

Real Estate Capstone

Course-end Project 1

The image displays three vertically stacked screenshots of a Jupyter Notebook interface, showing the process of importing and examining real estate data.

Screenshot 1: The first screenshot shows the initial steps of data import. In [4] contains the imports for time, random, math, operator, pandas, and numpy. In [5] shows df_train being read from "train.csv". In [6] shows df_test being read from "test.csv". In [7] displays the columns of df_train, which include various geographical and demographic features like UID, BLOCKID, SUMLEVEL, COUNTYID, STATEID, state, state_ab, city, place, type, primary, zip_code, area_code, lat, lon, Aland, Water, 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, universal_samples, used_samples, Al_mean, Al_median, hi_stdev, hi_mean, 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, home_equity_debt, debt, second_mortgage_cff, home_equity_cff, debt_cff, 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, and dtype='object'.

Screenshot 2: The second screenshot continues with df_train.columns. In [8] shows the columns of df_test, which are identical to df_train's columns. In [9] shows the len of df_train as 27321. In [10] shows the len of df_test as 11709.

Screenshot 3: The third screenshot shows the final output of df_train.columns and df_test.columns. Both outputs are identical to the ones shown in Screenshot 1.

```
In [4]: import time
import random
from math import *
import operator
import pandas as pd
import numpy as np

# Import plotting Libraries
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)

In [5]: df_train=pd.read_csv("train.csv")

In [6]: df_test=pd.read_csv("test.csv")

In [7]: df_train.columns
Out[7]: Index(['UID', 'BLOCKID', 'SUMLEVEL', 'COUNTYID', 'STATEID', 'state',
       'state_ab', 'city', 'place', 'type', 'primary', 'zip_code', 'area_code',
       'lat', 'lon', 'Aland', 'Water', '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',
       'universal_samples', 'used_samples', 'Al_mean', 'Al_median', 'hi_stdev',
       'hi_mean', '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_cff', 'home_equity_cff', 'debt_cff', '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 [8]: df_test.columns
Out[8]: Index(['UID', 'BLOCKID', 'SUMLEVEL', 'COUNTYID', 'STATEID', 'state',
       'state_ab', 'city', 'place', 'type', 'primary', 'zip_code', 'area_code',
       'lat', 'lon', 'Aland', 'Water', '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',
       'universal_samples', 'used_samples', 'Al_mean', 'Al_median', 'hi_stdev',
       'hi_mean', '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_cff', 'home_equity_cff', 'debt_cff', '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 [9]: len(df_train)
Out[9]: 27321

In [10]: len(df_test)
Out[10]: 11709
```

Jupyter Real Estate Project Last Checkpoint: Last Thursday at 11:37 AM (autosaved)

In [11]: df_train.head()

	UID	BLOCKID	SUMLEVEL	COUNTYID	STATEID	state	state_ab	city	place	type	...	female_age_mean	female_age_median	female_...
0	267822	NaN	140	53	36	New York	NY	Hamilton	Hamilton	City	...	44.48629	45.33333	22.5
1	246444	NaN	140	141	18	Indiana	IN	South Bend	Roseland	City	...	36.48391	37.58333	23.4
2	245683	NaN	140	63	18	Indiana	IN	Danville	Danville	City	...	42.15810	42.83333	23.5
3	279653	NaN	140	127	72	Puerto Rico	PR	San Juan	Guaynabo	Urban	...	47.77526	50.58333	24.1
4	247218	NaN	140	161	20	Kansas	KS	Manhattan	Manhattan	City	...	24.17693	21.58333	11.1

5 rows x 80 columns

In [12]: df_test.head()

	UID	BLOCKID	SUMLEVEL	COUNTYID	STATEID	state	state_ab	city	place	type	...	female_age_mean	female_age_median	female...
0	255504	NaN	140	163	26	Michigan	MI	Detroit	Dearborn Heights	CDP	...	34.78682	33.75000	...
1	252676	NaN	140	1	23	Maine	ME	Auburn	Auburn City	City	...	44.23451	46.66667	...
2	276314	NaN	140	15	42	Pennsylvania	PA	Pine City	Millerton	Borough	...	41.62426	44.50000	...

In [13]: df_train.describe()

	UID	BLOCKID	SUMLEVEL	COUNTYID	STATEID	zip_code	area_code	lat	lng	Aland	...	female_ag...
count	27321.000000	0.0	27321.000000	27321.000000	27321.000000	27321.000000	27321.000000	27321.000000	27321.000000	2.732100e+04	...	2711.
mean	257331.999303	NaN	140.0	85.646426	28.271806	50081.999524	598.507668	37.508813	-91.288394	1.29510e-08	...	40
std	21343.859725	NaN	0.0	98.333097	16.392846	29558.115660	232.497482	5.588268	16.343816	1.275531e-09	...	!
min	220342.000000	NaN	140.0	1.000000	1.000000	602.000000	201.000000	17.929085	165.453872	4.113400e-04	...	1
25%	238816.000000	NaN	140.0	29.000000	13.000000	26554.000000	405.000000	33.899064	-97.816067	1.799408e-06	...	3
50%	257220.000000	NaN	140.0	63.000000	28.000000	47715.000000	614.000000	38.755183	-86.554374	4.866940e+06	...	40
75%	275813.000000	NaN	140.0	109.000000	42.000000	77093.000000	801.000000	41.380606	-79.782503	3.359820e-07	...	4
max	294334.000000	NaN	140.0	840.000000	72.000000	99925.000000	989.000000	67.074017	-65.379332	1.039510e+11	...	71

8 rows x 74 columns

In [14]: df_test.describe()

	UID	BLOCKID	SUMLEVEL	COUNTYID	STATEID	zip_code	area_code	lat	lng	Aland	...	female_ag...
count	11709.000000	0.0	11709.000000	11709.000000	11709.000000	11709.000000	11709.000000	11709.000000	11709.000000	1.170900e+04	...	11613
mean	257525.004783	NaN	140.0	85.710650	28.489196	50123.418396	593.598514	37.405491	-91.340229	1.095500e+08	...	40
std	21466.372658	NaN	0.0	99.304334	16.607262	29775.134038	232.074263	5.625904	16.407818	7.624940e+08	...	5
min	220336.000000	NaN	140.0	1.000000	1.000000	601.000000	201.000000	17.965835	166.770979	8.290000e+03	...	15

In [15]: 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  zip_code          27321 non-null  object 
 11  area_code         27321 non-null  int64  
 12  lat               27321 non-null  int64  
 13  lon               27321 non-null  float64
 14  lng               27321 non-null  float64
 15  Aland            27321 non-null  float64
 16  Awater           27321 non-null  float64
 17  top               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_stddev         27007 non-null  float64
 23  rent_sample_weight  27007 non-null  float64

```

jupyter Real Estate Project Last Checkpoint: Last Thursday at 11:37 AM (autosaved)

```
In [16]: df_test.info()
<class 'pandas.core.frame.DataFrame'>
RangeIndex: 11709 entries, 0 to 11708
Data columns (total 88 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  long               11709 non-null   float64
 15  ALand              11709 non-null   int64  
 16  Akwater            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_stddev         11561 non-null   float64
 23  rent_sample_weight 11561 non-null   float64
```

jupyter Real Estate Project Last Checkpoint: Last Thursday at 11:37 AM (autosaved)

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

```
In [17]: #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 [18]: df_train.head(2)
Out[18]:
   BLOCKID  SUMLEVEL  COUNTYID  STATEID  state  state_ab  city  place  type  primary ...  female_age_mean  female_age_median  female_agr
   UID
267822    NaN        140       53     36  New York    NY  Hamilton  Hamilton  City  tract ...  44.48629  45.33333  2
246444    NaN        140      141     18  Indiana    IN  South Bend  Roseland  City  tract ...  36.48391  37.58333  2
2 rows x 79 columns
```

```
In [19]: df_test.head(2)
Out[19]:
   BLOCKID  SUMLEVEL  COUNTYID  STATEID  state  state_ab  city  place  type  primary ...  female_age_mean  female_age_median  female_agr
   UID
255504    NaN        140       163     26  Michigan    MI  Detroit  Dearborn Heights  CDP  tract ...  34.78682  33.75000  2
```

jupyter Real Estate Project Last Checkpoint: Last Thursday at 11:37 AM (autosaved)

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 [20]: #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]

Out[20]:
   Percentage of missing values
   BLOCKID          100.000000
   hc_samples        2.986113
   hc_mean           2.986113
   hc_median         2.986113
   hc_stdev          2.986113
   hc_sample_weight  2.986113
   hc_mortgage_mean  2.097288
   hc_mortgage_stdev 2.097288
   hc_mortgage_sample_weight  2.097288
   hc_mortgage_samples  2.097288
```

jupyter Real Estate Project Last Checkpoint: Last Thursday at 11:37 AM (autosaved)

```
In [21]: #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 436 missing values
```

