

Real Estate

Problem Statement:

A banking institution requires actionable insights from the perspective of Mortgage-Backed Securities, Geographic Business Investment and Real Estate Analysis.

1. The objective is to identify white spaces/potential business in the mortgage loan.
2. The mortgage bank would like to identify potential monthly mortgage expenses for each region based on monthly family income and rental of the real estate.
3. A statistical model needs to be created to predict the potential demand in dollars amount of loan for each of the region in the USA. Also, there is a need to create a dashboard which would refresh periodically post data retrieval from the agencies. This would help to monitor the key metrics and trends.
4. The dashboard must demonstrate relationships and trends for the key metrics as follows: number of loans, average rental income, monthly mortgage and owner's cost, family income vs mortgage cost comparison across different regions. The metrics are described not to limit the dashboard to these few only.

Dataset Description :-

Following are the themes the fields fall under Home Owner Costs: Sum of utilities, property taxes.

1. Second Mortgage: Households with a second mortgage statistics.
2. Home Equity Loan: Households with a Home equity Loan statistics.
3. Debt: Households with any type of debt statistics.
4. Mortgage Costs: Statistics regarding mortgage payments, home equity loans, utilities and property taxes
5. Home Owner Costs: Sum of utilities, property taxes statistics
6. Gross Rent: Contract rent plus the estimated average monthly cost of utility features
7. Gross Rent as Percent of Income Gross rent as the percent of income very interesting
8. High school Graduation: High school graduation statistics.
9. Population Demographics: Population demographic statistics.
10. Age Demographics: Age demographic statistics.
11. Household Income: Total income of people residing in the household.
12. Family Income: Total income of people related to the householder.

```
In [1]:  
1 import time  
2 import random  
3 from math import *  
4 import operator  
5 import pandas as pd  
6 import numpy as np  
7  
# import plotting libraries  
9 import matplotlib  
10 import matplotlib.pyplot as plt  
11 from pandas.plotting import scatter_matrix  
%matplotlib inline  
13  
14 import seaborn as sns  
15 sns.set(style="white", color_codes=True)  
16 sns.set(font_scale=1.5)
```

1. Import data

```
In [2]: 1 df_train=pd.read_csv("train.csv")
```

```
In [3]: 1 df_test=pd.read_csv("test.csv")
```

```
In [4]: 1 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]: 1 df_test.columns
```

```
Out[5]: Index(['UID', 'BLOCKID', 'SUMLEVEL', 'COUNTYID', 'STATEID', 'state',
       'state_ab', 'city', 'place', 'type', 'zip_code', 'area_code',
       'lat', 'lng', 'ALand', 'AWater', 'pop', 'male_pop', 'female_pop',
       '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',
       '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]: 1 len(df_train)
```

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

```
In [7]: 1 len(df_test)
```

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

```
In [8]: 1 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	female
0	267822	NaN	140	53	36	New York	NY	Hamilton	Hamilton	City	...	44.48629	45.33333	22.51276	
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]: 1 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	female
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]: 1 df_train.describe()
```

```
Out[10]:
```

	UID	BLOCKID	SUMLEVEL	COUNTYID	STATEID	zip_code	area_code	lat	lng	ALand	...	female_age_mean	female
count	27321.000000	0.0	27321.0	27321.000000	27321.000000	27321.000000	27321.000000	27321.000000	27321.000000	2.732100e+04	...	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	
std	21343.859725	NaN	0.0	98.333097	16.392846	29558.115660	232.497482	5.588268	16.343816	1.275531e+09	...	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	
25%	238816.000000	NaN	140.0	29.000000	13.000000	26554.000000	405.000000	33.899064	-97.816067	1.799408e+06	...	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	
75%	275818.000000	NaN	140.0	109.000000	42.000000	77093.000000	801.000000	41.380606	-79.782503	3.359820e+07	...	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	

8 rows × 74 columns

```
In [11]: 1 df_test.describe()
```

Out[11]:

	UID	BLOCKID	SUMLEVEL	COUNTYID	STATEID	zip_code	area_code	lat	lng	ALand	...	female_age_mean	fem:
count	11709.000000	0.0	11709.0	11709.000000	11709.000000	11709.000000	11709.000000	11709.000000	11709.000000	1.170900e+04	...	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	
std	21466.372658	NaN	0.0	99.304334	16.607262	29775.134038	232.074263	5.625904	16.407818	7.624940e+08	...	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	
25%	238819.000000	NaN	140.0	29.000000	13.000000	25570.000000	404.000000	33.919813	-97.816561	1.718660e+06	...	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	
75%	276300.000000	NaN	140.0	109.000000	42.000000	77406.000000	787.000000	41.232973	-79.697311	3.204540e+07	...	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	

8 rows × 74 columns

```
In [12]: 1 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 
 ...   ...
```

```
In [13]: 1 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 
 ...   ...
```

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

```
In [14]: 1 #UID is unique userID value in the train and test dataset. So an index can be created from the UID feature
 2 df_train.set_index(keys=['UID'], inplace=True) #Set the DataFrame index using existing columns.
 3 df_test.set_index(keys=['UID'], inplace=True)
```

```
In [15]: 1 df_train.head(2)
```

```
Out[15]:
```

UID	BLOCKID	SUMLEVEL	COUNTYID	STATEID	state	state_ab	city	place	type	primary	...	female_age_mean	female_age_median	female_age_stdev	fem:
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]: 1 df_test.head(2)
```

Out[16]:

UID	BLOCKID	SUMLEVEL	COUNTYID	STATEID	state	state_ab	city	place	type	primary	...	female_age_mean	female_age_median	female_age_stdev	fema
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]: 1 #percentage of missing values in train set
2 missing_list_train=df_train.isnull().sum()*100/len(df_train)
3 missing_values_df_train=pd.DataFrame(missing_list_train,columns=['Percentage of missing values'])
4 missing_values_df_train.sort_values(by=['Percentage of missing values'],inplace=True,ascending=False)
5 missing_values_df_train[missing_values_df_train['Percentage of missing values'] >0][:10]
6 #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]: 1 #percentage of missing values in test set
2 missing_list_test=df_test.isnull().sum()*100/len(df_train)
3 missing_values_df_test=pd.DataFrame(missing_list_test,columns=['Percentage of missing values'])
4 missing_values_df_test.sort_values(by=['Percentage of missing values'],inplace=True,ascending=False)
5 missing_values_df_test[missing_values_df_test['Percentage of missing values'] >0][:10]
6 #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]: 1 df_train .drop(columns=['BLOCKID','SUMLEVEL'],inplace=True) #SUMLEVEL doest not have any predictive power and no variance
```

```
In [20]: 1 df_test .drop(columns=['BLOCKID','SUMLEVEL'],inplace=True) #SUMLEVEL doest not have any predictive power
```

```
In [21]: 1 # Imputing missing values with mean
```

```
2 missing_train_cols=[]
3 for col in df_train.columns:
4     if df_train[col].isna().sum() !=0:
5         missing_train_cols.append(col)
6 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_std
ev', '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_sp', 'separated', 'divorced']
```

```
In [22]: 1 # Imputing missing values with mean
2 missing_test_cols=[]
3 for col in df_test.columns:
4     if df_test[col].isna().sum() !=0:
5         missing_test_cols.append(col)
6 print(missing_test_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_std
dev', '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_sn
p', 'separated', 'divorced']

In [23]: 1 # Missing cols are all numerical variables
2 for col in df_train.columns:
3     if col in (missing_train_cols):
4         df_train[col].replace(np.nan, df_train[col].mean(), inplace=True)

In [24]: 1 # Missing cols are all numerical variables
2 for col in df_test.columns:
3     if col in (missing_test_cols):
4         df_test[col].replace(np.nan, df_test[col].mean(), inplace=True)

In [25]: 1 df_train.isna().sum().sum()

Out[25]: 0

In [26]: 1 df_test.isna().sum().sum()

Out[26]: 0
```

