

```
1 #%% md
2 # Real Estate - PGP - Capstone Project
3
4 ## Problem Statement
5 - A banking institution requires actionable insights
  into mortgage-backed securities, geographic business
  investment, and real estate analysis.
6 - 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.
7 - 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.
8
9 ## Dataset Description
10
11 - Second mortgage: Households with a second mortgage
  statistics
12 - Home equity: Households with a home equity loan
  statistics
13 - Debt: Households with any type of debt statistics
14 - Mortgage Costs: Statistics regarding mortgage
  payments, home equity loans, utilities, and property
  taxes
15 - Home-Owner Costs: Sum of utilities, and property
  taxes statistics
16 - Gross Rent: Contract rent plus the estimated
  average monthly cost of utility features
17 - High school Graduation: High school graduation
  statistics
18 - Population Demographics: Population demographics
  statistics
19 - Age Demographics: Age demographic statistics
20 - Household Income: Total income of people residing
  in the household
21 - Family Income: Total income of people related to
```

```
21 the householder
22 #%% md
23 ### 1. Import Data
24 #%%
25 #IMPORT PACKAGES AND DATA
26
27 import numpy as np
28 import pandas as pd
29 import os
30 import warnings
31 warnings.filterwarnings('ignore')
32
33 import matplotlib.pyplot as plt
34 import seaborn as sns
35 sns.set(style= 'white', color_codes = True)
36 sns.set(font_scale = 1.5)
37 from itertools import cycle
38 #%%
39 train_df = pd.read_csv('train.csv')
40 test_df = pd.read_csv('test.csv')
41
42 train_df.head()
43 #%%
44 train_df.columns
45 #%%
46 test_df.columns
47 #%%
48 len(train_df)
49 #%%
50 train_df.describe()
51 #%%
52 train_df.info()
53 #%% md
54 ### 2. Find the primary key and locate the indexing
      requirement
55 #%%
56 #UID is the unique identifier in the test and
      training set data. So an index can be created from
```

```

56 the UID Feature
57 train_df.set_index(keys = ['UID'], inplace=True)
58 test_df.set_index(keys = ['UID'], inplace= True)
59 #%%
60 train_df.head()
61 #%%
62 test_df.head()
63 #%% md
64 ### 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.
65 #%%
66 #PERCENTAGE TOTALS OF MISSING VALUES IN TEST AND
       TRAINING DATASETS
67
68 mvt_train = train_df.isnull().sum() *100/len(train_df)
   )
69 mvt_train_df = pd.DataFrame(mvt_train, columns = [
   'Percentage of Values Missing'])
70 mvt_train_df.sort_values(by=['Percentage of Values
       Missing'],inplace=True, ascending=False)
71 mvt_train_df[mvt_train_df['Percentage of Values
       Missing'] > 0][:10]
72 #%%
73 mvt_test = test_df.isnull().sum() *100/len(train_df)
74 mvt_test_df = pd.DataFrame(mvt_test, columns = [
   'Percentage of Values Missing'])
75 mvt_test_df.sort_values(by=['Percentage of Values
       Missing'],inplace=True, ascending=False)
76 mvt_test_df[mvt_test_df['Percentage of Values Missing
       '] > 0][:10]
77 #%%
78 train_df['SUMLEVEL'].unique() #SUM LEVEL HAS ONLY ONE
       VALUE AND WILL NOT HELP WITH DIFFERENTIATING RESULTS
79 #%% md
80 **Based on the missing value analysis, we will drop
       the BLOCKID value entirely since the TRAIN_DF sample

```

```
80 has 100% NULL values. We will also impute the  
missing values for the other columns with the mean  
of the data since the values for the dataset are  
negligible. (Note: If the dataset were larger we  
would drop the missing queries as they would not  
impact the study from the models.**  
81 #%%  
82 train_df.drop(columns = ['BLOCKID', 'SUMLEVEL'],  
    inplace=True)  
83 test_df.drop(columns = ['BLOCKID', 'SUMLEVEL'],  
    inplace=True)  
84 #%%  
85 #IMPUTE THE MEAN  
86  
87 missing_train_cols = []  
88  
89 for col in train_df.columns:  
90     if train_df[col].isna().sum() !=0:  
91         missing_train_cols.append(col)  
92     print(missing_train_cols)  
93 #%%  
94 missing_test_cols = []  
95  
96 for col in test_df.columns:  
97     if test_df[col].isna().sum() !=0:  
98         missing_test_cols.append(col)  
99     print(missing_test_cols)  
100 #%%  
101 for col in train_df.columns:  
102     if col in (missing_train_cols):  
103         train_df[col].replace(np.nan, train_df[col].  