Out[21]: Percentage of missing values

	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	0.908930
hc_mortgage_stdev	0.908930
hc_mortgage_sample_weight	0.908930
hc_mortgage_samples	0.908930

```
In [22]: df_train .drop(columns=['BLOCKID','SUMLEVEL'],inplace=True) #SUMLEVEL does not have any predictive power and no variance
```

```
In [23]: df_test .drop(columns=['BLOCKID','SUMLEVEL'],inplace=True) #SUMLEVEL does not have any predictive power
```

jupyter Real Estate Project Last Checkpoint: Last Thursday at 11:37 AM (autosaved)

```
In [24]: # Imputing missing values 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_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', 'equity_cdf', 'debt_cdf', 'hs_degree', 'hs_degree_mean', 'hs_degree_stdev', '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_sample_weight', 'female_age_samples', 'pct_own', 'married', 'married_snp', 'separated', 'divorced']
```

```
In [25]: # Imputing missing values 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_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_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', 'equity_cdf', 'debt_cdf', 'hs_degree', 'hs_degree_mean', 'hs_degree_stdev', '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_sample_weight', 'female_age_samples', 'pct_own', 'married', 'married_snp', 'separated', 'divorced']
```

jupyter Real Estate Project Last Checkpoint: Last Thursday at 11:37 AM (autosaved)

```
In [26]: # Missing cols are all numerical variables
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 [27]: # Missing cols are all numerical variables
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 [28]: df_train.isna().sum().sum()
Out[28]: 0

In [29]: df_test.isna().sum().sum()
Out[29]: 0
```

Exploratory Data Analysis (EDA):

Perform debt analysis. You may take the following steps:

- 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.

localhost:8888/notebooks/Real%20Estate%20Project.ipynb#

jupyter Real Estate Project Last Checkpoint: Last Thursday at 11:37 AM (autosaved)

```
In [30]: import time
import random
from math import *
import operator
import pandas as pd
import numpy as np

# import plotting libraries
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)

In [31]: pip install pandasql

Requirement already satisfied: pandasql in c:\users\student_0002\anaconda3\lib\site-packages (0.7.3)
Requirement already satisfied: numpy in c:\users\student_0002\anaconda3\lib\site-packages (from pandasql) (1.23.5)
Requirement already satisfied: sqlalchemy in c:\users\student_0002\anaconda3\lib\site-packages (from pandasql) (1.4.39)
Requirement already satisfied: pandas in c:\users\student_0002\anaconda3\lib\site-packages (from pandasql) (1.5.3)
Requirement already satisfied: pytz>=2020.1 in c:\users\student_0002\anaconda3\lib\site-packages (from pandas->pandasql) (2022.7)
Requirement already satisfied: python-dateutil>=2.8.1 in c:\users\student_0002\anaconda3\lib\site-packages (from pandasql) (2.8.2)
Requirement already satisfied: greenlet<=0.4.17 in c:\users\student_0002\anaconda3\lib\site-packages (from sqlalchemy>pandasql) (2.0.1)
Requirement already satisfied: enum34 in c:\users\student_0002\anaconda3\lib\site-packages (from python-dateutil) (1.1.6)
Requirement already satisfied: six in c:\users\student_0002\anaconda3\lib\site-packages (from python-dateutil) (1.1.5)
Requirement already satisfied: pytz in c:\users\student_0002\anaconda3\lib\site-packages (from python-dateutil) (2022.7)
Requirement already satisfied: tzlocal in c:\users\student_0002\anaconda3\lib\site-packages (from python-dateutil) (4.1.0)
Requirement already satisfied: dateutil in c:\users\student_0002\anaconda3\lib\site-packages (from python-dateutil) (2.8.2)
Requirement already satisfied: six in c:\users\student_0002\anaconda3\lib\site-packages (from dateutil) (1.1.5)
Requirement already satisfied: pytz in c:\users\student_0002\anaconda3\lib\site-packages (from dateutil) (2022.7)
Requirement already satisfied: tzlocal in c:\users\student_0002\anaconda3\lib\site-packages (from dateutil) (4.1.0)
Requirement already satisfied: six in c:\users\student_0002\anaconda3\lib\site-packages (from tzlocal) (1.1.5)
Requirement already satisfied: pytz in c:\users\student_0002\anaconda3\lib\site-packages (from tzlocal) (2022.7)
Requirement already satisfied: tzlocal in c:\users\student_0002\anaconda3\lib\site-packages (from tzlocal) (4.1.0)

71°F Sunny 1:52 PM 8/29/2023
```

localhost:8888/notebooks/Real%20Estate%20Project.ipynb#

jupyter Real Estate Project Last Checkpoint: Last Thursday at 11:37 AM (autosaved)

```
In [32]: from pandasql import sqldf
q = "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"
psqldf = lambda q: sqldf(q, globals())
df_train_location_mort_pct = psqldf(q)

In [33]: df_train_location_mort_pct.head()

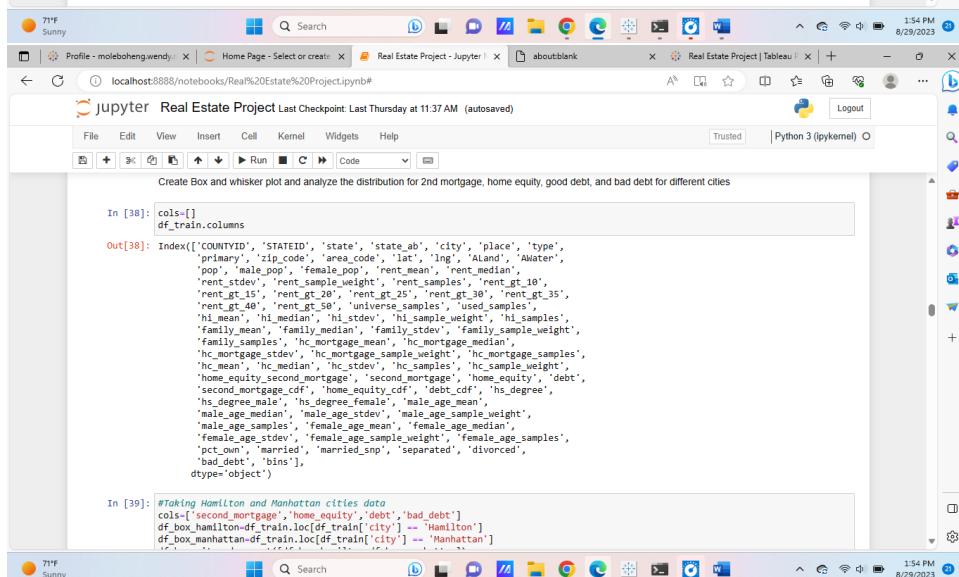
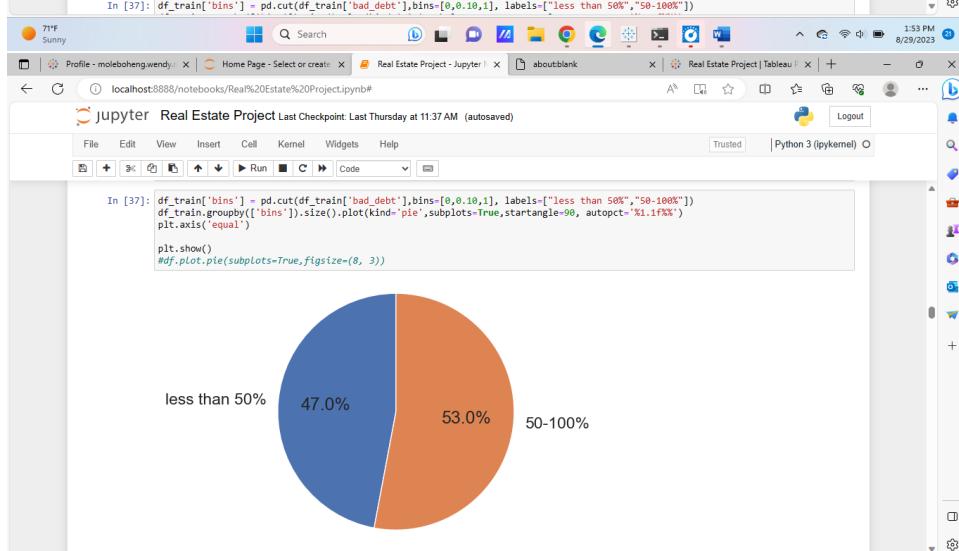
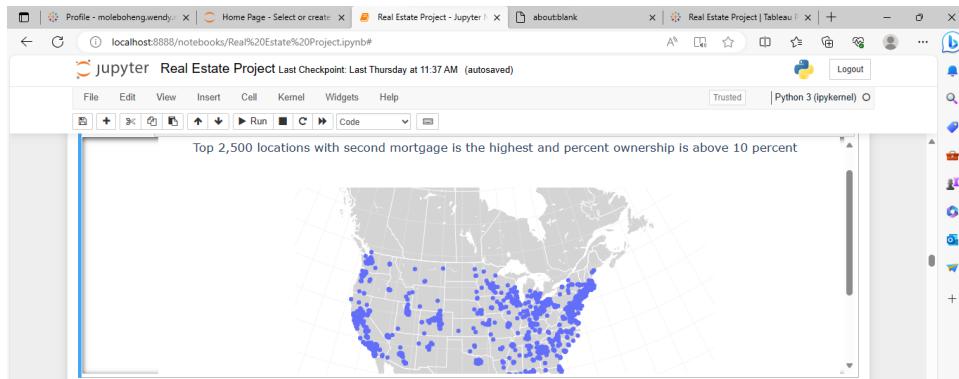
Out[33]:
   place  pct_own  second_mortgage      lat      lng
0  Worcester City  0.20247  0.43383  42.254262 -71.800347
1  Harbor Hills  0.15618  0.31818  40.751809 -73.853582
2  Glen Burnie  0.22380  0.30212  39.127273 -76.635265
3  Egypt Lake-leto  0.11618  0.28972  28.029063 -82.495395
4  Lincolnwood  0.14228  0.28899  41.967289 -87.652434

In [34]: import plotly.express as px
import plotly.graph_objects as go

In [35]: 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,
        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()
```