Exploratory Data Analysis (EDA):

Perform debt analysis. You may take the following steps:

- 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]: 1 from pandasql import sqldf
2 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
3 pysqldf = lambda q: sqldf(q, globals())
4 df_train_location_mort_pct=pysqldf(q1)
5

In [28]: 1 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-leto    0.11618       0.28972  28.029063 -82.495395
4   Lincolnwood    0.14228       0.28899  41.967289 -87.652434

In [29]: 1 import plotly.express as px
2 import plotly.graph_objects as go
```

In [30]:

```
1 fig = go.Figure(data=go.Scattergeo(
2     lat = df_train_location_mort_pct['lat'],
3     lon = df_train_location_mort_pct['lng'],
4 )
5 fig.update_layout(
6     geo=dict(
7         scope = 'north america',
8         showland = True,
9         landcolor = "rgb(212, 212, 212)",
10        subunitcolor = "rgb(255, 255, 255)",
11        countrycolor = "rgb(255, 255, 255)",
12        showlakes = True,
13        lakecolor = "rgb(255, 255, 255)",
14        showsubunits = True,
15        showcountries = True,
16        resolution = 50,
17        projection = dict(
18            type = 'conic conformal',
19            rotation_lon = -100
20        ),
21        lonaxis = dict(
22            showgrid = True,
23            gridwidth = 0.5,
24            range= [ -140.0, -55.0 ],
25            dtick = 5
26        ),
27        lataxis = dict (
28            showgrid = True,
29            gridwidth = 0.5,
30            range= [ 20.0, 60.0 ],
31            dtick = 5
32        )
33    ),
34    title='Top 2,500 locations with second mortgage is the highest and percent ownership is above 10 percent')
35 fig.show()
36
```

Top 2,500 locations with second mortgage is the highest and percent ownership is above 10 percent

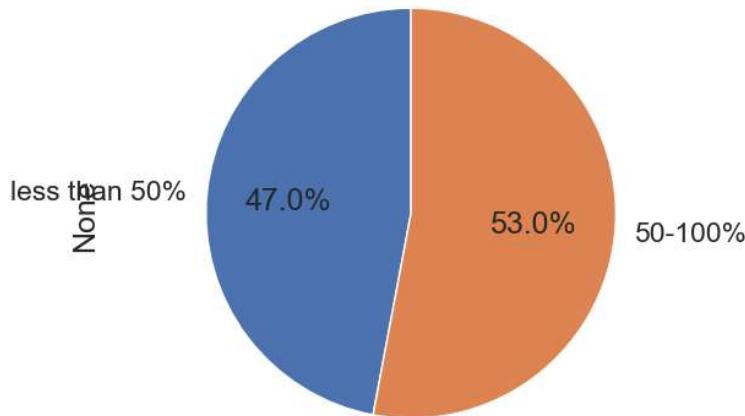


Use the following bad debt equation: $\text{Bad Debt} = P(\text{Second Mortgage} \cap \text{Home Equity Loan})$ $\text{Bad Debt} = \text{second_mortgage} + \text{home_equity} - \text{home_equity_second_mortgage}$ c) Create pie charts to show overall debt and bad debt

In [31]:

```
1 df_train['bad_debt']=df_train['second_mortgage']+df_train['home_equity']-df_train['home_equity_second_mortgage']
```

```
In [32]: 1 df_train['bins'] = pd.cut(df_train['bad_debt'], bins=[0,0.1,1], labels=["less than 50%","50-100%"])
2 df_train.groupby(['bins']).size().plot(kind='pie', subplots=True, startangle=90, autopct='%1.1f%%')
3 plt.axis('equal')
4
5 plt.show()
6 #df.plot.pie(subplots=True, figsize=(8, 3))
```



Create Box and whisker plot and analyze the distribution for 2nd mortgage, home equity, good debt, and bad debt for different cities

```
In [33]: 1 cols=[]
2 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')
```

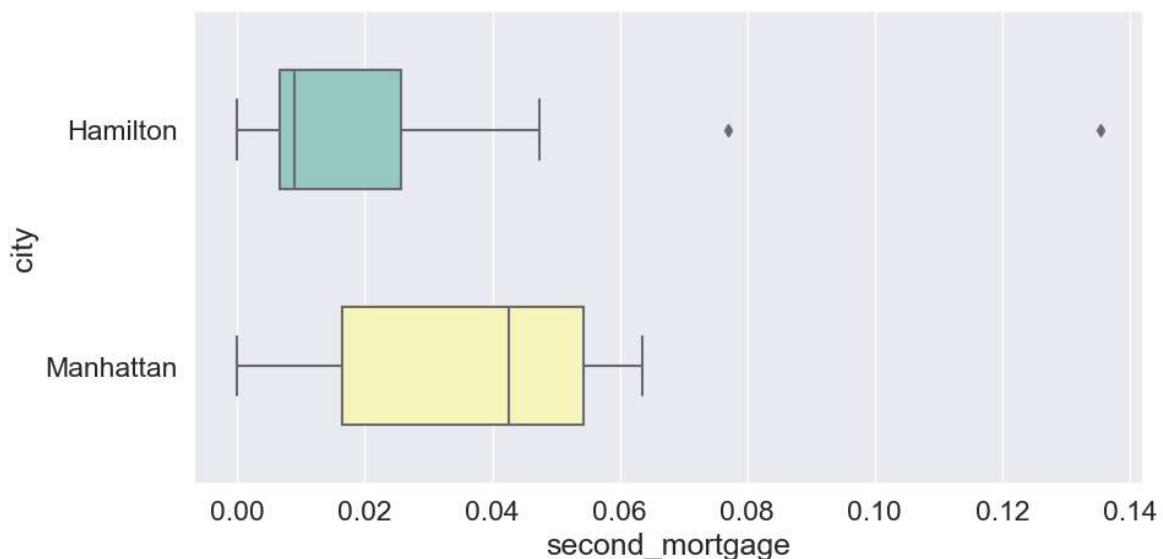
```
In [34]: 1 #Taking Hamilton and Manhattan cities data
2 cols=['second_mortgage','home_equity','debt','bad_debt']
3 df_box_hamilton=df_train.loc[df_train['city'] == 'Hamilton']
4 df_box_manhattan=df_train.loc[df_train['city'] == 'Manhattan']
5 df_box_city=pd.concat([df_box_hamilton,df_box_manhattan])
6 df_box_city.head(4)
```

