

Project 1 Report

Prepared by: Agathiyam K

Roll.No: 045005

Project Title: Demographic Patterns and Home Loan Clustering: A K-means Analysis in Washington State

1. Project Objectives | Problem Statements 1.1. PO1 | PS1: Segmentation of Consumer Data using Unsupervised Machine Learning Clustering Algorithms 1.2. PO2 | PS2: Identification of Appropriate Number of Segments or Clusters 1.3. PO3 | PS3: Determination of Segment or Cluster Characteristics
2. **Description of Data** 2.1. Data Source, Size, Shape 2.1.1. Data Source (Website Link) <https://www.kaggle.com/datasets/miker400/washington-state-home-mortgage-hdma2016>

** 2.1.2. Data Size (in KB | MB | GB ...)**

30 MB

2.1.3. Data Shape (Dimension: Number of Variables | Number of Records)

37 Variables (37 columns) Maximum number of rows 60,000 Total number of records: 60000
Total number of filled cells: 1918780 Missed cells: 301220

2.2. Description of Variables 2.2.1. Index Variable(s): I1, I2, ... The index is the sequence_number. This column has 39340 unique count of variables

2.2.2. Variables or Features having Categories | Categorical Variables or Features (CV)
Categorical Columns: ['state_name', 'state_abbr', 'respondent_id', 'purchaser_type_name', 'property_type_name', 'preapproval_name', 'owner_occupancy_name', 'msamd_name', 'loan_type_name', 'loan_purpose_name', 'lien_status_name', 'hoepa_status_name', 'edit_status_name', 'denial_reason_name_3', 'denial_reason_name_2', 'denial_reason_name_1', 'county_name', 'co_applicant_sex_name', 'co_applicant_ethnicity_name', 'applicant_sex_name', 'applicant_ethnicity_name', 'agency_name', 'agency_abbr', 'action_taken_name']

Non-Categorical Columns: ['tract_to_msamd_income', 'rate_spread', 'population', 'minority_population', 'number_of_owner_occupied_units', 'number_of_1_to_4_family_units', 'loan_amount_000s', 'hud_median_family_income', 'applicant_income_000s', 'sequence_number', 'census_tract_number', 'as_of_year', 'application_date_indicator']

2.2.2.1. Variables or Features having Nominal Categories | Categorical Variables or Features - Nominal Type: CNV1, CNV2, ... All the categorical variables available in the dataset for nominal variables

2.2.2.2. Variables or Features having Ordinal Categories | Categorical Variables or Features - Ordinal Type: COV1, COV2, ... No ordinal data available in the dataset

2.2.3. Non-Categorical Variables or Features: NCV1, NCV2, ... Non-Categorical Columns:
['tract_to_msamd_income', 'rate_spread', 'population', 'minority_population',
'number_of_owner_occupied_units', 'number_of_1_to_4_family_units', 'loan_amount_000s',
'hud_median_family_income', 'applicant_income_000s', 'sequence_number',
'census_tract_number', 'as_of_year', 'application_date_indicator']

2.3. Descriptive Statistics

2.3.1. Descriptive Statistics: Categorical Variables or Features

2.3.1.1. Count | Frequency Statistics

Variable	Count	Unique
state_name	60000	1
state_abbr	60000	1
respondent_id	60000	593
purchaser_type_name	60000	10
property_type_name	60000	3
preapproval_name	60000	3
owner_occupancy_name	60000	3
msamd_name	51619	14
loan_type_name	60000	4
loan_purpose_name	60000	3
lien_status_name	60000	4
hoepa_status_name	60000	2
county_name	59908	39
co_applicant_sex_name	60000	5
co_applicant_ethnicity_n ame	60000	5
applicant_sex_name	60000	4
applicant_ethnicity_nam e	60000	4
agency_name	60000	6
agency_abbr	60000	6
action_taken_name	60000	8

Variable	Top Value
state_name	Washington
state_abbr	WA
respondent_id	32489
purchaser_type_name	Loan was not originated or was not sold in calendar year covered by register
property_type_name	One-to-four family dwelling (other than manufactured housing)
preapproval_name	Not applicable

Variable	Top Value
owner_occupancy_name	Owner-occupied as a principal dwelling
msamd_name	Seattle, Bellevue, Everett - WA
loan_type_name	Conventional
loan_purpose_name	Refinancing
lien_status_name	Secured by a first lien
hoepa_status_name	Not a HOEPA loan
county_name	King County
co_applicant_sex_name	No co-applicant
co_applicant_ethnicity_na	No co-applicant
me	
applicant_sex_name	Male
applicant_ethnicity_name	Not Hispanic or Latino
agency_name	Department of Housing and Urban Development
agency_abbr	HUD
action_taken_name	Loan originated

Variable	Frequency
state_name	60000
state_abbr	60000
respondent_id	5006
purchaser_type_name	16112
property_type_name	57630
preapproval_name	47832
owner_occupancy_name	53940
msamd_name	17965
loan_type_name	42917
loan_purpose_name	28576
lien_status_name	57046
hoepa_status_name	59996
county_name	12915
co_applicant_sex_name	26987
co_applicant_ethnicity_na	26987
me	
applicant_sex_name	37070
applicant_ethnicity_name	44014
agency_name	27514
agency_abbr	27514

Variable	Frequency
action_taken_name	55815

In the context of catdf dataset: state_name and state_abbr: These columns have only one unique value, which is "Washington" for state_name and "WA" for state_abbr. This suggests that these columns may not provide much information for analysis as they have constant values for all rows.

respondent_id: This column has 593 unique values, and the most frequent respondent_id is "32489" with a frequency of 5006. This column likely identifies different respondents.

Other categorical columns: Each column represents a categorical variable, and the summary provides information about the number of unique categories, the most frequent category (top), and its frequency.

action_taken_name: This column represents the action taken for the loan application. It has 8 unique values, and "Loan originated" is the most frequent action with a frequency of 55815.

2.3.1.2. Proportion (Relative Frequency) Statistics

Variable	Percentage
state_name	Washington: 100.00%
state_abbr	WA: 100.00%
respondent_id	32489: 8.34%
purchaser_type_name	Loan was not originated or was not sold in calendar year covered by register: 79.72%
property_type_name	One-to-four family dwelling (other than manufactured housing): 100.00%
preapproval_name	Not applicable: 79.72%
owner_occupancy_name	Owner-occupied as a principal dwelling: 89.90%
msamd_name	Seattle, Bellevue, Everett - WA: 34.80%
loan_type_name	Conventional: 71.53%
loan_purpose_name	Refinancing: 47.63%
lien_status_name	Secured by a first lien: 95.08%
hoepa_status_name	Not a HOEPA loan: 99.99%
county_name	King County: 21.56%
co_applicant_sex_name	No co-applicant: 44.98%
co_applicant_ethnicity_name	No co-applicant: 44.98%
applicant_sex_name	Male: 61.78%
applicant_ethnicity_name	Not Hispanic or Latino: 73.36%
agency_name	Department of Housing and Urban Development: 41.19%
agency_abbr	HUD: 45.86%
action_taken_name	Loan originated: 93.03%

2.3.2. Descriptive Statistics: Non-Categorical Variables or Features

2.3.2.1. Measures of Central Tendency and Dispersion

Variable	Count/ Std	Variable	Count/ Std
tract_to_msamd_income	59878	number_of_1_to_4_family_units	59878
		Mean: 107.617351 Std: 28.233471 Min: 14.050000 25%: 88.970001 50%: 105.550003 75%: 123.330002 Max: 257.140015	Mean: 1873.281456 Std: 738.505184 Min: 27.000000 25%: 1414.000000 50%: 1770.000000 75%: 2249.000000 Max: 5893.000000
population	59878	loan_amount_000s	60000
		Mean: 5278.782157 Std: 1716.101490 Min: 98.000000 25%: 4070.000000 50%: 5145.000000 75%: 6382.000000 Max: 13025.000000	Mean: 291.358717 Std: 604.958183 Min: 1.000000 25%: 170.000000 50%: 242.000000 75%: 337.000000 Max: 55000.000000
minority_population	59878	hud_median_family_income	59878
		Mean: 23.244442 Std: 14.416209 Min: 2.040000	Mean: 73869.411136 Std: 12811.243390 Min:

Variable	Count/ Std	Variable	Count/ Std
			48700.000000
	25%: 12.950000		25%: 63100.000000
	50%: 19.420000		50%: 73300.000000
	75%: 29.680000		75%: 90300.000000
	Max: 94.790001		Max: 90300.000000
number_of_owner_occupie d_units	59876	applicant_income_000s	53630
	Mean: 1399.044375		Mean: 112.822301
	Std: 518.330561		Std: 122.862496
	Min: 15.000000		Min: 1.000000
	25%: 1034.000000		25%: 61.000000
	50%: 1359.000000		50%: 89.000000
	75%: 1722.000000		75%: 132.000000
	Max: 2997.000000		Max: 6161.000000
sequence_number	60000	census_tract_number	59878
	Mean: 77526.47		Mean: 1750.597252
	Std: 150515.7		Std: 3359.676740
	Min: 1.000000		Min: 1.000000
	25%: 3231.500000		25%: 114.020000
	50%: 16481.000000		50%: 403.020000
	75%: 72762.25		75%: 713.100000
	Max: 1241590.0		Max: 9757.000000
as_of_year	60000	application_date_indicator	60000

Variable	Count/ Mean/ Std	Variable	Count/ Mean/ Std
	Mean: 2016.0		Mean: 0.025867
	Std: 0.0		Std: 0.225976
	Min: 2016.0		Min: 0.000000
	25%: 2016.0		25%: 0.000000
	50%: 2016.0		50%: 0.000000
	75%: 2016.0		75%: 0.000000
	Max: 2016.0		Max: 2.000000

2.3.2.3. Correlation Statistics (with Test of Correlation)

3. Analysis of Data 3.1. Data Pre-Processing 3.1.1. Missing Data Statistics and Treatment

3.1.1.1.1. Missing Data Statistics: Records Number of rows with missing data: 60000 Number of rows with more than 50% missing data: 0

3.1.1.1.2. Missing Data Treatment: Records 3.1.1.1.2.1. Removal of Records with More Than 50% Missing Data: None | R1, R2, ... No rows with more than 50% missing values

3.1.1.2.1. Missing Data Statistics: Categorical Variables or Features

Feature	Missing_Records	Percentage_Missing
denial_reason_name_3	59999	99.998333
denial_reason_name_2	59957	99.928333
denial_reason_name_1	59882	99.803333
rate_spread	58134	96.890000
edit_status_name	47549	79.248333
msamd_name	8381	13.968333
applicant_income_000s	6370	10.616667
number_of_owner_occu pied_units	124	0.206667
census_tract_number	122	0.203333
tract_to_msamd_incom e	122	0.203333
hud_median_family_inc ome	122	0.203333
number_of_1_to_4_fami ly_units	122	0.203333
minority_population	122	0.203333
population	122	0.203333

Feature	Missing_Records	Percentage_Missing
county_name	92	0.153333

3.1.1.2.2. Missing Data Treatment: Categorical Variables or Features

3.1.1.2.2.1. Removal of Variables or Features with More Than 50% Missing Data: None | CV1, CV2, ... Removed the below columns as they have more than 50% data missing • denial_reason_name_3 • denial_reason_name_2 • denial_reason_name_1 • rate_spread (non-cat) • edit_status_name

3.1.1.2.2.2. Imputation of Missing Data using Descriptive Statistics: Mode

3.1.1.3.1. Missing Data Statistics: Non-Categorical Variables or Features

Feature	Missing_Records
tract_to_msamd_income	122
population	122
minority_population	122
number_of_owner_occupied_units	124
number_of_1_to_4_family_units	122
loan_amount_000s	0
hud_median_family_income	122
applicant_income_000s	6370
sequence_number	0
census_tract_number	122
as_of_year	0
application_date_indicator	0

3.1.1.3.2. Missing Data Treatment: Non-Categorical Variables or Features **3.1.1.3.2.1. Removal of Variables or Features with More Than 50% Missing** Data: None | NCV1, NCV2, ... • rate_spread

3.1.1.3.2.2. Imputation of Missing Data using Descriptive Statistics: Mean | Median • Imputing the missing values using mean

3.1.2. Numerical Encoding of Categorical Variables or Features (Encoding Schema - Alphanumeric Order)

Feature	Number_of_Uneque_Values
state_name	1
state_abbr	1
respondent_id	593
purchaser_type_name	10
property_type_name	3
preapproval_name	3

Feature	Number_of_Uneque_Values
owner_occupancy_name	3
msamd_name	14
loan_type_name	4
loan_purpose_name	3
lien_status_name	4
hoepa_status_name	2
county_name	39
co_applicant_sex_name	5
co_applicant_ethnicity_name	5
applicant_sex_name	4
applicant_ethnicity_name	4
agency_name	6
agency_abbr	6
action_taken_name	8

We are converting the above variables into numeric format in the alpha numeric order

3.1.3. Outlier Statistics and Treatment (Scaling | Transformation) 3.1.3.1.1. Outlier Statistics: Non-Categorical Variables or Features

Outliers count for the Non-Categorical Variables

Variable	Count
tract_to_msamd_income	1309
population	553
minority_population	2641
number_of_owner_occupied_units	478
number_of_1_to_4_family_units	2150
loan_amount_000s	2467
hud_median_family_income	0
applicant_income_000s	3765
sequence_number	7898
census_tract_number	9391
as_of_year	0
application_date_indicator	776

3.1.3.1.2. Outlier Treatment: Non-Categorical Variables or Features 3.1.3.1.2.1.

Standardization: OV1, OV2, ... 3.1.3.1.2.2. Normalization using Min-Max Scaler: OV3, OV4, ...

3.1.3.1.2.3. Log Transformation: OV5, OV6, ... I performed scaling using normalization using min-max scaler. But post the scaling, bubbles were still visible in the box plot. This signifies that

the outliers present in the non categorical dataset are not heavily influenced by the scaling method. The count of outliers seems consistent across different scaling methods.

3.1.4. Data Bifurcation: Training & Testing Sets (Not Required)

3.2. Data Analysis 3.2.1.1. PO1 | PS1:: Unsupervised Machine Learning Clustering Algorithm: K-Means (Base Model) | Metrics Used - Euclidean Distance 3.2.2.1.1. PO2 | PS2:: Clustering Model Performance Evaluation: Silhouette Score | Davies-Bouldin Score (Base Model: K-Mean)

K = 2 Davies-Bouldin Index (DBI): Interpretation: The Davies-Bouldin Index measures the compactness and separation between clusters. Lower DBI values are better: A lower DBI indicates better clustering. It is minimized when clusters are compact (small intra-cluster distances) and well-separated (large inter-cluster distances). Range: The DBI value is non-negative, and lower values are generally better, with 0 being the optimal score. A DBI of 0.462 suggests moderate compactness and separation in your clusters.

Silhouette Score (SS): Interpretation: The Silhouette Score measures how well-defined and separated the clusters are. Range: The SS ranges from -1 to 1. Higher values indicate better-defined clusters. Interpretation of values: Close to +1: Indicates that the object is well matched to its own cluster and poorly matched to neighboring clusters. Close to 0: Indicates overlapping clusters. Close to -1: Indicates that the object is poorly matched to its own cluster and may be a member of a neighboring cluster. An SS of 0.663 suggests well-defined clusters with relatively good separation.

The **cluster profiling** report provides a detailed overview of two clusters (Cluster 0 and Cluster 1) obtained through K-Means clustering with k=2. The clustering is based on various variables, and this report aims to characterize each cluster based on the descriptive statistics of significant variables.

Cluster 0: Overview: Number of Observations: 38,597 Dominant Characteristics: Lower values in key variables such as tract_to_msamd_income, minority_population, number_of_owner_occupied_units, number_of_1_to_4_family_units, loan_amount_000s, applicant_income_000s, sequence_number, census_tract_number, hud_median_family_income. Higher value in application_date_indicator.

Characteristic Variables: tract_to_msamd_income: Lower minority_population: Lower number_of_owner_occupied_units: Lower number_of_1_to_4_family_units: Lower loan_amount_000s: Lower applicant_income_000s: Lower sequence_number: Lower census_tract_number: Lower hud_median_family_income: Lower application_date_indicator: Higher

Cluster 1: Overview: Number of Observations: 21,403 Dominant Characteristics: Higher values in key variables such as tract_to_msamd_income, minority_population, number_of_owner_occupied_units, number_of_1_to_4_family_units, loan_amount_000s, applicant_income_000s, sequence_number, census_tract_number, hud_median_family_income. Lower value in application_date_indicator.

Characteristic Variables: tract_to_msamd_income: Higher minority_population: Higher number_of_owner_occupied_units: Higher number_of_1_to_4_family_units: Higher loan_amount_000s: Higher applicant_income_000s: Higher sequence_number: Higher

census_tract_number: Higher hud_median_family_income: Higher application_date_indicator: Lower

Notes: state_name and state_abbr: These variables have a chi-square p-value of 1.0, suggesting that they might not be strong differentiators between clusters. as_of_year: It seems to have a constant value for both clusters, possibly not providing useful information for clustering.

Conclusion: The clusters exhibit distinct characteristics in terms of socio-economic and loan-related variables. Further domain-specific analysis and business context consideration are recommended to enhance the understanding of the clusters' implications.

Cluster Profiling Report for K=3 Introduction: The **cluster profili**n report provides a detailed overview of three clusters (Cluster 0, Cluster 1, and Cluster 2) obtained through K-Means clustering with k=3. The clustering is based on various variables, and this report aims to characterize each cluster based on the descriptive statistics of significant variables.

Cluster 0: Overview: Number of Observations: 26,674 Dominant Characteristics: Lower values in key variables such as tract_to_msamd_income, population, minority_population, number_of_owner_occupied_units, number_of_1_to_4_family_units, loan_amount_000s, applicant_income_000s, sequence_number, census_tract_number, hud_median_family_income. Higher value in application_date_indicator.

Characteristic Variables: tract_to_msamd_income: Lower population: Lower minority_population: Lower number_of_owner_occupied_units: Lower number_of_1_to_4_family_units: Lower loan_amount_000s: Lower applicant_income_000s: Lower sequence_number: Lower census_tract_number: Lower hud_median_family_income: Lower application_date_indicator: Higher

Cluster 1: Overview: Number of Observations: 15,361 Dominant Characteristics: Higher values in key variables such as tract_to_msamd_income, population, minority_population, number_of_owner_occupied_units, number_of_1_to_4_family_units, loan_amount_000s, applicant_income_000s, sequence_number, census_tract_number, hud_median_family_income. Lower value in application_date_indicator.

Characteristic Variables: tract_to_msamd_income: Higher population: Higher minority_population: Higher number_of_owner_occupied_units: Higher number_of_1_to_4_family_units: Higher loan_amount_000s: Higher applicant_income_000s: Higher sequence_number: Higher census_tract_number: Higher hud_median_family_income: Higher application_date_indicator: Lower

Cluster 2: Overview: Number of Observations: 17,965 Dominant Characteristics: Moderate values in key variables. Comparable characteristics to both Cluster 0 and Cluster 1.

Characteristic Variables: The cluster exhibits characteristics between Cluster 0 and Cluster 1.

Notes: as_of_year: It seems to have a constant value for all clusters, possibly not providing useful information for clustering.

Conclusion: The clusters exhibit distinct characteristics in terms of socio-economic and loan-related variables. Cluster 2 serves as an intermediate group with features falling between Cluster 0 and Cluster 1. Further domain-specific analysis and business context consideration are recommended to enhance the understanding of the clusters' implications.

Cluster Profiling Report for K=4 Introduction: This **cluster profiling** report provides a detailed analysis of the characteristics of each cluster obtained through K-Means clustering with k=4. The clusters are defined based on various variables, and this report aims to highlight the distinctive features of each cluster.

Cluster 0: Overview: Number of Observations: 8,935 Dominant Characteristics: Moderate values in key variables. Comparable characteristics to Clusters 1 and 2. Lower values in hud_median_family_income compared to other clusters.

Characteristic Variables: tract_to_msamd_income: Moderate population: Moderate minority_population: Moderate number_of_owner_occupied_units: Moderate number_of_1_to_4_family_units: Moderate loan_amount_000s: Moderate applicant_income_000s: Moderate sequence_number: Moderate census_tract_number: Moderate hud_median_family_income: Lower application_date_indicator: Moderate

Cluster 1: Overview: Number of Observations: 10,205 Dominant Characteristics: Higher values in key variables compared to Cluster 0. Comparable characteristics to Clusters 2 and 3.

Characteristic Variables: tract_to_msamd_income: Higher population: Higher minority_population: Higher number_of_owner_occupied_units: Higher number_of_1_to_4_family_units: Higher loan_amount_000s: Higher applicant_income_000s: Higher sequence_number: Higher census_tract_number: Higher hud_median_family_income: Moderate application_date_indicator: Moderate

Cluster 2: Overview: Number of Observations: 17,965 Dominant Characteristics: Moderate values in key variables. Comparable characteristics to Clusters 0 and 1.

Characteristic Variables: tract_to_msamd_income: Moderate population: Moderate minority_population: Moderate number_of_owner_occupied_units: Moderate number_of_1_to_4_family_units: Moderate loan_amount_000s: Moderate applicant_income_000s: Moderate sequence_number: Moderate census_tract_number: Moderate hud_median_family_income: Moderate application_date_indicator: Moderate

Cluster 3: Overview: Number of Observations: 22,895 Dominant Characteristics: Higher values in key variables compared to Clusters 0 and 2. Higher values in hud_median_family_income.

Characteristic Variables: tract_to_msamd_income: Higher population: Higher minority_population: Higher number_of_owner_occupied_units: Higher number_of_1_to_4_family_units: Higher loan_amount_000s: Higher applicant_income_000s: Higher sequence_number: Higher census_tract_number: Higher hud_median_family_income: Higher application_date_indicator: Moderate

Notes: as_of_year: It seems to have a constant value for all clusters, possibly not providing useful information for clustering.

Conclusion: The clusters exhibit distinct characteristics in terms of socio-economic and loan-related variables. Each cluster represents a different profile of observations, providing valuable insights into the underlying patterns within the data. Further domain-specific analysis and business context consideration are recommended to enhance the understanding of the clusters' implications.

Cluster Profiling Report for K=5 Introduction: This **cluster profiling** report provides a comprehensive analysis of the characteristics of each cluster obtained through K-Means clustering with k=5. The clusters are defined based on various variables, and this report aims to highlight the distinctive features of each cluster.

Cluster 0: Overview: Number of Observations: 8,935 Dominant Characteristics: Moderate values in key variables. Lower values in hud_median_family_income. Comparable characteristics to Clusters 1 and 2.

Characteristic Variables: tract_to_msamd_income: Moderate population: Moderate minority_population: Moderate number_of_owner_occupied_units: Moderate number_of_1_to_4_family_units: Moderate loan_amount_000s: Moderate applicant_income_000s: Moderate sequence_number: Moderate census_tract_number: Moderate hud_median_family_income: Lower application_date_indicator: Moderate

Cluster 1: Overview: Number of Observations: 10,205 Dominant Characteristics: Higher values in key variables compared to Cluster 0. Comparable characteristics to Clusters 2 and 3.

Characteristic Variables: tract_to_msamd_income: Higher population: Higher minority_population: Higher number_of_owner_occupied_units: Higher number_of_1_to_4_family_units: Higher loan_amount_000s: Higher applicant_income_000s: Higher sequence_number: Higher census_tract_number: Higher hud_median_family_income: Moderate application_date_indicator: Moderate

Cluster 2: Overview: Number of Observations: 17,965 Dominant Characteristics: Moderate values in key variables. Comparable characteristics to Clusters 0 and 1.

Characteristic Variables: tract_to_msamd_income: Moderate population: Moderate minority_population: Moderate number_of_owner_occupied_units: Moderate number_of_1_to_4_family_units: Moderate loan_amount_000s: Moderate applicant_income_000s: Moderate sequence_number: Moderate census_tract_number: Moderate hud_median_family_income: Moderate application_date_indicator: Moderate

Cluster 3: Overview: Number of Observations: 19,457 Dominant Characteristics: Higher values in key variables compared to Clusters 0 and 2. Higher values in hud_median_family_income. Comparable characteristics to Cluster 4. Characteristic Variables: tract_to_msamd_income: Higher population: Higher minority_population: Higher number_of_owner_occupied_units: Higher number_of_1_to_4_family_units: Higher loan_amount_000s: Higher applicant_income_000s: Higher sequence_number: Higher census_tract_number: Higher hud_median_family_income: Higher application_date_indicator: Moderate

Cluster 4: Overview: Number of Observations: 3,438 Dominant Characteristics: Lower values in key variables compared to other clusters. Lower values in hud_median_family_income. Comparable characteristics to Cluster 0.

Characteristic Variables: tract_to_msamd_income: Lower population: Lower minority_population: Lower number_of_owner_occupied_units: Lower number_of_1_to_4_family_units: Lower loan_amount_000s: Lower applicant_income_000s: Lower sequence_number: Lower census_tract_number: Lower hud_median_family_income: Lower application_date_indicator: Lower

Notes: as_of_year: It seems to have a constant value for all clusters, possibly not providing useful information for clustering.

Conclusion: The clusters exhibit distinct characteristics in terms of socio-economic and loan-related variables. Each cluster represents a different profile of observations, providing valuable insights into the underlying patterns within the data. Further domain-specific analysis and business context consideration are recommended to enhance the understanding of the clusters' implications. The addition of Cluster 4 introduces a group with lower values across various features, offering a more diversified segmentation.

3.2.2.1.2. PO2 | PS2:: Clustering Model Performance Evaluation: Time Statistics | (CPU | GPU) Memory Statistics (Base Model: K-Mean)

K	Time	Peak Memory
2	67s	408.77 MiB
3	47.3s	432.41 MiB
4	56.9s	443.09 MiB
5	46.2s	447.80 MiB

3.2.3.1. PO3 | PS3:: Cluster Analysis: Base Model (K-Means) 3.2.3.1.1. Cluster Analysis with Categorical Variables or Features: Chi-Square Test of Independence

Categorical Variable	Chi-Sq Value
state_name	0
state_abbr	0
respondent_id	37667.00325
purchaser_type_name	3421.738007
property_type_name	1313.264545
preapproval_name	494.6907302
owner_occupancy_name	326.8773249
msamd_name	187202.447
loan_type_name	4193.389894
loan_purpose_name	209.9763256
lien_status_name	3332.387181
hoepa_status_name	2.423712105
county_name	239169.2752
co_applicant_sex_name	2320.820989
co_applicant_ethnicity_name	2472.609477
applicant_sex_name	988.6046171
applicant_ethnicity_name	1517.789761
agency_name	4859.037635
agency_abbr	4859.037635
action_taken_name	8254.230514

3.2.3.1.2. Cluster Analysis with Non-Categorical Variables or Features: Analysis of Variance (ANOVA)

Non-Categorical Variable	F-value	P-value
hud_median_family_income	1143184.92	0.0
census_tract_number	7928.80	0.0
number_of_1_to_4_family_units	1808.53	0.0
applicant_income_000s	466.13	0.0
loan_amount_000s	276.42	6.69e-236
minority_population	1543.64	0.0
population	120.34	1.86e-102
tract_to_msamid_income	93.49	2.06e-79
number_of_owner_occupied_units	184.42	2.28e-157
sequence_number	65.08	5.16e-55
application_date_indicator	859.68	0.0

4. Results | Observations 4.1. Appropriate Number of Segments | Clusters: Base Model (K-Means)

K	SS	DBS
2	0.663	0.462
3	0.8003	0.2801
4	0.8316	0.3288
5	0.8465	0.2879

4.2. Cluster Size (Base Model | Comparison Models)

K	C0	C1	C2	C3	C4
2	38597	21403	-	-	-
3	26674	15361	17965	-	-
4	8935	10205	17965	22895	-
5	8935	10205	17965	19457	3438

4.3. Clustering Model Performance: Time & Memory Statistics [Base Model (K-Means)]

K	time	peak memory
2	67s	408.77 MiB
3	47.3s	432.41 MiB
4	56.9s	443.09 MiB
5	46.2s	447.80 MiB

4.4. Cluster Analysis: Base Model (K-Means) 4.4.1. Categorical Variables or Features: Contributing or Significant | Non-Contributing or Non-Significant Significant Variable:**

Categorical Variable	Chi-Sq Value
state_name	0
state_abbr	0
respondent_id	37667.00325
purchaser_type_name	3421.738007
property_type_name	1313.264545
preapproval_name	494.6907302
owner_occupancy_name	326.8773249
msamd_name	187202.447
loan_type_name	4193.389894
loan_purpose_name	209.9763256
lien_status_name	3332.387181
hoepa_status_name	2.423712105
county_name	239169.2752
co_applicant_sex_name	2320.820989
co_applicant_ethnicity_name	2472.609477
applicant_sex_name	988.6046171
applicant_ethnicity_name	1517.789761
agency_name	4859.037635
agency_abbr	4859.037635
action_taken_name	8254.230514

Based on the Chi-Square value county_name is the most significant variable. Followed by msamd_name, respondent_id, action_taken_name and agency_name are the next most significant variables in the same order. Meanwhile, applicant_sex_name, preapproval_name, owner_occupancy_name, loan_purpose_name, hoepa_status_name, state_name, and state_abbr are least significant in the clusters

4.4.2. Non-Categorical Variables or Features: Contributing or Significant | Non-Contributing or Non-Significant

Non-Significant Variable:

Non-Categorical Variable	F-value	P-value
hud_median_family_inc	1143184.92	0.0
ome		
census_tract_number	7928.80	0.0
number_of_1_to_4_fami	1808.53	0.0
ly_units		
applicant_income_000s	466.13	0.0
loan_amount_000s	276.42	6.69e-236

Non-Categorical Variable	F-value	P-value
minority_population	1543.64	0.0
population	120.34	1.86e-102
tract_to_msamd_income	93.49	2.06e-79
number_of_owner_occupied_units	184.42	2.28e-157
sequence_number	65.08	5.16e-55
application_date_indicator	859.68	0.0

As per the f and p value, hud_median_family_income, census_tract_number, number_of_1_to_4_family_unit and minority_population are the most significant non-categorical variable

5. Managerial Insights

5.1. Appropriate Number of Segments | Clusters (Given the Appropriate Model) Silhouette Scores (ss):

The silhouette score measures how well-separated the clusters are. A higher silhouette score indicates better-defined clusters. Here, the silhouette score increases from k=2 to k=5, suggesting an improvement in cluster separation as k increases.

Davies-Bouldin Scores (dbs): The Davies-Bouldin index measures cluster compactness and separation. A lower Davies-Bouldin score indicates better clustering. Here, the Davies-Bouldin score decreases from k=2 to k=5, indicating improved cluster quality. Based on the ss and db scores for k=2,3,4 and 5, both the silhouette scores and Davies-Bouldin scores suggest that a higher value of k might be better. Given the increasing trend in silhouette scores and decreasing trend in Davies-Bouldin scores, we consider **k=5** for potentially better-defined and more compact clusters.

Cluster 0: Action Taken: The mean action taken by applicants is around 3.96, with a low standard deviation (0.28), indicating relatively consistent behavior within the cluster. Agency: The mean agency involvement is approximately 1.29, with a wide range of variation (standard deviation = 1.21). Loan Type: Predominantly type 0 and type 3 loans. Median Family Income: The median family income is around \$63,375, with relatively low variance (standard deviation = \$2,142). Number of 1 to 4 Family Units: The mean value is approximately 0.33, with some variation (standard deviation = 0.12).

Cluster 1: Action Taken: Similar mean action taken as Cluster 0, but with a higher standard deviation (0.85), indicating more variability. Agency: Mean agency involvement is lower compared to Cluster 0, with a similar standard deviation. Loan Type: Predominantly type 0 and type 3 loans, with some variability. Median Family Income: Lower median family income compared to Cluster 0, with slightly higher variance. Number of 1 to 4 Family Units: Higher mean value compared to Cluster 0, with higher variance.

Cluster 2: Action Taken: The mean action taken by applicants is around 3.99, with low standard deviation (0.20), indicating consistent behavior. Agency: Mean agency involvement is higher compared to Clusters 0 and 1, with low variability. Loan Type: Predominantly type 0 loans, with

lower variability compared to other clusters. Median Family Income: Highest median family income among clusters, with no variance. Number of 1 to 4 Family Units: Lower mean value compared to Clusters 0 and 1, with low variance.

Cluster 3: Action Taken: Mean action taken is slightly lower compared to Clusters 0 and 2, with moderate variability. Agency: Similar mean agency involvement as Cluster 1, but with slightly lower variability. Loan Type: Predominantly type 0 and type 3 loans, similar to other clusters. Median Family Income: Lower median family income compared to Clusters 0 and 2, with moderate variance. Number of 1 to 4 Family Units: Similar mean value as Cluster 0, with moderate variability.

Cluster 4: Action Taken: Similar mean action taken as Clusters 0 and 2, with low standard deviation. Agency: Mean agency involvement is lower compared to other clusters, with low variability. Loan Type: Predominantly type 1 loans, with lower variability compared to other clusters. Median Family Income: Similar median family income as Cluster 2, with no variance. Number of 1 to 4 Family Units: Lower mean value compared to other clusters, with low variance.

5.3. Identification of a 'Niche' Segment | Cluster (If Any)

Cluster 4 is a niche segment with 3438 records Cluster 4 emerges as a niche segment within the dataset due to its distinct characteristics and relatively smaller size compared to other clusters. This segment primarily comprises loans predominantly purchased by entities like "Ginnie Mae (GNMA)" and "Freddie Mac (FHLMC)". It is notable for its preference for conventional loans, particularly for one-to-four family dwellings, often utilized for refinancing or home purchases. Moreover, these loans are typically secured by a first lien and are mostly originated without the need for preapproval.

Demographically, the cluster represents properties that are primarily owner-occupied, with a relatively balanced distribution of male and female applicants, most of whom are non-Hispanic or Latino. Despite being a niche segment, the lending activities within this cluster are relatively successful, with a significant proportion resulting in loan origination. From a numerical standpoint, the cluster exhibits moderate economic indicators, such as tract-to-msamd-income ratios and median family incomes, suggesting stability within the localities represented. Overall, this cluster's unique composition and successful lending outcomes make it a distinctive niche within the broader lending landscape.

Managerial Insights: Tract to MSAMD Income and Population: Clusters 0, 1, and 2 have relatively similar average values for these variables, indicating similar socioeconomic conditions in these clusters. Clusters 3 and 4 have higher average values for both variables, suggesting potentially higher-income areas with larger populations.

Minority Population: Cluster 2 has the highest average minority population, indicating a cluster with a higher proportion of minority residents compared to other clusters.

Owner-Occupied Units and 1-4 Family Units: Clusters 3 and 4 have slightly higher average values for both variables, indicating potentially more stable housing conditions with higher rates of ownership.

Loan Amount and Applicant Income: Clusters 2 and 1 have slightly higher average loan amounts and applicant incomes compared to other clusters, indicating higher financial capability in these clusters.

Census Tract Number: Cluster 0 has the highest average census tract number, suggesting potentially more urbanized areas compared to other clusters.

HUD Median Family Income: Cluster 2 has the highest average HUD median family income, indicating potentially higher-income areas compared to other clusters.

Application Date Indicator: Cluster 0 has a significantly higher average application date indicator, suggesting a higher rate of applications within a certain time period compared to other clusters.

Categorical Variables: Each cluster has a dominant category for categorical variables such as purchaser type, property type, preapproval status, etc. These dominant categories can provide insights into the prevalent characteristics or preferences within each cluster.

Clusters with higher average incomes and housing stability (e.g., clusters 2 and 1) may present opportunities for targeting higher-value loan products or services. Understanding the demographics of each cluster (e.g., minority population percentage) can help tailor marketing strategies and product offerings to better serve specific communities. Clusters with higher application rates within a specific time period (e.g., cluster 0) may require additional attention in terms of processing efficiency and customer service. Analyzing dominant categories within each cluster can inform targeted marketing campaigns and product development initiatives tailored to the preferences and needs of each group.