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

71°F Sunny 1:53 PM 8/29/2023



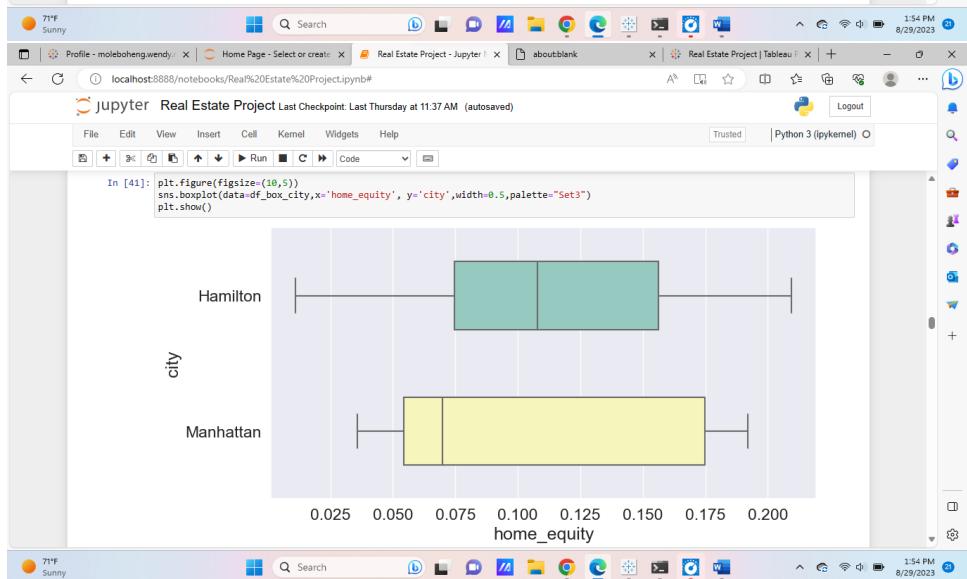
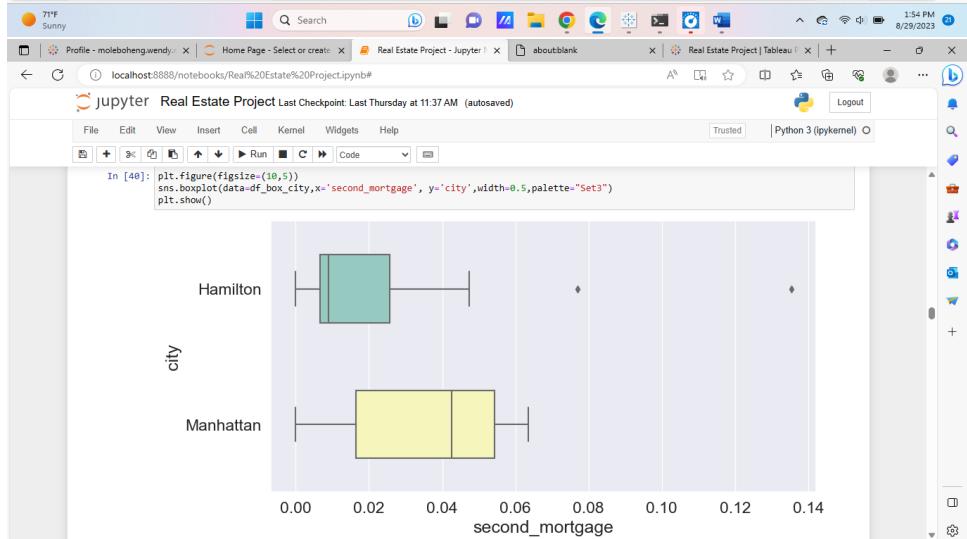
jupyter Real Estate Project Last Checkpoint: Last Thursday at 11:37 AM (autosaved)

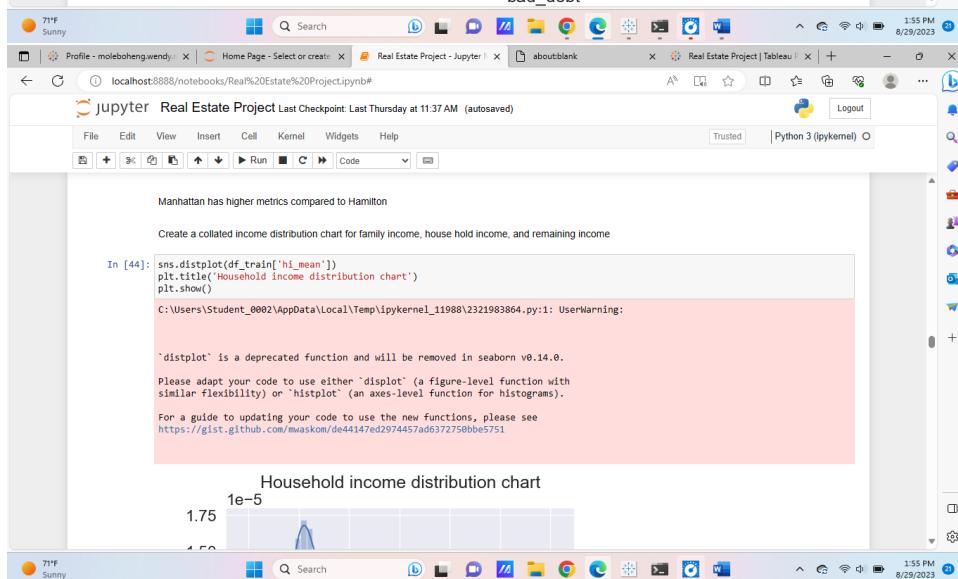
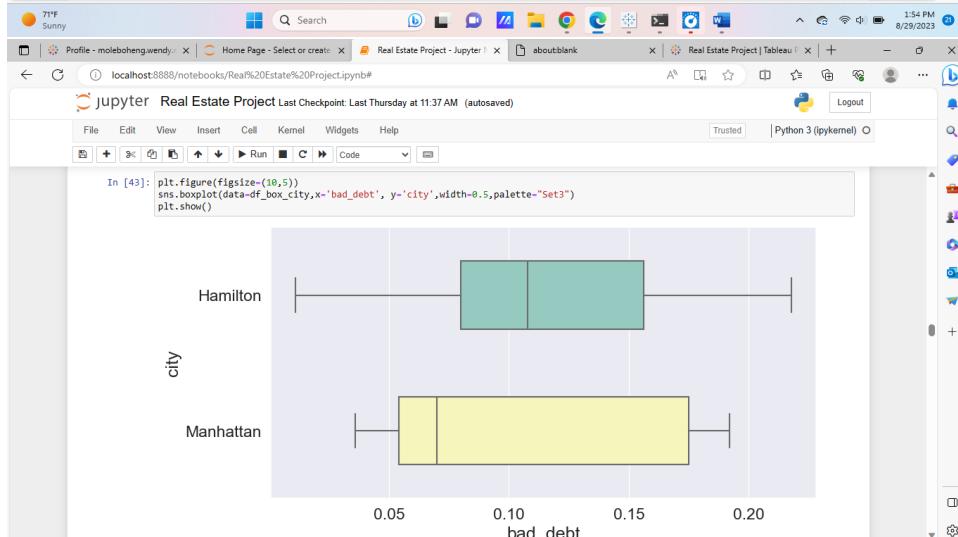
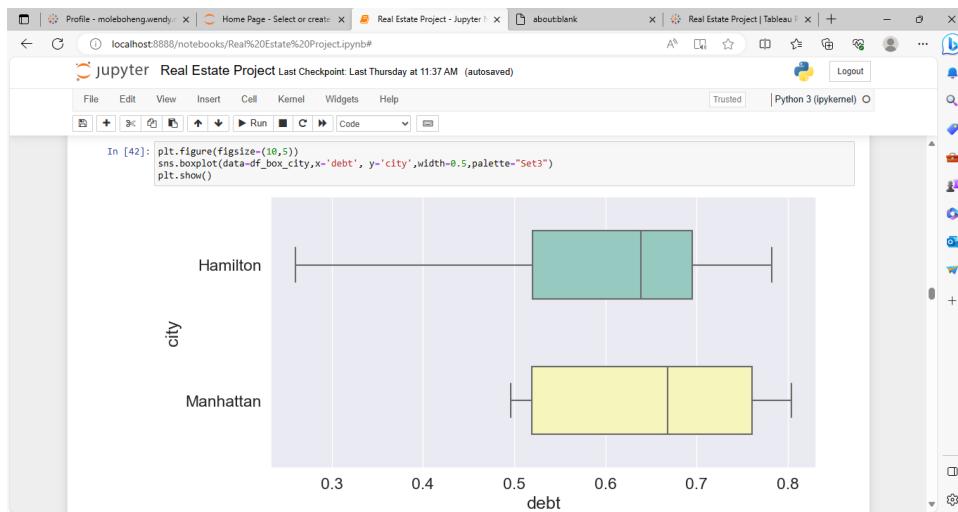
```
In [39]: #Taking Hamilton and Manhattan cities data
cols=['second_mortgage','home_equity','debt','bad_debt']
df_box_hamilton=df_train.loc[df_train['city'] == 'Hamilton']
df_box_manhattan=df_train.loc[df_train['city'] == 'Manhattan']
df_box_city=pd.concat([df_box_hamilton,df_box_manhattan])
df_box_city.head(4)
```