```
Out[34]:
```

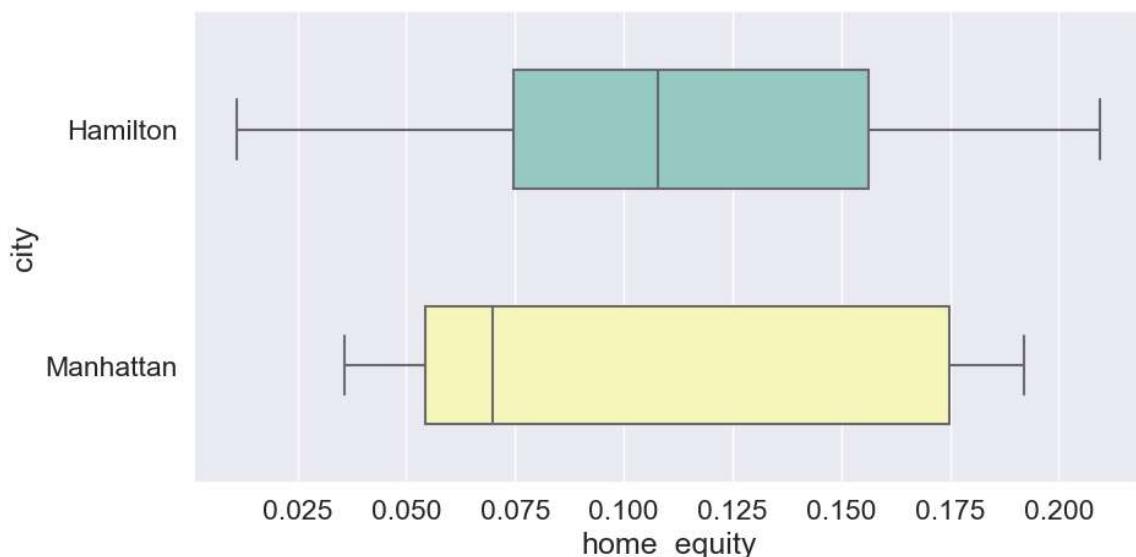
	COUNTYID	STATEID	state	state_ab	city	place	type	primary	zip_code	area_code	...	female_age_stdev	female_age_sample_weight	female_age_s
UID														
267822	53	36	New York	NY	Hamilton	Hamilton	City	tract	13346	315	...	22.51276	685.33845	
263797	21	34	New Jersey	NJ	Hamilton	Yardville	City	tract	8610	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

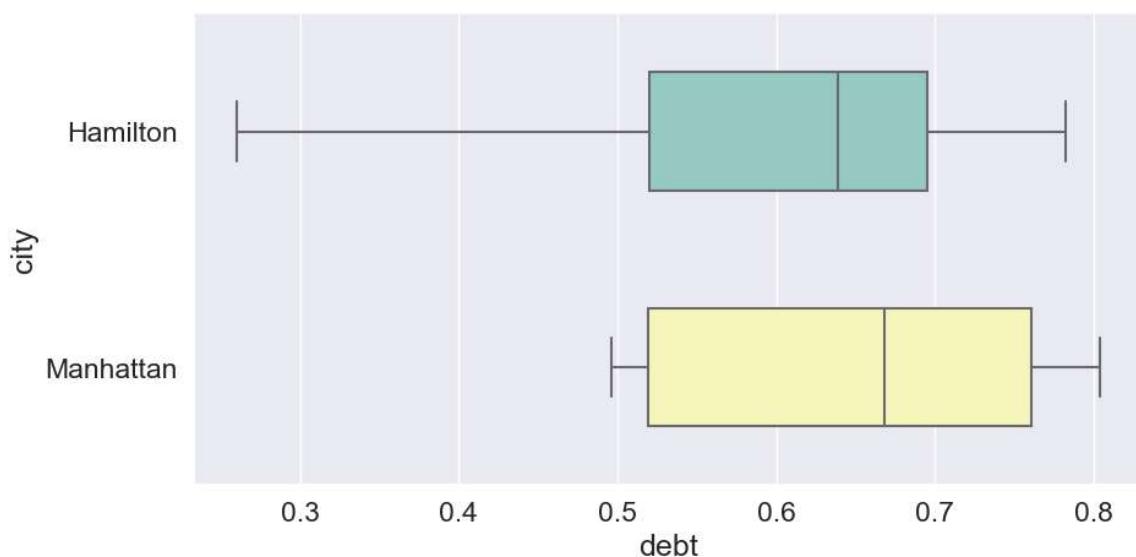
```
In [35]: 1 plt.figure(figsize=(10,5))
2 sns.boxplot(data=df_box_city,x='second_mortgage', y='city',width=0.5,palette="Set3")
3 plt.show()
```



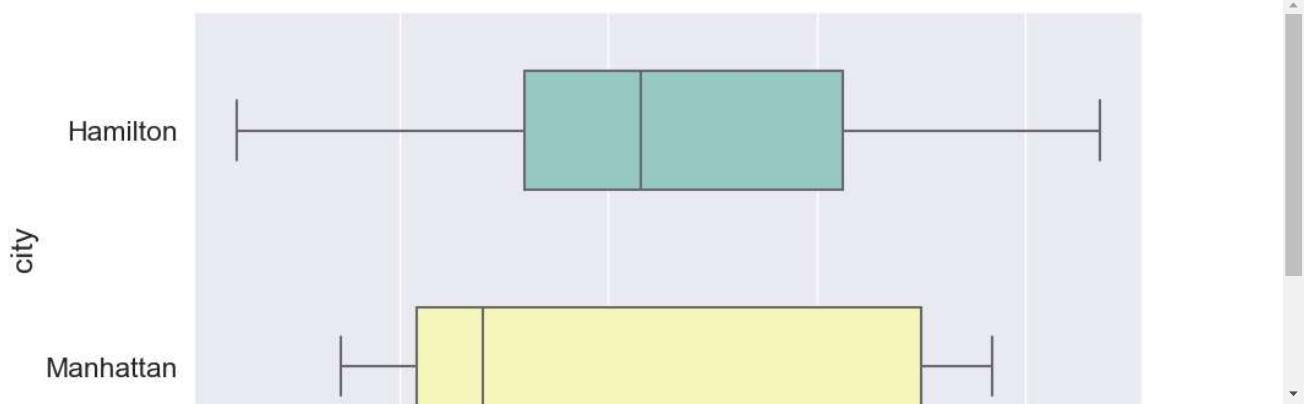
```
In [36]: 1 plt.figure(figsize=(10,5))
2 sns.boxplot(data=df_box_city,x='home_equity', y='city',width=0.5,palette="Set3")
3 plt.show()
```



