

# WASHINGTON STATE HOME LOANS 2016

## 1 Problem Statement

Factors affecting loan approval decisions?  
Is there any area, demographic or gender bias?  
Current trends?

## 2 Data Wrangling Report:

### 2.1 Dataset source

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

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

### 2.2 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

## 2.3 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
  - denial\_reason\_name\_3 1246 non-null object
  - denial\_reason\_name\_2 6746 non-null object
  - denial\_reason\_name\_1 34499 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
- Some columns missing (NaN) value
- Duplicate columns
  - “applicant\_race\_name\_1” duplicate of “applicant\_ethnicity\_name”
  - “co\_applicant\_race\_name\_1” duplicate of “co\_applicant\_ethnicity\_name”
  - “agency\_abbr” duplicate of “agency\_name”
  - “state\_abbr” duplicate of “state\_name”
- Columns with no significant information for statistical analysis
  - sequence\_number
  - respondent\_id

## 2.4 Drop Columns

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

denial_reason_name_3
denial_reason_name_2
denial_reason_name_1
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
applicant_race_name_5
applicant_race_name_4
applicant_race_name_3
applicant_race_name_2
applicant_race_name_1
agency_abbr
state_abbr
sequence_number
respondent_id

## 2.5 Column Data Types

Data type of all the columns with datatype “object” is changed to “category”.

## 2.6 Missing Data

### 2.6.1 Drop missing rows

Rows in column “**tract\_to\_msamd\_income**” with NaN are dropped, since the missing values are less than 20% of the data. Dropping of NaN values from “**tract\_to\_msamd\_income**”, also eliminated missing values from “**population**”, “**minority\_population**”, “**number\_of\_1\_to\_4\_family\_units**”, “**loan\_amount\_000s**” and “**hud\_median\_family\_income**”.

### 2.6.2 Impute missing values

Impute “number\_of\_owner\_occupied\_units” and “applicant\_income\_000s” with median value.

Substitute missing values in msamd\_name with value based on value in census\_tract\_number

## 2.7 Outliers

Outliers are visually inspected with box plots. Outliers from the following columns are filtered

Column Name
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

## 2.8 Duplicate Rows

Remove duplicate rows if any.

## 2.9 Save data

Save cleaned data to csv for further analysis.

## 3 Exploratory Data Analysis

### 3.1 Summary of Findings

#### 3.1.1 Trend Analysis

Seattle-Bellevue-Everett Metropolitan area trends-

1. Has largest population and highest income
2. Most number of applications
3. Among the lowest 1 to 4 family homes.
4. Most of the property are owner occupied

General trends –

1. Most of the loans applications are either for Home purchase or Refinancing
2. Most of the loan applications are conventional loans
3. Most of the loan applicants are Male

Action taken trends –

1. Most of the loan applications are accepted, there seems to be a good acceptance rate.
2. Loan applicants in each action taken category is “Not Hispanic or Latino”
3. Loan applicants in each action taken category is “Male”
4. Seattle-Bellevue-Everett Metropolitan area has the highest loan approval this could because of highest number of loan applications.
5. Most of the loans applications accepted are either for Home purchase or Refinancing
6. Conventional loans types are most accepted.

Distribution of loan amounts are between 150K to 300K across all the Metropolitan areas, loan purpose and action taken.

Please see a detailed report at this link.

[https://github.com/khmsa/Springboard/blob/master/Capstone\\_Project1/Data%20Story/WashingtonHomeLoans%20Data%20Story.ipynb](https://github.com/khmsa/Springboard/blob/master/Capstone_Project1/Data%20Story/WashingtonHomeLoans%20Data%20Story.ipynb)

#### 3.1.2 Gender Bias

A hypothesis test was conducted to determine Gender bias.

H0:  $p = 0$  There is no gender bias.

H1:  $p \neq 0$  There is a gender bias.

A t-statistic and P-value was calculated as per of hypothesis test and determined that there is a gender bias in loan approval.

#### 3.1.3 Racial Bias

Same procedure followed for determining gender bias was followed for determining racial bias.

A Hypothesis test was defined to evaluate if there is racial bias between ‘Not Hispanic or Latino’ and ‘Hispanic or Latino’.

H0:  $p = 0$  There is no gender bias.

H1:  $p \neq 0$  There is a gender bias.

After calculating t-statistic and p-value it was determined that there is no racial bias in loan approval.

### 3.1.4 Correlation Between Loan Amount and Applicant Income

Pearson correlation coefficient was calculated using `np.corrcoef()` and a hypothesis test was calculated to test the correlation coefficient.

H0:  $\rho = 0$  There is no correlation loan amount and applicant income.

H1:  $\rho \neq 0$  There is a correlation loan amount and applicant income.

A t-statistic, confidence interval and P-value was calculated and determined that there was a low chance of getting the observed correlation coefficient.

**From the statistical calculation there is no significance between loan amount and applicant income however scatter plot does show some positive correlation between loan amount and applicant income.**

Please see a detailed statistical report at this link.

[https://github.com/khmsha/Springboard/blob/master/Capstone\\_Project1/EDA/WashingtonHomeLoans%20-%20EDA.ipynb](https://github.com/khmsha/Springboard/blob/master/Capstone_Project1/EDA/WashingtonHomeLoans%20-%20EDA.ipynb)