Overall, leveraging these insights can enable financial institutions to better understand their customer base, optimize operations, and tailor their offerings to meet the diverse needs of different communities.

```
# Required Libraries
import pandas as pd, numpy as np # For Data Manipulation
from sklearn.preprocessing import LabelEncoder, OrdinalEncoder # For Encoding Categorical Data [Nominal | Ordinal]
from sklearn.preprocessing import OneHotEncoder # For Creating Dummy Variables of Categorical Data [Nominal]
from sklearn.impute import SimpleImputer, KNNImputer # For Imputation of Missing Data
from sklearn.preprocessing import StandardScaler, MinMaxScaler,
RobustScaler # For Rescaling Data
from sklearn.model_selection import train_test_split # For Splitting Data into Training & Testing Sets
import matplotlib.pyplot as plt
import numpy as np
from scipy.stats import pearsonr
from scipy import stats

# Required Libraries
import pandas as pd, numpy as np # For Data Manipulation
import matplotlib.pyplot as plt, seaborn as sns # For Data Visualization
import scipy.cluster.hierarchy as sch # For Hierarchical Clustering
from sklearn.cluster import AgglomerativeClustering as agclus, KMeans
as kmclus # For Agglomerative & K-Means Clustering
```

```

from sklearn.metrics import silhouette_score as sscore,
davies_bouldin_score as dbscore # For Clustering Model Evaluation

# @title load library { display-mode: "form" }
# Load IPython extension for measuring time
!pip install ipython-autotime
%reload_ext autotime

# Load IPython extension for memory profiling
!pip install memory-profiler
%reload_ext memory_profiler

# Your imports
import pandas as pd
import numpy as np
import matplotlib.pyplot as plt
import seaborn as sns
import scipy.cluster.hierarchy as sch
from sklearn.cluster import AgglomerativeClustering as agclus, KMeans
as kmclus
from sklearn.metrics import silhouette_score as sscore,
davies_bouldin_score as dbscore
from scipy.cluster.hierarchy import dendrogram, linkage
import plotly.graph_objects as go

# Load preprocessing libraries
from sklearn.preprocessing import LabelEncoder, OrdinalEncoder,
OneHotEncoder
from sklearn.impute import SimpleImputer, KNNImputer
from sklearn.preprocessing import StandardScaler, MinMaxScaler,
RobustScaler
from sklearn.model_selection import train_test_split

from scipy.stats import f_oneway

Collecting ipython-autotime
  Downloading ipython_autotime-0.3.2-py2.py3-none-any.whl (7.0 kB)
Requirement already satisfied: ipython in
/usr/local/lib/python3.10/dist-packages (from ipython-autotime)
(7.34.0)
Requirement already satisfied: setuptools>=18.5 in
/usr/local/lib/python3.10/dist-packages (from ipython->ipython-
autotime) (67.7.2)
Collecting jedi>=0.16 (from ipython->ipython-autotime)
  Downloading jedi-0.19.1-py2.py3-none-any.whl (1.6 MB)
1.6/1.6 MB 7.7 MB/s eta
0:00:00
Requirement already satisfied: decorator in /usr/local/lib/python3.10/dist-
packages (from ipython->ipython-autotime) (4.4.2)
Requirement already satisfied: pickleshare in

```

```
/usr/local/lib/python3.10/dist-packages (from ipython->ipython-autotime) (0.7.5)
Requirement already satisfied: traitlets>=4.2 in
/usr/local/lib/python3.10/dist-packages (from ipython->ipython-autotime) (5.7.1)
Requirement already satisfied: prompt-toolkit!=3.0.0,!
=3.0.1,<3.1.0,>=2.0.0 in /usr/local/lib/python3.10/dist-packages (from ipython->ipython-autotime) (3.0.43)
Requirement already satisfied: pygments in
/usr/local/lib/python3.10/dist-packages (from ipython->ipython-autotime) (2.16.1)
Requirement already satisfied: backcall in
/usr/local/lib/python3.10/dist-packages (from ipython->ipython-autotime) (0.2.0)
Requirement already satisfied: matplotlib-inline in
/usr/local/lib/python3.10/dist-packages (from ipython->ipython-autotime) (0.1.6)
Requirement already satisfied: pexpect>4.3 in
/usr/local/lib/python3.10/dist-packages (from ipython->ipython-autotime) (4.9.0)
Requirement already satisfied: parso<0.9.0,>=0.8.3 in
/usr/local/lib/python3.10/dist-packages (from jedi>=0.16->ipython->ipython-autotime) (0.8.3)
Requirement already satisfied: ptyprocess>=0.5 in
/usr/local/lib/python3.10/dist-packages (from pexpect>4.3->ipython->ipython-autotime) (0.7.0)
Requirement already satisfied: wcwidth in
/usr/local/lib/python3.10/dist-packages (from prompt-toolkit!=3.0.0,!
=3.0.1,<3.1.0,>=2.0.0->ipython->ipython-autotime) (0.2.13)
Installing collected packages: jedi, ipython-autotime
Successfully installed ipython-autotime-0.3.2 jedi-0.19.1
Collecting memory-profiler
  Downloading memory_profiler-0.61.0-py3-none-any.whl (31 kB)
Requirement already satisfied: psutil in
/usr/local/lib/python3.10/dist-packages (from memory-profiler) (5.9.5)
Installing collected packages: memory-profiler
Successfully installed memory-profiler-0.61.0
time: 9.78 s (started: 2024-03-19 18:26:21 +00:00)

import pandas as pd
import gdown

# Google Drive file ID
file_id = '1uCd0yJJa8yLNPvi4un6YXxyzPe9An_Fv'

# Downloading the CSV file from Google Drive
url = f'https://drive.google.com/uc?id={file_id}'
csv_file_path = 'Main_Washington_State_HDMA_2016.csv'
gdown.download(url, csv_file_path, quiet=False)
```

```

# Read the CSV file into a pandas DataFrame
df = pd.read_csv(csv_file_path)

# Display the first few rows of the DataFrame to verify the data
print(df.head())

Downloading...
From: https://drive.google.com/uc?id=1uCd0yJJa8yLNPvi4un6YXxyzPe9An_Fv
To: /content/Main_Washington_State_HDMA_2016.csv
100%|██████████| 31.6M/31.6M [00:00<00:00, 76.4MB/s]
<ipython-input-3-e36da72ebe94>:13: DtypeWarning: Columns (23,24,25)
have mixed types. Specify dtype option on import or set
low_memory=False.

df = pd.read_csv(csv_file_path)

      tract_to_msamd_income  rate_spread  population  minority_population
\0              121.690002        NaN       8381.0          23.790001
1              83.370003        NaN       4915.0          23.990000
2              91.129997        NaN       5075.0          11.820000
3              146.169998        NaN       5032.0          8.590000
4              162.470001        NaN       5183.0          10.500000

      number_of_owner_occupied_units  number_of_1_to_4_family_units \
0                  2175.0                      2660.0
1                  1268.0                      1777.0
2                  1136.0                      1838.0
3                  1525.0                      1820.0
4                  1705.0                      2104.0

      loan_amount_000s  hud_median_family_income
applicant_income_000s \
0                  227                      73300.0          116.0
1                  240                      57900.0          42.0
2                  241                      73300.0          117.0
3                  351                      73300.0          315.0
4                  417                      78100.0          114.0

      state_name  ... co_applicant_sex_name \
0 Washington  ...             Male
1 Washington  ...        No co-applicant

```

```

2 Washington ... Female
3 Washington ... Female
4 Washington ... Male

co_applicant_ethnicity_name
census_tract_number \
0 Not Hispanic or Latino
413.27
1 No co-applicant
9208.01
2 Not Hispanic or Latino
414.00
3 Information not provided by applicant in mail, ...
405.10
4 Not Hispanic or Latino
907.00

as_of_year application_date_indicator applicant_sex_name \
0 2016 0 Female
1 2016 0 Male
2 2016 0 Male
3 2016 0 Male
4 2016 0 Female

applicant_ethnicity_name \
0 Not Hispanic or Latino
1 Hispanic or Latino
2 Not Hispanic or Latino
3 Information not provided by applicant in mail, ...
4 Not Hispanic or Latino

agency_name agency_abbr
action_taken_name
0 Consumer Financial Protection Bureau CFPB Loan
originated
1 Department of Housing and Urban Development HUD Loan
originated
2 Department of Housing and Urban Development HUD Loan
originated
3 National Credit Union Administration NCUA Loan
originated
4 Federal Deposit Insurance Corporation FDIC Loan
originated

[5 rows x 37 columns]
time: 3.85 s (started: 2024-03-19 18:26:58 +00:00)

# Add a new column named "S.no" as the first column with serial numbers
df.insert(0, "S.no", range(1, len(df) + 1))

```



```

census_tract_number as_of_year application_date_indicator \
0           413.27      2016          0
1           9208.01      2016          0
2           414.00      2016          0
3           405.10      2016          0
4           907.00      2016          0

applicant_sex_name
applicant_ethnicity_name \
0                  Female           Not Hispanic or
Latino
1                  Male            Hispanic or
Latino
2                  Male           Not Hispanic or
Latino
3                  Male  Information not provided by applicant in
mail,...
4                  Female          Not Hispanic or
Latino

agency_name agency_abbr
action_taken_name
0           Consumer Financial Protection Bureau      CFPB   Loan
originated
1           Department of Housing and Urban Development      HUD   Loan
originated
2           Department of Housing and Urban Development      HUD   Loan
originated
3           National Credit Union Administration      NCUA   Loan
originated
4           Federal Deposit Insurance Corporation      FDIC   Loan
originated

[5 rows x 38 columns]
time: 33.2 ms (started: 2024-03-19 18:27:01 +00:00)

df.info()
list(df.columns)

# Assuming df is your original DataFrame
# Add your normalization or standardization code here

# Display summary statistics
df.describe()

total_records = len(df)
print(f"Total number of records: {total_records}")

# Calculate the total number of filled cells in each column

```

```

filled_cells_count = df.count()

# Sum up the counts to get the total number of filled cells in the DataFrame
total_filled_cells = filled_cells_count.sum()

print(f"Total number of filled cells: {total_filled_cells}")

# Assuming df is your DataFrame
unique_counts = df.nunique()

# Display the number of unique values in each column
print(unique_counts)

<class 'pandas.core.frame.DataFrame'>
RangeIndex: 60000 entries, 0 to 59999
Data columns (total 38 columns):
 #   Column           Non-Null Count  Dtype  
 ---  -- 
 0   S.no             60000 non-null   int64  
 1   tract_to_msamd_income 59878 non-null   float64 
 2   rate_spread       1866 non-null   float64 
 3   population        59878 non-null   float64 
 4   minority_population 59878 non-null   float64 
 5   number_of_owner_occupied_units 59876 non-null   float64 
 6   number_of_1_to_4_family_units 59878 non-null   float64 
 7   loan_amount_000s    60000 non-null   int64  
 8   hud_median_family_income 59878 non-null   float64 
 9   applicant_income_000s   53630 non-null   float64 
 10  state_name        60000 non-null   object  
 11  state_abbr        60000 non-null   object  
 12  sequence_number   60000 non-null   int64  
 13  respondent_id     60000 non-null   object  
 14  purchaser_type_name 60000 non-null   object  
 15  property_type_name 60000 non-null   object  
 16  preapproval_name   60000 non-null   object  
 17  owner_occupancy_name 60000 non-null   object  
 18  msamd_name        51619 non-null   object  
 19  loan_type_name     60000 non-null   object  
 20  loan_purpose_name  60000 non-null   object  
 21  lien_status_name   60000 non-null   object  
 22  hoepa_status_name  60000 non-null   object  
 23  edit_status_name   12451 non-null   object  
 24  denial_reason_name_3 1 non-null    object  
 25  denial_reason_name_2 43 non-null   object  
 26  denial_reason_name_1 118 non-null  object  
 27  county_name        59908 non-null  object  
 28  co_applicant_sex_name 60000 non-null  object  
 29  co_applicant_ethnicity_name 60000 non-null  object  
 30  census_tract_number 59878 non-null  float64

```

```

31 as_of_year           60000 non-null int64
32 application_date_indicator 60000 non-null int64
33 applicant_sex_name    60000 non-null object
34 applicant_ethnicity_name 60000 non-null object
35 agency_name           60000 non-null object
36 agency_abbr           60000 non-null object
37 action_taken_name     60000 non-null object
dtypes: float64(9), int64(5), object(24)
memory usage: 17.4+ MB
Total number of records: 60000
Total number of filled cells: 1978780
S.no          60000
tract_to_msamd_income 1327
rate_spread      320
population       1283
minority_population 1216
number_of_owner_occupied_units 996
number_of_1_to_4_family_units 1051
loan_amount_000s   1344
hud_median_family_income 15
applicant_income_000s 807
state_name        1
state_abbr         1
sequence_number    39340
respondent_id      593
purchaser_type_name 10
property_type_name 3
preapproval_name    3
owner_occupancy_name 3
msamd_name          14
loan_type_name      4
loan_purpose_name    3
lien_status_name     4
hoepa_status_name    2
edit_status_name     1
denial_reason_name_3 1
denial_reason_name_2 8
denial_reason_name_1 8
county_name          39
co_applicant_sex_name 5
co_applicant_ethnicity_name 5
census_tract_number 1108
as_of_year           1
application_date_indicator 2
applicant_sex_name    4
applicant_ethnicity_name 4
agency_name           6
agency_abbr           6
action_taken_name     8

```

```
dtype: int64
time: 2.06 s (started: 2024-03-19 18:27:02 +00:00)

# Assuming df is your DataFrame
columns_list = df.columns.tolist()
columns_list

list(df.columns)

['S.no',
 'tract_to_msamd_income',
 'rate_spread',
 'population',
 'minority_population',
 'number_of_owner_occupied_units',
 'number_of_1_to_4_family_units',
 'loan_amount_000s',
 'hud_median_family_income',
 'applicant_income_000s',
 'state_name',
 'state_abbr',
 'sequence_number',
 'respondent_id',
 'purchaser_type_name',
 'property_type_name',
 'preapproval_name',
 'owner_occupancy_name',
 'msamd_name',
 'loan_type_name',
 'loan_purpose_name',
 'lien_status_name',
 'hoepa_status_name',
 'edit_status_name',
 'denial_reason_name_3',
 'denial_reason_name_2',
 'denial_reason_name_1',
 'county_name',
 'co_applicant_sex_name',
 'co_applicant_ethnicity_name',
 'census_tract_number',
 'as_of_year',
 'application_date_indicator',
 'applicant_sex_name',
 'applicant_ethnicity_name',
 'agency_name',
 'agency_abbr',
 'action_taken_name']

time: 6.55 ms (started: 2024-03-19 18:27:04 +00:00)
```

```

# Nominal and Ordinal Columns

# Continuous and Non Continuous Columns

import pandas as pd

# Assuming df is your DataFrame
continuous_columns = df.select_dtypes(include=['float64',
'int64']).columns
non_continuous_columns = df.select_dtypes(exclude=['float64',
'int64']).columns

print("Continuous Columns:", list(continuous_columns))
print("Non-Continuous Columns:", list(non_continuous_columns))

# Assuming df is your DataFrame
categorical_columns = df.select_dtypes(include=['object',
'category']).columns
non_categorical_columns = df.select_dtypes(exclude=['object',
'category']).columns

print("Categorical Columns:", list(categorical_columns))
print("Non-Categorical Columns:", list(non_categorical_columns))

...
# Assuming df is your DataFrame

# Define a dictionary to store information about the nature of each variable
variable_info = {}

# Loop through columns in the DataFrame
for column in df.columns:
    # Check if the variable has a limited number of unique values
    unique_values = df[column].nunique()

    if unique_values <= 10:
        # If the number of unique values is small, it could be an ordinal variable
        variable_info[column] = 'Ordinal'
    else:
        # If the number of unique values is larger, it could be a nominal variable
        variable_info[column] = 'Nominal'

# Print the information about each variable
for column, nature in variable_info.items():

```

```

    print(f"{column}: {nature}")
...
Continuous Columns: ['S.no', 'tract_to_msamd_income', 'rate_spread',
'population', 'minority_population', 'number_of_owner_occupied_units',
'number_of_1_to_4_family_units', 'loan_amount_000s',
'hud_median_family_income', 'applicant_income_000s',
'sequence_number', 'census_tract_number', 'as_of_year',
'application_date_indicator']
Non-Continuous Columns: ['state_name', 'state_abbr', 'respondent_id',
'purchaser_type_name', 'property_type_name', 'preapproval_name',
'owner_occupancy_name', 'msamd_name', 'loan_type_name',
'loan_purpose_name', 'lien_status_name', 'hoepa_status_name',
'edit_status_name', 'denial_reason_name_3', 'denial_reason_name_2',
'denial_reason_name_1', 'county_name', 'co_applicant_sex_name',
'co_applicant_ethnicity_name', 'applicant_sex_name',
'applicant_ethnicity_name', 'agency_name', 'agency_abbr',
'action_taken_name']
Categorical Columns: ['state_name', 'state_abbr', 'respondent_id',
'purchaser_type_name', 'property_type_name', 'preapproval_name',
'owner_occupancy_name', 'msamd_name', 'loan_type_name',
'loan_purpose_name', 'lien_status_name', 'hoepa_status_name',
'edit_status_name', 'denial_reason_name_3', 'denial_reason_name_2',
'denial_reason_name_1', 'county_name', 'co_applicant_sex_name',
'co_applicant_ethnicity_name', 'applicant_sex_name',
'applicant_ethnicity_name', 'agency_name', 'agency_abbr',
'action_taken_name']
Non-Categorical Columns: ['S.no', 'tract_to_msamd_income',
'rate_spread', 'population', 'minority_population',
'number_of_owner_occupied_units', 'number_of_1_to_4_family_units',
'loan_amount_000s', 'hud_median_family_income',
'applicant_income_000s', 'sequence_number', 'census_tract_number',
'as_of_year', 'application_date_indicator']

{"type": "string"}

time: 114 ms (started: 2024-03-19 18:27:04 +00:00)

### Missing Data Statistics and Treatment
### Missing Data Statistics: Records

# Assuming df is your DataFrame

# Count the missing values in each column
missing_data = df.isnull().sum()

# Create a DataFrame to display missing data statistics
missing_data_stats = pd.DataFrame({
    'Column': missing_data.index,
    'Missing Records': missing_data.values,
})

```

```

    'Percentage Missing': (missing_data / len(df)) * 100
})

# Sort the DataFrame by the percentage of missing values in descending
# order
missing_data_stats = missing_data_stats.sort_values(by='Percentage
Missing', ascending=False)

# Print the missing data statistics
print(missing_data_stats)

```

	Column \
denial_reason_name_3	denial_reason_name_3
denial_reason_name_2	denial_reason_name_2
denial_reason_name_1	denial_reason_name_1
rate_spread	rate_spread
edit_status_name	edit_status_name
msamd_name	msamd_name
applicant_income_000s	applicant_income_000s
number_of_owner_occupied_units	number_of_owner_occupied_units
population	population
minority_population	minority_population
number_of_1_to_4_family_units	number_of_1_to_4_family_units
hud_median_family_income	hud_median_family_income
tract_to_msamd_income	tract_to_msamd_income
census_tract_number	census_tract_number
county_name	county_name
co_applicant_sex_name	co_applicant_sex_name
S.no	S.no
co_applicant_ethnicity_name	co_applicant_ethnicity_name
as_of_year	as_of_year
application_date_indicator	application_date_indicator
applicant_ethnicity_name	applicant_ethnicity_name
agency_name	agency_name
agency_abbr	agency_abbr
applicant_sex_name	applicant_sex_name
loan_type_name	loan_type_name
hoepa_status_name	hoepa_status_name
lien_status_name	lien_status_name
loan_purpose_name	loan_purpose_name
owner_occupancy_name	owner_occupancy_name
preapproval_name	preapproval_name
property_type_name	property_type_name
purchaser_type_name	purchaser_type_name
respondent_id	respondent_id
sequence_number	sequence_number
state_abbr	state_abbr
state_name	state_name
loan_amount_000s	loan_amount_000s
action_taken_name	action_taken_name

	Missing Records	Percentage Missing
denial_reason_name_3	59999	99.998333
denial_reason_name_2	59957	99.928333
denial_reason_name_1	59882	99.803333
rate_spread	58134	96.890000
edit_status_name	47549	79.248333
msamd_name	8381	13.968333
applicant_income_000s	6370	10.616667
number_of_owner_occupied_units	124	0.206667
population	122	0.203333
minority_population	122	0.203333
number_of_1_to_4_family_units	122	0.203333
hud_median_family_income	122	0.203333
tract_to_msamd_income	122	0.203333
census_tract_number	122	0.203333
county_name	92	0.153333
co_applicant_sex_name	0	0.000000
S.no	0	0.000000
co_applicant_ethnicity_name	0	0.000000
as_of_year	0	0.000000
application_date_indicator	0	0.000000
applicant_ethnicity_name	0	0.000000
agency_name	0	0.000000
agency_abbr	0	0.000000
applicant_sex_name	0	0.000000
loan_type_name	0	0.000000
hoepa_status_name	0	0.000000
lien_status_name	0	0.000000
loan_purpose_name	0	0.000000
owner_occupancy_name	0	0.000000
preapproval_name	0	0.000000
property_type_name	0	0.000000
purchaser_type_name	0	0.000000
respondent_id	0	0.000000
sequence_number	0	0.000000
state_abbr	0	0.000000
state_name	0	0.000000
loan_amount_000s	0	0.000000
action_taken_name	0	0.000000

time: 1.13 s (started: 2024-03-19 18:27:07 +00:00)

```
# List of columns to drop
columns_to_drop = ['denial_reason_name_3', 'denial_reason_name_2',
'denial_reason_name_1', 'rate_spread', 'edit_status_name']

# Drop columns with more than 50% missing values
df_cleaned = df.drop(columns=columns_to_drop)
```

```

# Print the cleaned DataFrame
df1 = df_cleaned

# Count the missing values in each column
missing_data = df1.isnull().sum()

# Create a DataFrame to display missing data statistics
missing_data_stats = pd.DataFrame({
    'Column': missing_data.index,
    'Missing Records': missing_data.values,
    'Percentage Missing': (missing_data / len(df)) * 100
})
# Sort the DataFrame by the percentage of missing values in descending order
missing_data_stats = missing_data_stats.sort_values(by='Percentage Missing', ascending=False)

# Print the missing data statistics
print(missing_data_stats)

```

	Column \
msamd_name	msamd_name
applicant_income_000s	applicant_income_000s
number_of_owner_occupied_units	number_of_owner_occupied_units
population	population
minority_population	minority_population
number_of_1_to_4_family_units	number_of_1_to_4_family_units
hud_median_family_income	hud_median_family_income
tract_to_msamd_income	tract_to_msamd_income
census_tract_number	census_tract_number
county_name	county_name
as_of_year	as_of_year
co_applicant_sex_name	co_applicant_sex_name
co_applicant_ethnicity_name	co_applicant_ethnicity_name
S.no	S.no
application_date_indicator	application_date_indicator
applicant_sex_name	applicant_sex_name
hoepa_status_name	hoepa_status_name
agency_name	agency_name
agency_abbr	agency_abbr
applicant_ethnicity_name	applicant_ethnicity_name
owner_occupancy_name	owner_occupancy_name
lien_status_name	lien_status_name
loan_purpose_name	loan_purpose_name
loan_type_name	loan_type_name
preapproval_name	preapproval_name
property_type_name	property_type_name
purchaser_type_name	purchaser_type_name
respondent_id	respondent_id
sequence_number	sequence_number

	state_abbr	state_name	loan_amount_000s	action_taken_name	state_abbr	state_name	loan_amount_000s	action_taken_name
					Missing	Records	Percentage	Missing
msamd_name					8381		13.968333	
applicant_income_000s					6370		10.616667	
number_of_owner_occupied_units					124		0.206667	
population					122		0.203333	
minority_population					122		0.203333	
number_of_1_to_4_family_units					122		0.203333	
hud_median_family_income					122		0.203333	
tract_to_msamd_income					122		0.203333	
census_tract_number					122		0.203333	
county_name					92		0.153333	
as_of_year					0		0.000000	
co_applicant_sex_name					0		0.000000	
co_applicant_ethnicity_name					0		0.000000	
S.no					0		0.000000	
application_date_indicator					0		0.000000	
applicant_sex_name					0		0.000000	
hoepa_status_name					0		0.000000	
agency_name					0		0.000000	
agency_abbr					0		0.000000	
applicant_ethnicity_name					0		0.000000	
owner_occupancy_name					0		0.000000	
lien_status_name					0		0.000000	
loan_purpose_name					0		0.000000	
loan_type_name					0		0.000000	
preapproval_name					0		0.000000	
property_type_name					0		0.000000	
purchaser_type_name					0		0.000000	
respondent_id					0		0.000000	
sequence_number					0		0.000000	
state_abbr					0		0.000000	
state_name					0		0.000000	
loan_amount_000s					0		0.000000	
action_taken_name					0		0.000000	

time: 758 ms (started: 2024-03-19 18:27:10 +00:00)

Missing Records (ROWS)

```
# Count the missing values in each row
missing_rows = df1.isnull().sum(axis=1)
```

```
# Count the number of rows with at least one missing value
num_rows_with_missing = len(missing_rows[missing_rows > 0])
```

Print the number of rows with missing data

```

print("Number of rows with missing data:", num_rows_with_missing)

# Calculate the percentage of missing values in each row
missing_percentage_rows = (df1.isnull().sum(axis=1) /
len(df1.columns)) * 100

# Count the number of rows with more than 50% missing data
num_rows_more_than_50_percent_missing =
len(missing_percentage_rows[missing_percentage_rows > 50])

# Print the number of rows with more than 50% missing data
print("Number of rows with more than 50% missing data:",
num_rows_more_than_50_percent_missing)

Number of rows with missing data: 13629
Number of rows with more than 50% missing data: 0
time: 911 ms (started: 2024-03-19 18:27:13 +00:00)

# DIVIDING DF1 into Cat and Non Cat

# Assuming df1 is your DataFrame
cat_columns = df1.select_dtypes(include=['object']).columns
noncat_columns = df1.select_dtypes(exclude=['object']).columns

# Creating categorical and non-categorical DataFrames
catdf1 = df1[cat_columns]
noncatdf1 = df1[noncat_columns]

#print(list(catdf.columns))
#print(list(noncatdf.columns))
print(list(catdf1.columns))
print(list(noncatdf1.columns))

#20
#list(noncatdf.columns)

['state_name', 'state_abbr', 'respondent_id', 'purchaser_type_name',
'property_type_name', 'preapproval_name', 'owner_occupancy_name',
'msamd_name', 'loan_type_name', 'loan_purpose_name',
'lien_status_name', 'hoepa_status_name', 'county_name',
'co_applicant_sex_name', 'co_applicant_ethnicity_name',
'applicant_sex_name', 'applicant_ethnicity_name', 'agency_name',
'agency_abbr', 'action_taken_name']
['S.no', 'tract_to_msamd_income', 'population', 'minority_population',
'number_of_owner_occupied_units', 'number_of_1_to_4_family_units',
'loan_amount_000s', 'hud_median_family_income',
'applicant_income_000s', 'sequence_number', 'census_tract_number',
'as_of_year', 'application_date_indicator']
time: 35.5 ms (started: 2024-03-19 18:27:27 +00:00)

```

```

# Data Bifurcation
catdf = df[['S.no', 'state_name', 'state_abbr', 'respondent_id',
'purchaser_type_name', 'property_type_name', 'preapproval_name',
'owner_occupancy_name', 'msamd_name', 'loan_type_name',
'loan_purpose_name',
        'lien_status_name', 'hoepa_status_name', 'county_name',
'co_applicant_sex_name', 'co_applicant_ethnicity_name',
'applicant_sex_name', 'applicant_ethnicity_name', 'agency_name',
        'agency_abbr', 'action_taken_name']] # Categorical Data
[Nominal | Ordinal]

noncatdf = df[['S.no', 'tract_to_msamd_income', 'population',
'minority_population', 'number_of_owner_occupied_units',
'number_of_1_to_4_family_units', 'loan_amount_000s',
        'hud_median_family_income', 'applicant_income_000s',
'sequence_number', 'census_tract_number', 'as_of_year',
'application_date_indicator']] # Non-Categorical Data

time: 25.2 ms (started: 2024-03-19 18:27:30 +00:00)

time: 2.34 s (started: 2024-03-19 18:27:30 +00:00)

#### STATISTICS OF CAT DATASET

# Count and frequency statistics for each column in catdf
catdf_stats = pd.DataFrame()

for column in catdf.columns:
    col_count = catdf[column].value_counts().reset_index()
    col_count.columns = [column, 'Frequency']
    catdf_stats = pd.concat([catdf_stats, col_count], axis=1)

# Display the count and frequency statistics
#print(catdf_stats)

# Summary for each column in catdf
catdf_summary = catdf.describe(include='all').transpose()

# Display the summary
print(catdf_summary)

# Calculate the proportion (relative frequency) for each categorical
# column
#proportion_stats = catdf.apply(lambda x:
#x.value_counts(normalize=True).idxmax() + ': ' +
#"{: .2%}".format(x.value_counts(normalize=True).max())))

```

```

# Display the proportion statistics
#print(proportion_stats)

          count unique \
S.no
state_name      60000.0    NaN
state_abbr       60000     1
respondent_id    60000     1
purchaser_type_name  60000    10
property_type_name  60000     3
preapproval_name   60000     3
owner_occupancy_name  60000     3
msamd_name        51619    14
loan_type_name    60000     4
loan_purpose_name  60000     3
lien_status_name   60000     4
hoepa_status_name  60000     2
county_name        59908    39
co_applicant_sex_name  60000     5
co_applicant_ethnicity_name  60000     5
applicant_sex_name  60000     4
applicant_ethnicity_name  60000     4
agency_name         60000     6
agency_abbr         60000     6
action_taken_name   60000     8

top \
S.no
NaN
state_name
Washington
state_abbr
WA
respondent_id
32489
purchaser_type_name      Loan was not originated or was not sold
in cal...
property_type_name        One-to-four family dwelling (other than
manufa...
preapproval_name           Not
applicable
owner_occupancy_name       Owner-occupied as a principal
dwelling
msamd_name                 Seattle, Bellevue,
Everett - WA
loan_type_name
Conventional
loan_purpose_name
Refinancing

```

lien_status_name					Secured by a
first_lien					
hoepa_status_name					Not a
HOEPA_loan					
county_name					
King_County					
co_applicant_sex_name					No co-
applicant					
co_applicant_ethnicity_name					No co-
applicant					
applicant_sex_name					
Male					
applicant_ethnicity_name					Not Hispanic
or_Latino					
agency_name					Department of Housing and Urban
Development					
agency_abbr					
HUD					
action_taken_name					Loan
originated					
		freq	mean	std	min
25% \					
S.no		NaN	30000.5	17320.652413	1.0
15000.75					
state_name	60000		NaN		NaN
NaN					
state_abbr	60000		NaN		NaN
NaN					
respondent_id	5006		NaN		NaN
NaN					
purchaser_type_name	16112		NaN		NaN
NaN					
property_type_name	57630		NaN		NaN
NaN					
preapproval_name	47832		NaN		NaN
NaN					
owner_occupancy_name	53940		NaN		NaN
NaN					
msamd_name	17965		NaN		NaN
NaN					
loan_type_name	42917		NaN		NaN
NaN					
loan_purpose_name	28576		NaN		NaN
NaN					
lien_status_name	57046		NaN		NaN
NaN					
hoepa_status_name	59996		NaN		NaN
NaN					

county_name	12915	NaN	NaN	NaN
NaN				
co_applicant_sex_name	26987	NaN	NaN	NaN
NaN				
co_applicant_ethnicity_name	26987	NaN	NaN	NaN
NaN				
applicant_sex_name	37070	NaN	NaN	NaN
NaN				
applicant_ethnicity_name	44014	NaN	NaN	NaN
NaN				
agency_name	27514	NaN	NaN	NaN
NaN				
agency_abbr	27514	NaN	NaN	NaN
NaN				
action_taken_name	55815	NaN	NaN	NaN
NaN				
S.no	50%	75%	max	
	30000.5	45000.25	60000.0	
state_name	NaN	NaN	NaN	
state_abbr	NaN	NaN	NaN	
respondent_id	NaN	NaN	NaN	
purchaser_type_name	NaN	NaN	NaN	
property_type_name	NaN	NaN	NaN	
preapproval_name	NaN	NaN	NaN	
owner_occupancy_name	NaN	NaN	NaN	
msamd_name	NaN	NaN	NaN	
loan_type_name	NaN	NaN	NaN	
loan_purpose_name	NaN	NaN	NaN	
lien_status_name	NaN	NaN	NaN	
hoepa_status_name	NaN	NaN	NaN	
county_name	NaN	NaN	NaN	
co_applicant_sex_name	NaN	NaN	NaN	
co_applicant_ethnicity_name	NaN	NaN	NaN	
applicant_sex_name	NaN	NaN	NaN	
applicant_ethnicity_name	NaN	NaN	NaN	
agency_name	NaN	NaN	NaN	
agency_abbr	NaN	NaN	NaN	
action_taken_name	NaN	NaN	NaN	

time: 1.06 s (started: 2024-03-19 18:27:35 +00:00)

STATISTICS OF NONCAT DATASET

```
# Display descriptive statistics for non-categorical variables
noncatdf_descriptive_stats = noncatdf.describe()

# Print the descriptive statistics
print(noncatdf_descriptive_stats)
```

	S.no	tract_to_msam_d_income	population
minority_population	\		
count	60000.000000	59878.000000	59878.000000
59878.000000			
mean	30000.500000	107.617351	5278.782157
23.244442			
std	17320.652413	28.233471	1716.101490
14.416209			
min	1.000000	14.050000	98.000000
2.040000			
25%	15000.750000	88.970001	4070.000000
12.950000			
50%	30000.500000	105.550003	5145.000000
19.420000			
75%	45000.250000	123.330002	6382.000000
29.680000			
max	60000.000000	257.140015	13025.000000
94.790001			
number_of_owner_occupied_units			
number_of_1_to_4_family_units	\		
count	59876.000000		59878.000000
mean	1399.044375		1873.281456
std	518.330561		738.505184
min	15.000000		27.000000
25%	1034.000000		1414.000000
50%	1359.000000		1770.000000
75%	1722.000000		2249.000000
max	2997.000000		5893.000000
loan_amount_000s hud_median_family_income			
applicant_income_000s	\		
count	60000.000000	59878.000000	
53630.000000			
mean	291.358717	73869.411136	
112.822301			
std	604.958183	12811.243390	
122.862496			
min	1.000000	48700.000000	
1.000000			
25%	170.000000	63100.000000	
61.000000			
50%	242.000000	73300.000000	

```

89.000000
75%      337.000000
132.000000
max      55000.000000
6161.000000

      sequence_number  census_tract_number  as_of_year \
count    6.000000e+04      59878.000000   60000.0
mean     7.752647e+04      1750.597252   2016.0
std      1.505157e+05      3359.676740      0.0
min     1.000000e+00      1.000000   2016.0
25%     3.231500e+03      114.020000   2016.0
50%     1.648100e+04      403.020000   2016.0
75%     7.276225e+04      713.100000   2016.0
max     1.241590e+06      9757.000000  2016.0

      application_date_indicator
count      60000.000000
mean       0.025867
std        0.225976
min        0.000000
25%        0.000000
50%        0.000000
75%        0.000000
max        2.000000
time: 122 ms (started: 2024-03-19 18:27:38 +00:00)

...
correlation_matrix = noncatdf.corr()

# Initialize empty matrices for p-values
p_values = np.zeros_like(correlation_matrix)
n = len(noncatdf)

# Calculate correlation coefficients and p-values
for i in range(len(noncatdf.columns)):
    for j in range(i, len(noncatdf.columns)):
        corr, p_value = pearsonr(noncatdf.iloc[:, i], noncatdf.iloc[:, j])
        correlation_matrix.iloc[i, j] = corr
        correlation_matrix.iloc[j, i] = corr
        p_values[i, j] = p_value
        p_values[j, i] = p_value

# Create DataFrames for correlation coefficients and p-values
correlation_df = pd.DataFrame(correlation_matrix,
columns=noncatdf.columns, index=noncatdf.columns)
p_values_df = pd.DataFrame(p_values, columns=noncatdf.columns,
index=noncatdf.columns)
...

```

```
{"type": "string"}
```

time: 11.5 ms (started: 2024-03-19 18:27:41 +00:00)

```
# Missing Data Statistics: Non-Categorical Variables or Features
```

```
# Calculate missing data statistics for non-categorical columns
```

```
missing_data_non_categorical = noncatdf.isnull().sum().reset_index()
```

```
missing_data_non_categorical.columns = ['Feature', 'Missing_Records']
```

```
# Display the missing data statistics
```

```
print(missing_data_non_categorical)
```

	Feature	Missing_Records
0	S.no	0
1	tract_to_msamd_income	122
2	population	122
3	minority_population	122
4	number_of_owner_occupied_units	124
5	number_of_1_to_4_family_units	122
6	loan_amount_000s	0
7	hud_median_family_income	122
8	applicant_income_000s	6370
9	sequence_number	0
10	census_tract_number	122
11	as_of_year	0
12	application_date_indicator	0

time: 9.94 ms (started: 2024-03-19 18:28:08 +00:00)

```
# Missing Data Treatment: Non-Categorical Variables or Features
```

```
# Dataset Used : df_noncat
```

```
si_noncat = SimpleImputer(missing_values=np.nan, strategy='mean') # Other Strategy : mean | median | most_frequent | constant
```

```
si_noncat_fit = si_noncat.fit_transform(noncatdf)
```

```
imputed_data_non_categorical = pd.DataFrame(si_noncat_fit,
```

```
columns=noncatdf.columns); # Missing Non-Categorical Data Imputed
```

Subset using Simple Imputer

```
imputed_data_non_categorical.info()
```

<class 'pandas.core.frame.DataFrame'>

RangeIndex: 60000 entries, 0 to 59999

Data columns (total 13 columns):

#	Column	Non-Null Count	Dtype	
0	S.no	60000	non-null	float64
1	tract_to_msamd_income	60000	non-null	float64
2	population	60000	non-null	float64
3	minority_population	60000	non-null	float64

```

4   number_of_owner_occupied_units    60000 non-null float64
5   number_of_1_to_4_family_units    60000 non-null float64
6   loan_amount_000s                60000 non-null float64
7   hud_median_family_income       60000 non-null float64
8   applicant_income_000s          60000 non-null float64
9   sequence_number                 60000 non-null float64
10  census_tract_number            60000 non-null float64
11  as_of_year                     60000 non-null float64
12  application_date_indicator    60000 non-null float64
dtypes: float64(13)
memory usage: 6.0 MB
time: 28.2 ms (started: 2024-03-19 18:28:12 +00:00)

# Calculate standard deviation for non-categorical columns
std_deviation_non_categorical = imputed_data_non_categorical.std()

# Creating a DataFrame to display the results
dispersion_non_categorical_df = pd.DataFrame({
    'Variable': imputed_data_non_categorical.columns,
    'Standard Deviation': std_deviation_non_categorical.values
})

print(dispersion_non_categorical_df)

      Variable  Standard Deviation
0           S.no        17320.652413
1  tract_to_msamid_income         28.204752
2           population        1714.355870
3      minority_population        14.401545
4  number_of_owner_occupied_units        517.794667
5  number_of_1_to_4_family_units        737.753976
6      loan_amount_000s        604.958183
7  hud_median_family_income        12798.211781
8      applicant_income_000s        116.157479
9      sequence_number        150515.678152
10  census_tract_number        3356.259273
11      as_of_year        0.000000
12  application_date_indicator        0.225976
time: 16 ms (started: 2024-03-19 18:28:16 +00:00)

# Compute the correlation matrix for non-categorical variables
correlation_matrix_non_categorical =
imputed_data_non_categorical.corr()

# Print correlation matrix
print("Correlation Matrix (Non-Categorical Variables):")
print(correlation_matrix_non_categorical)

# Compute correlation coefficients for non-categorical variables
correlation_coefficients_non_categorical =

```

```

correlation_matrix_non_categorical.unstack().sort_values(ascending=False)

# Print correlation coefficients
print("\nCorrelation Coefficients (Non-Categorical Variables):")
print(correlation_coefficients_non_categorical)

# Test of correlation (using Pearson correlation coefficient) for specific variables
corr_tract_income, p_value_tract_income =
stats.pearsonr(imputed_data_non_categorical['tract_to_msamd_income'],
imputed_data_non_categorical['hud_median_family_income'])
corr_population_minority, p_value_population_minority =
stats.pearsonr(imputed_data_non_categorical['population'],
imputed_data_non_categorical['minority_population'])
corr_owner_units_loan, p_value_owner_units_loan =
stats.pearsonr(imputed_data_non_categorical['number_of_owner_occupied_units'],
imputed_data_non_categorical['loan_amount_000s'])

# Print correlation coefficients and p-values for specific variables
print("\nCorrelation Coefficient (tract_to_msamd_income vs
hud_median_family_income):", corr_tract_income)
print("P-value:", p_value_tract_income)

print("\nCorrelation Coefficient (population vs
minority_population):", corr_population_minority)
print("P-value:", p_value_population_minority)

print("\nCorrelation Coefficient (number_of_owner_occupied_units vs
loan_amount_000s):", corr_owner_units_loan)
print("P-value:", p_value_owner_units_loan)

Correlation Matrix (Non-Categorical Variables):
S.no tract_to_msamd_income
population \
S.no 1.000000 0.003807 -
0.008063
tract_to_msamd_income 0.003807 1.000000
0.048580
population -0.008063 0.048580
1.000000
minority_population 0.010181 -0.403257
0.184727
number_of_owner_occupied_units -0.032293 0.370891
0.777449
number_of_1_to_4_family_units -0.116825 0.132696
0.709542
loan_amount_000s 0.021756 0.081783
0.002061
hud_median_family_income 0.208449 -0.009365 -

```

0.012832		
applicant_income_000s	0.018066	0.184541
0.002256		
sequence_number	-0.038374	0.012611
0.004819		
census_tract_number	-0.190398	-0.022382
0.105298		
as_of_year	NaN	NaN
Nan		
application_date_indicator	-0.018438	0.011541
0.015836		
minority_population \		
S.no	0.010181	
tract_to_msamd_income	-0.403257	
population	0.184727	
minority_population	1.000000	
number_of_owner_occupied_units	-0.202225	
number_of_1_to_4_family_units	-0.175571	
loan_amount_000s	0.012030	
hud_median_family_income	0.242211	
applicant_income_000s	-0.031117	
sequence_number	-0.022516	
census_tract_number	-0.143931	
as_of_year	NaN	
application_date_indicator	-0.042659	
number_of_owner_occupied_units \		
S.no	-0.032293	
tract_to_msamd_income	0.370891	
population	0.777449	
minority_population	-0.202225	
number_of_owner_occupied_units	1.000000	
number_of_1_to_4_family_units	0.849083	
loan_amount_000s	0.002045	
hud_median_family_income	-0.079037	
applicant_income_000s	0.037969	
sequence_number	0.015118	
census_tract_number	0.019509	
as_of_year	NaN	
application_date_indicator	0.016227	
number_of_1_to_4_family_units \		
S.no	-0.116825	
tract_to_msamd_income	0.132696	
population	0.709542	
minority_population	-0.175571	
number_of_owner_occupied_units	0.849083	
number_of_1_to_4_family_units	1.000000	
loan_amount_000s	-0.032803	

hud_median_family_income	-0.298878	
applicant_income_000s	0.002561	
sequence_number	0.024564	
census_tract_number	0.205901	
as_of_year	NaN	
application_date_indicator	0.070299	
loan_amount_000s		
hud_median_family_income \ S.no	0.021756	
0.208449		
tract_to_msamd_income	0.081783	
0.009365		
population	0.002061	
0.012832		
minority_population	0.012030	
0.242211		
number_of_owner_occupied_units	0.002045	
0.079037		
number_of_1_to_4_family_units	-0.032803	
0.298878		
loan_amount_000s	1.000000	
0.128613		
hud_median_family_income	0.128613	
1.000000		
applicant_income_000s	0.146546	
0.153824		
sequence_number	0.011604	
0.060147		
census_tract_number	-0.040378	
0.493106		
as_of_year	NaN	
Nan		
application_date_indicator	-0.000406	
0.153317		
applicant_income_000s sequence_number		
\ S.no	0.018066	-0.038374
tract_to_msamd_income	0.184541	0.012611
population	0.002256	0.004819
minority_population	-0.031117	-0.022516
number_of_owner_occupied_units	0.037969	0.015118
number_of_1_to_4_family_units	0.002561	0.024564

loan_amount_000s	0.146546	0.011604
hud_median_family_income	0.153824	-0.060147
applicant_income_000s	1.000000	-0.002903
sequence_number	-0.002903	1.000000
census_tract_number	-0.033578	0.028415
as_of_year	NaN	NaN
application_date_indicator	-0.010670	0.124127
S.no	census_tract_number	as_of_year \
tract_to_msamd_income	-0.190398	NaN
population	-0.022382	NaN
minority_population	-0.105298	NaN
number_of_owner_occupied_units	0.019509	NaN
number_of_1_to_4_family_units	0.205901	NaN
loan_amount_000s	-0.040378	NaN
hud_median_family_income	-0.493106	NaN
applicant_income_000s	-0.033578	NaN
sequence_number	0.028415	NaN
census_tract_number	1.000000	NaN
as_of_year	NaN	NaN
application_date_indicator	0.149146	NaN
S.no	application_date_indicator	
tract_to_msamd_income	-0.018438	
population	0.011541	
minority_population	-0.015836	
number_of_owner_occupied_units	-0.042659	
number_of_1_to_4_family_units	0.016227	
loan_amount_000s	0.070299	
hud_median_family_income	-0.000406	
applicant_income_000s	-0.153317	
sequence_number	-0.010670	
census_tract_number	0.124127	
as_of_year	0.149146	
application_date_indicator	NaN	
application_date_indicator	1.000000	
Correlation Coefficients (Non-Categorical Variables):		
S.no	S.no	1.0
tract_to_msamd_income	tract_to_msamd_income	1.0
census_tract_number	census_tract_number	1.0
sequence_number	sequence_number	1.0

```
applicant_income_000s      applicant_income_000s      1.0
as_of_year                   sequence_number          NaN
                             census_tract_number    NaN
                             as_of_year            NaN
                             application_date_indicator  NaN
application_date_indicator  as_of_year            NaN
Length: 169, dtype: float64
```

```
Correlation Coefficient (tract_to_msamd_income vs
hud_median_family_income): -0.00936470130537446
P-value: 0.02179759460285561
```

```
Correlation Coefficient (population vs minority_population):
0.18472678320300429
P-value: 0.0
```

```
Correlation Coefficient (number_of_owner_occupied_units vs
loan_amount_000s): 0.0020453682357100228
P-value: 0.6163704631114627
time: 103 ms (started: 2024-03-19 18:28:20 +00:00)
```

```
...
```