```
In [37]: 1 plt.figure(figsize=(10,5))
2 sns.boxplot(data=df_box_city,x='debt', y='city',width=0.5,palette="Set3")
3 plt.show()
```



```
In [38]: 1 plt.figure(figsize=(10,5))
2 sns.boxplot(data=df_box_city,x='bad_debt', y='city',width=0.5,palette="Set3")
3 plt.show()
```



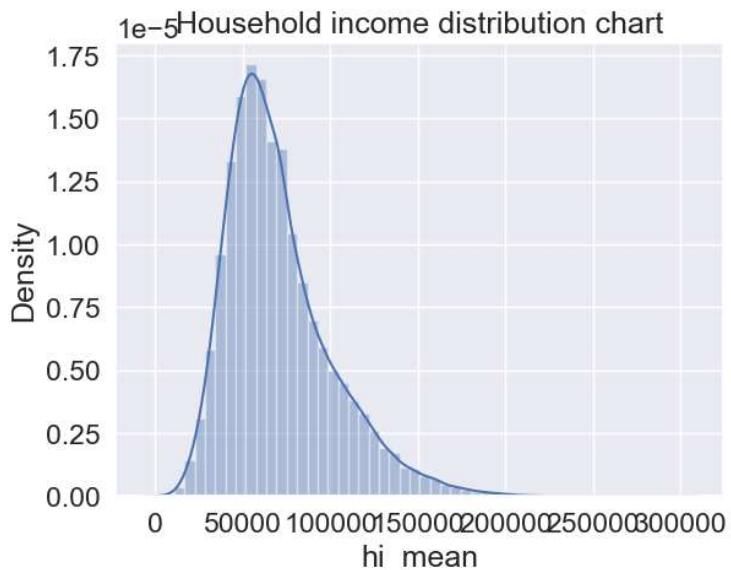
Manhattan has higher metrics compared to Hamilton

Create a collated income distribution chart for family income, house hold income, and remaining income

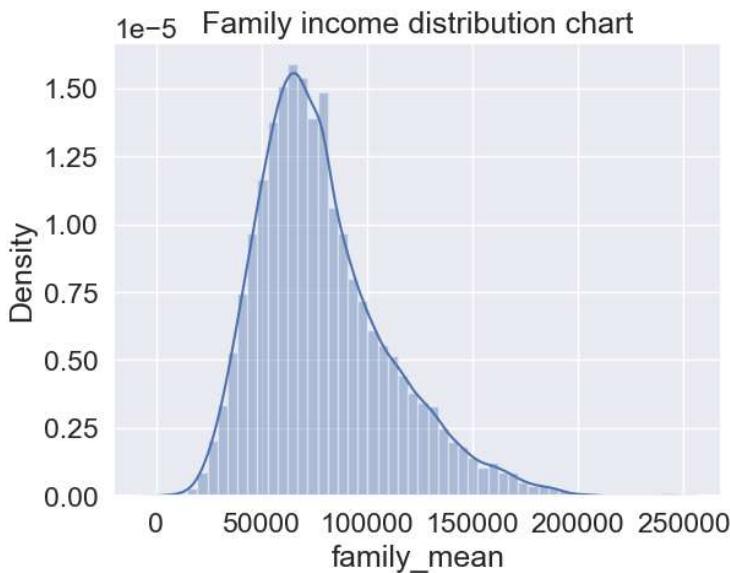
```
In [39]: 1 sns.distplot(df_train['hi_mean'])
2 plt.title('Household income distribution chart')
3 plt.show()
```

C:\Users\ANANT KALEKAR\anaconda3\lib\site-packages\seaborn\distributions.py:2619: FutureWarning:

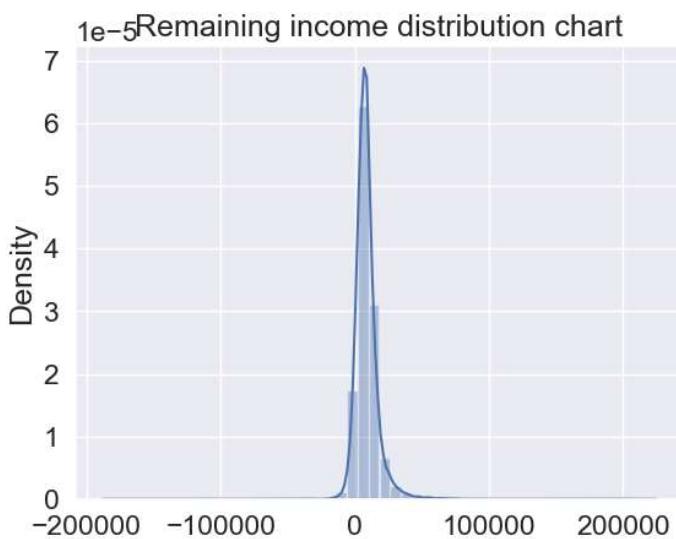
`distplot` is a deprecated function and will be removed in a future version. Please adapt your code to use either `displot` (a figure-level function with similar flexibility) or `histplot` (an axes-level function for histograms).



```
In [40]: 1 sns.distplot(df_train['family_mean'])
2 plt.title('Family income distribution chart')
3 plt.show()
C:\Users\ANANT KALEKAR\anaconda3\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]: 1 sns.distplot(df_train['family_mean']-df_train['hi_mean'])
2 plt.title('Remaining income distribution chart')
3 plt.show()
C:\Users\ANANT KALEKAR\anaconda3\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]:

```
1 #plt.figure(figsize=(25,10))
2 fig,(ax1,ax2,ax3)=plt.subplots(3,1)
3 sns.distplot(df_train['pop'],ax=ax1)
4 sns.distplot(df_train['male_pop'],ax=ax2)
5 sns.distplot(df_train['female_pop'],ax=ax3)
6 plt.subplots_adjust(wspace=0.8,hspace=0.8)
7 plt.tight_layout()
8 plt.show()
```

C:\Users\ANANT KALEKAR\anaconda3\lib\site-packages\seaborn\distributions.py:2619: FutureWarning:

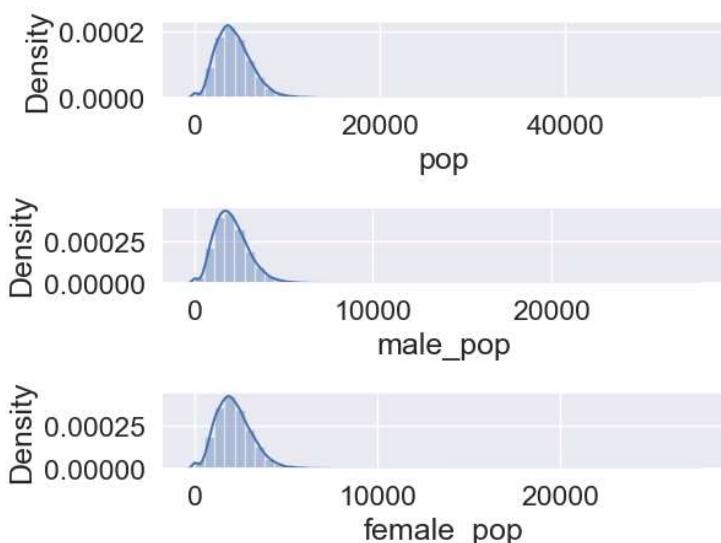
`distplot` is a deprecated function and will be removed in a future version. Please adapt your code to use either `displot` (a figure-level function with similar flexibility) or `histplot` (an axes-level function for histograms).