            mean(), inplace=True)  
104 #%%  
105 for col in test_df.columns:  
106     if col in (missing_test_cols):  
107         test_df[col].replace(np.nan, test_df[col].  
            mean(), inplace=True)  
108 #%%
```

```
109 train_df.isna().sum().sum()
110 #%%
111 test_df.isna().sum().sum()
112 #%% md
113 **ALL NULL VALUES HAVE BEEN REMOVED**
114 #%% md
115 ## Exploratory Data Analysis:
116
117 ### 4. Perform Debt analysis and Study Population
Density
118
119 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
120 #%%
121 #GETTING A SMALLER QUERIED DATASET WITH REQUIRED
ELEMENTS TO EXPLORE ITEM (A)
122
123 from pandasql import sqldf
124 q1 = "select place,pct_own,second_mortgage,lat,lng
from train_df where pct_own >0.10 and
second_mortgage <0.5 order by second_mortgage DESC
LIMIT 2500;"
125 pymysql_df = lambda q: sqldf(q, globals())
126 train_df_location_mort_pct=pymysql_df(q1)
127 #%%
128 train_df_location_mort_pct.head()
129 #%%
130 import plotly.express as px
131 import plotly.graph_objects as go
132 #%%
133 #VISUALIZE USING PLOTLY GEO-MAPPING PACKAGE. <USED
REFERENCE CODE FROM GITHUB>
134
135 fig = go.Figure(data = go.Scattergeo(
```

```
136     lat = train_df_location_mort_pct['lat'],
137     lon = df_train_location_mort_pct['lng']
138 ),)
139
140 fig.update_layout(
141     geo=dict(
142         scope = 'north america',
143         showland = True,
144         landcolor= 'rgb(212,212,212)',
145         subunitcolor ='rgb(255,255,255)',
146         showlakes = True,
147         lakecolor = 'rgb(255,255,255)',
148         showsubunits = True,
149         showcountries = True,
150         resolution = 50,
151         projection = dict(
152             type = 'conic conformal',
153             rotation_lon = -100
154         ),
155         lonaxis = dict(
156             showgrid = True,
157             gridwidth = 0.5,
158             range = [-140.0, -55.0],
159             dtick = 5
160         ),
161         lataxis = dict(
162             showgrid = True,
163             gridwidth = 0.5,
164             range = [20.0, 60.0],
165             dtick = 5
166         )
167     ),
168     title = 'Top 2,500 locations with Highest rate
169       of Second Mortgages and Percent Ownership is Above
170       10%'
171 )
170 fig.show()
171 #%% md
```

```

172 **b) Use the following bad debt equation: Bad Debt
    = P (Second Mortgage ∩ Home Equity Loan) Bad Debt
    = second_mortgage + home_equity -
    home_equity_second_mortgage c) Create pie charts to
    show overall debt and bad debt**
173 #%%
174 train_df['bad_debt'] = train_df['second_mortgage']
    ] + train_df['home_equity'] - train_df['
    home_equity_second_mortgage']
175 #%%
176 train_df['bins'] = pd.cut(train_df['bad_debt'], bins
    = [0,0.1,1], labels = ['less than 50%', 'greater
    than 50%'])
177 train_df.groupby(['bins']).size().plot(kind='pie',
    subplots=True, startangle =90, autopct = '%1.1f%%')
178 plt.axis('equal')
179 plt.title('Bad Debt Distribution')
180
181 plt.show()
182 #%% md
183 **d) Create Box and whisker plot and analyze the
    distribution for 2nd mortgage, home equity, good
    debt, and bad debt for different cities**
184 #%%
185 #SELECT TWO CITIES TO COMPARE
186 train_df['city'].unique()
187
188 #PANTEGO and STANTON CITY
189 #%%
190 cols = ['second_mortgage', 'home_equity', 'debt', ' '
    'bad_debt']
191 box_pantego = train_df.loc[train_df['city'] == ' '
    'Trinity']
192 box_stantoncity = train_df.loc[train_df['city'] == ' '
    'Nocona']
193 box_cities = pd.concat([box_pantego,box_stantoncity
    ])
194 box_cities.head()

```

```
195 #%%
196 plt.figure(figsize=(7,5))
197 sns.boxplot(data=box_cities,x='second_mortgage', y='city',width=0.5,palette="Set3")
198 plt.show()
199
200 plt.figure(figsize=(7,5))
201 sns.boxplot(data=box_cities,x='home_equity', y='city',width=0.5,palette="Set3")
202 plt.show()
203
204 plt.figure(figsize=(7,5))
205 sns.boxplot(data=box_cities,x='debt', y='city',width=0.5,palette="Set3")
206 plt.show()
207
208 plt.figure(figsize=(7,5))
209 sns.boxplot(data=box_cities,x='bad_debt', y='city',width=0.5,palette="Set3")