Correlation Statistics (with Test of Correlation)

```
# Initialize correlation matrix and p-values matrix
correlation_matrix =
pd.DataFrame(index=imputed_data_non_categorical.columns,
columns=imputed_data_non_categorical.columns)
p_values = np.zeros((len(imputed_data_non_categorical.columns),
len(imputed_data_non_categorical.columns)))

# Calculate correlations and p-values
for i in range(len(imputed_data_non_categorical.columns)):
    for j in range(i, len(imputed_data_non_categorical.columns)):
        if imputed_data_non_categorical.iloc[:, i].nunique() > 1 and
imputed_data_non_categorical.iloc[:, j].nunique() > 1:
            corr, p_value =
pearsonr(imputed_data_non_categorical.iloc[:, i],
imputed_data_non_categorical.iloc[:, j])
            correlation_matrix.iloc[i, j] = corr
            correlation_matrix.iloc[j, i] = corr
            p_values[i, j] = p_value
            p_values[j, i] = p_value

# Create DataFrames for correlation coefficients and p-values
correlation_df = pd.DataFrame(correlation_matrix,
columns=imputed_data_non_categorical.columns,
index=imputed_data_non_categorical.columns)
```

```
p_values_df = pd.DataFrame(p_values,
columns=imputed_data_non_categorical.columns,
index=imputed_data_non_categorical.columns)
correlation_df
p_values_df
...
{"type":"string"}  
time: 13.2 ms (started: 2024-03-16 13:28:13 +00:00)  
...
# Missing Data Treatment: Non-Categorical Variables or Features  
  
# Impute missing values with the mean
df_non_categorical_imputed = noncatdf.fillna(noncatdf.mean())  
  
# Calculate missing data statistics for non-categorical columns
imputed_data_non_categorical =
df_non_categorical_imputed.isnull().sum().reset_index()
imputed_data_non_categorical.columns = ['Feature', 'Missing_Records']  
  
# Display the missing data statistics
print(imputed_data_non_categorical)
...  
  
{"type":"string"}  
time: 4.29 ms (started: 2024-03-16 13:28:13 +00:00)  
...
# Missing Data Treatment: Categorical Variables or Features  
  
# Impute missing values with the mode
df_categorical_imputed = catdf.apply(lambda x: x.fillna(x.mode()[0]))  
  
# Calculate missing data statistics for categorical columns
imputed_data_categorical =
df_categorical_imputed.isnull().sum().reset_index()
imputed_data_categorical.columns = ['Feature', 'Missing_Records']  
  
# Display the missing data statistics
print(imputed_data_categorical)
...  
  
{"type":"string"}  
time: 8.07 ms (started: 2024-03-16 13:28:13 +00:00)  
# Dataset Used : df_cat
```

```

si_cat = SimpleImputer(missing_values=np.nan,
strategy='most_frequent') # Strategy = median [When Odd Number of
Categories Exists]
si_cat_fit = si_cat.fit_transform(catdf)
imputed_data_categorical = pd.DataFrame(si_cat_fit,
columns=catdf.columns); # Missing Categorical Data Imputed Subset
imputed_data_categorical.info()
imputed_data_categorical.head()

<class 'pandas.core.frame.DataFrame'>
RangeIndex: 60000 entries, 0 to 59999
Data columns (total 21 columns):
 #   Column           Non-Null Count  Dtype  
--- 
 0   S.no             60000 non-null   object  
 1   state_name       60000 non-null   object  
 2   state_abbr       60000 non-null   object  
 3   respondent_id   60000 non-null   object  
 4   purchaser_type_name  60000 non-null   object  
 5   property_type_name  60000 non-null   object  
 6   preapproval_name  60000 non-null   object  
 7   owner_occupancy_name  60000 non-null   object  
 8   msamd_name       60000 non-null   object  
 9   loan_type_name   60000 non-null   object  
 10  loan_purpose_name  60000 non-null   object  
 11  lien_status_name  60000 non-null   object  
 12  hoepa_status_name  60000 non-null   object  
 13  county_name      60000 non-null   object  
 14  co_applicant_sex_name  60000 non-null   object  
 15  co_applicant_ethnicity_name  60000 non-null   object  
 16  applicant_sex_name  60000 non-null   object  
 17  applicant_ethnicity_name  60000 non-null   object  
 18  agency_name       60000 non-null   object  
 19  agency_abbr       60000 non-null   object  
 20  action_taken_name  60000 non-null   object  
dtypes: object(21)
memory usage: 9.6+ MB

{"type": "dataframe", "variable_name": "imputed_data_categorical"}

time: 463 ms (started: 2024-03-19 18:28:28 +00:00)

# ENCODING
# Converting Categorical Variable into Numeric

# Calculate the number of unique values in each column
unique_values_categorical =
imputed_data_categorical.nunique().reset_index()
unique_values_categorical.columns = ['Feature',
'Number_of_Unique_Values']

```

```

# Display the number of unique values
print(unique_values_categorical)

# Initialize LabelEncoder
label_encoder = LabelEncoder()

# Create a copy of the imputed_data_categorical dataframe to avoid
# modifying the original
encoded_data_categorical = imputed_data_categorical.copy()

# Iterate through each column in the dataframe
mapping = {} # To store the mapping of variable names to numeric
# representation

for column in encoded_data_categorical.columns:
    if column != "S.no": # Check if the column is not "S.no"
        # Perform numerical encoding
        encoded_data_categorical[column] =
label_encoder.fit_transform(encoded_data_categorical[column])

        # Store the mapping information
        mapping[column] = dict(zip(label_encoder.classes_,
label_encoder.transform(label_encoder.classes_)))

# Display the mapping
for variable, variable_mapping in mapping.items():
    print(f"\nMapping for {variable}:")
    print(variable_mapping)

# Display the encoded data
print(encoded_data_categorical)

      Feature  Number_of_Unique_Values
0          S.no                      60000
1      state_name                         1
2      state_abbr                         1
3  respondent_id                        593
4  purchaser_type_name                   10
5  property_type_name                   3
6  preapproval_name                     3
7  owner_occupancy_name                  3
8      msamd_name                       14
9      loan_type_name                     4
10     loan_purpose_name                  3
11     lien_status_name                   4
12     hoepa_status_name                  2
13      county_name                      39

```

14	co_applicant_sex_name	5
15	co_applicant_ethnicity_name	5
16	applicant_sex_name	4
17	applicant_ethnicity_name	4
18	agency_name	6
19	agency_abbr	6
20	action_taken_name	8

Mapping for state_name:

```
{'Washington': 0}
```

Mapping for state_abbr:

```
{'WA': 0}
```

Mapping for respondent_id:

```
{'01-0681100': 0, '01-0726495': 1, '02-0793125': 2, '03-0488052': 3,
'04-3212636': 4, '04-3568208': 5, '04-3660901': 6, '04-7534967': 7,
'05-0402708': 8, '06-1016329': 9, '1015560': 10, '10257': 11, '1047': 12,
'1097500000': 13, '1099800006': 14, '11-3399725': 15, '11-3412303': 16,
'11-3714032': 17, '11162': 18, '112837': 19, '11443': 20,
'1146500007': 21, '11729': 22, '11734': 23, '12135': 24,
'1216826': 25, '12219': 26, '1227300009': 27, '12311': 28, '1265': 29,
'1281': 30, '13-3222578': 31, '13-3602661': 32, '13-3753941': 33, '13-4225190': 34, '13-4362989': 35, '13-6131491': 36, '13232': 37,
'13303': 38, '13392': 39, '13778': 40, '14-1841762': 41, '14206': 42,
'14252': 43, '143662': 44, '1461700004': 45, '14662': 46, '146672': 47,
'14740': 48, '14843': 49, '151': 50, '15219': 51, '15732': 52,
'16-1686740': 53, '16243': 54, '1635900004': 55, '16450': 56, '16814': 57,
'169653': 58, '17587': 59, '17672': 60, '177957': 61, '17874': 62,
'17884': 63, '1842065': 64, '19307': 65, '19628': 66, '197478': 67,
'19899': 68, '19976': 69, '20-0142846': 70, '20-0192872': 71, '20-0304793': 72, '20-0640473': 73, '20-0740151': 74, '20-1255434': 75,
'20-1832276': 76, '20-2053401': 77, '20-2355296': 78, '20-2470783': 79,
'20-2471369': 80, '20-2485875': 81, '20-2693054': 82, '20-2718340': 83, '20-2752826': 84, '20-2928975': 85, '20-3702275': 86,
'20-3828708': 87, '20-4136310': 88, '20-4224234': 89, '20-4255880': 90,
'20-4866754': 91, '20-5238443': 92, '20-5239910': 93, '20-5741925': 94, '20-8006279': 95, '20-8083209': 96, '20-8544905': 97,
'20-8745846': 98, '20-8803449': 99, '20-8921389': 100, '2003500009': 101,
'20061': 102, '20068': 103, '20214': 104, '20516': 105, '20624': 106,
'20774': 107, '210434': 108, '21122': 109, '212465': 110,
'2137100009': 111, '2149009991': 112, '21717': 113, '2191': 114,
'2193616': 115, '22-3039688': 116, '22-3470404': 117, '22-3554558': 118,
'22-3626426': 119, '22-3747694': 120, '22-3887207': 121, '22134': 122,
'22157': 123, '22407': 124, '22444': 125, '22637': 126, '2285': 127,
'22939': 128, '23-2470039': 129, '23-2769131': 130, '23041': 131,
'2317700005': 132, '23216': 133, '23416': 134, '23521': 135, '23850': 136,
'23922': 137, '23957': 138, '24077': 139, '24080': 140, '24107': 141,
'24169': 142, '24224': 143, '24235': 144, '24326': 145, '24382': 146,
'24671': 147, '24708': 148, '24713': 149, '24719': 150, '24753': 151}
```

151, '24760': 152, '24831': 153, '24849': 154, '2489805': 155, '25080': 156, '25093': 157, '25103': 158, '2562164': 159, '2590037': 160, '26-0012825': 161, '26-0021318': 162, '26-0335190': 163, '26-0360466': 164, '26-0362771': 165, '26-0423240': 166, '26-0455770': 167, '26-0508430': 168, '26-0595342': 169, '26-0707492': 170, '26-1242154': 171, '26-1334020': 172, '26-1589507': 173, '26-1773722': 174, '26-2049351': 175, '26-2261031': 176, '26-2593704': 177, '26-2689428': 178, '26-2916887': 179, '26-3264687': 180, '26-3416474': 181, '26-3780954': 182, '26-4193875': 183, '26-4461592': 184, '26-4558390': 185, '26-4599244': 186, '2618780': 187, '26610': 188, '267100004': 189, '27-0222046': 190, '27-0267182': 191, '27-0684906': 192, '27-0812934': 193, '27-1190043': 194, '27-1438405': 195, '27-2389039': 196, '27-3459350': 197, '27-4023565': 198, '27-4626537': 199, '27280': 200, '2734': 201, '2735146': 202, '27601': 203, '276579': 204, '280110': 205, '28116': 206, '28151': 207, '28316': 208, '28405': 209, '28453': 210, '28454': 211, '28489': 212, '28599': 213, '2888798': 214, '29012': 215, '29058': 216, '29209': 217, '2945': 218, '3027509990': 219, '3076248': 220, '30788': 221, '30810': 222, '31-1197926': 223, '31-1690008': 224, '31-1712553': 225, '311845': 226, '3121': 227, '31286': 228, '3150447': 229, '319': 230, '32-0016270': 231, '3212149': 232, '32178': 233, '32489': 234, '3284070': 235, '33-0231744': 236, '33-0397503': 237, '33-0419992': 238, '33-0594693': 239, '33-0750812': 240, '33-0816610': 241, '33-0828099': 242, '33-0941669': 243, '33-0962918': 244, '33-0975529': 245, '33183': 246, '33508': 247, '3374412': 248, '33806': 249, '33826': 250, '339858': 251, '34-1194858': 252, '34-1633105': 253, '34-1716542': 254, '34-1719615': 255, '34-2000096': 256, '34106': 257, '34214': 258, '34585': 259, '34607': 260, '34627': 261, '34953': 262, '35-2486440': 263, '35013': 264, '35014': 265, '35139': 266, '35261': 267, '35355': 268, '3537897': 269, '35406': 270, '36-4327855': 271, '365325': 272, '37-1493496': 273, '37-1542226': 274, '38-2434249': 275, '38-2749215': 276, '38-2750395': 277, '38-2799035': 278, '38-3564305': 279, '39-2001010': 280, '3918898': 281, '3938186': 282, '3956': 283, '4004574': 284, '41-1795868': 285, '41-1914032': 286, '41-2277737': 287, '4114567': 288, '413208': 289, '4142': 290, '42-1554181': 291, '42-1739728': 292, '4239': 293, '43-0951349': 294, '43-1965151': 295, '435': 296, '45-3143795': 297, '45-3508295': 298, '451965': 299, '46-0530020': 300, '46-1728831': 301, '46-1836968': 302, '46-2477192': 303, '46-2888046': 304, '46-3435079': 305, '46-3528270': 306, '46-3676810': 307, '46-5671661': 308, '47-0873092': 309, '47-0912342': 310, '47-0933090': 311, '47-3632618': 312, '471809999': 313, '476810': 314, '48-1148159': 315, '48-1236121': 316, '480228': 317, '4878': 318, '491224': 319, '501105': 320, '504713': 321, '51-0488301': 322, '51-0517525': 323, '52-2091594': 324, '52-2276553': 325, '52-2321476': 326, '52-2323186': 327, '5380': 328, '54-0259290': 329, '54-1094297': 330, '54-1994393': 331, '54-2070914': 332, '542409990': 333, '542649': 334, '546571': 335, '55130': 336, '5556209999': 337, '5582': 338, '5588': 339, '56-2103469': 340, '56-2237729': 341, '56-2471041': 342, '566': 343, '57-1175755': 344, '57033': 345, '57071': 346, '57074':

347, '57167': 348, '57451': 349, '57542': 350, '57607': 351, '57614': 352, '57633': 353, '57776': 354, '57777': 355, '57949': 356, '57955': 357, '57978': 358, '58-1865166': 359, '58305': 360, '58322': 361, '58341': 362, '58380': 363, '58778': 364, '59-3378746': 365, '5912': 366, '5913': 367, '592448': 368, '595270': 369, '595869': 370, '60042': 371, '60059': 372, '60079': 373, '601050': 374, '60143': 375, '60438': 376, '606046': 377, '60613': 378, '6158': 379, '6161': 380, '617677': 381, '619877': 382, '62-0997810': 383, '62-1494087': 384, '62-1532940': 385, '624': 386, '6248': 387, '62659': 388, '62665': 389, '62745': 390, '6288': 391, '63-1052225': 392, '63069': 393, '63194': 394, '63196': 395, '6328': 396, '63315': 397, '63440': 398, '63799': 399, '64103': 400, '644': 401, '6443809990': 402, '64482': 403, '64546': 404, '64552': 405, '65595': 406, '656377': 407, '65644': 408, '656733': 409, '66157': 410, '66328': 411, '66331': 412, '66337': 413, '66349': 414, '66373': 415, '66399': 416, '665': 417, '66734': 418, '66751': 419, '66841': 420, '67201': 421, '67262': 422, '67264': 423, '67389': 424, '675332': 425, '67911': 426, '67955': 427, '68-0151632': 428, '68-0295876': 429, '68-0309242': 430, '68061': 431, '68095': 432, '68186': 433, '68187': 434, '68196': 435, '68203': 436, '68205': 437, '68222': 438, '68223': 439, '68237': 440, '68239': 441, '68253': 442, '68255': 443, '68271': 444, '68278': 445, '68284': 446, '68293': 447, '68298': 448, '68304': 449, '68315': 450, '68362': 451, '68375': 452, '68394': 453, '68423': 454, '68457': 455, '68465': 456, '68517': 457, '68530': 458, '68546': 459, '685676': 460, '68576': 461, '68598': 462, '68601': 463, '68613': 464, '68622': 465, '694904': 466, '697633': 467, '702825': 468, '702889': 469, '703927': 470, '704347': 471, '7056000000': 472, '706707': 473, '706809': 474, '707674': 475, '708146': 476, '708186': 477, '7101': 478, '712504': 479, '713570': 480, '713964': 481, '714970': 482, '715018': 483, '716195': 484, '7162800002': 485, '716456': 486, '717936': 487, '7197000003': 488, '72-1545376': 489, '7257500009': 490, '7272800006': 491, '7275700004': 492, '73-1374559': 493, '73-1545233': 494, '73-1577221': 495, '74-2508160': 496, '75-2585326': 497, '75-2695327': 498, '75-2921540': 499, '75-3170028': 500, '7505400005': 501, '7516800003': 502, '75633': 503, '76-0236067': 504, '76-0503625': 505, '76-0561995': 506, '76-0629353': 507, '7635500004': 508, '7638200000': 509, '7667200009': 510, '7674000006': 511, '77-0158990': 512, '77-0605392': 513, '77-0672274': 514, '77-0717225': 515, '7748': 516, '7756300009': 517, '7810600004': 518, '7811300008': 519, '7927200007': 520, '7983500003': 521, '7992700007': 522, '80-0233937': 523, '80-0312140': 524, '80-0860209': 525, '80122': 526, '804963': 527, '8100': 528, '817824': 529, '83-0171636': 530, '83-0368862': 531, '84-0927358': 532, '84-1040263': 533, '84-1412422': 534, '84-1496821': 535, '84-1564935': 536, '84-1594306': 537, '85-0260899': 538, '852218': 539, '852320': 540, '857': 541, '86-0415227': 542, '86-0431588': 543, '86-0634557': 544, '86-0860478': 545, '8663': 546, '87-0623581': 547, '87-0675992': 548, '87-0682600': 549, '87-0691650': 550, '8796': 551, '8797': 552, '88-0209429': 553, '88-0508228': 554, '8854': 555, '90-0790926': 556, '91-1374387': 557, '91-1395192': 558, '91-1441009': 559, '91-1465333':

```
560, '91-1529683': 561, '91-1569077': 562, '91-1780488': 563, '91-1841798': 564, '91-1913382': 565, '91-2006136': 566, '915878': 567, '9289': 568, '93-1231049': 569, '93-1248952': 570, '93-1296762': 571, '93-1301081': 572, '934329': 573, '9366': 574, '936855': 575, '9373': 576, '94-3195577': 577, '9483': 578, '9486': 579, '95-3821253': 580, '95-3990375': 581, '95-4196389': 582, '95-4234730': 583, '95-4267987': 584, '95-4462959': 585, '95-4482547': 586, '95-4523866': 587, '95-4623407': 588, '95-4762204': 589, '95-4769926': 590, '95-4866828': 591, '972590': 592}
```

Mapping for purchaser_type_name:

```
{'Affiliate institution': 0, 'Commercial bank, savings bank or savings association': 1, 'Fannie Mae (FNMA)': 2, 'Farmer Mac (FAMC)': 3, 'Freddie Mac (FHLMC)': 4, 'Ginnie Mae (GNMA)': 5, 'Life insurance company, credit union, mortgage bank, or finance company': 6, 'Loan was not originated or was not sold in calendar year covered by register': 7, 'Other type of purchaser': 8, 'Private securitization': 9}
```

Mapping for property_type_name:

```
{'Manufactured housing': 0, 'Multifamily dwelling': 1, 'One-to-four family dwelling (other than manufactured housing)': 2}
```

Mapping for preapproval_name:

```
{'Not applicable': 0, 'Preapproval was not requested': 1, 'Preapproval was requested': 2}
```

Mapping for owner_occupancy_name:

```
{'Not applicable': 0, 'Not owner-occupied as a principal dwelling': 1, 'Owner-occupied as a principal dwelling': 2}
```

Mapping for msamid_name:

```
{'Bellingham - WA': 0, 'Bremerton, Silverdale - WA': 1, 'Kennewick, Richland - WA': 2, 'Lewiston - ID, WA': 3, 'Longview - WA': 4, 'Mount Vernon, Anacortes - WA': 5, 'Olympia, Tumwater - WA': 6, 'Portland, Vancouver, Hillsboro - OR, WA': 7, 'Seattle, Bellevue, Everett - WA': 8, 'Spokane, Spokane Valley - WA': 9, 'Tacoma, Lakewood - WA': 10, 'Walla Walla - WA': 11, 'Wenatchee - WA': 12, 'Yakima - WA': 13}
```

Mapping for loan_type_name:

```
{'Conventional': 0, 'FHA-insured': 1, 'FSA/RHS-guaranteed': 2, 'VA-guaranteed': 3}
```

Mapping for loan_purpose_name:

```
{'Home improvement': 0, 'Home purchase': 1, 'Refinancing': 2}
```

Mapping for lien_status_name:

```
{'Not applicable': 0, 'Not secured by a lien': 1, 'Secured by a first lien': 2, 'Secured by a subordinate lien': 3}
```

```
Mapping for hoepa_status_name:  
{'HOEPA loan': 0, 'Not a HOEPA loan': 1}  
  
Mapping for county_name:  
{'Adams County': 0, 'Asotin County': 1, 'Benton County': 2, 'Chelan  
County': 3, 'Clallam County': 4, 'Clark County': 5, 'Columbia County':  
6, 'Cowlitz County': 7, 'Douglas County': 8, 'Ferry County': 9,  
'Franklin County': 10, 'Garfield County': 11, 'Grant County': 12,  
'Grays Harbor County': 13, 'Island County': 14, 'Jefferson County':  
15, 'King County': 16, 'Kitsap County': 17, 'Kittitas County': 18,  
'Klickitat County': 19, 'Lewis County': 20, 'Lincoln County': 21,  
'Mason County': 22, 'Okanogan County': 23, 'Pacific County': 24, 'Pend  
Oreille County': 25, 'Pierce County': 26, 'San Juan County': 27,  
'Skagit County': 28, 'Skamania County': 29, 'Snohomish County': 30,  
'Spokane County': 31, 'Stevens County': 32, 'Thurston County': 33,  
'Wahkiakum County': 34, 'Walla Walla County': 35, 'Whatcom County':  
36, 'Whitman County': 37, 'Yakima County': 38}  
  
Mapping for co_applicant_sex_name:  
{'Female': 0, 'Information not provided by applicant in mail,  
Internet, or telephone application': 1, 'Male': 2, 'No co-applicant':  
3, 'Not applicable': 4}  
  
Mapping for co_applicant_ethnicity_name:  
{'Hispanic or Latino': 0, 'Information not provided by applicant in  
mail, Internet, or telephone application': 1, 'No co-applicant': 2,  
'Not Hispanic or Latino': 3, 'Not applicable': 4}  
  
Mapping for applicant_sex_name:  
{'Female': 0, 'Information not provided by applicant in mail,  
Internet, or telephone application': 1, 'Male': 2, 'Not applicable':  
3}  
  
Mapping for applicant_ethnicity_name:  
{'Hispanic or Latino': 0, 'Information not provided by applicant in  
mail, Internet, or telephone application': 1, 'Not Hispanic or  
Latino': 2, 'Not applicable': 3}  
  
Mapping for agency_name:  
{'Consumer Financial Protection Bureau': 0, 'Department of Housing and  
Urban Development': 1, 'Federal Deposit Insurance Corporation': 2,  
'Federal Reserve System': 3, 'National Credit Union Administration':  
4, 'Office of the Comptroller of the Currency': 5}  
  
Mapping for agency_abbr:  
{'CFPB': 0, 'FDIC': 1, 'FRS': 2, 'HUD': 3, 'NCUA': 4, 'OCC': 5}  
  
Mapping for action_taken_name:  
{'Application approved but not accepted': 0, 'Application denied by  
financial institution': 1, 'Application withdrawn by applicant': 2,
```

```
'File closed for incompleteness': 3, 'Loan originated': 4, 'Loan purchased by the institution': 5, 'Preapproval request approved but not accepted': 6, 'Preapproval request denied by financial institution': 7}
    S.no state_name state_abbr respondent_id
purchaser_type_name \
0      1          0          0        317
4
1      2          0          0        490
6
2      3          0          0        489
7
3      4          0          0        318
7
4      5          0          0        234
4
4
...
59995  59996          0          0        472
8
59996  59997          0          0        488
4
59997  59998          0          0        55
5
59998  59999          0          0        116
5
59999  60000          0          0         47
2

    property_type_name preapproval_name owner_occupancy_name
msamd_name \
0                  2          0          2
7
1                  2          0          2
11
2                  2          0          2
7
3                  2          0          2
7
4                  2          0          2
1
1
...
59995              2          0          2
7
59996              2          0          1
5
59997              2          0          2
0
```

59998	2	0	2
1			
59999	2	0	2
7			
loan_type_name ... lien_status_name hoepa_status_name			
county_name \	0 ...	2	1
0	0 ...	2	1
5			
1	1 ...	2	1
35			
2	0 ...	2	1
5			
3	0 ...	2	1
5			
4	0 ...	2	1
17			
...
...			
59995	0 ...	2	1
5			
59996	0 ...	2	1
28			
59997	1 ...	2	1
36			
59998	3 ...	2	1
17			
59999	0 ...	2	1
5			
co_applicant_sex_name co_applicant_ethnicity_name			
applicant_sex_name \	2	3	
0	2	3	
0			
1	3	2	
2			
2	0	3	
2			
3	0	1	
2			
4	2	3	
0			
...
...			
59995	0	0	
2			
59996	3	2	
0			
59997	0	3	

```
2  
59998 0 3  
2  
59999 0 3  
2  
  
    applicant_ethnicity_name agency_name agency_abbr  
action_taken_name  
0 2 0 0  
4  
1 0 1 3  
4  
2 2 1 3  
4  
3 1 4 4  
4  
4 2 2 1  
4  
4  
... ... ... ...  
...  
59995 2 1 3  
4  
59996 1 1 3  
4  
59997 2 1 3  
4  
59998 2 0 0  
4  
59999 2 0 0  
4  
  
[60000 rows x 21 columns]  
time: 1 s (started: 2024-03-19 18:29:23 +00:00)  
  
print(imputed_data_non_categorical.columns)  
  
Index(['S.no', 'tract_to_msamd_income', 'population',  
'minority_population',  
       'number_of_owner_occupied_units',  
'number_of_1_to_4_family_units',  
       'loan_amount_000s', 'hud_median_family_income',  
'applicant_income_000s',  
       'sequence_number', 'census_tract_number', 'as_of_year',  
       'application_date_indicator'],  
      dtype='object')  
time: 801 µs (started: 2024-03-19 18:29:26 +00:00)  
  
def identify_outliers(column):  
    Q1 = np.percentile(column, 25)  
    Q3 = np.percentile(column, 75)
```

```

IQR = Q3 - Q1
lower_bound = Q1 - 1.5 * IQR
upper_bound = Q3 + 1.5 * IQR
outliers = (column < lower_bound) | (column > upper_bound)
return outliers

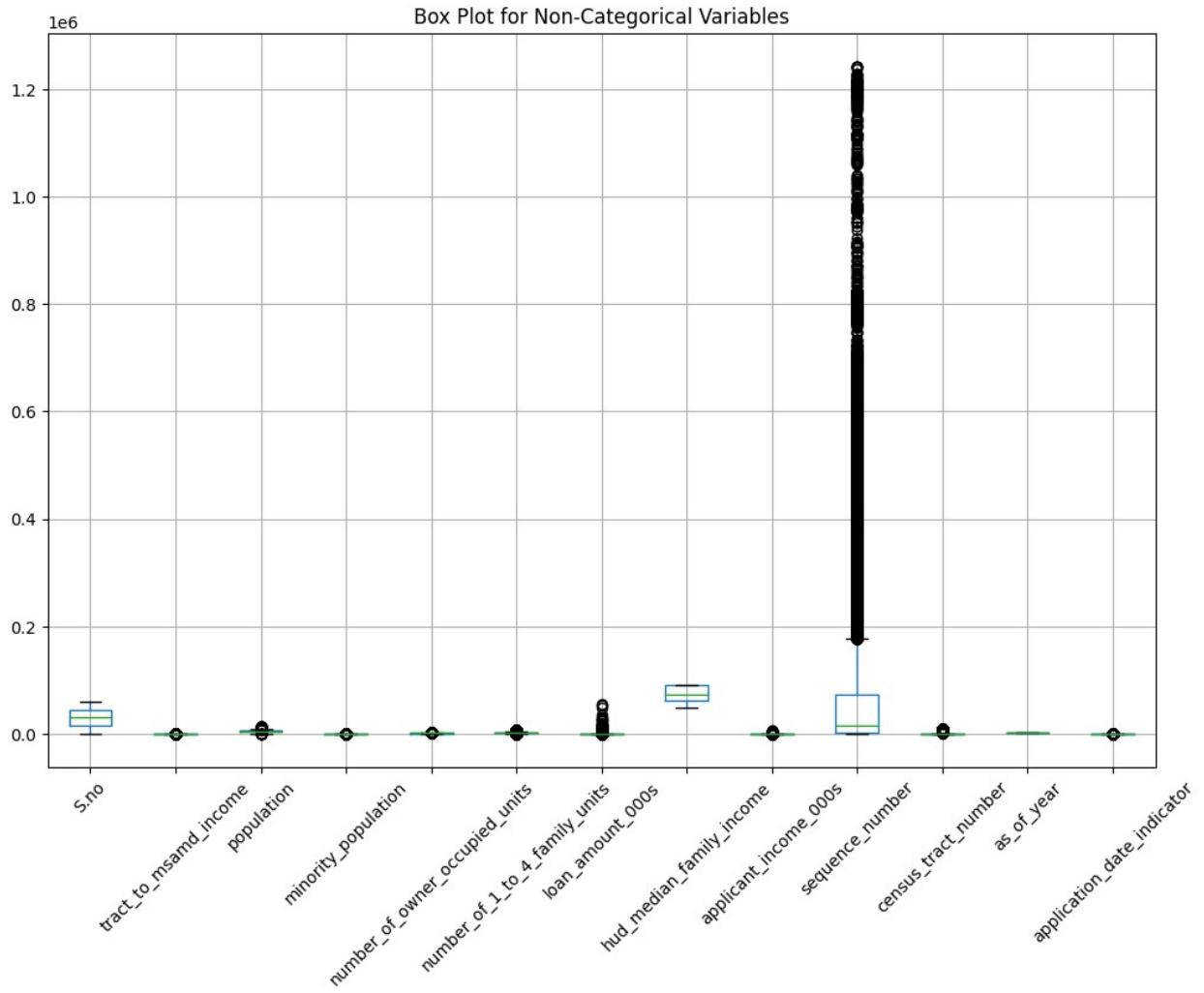
# Apply the function to each column to get a DataFrame of True/False values
outliers = imputed_data_non_categorical.apply(identify_outliers)

# Display the number of outliers for each column
outlier_counts = outliers.sum()
print(outlier_counts)

S.no                      0
tract_to_msamd_income     1309
population                 553
minority_population        2641
number_of_owner_occupied_units   478
number_of_1_to_4_family_units    2150
loan_amount_000s            2467
hud_median_family_income      0
applicant_income_000s         3765
sequence_number              7898
census_tract_number          9391
as_of_year                   0
application_date_indicator    776
dtype: int64
time: 40.5 ms (started: 2024-03-19 18:29:38 +00:00)

# Create a box plot for non-categorical variables
imputed_data_non_categorical.boxplot(rot=45, figsize=(12, 8))
plt.title('Box Plot for Non-Categorical Variables')
plt.show()

```



```
time: 1.04 s (started: 2024-03-19 18:29:46 +00:00)
```

```
# Iterate through each column and print count of unique values
for column in imputed_data_non_categorical.columns:
    unique_count = imputed_data_non_categorical[column].nunique()
    print(f"Count of unique values in {column} column: {unique_count}")

Count of unique values in S.no column: 60000
Count of unique values in tract_to_msamd_income column: 1328
Count of unique values in population column: 1284
Count of unique values in minority_population column: 1217
Count of unique values in number_of_owner_occupied_units column: 997
Count of unique values in number_of_1_to_4_family_units column: 1052
Count of unique values in loan_amount_000s column: 1344
Count of unique values in hud_median_family_income column: 16
Count of unique values in applicant_income_000s column: 808
Count of unique values in sequence_number column: 39340
Count of unique values in census_tract_number column: 1109
```

```

Count of unique values in as_of_year column: 1
Count of unique values in application_date_indicator column: 2
time: 36.7 ms (started: 2024-03-19 18:29:54 +00:00)

# Initialize the StandardScaler
scaler = StandardScaler()

# Apply Standard Scaling to your dataset
scaled_data = scaler.fit_transform(imputed_data_non_categorical)

def identify_outliers(column):
    Q1 = np.percentile(column, 25)
    Q3 = np.percentile(column, 75)
    IQR = Q3 - Q1
    lower_bound = Q1 - 1.5 * IQR
    upper_bound = Q3 + 1.5 * IQR
    outliers = (column < lower_bound) | (column > upper_bound)
    return outliers

# Apply the function to each column in the scaled dataset
outliers_scaled = pd.DataFrame(scaled_data,
columns=imputed_data_non_categorical.columns).apply(identify_outliers)

# Display the number of outliers for each column in the scaled dataset
outlier_counts_scaled = outliers_scaled.sum()
print(outlier_counts_scaled)

S.no          0
tract_to_msamd_income 1309
population      553
minority_population 2641
number_of_owner_occupied_units 478
number_of_1_to_4_family_units 2150
loan_amount_000s 2467
hud_median_family_income 0
applicant_income_000s 3765
sequence_number 7898
census_tract_number 9391
as_of_year        0
application_date_indicator 776
dtype: int64
time: 51 ms (started: 2024-03-19 18:30:03 +00:00)

# Initialize the RobustScaler
scaler = RobustScaler()

# Apply Robust Scaling to your dataset
scaled_data_robust =
scaler.fit_transform(imputed_data_non_categorical)

# Check for outliers in the scaled dataset

```

```

outliers_robust = pd.DataFrame(scaled_data_robust,
columns=imputed_data_non_categorical.columns).apply(identify_outliers)

# Display the number of outliers for each column in the scaled dataset
outlier_counts_robust = outliers_robust.sum()
print(outlier_counts_robust)

S.no          0
tract_to_msamd_income    1309
population      553
minority_population  2641
number_of_owner_occupied_units 478
number_of_1_to_4_family_units   2150
loan_amount_000s        2467
hud_median_family_income 0
applicant_income_000s     3765
sequence_number       7898
census_tract_number     9391
as_of_year           0
application_date_indicator 776
dtype: int64
time: 80.4 ms (started: 2024-03-19 18:30:21 +00:00)

# Define columns to exclude from normalization
columns_to_exclude = ['S.no', 'hud_median_family_income', 'as_of_year',
'application_date_indicator']

# Create a copy of the DataFrame with excluded columns
data_to_scale =
imputed_data_non_categorical.drop(columns=columns_to_exclude)

# Initialize the MinMaxScaler
scaler = MinMaxScaler()

# Apply Min-Max Scaling to the selected columns
scaled_data = scaler.fit_transform(data_to_scale)

# Create a DataFrame with scaled data and original column names
scaled_df = pd.DataFrame(scaled_data, columns=data_to_scale.columns)

# Add back the excluded columns to the scaled DataFrame
scaled_df[columns_to_exclude] =
imputed_data_non_categorical[columns_to_exclude]

# Display the scaled DataFrame
print(scaled_df)

      tract_to_msamd_income  population  minority_population \
0            0.442799      0.640752        0.234501
1            0.285162      0.372631        0.236658
2            0.317084      0.385008        0.105445

```

3	0.543502	0.381682	0.070620
4	0.610556	0.393363	0.091213
...
59995	0.394915	0.299373	0.149434
59996	0.445679	0.369923	0.060162
59997	0.418734	0.609964	0.105553
59998	0.336419	0.179392	0.322803
59999	0.442799	0.640752	0.234501
number_of_owner_occupied_units			
number_of_1_to_4_family_units \			
0	0.724346		0.448858
1	0.420188		0.298329
2	0.375922		0.308728
3	0.506372		0.305660
4	0.566734		0.354074
...
59995	0.370557		0.240709
59996	0.556673		0.383396
59997	0.783032		0.558814
59998	0.195171		0.116263
59999	0.724346		0.448858
loan_amount_000s applicant_income_000s sequence_number \			
0	0.004109	0.018669	0.096625
1	0.004346	0.006656	0.042368
2	0.004364	0.018831	0.005001
3	0.006364	0.050974	0.000158
4	0.007564	0.018344	0.026241
...
59995	0.002927	0.012013	0.024452
59996	0.002564	0.007955	0.268165
59997	0.004546	0.014123	0.020504
59998	0.004655	0.018153	0.289037
59999	0.005309	0.013149	0.032247
census_tract_number S.no hud_median_family_income			
as_of_year \			
0	0.042258	1.0	73300.0
2016.0			

1	0.943728	2.0	57900.0
2016.0			
2	0.042333	3.0	73300.0
2016.0			
3	0.041421	4.0	73300.0
2016.0			
4	0.092866	5.0	78100.0
2016.0			
...
...			
59995	0.041926	59996.0	73300.0
2016.0			
59996	0.963715	59997.0	61400.0
2016.0			
59997	0.000724	59998.0	69900.0
2016.0			
59998	0.082002	59999.0	78100.0
2016.0			
59999	0.042258	60000.0	73300.0
2016.0			
	application_date_indicator		
0	0.0		
1	0.0		
2	0.0		
3	0.0		
4	0.0		
...	...		
59995	0.0		
59996	0.0		
59997	0.0		
59998	0.0		
59999	0.0		

[60000 rows x 13 columns]
time: 26.4 ms (started: 2024-03-19 18:30:44 +00:00)

```
def identify_outliers(column):
    Q1 = np.percentile(column, 25)
    Q3 = np.percentile(column, 75)
    IQR = Q3 - Q1
    lower_bound = Q1 - 1.5 * IQR
    upper_bound = Q3 + 1.5 * IQR
    outliers = (column < lower_bound) | (column > upper_bound)
    return outliers

# Apply the function to each column in the scaled dataset
scaled_outliers = scaled_df.apply(identify_outliers)

# Display the number of outliers for each column in the scaled dataset
```

```

scaled_outlier_counts = scaled_outliers.sum()
print(scaled_outlier_counts)

tract_to_msamd_income           1309
population                      553
minority_population              2641
number_of_owner_occupied_units   478
number_of_1_to_4_family_units    2150
loan_amount_000s                 2467
applicant_income_000s             3765
sequence_number                  7898
census_tract_number               9391
S.no                            0
hud_median_family_income         0
as_of_year                       0
application_date_indicator      776
dtype: int64
time: 41.3 ms (started: 2024-03-19 18:30:53 +00:00)

def identify_outliers(column):
    Q1 = np.percentile(column, 25)
    Q3 = np.percentile(column, 75)
    IQR = Q3 - Q1
    lower_bound = Q1 - 1.5 * IQR
    upper_bound = Q3 + 1.5 * IQR
    outliers1 = (column < lower_bound) | (column > upper_bound)
    return outliers1

# Apply the function to each column to get a DataFrame of True/False values
outliers1 = scaled_df.apply(identify_outliers)

# Display the number of outliers for each column
outlier_counts1 = outliers1.sum()
print(outlier_counts1)

# Create a box plot for non-categorical variables
scaled_df.boxplot(rot=45, figsize=(12, 8))
plt.title('Box Plot for Non-Categorical Variables')
plt.show()

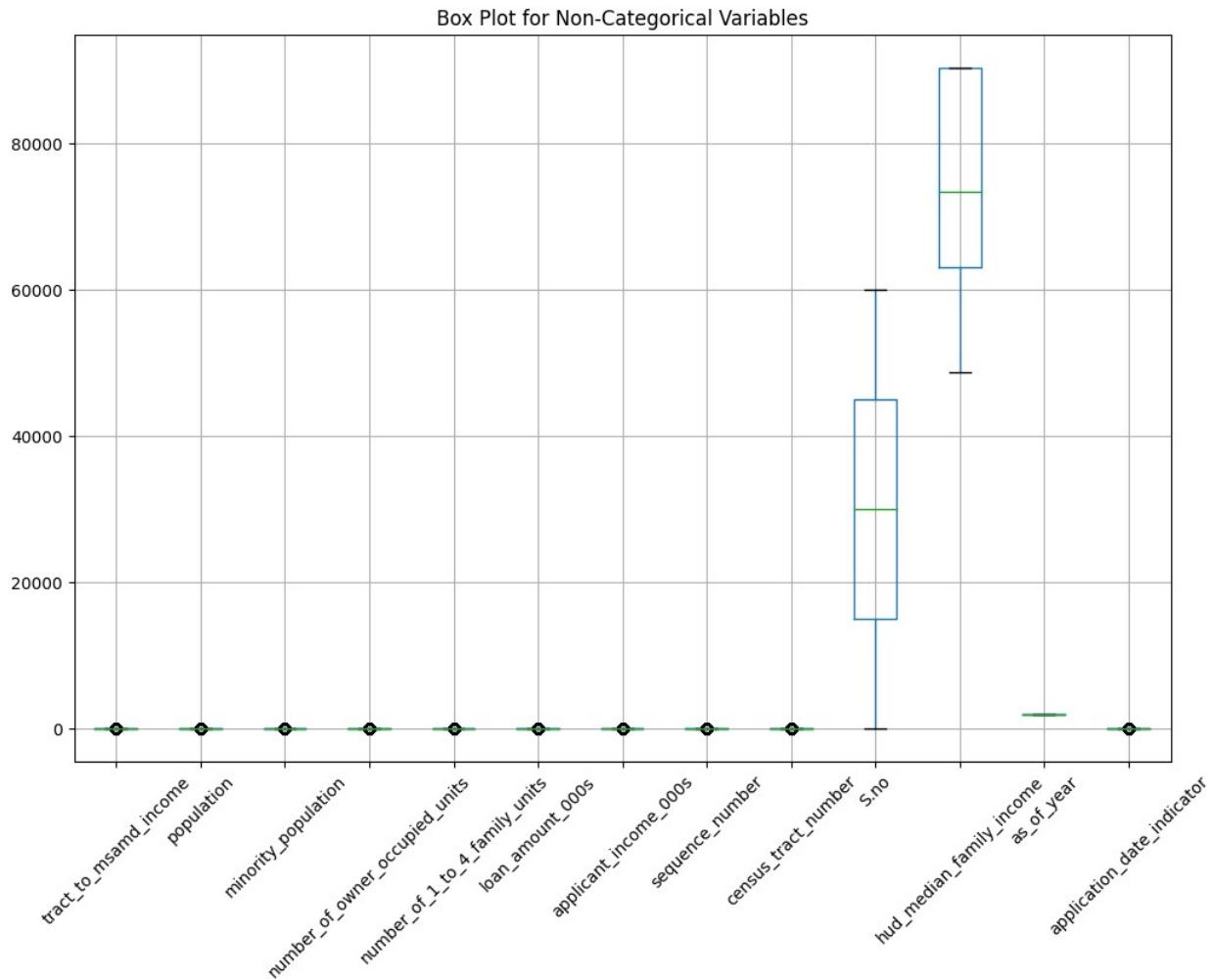
tract_to_msamd_income           1309
population                      553
minority_population              2641
number_of_owner_occupied_units   478
number_of_1_to_4_family_units    2150
loan_amount_000s                 2467
applicant_income_000s             3765
sequence_number                  7898