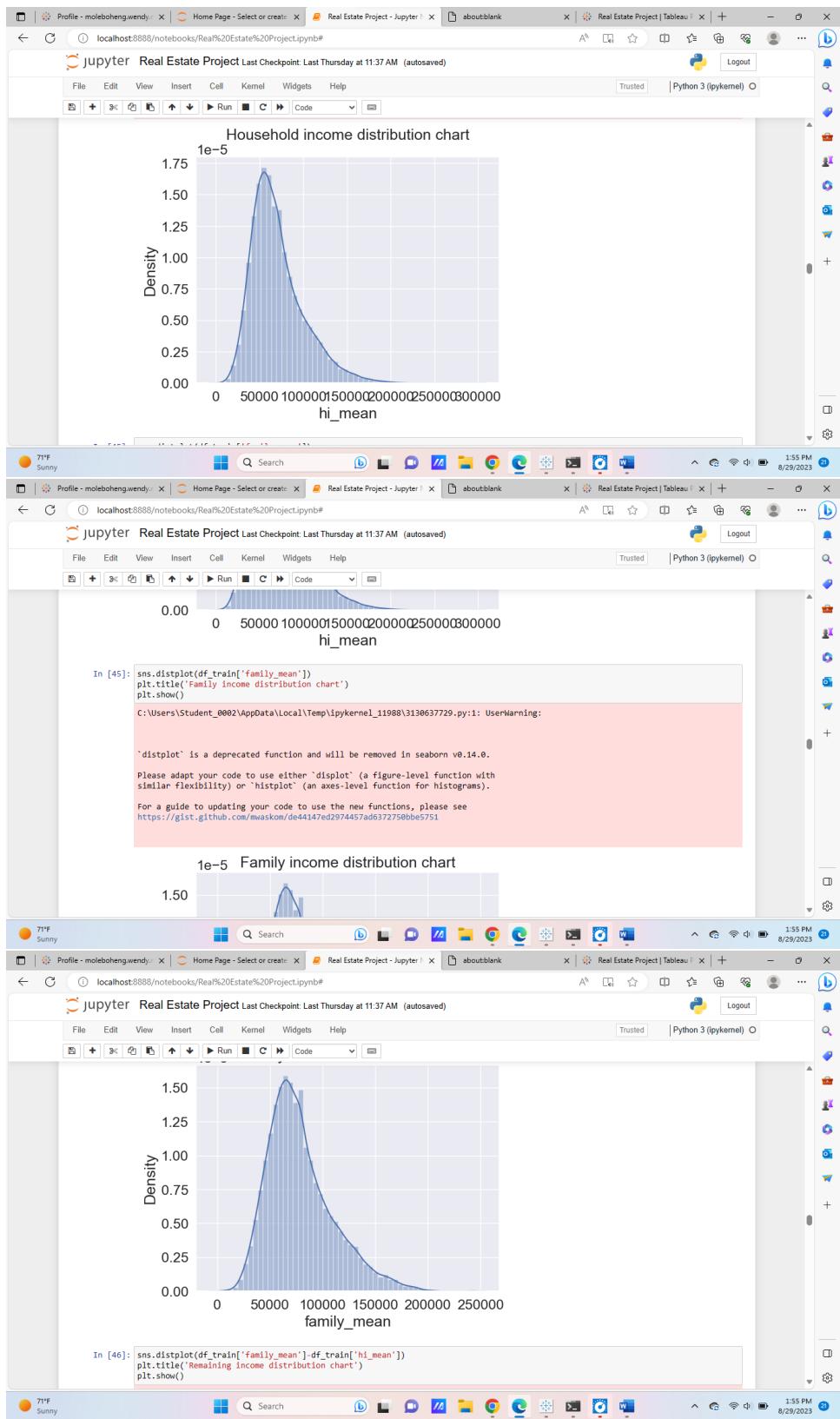
UID	COUNTYID	STATEID	state	state_ab	city	place	type	primary	zip_code	area_code	female_age_stddev	female_age_sample_weight
267822	53	36	New York	NY	Hamilton	Hamilton	City	tract	13346	315 ...	22.51276	685.33845
263797	21	34	New Jersey	NJ	Hamilton	Yardville	City	tract	8610	609 ...	24.05831	732.58443
270879	17	39	Ohio	OH	Hamilton	Hamilton City	Village	tract	45015	513 ...	22.66500	565.32725
259028	95	28	Mississippi	MS	Hamilton	Hamilton	CDP	tract	39748	662 ...	22.79602	483.01311

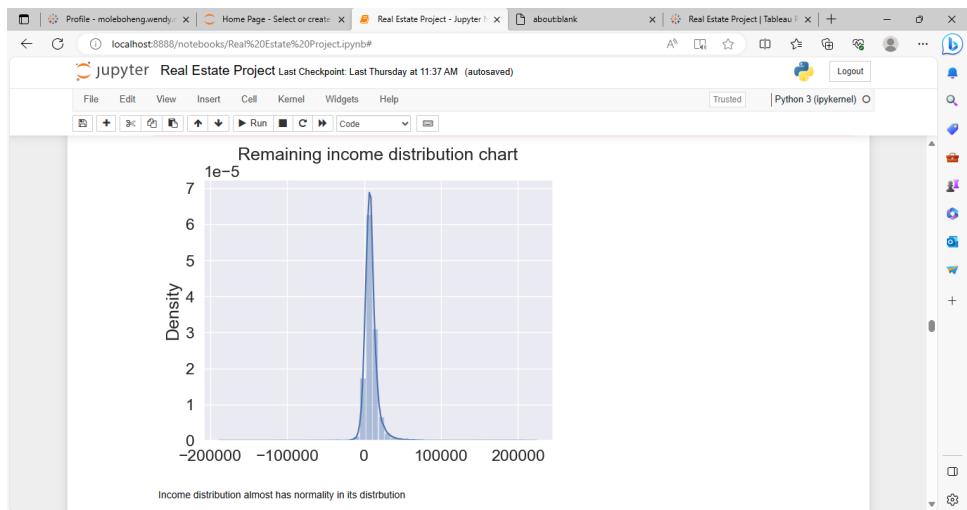
4 rows × 9 columns

```
In [40]: plt.figure(figsize=(10,5))
sns.boxplot(data=df_box_city,x='second_mortgage', y='city',width=0.5,palette="Set3")
plt.show()
```









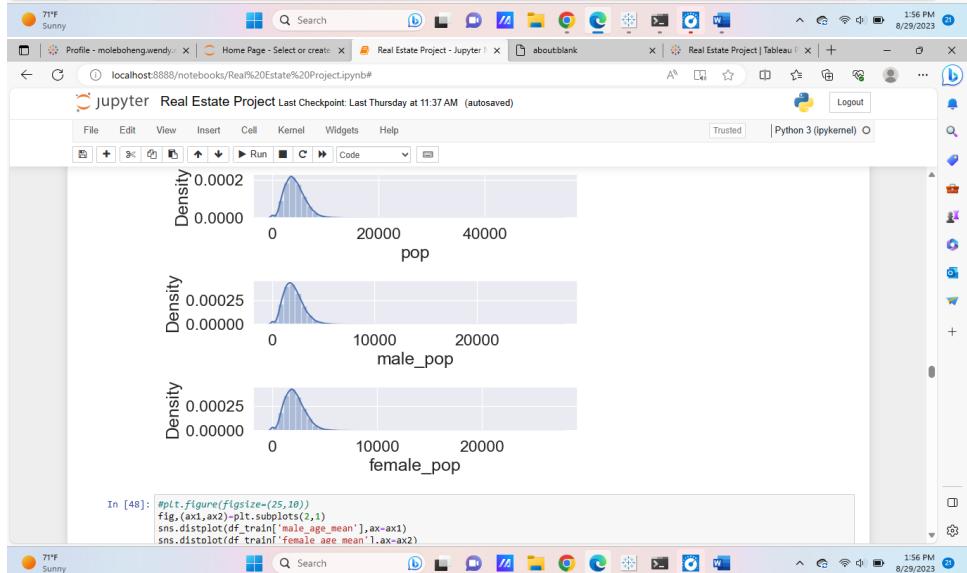
```
In [47]: plt.figure(figsize=(25,10))
fig,(ax1,ax2,ax3)=plt.subplots(3,1)
sns.distplot(df_train['pop'],ax=ax1)
sns.distplot(df_train['male_pop'],ax=ax2)
sns.distplot(df_train['female_pop'],ax=ax3)
plt.subplots_adjust(wspace=0.8,hspace=0.8)
plt.tight_layout()
plt.show()

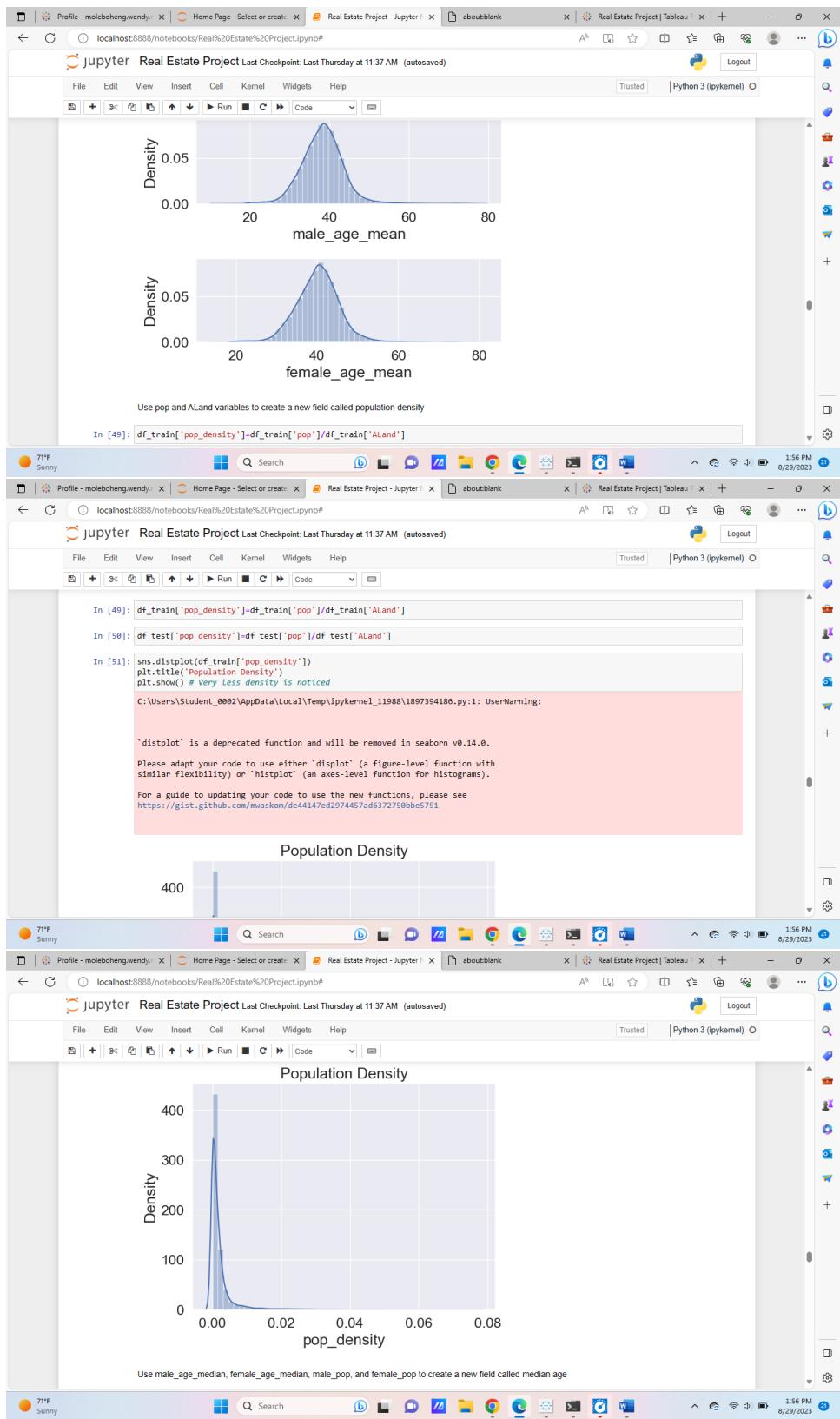
C:\Users\Student_0002\AppData\Local\Temp\ipykernel_11988\222623768.py:3: UserWarning:
```