C:\Users\ANANT KALEKAR\anaconda3\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).

C:\Users\ANANT KALEKAR\anaconda3\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).



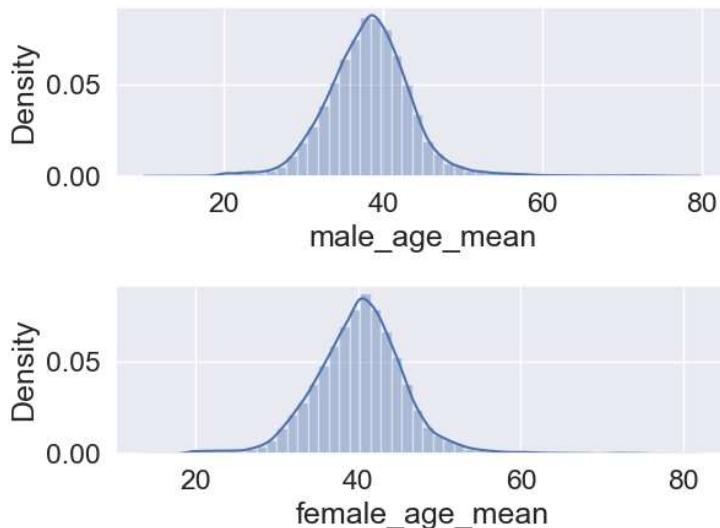
In []:

1

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

C:\Users\ANANT KALEKAR\anaconda3\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).

C:\Users\ANANT KALEKAR\anaconda3\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).
```



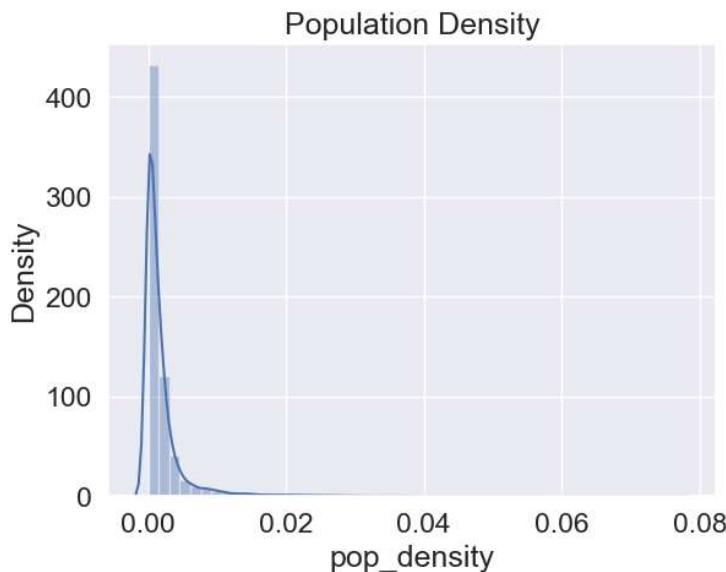
a) Use pop and ALand variables to create a new field called population density

```
In [44]: 1 df_train['pop_density']=df_train['pop']/df_train['ALand']

In [45]: 1 df_test['pop_density']=df_test['pop']/df_test['ALand']

In [46]: 1 sns.distplot(df_train['pop_density'])
2 plt.title('Population Density')
3 plt.show() # Very less density is noticed
```

C:\Users\ANANT KALEKAR\anaconda3\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]: 1 df_train['age_median']=(df_train['male_age_median']+df_train['female_age_median'])/2  
2 df_test['age_median']=(df_test['male_age_median']+df_test['female_age_median'])/2
```

```
In [48]: 1 df_train[['male_age_median','female_age_median','male_pop','female_pop','age_median']].head()
```

```
Out[48]:
```

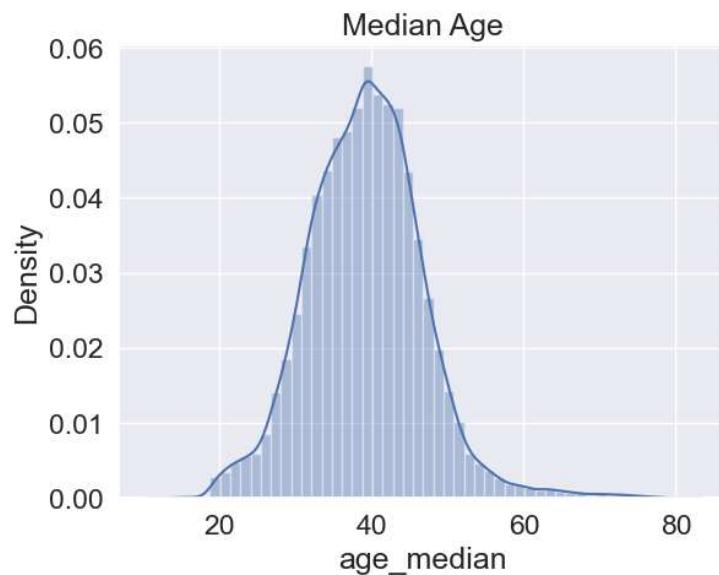
UID	male_age_median	female_age_median	male_pop	female_pop	age_median
267822	44.00000	45.33333	2612	2618	44.666665
246444	32.00000	37.58333	1349	1284	34.791665
245683	40.83333	42.83333	3643	3238	41.833330
279653	48.91667	50.58333	1141	1559	49.750000
247218	22.41667	21.58333	2586	3051	22.000000

```
In [49]: 1 sns.distplot(df_train['age_median'])  
2 plt.title('Median Age')  
3 plt.show()
```

```
4 # Age of population is mostly between 20 and 60  
5 # Majority are of age around 40  
6 # Median age distribution has a gaussian distribution  
7 # Some right skewness is noticed
```

```
C:\Users\ANANT KALEKAR\anaconda3\lib\site-packages\seaborn\distributions.py:2619: FutureWarning:
```