```

```

census_tract_number          9391
S.no                         0
hud_median_family_income     0
as_of_year                    0
application_date_indicator   776
dtype: int64

```



```
time: 1.03 s (started: 2024-03-19 18:30:58 +00:00)
```

```

# Calculate standard deviation for non-categorical columns
std_deviation_non_categorical1 = scaled_df.std()

# Creating a DataFrame to display the results
dispersion_non_categorical_df1 = pd.DataFrame({
    'Variable': scaled_df.columns,
    'Standard Deviation': std_deviation_non_categorical1.values
})

print(dispersion_non_categorical_df1)

```

```

      Variable Standard Deviation
0    tract_to_msamd_income      0.116026
1           population        0.132618
2   minority_population       0.155273
3 number_of_owner_occupied_units 0.173640
4 number_of_1_to_4_family_units 0.125768
5          loan_amount_000s     0.010999
6 applicant_income_000s        0.018857
7         sequence_number      0.121228
8 census_tract_number        0.344020
9                 S.no        17320.652413
10      hud_median_family_income 12798.211781
11           as_of_year        0.000000
12 application_date_indicator 0.225976
time: 15.1 ms (started: 2024-03-19 18:31:06 +00:00)

# Pre-Processed Dataset
combined_data = pd.merge(encoded_data_categorical, scaled_df,
on='S.no')

# Display the Pre-Processed Dataset
%memit
combined_data

peak memory: 454.77 MiB, increment: 0.02 MiB
{"type": "dataframe", "variable_name": "combined_data"}

time: 346 ms (started: 2024-03-19 18:31:11 +00:00)
...
# Assuming encoded_data_categorical and scaled_df have the same number of rows
combined_data = pd.concat([encoded_data_categorical, scaled_df],
axis=1)

# Display the combined DataFrame
print(combined_data)

list(combined_data.columns)
...
{"type": "string"}

time: 8.32 ms (started: 2024-03-19 12:37:54 +00:00)
df_ppd_subset = combined_data.copy()

time: 5.43 ms (started: 2024-03-19 18:31:46 +00:00)

# Create K-Means Clusters [K=2]
km_2cluster = kmclus(n_clusters=2, init='random', random_state=333)

```

```

df_ppd_subset['Cluster_Label'] =
km_2cluster.fit_predict(df_ppd_subset)
km_2cluster_model = df_ppd_subset['Cluster_Label'].values
km_2cluster_model

/usr/local/lib/python3.10/dist-packages/sklearn/cluster/
_kmeans.py:870: FutureWarning: The default value of `n_init` will
change from 10 to 'auto' in 1.4. Set the value of `n_init` explicitly
to suppress the warning
    warnings.warn(
array([0, 0, 0, ..., 1, 1, 1], dtype=int32)

time: 745 ms (started: 2024-03-19 18:31:51 +00:00)

# K-Means Clustering Model Evaluation [K=2]
# -----

sscore_km_2cluster = sscore(df_ppd_subset, km_2cluster_model)
dbscore_km_2cluster = dbscore(df_ppd_subset, km_2cluster_model);
%memit
print(f"Davies-Bouldin Index for 2 clusters: {dbscore_km_2cluster}")
print(f"Silhouette Score for 2 clusters: {sscore_km_2cluster}")

peak memory: 486.98 MiB, increment: 0.07 MiB
Davies-Bouldin Index for 2 clusters: 0.8948706111369534
Silhouette Score for 2 clusters: 0.456533778535065
time: 53 s (started: 2024-03-19 18:31:59 +00:00)

list(scaled_df.columns)

['tract_to_msamd_income',
 'population',
 'minority_population',
 'number_of_owner_occupied_units',
 'number_of_1_to_4_family_units',
 'loan_amount_000s',
 'applicant_income_000s',
 'sequence_number',
 'census_tract_number',
 'S.no',
 'hud_median_family_income',
 'as_of_year',
 'application_date_indicator']

time: 4.39 ms (started: 2024-03-19 18:32:51 +00:00)

# Joining cluster labels with the original dataset
df_with_clusters = df_ppd_subset.copy()
df_with_clusters['Cluster_Label'] = km_2cluster_model

```

```

# Extracting non-categorical variables
non_cat_variables = ['tract_to_msamd_income',
    'population',
    'minority_population',
    'number_of_owner_occupied_units',
    'number_of_1_to_4_family_units',
    'loan_amount_000s',
    'applicant_income_000s',
    'sequence_number',
    'census_tract_number',
    'hud_median_family_income',
    'as_of_year',
    'application_date_indicator']

# Grouping variables by cluster label
cluster_groups = df_with_clusters.groupby('Cluster_Label')

# Perform ANOVA for each non-categorical variable
anova_results_non_cat = {}
for column in non_cat_variables:
    anova_results_non_cat[column] = f_oneway(*[group[column] for name,
group in cluster_groups])

# Print ANOVA results for non-categorical variables
for column, result in anova_results_non_cat.items():
    print(f"Variable: {column}")
    print(f"F-value: {result.statistic}")
    print(f"P-value: {result.pvalue}")
    print()

Variable: tract_to_msamd_income
F-value: 0.4899305755891777
P-value: 0.4839609885487106

Variable: population
F-value: 0.36860240058865373
P-value: 0.5437693642942767

Variable: minority_population
F-value: 53.51127032687608
P-value: 2.603129016437262e-13

Variable: number_of_owner_occupied_units
F-value: 79.20665991523715
P-value: 5.745919129964301e-19

Variable: number_of_1_to_4_family_units
F-value: 1447.4507108159823
P-value: 5.59822836275e-313

```

```
Variable: loan_amount_000s
F-value: 66.8511172102048
P-value: 2.9847080460709284e-16

Variable: applicant_income_000s
F-value: 74.4120196977382
P-value: 6.492036740328373e-18

Variable: sequence_number
F-value: 208.91601920080888
P-value: 2.843906136761009e-47

Variable: census_tract_number
F-value: 4396.470592814403
P-value: 0.0

Variable: hud_median_family_income
F-value: 7886.14098144448
P-value: 0.0

Variable: as_of_year
F-value: nan
P-value: nan

Variable: application_date_indicator
F-value: 789.3554434554261
P-value: 1.468872701421905e-172

time: 136 ms (started: 2024-03-19 18:33:23 +00:00)

/usr/local/lib/python3.10/dist-packages/scipy/stats/_stats_py.py:4167:
ConstantInputWarning: Each of the input arrays is constant; the F
statistic is not defined or infinite
    warnings.warn(stats.ConstantInputWarning(msg))

list(encoded_data_categorical.columns)

['S.no',
 'state_name',
 'state_abbr',
 'respondent_id',
 'purchaser_type_name',
 'property_type_name',
 'preapproval_name',
 'owner_occupancy_name',
 'msamd_name',
 'loan_type_name',
 'loan_purpose_name',
 'lien_status_name',
 'hoepa_status_name',
 'county_name',
```

```

'co_applicant_sex_name',
'co_applicant_ethnicity_name',
'applicant_sex_name',
'applicant_ethnicity_name',
'agency_name',
'agency_abbr',
'action_taken_name']

time: 3.65 ms (started: 2024-03-19 18:33:29 +00:00)

from scipy.stats import chi2_contingency

# Extracting categorical variables
cat_variables = ['state_name',
    'state_abbr',
    'respondent_id',
    'purchaser_type_name',
    'property_type_name',
    'preapproval_name',
    'owner_occupancy_name',
    'msamd_name',
    'loan_type_name',
    'loan_purpose_name',
    'lien_status_name',
    'hoepa_status_name',
    'county_name',
    'co_applicant_sex_name',
    'co_applicant_ethnicity_name',
    'applicant_sex_name',
    'applicant_ethnicity_name',
    'agency_name',
    'agency_abbr',
    'action_taken_name']

# Perform Chi-square test for each categorical variable
chi2_results_cat = {}
for column in cat_variables:
    contingency_table = pd.crosstab(df_with_clusters[column],
    km_2cluster_model)
    chi2, p_value, _, _ = chi2_contingency(contingency_table)
    chi2_results_cat[column] = {'Chi-square': chi2, 'P-value': p_value}

# Print Chi-square test results for categorical variables
for column, result in chi2_results_cat.items():
    print(f"Variable: {column}")
    print(f"Chi-square: {result['Chi-square']}")
    print(f"P-value: {result['P-value']}")
    print()

```

```
Variable: state_name
Chi-square: 0.0
P-value: 1.0

Variable: state_abbr
Chi-square: 0.0
P-value: 1.0

Variable: respondent_id
Chi-square: 8854.178279021577
P-value: 0.0

Variable: purchaser_type_name
Chi-square: 498.35007306683207
P-value: 1.296870862833355e-101

Variable: property_type_name
Chi-square: 198.0655253931781
P-value: 9.786279905934946e-44

Variable: preapproval_name
Chi-square: 52.32061858505
P-value: 4.352332130547724e-12

Variable: owner_occupancy_name
Chi-square: 68.33880256812851
P-value: 1.446830223606883e-15

Variable: msamd_name
Chi-square: 10170.372989395928
P-value: 0.0

Variable: loan_type_name
Chi-square: 418.8003280951333
P-value: 1.8734457444207908e-90

Variable: loan_purpose_name
Chi-square: 118.33618340196239
P-value: 2.0119828269018097e-26

Variable: lien_status_name
Chi-square: 786.5195973265372
P-value: 3.6289872720942084e-170

Variable: hoepa_status_name
Chi-square: 0.33947432487372337
P-value: 0.5601328088115435

Variable: county_name
Chi-square: 18426.294723230283
P-value: 0.0
```

```
Variable: co_applicant_sex_name
Chi-square: 642.3465051284185
P-value: 1.0575920026155875e-137

Variable: co_applicant_ethnicity_name
Chi-square: 819.0900740581233
P-value: 5.626147759343297e-176

Variable: applicant_sex_name
Chi-square: 388.1480742650455
P-value: 8.170901376968501e-84

Variable: applicant_ethnicity_name
Chi-square: 702.3318873173232
P-value: 6.5523447107087534e-152

Variable: agency_name
Chi-square: 3432.0041383942594
P-value: 0.0

Variable: agency_abbr
Chi-square: 3432.0041383942594
P-value: 0.0

Variable: action_taken_name
Chi-square: 4599.781160569204
P-value: 0.0

time: 284 ms (started: 2024-03-19 18:33:47 +00:00)

# Get the cluster centers
cluster_centers = km_2cluster.cluster_centers_

# Convert cluster_centers to a DataFrame
centroids_df = pd.DataFrame(cluster_centers,
columns=df_ppd_subset.columns[:-1])

# Display the centroids of clusters
centroids_df

{"type": "dataframe", "variable_name": "centroids_df"}

time: 31.4 ms (started: 2024-03-19 18:34:39 +00:00)

df_ppd_subset['Cluster_Label'] = km_2cluster_model

# Calculate cluster-wise descriptive statistics
cluster_stats = df_ppd_subset.groupby('Cluster_Label').describe()

# Variables of interest
```

```

variables_of_interest = ['tract_to_msamd_income',
'minority_population',
'number_of_owner_occupied_units',
'number_of_1_to_4_family_units',
'loan_amount_000s',
'applicant_income_000s',
'sequence_number',
'census_tract_number',
'hud_median_family_income',
'application_date_indicator']

# Print descriptive statistics for each variable
for variable in variables_of_interest:
    print(f"Descriptive statistics for variable: {variable}")
    print(cluster_stats[variable])
    print()

Descriptive statistics for variable: tract_to_msamd_income
      count        mean         std       min       25%       50%
75% \
Cluster_Label
0          28768.0  0.385254  0.113063  0.0  0.312929  0.378090
0.452343
1          31232.0  0.384590  0.118690  0.0  0.307047  0.374964
0.445679

      max
Cluster_Label
0          1.0
1          1.0

Descriptive statistics for variable: minority_population
      count        mean         std       min       25%
50% \
Cluster_Label
0          28768.0  0.223790  0.165685  0.030296  0.107278
0.176927
1          31232.0  0.233068  0.144882  0.000000  0.130135
0.196119

      75%   max
Cluster_Label
0          0.292183  1.0
1          0.302749  1.0

Descriptive statistics for variable: number_of_owner_occupied_units
      count        mean         std       min       25%
50% \

```

```

Cluster_Label

0          28768.0  0.470702  0.178937  0.007378  0.343729
0.457411
1          31232.0  0.458082  0.168391  0.000000  0.340376
0.447686

    75%   max
Cluster_Label
0          0.585178  1.0
1          0.566734  1.0

Descriptive statistics for variable: number_of_1_to_4_family_units
      count      mean       std     min      25%      50%
75% \
Cluster_Label

0          28768.0  0.334855  0.143088  0.0  0.242925  0.311626
0.397716
1          31232.0  0.296217  0.103985  0.0  0.231504  0.288442
0.360552

      max
Cluster_Label
0          1.0
1          1.0

Descriptive statistics for variable: loan_amount_000s
      count      mean       std     min      25%      50%
75% \
Cluster_Label

0          28768.0  0.004897  0.012141  0.0  0.002800  0.004091
0.005691
1          31232.0  0.005632  0.009818  0.0  0.003346  0.004618
0.006491

      max
Cluster_Label
0          1.000000
1          0.887871

Descriptive statistics for variable: applicant_income_000s
      count      mean       std     min      25%
50% \
Cluster_Label

0          28768.0  0.017461  0.016306  0.000000  0.009740
0.015179
1          31232.0  0.018790  0.020913  0.000162  0.010714

```

```

0.016071

    75%      max
Cluster_Label
0          0.019156  0.635877
1          0.021104  1.000000

Descriptive statistics for variable: sequence_number
      count      mean      std   min      25%      50%
75% \
Cluster_Label

0          28768.0  0.069881  0.130025  0.0  0.003121  0.015698
0.063753
1          31232.0  0.055587  0.112085  0.0  0.002261  0.011350
0.048852

      max
Cluster_Label
0          0.999993
1          1.000000

Descriptive statistics for variable: census_tract_number
      count      mean      std   min      25%      50%
75% \
Cluster_Label

0          28768.0  0.272995  0.414842  0.0  0.010873  0.041934
0.943214
1          31232.0  0.093065  0.230925  0.0  0.011900  0.033021
0.043563

      max
Cluster_Label
0          1.0
1          1.0

Descriptive statistics for variable: hud_median_family_income
      count      mean      std   min      25%
50% \
Cluster_Label

0          28768.0  69324.359294  13639.895277  48700.0  55600.0
69900.0
1          31232.0  78055.888192  10331.847326  48700.0  73300.0
73600.0

      75%      max
Cluster_Label
0          78100.0  90300.0

```

```

1          90300.0  90300.0

Descriptive statistics for variable: application_date_indicator
   count      mean       std    min   25%   50%   75%   max
Cluster_Label
0        28768.0  0.052697  0.320346  0.0   0.0   0.0   0.0   2.0
1        31232.0  0.001153  0.048001  0.0   0.0   0.0   0.0   2.0

time: 300 ms (started: 2024-03-19 18:35:11 +00:00)

# Assuming km_2cluster_model contains the cluster labels

# Grouping the data by cluster labels
cluster_groups = df_with_clusters.groupby(km_2cluster_model)

# Counting the number of variables in each cluster
num_variables_in_clusters = cluster_groups.size()

# Displaying the results
print("Number of Variables in Each Cluster:")
print(num_variables_in_clusters)

Number of Variables in Each Cluster:
0      28768
1      31232
dtype: int64
time: 7.71 ms (started: 2024-03-19 18:35:21 +00:00)

# Joining cluster labels with the original dataset
df_with_clusters = df_ppd_subset.copy()
df_with_clusters['Cluster_Label'] = km_2cluster_model

# Extracting non-categorical variables
non_cat_variables = ['tract_to_msamd_income',
 'population',
 'minority_population',
 'number_of_owner_occupied_units',
 'number_of_1_to_4_family_units',
 'loan_amount_000s',
 'applicant_income_000s',
 'sequence_number',
 'census_tract_number',
 'hud_median_family_income',
 'as_of_year',
 'application_date_indicator']

# Grouping variables by cluster label
cluster_groups = df_with_clusters.groupby('Cluster_Label')

# Perform ANOVA for each non-categorical variable
anova_results_non_cat = {}

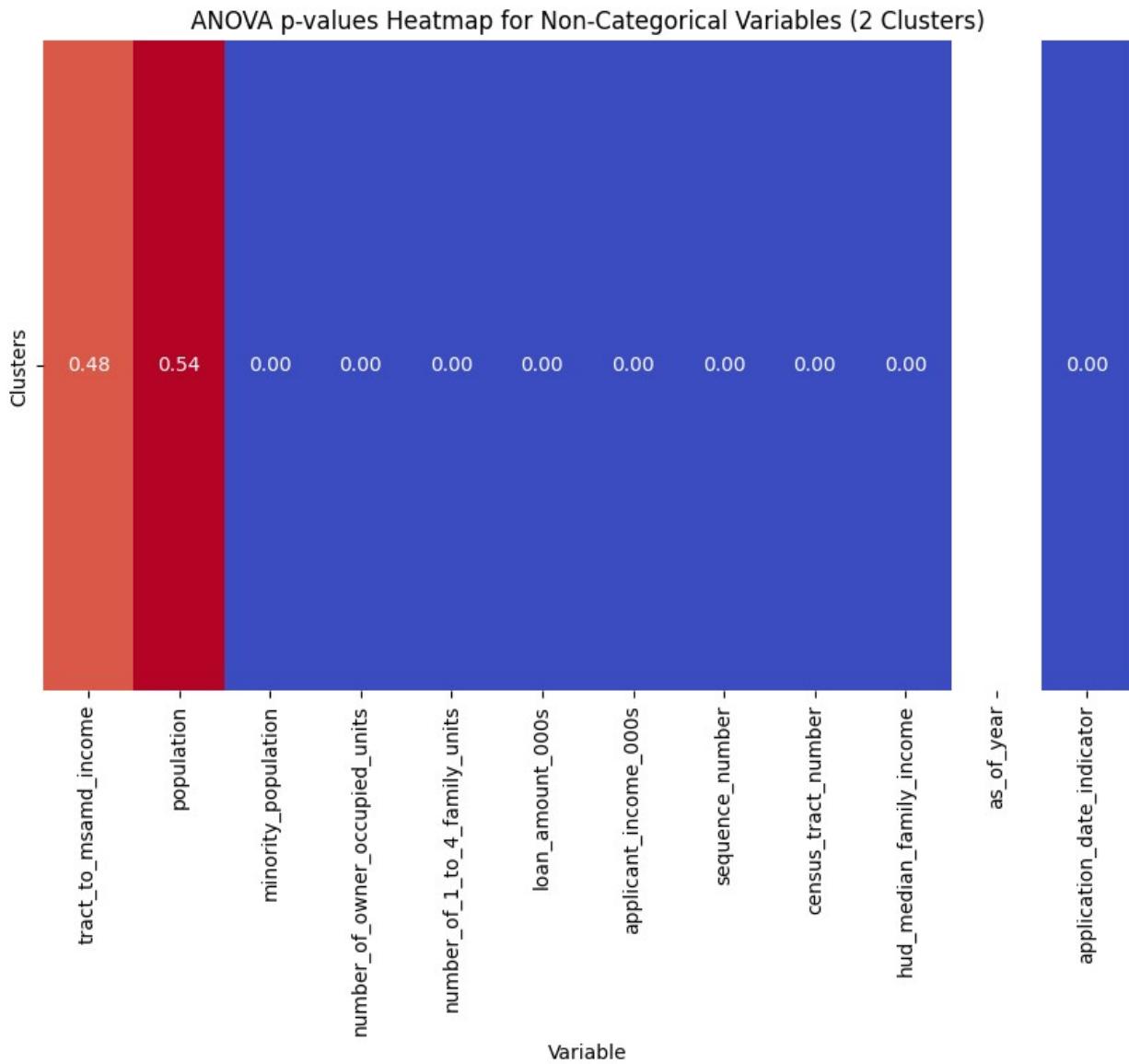
```

```
for column in non_cat_variables:
    anova_results_non_cat[column] = f_oneway(*[group[column] for name,
group in cluster_groups])

# Extract p-values for non-categorical variables
non_cat_p_values = [result.pvalue for result in
anova_results_non_cat.values()]

# Plotting heatmap for non-categorical variables
plt.figure(figsize=(10, 6))
sns.heatmap([non_cat_p_values], cmap='coolwarm', annot=True,
fmt='.2f', xticklabels=non_cat_variables, yticklabels=['Clusters'],
cbar=False)
plt.xlabel('Variable')
plt.title('ANOVA p-values Heatmap for Non-Categorical Variables (2
Clusters)')
plt.show()

/usr/local/lib/python3.10/dist-packages/scipy/stats/_stats_py.py:4167:
ConstantInputWarning: Each of the input arrays is constant; the F
statistic is not defined or infinite
warnings.warn(stats.ConstantInputWarning(msg))
```



```
time: 966 ms (started: 2024-03-19 18:36:04 +00:00)
```

```
cat_variables = ['state_name',
'state_abbr',
'respondent_id',
'purchaser_type_name',
'property_type_name',
'preapproval_name',
'owner_occupancy_name',
'msamd_name',
'loan_type_name',
'loan_purpose_name',
'lien_status_name',
'hoepa_status_name',
```

```

'county_name',
'co_applicant_sex_name',
'co_applicant_ethnicity_name',
'applicant_sex_name',
'applicant_ethnicity_name',
'agency_name',
'agency_abbr',
'action_taken_name']

# Initialize dictionary to store chi-square results
chi2_results_cat = {}

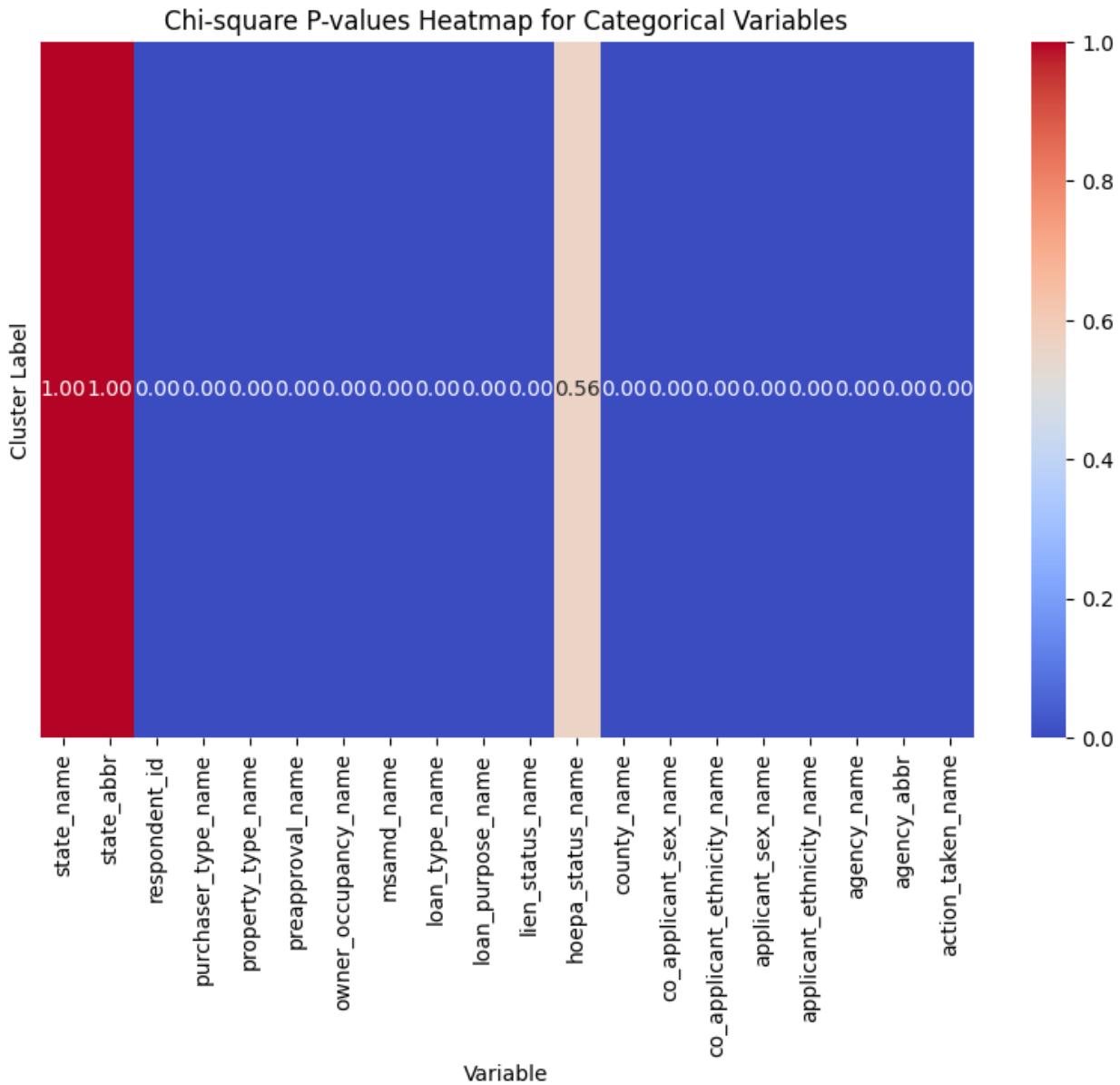
# Calculate chi-square for each categorical variable
for column in cat_variabels:
    contingency_table = pd.crosstab(df_with_clusters[column],
km_2cluster_model)
    chi2, p_value, _, _ = chi2_contingency(contingency_table)
    chi2_results_cat[column] = {'P-value': p_value}

# Extract p-values
p_values = [[result['P-value']] for result in
chi2_results_cat.values()]

# Get variable names
variables = list(chi2_results_cat.keys())

# Plotting heatmap for p-values
plt.figure(figsize=(10, 6))
sns.heatmap(p_values, cmap='coolwarm', annot=True, fmt='.2f',
xticklabels=variables, yticklabels=False)
plt.xlabel('Variable')
plt.ylabel('Cluster Label')
plt.title('Chi-square P-values Heatmap for Categorical Variables')
plt.show()

```



```
time: 1.15 s (started: 2024-03-19 18:36:10 +00:00)

# Assign cluster labels to the DataFrame
df_ppd_subset['Cluster_Label'] = km_2cluster_model

# Plot the scatter plot with clusters and centroids
plt.figure(figsize=(10, 6))
sns.scatterplot(data=df_ppd_subset, x='tract_to_msamd_income',
y='property_type_name', hue='Cluster_Label', legend='full')

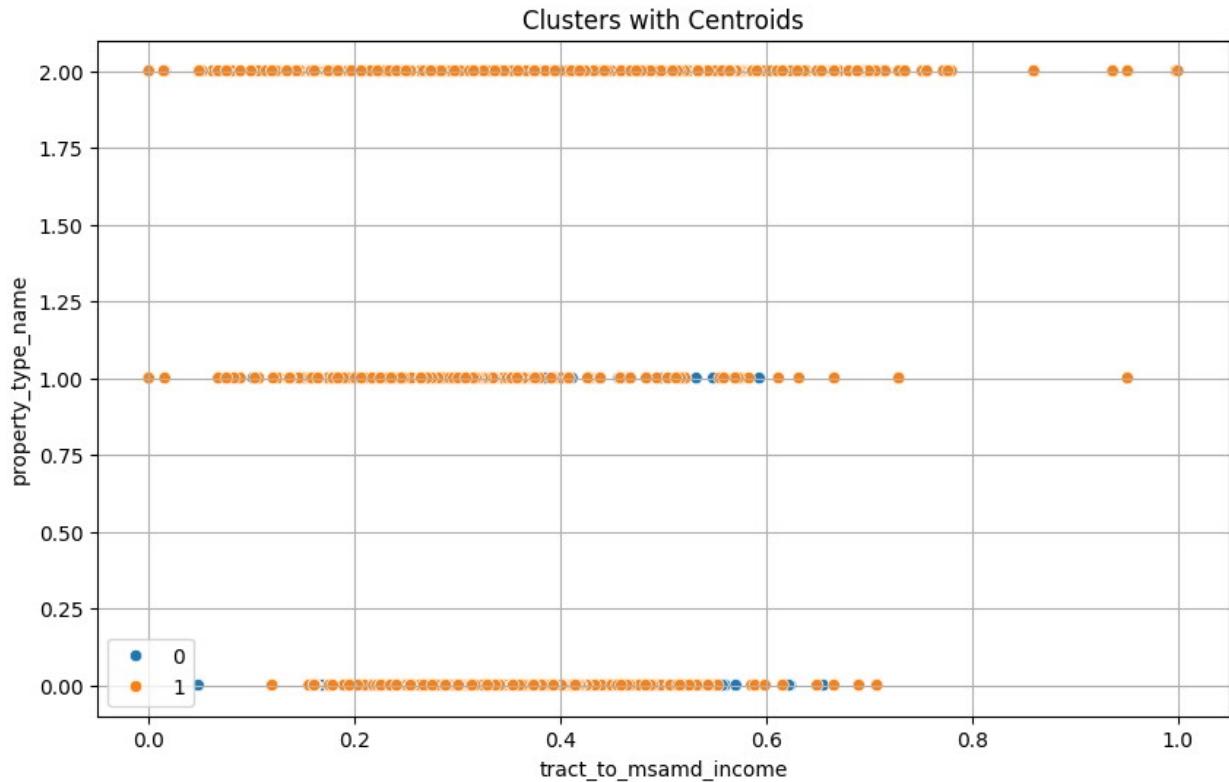
# Set plot title and labels
plt.title('Clusters with Centroids')
plt.xlabel('tract_to_msamd_income')
```

```

plt.ylabel('property_type_name')
plt.legend()
plt.grid(True)
%memit
plt.show()

peak memory: 487.18 MiB, increment: 0.00 MiB

```



```

time: 3.79 s (started: 2024-03-19 18:36:14 +00:00)

import plotly.graph_objs as go

# Create 3D scatter plot
fig = go.Figure()

# Add traces for each cluster
for cluster_label in df_ppd_subset['Cluster_Label'].unique():
    cluster_data = df_ppd_subset[df_ppd_subset['Cluster_Label'] == cluster_label]
    fig.add_trace(go.Scatter3d(
        x=cluster_data['tract_to_msamd_income'],
        y=cluster_data['loan_amount_000s'],
        z=cluster_data['hud_median_family_income'],
        mode='markers',
        marker=dict(
            size=10,
            color=cluster_label
        )
    ))

```

```

        size=5,
        color=cluster_label,
        colorscale='Viridis', # Adjust colorscale if needed
        opacity=0.8
    ),
    name=f'Cluster {cluster_label}'
))

# Set layout
fig.update_layout(
    title='Clusters with Centroids (3D)',
    scene=dict(
        xaxis=dict(title='tract_to_msamdi_income'),
        yaxis=dict(title='loan_amount_000s'),
        zaxis=dict(title='hud_median_family_income'),
    ),
    margin=dict(l=0, r=0, b=0, t=40)
)

# Show interactive 3D scatter plot
fig.show()

time: 436 ms (started: 2024-03-19 18:36:21 +00:00)

# Create subplots for each cluster
fig, axes = plt.subplots(1, 2, figsize=(20, 4), sharex=True,
sharey=True)
fig.suptitle('Clusters with Centroids')

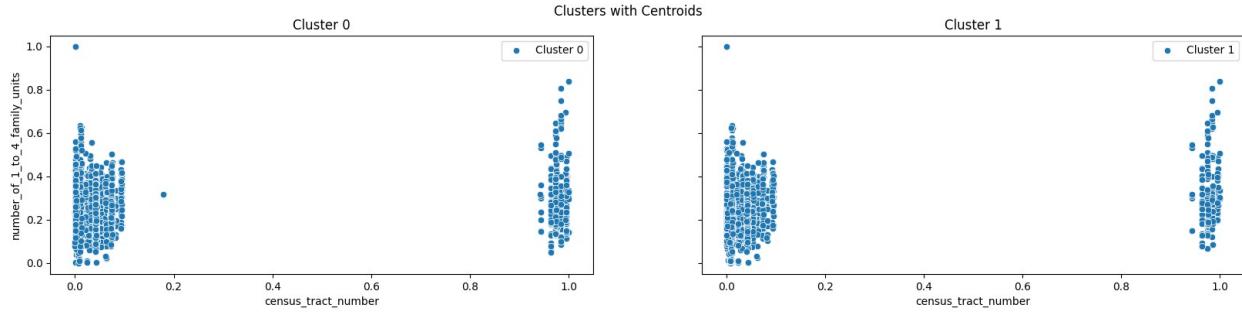
# Iterate through each cluster
for i in range(2):
    # Filter data points belonging to the current cluster
    cluster_data = df_ppd_subset[df_ppd_subset['Cluster_Label'] == i]

    # Scatter plot of data points
    sns.scatterplot(data=cluster_data, x='census_tract_number',
y='number_of_1_to_4_family_units', ax=axes[i], label=f'Cluster {i}')

    # Set title and labels for each subplot
    axes[i].set_title(f'Cluster {i}')
    axes[i].set_xlabel('census_tract_number')
    axes[i].set_ylabel('number_of_1_to_4_family_units')
    axes[i].legend()

plt.show()
%memit

```



```
peak memory: 488.56 MiB, increment: 0.00 MiB
time: 1.21 s (started: 2024-03-19 18:39:52 +00:00)
```

```
##### KKKKKKKKKK ====== 333333333333
```

```
time: 321 µs (started: 2024-03-19 18:39:54 +00:00)
```

```
# Create K-Means Clusters [K=3]
```

```
km_3cluster = kmclust(n_clusters=3, init='random', random_state=333)
df_ppd_subset['Cluster_Label'] =
km_3cluster.fit_predict(df_ppd_subset)
km_3cluster_model = df_ppd_subset['Cluster_Label'].values
km_3cluster_model
```

```
/usr/local/lib/python3.10/dist-packages/sklearn/cluster/
_kmeans.py:870: FutureWarning:
```

```
The default value of `n_init` will change from 10 to 'auto' in 1.4.
Set the value of `n_init` explicitly to suppress the warning
```

```
array([0, 0, 0, ..., 1, 1, 1], dtype=int32)
```

```
time: 2.11 s (started: 2024-03-19 18:39:55 +00:00)
```

```
# K-Means Clustering Model Evaluation [K=3]
```

```
# -----
```

```
sscore_km_3cluster = sscore(df_ppd_subset, km_3cluster_model)
dbscore_km_3cluster = dbscore(df_ppd_subset, km_3cluster_model);
%memit
print(f"Davies-Bouldin Index for 3 clusters: {dbscore_km_3cluster}")
print(f"Silhouette Score for 3 clusters: {sscore_km_3cluster}")
```

```
peak memory: 486.86 MiB, increment: 0.00 MiB
Davies-Bouldin Index for 3 clusters: 0.8837455599506786
Silhouette Score for 3 clusters: 0.4312014782480616
time: 53.5 s (started: 2024-03-19 18:39:58 +00:00)
```

```
# Joining cluster labels with the original dataset
```

```
df_with_clusters = df_ppd_subset.copy()
df_with_clusters['Cluster_Label'] = km_3cluster_model
```

```

# Extracting non-categorical variables
non_cat_variables = ['tract_to_msamd_income',
    'population',
    'minority_population',
    'number_of_owner_occupied_units',
    'number_of_1_to_4_family_units',
    'loan_amount_000s',
    'applicant_income_000s',
    'sequence_number',
    'census_tract_number',
    'hud_median_family_income',
    'as_of_year',
    'application_date_indicator']

# Grouping variables by cluster label
cluster_groups = df_with_clusters.groupby('Cluster_Label')

# Perform ANOVA for each non-categorical variable
anova_results_non_cat = {}
for column in non_cat_variables:
    anova_results_non_cat[column] = f_oneway(*[group[column] for name,
group in cluster_groups])

# Print ANOVA results for non-categorical variables
for column, result in anova_results_non_cat.items():
    print(f"Variable: {column}")
    print(f"F-value: {result.statistic}")
    print(f"P-value: {result.pvalue}")
    print()

Variable: tract_to_msamd_income
F-value: 12.703560383256072
P-value: 3.0484681673235144e-06

Variable: population
F-value: 29.2373499686067
P-value: 2.034997816482065e-13

Variable: minority_population
F-value: 1174.0980790615254
P-value: 0.0

Variable: number_of_owner_occupied_units
F-value: 228.53409104816646
P-value: 1.3336833013916693e-99

Variable: number_of_1_to_4_family_units
F-value: 2429.572179442173
P-value: 0.0

```

```
Variable: loan_amount_000s
F-value: 231.0255989132163
P-value: 1.1252001181830643e-100

Variable: applicant_income_000s
F-value: 238.57253758632996
P-value: 6.296373939673157e-104

Variable: sequence_number
F-value: 76.67894384932364
P-value: 5.511197191663442e-34

Variable: census_tract_number
F-value: 6504.82913373703
P-value: 0.0

Variable: hud_median_family_income
F-value: 36672.72153507024
P-value: 0.0

Variable: as_of_year
F-value: nan
P-value: nan

Variable: application_date_indicator
F-value: 797.7152901599333
P-value: 0.0

time: 420 ms (started: 2024-03-19 18:40:51 +00:00)
/usr/local/lib/python3.10/dist-packages/scipy/stats/_stats_py.py:4167:
ConstantInputWarning:

Each of the input arrays is constant; the F statistic is not defined or
infinite

from scipy.stats import chi2_contingency

# Extracting categorical variables
cat_variables = ['state_name',
                 'state_abbr',
                 'respondent_id',
                 'purchaser_type_name',
                 'property_type_name',
                 'preapproval_name',
                 'owner_occupancy_name',
                 'msamd_name',
                 'loan_type_name',
                 'loan_purpose_name',
```

```

'lien_status_name',
'hoepa_status_name',
'county_name',
'co_applicant_sex_name',
'co_applicant_ethnicity_name',
'applicant_sex_name',
'applicant_ethnicity_name',
'agency_name',
'agency_abbr',
'action_taken_name']

# Perform Chi-square test for each categorical variable
chi2_results_cat = {}
for column in cat_variables:
    contingency_table = pd.crosstab(df_with_clusters[column],
km_3cluster_model)
    chi2, p_value, _, _ = chi2_contingency(contingency_table)
    chi2_results_cat[column] = {'Chi-square': chi2, 'P-value': p_value}

# Print Chi-square test results for categorical variables
for column, result in chi2_results_cat.items():
    print(f"Variable: {column}")
    print(f"Chi-square: {result['Chi-square']}") 
    print(f"P-value: {result['P-value']}") 
    print()

```

Variable: state_name
Chi-square: 0.0
P-value: 1.0

Variable: state_abbr
Chi-square: 0.0
P-value: 1.0

Variable: respondent_id
Chi-square: 12446.017333693797
P-value: 0.0

Variable: purchaser_type_name
Chi-square: 1398.5495057897674
P-value: 2.9206912601191758e-286

