

# Documentation for the Urban Institute's Census Tract–Level Summary Files of the Home Mortgage Disclosure Act (HMDA)

The Home Mortgage Disclosure Act (HMDA) is a federal act that requires mortgage lenders to keep records of loans and lending practices and submit data to regulatory authorities. HMDA reporting allows government regulators to analyze information on mortgage loans and mortgage lending trends. HMDA data are used to understand credit accessibility, fair lending, and the mortgage market. For example, Urban researchers have used HMDA data to analyze the [state of the mortgage market](#) and understand how well [GSEs serve minority borrowers](#). HMDA data are also an important research tool understanding housing market dynamics and neighborhood change. The data are of special interest to the [National Neighborhood Indicators Partnership](#) (NNIP) because the data includes information about the location (census tract) of the property for each individual loan application. Individual loan application data can be aggregated to the census-tract level to provide insights on the demographic and economic characteristics of people purchasing homes.

These data files are very large and cumbersome. The 2018 file has over 15 million observations and is 5.7 gigabytes. To make these data accessible to more analysts, the Urban Institute has published a public use census-tract level files with selected indicators that focus on monitoring neighborhood change. This is a brief description of the source and data provided by Urban. If we see field demand for this data, we will continue to update as additional years of HMDA are released.

## Source data

[HMDA's Snapshot National Loan Level Dataset](#) contains the applications for a calendar year as of a fixed date for all HMDA reporters, as modified by the Consumer Financial Protection Bureau to protect applicant and borrower privacy. For example, the 2018 Snapshot National Loan Level Dataset reflects all the loan applications reported for the calendar year 2018 as of August 7th, 2019. For more information on the source data, see the [full documentation from the Consumer Financial Protection Bureau](#).

## Notes on definitions and creation of variables

The indicators included in this dataset are restricted to loans that meet the following criteria.

1. Loan completed the origination process (action\_taken=1 - Loan originated)
2. Loan is for a home purchase (loan\_purpose=1 – Home Purchase)
3. Loan is first lien (lien\_status=1 – Secured by a first lien)
4. Home is single family or 1-4 unit building (derived\_dwelling\_category=Single Family (1-4 Units):Site-Built or Manufactured)
5. Home is owner occupied (occupancy\_type=1 – Principal Residence)

Within this universe, the Urban Institute uses a three-step process to create HMDA neighborhood variables. First, we use information from individual loan applications to flag characteristics of that application based on several key measures described below. Second, we combine those characteristic flags to indicate whether a loan application meets all the criteria for that indicator. And third, loan applications that meet the criteria are counted for each census tract.

Counts are then created for the number of loans:

1. By race/ethnicity
2. By income level (relative to surrounding geographic area)
3. By race/ethnicity AND relative income

The median loan amount and median borrower income are also calculated using the same universe as described above.

## Geography

This data set is published at the census-tract level. Census tract codes follow the format: 2-digit state FIPS code, 3-digit county 6-digit FIPS code, census tract FIPS code (ie: sssccttttt).

However, approximately one percent of records on the Snapshot National Loan Level Dataset are missing census-tract information. For that reason, attempting to aggregate all of the census tracts in a county or state will not add up to the total number of records in that geography.

Rather than drop these records, we created special census tract codes for them to facilitate accurate summaries of larger geographies. For partial geographic matches we use the following format:

- If a record has a known county it is assigned: scccc000000.
- If a record only has a known state, it is assigned: ss000000000.
- If a record has no geography details at all, it is assigned 99999999999.

## Race/Ethnicity

To create the indicators of race ethnicity, we evaluate the race and ethnicity variables for both the applicant and co-applicant. A record may have up to five race/ethnicities each for the applicant and co-applicant. In the source data, the field series we use are applicant\_race\_1 - applicant\_race\_5; applicant\_ethnicity\_1 - applicant\_ethnicity\_5; co\_applicant\_race\_1 - co\_applicant\_race\_5; and co\_applicant\_ethnicity\_1 - co\_applicant\_ethnicity\_5. Note, the source data file includes two other variable series about race on the data file that we do not use: derived race (Derived\_race) and race based on observation (applicant\_race\_observed).

The following steps are used to create the “Household race/ethnicity” indicator:

First, we determine a race/ethnicity category for the applicant:

- IF applicant\_ethnicity is Hispanic then they are coded: Hispanic
- ELSE IF their race is missing or is marked as “Not applicable” or “Not provided” then they are coded as: Race not available
- ELSE IF the applicant indicated identifying with multiple races (using fields applicant\_race\_2 - applicant\_race\_5) then they are coded as: Multiple races

- ELSE they are coded as the race indicated in applicant\_race\_1

Second, the same logic is followed to determine a race/ethnicity category for the co-applicant (if there is a co-applicant).

Third, the values determined in steps one and two are used to determine the household race/ethnicity category for each record:

- IF there is a co-applicant AND both the applicant and co-applicant race are available AND they are not the same, then the record is coded as: Multiple races
- ELSE IF the applicant and co-applicant's race/ethnicity are the same OR there is no co-applicant, then the record is coded as the applicant's race/ethnicity

## Relative Income

We provide indicators based relative income categories compared to the median area income for metropolitan areas (for counties in those areas) or counties (for non-metropolitan areas). Using the relative income level is helpful in comparing across areas of the country with different housing costs and other costs of living. Income limit data are published by the Department of Housing and Urban Development (HUD) and can be [acquired from the HUD USER website](#). We merged the dataset with the Snapshot National Loan Level Dataset using county identifier.

We compare the borrower's income to the Area Median Income to determine whether the applicant falls into the category of:

- Very low income (0-50 percent of Area Median Income)
- Low income (50.1-80 percent of Area Median Income)
- Moderate income (80.1-120 percent of Area Median Income)
- High income (120.1 percent+ of Area Median Income)

## Citation and license

These data are published under an ODC-BY 1.0 license. You are free to share these data, produce works from these data, and adapt the files as long as you attribute any public use of the database or works produced from the database. The citation should be:

Urban Institute. 2020. Census Tract–Level Home Mortgage Disclosure Act Indicators for Neighborhood Change. Accessible from <https://datacatalog.urban.org/dataset/home-mortgage-disclosure-act-hmda-neighborhood-summary-files-census-tract-level>

Data originally sourced from the Consumer Financial Protection Bureau, developed at the Urban Institute, and made available under the ODC-BY 1.0 Attribution License. For full details on the license, please see <https://opendatacommons.org/licenses/by/summary/index.html>. Consumer Financial Protection Bureau (2018) [computer file]. Washington, DC: Snapshot National Loan Level Dataset, accessed May 1, 2020 at <https://ffiec.cfpb.gov/data-publication/>