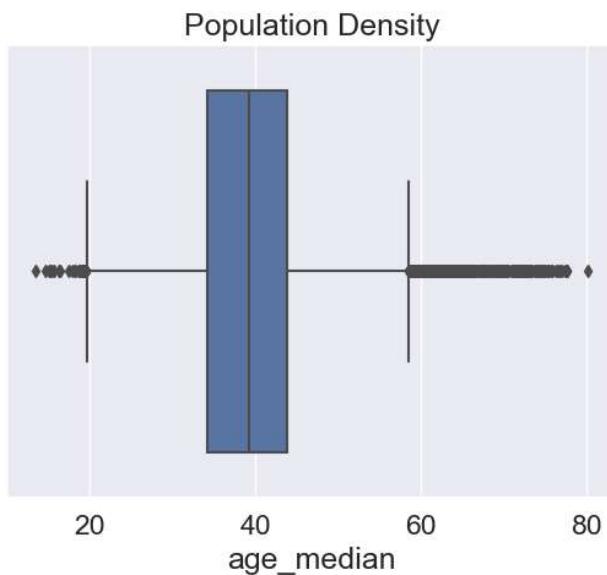
```
`distplot` is a deprecated function and will be removed in a future version. Please adapt your code to use either `displot` (a figure-level function with similar flexibility) or `histplot` (an axes-level function for histograms).
```



```
In [50]: 1 sns.boxplot(df_train['age_median'])
2 plt.title('Population Density')
3 plt.show()
```

C:\Users\ANANT KALEKAR\anaconda3\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]: 1 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]: 1 df_train['pop_bins']=pd.cut(df_train['pop'],bins=5,labels=['very low','low','medium','high','very high'])
```

```
In [53]: 1 df_train[['pop','pop_bins']]
```

```
Out[53]:
   pop    pop_bins
  UID
267822  5230  very low
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]: 1 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]: 1 df_train.groupby(by='pop_bins')[['married', 'separated', 'divorced']].count()
```

```
Out[55]:
```

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

```
In [56]: 1 df_train.groupby(by='pop_bins')[['married', 'separated', 'divorced']].agg(["mean", "median"])
```

```
Out[56]:
```

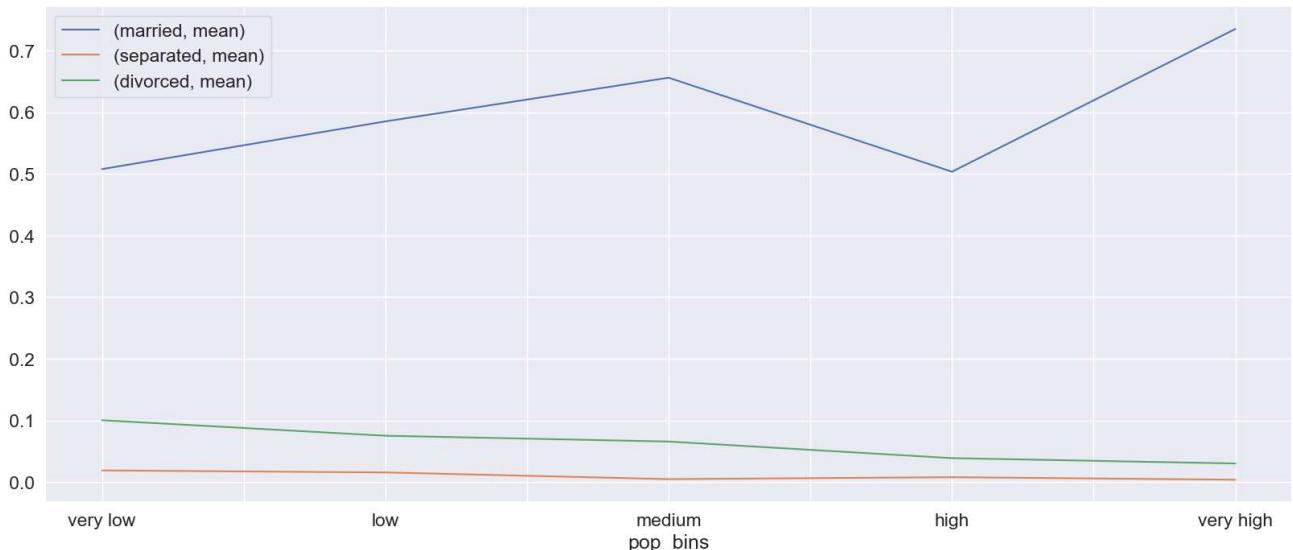
pop_bins	married		separated		divorced	
	mean	median	mean	median	mean	median
very low	0.507548	0.524680	0.019126	0.013650	0.100504	0.096020
low	0.584894	0.593135	0.015833	0.011195	0.075348	0.070045
medium	0.655737	0.618710	0.005003	0.004120	0.065927	0.064890
high	0.503359	0.335660	0.008141	0.002500	0.039030	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]: 1 plt.figure(figsize=(10,5))  
2 pop_bin_married=df_train.groupby(by='pop_bins')[['married', 'separated', 'divorced']].agg(["mean"])  
3 pop_bin_married.plot(figsize=(20,8))  
4 plt.legend(loc='best')  
5 plt.show()
```

<Figure size 1000x500 with 0 Axes>



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

```
In [58]: 1 rent_state_mean=df_train.groupby(by='state')['rent_mean'].agg(["mean"])  
2 rent_state_mean.head()
```

```
Out[58]:
```

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

```
In [59]: 1 income_state_mean=df_train.groupby(by='state')['family_mean'].agg(["mean"])
2 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]: 1 rent_perc_of_income=rent_state_mean['mean']/income_state_mean['mean']
2 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]: 1 #overall level rent as a percentage of income
2 sum(df_train['rent_mean'])/sum(df_train['family_mean'])
```