'distplot' is a deprecated function and will be removed in seaborn v0.14.0.
Please adapt your code to use either 'displot' (a figure-level function with similar flexibility) or 'histplot' (an axes-level function for histograms).
For a guide to updating your code to use the new functions, please see
<https://gist.github.com/mwaskom/de44147ed2974457ad6372750bbe5751>

C:\Users\Student_0002\AppData\Local\Temp\ipykernel_11988\222623768.py:4: UserWarning:

'distplot' is a deprecated function and will be removed in seaborn v0.14.0.
Please adapt your code to use either 'displot' (a figure-level function with similar flexibility) or 'histplot' (an axes-level function for histograms).



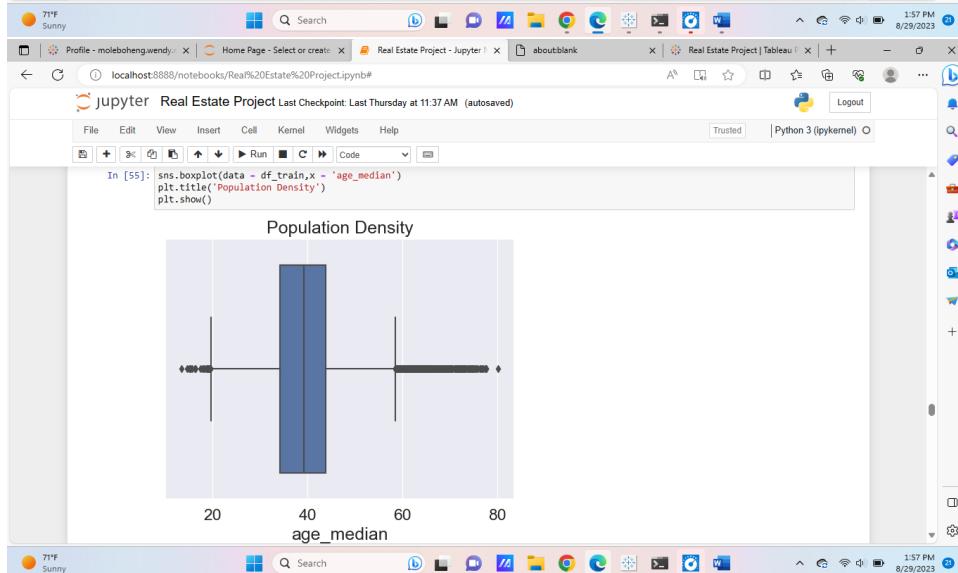
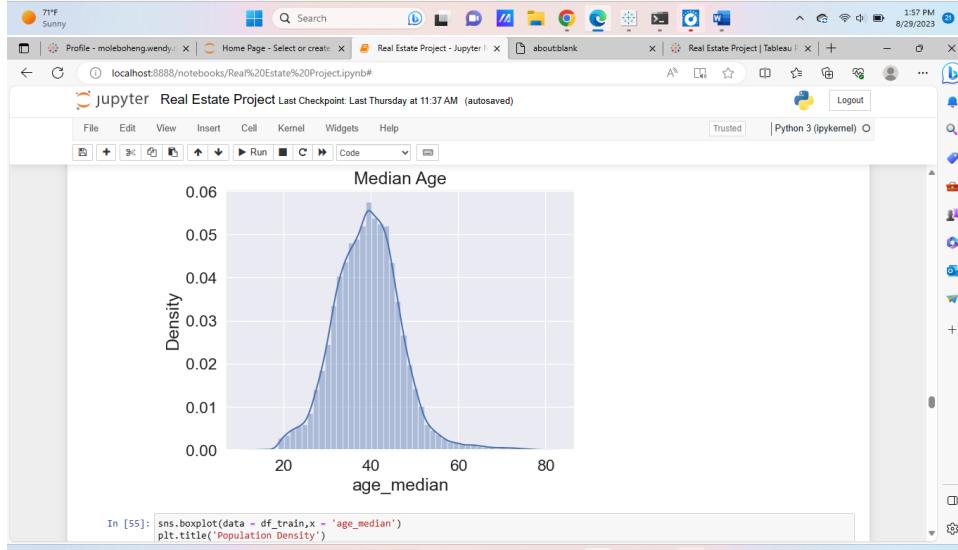


jupyter Real Estate Project Last Checkpoint: Last Thursday at 11:37 AM (autosaved)

```
In [52]: df_train['age_median']=df_train['male_age_median']+df_train['female_age_median'])/2  
df_test['age_median']=(df_train['male_age_median']+df_test['female_age_median'])/2  
  
In [53]: df_train[['male_age_median','female_age_median','male_pop','female_pop','age_median']].head()  
Out[53]:
```

UID	male_age_median	female_age_median	male_pop	female_pop	age_median
267822	44.00000	45.33333	2612	2618	44.66665
246444	32.00000	37.58333	1349	1284	34.79165
245683	40.83333	42.83333	3643	3238	41.83330
279653	48.91667	50.58333	1141	1559	49.75000
247218	22.41667	21.58333	2586	3051	22.00000

```
In [54]: 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  
C:\Users\Student_002\AppData\Local\Temp\ipykernel_11988\195219963.py:1: UserWarning:  
  
'distplot' is a deprecated function and will be removed in seaborn v0.14.0.
```



jupyter Real Estate Project Last Checkpoint: Last Thursday at 11:37 AM (autosaved)

```
In [56]: df_train['pop'].describe()
Out[56]:
count    27321.000000
mean     4316.012685
std      2169.226173
min      0.000000
25%    2885.000000
50%    4042.000000
75%    5400.000000
max    53812.000000
Name: pop, dtype: float64
```

```
In [57]: df_train['pop_bins']=pd.cut(df_train['pop'],bins=5,labels=['very low','low','medium','high','very high'])
In [58]: df_train[['pop','pop_bins']]
```

```
Out[58]:
   pop    pop_bins
   UID
26722  5230  very low
246444 2633  very low
245683 3881  very low
279653 2706  very low
247218 5637  very low
```

71°F Sunny 1:57 PM 8/29/2023

jupyter Real Estate Project Last Checkpoint: Last Thursday at 11:37 AM (autosaved)

```
In [59]: df_train['pop_bins'].value_counts()
Out[59]:
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 [60]: df_train.groupby(by='pop_bins')[['married','separated','divorced']].count()
Out[60]:
   married  separated  divorced
   pop_bins
   very low    27058    27058    27058
   low        246      246      246
```

71°F Sunny 1:57 PM 8/29/2023

jupyter Real Estate Project Last Checkpoint: Last Thursday at 11:37 AM (autosaved)