210 plt.show()
211 #%% md
212 **e) Create a collated income distribution chart for
      family income, house hold income, and remaining
      income**
213 #%%
214 sns.distplot(train_df['hi_mean'])
215 plt.title('Household Income Distribution Chart')
216 plt.show()
217
218 sns.distplot(train_df['family_mean'])
219 plt.title('Family Income Distribution Chart')
220 plt.show()
221
222 sns.distplot(train_df['family_mean']-train_df['hi_mean'])
223 plt.title('Remaining Income Distribution Chart')
224 plt.show()
225 #%% md
```

```
226 **f) Use pop and ALand variables to create a new  
field called population density**  
227  
228 ALand - Land Square Footage  
229 pop - Male and Female Population  
230 #%%  
231 train_df['pop_density'] = train_df['pop']/train_df['  
ALand']  
232 test_df['pop_density'] = test_df['pop']/test_df['  
ALand']  
233 #%%  
234 sns.displot(train_df['pop_density'])  
235 plt.title('Population Density')  
236 plt.show()  
237 #%% md  
238 **b) 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**  
239 #%%  
240 train_df['age_median']=(train_df['male_age_median']+  
train_df['female_age_median'])/2  
241  
242 test_df['age_median']=(test_df['male_age_median']+  
test_df['female_age_median'])/2  
243  
244 train_df[['male_age_median','female_age_median', '  
male_pop','female_pop','age_median']]  
245 #%%  
246 sns.displot([train_df['age_median']])  
247 plt.title('Median Age')  
248 plt.show()  
249 #%% md  
250 **FINDINGS:**  
251  
252 - Median Age is concentrated between ages 20 and 60  
, with a majority around the age of 40.  
253 - Median age distribution has a Gaussian
```

```
253 Distribution
254 - Slight Right Skewness is noted
255 #%%
256 train_df['pop_bins'] = pd.cut(train_df['pop'], bins
257     = 5, labels = ['lowest','low','medium', 'high', 'highest'])
258 #%%
259 train_df[['pop', 'pop_bins']]
260 #%%
261 train_df['pop_bins'].value_counts()
262 #%% md
263 **Analyze the married, separated, and divorced
264 population for these population brackets. Visualize
265 using appropriate chart type**
266 #%%
267 train_df.groupby(by='pop_bins')[['married', 'separated', 'divorced']].count()
268 #%%
269 plt.figure(figsize=(20,10))
270 pop_bin_married=train_df.groupby(by='pop_bins')[['married', 'separated', 'divorced']].agg(["mean"])
271 pop_bin_married.plot(figsize = (20,10))
272 plt.legend(loc = 'best')
273 plt.show()
274 #%%
275 state_rent_mean = train_df.groupby(by = 'state')[['rent_mean']].agg(['mean'])
276 state_income_mean = train_df.groupby(by = 'state')[['family_mean']].agg(['mean'])
277
278 percentage_rent_income = state_rent_mean['mean']/state_income_mean['mean']
279 #%%
280 percentage_rent_income.head(10)
```

```
281 #%%
282 #DEPICT LEVEL OF RENT AS AN OVERALL PERCENTAGE OF
INCOME
283 print(round(sum(train_df['rent_mean']) / sum(
train_df['family_mean']) *100),'%')
284 %% md
285 ##### 5. Perform correlation analysis for all the
relevant variables by creating a heatmap. Describe
your findings.
286 %%
287 train_df.columns
288 %%
289 #SELECT DATA TO CORRELATE USING HEATMAP
290 co = train_df[['COUNTYID','STATEID','zip_code','type',
'pop', 'family_mean',
291           'second_mortgage', 'home_equity', ' '
debt','rent_median',
292           'age_median','hs_degree','pct_own', ' '
married','separated', 'divorced']].corr()
293 %%
294 plt.figure(figsize= (20,10))
295 sns.heatmap(co, annot= True, cmap = 'plasma')
296 plt.show()
297
298 %% md
299 **NOTES FROM HEATMAP**
300 - High correlation seen between hs_degree and
family_mean
301 - High positive correlation seen between pct_down
and married
302 %% md
303 ## DATA PRE-PROCESSING
304
305 ### 1. The multivariate data has a significant
number of measured variables, not all of which have
an impact on the selected predicted variable. The
goal is to find the measured variables that would be
dependent on smaller unobserved common factors of
```

```
305 latent variables.  
306  
307 ### See the common factors below and plotted  
loadings below. Use factor analysis to find the  
latent variables within the dataset and gain  
insights on the linear relationships within the data  
. .  