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

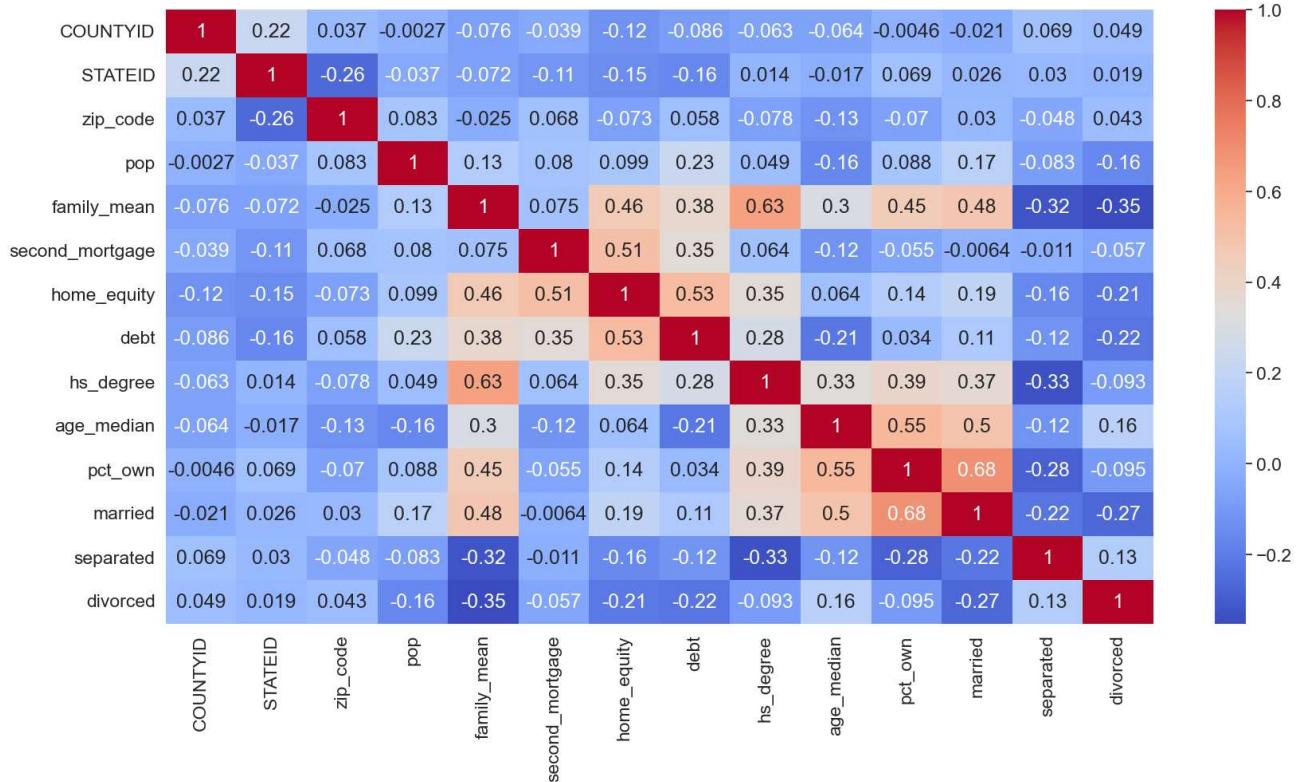
Perform correlation analysis for all the relevant variables by creating a heatmap. Describe your findings.

```
In [62]: 1 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_snip', 'separated', 'divorced',
 'bad_debt', 'bins', 'pop_density', 'age_median', 'pop_bins'],
 dtype='object')
```

```
In [63]: 1 cor=df_train[['COUNTYID','STATEID','zip_code','type','pop', 'family_mean',
2           'second_mortgage', 'home_equity', 'debt', 'hs_degree',
3           'age_median','pct_own', 'married', 'separated', 'divorced']].corr()
```

```
In [64]: 1 plt.figure(figsize=(20,10))
2 sns.heatmap(cor, annot=True, cmap='coolwarm')
3 plt.show()
```



1. High positive correlation is noticed between pop, male_pop and female_pop
2. High positive correlation 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]: 1 from sklearn.decomposition import FactorAnalysis
2 from factor_analyzer import FactorAnalyzer
```

```
In [66]: 1 #pip install factor_analyzer
```

```
In [67]: 1 fa=FactorAnalyzer(n_factors=5)
2 fa.fit_transform(df_train.select_dtypes(exclude= ('object','category')))
3 fa.loadings_
```

```
Out[67]: array([[-1.12589168e-01,  1.95646474e-02, -2.39331088e-02,
   -6.27632645e-02,  4.23474760e-02],
  [-1.10186763e-01,  1.33506219e-02,  2.79651249e-02,
   -1.49825865e-01,  1.10838809e-01],
  [-8.28678669e-02,  5.16372375e-02, -1.36451871e-01,
   -4.98918635e-02, -1.04824845e-01],
  [ 1.80961156e-02,  1.92013756e-02,  5.81329872e-03,
   2.64842748e-02, -6.12442393e-03],
  [ 9.02324703e-02, -9.72544309e-02, -6.54601319e-02,
   -1.33145902e-01, -1.48594605e-01],
  [-1.07335681e-02, -4.12376814e-02,  1.45853485e-01,
   8.80433354e-03,  1.08227569e-01],
  [-4.28796972e-02, -2.09780212e-02,  3.66726852e-02,
   -9.45597372e-02,  5.91380515e-02],
  [-2.44243021e-03, -1.53245408e-02, -2.68300881e-03,
   -4.52473039e-02,  2.37240654e-02],
  [ 7.92164326e-02,  9.57453312e-01, -8.71151658e-02,
   -6.59923876e-03, -3.97273205e-02],
  [ 7.39808208e-02,  9.18750523e-01, -1.08834842e-01,
```

Data Modeling : Linear Regression

1. Build a linear Regression model to predict the total monthly expenditure for home mortgages loan. Please refer 'deploment_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 [68]: 1 df_train.columns
```

```
Out[68]: 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_sn', 'separated', 'divorced',
       'bad_debt', 'bins', 'pop_density', 'age_median', 'pop_bins'],
      dtype='object')
```

```
In [69]: 1 df_train['type'].unique()
2 type_dict={'type':{1:
3     'Urban':2,
4     'Town':3,
5     'CDP':4,
6     'Village':5,
7     'Borough':6}
8   }
9 df_train.replace(type_dict,inplace=True)
```

```
In [70]: 1 df_train['type'].unique()
```

```
Out[70]: array([1, 2, 3, 4, 5, 6], dtype=int64)
```

```
In [71]: 1 df_test.replace(type_dict,inplace=True)
```

```
In [72]: 1 df_test['type'].unique()
```

```
Out[72]: array([4, 1, 6, 3, 5, 2], dtype=int64)
```

```
In [73]: 1 feature_cols=['COUNTYID','STATEID','zip_code','type','pop', 'family_mean',
2       'second_mortgage', 'home_equity', 'debt', 'hs_degree',
3       'age_median', 'pct_own', 'married', 'separated', 'divorced']
```

```
In [74]: 1 x_train=df_train[feature_cols]
2 y_train=df_train['hc_mortgage_mean']
```

```
In [75]: 1 x_test=df_test[feature_cols]
2 y_test=df_test['hc_mortgage_mean']
```

```
In [76]: 1 from sklearn.preprocessing import StandardScaler
2 from sklearn.linear_model import LinearRegression
3 from sklearn.metrics import r2_score, mean_absolute_error, mean_squared_error, accuracy_score
```

```
In [77]: 1 x_train.head()
```

```
Out[77]:
   COUNTYID STATEID zip_code type pop family_mean second_mortgage home_equity debt hs_degree age_median pct_own married separated divorced
    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  0.08770
246444      141      18  46616    1  2633  50670.10337      0.02222   0.04274  0.60855  0.90487  34.791665  0.52483  0.34886  0.01426  0.09030
245683       63      18  46122    1  6881  95262.51431      0.00000   0.09512  0.73484  0.94288  41.833330  0.85331  0.64745  0.01607  0.10657
279653      127      72     927    2  2700  56401.68133      0.01086   0.01086  0.52714  0.91500  49.750000  0.65037  0.47257  0.02021  0.10106
247218      161      20  66502    1  5637  54053.42396      0.05426   0.05426  0.51938  1.00000  22.000000  0.13046  0.12356  0.00000  0.03109
```

```
In [78]: 1 sc=StandardScaler()
2 x_train_scaled=sc.fit_transform(x_train)
3 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 [79]: 1 linereg=LinearRegression()
2 linereg.fit(x_train_scaled,y_train)
```

```
Out[79]: LinearRegression()
```

```
In [80]: 1 y_pred=linereg.predict(x_test_scaled)
```

```
In [81]: 1 print("Overall R2 score of linear regression model", r2_score(y_test,y_pred))
2 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 323.1018894984635
```

