	50670.10337	0.02222	0.04274	0.60855	0.90487	34.791665	0.52483	0.34886	0.01426	0.090
245683	63	18	46122	1	6881	95262.51431	0.00000	0.09512	0.73484	0.94288	41.833330	0.85331	0.64745	0.01607	0.106
279653	127	72	927	2	2700	56401.68133	0.01086	0.01086	0.52714	0.91500	49.750000	0.65037	0.47257	0.02021	0.101
247218	161	20	66502	1	5637	54053.42396	0.05426	0.05426	0.51938	1.00000	22.000000	0.13046	0.12356	0.00000	0.031

In [77]: sc=StandardScaler()  
         x\_train\_scaled=sc.fit\_transform(x\_train)  
         x\_test\_scaled=sc.fit\_transform(x\_test)

Run a model at a Nation level. If the accuracy levels and R square are not satisfactory proceed to below step.

In [78]: linereg=LinearRegression()  
         linereg.fit(x\_train\_scaled,y\_train)

Out[78]: LinearRegression()

In [79]: y\_pred=linereg.predict(x\_test\_scaled)

In [80]: print("Overall R2 score of linear regression model", r2\_score(y\_test,y\_pred))  
         print("Overall RMSE of linear regression model", np.sqrt(mean\_squared\_error(y\_test,y\_pred)))

Overall R2 score of linear regression model 0.7348210754610929  
Overall RMSE of linear regression model 323.10188949846344

The Accuracy and R2 score are good, but still will investigate the model performance at state level

Run another model at State level. There are 52 states in USA.

In [82]: state=df\_train['STATEID'].unique()  
         state[0:5]

Out[82]: array([36, 18, 72, 20, 1], dtype=int64)

In [83]: for i in [20,1,45]:  
         print("State ID-",i)  
  
         x\_train\_nation=df\_train[df\_train['COUNTYID']==i][feature\_cols]  
         y\_train\_nation=df\_train[df\_train['COUNTYID']==i]['hc\_mortgage\_mean']  
  
         x\_test\_nation=df\_test[df\_test['COUNTYID']==i][feature\_cols]  
         y\_test\_nation=df\_test[df\_test['COUNTYID']==i]['hc\_mortgage\_mean']  
  
         x\_train\_scaled\_nation=sc.fit\_transform(x\_train\_nation)  
         x\_test\_scaled\_nation=sc.fit\_transform(x\_test\_nation)  
  
         linereg.fit(x\_train\_scaled\_nation,y\_train\_nation)  
         y\_pred\_nation=linereg.predict(x\_test\_scaled\_nation)  
  
         print("Overall R2 score of linear regression model for state",i,":-" ,r2\_score(y\_test\_nation,y\_pred\_nation))  
         print("Overall RMSE of linear regression model for state",i,":-" ,np.sqrt(mean\_squared\_error(y\_test\_nation,y\_pred\_nation)))  
         print("\n")

State ID- 20  
Overall R2 score of linear regression model for state, 20 :- 0.6046603766461807  
Overall RMSE of linear regression model for state, 20 :- 307.97188999314733

State ID- 1  
Overall R2 score of linear regression model for state, 1 :- 0.8104382475484616  
Overall RMSE of linear regression model for state, 1 :- 307.8275861848435

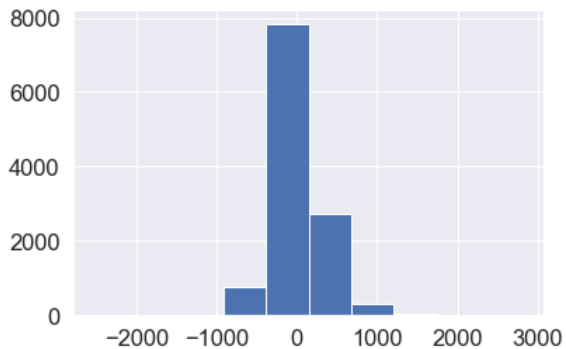
State ID- 45  
Overall R2 score of linear regression model for state, 45 :- 0.7887446497855253  
Overall RMSE of linear regression model for state, 45 :- 225.69615420724128

```
In [84]: residuals=y_test-y_pred  
residuals
```

```
Out[84]: UID  
255504    281.969088  
252676    -69.935775  
276314    190.761969  
248614   -157.290627  
286865    -9.887017  
...  
238088   -67.541646  
242811   -41.578757  
250127   -127.427569  
241096   -330.820475  
287763    217.760642  
Name: hc_mortgage_mean, Length: 11709, dtype: float64
```

```
In [85]: plt.hist(residuals) # Normal distribution of residuals
```

```
Out[85]: (array([6.000e+00, 3.000e+00, 2.900e+01, 7.670e+02, 7.823e+03, 2.716e+03,  
        3.010e+02, 4.900e+01, 1.200e+01, 3.000e+00]),  
array([-2515.04284233, -1982.92661329, -1450.81038425, -918.69415521,  
       -386.57792617,  145.53830287,  677.65453191, 1209.77076095,  
       1741.88698999, 2274.00321903, 2806.11944807]),  
<BarContainer object of 10 artists>)
```



```
In [86]: sns.distplot(residuals)
```

E:\Programs\Anaconda\lib\site-packages\seaborn\distributions.py:2619: FutureWarning:

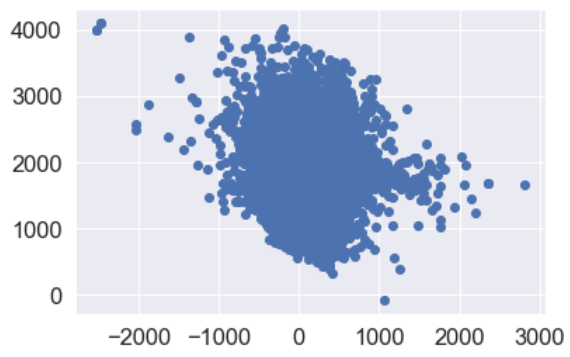
`distplot` is a deprecated function and will be removed in a future version. Please adapt your code to use either `displot` (a figure-level function with similar flexibility) or `histplot` (an axes-level function for histograms).

```
Out[86]: <AxesSubplot:xlabel='hc_mortgage_mean', ylabel='Density'>
```



```
In [87]: plt.scatter(residuals,y_pred)
```

```
Out[87]: <matplotlib.collections.PathCollection at 0x20704929400>
```



In [ ]: