

WASHINGTON STATE HOME LOANS 2016

Data Wrangling Report:

Washington state Home Loans dataset is obtained from Kaggle in CSV format.

<https://www.kaggle.com/miker400/washington-state-home-mortgage-hdma2016>

1. Dataset details –

Dataset contains following 47 columns and 466,566 records

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_race_name_5
co_applicant_race_name_4
co_applicant_race_name_3
co_applicant_race_name_2
co_applicant_race_name_1
co_applicant_ethnicity_name
census_tract_number
as_of_year
application_date_indicator
applicant_sex_name
applicant_race_name_5
applicant_race_name_4

applicant_race_name_3
applicant_race_name_2
applicant_race_name_1
applicant_ethnicity_name
agency_name
agency_abbr
action_taken_name

2. Data Cleaning Methods

All the columns in the dataset are imported as object type. After inspecting data following observations are found

- Some columns contain very few non-null values
- All the columns have object data type
- Some columns missing (NaN) value

2.1 Drop Columns

Columns with few data elements are dropped. Following Data columns are dropped.

denial_reason_name_3	1246 non-null object
denial_reason_name_2	6746 non-null object
co_applicant_race_name_5	14 non-null object
co_applicant_race_name_4	21 non-null object
co_applicant_race_name_3	105 non-null object
co_applicant_race_name_2	1862 non-null object
applicant_race_name_5	46 non-null object
applicant_race_name_4	68 non-null object
applicant_race_name_3	297 non-null object
applicant_race_name_2	4478 non-null object

2.2 Column Data Types

Data type of following columns are changed based on the data present in columns.

Column Name	Changed Data Type
rate_spread	float
population	float
minority_population	float
number_of_owner_occupied_units	float
number_of_1_to_4_family_units	float
loan_amount_000s	float
hud_median_family_income	float
applicant_income_000s	float
co_applicant_sex_name	Category
application_date_indicator	int
applicant_sex_name	Category

2.3 Missing Data

Numeric columns with missing data (NaN value) are substituted with median value.

Column Name
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

2.4 Outliers

Outliers from the following columns are filtered

Column Name
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

2.5 Duplicate Rows

Remove duplicate rows if any.