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]: 1 state=df_train['STATEID'].unique()
2 state[0:5]
3 #Picking a few IDs 20,1,45,6
```

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

```
In [83]: 1 for i in [20,1,45]:
2     print("State ID-",i)
3
4     x_train_nation=df_train[df_train['COUNTYID']==i][feature_cols]
5     y_train_nation=df_train[df_train['COUNTYID']==i]['hc_mortgage_mean']
6
7     x_test_nation=df_test[df_test['COUNTYID']==i][feature_cols]
8     y_test_nation=df_test[df_test['COUNTYID']==i]['hc_mortgage_mean']
9
10    x_train_scaled_nation=sc.fit_transform(x_train_nation)
11    x_test_scaled_nation=sc.fit_transform(x_test_nation)
12
13    linereg.fit(x_train_scaled_nation,y_train_nation)
14    y_pred_nation=linereg.predict(x_test_scaled_nation)
15
16    print("Overall R2 score of linear regression model for state," ,i,":-" ,r2_score(y_test_nation,y_pred_nation))
17    print("Overall RMSE of linear regression model for state," ,i,":-" ,np.sqrt(mean_squared_error(y_test_nation,y_pred_nation)))
18    print("\n")
```

State ID- 20

Overall R2 score of linear regression model for state, 20 :- 0.6046603766461809

Overall RMSE of linear regression model for state, 20 :- 307.97188999314716

State ID- 1

Overall R2 score of linear regression model for state, 1 :- 0.8104382475484616

Overall RMSE of linear regression model for state, 1 :- 307.82758618484354

State ID- 45

Overall R2 score of linear regression model for state, 45 :- 0.7887446497855252

Overall RMSE of linear regression model for state, 45 :- 225.69615420724134

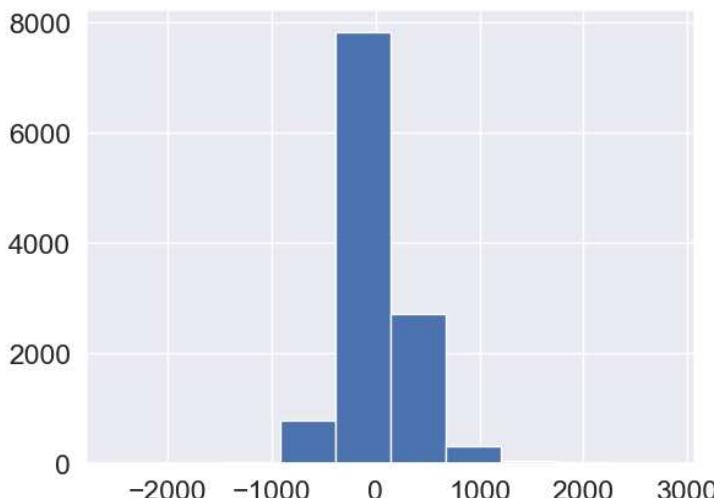
```
In [84]: 1 # To check the residuals
```

```
In [85]: 1 residuals=y_test-y_pred
2 residuals
```

```
Out[85]: 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
```

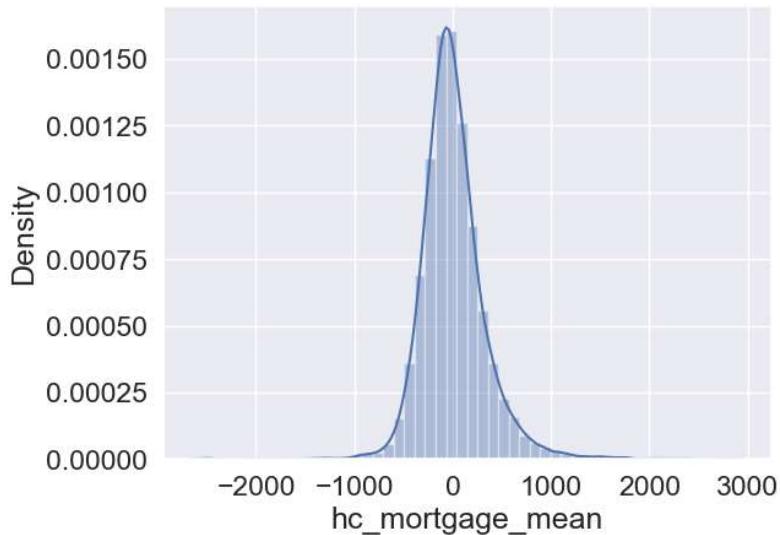
```
In [86]: 1 plt.hist(residuals) # Normal distribution of residuals
```

```
Out[86]: (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>)
```



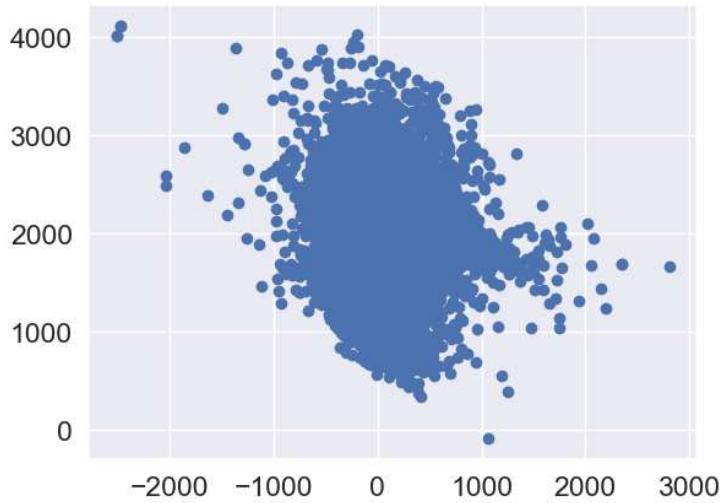
```
In [87]: 1 sns.distplot(residuals)
C:\Users\ANANT KALEKAR\anaconda3\lib\site-packages\seaborn\distributions.py:2619: FutureWarning:
`distplot` is a deprecated function and will be removed in a future version. Please adapt your code to use either `displot` (a figure-level function with similar flexibility) or `histplot` (an axes-level function for histograms).
```

```
Out[87]: <AxesSubplot:xlabel='hc_mortgage_mean', ylabel='Density'>
```



```
In [88]: 1 plt.scatter(residuals,y_pred) # Same variance and residuals does not have correlation with predictor
2 # Independence of residuals
```

```
Out[88]: <matplotlib.collections.PathCollection at 0x15727c9e3d0>
```



```
In [ ]: 1
```