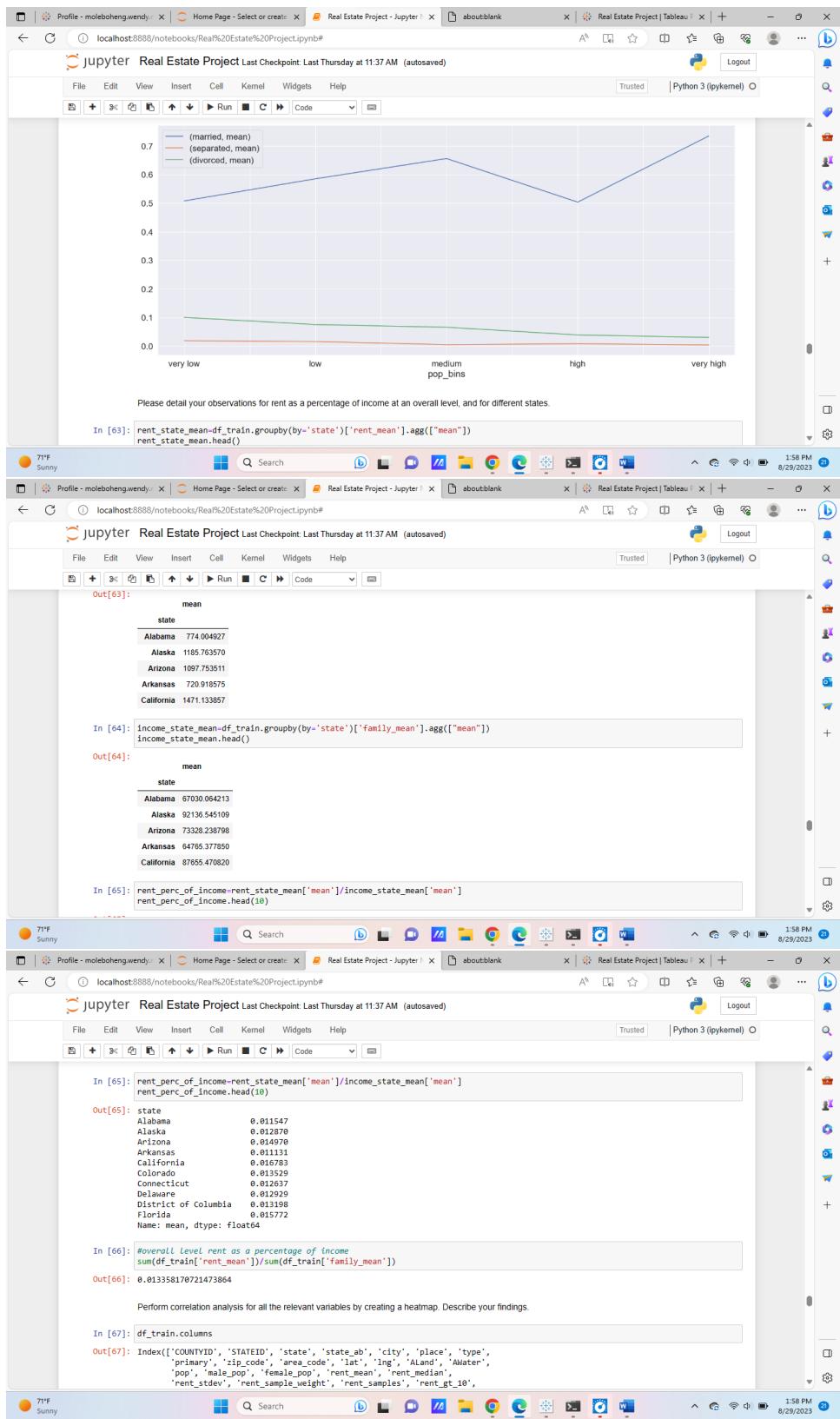
```
In [61]: df_train.groupby(by='pop_bins')[['married','separated','divorced']].agg(['mean', 'median'])
Out[61]:
   married           separated           divorced
   mean    median    mean    median    mean    median
   pop_bins
   very low  0.507548  0.524680  0.019126  0.013650  0.100504  0.096020
   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  0.030380  0.030380
```

Visualize using appropriate chart type

```
In [62]: plt.figure(figsize=(10,5))
pop_bin_married=df_train.groupby(by='pop_bins')[['married','separated','divorced']].agg(['mean'])
pop_bin_married.plot(figsize=(28,8))
plt.legend(loc='best')
plt.show()
```

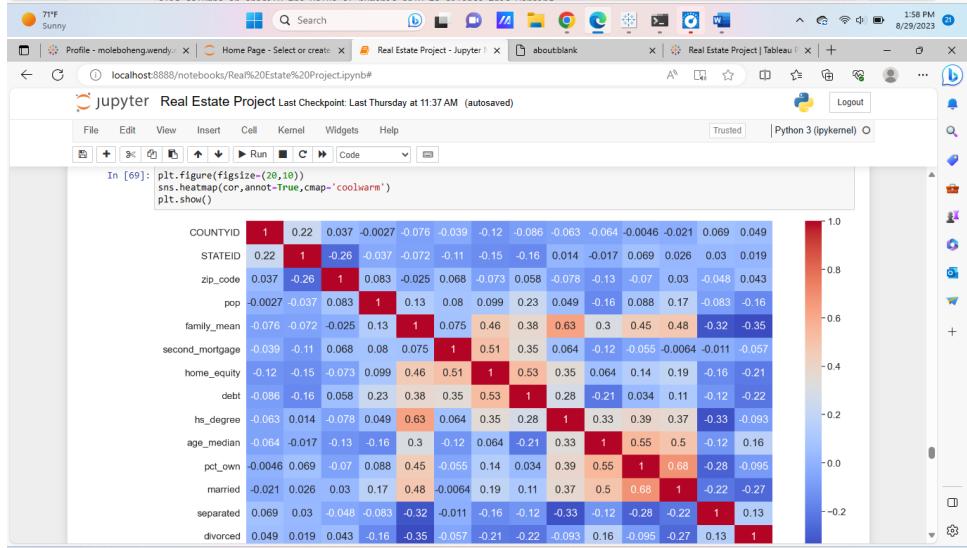
<Figure size 1000x500 with 0 Axes>

71°F Sunny 1:58 PM 8/29/2023



```
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_mean_median',
       '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_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_ag_mean',
       'male_ag_stdev', 'male_ag_sample_weight',
       'male_ag_samples', 'female_ag_mean', 'female_ag_median',
       'female_ag_stdev', 'female_ag_sample_weight', 'female_ag_samples',
       'married', 'separated', 'married_sep', 'separated_sep',
       'divorced', 'bad_debt', 'bins', 'pop_density', 'age_median', 'pop_bins'],
      dtype='object')

In [68]: cor_df_train[['COUNTYID','STATEID','zip_code','type','pop','family_mean',
       'second_mortgage','home_equity','debt','hs_degree',
       'age_mean','pct_own','married','separated']].corr()
C:\Users\Student_0002\AppData\Local\Temp\ipykernel_11988\3214557709.py:3: FutureWarning:
The default value of numeric_only in DataFrame.corr is deprecated. In a future version, it will default to False. Select only v
```



High positive correlation is noticed between pop, male_pop and female_pop
 High positive correlation is noticed between rent_mean_hi_mean, family_mean_hc_mean
 Project Task: Week 2
 Data Pre-processing:
 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.
 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.

```
In [70]: !pip install factor_analyzer
Requirement already satisfied: factor_analyzer in c:\users\student_0002\anaconda3\lib\site-packages (0.5.0)
```

In [70]:

```
pip install factor_analyzer
Requirement already satisfied: factor_analyzer in c:\users\student_0002\anaconda3\lib\site-packages (0.5.0)
Requirement already satisfied: numpy in c:\users\student_0002\anaconda3\lib\site-packages (from factor_analyzer) (1.23.5)
Requirement already satisfied: pre-commit in c:\users\student_0002\anaconda3\lib\site-packages (from factor_analyzer) (3.3.3)
Requirement already satisfied: scipy in c:\users\student_0002\anaconda3\lib\site-packages (from factor_analyzer) (1.10.0)
Requirement already satisfied: pandas in c:\users\student_0002\anaconda3\lib\site-packages (from factor_analyzer) (1.5.3)
Requirement already satisfied: scikit-learn in c:\users\student_0002\anaconda3\lib\site-packages (from factor_analyzer) (1.2.1)
Requirement already satisfied: pytz>=2020.1 in c:\users\student_0002\anaconda3\lib\site-packages (from pandas->factor_analyzer)
Requirement already satisfied: python-dateutil>=2.8.1 in c:\users\student_0002\anaconda3\lib\site-packages (from pandas->factor_analyzer) (2.8.2)
Requirement already satisfied: cfgv>=2.0.0 in c:\users\student_0002\anaconda3\lib\site-packages (from pre-commit->factor_analyzer) (3.4.0)
Requirement already satisfied: virtualenv>=20.10.0 in c:\users\student_0002\anaconda3\lib\site-packages (from pre-commit->factor_analyzer) (20.24.3)
Requirement already satisfied: identify>=1.0.0 in c:\users\student_0002\anaconda3\lib\site-packages (from pre-commit->factor_analyzer) (2.5.27)
Requirement already satisfied: nodeenv>=0.11.1 in c:\users\student_0002\anaconda3\lib\site-packages (from pre-commit->factor_analyzer) (1.8.0)
Requirement already satisfied: pyyaml>=5.1 in c:\users\student_0002\anaconda3\lib\site-packages (from pre-commit->factor_analyzer) (6.0)
Requirement already satisfied: joblib>=1.1.1 in c:\users\student_0002\anaconda3\lib\site-packages (from scikit-learn->factor_analyzer) (1.1.1)
Requirement already satisfied: threadpoolctl>=2.0.0 in c:\users\student_0002\anaconda3\lib\site-packages (from scikit-learn->factor_analyzer) (2.2.0)
Requirement already satisfied: setuptools in c:\users\student_0002\anaconda3\lib\site-packages (from nodeenv>=0.11.1->pre-commit->factor_analyzer) (65.6.3)
Requirement already satisfied: six>=1.5 in c:\users\student_0002\anaconda3\lib\site-packages (from python-dateutil>=2.8.1->pandas->factor_analyzer) (1.16.0)
Requirement already satisfied: distlib<1,>=0.3.7 in c:\users\student_0002\anaconda3\lib\site-packages (from virtualenv>=20.10.0)
```