Variable: property_type_name
Chi-square: 633.4754279952812
P-value: 8.802932399978627e-136

Variable: preapproval_name
Chi-square: 281.8678870502726
P-value: 8.81560826993418e-60

Variable: owner_occupancy_name
Chi-square: 186.68366050051145
P-value: 2.7344027928516775e-39

Variable: msamd_name
Chi-square: 22453.668788654228
P-value: 0.0

Variable: loan_type_name
Chi-square: 1212.769039544149
P-value: 8.249882649431631e-259

Variable: loan_purpose_name
Chi-square: 281.5374598610049
P-value: 1.038718377177558e-59

Variable: lien_status_name
Chi-square: 1642.7507167875815
P-value: 0.0

Variable: hoepa_status_name
Chi-square: 0.7738482595695989
P-value: 0.6791426196687216

Variable: county_name
Chi-square: 50095.626895059424
P-value: 0.0

Variable: co_applicant_sex_name
Chi-square: 1075.5082905841657
P-value: 7.453619130583919e-227

Variable: co_applicant_ethnicity_name
Chi-square: 1221.2836019707781
P-value: 2.415344728509488e-258

Variable: applicant_sex_name
Chi-square: 509.45934888722127
P-value: 7.706317691417523e-107

Variable: applicant_ethnicity_name
Chi-square: 819.3616950940085
P-value: 1.008904472942779e-173

Variable: agency_name
Chi-square: 1045.9689304422648
P-value: 2.332408160942132e-218

Variable: agency_abbr

```
Chi-square: 1045.9689304422648
P-value: 2.332408160942132e-218

Variable: action_taken_name
Chi-square: 4823.794971357527
P-value: 0.0

time: 887 ms (started: 2024-03-19 18:40:52 +00:00)

# Get the cluster centers
cluster_centers = km_3cluster.cluster_centers_

# Convert cluster_centers to a DataFrame
centroids_df = pd.DataFrame(cluster_centers,
columns=df_ppd_subset.columns)

# Display the centroids of clusters
centroids_df

{"type": "dataframe", "variable_name": "centroids_df"}

time: 55 ms (started: 2024-03-19 18:40:52 +00:00)

# Calculate cluster-wise descriptive statistics
cluster_stats = df_ppd_subset.groupby('Cluster_Label').describe()

# Print the descriptive statistics for each cluster
print("Cluster-wise Descriptive Statistics:")
cluster_stats

Cluster-wise Descriptive Statistics:

{"type": "dataframe", "variable_name": "cluster_stats"}

time: 730 ms (started: 2024-03-19 18:40:53 +00:00)

# Grouping the data by cluster labels
cluster_groups = df_with_clusters.groupby(km_3cluster_model)

# Counting the number of variables in each cluster
num_variables_in_clusters = cluster_groups.size()

# Displaying the results
print("Number of Variables in Each Cluster:")
print(num_variables_in_clusters)

Number of Variables in Each Cluster:
0    18138
1    23766
2    18096
dtype: int64
time: 8.84 ms (started: 2024-03-19 18:40:53 +00:00)
```

```

# Joining cluster labels with the original dataset
df_with_clusters = df_ppd_subset.copy()
df_with_clusters['Cluster_Label'] = km_3cluster_model

# Extracting non-categorical variables
non_cat_variables = ['tract_to_msamd_income',
    'population',
    'minority_population',
    'number_of_owner_occupied_units',
    'number_of_1_to_4_family_units',
    'loan_amount_000s',
    'applicant_income_000s',
    'sequence_number',
    'census_tract_number',
    'hud_median_family_income',
    'as_of_year',
    'application_date_indicator']

# Grouping variables by cluster label
cluster_groups = df_with_clusters.groupby('Cluster_Label')

# Perform ANOVA for each non-categorical variable
anova_results_non_cat = {}
for column in non_cat_variables:
    anova_results_non_cat[column] = f_oneway(*[group[column] for name, group in cluster_groups])

# Extract p-values for non-categorical variables
non_cat_p_values = [result.pvalue for result in anova_results_non_cat.values()]

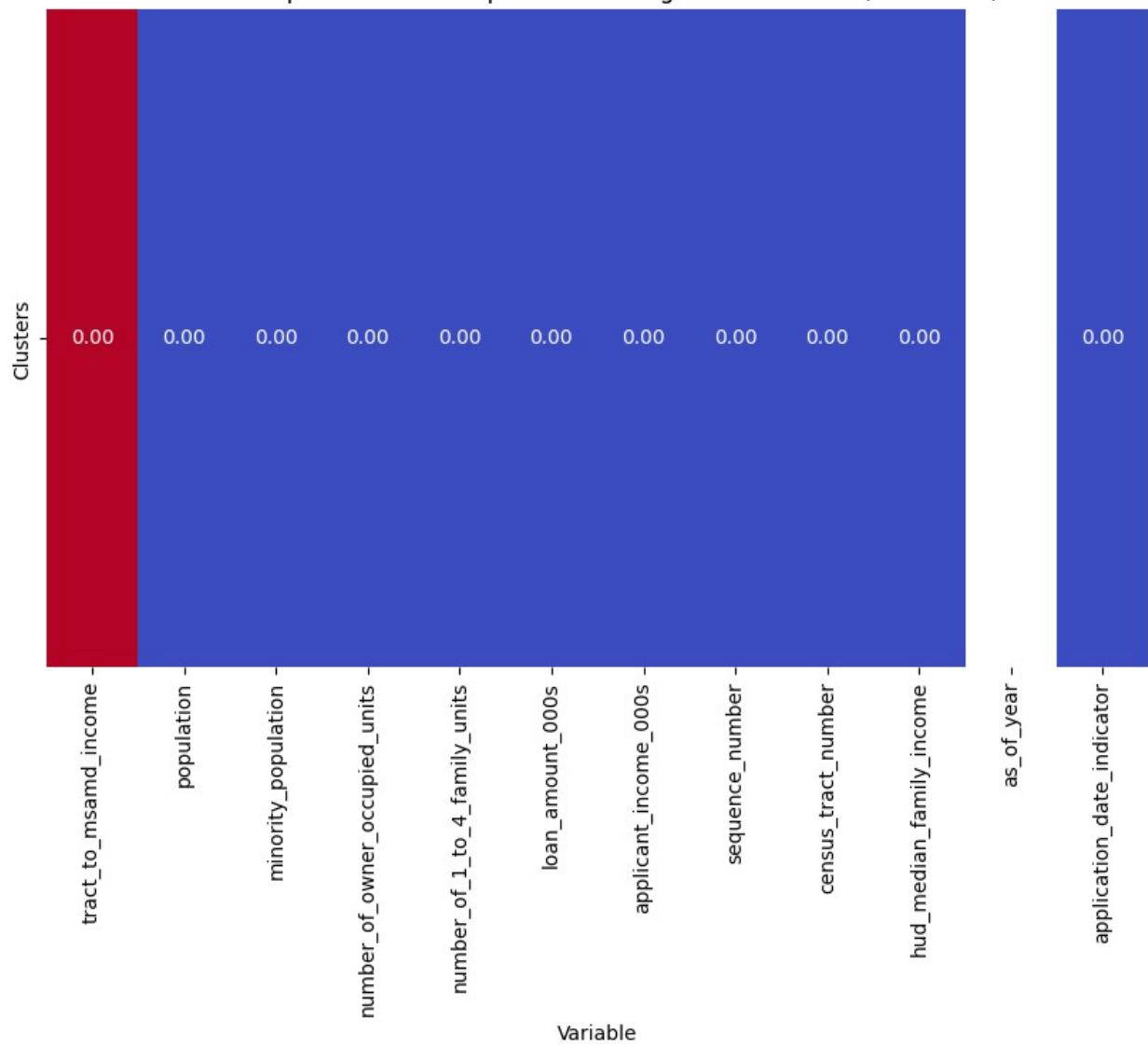
# Plotting heatmap for non-categorical variables
plt.figure(figsize=(10, 6))
sns.heatmap([non_cat_p_values], cmap='coolwarm', annot=True,
fmt='.2f', xticklabels=non_cat_variables, yticklabels=['Clusters'],
cbar=False)
plt.xlabel('Variable')
plt.title('ANOVA p-values Heatmap for Non-Categorical Variables (3 Clusters)')
plt.show()

/usr/local/lib/python3.10/dist-packages/scipy/stats/_stats_py.py:4167:
ConstantInputWarning:

Each of the input arrays is constant; the F statistic is not defined or infinite

```

ANOVA p-values Heatmap for Non-Categorical Variables (3 Clusters)



```
time: 744 ms (started: 2024-03-19 18:40:53 +00:00)
```

```
cat_variables = ['state_name',
 'state_abbr',
 'respondent_id',
 'purchaser_type_name',
 'property_type_name',
 'preapproval_name',
 'owner_occupancy_name',
 'msamd_name',
 'loan_type_name',
 'loan_purpose_name',
 'lien_status_name',
 'hoepa_status_name',
```

```

'county_name',
'co_applicant_sex_name',
'co_applicant_ethnicity_name',
'applicant_sex_name',
'applicant_ethnicity_name',
'agency_name',
'agency_abbr',
'action_taken_name']

# Initialize dictionary to store chi-square results
chi2_results_cat = {}

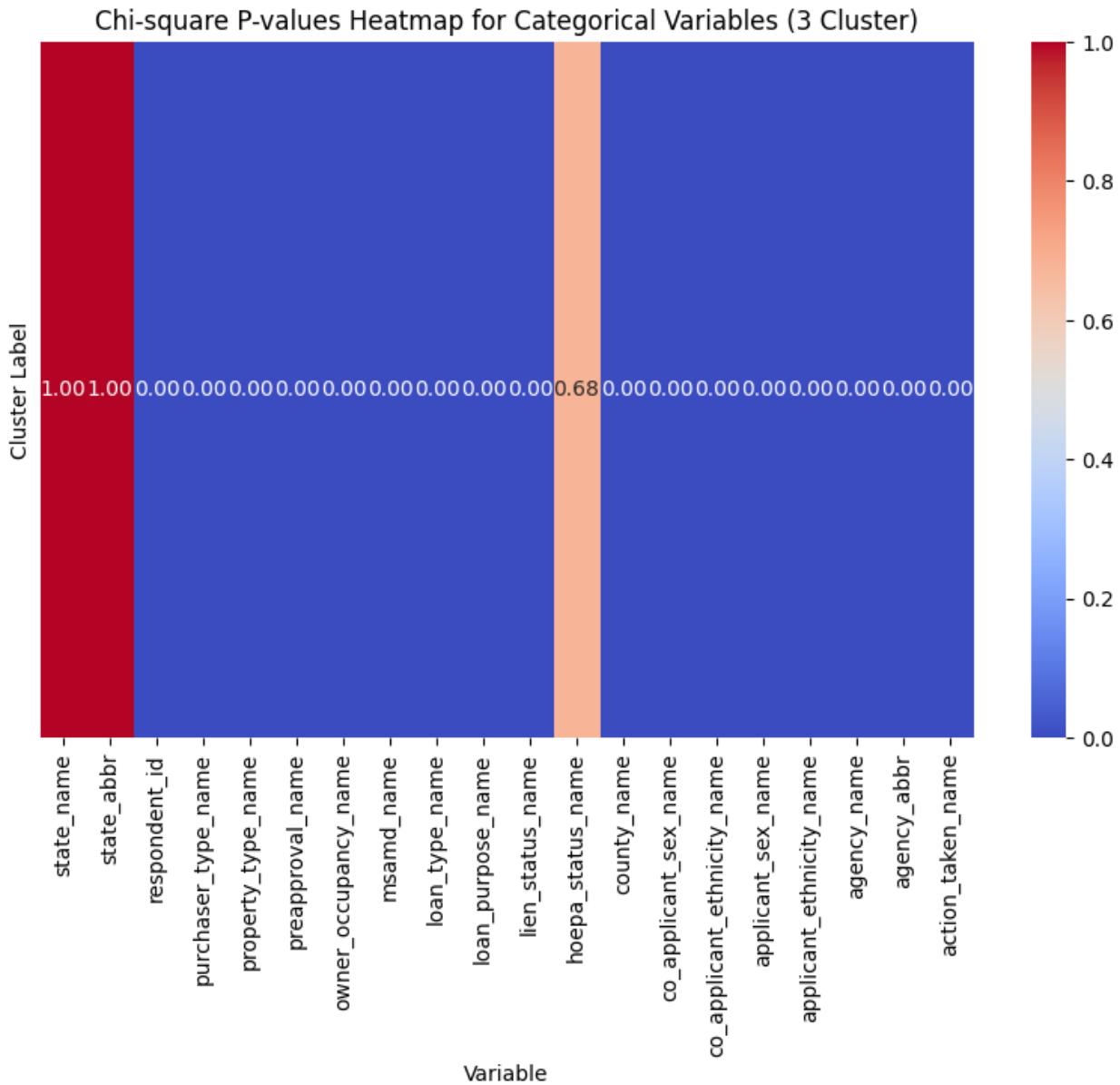
# Calculate chi-square for each categorical variable
for column in cat_variabels:
    contingency_table = pd.crosstab(df_with_clusters[column],
km_3cluster_model)
    chi2, p_value, _, _ = chi2_contingency(contingency_table)
    chi2_results_cat[column] = {'P-value': p_value}

# Extract p-values
p_values = [[result['P-value']] for result in
chi2_results_cat.values()]

# Get variable names
variables = list(chi2_results_cat.keys())

# Plotting heatmap for p-values
plt.figure(figsize=(10, 6))
sns.heatmap(p_values, cmap='coolwarm', annot=True, fmt='.2f',
xticklabels=variables, yticklabels=False)
plt.xlabel('Variable')
plt.ylabel('Cluster Label')
plt.title('Chi-square P-values Heatmap for Categorical Variables (3
Cluster)')
plt.show()

```



```
time: 1.37 s (started: 2024-03-19 18:40:54 +00:00)

# Assign cluster labels to the DataFrame
df_ppd_subset['Cluster_Label'] = km_3cluster_model

# Plot the scatter plot with clusters and centroids
plt.figure(figsize=(10, 6))
sns.scatterplot(data=df_ppd_subset, x='number_of_1_to_4_family_units',
y='hud_median_family_income', hue='Cluster_Label', legend='full')

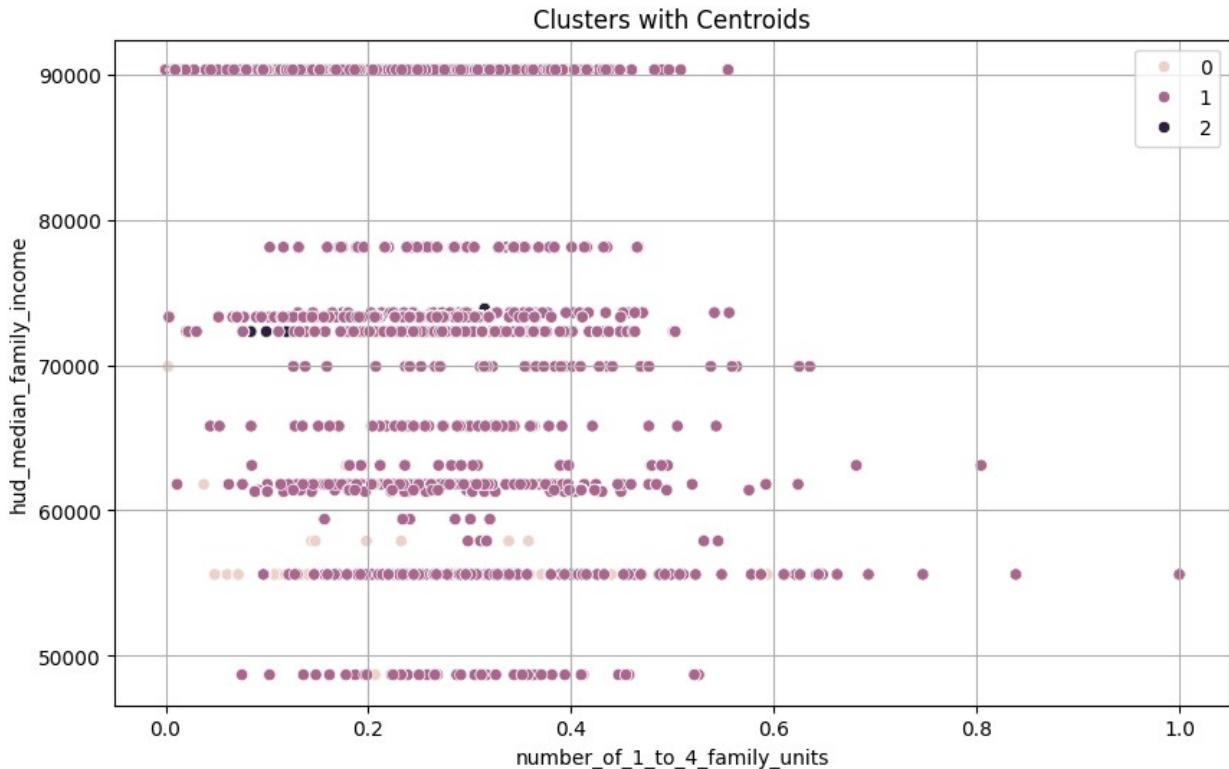
# Set plot title and labels
plt.title('Clusters with Centroids')
plt.xlabel('number_of_1_to_4_family_units')
```

```

plt.ylabel('hud_median_family_income')
plt.legend()
plt.grid(True)
%memit
plt.show()

peak memory: 512.51 MiB, increment: 0.00 MiB

```



```

time: 8.4 s (started: 2024-03-19 18:40:56 +00:00)

import plotly.graph_objs as go
import numpy as np

# Create 3D scatter plot
fig = go.Figure()

# Get unique cluster labels
unique_clusters = np.unique(km_3cluster_model)

# Add traces for each cluster
for cluster_label in unique_clusters:
    cluster_data = df_ppd_subset[km_3cluster_model == cluster_label]
    fig.add_trace(go.Scatter3d(
        x=cluster_data['hud_median_family_income'],
        y=cluster_data['applicant_income_000s'],
        z=cluster_data['minority_population'],

```

```

        mode='markers',
        marker=dict(
            size=5,
            color=cluster_label,
            colorscale='Viridis', # Adjust colorscale if needed
            opacity=0.8
        ),
        name=f'Cluster {cluster_label}'
    ))

# Set layout
fig.update_layout(
    title='Clusters with Centroids (3D)',
    scene=dict(
        xaxis=dict(title='hud_median_family_income'),
        yaxis=dict(title='applicant_income_000s'),
        zaxis=dict(title='minority_population'),
    ),
    margin=dict(l=0, r=0, b=0, t=40)
)

# Show interactive 3D scatter plot
fig.show()

time: 184 ms (started: 2024-03-19 18:41:04 +00:00)

# Create subplots for each cluster
fig, axes = plt.subplots(1, 3, figsize=(20, 4), sharex=True,
sharey=True)
fig.suptitle('Clusters with Centroids')

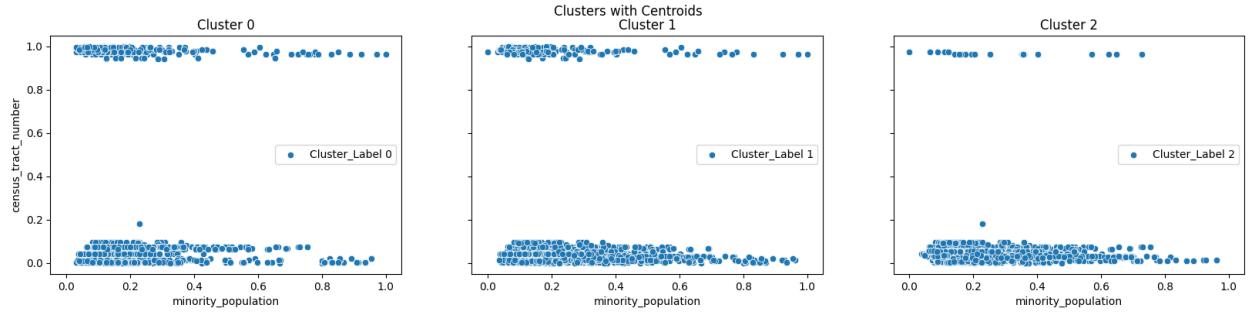
# Iterate through each cluster
for i in range(3):
    # Filter data points belonging to the current cluster
    cluster_data = df_ppd_subset[df_ppd_subset['Cluster_Label'] == i]

    # Scatter plot of data points
    sns.scatterplot(data=cluster_data, x='minority_population',
y='census_tract_number', ax=axes[i], label=f'Cluster_Label {i}')

    # Set title and labels for each subplot
    axes[i].set_title(f'Cluster {i}')
    axes[i].set_xlabel('minority_population')
    axes[i].set_ylabel('census_tract_number')
    axes[i].legend()

plt.show()
%memit

```



```
peak memory: 490.31 MiB, increment: 0.00 MiB
time: 2.37 s (started: 2024-03-19 18:41:04 +00:00)
```

```
# K=4
```

```
time: 393 µs (started: 2024-03-19 18:41:07 +00:00)
```

```
# Create K-Means Clusters [K=4]
km_4cluster = kmclus(n_clusters=4, init='random', random_state=333)
df_ppd_subset['Cluster_Label'] =
km_4cluster.fit_predict(df_ppd_subset)
km_4cluster_model = df_ppd_subset['Cluster_Label'].values
km_4cluster_model
```

```
/usr/local/lib/python3.10/dist-packages/sklearn/cluster/
_kmeans.py:870: FutureWarning:
```

```
The default value of `n_init` will change from 10 to 'auto' in 1.4.
Set the value of `n_init` explicitly to suppress the warning
```

```
array([3, 0, 3, ..., 1, 1, 1], dtype=int32)
```

```
time: 1.63 s (started: 2024-03-19 18:41:07 +00:00)
```

```
#####
KKKKKKKKKKKKKKKK
=====
4444444444444444
```

```
time: 362 µs (started: 2024-03-19 18:41:08 +00:00)
```

```
# K-Means Clustering Model Evaluation [K=4]
```

```
# -----
```

```
sscore_km_4cluster = sscore(df_ppd_subset, km_4cluster_model)
dbscore_km_4cluster = dbscore(df_ppd_subset, km_4cluster_model);
%memit
print(f"Davies-Bouldin Index for 4 clusters: {dbscore_km_4cluster}")
print(f"Silhouette Score for 4 clusters: {sscore_km_4cluster}")
```

```
peak memory: 491.44 MiB, increment: 0.00 MiB
```

```
Davies-Bouldin Index for 4 clusters: 0.871140449105725
```

```

Silhouette Score for 4 clusters: 0.44321886396320287
time: 49.4 s (started: 2024-03-19 18:41:08 +00:00)

# Joining cluster labels with the original dataset
df_with_clusters = df_ppd_subset.copy()
df_with_clusters['Cluster_Label'] = km_4cluster_model

# Extracting non-categorical variables
non_cat_variables = ['tract_to_msamd_income',
    'population',
    'minority_population',
    'number_of_owner_occupied_units',
    'number_of_1_to_4_family_units',
    'loan_amount_000s',
    'applicant_income_000s',
    'sequence_number',
    'census_tract_number',
    'hud_median_family_income',
    'as_of_year',
    'application_date_indicator']

# Grouping variables by cluster label
cluster_groups = df_with_clusters.groupby('Cluster_Label')

# Perform ANOVA for each non-categorical variable
anova_results_non_cat = {}
for column in non_cat_variables:
    anova_results_non_cat[column] = f_oneway(*[group[column] for name, group in cluster_groups])

# Print ANOVA results for non-categorical variables
for column, result in anova_results_non_cat.items():
    print(f"Variable: {column}")
    print(f"F-value: {result.statistic}")
    print(f"P-value: {result.pvalue}")
    print()

Variable: tract_to_msamd_income
F-value: 51.50331637112778
P-value: 3.0929224414263875e-33

Variable: population
F-value: 3.477308618427996
P-value: 0.015235602431356474

Variable: minority_population
F-value: 1700.228610153462
P-value: 0.0

Variable: number_of_owner_occupied_units

```

```
F-value: 139.13153059321243
P-value: 7.753377593427997e-90

Variable: number_of_1_to_4_family_units
F-value: 1878.2704837944307
P-value: 0.0

Variable: loan_amount_000s
F-value: 304.41370121181023
P-value: 3.6795449694673084e-196

Variable: applicant_income_000s
F-value: 513.6916989836737
P-value: 0.0

Variable: sequence_number
F-value: 179.56443922479943
P-value: 6.497425078355844e-116

Variable: census_tract_number
F-value: 7348.202916930307
P-value: 0.0

Variable: hud_median_family_income
F-value: 76468.96718372905
P-value: 0.0

Variable: as_of_year
F-value: nan
P-value: nan

Variable: application_date_indicator
F-value: 803.6973675994069
P-value: 0.0

time: 136 ms (started: 2024-03-19 18:41:58 +00:00)

/usr/local/lib/python3.10/dist-packages/scipy/stats/_stats_py.py:4167:
ConstantInputWarning:

Each of the input arrays is constant; the F statistic is not defined or
infinite

from scipy.stats import chi2_contingency

# Extracting categorical variables
cat_variables = ['state_name',
                 'state_abbr',
                 'respondent_id',
                 'purchaser_type_name',
```

```

'property_type_name',
'preapproval_name',
'owner_occupancy_name',
'msamd_name',
'loan_type_name',
'loan_purpose_name',
'lien_status_name',
'hoepa_status_name',
'county_name',
'co_applicant_sex_name',
'co_applicant_ethnicity_name',
'applicant_sex_name',
'applicant_ethnicity_name',
'agency_name',
'agency_abbr',
'action_taken_name']

# Perform Chi-square test for each categorical variable
chi2_results_cat = {}
for column in cat_variables:
    contingency_table = pd.crosstab(df_with_clusters[column],
km_4cluster_model)
    chi2, p_value, _, _ = chi2_contingency(contingency_table)
    chi2_results_cat[column] = {'Chi-square': chi2, 'P-value': p_value}

# Print Chi-square test results for categorical variables
for column, result in chi2_results_cat.items():
    print(f"Variable: {column}")
    print(f"Chi-square: {result['Chi-square']}")  

    print(f"P-value: {result['P-value']}")  

    print()
  

Variable: state_name
Chi-square: 0.0
P-value: 1.0
  

Variable: state_abbr
Chi-square: 0.0
P-value: 1.0
  

Variable: respondent_id
Chi-square: 32172.648842023515
P-value: 0.0
  

Variable: purchaser_type_name
Chi-square: 2898.260429675344
P-value: 0.0
  

Variable: property_type_name

```

Chi-square: 1040.919915658263
P-value: 1.2604505005053296e-221

Variable: preapproval_name
Chi-square: 205.038222928886
P-value: 1.6053406466766635e-41

Variable: owner_occupancy_name
Chi-square: 291.33769227256494
P-value: 5.867996971928788e-60

Variable: msamd_name
Chi-square: 58401.514737340054
P-value: 0.0

Variable: loan_type_name
Chi-square: 2591.6086431011804
P-value: 0.0

Variable: loan_purpose_name
Chi-square: 322.03678555239844
P-value: 1.5442432805374943e-66

Variable: lien_status_name
Chi-square: 2379.808079550548
P-value: 0.0

Variable: hoepa_status_name
Chi-square: 5.829280631440727
P-value: 0.1202180355419859

Variable: county_name
Chi-square: 96731.04367338086
P-value: 0.0

Variable: co_applicant_sex_name
Chi-square: 2055.1818571467297
P-value: 0.0

Variable: co_applicant_ethnicity_name
Chi-square: 2378.097969036845
P-value: 0.0

Variable: applicant_sex_name
Chi-square: 1281.3555668483068
P-value: 3.2896058159064975e-270

Variable: applicant_ethnicity_name
Chi-square: 1954.6346108272146
P-value: 0.0

```
Variable: agency_name
Chi-square: 10609.899635043503
P-value: 0.0

Variable: agency_abbr
Chi-square: 10609.899635043503
P-value: 0.0

Variable: action_taken_name
Chi-square: 7450.482100068883
P-value: 0.0

time: 303 ms (started: 2024-03-19 18:41:58 +00:00)
cluster_centers = km_4cluster.cluster_centers_

# Convert centroids to DataFrame for better visualization
centroids_df = pd.DataFrame(cluster_centers,
columns=df_ppd_subset.columns)
%memit
print("Centroids of Clusters:")
centroids_df

peak memory: 520.61 MiB, increment: 0.00 MiB
Centroids of Clusters:

{"type":"dataframe","variable_name":"centroids_df"}

time: 278 ms (started: 2024-03-19 18:41:58 +00:00)

# Calculate cluster-wise descriptive statistics
cluster_stats = df_ppd_subset.groupby('Cluster_Label').describe()

# Print the descriptive statistics for each cluster
print("Cluster-wise Descriptive Statistics:")
cluster_stats

Cluster-wise Descriptive Statistics:

{"type":"dataframe","variable_name":"cluster_stats"}

time: 358 ms (started: 2024-03-19 18:41:58 +00:00)

# Grouping the data by cluster labels
cluster_groups = df_with_clusters.groupby(km_4cluster_model)

# Counting the number of variables in each cluster
num_variables_in_clusters = cluster_groups.size()

# Displaying the results
```

```

print("Number of Variables in Each Cluster:")
print(num_variables_in_clusters)

Number of Variables in Each Cluster:
0    13654
1    20027
2    11451
3    14868
dtype: int64
time: 6.06 ms (started: 2024-03-19 18:41:59 +00:00)

# Joining cluster labels with the original dataset
df_with_clusters = df_ppd_subset.copy()
df_with_clusters['Cluster_Label'] = km_4cluster_model

# Extracting non-categorical variables
non_cat_variables = ['tract_to_msamd_income',
 'population',
 'minority_population',
 'number_of_owner_occupied_units',
 'number_of_1_to_4_family_units',
 'loan_amount_000s',
 'applicant_income_000s',
 'sequence_number',
 'census_tract_number',
 'hud_median_family_income',
 'as_of_year',
 'application_date_indicator']

# Grouping variables by cluster label
cluster_groups = df_with_clusters.groupby('Cluster_Label')

# Perform ANOVA for each non-categorical variable
anova_results_non_cat = {}
for column in non_cat_variables:
    anova_results_non_cat[column] = f_oneway(*[group[column] for name,
group in cluster_groups])

# Extract p-values for non-categorical variables
non_cat_p_values = [result.pvalue for result in
anova_results_non_cat.values()]

# Plotting heatmap for non-categorical variables
plt.figure(figsize=(10, 6))
sns.heatmap([non_cat_p_values], cmap='coolwarm', annot=True,
fmt='.2f', xticklabels=non_cat_variables, yticklabels=['Clusters'],
cbar=False)
plt.xlabel('Variable')
plt.title('ANOVA p-values Heatmap for Non-Categorical Variables (4

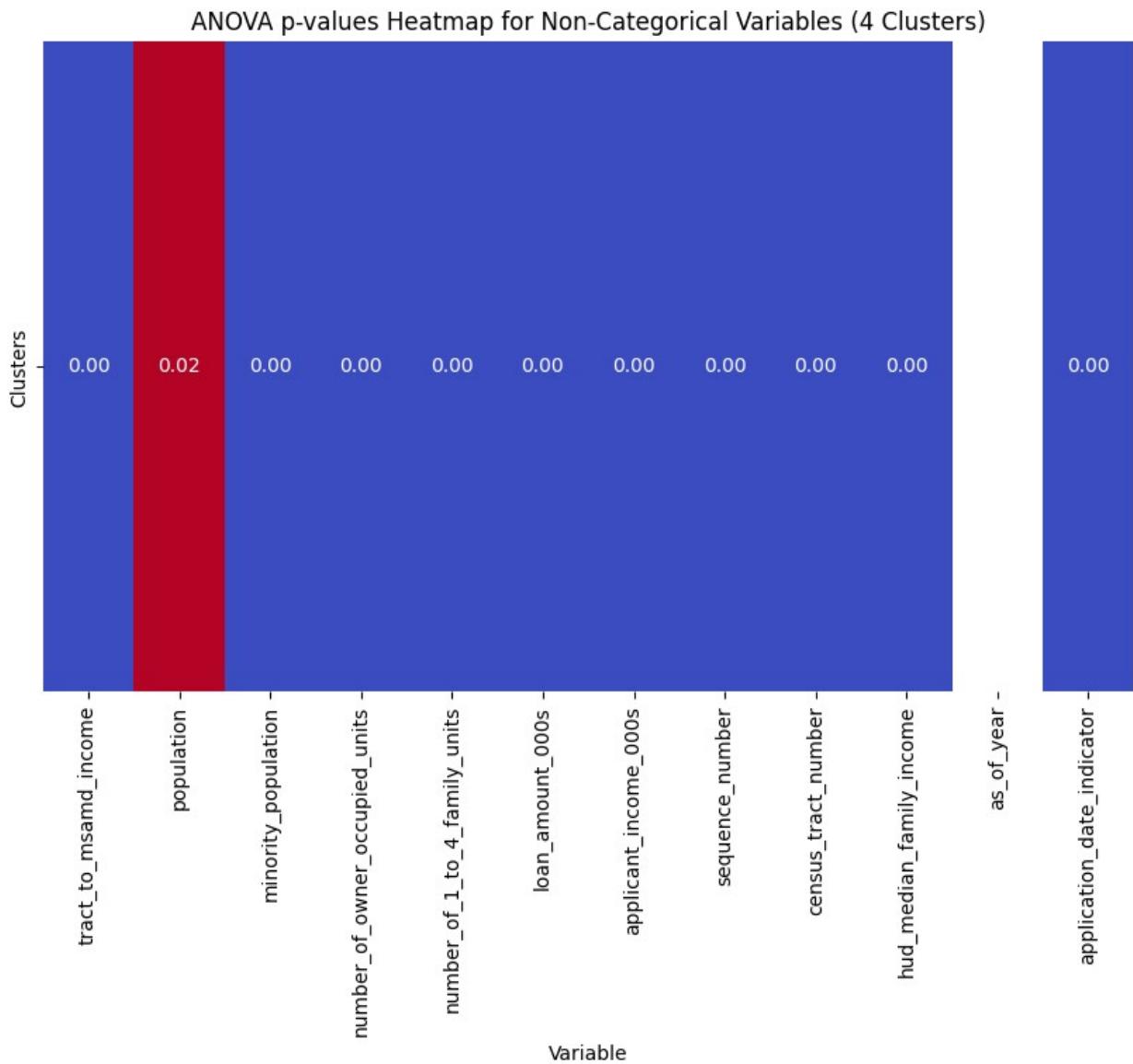
```

```

Clusters)
plt.show()

/usr/local/lib/python3.10/dist-packages/scipy/stats/_stats_py.py:4167:
ConstantInputWarning:
Each of the input arrays is constant;the F statistic is not defined or
infinite

```



```

time: 398 ms (started: 2024-03-19 18:41:59 +00:00)

cat_variables = ['state_name',
 'state_abbr',

```

```

'respondent_id',
'purchaser_type_name',
'property_type_name',
'preapproval_name',
'owner_occupancy_name',
'msamd_name',
'loan_type_name',
'loan_purpose_name',
'lien_status_name',
'hoepa_status_name',
'county_name',
'co_applicant_sex_name',
'co_applicant_ethnicity_name',
'applicant_sex_name',
'applicant_ethnicity_name',
'agency_name',
'agency_abbr',
'action_taken_name']

# Initialize dictionary to store chi-square results
chi2_results_cat = {}

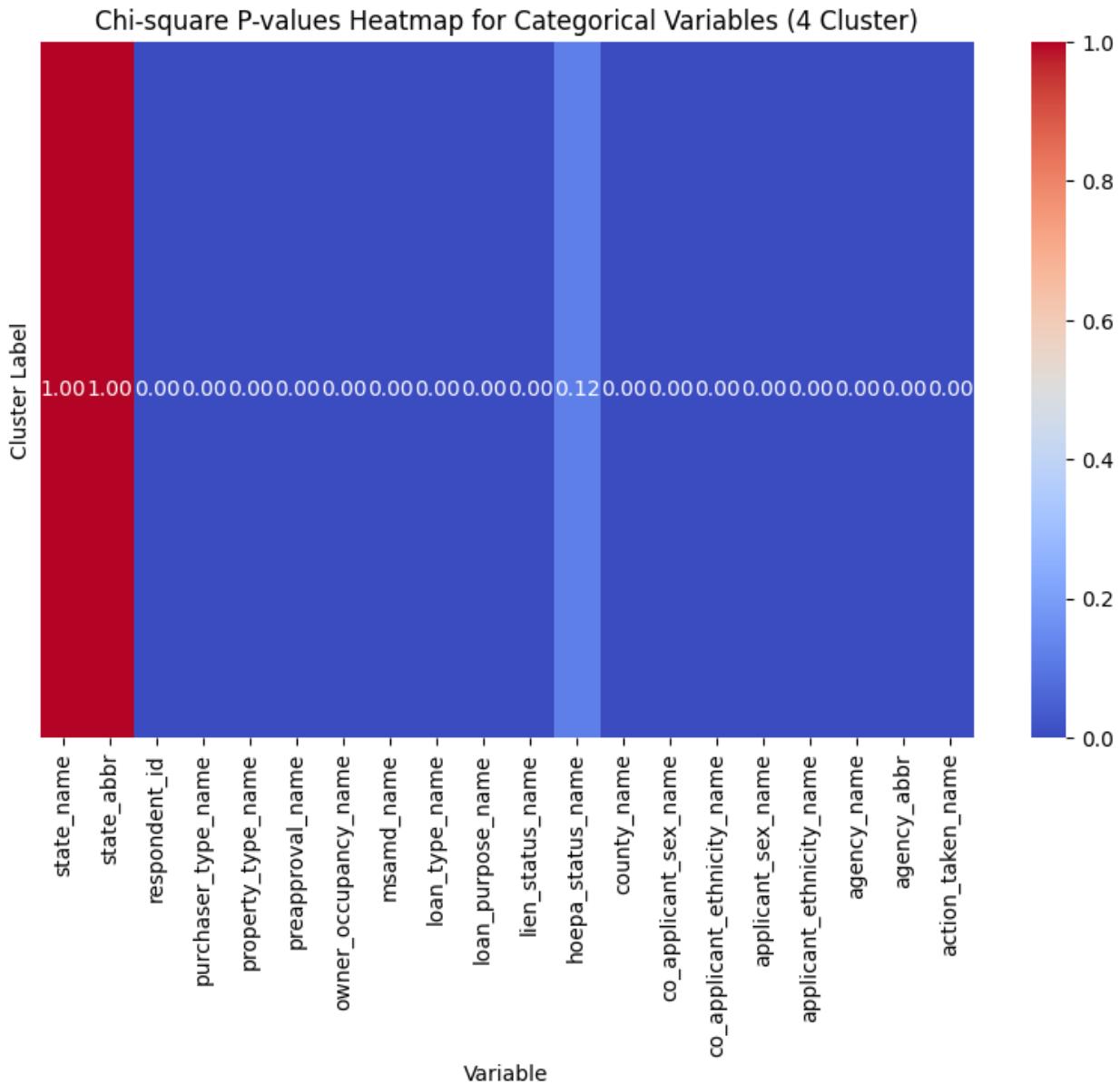
# Calculate chi-square for each categorical variable
for column in cat_variables:
    contingency_table = pd.crosstab(df_with_clusters[column],
    km_4cluster_model)
    chi2, p_value, _, _ = chi2_contingency(contingency_table)
    chi2_results_cat[column] = {'P-value': p_value}

# Extract p-values
p_values = [[result['P-value']] for result in
chi2_results_cat.values()]

# Get variable names
variables = list(chi2_results_cat.keys())

# Plotting heatmap for p-values
plt.figure(figsize=(10, 6))
sns.heatmap(p_values, cmap='coolwarm', annot=True, fmt='.2f',
xticklabels=variables, yticklabels=False)
plt.xlabel('Variable')
plt.ylabel('Cluster Label')
plt.title('Chi-square P-values Heatmap for Categorical Variables (4
Cluster)')
plt.show()