308  
309 **Latent Variables:**  
310 - High School graduation rates  
311 - Median population age  
312 - Second mortgage stats  
313 - Percent own  
314 - Bad debt expense  
315 #%%  
316 from sklearn.decomposition import _factor_analysis  
317 from factor_analyzer import FactorAnalyzer  
318 #%%  
319 fa = FactorAnalyzer(n_factors=5)  
320 fa.fit_transform(train_df.select_dtypes(exclude=[  
    'object', 'category']))  
321 fa.loadings_  
322 #%% md  
323 ## Linear Regression Model  
324  
325 ### 1. This linear regression model serves to  
predict the total monthly expenditure for the home  
mortgages loan. This is the mean monthly mortgage  
and owner costs of a specified geographical location  
. .  
326  
327 **NOTE: All NaN values have been removed from the  
dataset.**  
328 #%%  
329 train_df.columns  
330 #%%  
331 train_df['type'].unique()  
332 type_dict = {'type_2':{'City':1},
```

```
333             'Urban': 2,
334             'Town': 3,
335             'CDP': 4,
336             'Village': 5,
337             'Borough': 6}
338         }
339 train_df['type_2'] = train_df[type_dict,inplace =
340 True]
340 #%%
341 train_df['type_2'].unique()
342 #%%
343 test_df.replace(type_dict,inplace = True)
344 #%%
345 test_df['type'].unique()
346 #%%
347 feature_cols = ['COUNTYID','STATEID','zip_code','
348 type','pop','family_mean','second_mortgage','
349 home_equity','debt','rent_median','age_median','
350 hs_degree','pct_own','married','separated','
351 divorced']
352 X_train = train_df[feature_cols]
353 y_train = train_df['hc_mortgage_mean']
354
355 from sklearn.preprocessing import StandardScaler
356 from sklearn.linear_model import LinearRegression
357 from sklearn.metrics import r2_score,
358     mean_absolute_error,mean_squared_error,
359     accuracy_score
360 #%%
361 X_train.head()
362 #%%
363 scaler = StandardScaler()
364 X_train_scaled = scaler.fit_transform(X_train)
365 X_test_scaled = scaler.fit_transform(X_test)
```

```
364 #%% md
365 **a) Run a model at a Nation level. If the accuracy
   levels and R square are not satisfactory proceed to
   below step.**
366 #%%
367 lin_reg = LinearRegression()
368 lin_reg.fit(X_train_scaled , y_train)
369 #%%
370 y_pred = lin_reg.predict(X_test_scaled)
371 #%%
372 print('Overall R2 Score of linear regression model:'
      , r2_score(y_test,y_pred))
373 print("Overall RMSE of linear regression model", np.
      sqrt(mean_squared_error(y_test,y_pred)))
374 #%% md
375 Model has run successfully and has good results but
   we will further investigate the model's performance
   at the state level.
376
377 **b) Run another model at State level.**
378 #%%
379 state = train_df['STATEID'].unique()
380 state[0:5]
381 #%%
382 for i in [20,1,45]:
383     print('State-ID:',i)
384
385     X_train_state = train_df[train_df['COUNTYID'
] == i][feature_cols]
386     y_train_state = train_df[train_df['COUNTYID'
] == i]['hc_mortgage_mean']
387
388     x_test_state = test_df[test_df['COUNTYID'] == i
][feature_cols]
389     y_test_state = test_df[test_df['COUNTYID'] == i
]['hc_mortgage_mean']
390
391     x_train_state_scaled = scaler.fit_transform(
```

```
391 X_train_state)
392     x_test_state_scaled = scaler.fit_transform(
393         x_test_state)
394     lin_reg.fit(x_train_state_scaled, y_train_state)
395     y_pred_nation = lin_reg.predict(
396         x_test_state_scaled)
397     print('Overall R2 Score of linear regression
398 model for state, ', i, ':', r2_score(y_test_state,
399 y_pred_nation) )
400     print('Overall RMSE of linear regression model
401 for state, ', i, ':', np.sqrt(mean_squared_error(
402 y_test_state,y_pred_nation)))
403     print('\n')
404 #%%
405 #checking error (residuals)
406 residuals = y_test - y_pred
407 residuals
408 #%%
409 plt.hist(residuals)
410 #%%
411 sns.displot(residuals)
412 #%%
413 plt.scatter(residuals, y_pred)
414 #%%
415 train_df.to_excel (r'/Users/Dishant/Documents/DATA
416 SCIENCE/Machine Learning Projects/Real Estate Data.
417 xlsx', index = None, header=True)
418 #%%
419 %pip install xgboost
420 #%%
421
```