In [71]:

```
from sklearn.decomposition import FactorAnalysis
from factor_analyzer import FactorAnalyzer
```

In [72]:

```
fa=FactorAnalyzer(n_factors=5)
fa.fit_transform(df_train.select_dtypes(exclude= ('object','category')))
fa.loadings_
```

Out[72]:

```
array([[-1.12589166e-01, 1.95646474e-02, -2.39331091e-02,
       -6.27632651e-02, 4.23474792e-02],
      [-1.10186761e-01, 1.3356226e-02, 2.79651227e-02,
       -1.10186761e-01, 1.3356226e-02],
      [-8.2087652e-02, 5.16372376e-02, -1.36451866e-01,
      -4.98918619e-02, -1.04024837e-01],
      [-1.80961179e-02, 1.92013761e-02, 5.81329907e-03,
      2.64842763e-02, -6.12442299e-03],
      [ 9.02324706e-02, -9.72542429e-02, -6.54681291e-02,
      -1.33116561e-01, -1.04024837e-01],
      [-1.47323561e-02, -1.13376892e-02, -1.45853491e-01,
      8.80433440e-03, 1.003217577e-03],
      [-4.28796967e-02, -2.09780215e-02, 3.66726850e-02,
      -9.45597448e-02, 5.91380552e-02],
      [-2.44242051e-03, -1.53245410e-02, -2.68308947e-03,
      -4.52473058e-02, 2.37240673e-02],
      [ 7.92164335e-02, 9.57453338e-01, -8.71151690e-02,
      -6.59923712e-03, -3.97273189e-02],
      [ 7.39808227e-02, 9.18790522e-01, -1.08834844e-01,
```

Data Modeling :

In [73]:

```
df_train.columns
```

Out[73]:

```
Index(['COUNTYID', 'STATEID', 'state', 'state_ab', 'city', 'place', 'type',
       'primary', 'zip_code', 'area_code', 'lat', 'long', 'Aland', 'Water',
       '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',
       'fam_mean', 'fam_median', 'fam_stdev', 'fam_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_mean', 'hs_degree_median', 'male_age_mean',
       'male_age_stdev', 'male_age_sample_weight', 'female_age_mean',
       'female_age_stdev', 'female_age_sample_weight', 'female_age_samples',
       'pct_own', 'married', 'married_sm', 'separated', 'divorced',
       'bad_debt', 'bins', 'pop_density', 'age_median', 'pop_bins'],
      dtype='object')
```

jupyter Real Estate Project Last Checkpoint: Last Thursday at 11:37 AM (autosaved)

```
In [74]: df_train['type'].unique()
type_dict={ 'City':1,
            'Urban':2,
            'Town':3,
            'CDP':4,
            'Village':5,
            'Borough':6}
df_train.replace(type_dict,inplace=True)

In [75]: df_train['type'].unique()
Out[75]: array([1, 2, 3, 4, 5, 6], dtype=int64)

In [76]: df_test.replace(type_dict,inplace=True)

In [77]: df_test['type'].unique()
Out[77]: array([4, 1, 6, 3, 5, 2], dtype=int64)

In [78]: feature_cols=['COUNTYID','STATEID','zip_code','pop','family_mean',
                     'second_mortgage','home_equity','debt','hs_degree',
                     'age_median','pct_own','married','separated','divorced']

In [79]: x_train=df_train[feature_cols]
y_train=df_train['hc_mortgage_mean']

In [80]: x_test=df_test[feature_cols]
```

71°F Sunny 1:59 PM 8/29/2023

jupyter Real Estate Project Last Checkpoint: Last Thursday at 11:37 AM (autosaved)

```
In [79]: x_train=df_train[feature_cols]
y_train=df_train['hc_mortgage_mean']

In [80]: x_test=df_test[feature_cols]
y_test=df_test['hc_mortgage_mean']

In [81]: 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 [82]: x_train.head()
Out[82]:
   COUNTYID STATEID zip_code type pop family_mean second_mortgage home_equity debt hs_degree age_median pct_own married separate
   UID
267822    53     36  13346    1  5230  67994.14790    0.02077  0.08919  0.52963  0.89288  44.666665  0.79046  0.57851  0.012
246444    141    18  46616    1  2633  50670.10337    0.02222  0.04274  0.60855  0.90487  34.791665  0.52483  0.34886  0.014
245683    63     18  46122    1  6881  95262.51431    0.00000  0.09512  0.73484  0.94288  41.833330  0.85331  0.64745  0.016
279653    127    72  927    2  2700  56401.88133    0.01086  0.01086  0.52714  0.91500  49.750000  0.65037  0.47257  0.020
247218    161    20  66502    1  5637  54053.42396    0.05426  0.05426  0.51938  1.00000  22.000000  0.13046  0.12356  0.000
```

In [83]: sc=StandardScaler()
x_train_scaled=sc.fit_transform(x_train)
x_test_scaled=sc.fit_transform(x_test)

71°F Sunny 1:59 PM 8/29/2023

jupyter Real Estate Project Last Checkpoint: Last Thursday at 11:37 AM (autosaved)

```
In [83]: 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 [84]: linerreg=LinearRegression()
linerreg.fit(x_train_scaled,y_train)
Out[84]: <LinearRegression>
LinearRegression()

In [85]: y_pred=linerreg.predict(x_test_scaled)

In [86]: 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 323.1018894984635

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

In [87]: state=df_train['STATEID'].unique()
state[0:5]
#Picking a few IDs 20,1,45,6
```

71°F Sunny 2:00 PM 8/29/2023

jupyter Real Estate Project Last Checkpoint: Last Thursday at 11:37 AM (autosaved)

```

Out[87]: array([36, 18, 72, 20, 1], dtype=int64)

In [88]:
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 :- 0.6046603766461811
Overall RMSE of linear regression model for state, 20 :- 307.9718899931475

State ID- 1
Overall R2 score of linear regression model for state, 1 :- 0.8104382475484616
Overall RMSE of linear regression model for state, 1 :- 307.82758168484354

```

To check the residuals

```

In [89]:
residuals=y_test - y_pred
residuals

Out[89]:
UID
255504   281.969888
252676   -69.935775
252114   197.758949
248814   -157.209637
280805   -9.887017
...
238888   -67.541646
242811   -41.57857
242127   -127.44569
241996   -310.68845
287763   217.766642
Name: hc_mortgage_mean, Length: 11700, dtype: float64

In [90]:
plt.hist(residuals) # Normal distribution of residuals

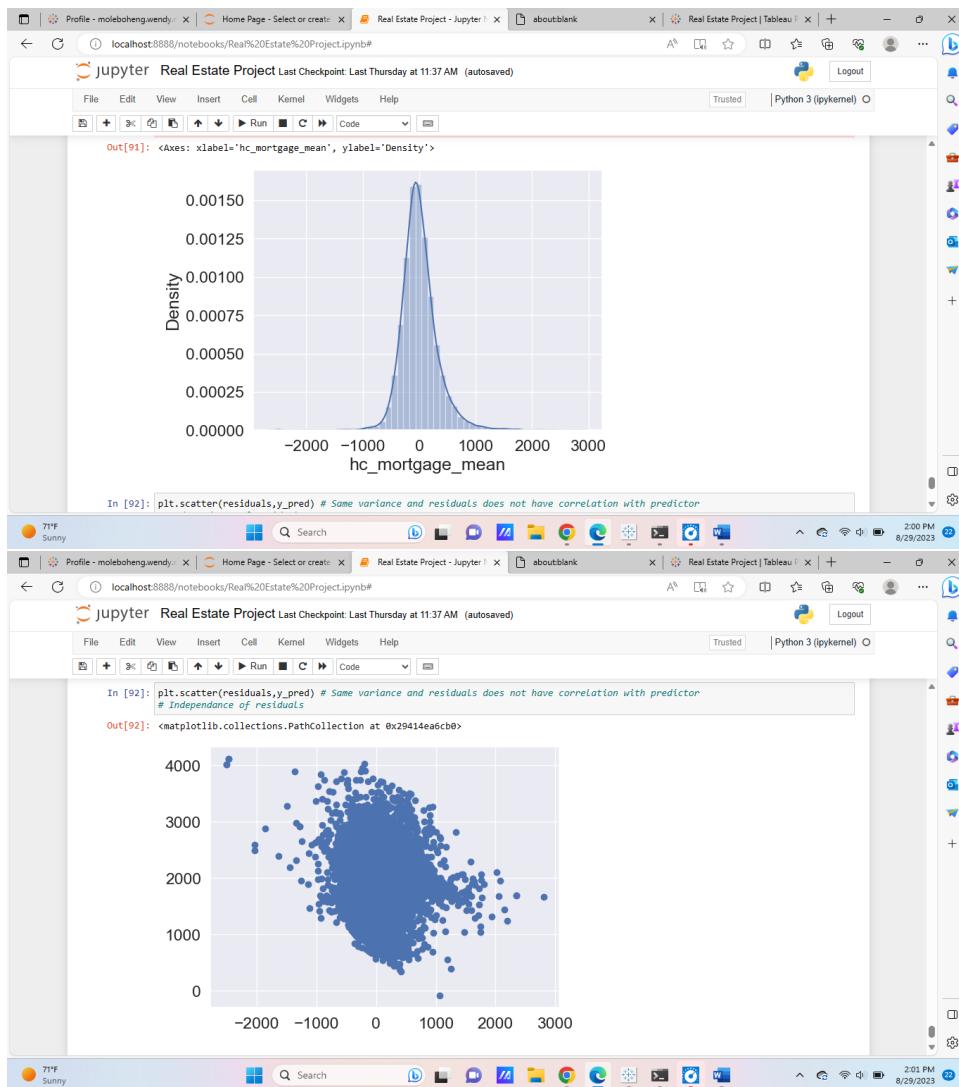
```

```

Out[90]:
(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([-2515.37284513, -1910.92311329, -1450.81038425, -918.69415521,
        398.57961217,  245.8389367,  677.65453151, 1289.77076695,
       1741.88698999,  2274.08321983,  2888.11944807]),

<BarContainer object of 10 artifacts>