```



```
time: 685 ms (started: 2024-03-19 18:41:59 +00:00)
```

```
import plotly.graph_objs as go
import numpy as np

# Create 3D scatter plot
fig = go.Figure()

# Get unique cluster labels
unique_clusters = np.unique(km_4cluster_model)

# Add traces for each cluster
for cluster_label in unique_clusters:
```

```

cluster_data = df_ppd_subset[km_4cluster_model == cluster_label]
fig.add_trace(go.Scatter3d(
    x=cluster_data['loan_amount_000s'],
    y=cluster_data['hud_median_family_income'],
    z=cluster_data['number_of_owner_occupied_units'],
    mode='markers',
    marker=dict(
        size=5,
        color=cluster_label,
        colorscale='Viridis', # Adjust colorscale if needed
        opacity=0.8
    ),
    name=f'Cluster {cluster_label}'
))

# Set layout
fig.update_layout(
    title='Clusters with Centroids (3D)',
    scene=dict(
        xaxis=dict(title='loan_amount_000s'),
        yaxis=dict(title='hud_median_family_income'),
        zaxis=dict(title='number_of_owner_occupied_units'),
    ),
    margin=dict(l=0, r=0, b=0, t=40)
)

# Show interactive 3D scatter plot
fig.show()

time: 176 ms (started: 2024-03-19 18:42:00 +00:00)

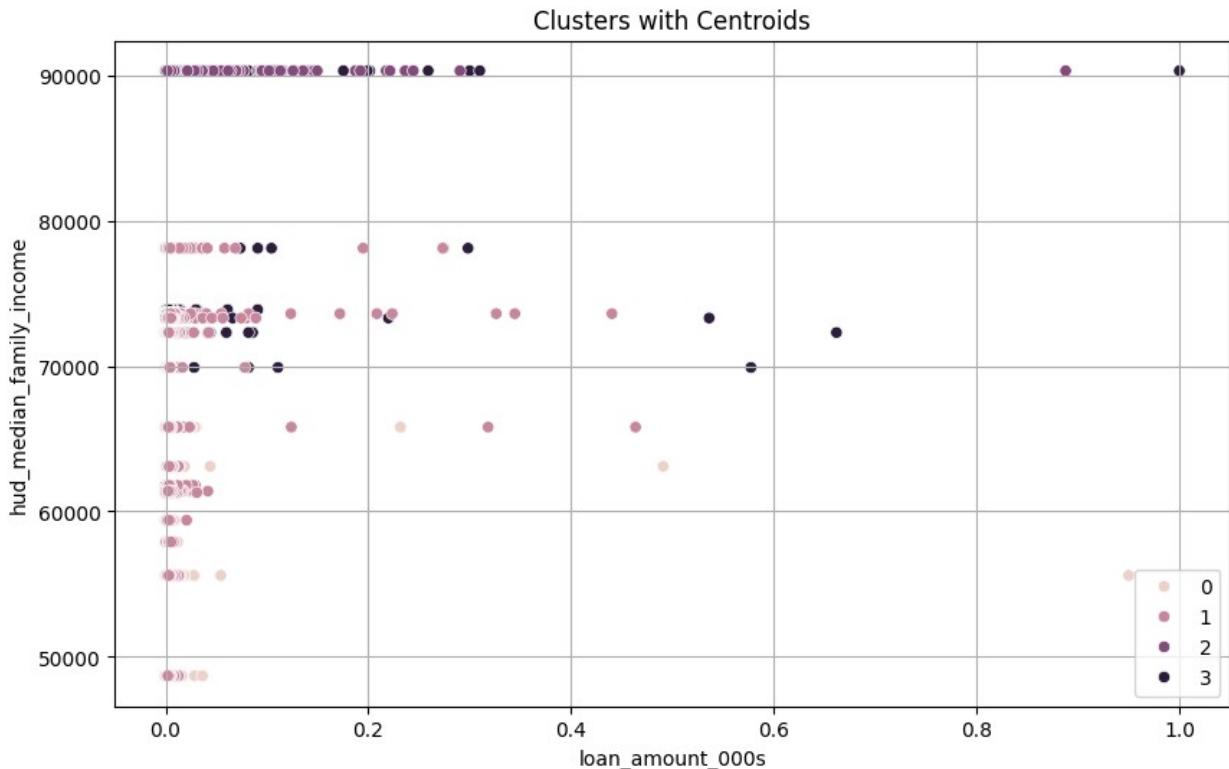
# Assign cluster labels to the DataFrame
df_ppd_subset['Cluster_Label'] = km_4cluster_model

# Plot the scatter plot with clusters and centroids
plt.figure(figsize=(10, 6))
sns.scatterplot(data=df_ppd_subset, x='loan_amount_000s',
y='hud_median_family_income', hue='Cluster_Label', legend='full')

# Set plot title and labels
plt.title('Clusters with Centroids')
plt.xlabel('loan_amount_000s')
plt.ylabel('hud_median_family_income')
plt.legend()
plt.grid(True)
%memit
plt.show()

peak memory: 488.68 MiB, increment: 0.00 MiB

```



```

time: 4.8 s (started: 2024-03-19 18:42:00 +00:00)

# Create subplots for each cluster
fig, axes = plt.subplots(1, 4, figsize=(20, 4), sharex=True,
sharey=True)
fig.suptitle('Clusters with Centroids')

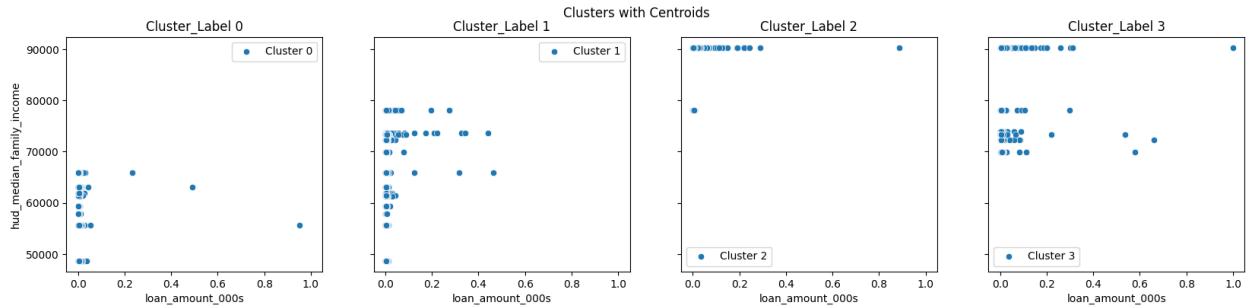
# Iterate through each cluster
for i in range(4):
    # Filter data points belonging to the current cluster
    cluster_data = df_ppd_subset[df_ppd_subset['Cluster_Label'] == i]

    # Scatter plot of data points
    sns.scatterplot(data=cluster_data, x='loan_amount_000s',
y='hud_median_family_income', ax=axes[i], label=f'Cluster {i}')

    # Set title and labels for each subplot
    axes[i].set_title(f'Cluster_Label {i}')
    axes[i].set_xlabel('loan_amount_000s')
    axes[i].set_ylabel('hud_median_family_income')
    axes[i].legend()

plt.show()
%memit

```



```
peak memory: 493.71 MiB, increment: 0.00 MiB
time: 2.08 s (started: 2024-03-19 18:42:05 +00:00)
```

```
#####
KKKKKKKKKK
=====
5555555555555555
```

```
time: 545 µs (started: 2024-03-19 18:42:07 +00:00)
```

```
# Create K-Means Clusters [K=5]
km_5cluster = kmclus(n_clusters=5, init='random', random_state=333)
df_ppd_subset['Cluster_Label'] =
km_5cluster.fit_predict(df_ppd_subset)
km_5cluster_model = df_ppd_subset['Cluster_Label'].values
km_5cluster_model

/usr/local/lib/python3.10/dist-packages/sklearn/cluster/
_kmeans.py:870: FutureWarning:
```

```
The default value of `n_init` will change from 10 to 'auto' in 1.4.
Set the value of `n_init` explicitly to suppress the warning
```

```
array([4, 3, 4, ..., 1, 1, 1], dtype=int32)
```

```
time: 1.57 s (started: 2024-03-19 18:42:07 +00:00)
```

```
# Create subsets for each cluster label
dfk0 = df_ppd_subset[df_ppd_subset['Cluster_Label'] == 0]
dfk1 = df_ppd_subset[df_ppd_subset['Cluster_Label'] == 1]
dfk2 = df_ppd_subset[df_ppd_subset['Cluster_Label'] == 2]
dfk3 = df_ppd_subset[df_ppd_subset['Cluster_Label'] == 3]
dfk4 = df_ppd_subset[df_ppd_subset['Cluster_Label'] == 4]
```

```
# Find the number of rows for each subset
```

```
num_rows_dfk0 = dfk0.shape[0]
num_rows_dfk1 = dfk1.shape[0]
num_rows_dfk2 = dfk2.shape[0]
num_rows_dfk3 = dfk3.shape[0]
num_rows_dfk4 = dfk4.shape[0]
```

```
print("Number of rows in dfk0:", num_rows_dfk0)
```

```

print("Number of rows in dfk1:", num_rows_dfk1)
print("Number of rows in dfk2:", num_rows_dfk2)
print("Number of rows in dfk3:", num_rows_dfk3)
print("Number of rows in dfk4:", num_rows_dfk4)

Number of rows in dfk0: 11435
Number of rows in dfk1: 14049
Number of rows in dfk2: 9611
Number of rows in dfk3: 14162
Number of rows in dfk4: 10743
time: 21.4 ms (started: 2024-03-19 18:42:09 +00:00)

# K-Means Clustering Model Evaluation [K=5]
# -----

sscore_km_5cluster = sscore(df_ppd_subset, km_5cluster_model)
dbscore_km_5cluster = dbscore(df_ppd_subset, km_5cluster_model);
%memit
print(f"Davies-Bouldin Index for 5 clusters: {dbscore_km_5cluster}")
print(f"Silhouette Score for 5 clusters: {sscore_km_5cluster}")

peak memory: 491.28 MiB, increment: 0.00 MiB
Davies-Bouldin Index for 5 clusters: 0.7202139693085351
Silhouette Score for 5 clusters: 0.4575747033407218
time: 50.7 s (started: 2024-03-19 18:42:09 +00:00)

# Joining cluster labels with the original dataset
df_with_clusters = df_ppd_subset.copy()
df_with_clusters['Cluster_Label'] = km_5cluster_model

# Extracting non-categorical variables
non_cat_variables = ['tract_to_msamd_income',
    'population',
    'minority_population',
    'number_of_owner_occupied_units',
    'number_of_1_to_4_family_units',
    'loan_amount_000s',
    'applicant_income_000s',
    'sequence_number',
    'census_tract_number',
    'hud_median_family_income',
    'as_of_year',
    'application_date_indicator']

# Grouping variables by cluster label
cluster_groups = df_with_clusters.groupby('Cluster_Label')

# Perform ANOVA for each non-categorical variable
anova_results_non_cat = {}
for column in non_cat_variables:
    anova_results_non_cat[column] = f_oneway(*[group[column] for name,

```

```
group in cluster_groups])

# Print ANOVA results for non-categorical variables
for column, result in anova_results_non_cat.items():
    print(f"Variable: {column}")
    print(f"F-value: {result.statistic}")
    print(f"P-value: {result.pvalue}")
    print()

Variable: tract_to_msamd_income
F-value: 36.735275944084634
P-value: 1.0048291798338516e-30

Variable: population
F-value: 24.32915120251411
P-value: 3.805312517512312e-20

Variable: minority_population
F-value: 1328.7527886289538
P-value: 0.0

Variable: number_of_owner_occupied_units
F-value: 155.05847521898747
P-value: 3.146355571197373e-132

Variable: number_of_1_to_4_family_units
F-value: 1583.7775363697897
P-value: 0.0

Variable: loan_amount_000s
F-value: 236.4565886026594
P-value: 7.717295085016942e-202

Variable: applicant_income_000s
F-value: 426.3789442097002
P-value: 0.0

Variable: sequence_number
F-value: 135.53709539075538
P-value: 1.7120081163319907e-115

Variable: census_tract_number
F-value: 4666.008424041755
P-value: 0.0

Variable: hud_median_family_income
F-value: 69396.54958000354
P-value: 0.0

Variable: as_of_year
F-value: nan
```

```
P-value: nan

Variable: application_date_indicator
F-value: 552.6604042303275
P-value: 0.0

time: 135 ms (started: 2024-03-19 18:42:59 +00:00)

/usr/local/lib/python3.10/dist-packages/scipy/stats/_stats_py.py:4167:
ConstantInputWarning:

Each of the input arrays is constant; the F statistic is not defined or
infinite

from scipy.stats import chi2_contingency

# Extracting categorical variables
cat_variables = ['state_name',
    'state_abbr',
    'respondent_id',
    'purchaser_type_name',
    'property_type_name',
    'preapproval_name',
    'owner_occupancy_name',
    'msamd_name',
    'loan_type_name',
    'loan_purpose_name',
    'lien_status_name',
    'hoepa_status_name',
    'county_name',
    'co_applicant_sex_name',
    'co_applicant_ethnicity_name',
    'applicant_sex_name',
    'applicant_ethnicity_name',
    'agency_name',
    'agency_abbr',
    'action_taken_name']

# Perform Chi-square test for each categorical variable
chi2_results_cat = {}
for column in cat_variables:
    contingency_table = pd.crosstab(df_with_clusters[column],
    km_5cluster_model)
    chi2, p_value, _, _ = chi2_contingency(contingency_table)
    chi2_results_cat[column] = {'Chi-square': chi2, 'P-value': p_value}

# Print Chi-square test results for categorical variables
for column, result in chi2_results_cat.items():
```

```
print(f"Variable: {column}")
print(f"Chi-square: {result['Chi-square']}") 
print(f"P-value: {result['P-value']}") 
print()

Variable: state_name
Chi-square: 0.0
P-value: 1.0

Variable: state_abbr
Chi-square: 0.0
P-value: 1.0

Variable: respondent_id
Chi-square: 36123.283985643044
P-value: 0.0

Variable: purchaser_type_name
Chi-square: 3251.6830606356034
P-value: 0.0

Variable: property_type_name
Chi-square: 1056.6104562325181
P-value: 8.972771855445554e-223

Variable: preapproval_name
Chi-square: 225.31724986883384
P-value: 2.895597746443831e-44

Variable: owner_occupancy_name
Chi-square: 294.20914989700015
P-value: 7.029244824919416e-59

Variable: msamd_name
Chi-square: 61680.386313027695
P-value: 0.0

Variable: loan_type_name
Chi-square: 2889.99453116356
P-value: 0.0

Variable: loan_purpose_name
Chi-square: 260.7401082334816
P-value: 9.086960668553878e-52

Variable: lien_status_name
Chi-square: 2163.0211978836865
P-value: 0.0

Variable: hoepa_status_name
Chi-square: 9.634656217418764
```

```
P-value: 0.04705269946396013

Variable: county_name
Chi-square: 100216.46178291588
P-value: 0.0

Variable: co_applicant_sex_name
Chi-square: 1920.4657939451813
P-value: 0.0

Variable: co_applicant_ethnicity_name
Chi-square: 2162.6201219838144
P-value: 0.0

Variable: applicant_sex_name
Chi-square: 1231.0024231266855
P-value: 3.645189435743338e-256

Variable: applicant_ethnicity_name
Chi-square: 1792.6960451546777
P-value: 0.0

Variable: agency_name
Chi-square: 10247.263124936819
P-value: 0.0

Variable: agency_abbr
Chi-square: 10247.263124936819
P-value: 0.0

Variable: action_taken_name
Chi-square: 8141.653891437003
P-value: 0.0

time: 301 ms (started: 2024-03-19 18:43:00 +00:00)
cluster_centers = km_5cluster.cluster_centers_
# Convert centroids to DataFrame for better visualization
centroids_df = pd.DataFrame(cluster_centers,
columns=df_ppd_subset.columns)
%memit
print("Centroids of Clusters:")
centroids_df

peak memory: 520.45 MiB, increment: 0.00 MiB
Centroids of Clusters:

{"type":"dataframe","variable_name":"centroids_df"}

time: 299 ms (started: 2024-03-19 18:43:00 +00:00)
```

```
# Calculate cluster-wise descriptive statistics
cluster_stats = df_ppd_subset.groupby('Cluster_Label').describe()

# Print the descriptive statistics for each cluster
print("Cluster-wise Descriptive Statistics:")
cluster_stats

Cluster-wise Descriptive Statistics:
{"type": "dataframe", "variable_name": "cluster_stats"}
time: 417 ms (started: 2024-03-19 18:43:00 +00:00)

# Grouping the data by cluster labels
cluster_groups = df_with_clusters.groupby(km_5cluster_model)

# Counting the number of variables in each cluster
num_variables_in_clusters = cluster_groups.size()

# Displaying the results
print("Number of Variables in Each Cluster:")
print(num_variables_in_clusters)

Number of Variables in Each Cluster:
0    11435
1    14049
2     9611
3    14162
4    10743
dtype: int64
time: 6.22 ms (started: 2024-03-19 18:43:01 +00:00)

# Joining cluster labels with the original dataset
df_with_clusters = df_ppd_subset.copy()
df_with_clusters['Cluster_Label'] = km_5cluster_model

# Extracting non-categorical variables
non_cat_variables = ['tract_to_msamd_income',
 'population',
 'minority_population',
 'number_of_owner_occupied_units',
 'number_of_1_to_4_family_units',
 'loan_amount_000s',
 'applicant_income_000s',
 'sequence_number',
 'census_tract_number',
 'hud_median_family_income',
 'as_of_year',
 'application_date_indicator']

# Grouping variables by cluster label
```

```
cluster_groups = df_with_clusters.groupby('Cluster_Label')

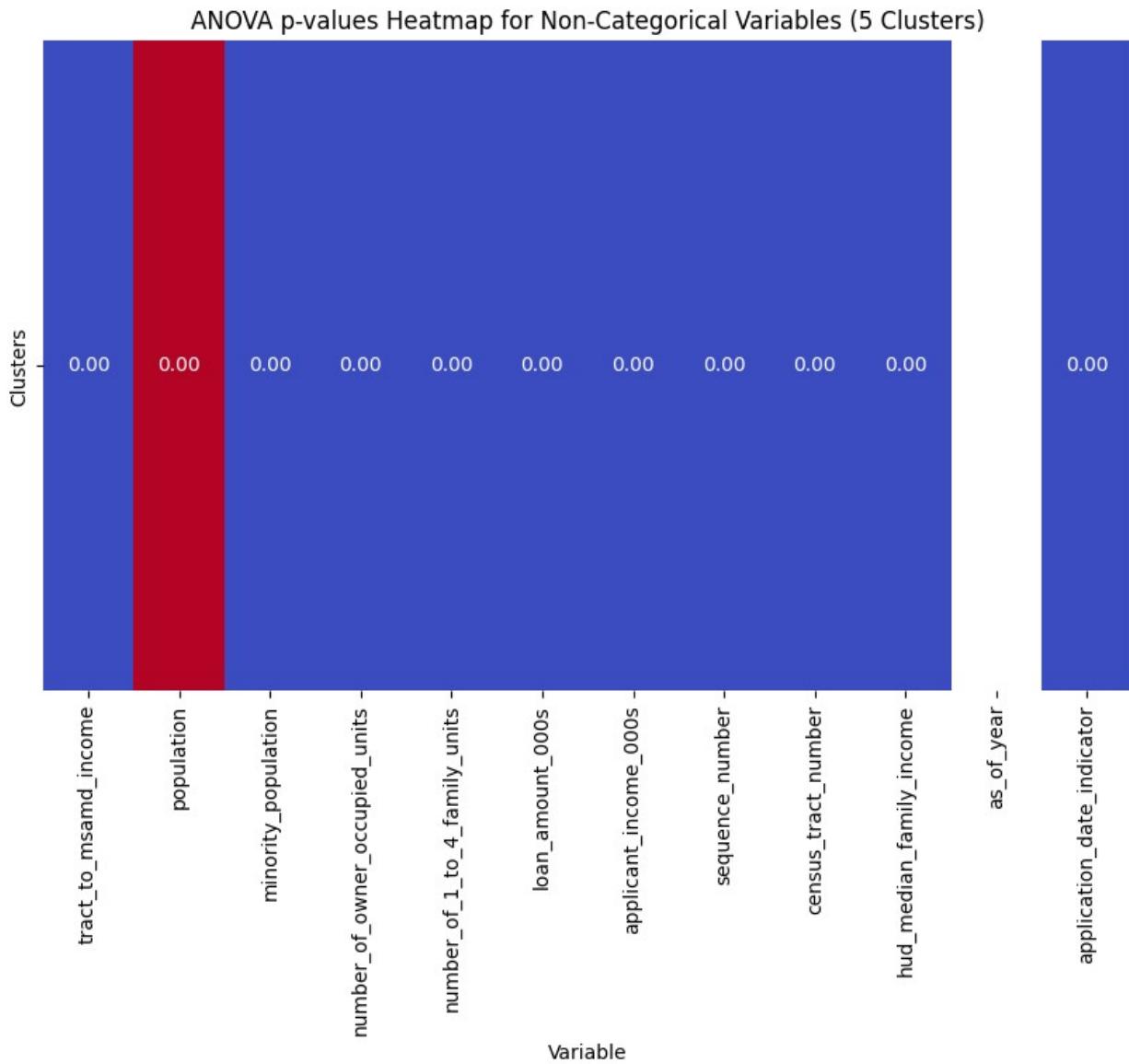
# Perform ANOVA for each non-categorical variable
anova_results_non_cat = {}
for column in non_cat_variables:
    anova_results_non_cat[column] = f_oneway(*[group[column] for name,
group in cluster_groups])

# Extract p-values for non-categorical variables
non_cat_p_values = [result.pvalue for result in
anova_results_non_cat.values()]

# Plotting heatmap for non-categorical variables
plt.figure(figsize=(10, 6))
sns.heatmap([non_cat_p_values], cmap='coolwarm', annot=True,
fmt='.2f', xticklabels=non_cat_variables, yticklabels=['Clusters'],
cbar=False)
plt.xlabel('Variable')
plt.title('ANOVA p-values Heatmap for Non-Categorical Variables (5
Clusters)')
plt.show()

/usr/local/lib/python3.10/dist-packages/scipy/stats/_stats_py.py:4167:
ConstantInputWarning:

Each of the input arrays is constant; the F statistic is not defined or
infinite
```



```
time: 383 ms (started: 2024-03-19 18:43:01 +00:00)
```

```
cat_variables = ['state_name',
 'state_abbr',
 'respondent_id',
 'purchaser_type_name',
 'property_type_name',
 'preapproval_name',
 'owner_occupancy_name',
 'msamd_name',
 'loan_type_name',
 'loan_purpose_name',
 'lien_status_name',
 'hoepa_status_name',
```

```

'county_name',
'co_applicant_sex_name',
'co_applicant_ethnicity_name',
'applicant_sex_name',
'applicant_ethnicity_name',
'agency_name',
'agency_abbr',
'action_taken_name']

# Initialize dictionary to store chi-square results
chi2_results_cat = {}

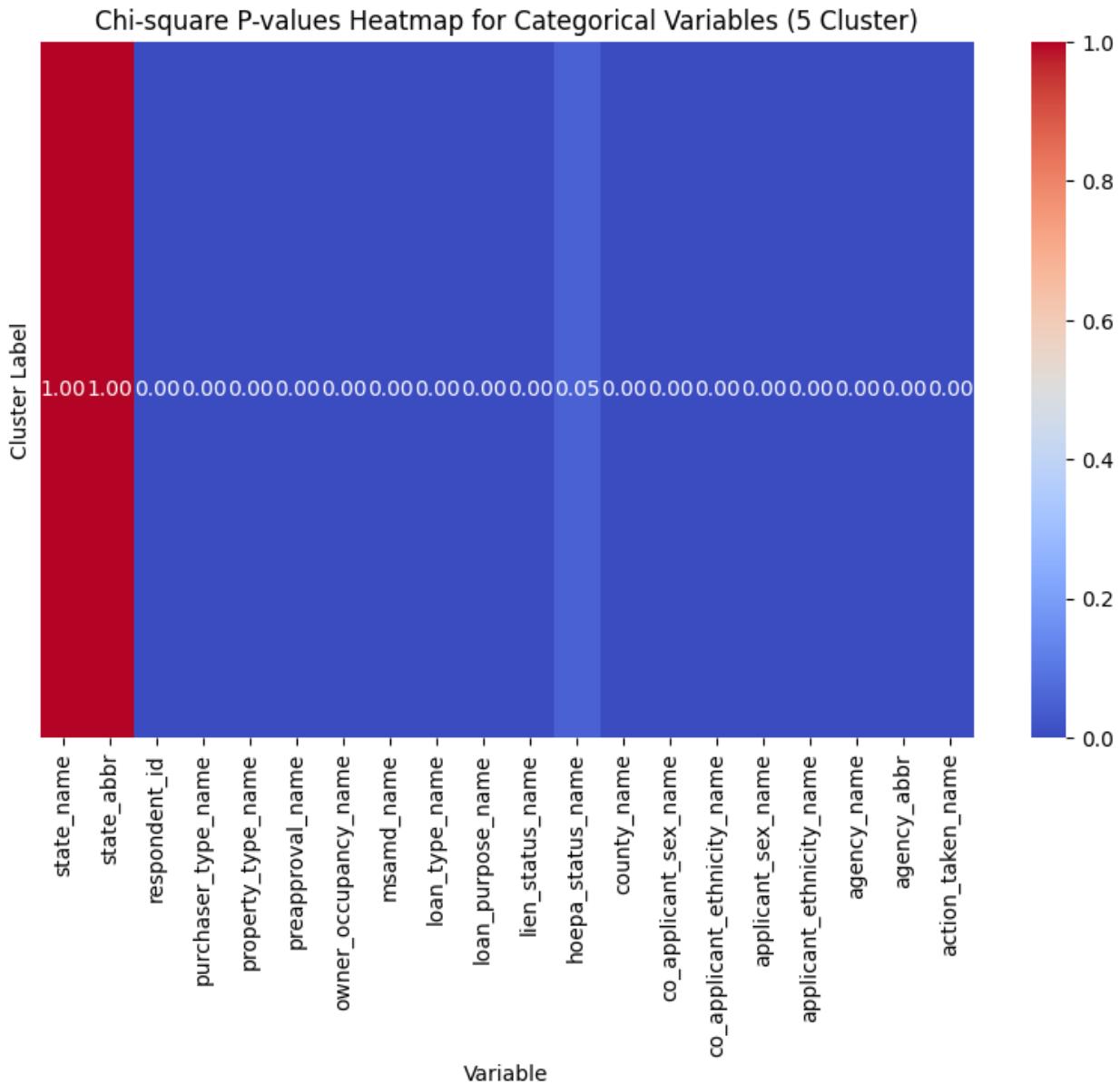
# Calculate chi-square for each categorical variable
for column in cat_variabels:
    contingency_table = pd.crosstab(df_with_clusters[column],
km_5cluster_model)
    chi2, p_value, _, _ = chi2_contingency(contingency_table)
    chi2_results_cat[column] = {'P-value': p_value}

# Extract p-values
p_values = [[result['P-value']] for result in
chi2_results_cat.values()]

# Get variable names
variables = list(chi2_results_cat.keys())

# Plotting heatmap for p-values
plt.figure(figsize=(10, 6))
sns.heatmap(p_values, cmap='coolwarm', annot=True, fmt='.2f',
xticklabels=variables, yticklabels=False)
plt.xlabel('Variable')
plt.ylabel('Cluster Label')
plt.title('Chi-square P-values Heatmap for Categorical Variables (5
Cluster)')
plt.show()

```



```
time: 685 ms (started: 2024-03-19 18:43:01 +00:00)
```

```
import plotly.graph_objs as go
import numpy as np

# Create 3D scatter plot
fig = go.Figure()

# Get unique cluster labels
unique_clusters = np.unique(km_5cluster_model)

# Add traces for each cluster
for cluster_label in unique_clusters:
```

```

cluster_data = df_ppd_subset[km_5cluster_model == cluster_label]
fig.add_trace(go.Scatter3d(
    x=cluster_data['number_of_owner_occupied_units'],
    y=cluster_data['number_of_1_to_4_family_units'],
    z=cluster_data['loan_amount_000s'],
    mode='markers',
    marker=dict(
        size=5,
        color=cluster_label,
        colorscale='Viridis', # Adjust colorscale if needed
        opacity=0.8
    ),
    name=f'Cluster {cluster_label}'
))

# Set layout
fig.update_layout(
    title='Clusters with Centroids (3D)',
    scene=dict(
        xaxis=dict(title='number_of_owner_occupied_units'),
        yaxis=dict(title='number_of_1_to_4_family_units'),
        zaxis=dict(title='loan_amount_000s'),
    ),
    margin=dict(l=0, r=0, b=0, t=40)
)

# Show interactive 3D scatter plot
fig.show()

time: 219 ms (started: 2024-03-19 18:43:02 +00:00)

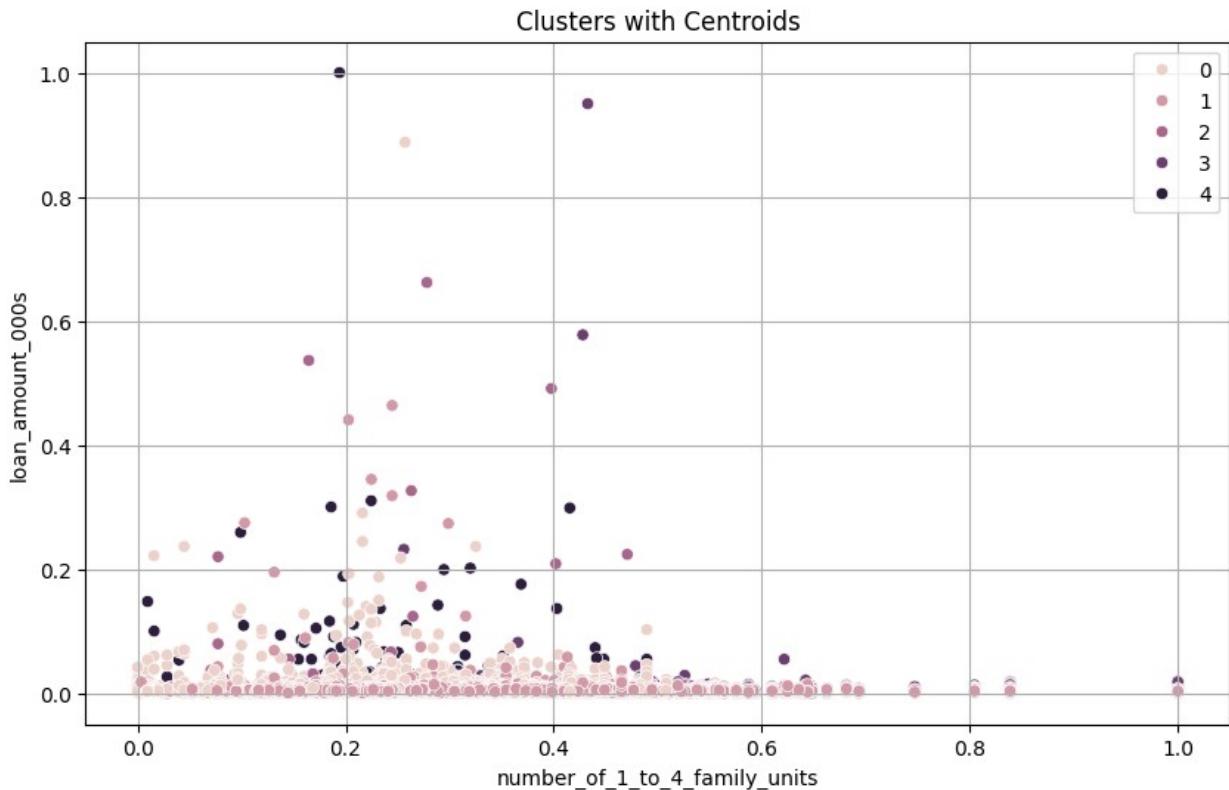
# Assign cluster labels to the DataFrame
df_ppd_subset['Cluster_Label'] = km_5cluster_model

# Plot the scatter plot with clusters and centroids
plt.figure(figsize=(10, 6))
sns.scatterplot(data=df_ppd_subset, x='number_of_1_to_4_family_units',
y='loan_amount_000s', hue='Cluster_Label', legend='full')

# Set plot title and labels
plt.title('Clusters with Centroids')
plt.xlabel('number_of_1_to_4_family_units')
plt.ylabel('loan_amount_000s')
plt.legend()
plt.grid(True)
%memit
plt.show()

peak memory: 513.79 MiB, increment: 0.00 MiB

```



```

time: 4.08 s (started: 2024-03-19 18:43:02 +00:00)

# Create subplots for each cluster
fig, axes = plt.subplots(1, 5, figsize=(20, 4), sharex=True,
sharey=True)
fig.suptitle('Clusters with Centroids')

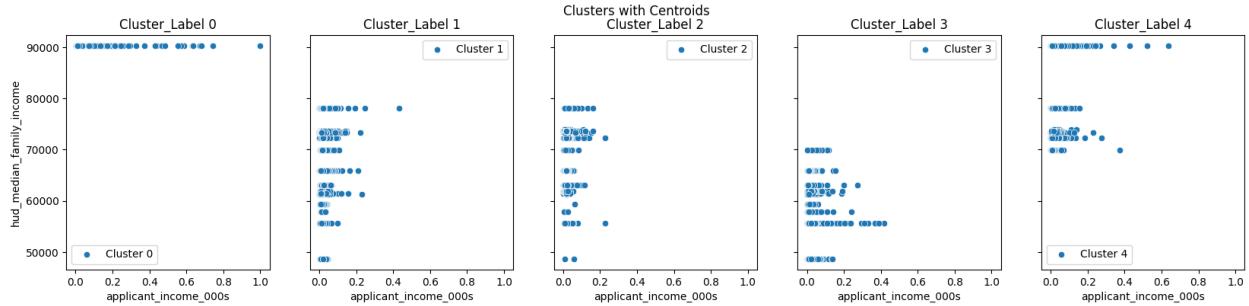
# Iterate through each cluster
for i in range(5):
    # Filter data points belonging to the current cluster
    cluster_data = df_ppd_subset[df_ppd_subset['Cluster_Label'] == i]

    # Scatter plot of data points
    sns.scatterplot(data=cluster_data, x='applicant_income_000s',
y='hud_median_family_income', ax=axes[i], label=f'Cluster {i}')

    # Set title and labels for each subplot
    axes[i].set_title(f'Cluster_Label {i}')
    axes[i].set_xlabel('applicant_income_000s')
    axes[i].set_ylabel('hud_median_family_income')
    axes[i].legend()

plt.show()
%memit

```



```
peak memory: 513.81 MiB, increment: 0.00 MiB
time: 2.08 s (started: 2024-03-19 18:43:06 +00:00)
```

```
# Example code for Elbow Method
import matplotlib.pyplot as plt
from sklearn.cluster import KMeans

# Assuming 'data' is your feature matrix
ssd = []
k_values = range(2, 11)

for k in k_values:
    kmeans = KMeans(n_clusters=k, random_state=42)
    kmeans.fit(df_ppd_subset)
    ssd.append(kmeans.inertia_)

# Plotting the Elbow curve
plt.plot(k_values, ssd, marker='o')
plt.title('Elbow Method for Optimal k')
plt.xlabel('Number of Clusters (k)')
plt.ylabel('Sum of Squared Distances (SSD)')
plt.show()

/usr/local/lib/python3.10/dist-packages/sklearn/cluster/_kmeans.py:870: FutureWarning:
The default value of `n_init` will change from 10 to 'auto' in 1.4.
Set the value of `n_init` explicitly to suppress the warning

/usr/local/lib/python3.10/dist-packages/sklearn/cluster/_kmeans.py:870
: FutureWarning:
The default value of `n_init` will change from 10 to 'auto' in 1.4.
Set the value of `n_init` explicitly to suppress the warning

/usr/local/lib/python3.10/dist-packages/sklearn/cluster/_kmeans.py:870
: FutureWarning:
The default value of `n_init` will change from 10 to 'auto' in 1.4.
Set the value of `n_init` explicitly to suppress the warning
```

```
/usr/local/lib/python3.10/dist-packages/sklearn/cluster/_kmeans.py:870
: FutureWarning:
```

```
The default value of `n_init` will change from 10 to 'auto' in 1.4.
Set the value of `n_init` explicitly to suppress the warning
```

```
/usr/local/lib/python3.10/dist-packages/sklearn/cluster/_kmeans.py:870
: FutureWarning:
```

```
The default value of `n_init` will change from 10 to 'auto' in 1.4.
Set the value of `n_init` explicitly to suppress the warning
```

```
/usr/local/lib/python3.10/dist-packages/sklearn/cluster/_kmeans.py:870
: FutureWarning:
```

```
The default value of `n_init` will change from 10 to 'auto' in 1.4.
Set the value of `n_init` explicitly to suppress the warning
```

```
/usr/local/lib/python3.10/dist-packages/sklearn/cluster/_kmeans.py:870
: FutureWarning:
```

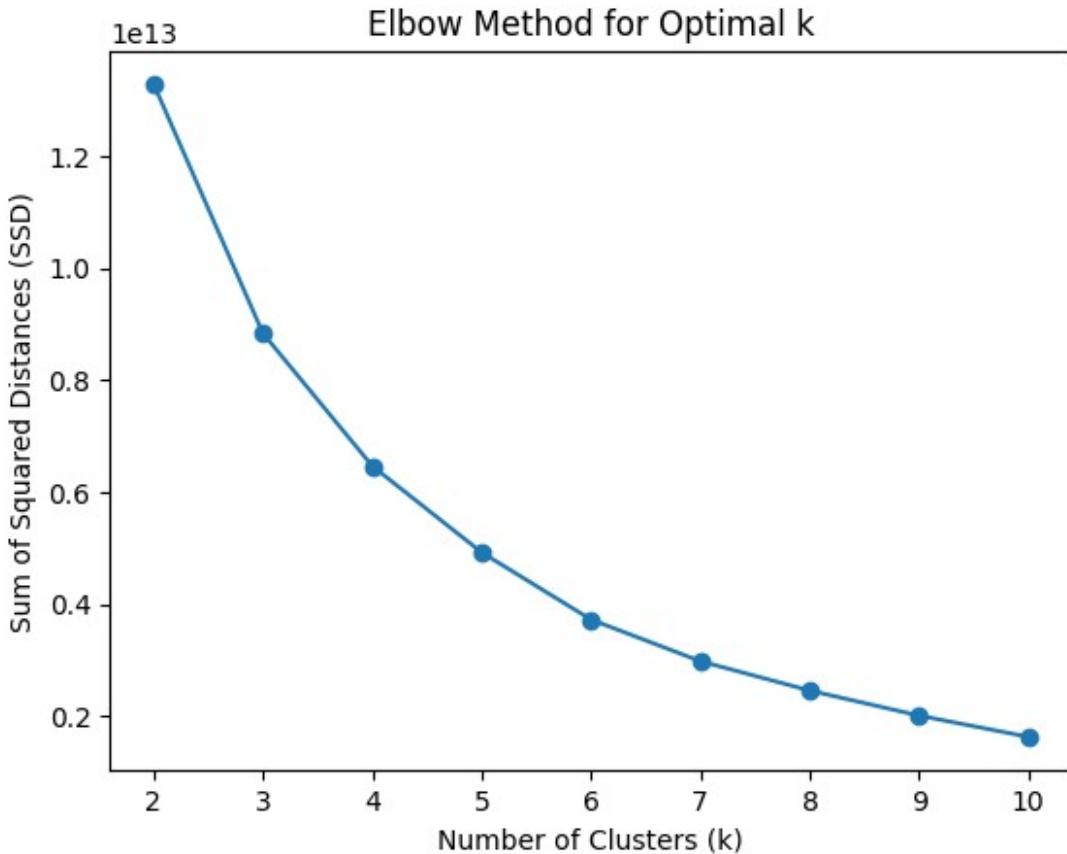
```
The default value of `n_init` will change from 10 to 'auto' in 1.4.
Set the value of `n_init` explicitly to suppress the warning
```

```
/usr/local/lib/python3.10/dist-packages/sklearn/cluster/_kmeans.py:870
: FutureWarning:
```

```
The default value of `n_init` will change from 10 to 'auto' in 1.4.
Set the value of `n_init` explicitly to suppress the warning
```

```
/usr/local/lib/python3.10/dist-packages/sklearn/cluster/_kmeans.py:870
: FutureWarning:
```

```
The default value of `n_init` will change from 10 to 'auto' in 1.4.
Set the value of `n_init` explicitly to suppress the warning
```



```

time: 23.1 s (started: 2024-03-19 18:43:08 +00:00)

import numpy as np

def cramers_v(confusion_matrix):
    chi2 = chi2_contingency(confusion_matrix)[0]
    n = confusion_matrix.sum()
    if n == 0:
        return np.nan
    phi2 = chi2 / n
    r, k = confusion_matrix.shape
    phi2corr = max(0, phi2 - ((k - 1) * (r - 1)) / (n - 1))
    rcorr = r - ((r - 1)**2) / (n - 1)
    kcorr = k - ((k - 1)**2) / (n - 1)
    return np.sqrt(phi2corr / min((kcorr - 1), (rcorr - 1)))

# Calculate Cramér's V for each significant categorical variable
cramers_v_results = {}
for column in chi2_results_cat.keys():
    contingency_table = pd.crosstab(df_with_clusters[column],
    km_5cluster_model)
    cramers_v_value = cramers_v(contingency_table.values)
    cramers_v_results[column] = cramers_v_value