```



Screenshots Of Tableau Public Real Estate Project

Tableau Public - Real Estate Project

File Data Window Help

Connections Add

train Text file

Files

Use Data Interpreter

Data Interpreter might be able to clean your Text file workbook.

test.csv
train.csv

New Union
New Table Extension

train

train.csv

Need more data?
Drag tables here to relate them. [Learn more](#)

train.csv 82 fields 27321 rows

100 rows

Name train.csv

Fields

Type	Field Name	Physical Table	Remote Field Name
#	UID	train.csv	UID

train.csv train.csv train.csv train.csv train.csv train.csv

UID Blockid Sumlevel Countyid Stateid train.csv State Ab

267822 null 140 53 36 New York NY
246444 null 140 141 18 Indiana IN
245683 null 140 63 18 Indiana IN
279653 null 140 127 72 Puerto Rico PR
247718 null 140 161 29 Kansas KS

Data Source Box Plot Debt vs Bad Credit Geo Map HeatMap Population Distribution By Type Real Estate Dashboard

Tableau Public - Real Estate Project

File Data Worksheet Dashboard Story Analysis Map Format Server Window Help

Data Analytics

Pages

Entire View

Box Plot

AVG(Rent Mean)

Box Plot

Type

Avg Rent Mean

Borough CDP City Town Urban Village

25 nulls

Box Plot

Sum of AVG(Rent Mean): 7,825.509

Tableau Public - Real Estate Project

File Data Worksheet Dashboard Story Analysis Map Format Server Window Help

Data Analytics

Pages

Entire View

Debt vs Bad Credit

Measure Names

Bad Debt Debt

Measure Values

2,923 16,903 Bad Debt
2,923 16,903 Debt

Pie

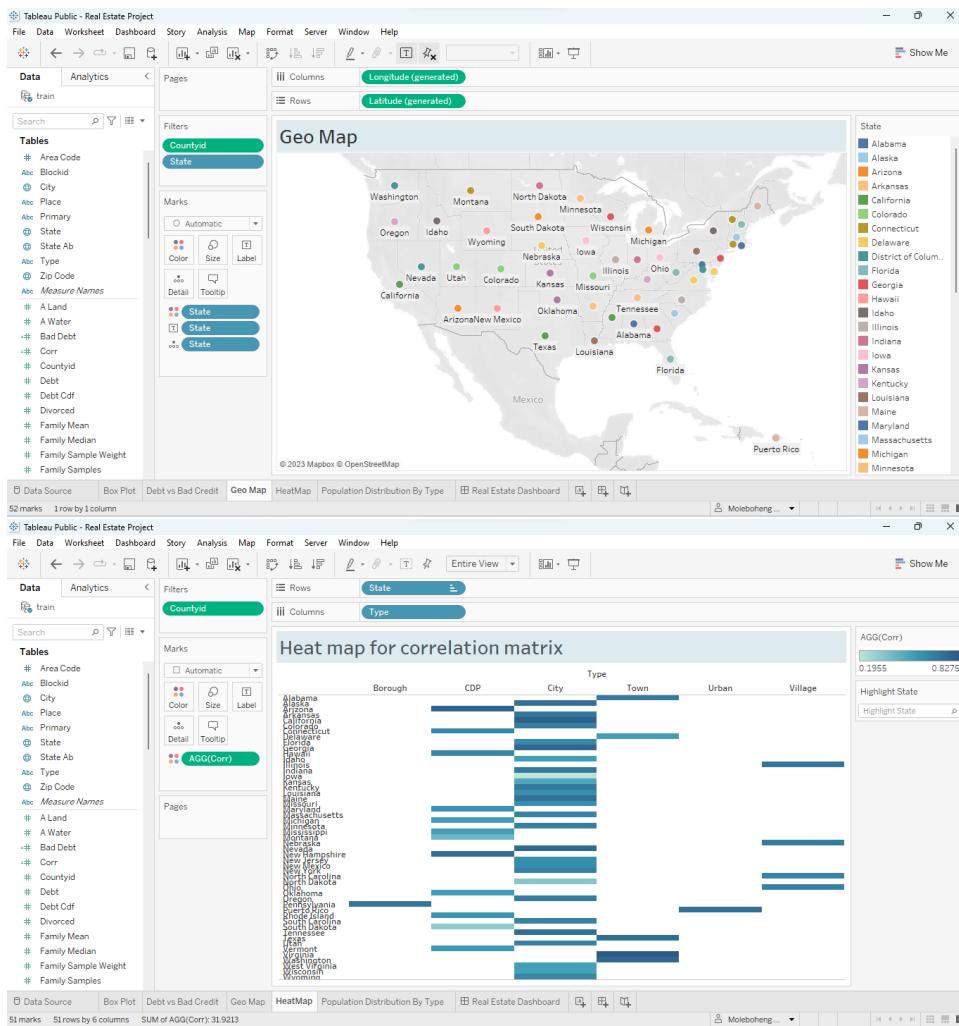
Measure Names

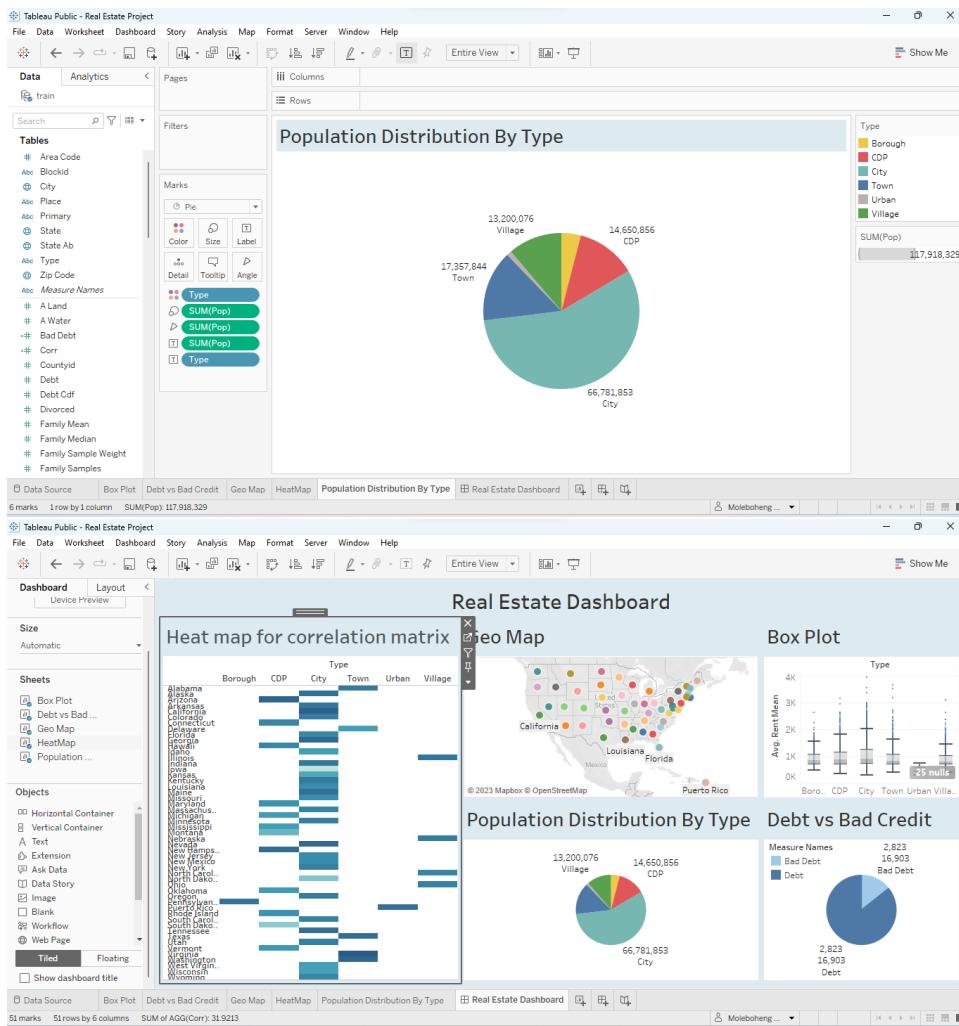
Bad Debt Debt

Measure Values

2,923 16,903 Bad Debt
2,923 16,903 Debt

2 marks 1 row by 1 column SUM(Bad Debt): 2,923





Link To Tableau Real Estate Project:

https://public.tableau.com/app/profile/moleboheng.wendy.mokoena/viz/RealEstateProject_16933092279690/BoxPlot