```

```
# Print Cramér's V results for significant categorical variables
for column, result in cramers_v_results.items():
    print(f"Variable: {column}")
    print(f"Cramér's V: {result}")
    print()

<ipython-input-91-82ce08ede275>:13: RuntimeWarning:
invalid value encountered in scalar divide
<ipython-input-91-82ce08ede275>:13: RuntimeWarning:
invalid value encountered in scalar divide

Variable: state_name
Cramér's V: nan

Variable: state_abbr
Cramér's V: nan

Variable: respondent_id
Cramér's V: 0.37504163688570363

Variable: purchaser_type_name
Cramér's V: 0.11575651272755497

Variable: property_type_name
Cramér's V: 0.09348107109294747

Variable: preapproval_name
Cramér's V: 0.042556272628072114

Variable: owner_occupancy_name
Cramér's V: 0.048838045051905615

Variable: msamd_name
Cramér's V: 0.506756410696931

Variable: loan_type_name
Cramér's V: 0.1264502150710077

Variable: loan_purpose_name
Cramér's V: 0.04589375411253423

Variable: lien_status_name
Cramér's V: 0.10931932217389556

Variable: hoepa_status_name
Cramér's V: 0.009690788979401373
```

```

Variable: county_name
Cramér's V: 0.6457267571292904

Variable: co_applicant_sex_name
Cramér's V: 0.08908311901027698

Variable: co_applicant_ethnicity_name
Cramér's V: 0.09457719455621329

Variable: applicant_sex_name
Cramér's V: 0.08229564512908187

Variable: applicant_ethnicity_name
Cramér's V: 0.09946481471353605

Variable: agency_name
Cramér's V: 0.2064374871676273

Variable: agency_abbr
Cramér's V: 0.2064374871676273

Variable: action_taken_name
Cramér's V: 0.18387262778971122

time: 310 ms (started: 2024-03-19 18:43:31 +00:00)

# Aggregate statistics across clusters
cluster_agg_stats =
df_ppd_subset.groupby('Cluster_Label').agg(['mean', 'std', 'min',
'max'])

# Print the aggregated statistics
print(cluster_agg_stats)

\

S.no                                     state_name
\\
mean          std      min      max    mean    std
min
Cluster_Label

0           41320.884040  9123.811002  26151  59994       0.0    0.0
0
1           50466.695352  5780.027202  38849  60000       0.0    0.0
0
2           31307.010717  5947.863420  17963  41676       0.0    0.0
0
3           13936.859130  8374.330385      2  36354       0.0    0.0
0
4           11193.671135  6291.511155      1  23783       0.0    0.0
0

```

```

state_abbr      ... hud_median_family_income
\             max     mean    std   ...
max
Cluster_Label      ... 

0           0     0.0   0.0   ...
90300.0      90300.0
1           0     0.0   0.0   ...
78100.0      78100.0
2           0     0.0   0.0   ...
78100.0      78100.0
3           0     0.0   0.0   ...
69900.0      69900.0
4           0     0.0   0.0   ...
90300.0      90300.0

as_of_year
application_date_indicator \
mean    std      min      max
mean
Cluster_Label

0           2016.0  0.0   2016.0  2016.0
0.003148      0.003148
1           2016.0  0.0   2016.0  2016.0
0.000000      0.000000
2           2016.0  0.0   2016.0  2016.0
0.006867      0.006867
3           2016.0  0.0   2016.0  2016.0
0.102387      0.102387
4           2016.0  0.0   2016.0  2016.0
0.000000      0.000000

std   min   max
Cluster_Label
0       0.079291  0.0   2.0
1       0.000000  0.0   0.0
2       0.116998  0.0   2.0
3       0.440799  0.0   2.0
4       0.000000  0.0   0.0

```

[5 rows x 132 columns]
time: 157 ms (started: 2024-03-19 18:43:32 +00:00)

```
# Select relevant columns
relevant_columns = ['tract_to_msamr_income',
 'population',
```

```

'minority_population',
'number_of_owner_occupied_units',
'number_of_1_to_4_family_units',
'loan_amount_000s',
'applicant_income_000s',
'sequence_number',
'census_tract_number',
'hud_median_family_income',
'as_of_year',
'application_date_indicator']

# Aggregate statistics across clusters for relevant columns
#cluster_agg_stats_relevant = df_ppd_subset.groupby('Cluster_Label')[relevant_columns].agg(['mean', 'std', 'min', 'max'])
cluster_agg_stats_relevant = df_ppd_subset.groupby('Cluster_Label')[relevant_columns].agg(['mean'])
# Print the aggregated statistics for relevant columns
print(cluster_agg_stats_relevant)



|               | tract_to_msamd_income | population | minority_population | \ |
|---------------|-----------------------|------------|---------------------|---|
|               | mean                  | mean       | mean                |   |
| Cluster_Label |                       |            |                     |   |
| 0             | 0.392587              | 0.403339   | 0.301487            |   |
| 1             | 0.382383              | 0.401038   | 0.196050            |   |
| 2             | 0.375166              | 0.391522   | 0.205855            |   |
| 3             | 0.389206              | 0.407582   | 0.189144            |   |
| 4             | 0.383087              | 0.396991   | 0.266054            |   |


|               | number_of_owner_occupied_units | number_of_1_to_4_family_units | \ |
|---------------|--------------------------------|-------------------------------|---|
|               | mean                           | mean                          |   |
| mean          |                                |                               |   |
| Cluster_Label |                                |                               |   |
| 0             | 0.453167                       | 0.274231                      |   |
| 1             | 0.463540                       | 0.310887                      |   |
| 2             | 0.447393                       | 0.302378                      |   |
| 3             | 0.494227                       | 0.380824                      |   |
| 4             | 0.451885                       | 0.286856                      |   |


|               | loan_amount_000s | applicant_income_000s |      |
|---------------|------------------|-----------------------|------|
|               | sequence_number  | \                     |      |
|               | mean             | mean                  | mean |
| Cluster_Label |                  |                       |      |


```

```

0          0.007573        0.023726        0.039357
1          0.004506        0.015841        0.064601
2          0.004702        0.016145        0.065780
3          0.003823        0.015801        0.070423
4          0.006285        0.020141        0.070676

census_tract_number  hud_median_family_income  as_of_year
\               mean           mean           mean

Cluster_Label
0          0.030713        90300.000000      2016.0
1          0.158277        70189.664745      2016.0
2          0.095602        72405.260387      2016.0
3          0.470078        57676.818246      2016.0
4          0.056708        83848.413905      2016.0

application_date_indicator
                           mean
Cluster_Label
0          0.003148
1          0.000000
2          0.006867
3          0.102387
4          0.000000
time: 29 ms (started: 2024-03-19 18:43:32 +00:00)

# Define the non-categorical variables
non_cat_variables = ['tract_to_msamd_income', 'population',
'minority_population',
'number_of_owner_occupied_units',
'number_of_1_to_4_family_units',
'loan_amount_000s', 'applicant_income_000s',
'sequence_number',
'census_tract_number',
'hud_median_family_income', 'as_of_year',
'application_date_indicator']

# Create a dictionary to store the average values for each non-

```

```

categorical variable in each cluster
cluster_profiles_non_cat = {}

# Calculate the average values for each non-categorical variable in each cluster
for variable in non_cat_variables:
    cluster_profiles_non_cat[variable] =
df_ppd_subset.groupby('Cluster_Label')[variable].mean()

# Print the cluster profiling for non-categorical variables
for variable, averages in cluster_profiles_non_cat.items():
    print(f"Variable: {variable}")
    for cluster, avg_value in averages.items():
        print(f"Cluster {cluster}: Average {variable} - {avg_value:.2f}")
    print()

```

Variable: tract_to_msamd_income

Cluster 0: Average tract_to_msamd_income - 0.39
 Cluster 1: Average tract_to_msamd_income - 0.38
 Cluster 2: Average tract_to_msamd_income - 0.38
 Cluster 3: Average tract_to_msamd_income - 0.39
 Cluster 4: Average tract_to_msamd_income - 0.38

Variable: population

Cluster 0: Average population - 0.40
 Cluster 1: Average population - 0.40
 Cluster 2: Average population - 0.39
 Cluster 3: Average population - 0.41
 Cluster 4: Average population - 0.40

Variable: minority_population

Cluster 0: Average minority_population - 0.30
 Cluster 1: Average minority_population - 0.20
 Cluster 2: Average minority_population - 0.21
 Cluster 3: Average minority_population - 0.19
 Cluster 4: Average minority_population - 0.27

Variable: number_of_owner_occupied_units

Cluster 0: Average number_of_owner_occupied_units - 0.45
 Cluster 1: Average number_of_owner_occupied_units - 0.46
 Cluster 2: Average number_of_owner_occupied_units - 0.45
 Cluster 3: Average number_of_owner_occupied_units - 0.49
 Cluster 4: Average number_of_owner_occupied_units - 0.45

Variable: number_of_1_to_4_family_units

Cluster 0: Average number_of_1_to_4_family_units - 0.27
 Cluster 1: Average number_of_1_to_4_family_units - 0.31
 Cluster 2: Average number_of_1_to_4_family_units - 0.30
 Cluster 3: Average number_of_1_to_4_family_units - 0.38

```
Cluster 4: Average number_of_1_to_4_family_units - 0.29
```

```
Variable: loan_amount_000s
```

```
Cluster 0: Average loan_amount_000s - 0.01
Cluster 1: Average loan_amount_000s - 0.00
Cluster 2: Average loan_amount_000s - 0.00
Cluster 3: Average loan_amount_000s - 0.00
Cluster 4: Average loan_amount_000s - 0.01
```

```
Variable: applicant_income_000s
```

```
Cluster 0: Average applicant_income_000s - 0.02
Cluster 1: Average applicant_income_000s - 0.02
Cluster 2: Average applicant_income_000s - 0.02
Cluster 3: Average applicant_income_000s - 0.02
Cluster 4: Average applicant_income_000s - 0.02
```

```
Variable: sequence_number
```

```
Cluster 0: Average sequence_number - 0.04
Cluster 1: Average sequence_number - 0.06
Cluster 2: Average sequence_number - 0.07
Cluster 3: Average sequence_number - 0.07
Cluster 4: Average sequence_number - 0.07
```

```
Variable: census_tract_number
```

```
Cluster 0: Average census_tract_number - 0.03
Cluster 1: Average census_tract_number - 0.16
Cluster 2: Average census_tract_number - 0.10
Cluster 3: Average census_tract_number - 0.47
Cluster 4: Average census_tract_number - 0.06
```

```
Variable: hud_median_family_income
```

```
Cluster 0: Average hud_median_family_income - 90300.00
Cluster 1: Average hud_median_family_income - 70189.66
Cluster 2: Average hud_median_family_income - 72405.26
Cluster 3: Average hud_median_family_income - 57676.82
Cluster 4: Average hud_median_family_income - 83848.41
```

```
Variable: as_of_year
```

```
Cluster 0: Average as_of_year - 2016.00
Cluster 1: Average as_of_year - 2016.00
Cluster 2: Average as_of_year - 2016.00
Cluster 3: Average as_of_year - 2016.00
Cluster 4: Average as_of_year - 2016.00
```

```
Variable: application_date_indicator
```

```
Cluster 0: Average application_date_indicator - 0.00
Cluster 1: Average application_date_indicator - 0.00
Cluster 2: Average application_date_indicator - 0.01
Cluster 3: Average application_date_indicator - 0.10
Cluster 4: Average application_date_indicator - 0.00
```

```

time: 48.1 ms (started: 2024-03-19 18:43:32 +00:00)

# Select relevant columns
relevant_columns = [ 'purchaser_type_name',
'property_type_name',
'preapproval_name',
'owner_occupancy_name',
'msamd_name',
'loan_type_name',
'loan_purpose_name',
'lien_status_name',
'hoepa_status_name',
'county_name',
'co_applicant_sex_name',
'co_applicant_ethnicity_name',
'applicant_sex_name',
'applicant_ethnicity_name',
'agency_name',
'agency_abbr',
'action_taken_name']

# Create a dictionary to store frequencies for each column and cluster
cluster_freq_stats_relevant = {}

# Iterate over relevant columns
for column in relevant_columns:
    # Group the data by 'Cluster_Label' and get value counts for each cluster
    cluster_freq_stats_relevant[column] =
df_ppd_subset.groupby('Cluster_Label')[column].value_counts()

# Print the frequency statistics for relevant columns
for column, value_counts in cluster_freq_stats_relevant.items():
    print(f"Variable: {column}")
    print(value_counts)
    print()

Variable: purchaser_type_name
Cluster_Label  purchaser_type_name
0              2                  4048
               7                  3521
               4                  1639
               5                  815
               6                  626
               1                  417
               8                  268
               9                  65
               0                  36
1              2                  3027

```

	5	2961
	7	2696
	4	2090
	6	1395
	1	961
	8	628
	0	168
	9	122
	3	1
2	7	3115
	2	1872
	5	1710
	4	1192
	6	814
	1	474
	8	281
	0	87
	9	66
3	7	4479
	5	2738
	2	2673
	4	1776
	6	946
	1	669
	8	652
	0	155
	9	74
4	2	2783
	7	2301
	4	1592
	5	1557
	6	1084
	1	797
	8	466
	0	107
	9	56

Name: purchaser_type_name, dtype: int64

Variable: property_type_name

Cluster_Label	property_type_name	
0	2	11170
	1	140
	0	125
1	2	13509
	0	498
	1	42
2	2	9324
	0	264
	1	23

```
3          2      13135
          0      1003
          1       24
4          2      10492
          0      171
          1       80
```

Name: property_type_name, dtype: int64

Variable: preapproval_name

```
Cluster_Label  preapproval_name
0            0      9248
            1      1559
            2       628
1            0     10959
            1      2444
            2       646
2            0      7719
            1      1558
            2       334
3            0     11156
            1      2564
            2       442
4            0      8750
            1      1573
            2       420
```

Name: preapproval_name, dtype: int64

Variable: owner_occupancy_name

```
Cluster_Label  owner_occupancy_name
0            2      10286
            1      1011
            0       138
1            2     12866
            1      1144
            0       39
2            2      8708
            1      872
            0       31
3            2     12391
            1      1725
            0       46
4            2      9689
            1      972
            0       82
```

Name: owner_occupancy_name, dtype: int64

Variable: msamd_name

```
Cluster_Label  msamd_name
0            8      11435
1            7      4143
```

	2	2301
	1	2257
	6	1674
	5	1198
	4	666
	10	527
	0	426
	9	238
	8	218
	13	169
	12	120
	3	100
	11	12
2	7	3653
	10	2406
	6	2175
	1	325
	2	279
	12	247
	0	175
	8	169
	9	158
	5	15
	11	4
	4	2
	13	2
	3	1
3	8	7978
	13	1414
	2	1199
	9	957
	0	815
	4	696
	5	446
	11	345
	12	221
	3	91
4	8	6546
	10	1676
	7	977
	1	856
	6	597
	0	91

Name: msamd_name, dtype: int64

Variable: loan_type_name
 Cluster_Label loan_type_name
 0 0 10169
 1 667

	3	578
	2	21
1	0	9230
	3	2720
	1	1931
	2	168
2	0	6233
	3	1894
	1	1414
	2	70
3	0	9135
	3	2412
	1	2277
	2	338
4	0	8150
	3	1339
	1	1208
	2	46

Name: loan_type_name, dtype: int64

Variable: loan_purpose_name

Cluster_Label	loan_purpose_name	
0	2	5373
	1	5291
	0	771
1	1	6964
	2	6300
	0	785
2	2	4718
	1	4371
	0	522
3	1	6991
	2	6463
	0	708
4	2	5722
	1	4490
	0	531

Name: loan_purpose_name, dtype: int64

Variable: lien_status_name

Cluster_Label	lien_status_name	
0	2	10996
	3	303
	1	118
	0	18
1	2	13528
	3	350
	1	171
2	2	9235

```
3          233
1          110
0           33
3          12879
2          725
0          348
3          210
1          10408
2          237
3           98
Name: lien_status_name, dtype: int64
```

```
Variable: hoepa_status_name
Cluster_Label  hoepa_status_name
0            1           11435
1            1           14048
0             1            1
2            1           9611
3            1           14162
4            1           10740
0             3            3
Name: hoepa_status_name, dtype: int64
```

```
Variable: county_name
Cluster_Label  county_name
0            16          8312
30           3123
1            5           4064
17           2257
2            1674
...
4            36           91
29           20
0             1
27           1
31           1
Name: county_name, Length: 114, dtype: int64
```

```
Variable: co_applicant_sex_name
Cluster_Label  co_applicant_sex_name
0            3           5577
0           3717
1           1154
2           968
4            19
1            3           6256
0           5891
2           1210
1           685
4             7
Name: co_applicant_sex_name, dtype: int64
```

2	3	4361
	0	3891
	2	819
	1	517
	4	23
3	3	5900
	0	5365
	1	1290
	2	1165
	4	442
4	3	4893
	0	3686
	1	1190
	2	968
	4	6

Name: co_applicant_sex_name, dtype: int64

Variable: co_applicant_ethnicity_name

Cluster_Label	co_applicant_ethnicity_name	
0	2	5577
	3	4065
	1	1619
	0	155
	4	19
1	3	6298
	2	6256
	1	1050
	0	435
	4	10
2	2	4361
	3	4132
	1	890
	0	201
	4	27
3	2	5900
	3	5445
	1	1907
	0	467
	4	443
4	2	4893
	3	3809
	1	1819
	0	215
	4	7

Name: co_applicant_ethnicity_name, dtype: int64

Variable: applicant_sex_name

Cluster_Label	applicant_sex_name	
0	2	6747

	0	2521
	1	1943
	3	224
1	2	9425
	0	3329
	1	1071
	3	224
2	2	6249
	0	2324
	1	841
	3	197
3	2	8462
	0	3114
	1	1955
	3	631
4	2	6187
	0	2405
	1	1913
	3	238

Name: applicant_sex_name, dtype: int64

Variable: applicant_ethnicity_name

Cluster_Label	applicant_ethnicity_name	
0	2	8194
	1	2670
	0	345
	3	226
1	2	11323
	1	1600
	0	898
	3	228
2	2	7564
	1	1441
	0	407
	3	199
3	2	9729
	1	2846
	0	949
	3	638
4	2	7204
	1	2901
	0	397
	3	241

Name: applicant_ethnicity_name, dtype: int64

Variable: agency_name

Cluster_Label	agency_name	
0	2	4837
	0	3188

	4	1865
	1	1333
	5	192
	3	20
1	1	7372
	0	3876
	2	1466
	4	986
	5	244
	3	105
2	1	5068
	0	2216
	2	1198
	4	927
	5	169
	3	33
3	1	7412
	0	3752
	2	1544
	4	1194
	5	161
	3	99
4	1	6329
	0	2580
	2	943
	4	685
	5	157
	3	49

Name: agency_name, dtype: int64

Variable:	agency_abbr	
Cluster_Label	agency_abbr	
0	1	4837
	0	3188
	4	1865
	3	1333
	5	192
	2	20
1	3	7372
	0	3876
	1	1466
	4	986
	5	244
	2	105
2	3	5068
	0	2216
	1	1198
	4	927
	5	169

```
      2          33
3      3         7412
      0         3752
      1        1544
      4        1194
      5         161
      2         99
4      3        6329
      0        2580
      1         943
      4        685
      5        157
      2         49
```

Name: agency_abbr, dtype: int64

Variable: action_taken_name
Cluster_Label action_taken_name

```
0      4        11368
      2         30
      5         18
      3         11
      6          7
      7          1
1      4        14049
2      4        8109
      2        1333
      3         49
      1         37
      7         34
      5         33
      6         10
      0          6
3      4        11620
      2        1267
      5         725
      3         386
      1         89
      0         75
4      4        10669
      2         36
      1         35
      0          3
```

Name: action_taken_name, dtype: int64

time: 159 ms (started: 2024-03-19 18:43:32 +00:00)

```
# Define a function to get the dominant category for each cluster
def get_dominant_category(cluster_data):
    return cluster_data.idxmax()
```

```

# Create a dictionary to store the dominant categories for each
# cluster and variable
cluster_profiles = {}

# Iterate over relevant columns
for column in relevant_columns:
    # Group the data by 'Cluster_Label' and get value counts for each
    # cluster
    cluster_freq = df_ppd_subset.groupby('Cluster_Label')
    [column].value_counts()

    # Find the dominant category for each cluster
    dominant_categories =
    cluster_freq.groupby(level=0).apply(get_dominant_category)

    # Store the dominant categories in the dictionary
    cluster_profiles[column] = dominant_categories

# Print the cluster profiling
for column, dominant_categories in cluster_profiles.items():
    print(f"Variable: {column}")
    for cluster, category in dominant_categories.items():
        print(f"Cluster {cluster}: Dominant category - {category}")
    print()

```

Variable: purchaser_type_name

Cluster 0: Dominant category - (0, 2)

Cluster 1: Dominant category - (1, 2)

Cluster 2: Dominant category - (2, 7)

Cluster 3: Dominant category - (3, 7)

Cluster 4: Dominant category - (4, 2)

Variable: property_type_name

Cluster 0: Dominant category - (0, 2)

Cluster 1: Dominant category - (1, 2)

Cluster 2: Dominant category - (2, 2)

Cluster 3: Dominant category - (3, 2)

Cluster 4: Dominant category - (4, 2)

Variable: preapproval_name

Cluster 0: Dominant category - (0, 0)

Cluster 1: Dominant category - (1, 0)

Cluster 2: Dominant category - (2, 0)

Cluster 3: Dominant category - (3, 0)

Cluster 4: Dominant category - (4, 0)

Variable: owner_occupancy_name

Cluster 0: Dominant category - (0, 2)

Cluster 1: Dominant category - (1, 2)

Cluster 2: Dominant category - (2, 2)

```
Cluster 3: Dominant category - (3, 2)
Cluster 4: Dominant category - (4, 2)
```

```
Variable: msamd_name
Cluster 0: Dominant category - (0, 8)
Cluster 1: Dominant category - (1, 7)
Cluster 2: Dominant category - (2, 7)
Cluster 3: Dominant category - (3, 8)
Cluster 4: Dominant category - (4, 8)
```

```
Variable: loan_type_name
Cluster 0: Dominant category - (0, 0)
Cluster 1: Dominant category - (1, 0)
Cluster 2: Dominant category - (2, 0)
Cluster 3: Dominant category - (3, 0)
Cluster 4: Dominant category - (4, 0)
```

```
Variable: loan_purpose_name
Cluster 0: Dominant category - (0, 2)
Cluster 1: Dominant category - (1, 1)
Cluster 2: Dominant category - (2, 2)
Cluster 3: Dominant category - (3, 1)
Cluster 4: Dominant category - (4, 2)
```

```
Variable: lien_status_name
Cluster 0: Dominant category - (0, 2)
Cluster 1: Dominant category - (1, 2)
Cluster 2: Dominant category - (2, 2)
Cluster 3: Dominant category - (3, 2)
Cluster 4: Dominant category - (4, 2)
```

```
Variable: hoepa_status_name
Cluster 0: Dominant category - (0, 1)
Cluster 1: Dominant category - (1, 1)
Cluster 2: Dominant category - (2, 1)
Cluster 3: Dominant category - (3, 1)
Cluster 4: Dominant category - (4, 1)
```

```
Variable: county_name
Cluster 0: Dominant category - (0, 16)
Cluster 1: Dominant category - (1, 5)
Cluster 2: Dominant category - (2, 5)
Cluster 3: Dominant category - (3, 14)
Cluster 4: Dominant category - (4, 16)
```

```
Variable: co_applicant_sex_name
Cluster 0: Dominant category - (0, 3)
Cluster 1: Dominant category - (1, 3)
Cluster 2: Dominant category - (2, 3)
Cluster 3: Dominant category - (3, 3)
```

```
Cluster 4: Dominant category - (4, 3)
```

```
Variable: co_applicant_ethnicity_name
Cluster 0: Dominant category - (0, 2)
Cluster 1: Dominant category - (1, 3)
Cluster 2: Dominant category - (2, 2)
Cluster 3: Dominant category - (3, 2)
Cluster 4: Dominant category - (4, 2)
```

```
Variable: applicant_sex_name
Cluster 0: Dominant category - (0, 2)
Cluster 1: Dominant category - (1, 2)
Cluster 2: Dominant category - (2, 2)
Cluster 3: Dominant category - (3, 2)
Cluster 4: Dominant category - (4, 2)
```

```
Variable: applicant_ethnicity_name
Cluster 0: Dominant category - (0, 2)
Cluster 1: Dominant category - (1, 2)
Cluster 2: Dominant category - (2, 2)
Cluster 3: Dominant category - (3, 2)
Cluster 4: Dominant category - (4, 2)
```

```
Variable: agency_name
Cluster 0: Dominant category - (0, 2)
Cluster 1: Dominant category - (1, 1)
Cluster 2: Dominant category - (2, 1)
Cluster 3: Dominant category - (3, 1)
Cluster 4: Dominant category - (4, 1)
```

```
Variable: agency_abbr
Cluster 0: Dominant category - (0, 1)
Cluster 1: Dominant category - (1, 3)
Cluster 2: Dominant category - (2, 3)
Cluster 3: Dominant category - (3, 3)
Cluster 4: Dominant category - (4, 3)
```

```
Variable: action_taken_name
Cluster 0: Dominant category - (0, 4)
Cluster 1: Dominant category - (1, 4)
Cluster 2: Dominant category - (2, 4)
Cluster 3: Dominant category - (3, 4)
Cluster 4: Dominant category - (4, 4)
```

```
time: 149 ms (started: 2024-03-19 18:43:32 +00:00)
```

```
# Filter the dataframe to include only data points from Cluster 4
cluster_4_data = df_ppd_subset[df_ppd_subset['Cluster_Label'] == 4]
```

```
# List of categorical variables
```

```

categorical_variables = ['purchaser_type_name', 'property_type_name',
'preapproval_name', 'owner_occupancy_name', 'msamd_name',
'loan_type_name', 'loan_purpose_name', 'lien_status_name',
'hoepa_status_name', 'county_name', 'co_applicant_sex_name',
'co_applicant_ethnicity_name', 'applicant_sex_name',
'applicant_ethnicity_name', 'agency_name',
'action_taken_name']

# List of numerical variables
numerical_variables = ['tract_to_msamd_income', 'population',
'minority_population',
'number_of_owner_occupied_units',
'number_of_1_to_4_family_units',
'loan_amount_000s', 'applicant_income_000s',
'census_tract_number',
'hud_median_family_income',
'application_date_indicator']

# Frequency calculation for categorical variables
categorical_freq =
cluster_4_data[categorical_variables].apply(pd.Series.value_counts)

# Mean calculation for numerical variables
numerical_mean = cluster_4_data[numerical_variables].mean()

# Display the frequencies and mean values
print("Frequency of Categorical Variables in Cluster 4:")
print(categorical_freq)
print("\nMean of Numerical Variables in Cluster 4:")
print(numerical_mean)

Frequency of Categorical Variables in Cluster 4:
  purchaser_type_name  property_type_name  preapproval_name \
0              107.0            171.0        8750.0
1              797.0             80.0        1573.0
2             2783.0           10492.0         420.0
3                NaN               NaN            NaN
4              1592.0             NaN            NaN
5              1557.0             NaN            NaN
6              1084.0             NaN            NaN
7              2301.0             NaN            NaN
8               466.0             NaN            NaN
9                 56.0             NaN            NaN
10                NaN             NaN            NaN
16                NaN             NaN            NaN
17                NaN             NaN            NaN

```

26	NaN	NaN	NaN
27	NaN	NaN	NaN
29	NaN	NaN	NaN
30	NaN	NaN	NaN
31	NaN	NaN	NaN
33	NaN	NaN	NaN
36	NaN	NaN	NaN
loan_purpose_name	owner_occupancy_name	msamid_name	loan_type_name
531.0	82.0	91.0	8150.0
4490.0	972.0	856.0	1208.0
5722.0	9689.0	NaN	46.0
3	NaN	NaN	1339.0
NaN	NaN	NaN	NaN
4	NaN	NaN	NaN
NaN	NaN	NaN	NaN
5	NaN	NaN	NaN
NaN	NaN	NaN	NaN
6	NaN	597.0	NaN
NaN	NaN	977.0	NaN
7	NaN	6546.0	NaN
NaN	NaN	1676.0	NaN
8	NaN	NaN	NaN
NaN	NaN	NaN	NaN
9	NaN	NaN	NaN
NaN	NaN	NaN	NaN
10	NaN	NaN	NaN
NaN	NaN	NaN	NaN
16	NaN	NaN	NaN
NaN	NaN	NaN	NaN
17	NaN	NaN	NaN
NaN	NaN	NaN	NaN
26	NaN	NaN	NaN
NaN	NaN	NaN	NaN
27	NaN	NaN	NaN
NaN	NaN	NaN	NaN
29	NaN	NaN	NaN
NaN	NaN	NaN	NaN
30	NaN	NaN	NaN
NaN	NaN	NaN	NaN
31	NaN	NaN	NaN
NaN	NaN	NaN	NaN
33	NaN	NaN	NaN
NaN	NaN	NaN	NaN
36	NaN	NaN	NaN

NaN

	lien_status_name	hoepa_status_name	county_name
co_applicant_sex_name	\\		
0	NaN	3.0	1.0
3686.0			
1	98.0	10740.0	NaN
1190.0			
2	10408.0	NaN	NaN
968.0			
3	237.0	NaN	NaN
4893.0			
4	NaN	NaN	NaN
6.0			
5	NaN	NaN	957.0
NaN			
6	NaN	NaN	NaN
NaN			
7	NaN	NaN	NaN
NaN			
8	NaN	NaN	NaN
NaN			
9	NaN	NaN	NaN
NaN			
10	NaN	NaN	NaN
NaN			
16	NaN	NaN	4615.0
NaN			
17	NaN	NaN	856.0
NaN			
26	NaN	NaN	1676.0
NaN			
27	NaN	NaN	1.0
NaN			
29	NaN	NaN	20.0
NaN			
30	NaN	NaN	1928.0
NaN			
31	NaN	NaN	1.0
NaN			
33	NaN	NaN	597.0
NaN			
36	NaN	NaN	91.0
NaN			

	co_applicant_ethnicity_name	applicant_sex_name
applicant_ethnicity_name	\\	
0	215.0	2405.0
397.0		

1	1819.0	1913.0
2901.0		
2	4893.0	6187.0
7204.0		
3	3809.0	238.0
241.0		
4	7.0	NaN
NaN		
5	NaN	NaN
NaN		
6	NaN	NaN
NaN		
7	NaN	NaN
NaN		
8	NaN	NaN
NaN		
9	NaN	NaN
NaN		
10	NaN	NaN
NaN		
16	NaN	NaN
NaN		
17	NaN	NaN
NaN		
26	NaN	NaN
NaN		
27	NaN	NaN
NaN		
29	NaN	NaN
NaN		
30	NaN	NaN
NaN		
31	NaN	NaN
NaN		
33	NaN	NaN
NaN		
36	NaN	NaN
NaN		
agency_name	action_taken_name	
0	2580.0	3.0
1	6329.0	35.0
2	943.0	36.0
3	49.0	NaN
4	685.0	10669.0
5	157.0	NaN
6	NaN	NaN
7	NaN	NaN
8	NaN	NaN

```

9      NaN      NaN
10     NaN      NaN
16     NaN      NaN
17     NaN      NaN
26     NaN      NaN
27     NaN      NaN
29     NaN      NaN
30     NaN      NaN
31     NaN      NaN
33     NaN      NaN
36     NaN      NaN

Mean of Numerical Variables in Cluster 4:
tract_to_msamd_income          0.383087
population                      0.396991
minority_population              0.266054
number_of_owner_occupied_units  0.451885
number_of_1_to_4_family_units   0.286856
loan_amount_000s                 0.006285
applicant_income_000s            0.020141
census_tract_number              0.056708
hud_median_family_income        83848.413905
application_date_indicator      0.000000
dtype: float64
time: 42.4 ms (started: 2024-03-19 18:43:32 +00:00)

import matplotlib.pyplot as plt
from sklearn.decomposition import PCA
from sklearn.manifold import TSNE

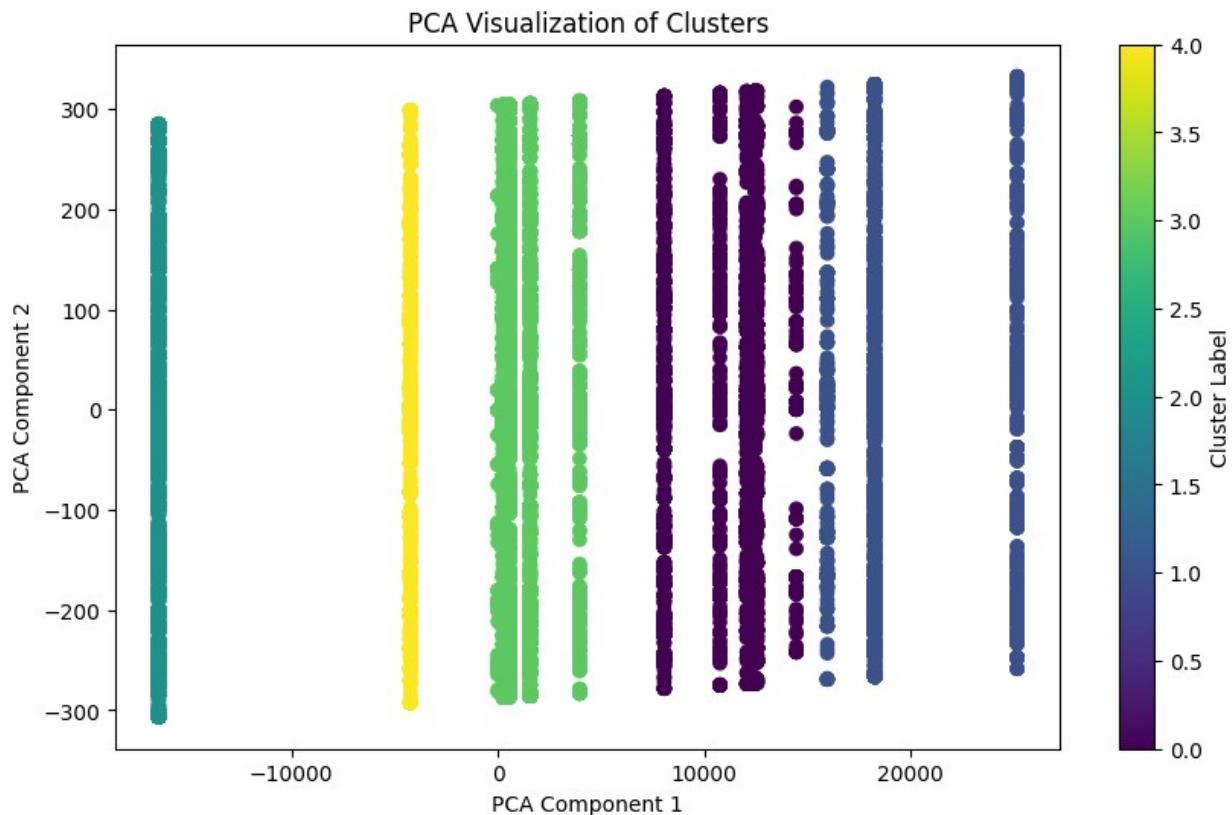
# Reduce dimensionality using PCA
pca = PCA(n_components=2)
pca_result =
pca.fit_transform(df_ppd_subset.drop(columns=['Cluster_Label']))

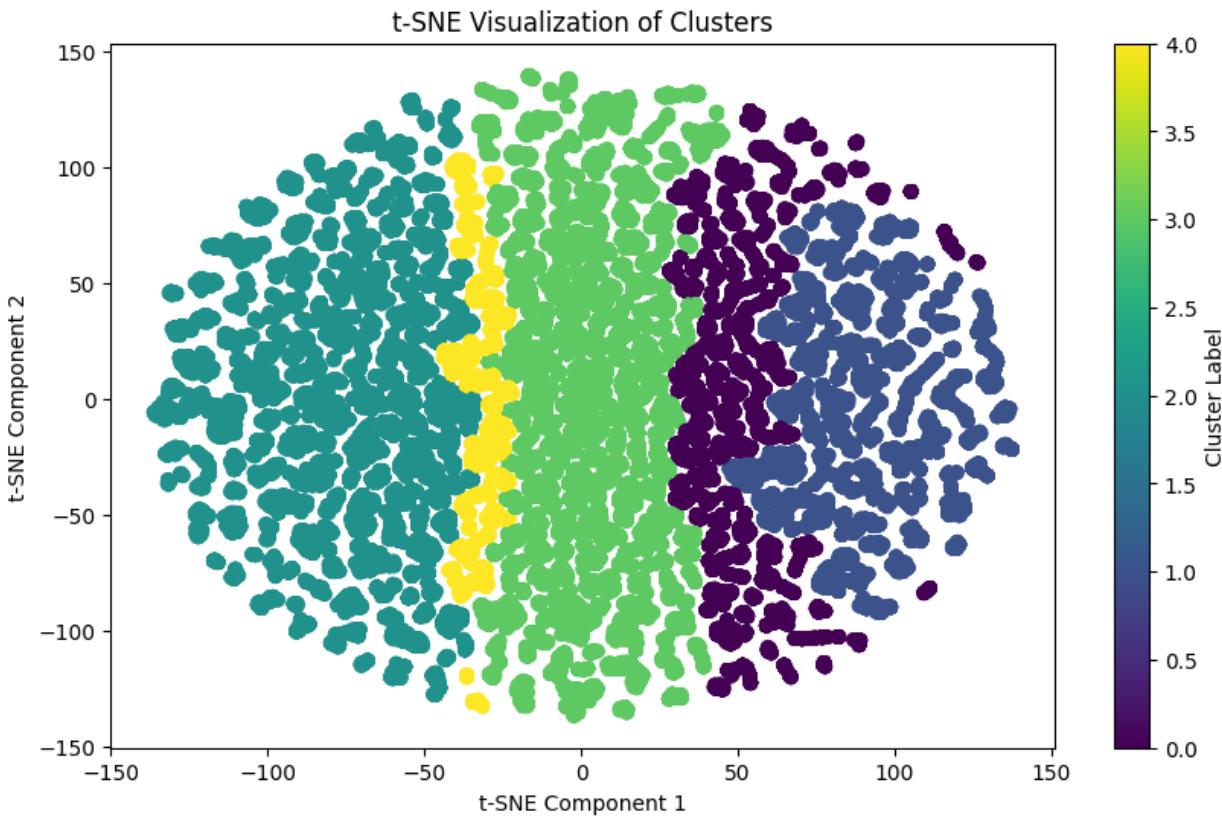
# Visualize clusters using PCA
plt.figure(figsize=(10, 6))
plt.scatter(pca_result[:, 0], pca_result[:, 1],
c=df_ppd_subset['Cluster_Label'], cmap='viridis')
plt.title('PCA Visualization of Clusters')
plt.xlabel('PCA Component 1')
plt.ylabel('PCA Component 2')
plt.colorbar(label='Cluster Label')
plt.show()

# Reduce dimensionality using t-SNE
tsne = TSNE(n_components=2, perplexity=30, random_state=42)
tsne_result =
tsne.fit_transform(df_ppd_subset.drop(columns=['Cluster_Label']))

```

```
# Visualize clusters using t-SNE
plt.figure(figsize=(10, 6))
plt.scatter(tsne_result[:, 0], tsne_result[:, 1],
c=df_ppd_subset['Cluster_Label'], cmap='viridis')
plt.title('t-SNE Visualization of Clusters')
plt.xlabel('t-SNE Component 1')
plt.ylabel('t-SNE Component 2')
plt.colorbar(label='Cluster Label')
plt.show()
```





```

time: 15min 41s (started: 2024-03-16 13:34:25 +00:00)

from sklearn.cluster import DBSCAN
from sklearn.metrics import silhouette_score

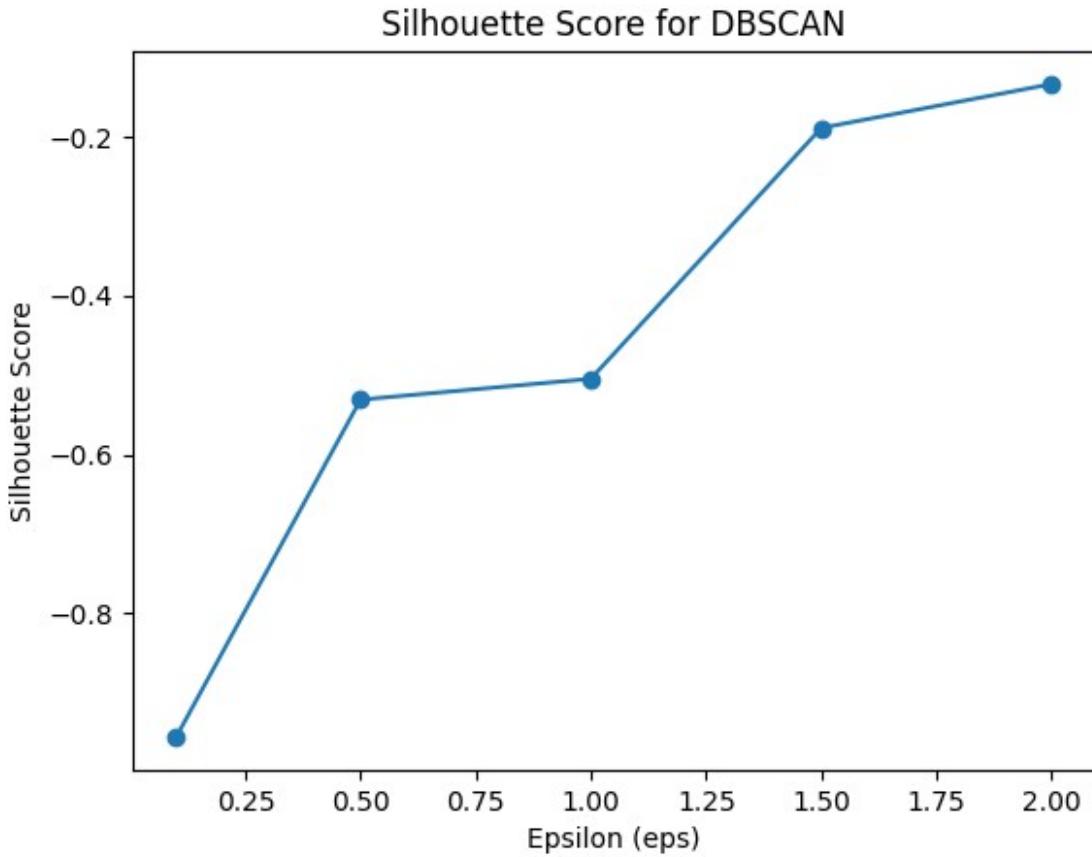
# Assuming 'data' is your feature matrix
silhouette_scores = []
eps_values = [0.1, 0.5, 1.0, 1.5, 2.0] # Example values for epsilon

for eps in eps_values:
    dbscan = DBSCAN(eps=eps)
    dbscan.fit(df_ppd_subset)
    if len(set(dbscan.labels_)) > 1: # Silhouette score requires at least 2 clusters
        silhouette_scores.append(silhouette_score(df_ppd_subset,
dbscan.labels_))
    else:
        silhouette_scores.append(-1) # Silhouette score is undefined for 1 cluster

# Plotting the silhouette scores
plt.plot(eps_values, silhouette_scores, marker='o')
plt.title('Silhouette Score for DBSCAN')
plt.xlabel('Epsilon (eps)')

```

```
plt.ylabel('Silhouette Score')
plt.show()
```



```
time: 5min 59s (started: 2024-03-16 15:21:31 +00:00)

import seaborn as sns
import matplotlib.pyplot as plt

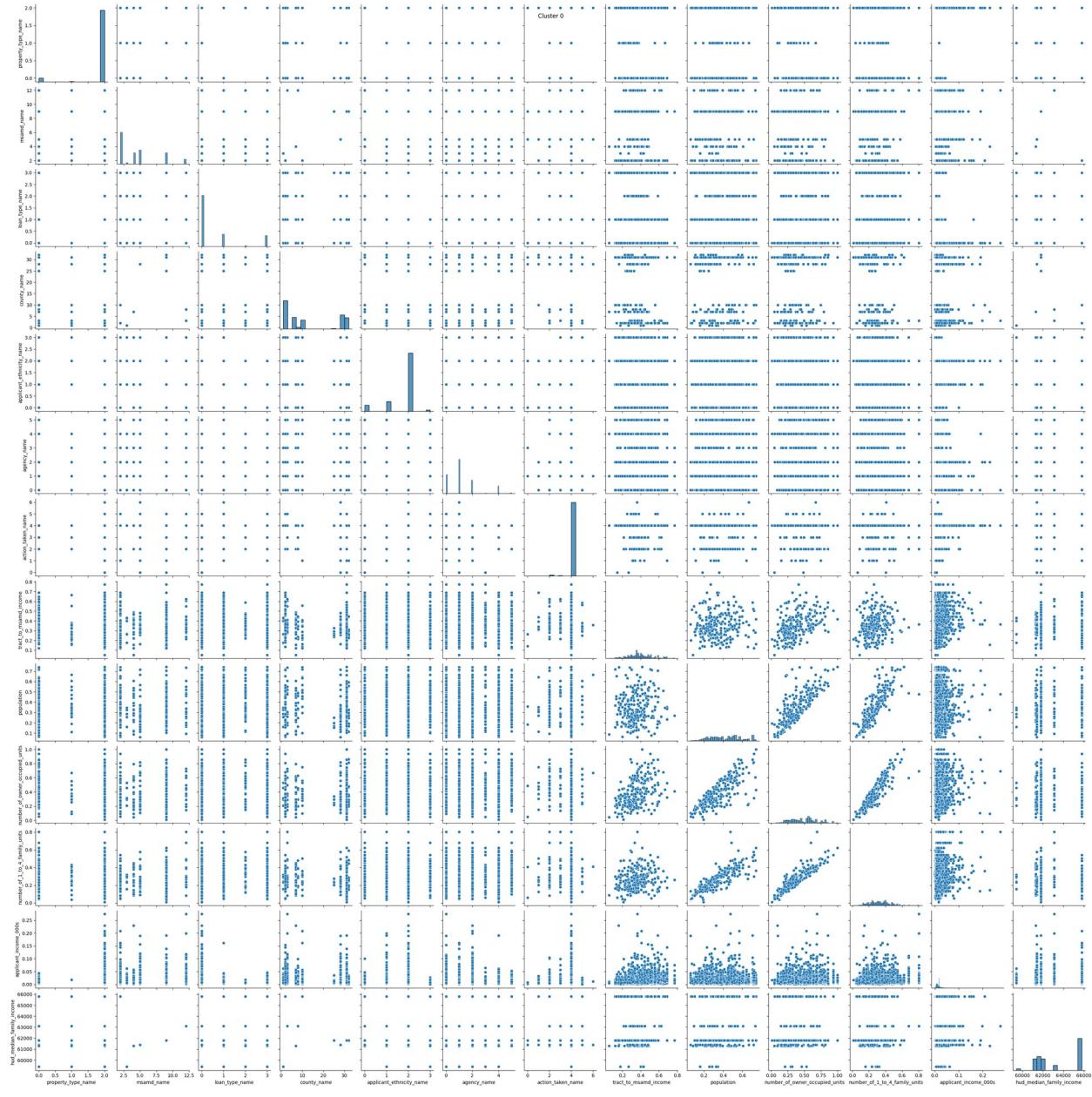
# Select top categorical and non-categorical significant variables
top_cat_variables = ['property_type_name', 'msamd_name',
'loan_type_name',
'agency_name',
'county_name', 'applicant_ethnicity_name',
'action_taken_name']
top_non_cat_variables = ['tract_to_msamd_income', 'population',
'number_of_owner_occupied_units',
'number_of_1_to_4_family_units',
'applicant_income_000s',
'hud_median_family_income']

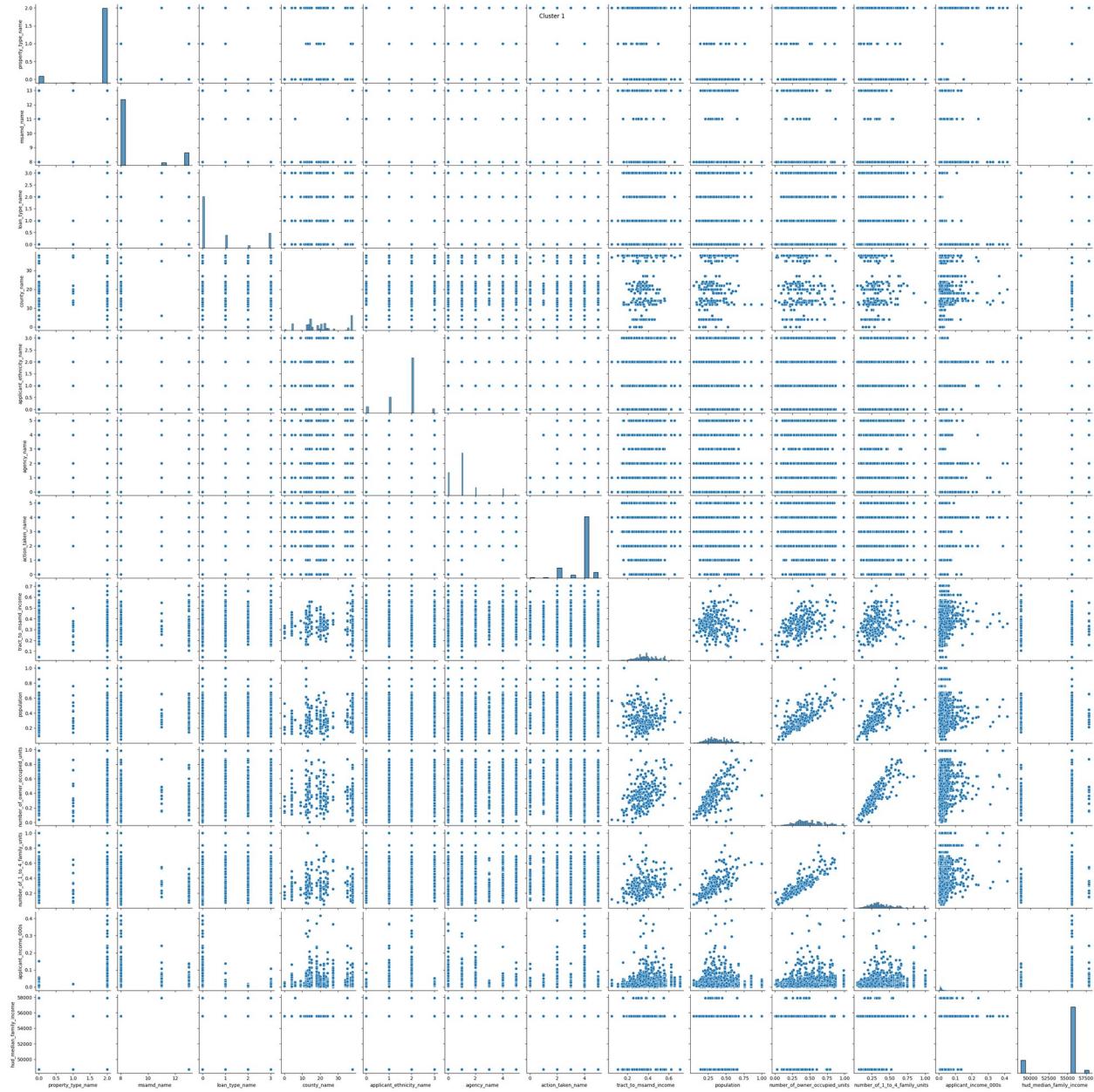
# Combine all significant variables
significant_variables = top_cat_variables + top_non_cat_variables
```

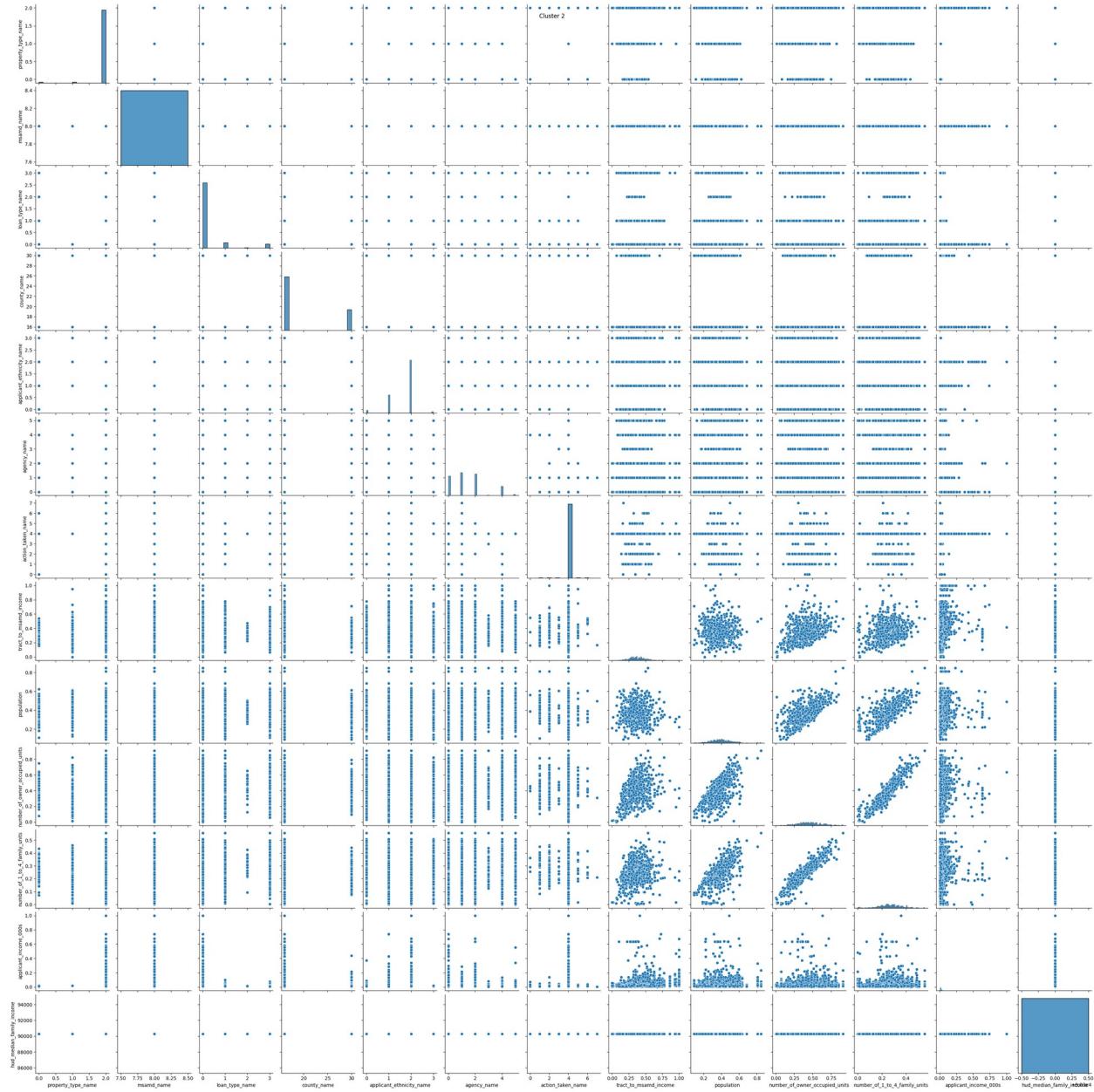
```
# Separate data for each cluster
cluster_data = {}
for cluster_label in range(5):
    cluster_data[cluster_label] =
df_ppd_subset[df_ppd_subset['Cluster_Label'] == cluster_label]

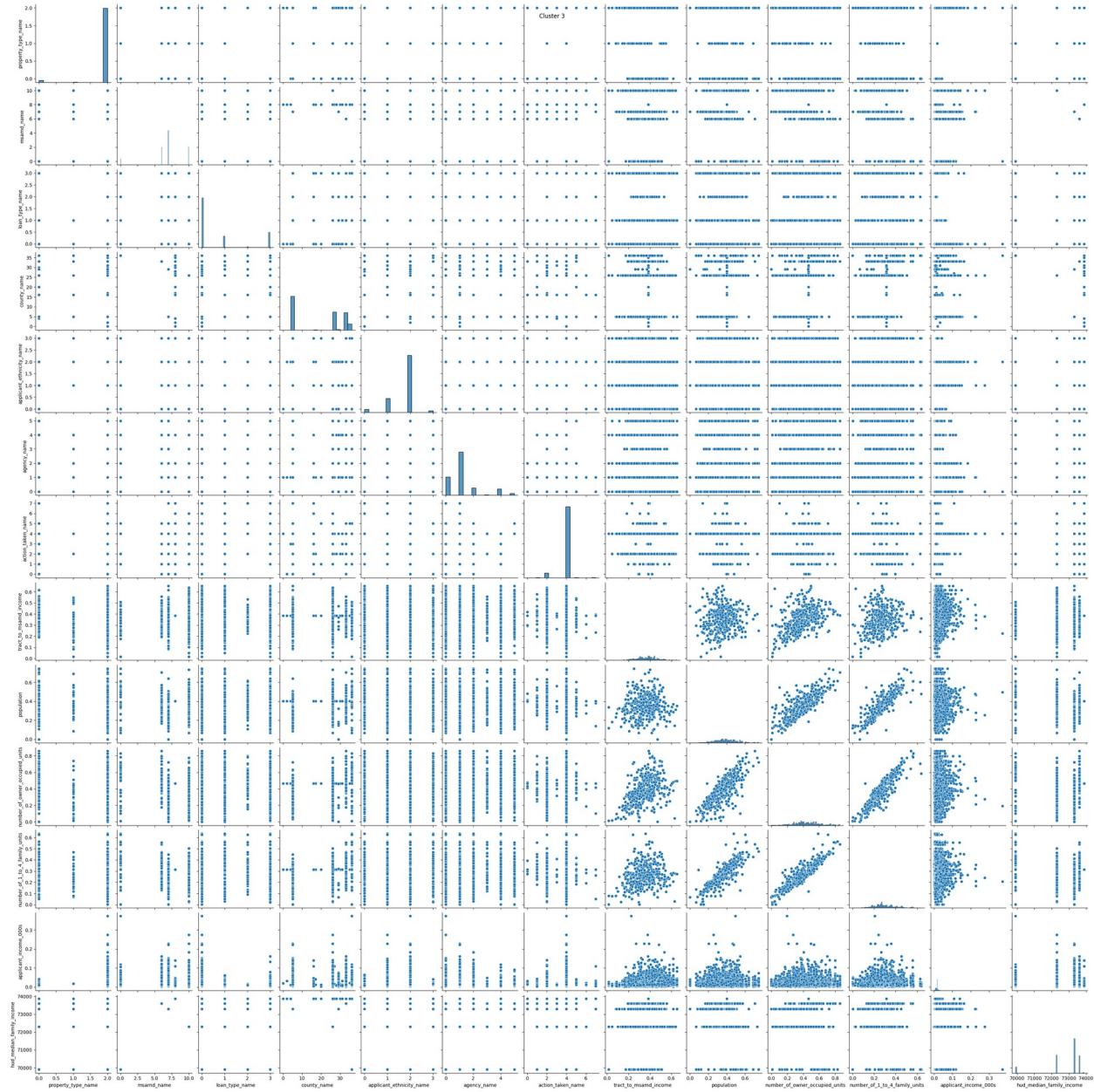
# Create pairplot for all clusters using only significant variables
for cluster_label, data in cluster_data.items():
    # Filter data for significant variables
    filtered_data = data[significant_variables]

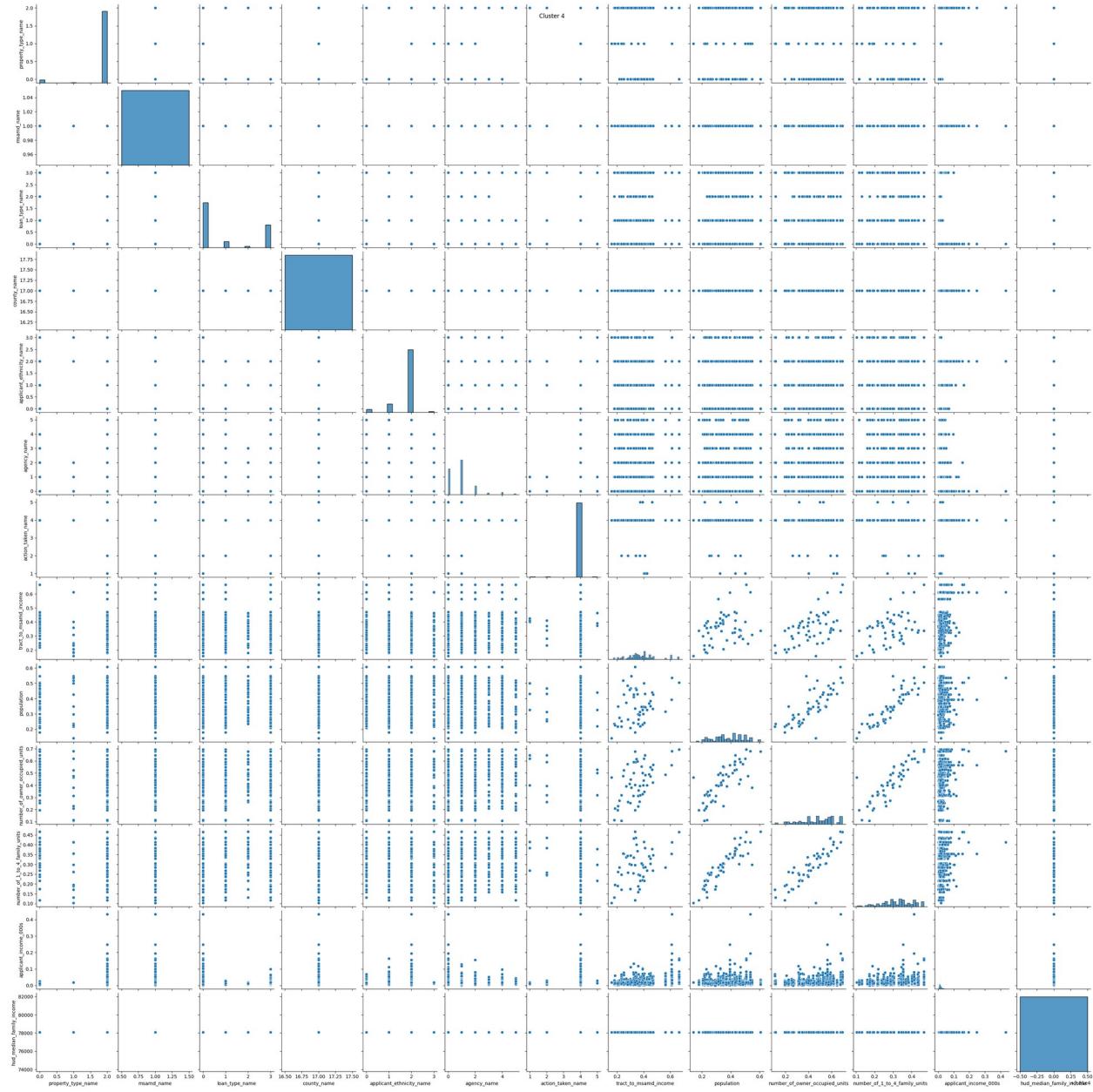
    # Plot pairplot
    sns.pairplot(filtered_data)
    plt.suptitle(f'Cluster {cluster_label}')
    plt.show()
```











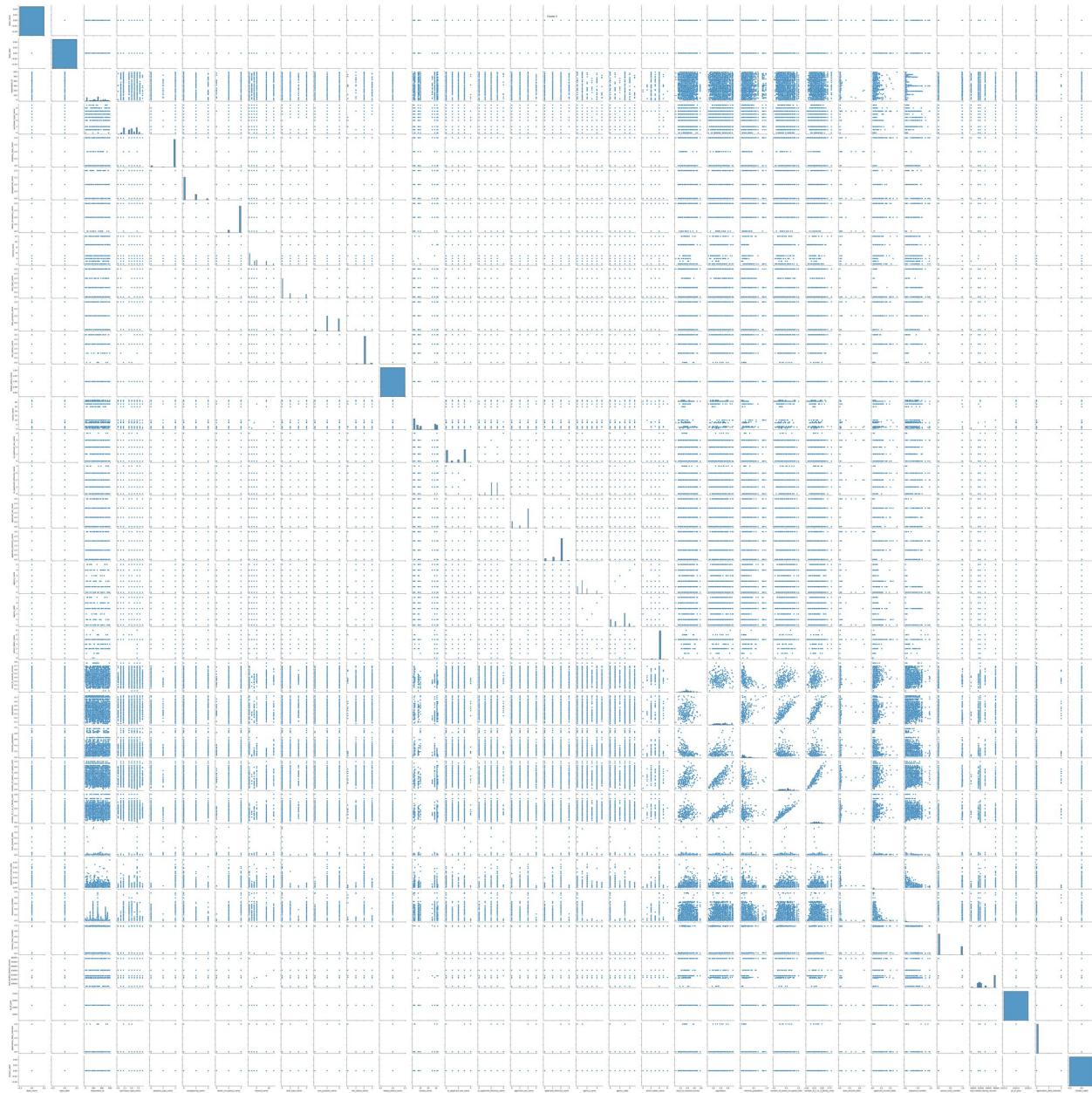
time: 9min 34s (started: 2024-03-16 15:11:56 +00:00)

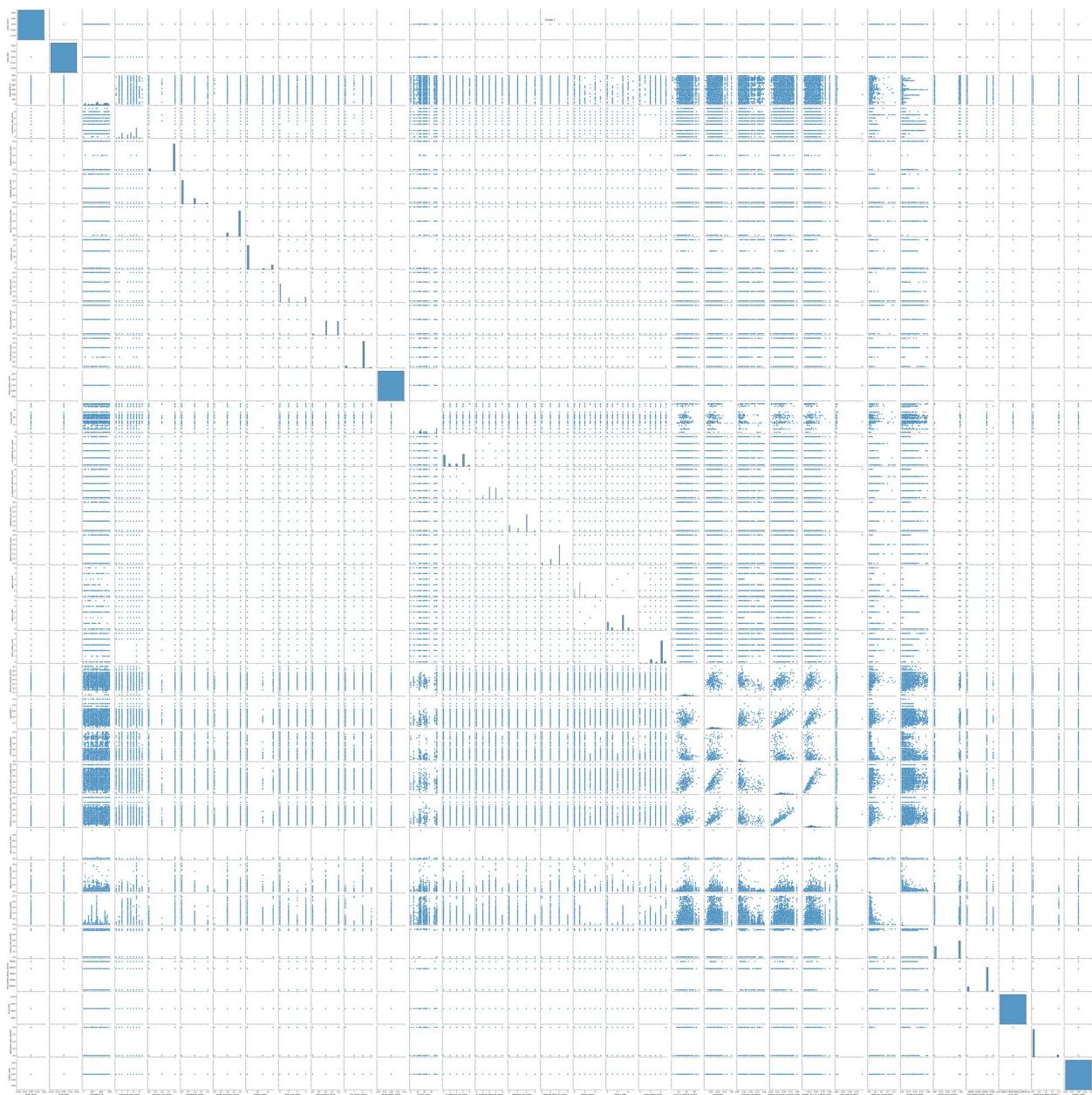
..

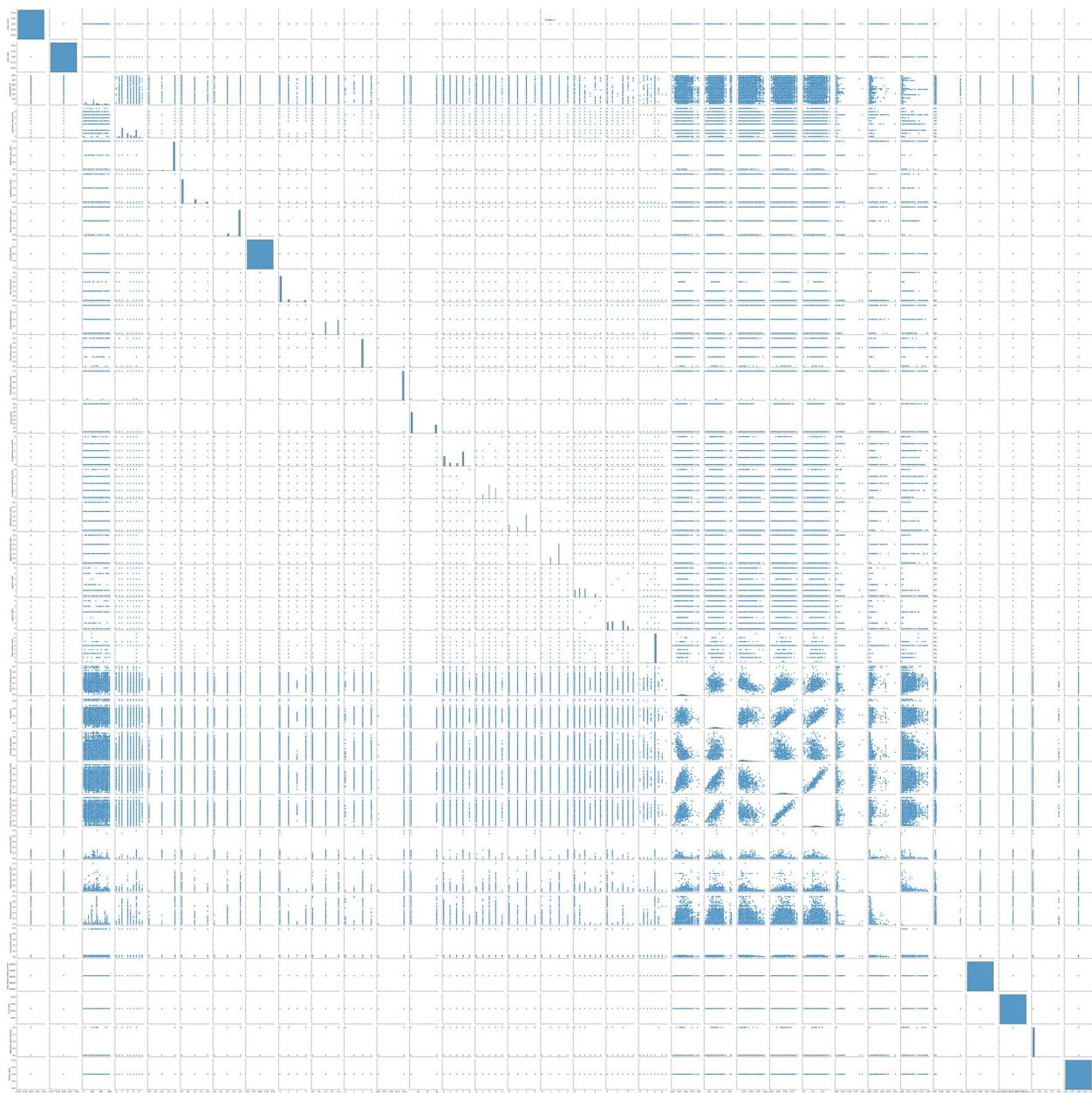
```
import seaborn as sns
import matplotlib.pyplot as plt

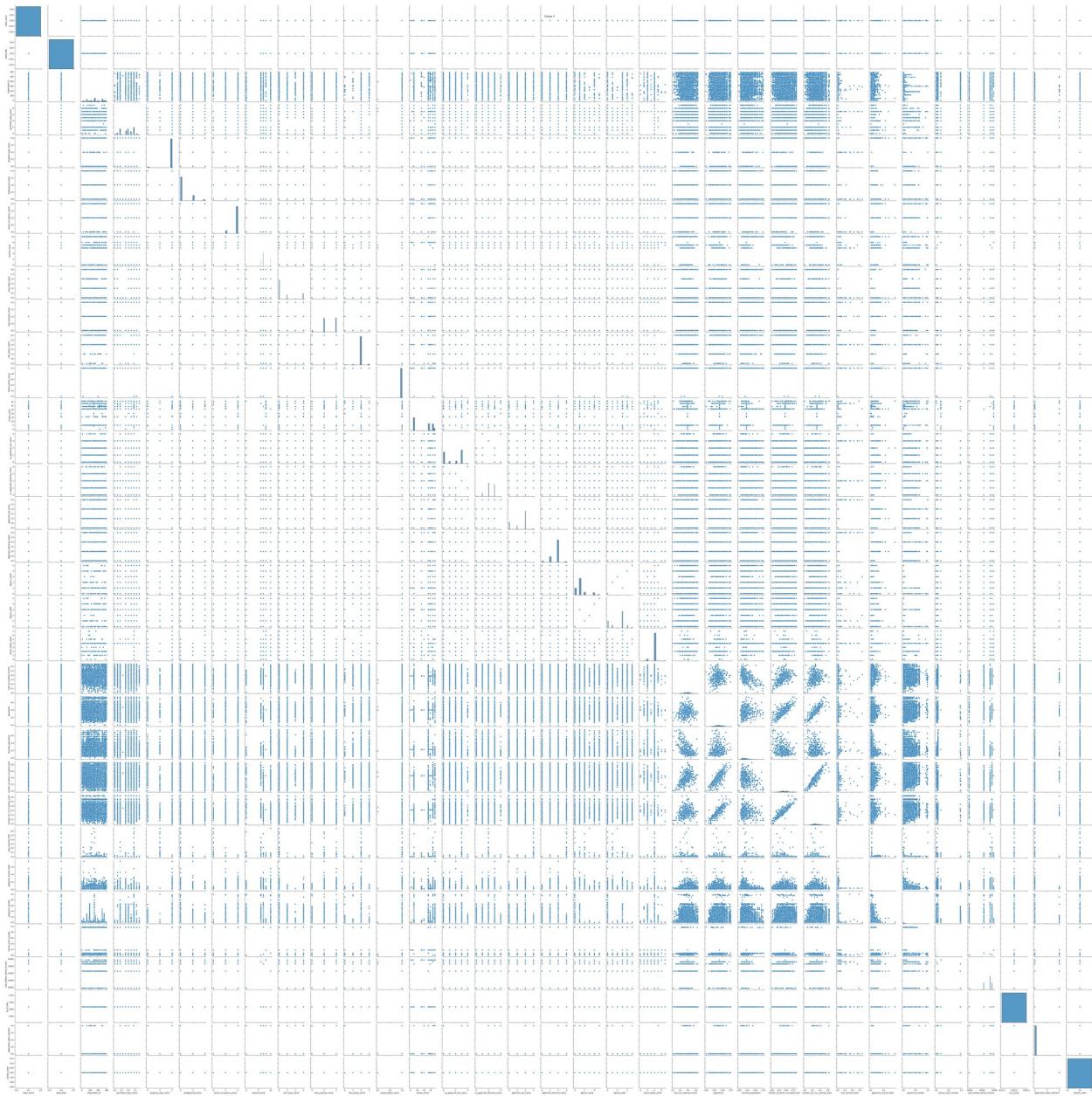
# Separate data for each cluster
cluster_data = {}
for cluster_label in range(5):
    cluster_data[cluster_label] =
df_ppd_subset[df_ppd_subset['Cluster_Label'] == cluster_label]
```

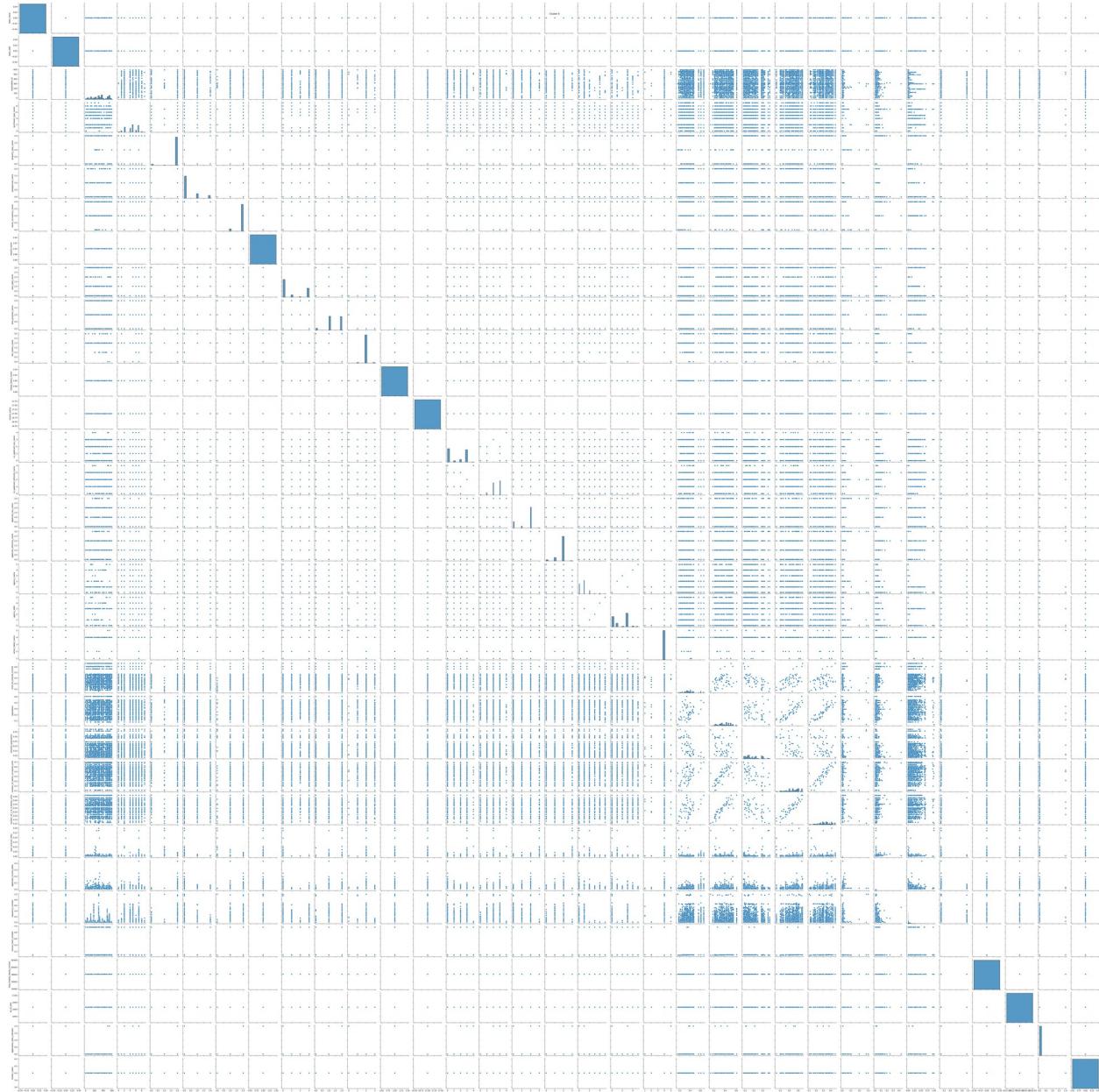
```
# Create pairplot for all clusters
for cluster_label, data in cluster_data.items():
    sns.pairplot(data)
    plt.suptitle(f'Cluster {cluster_label}')
    plt.show()
...  
...
```











time: 1h 20min 23s (started: 2024-03-16 13:50:07 +00:00)