

Link to dataset:

<https://www.kaggle.com/datasets/saravananselvamohan/freddie-mac-singlefamily-loanlevel-dataset>

This data set contains information on single family loans purchased by Freddie Mac from 1999 to 2018. There are 27 million rows in this dataset, and that could be too much for R to handle so I will just take the first 3,000 mortgages in 1999 for the analysis.

Uploading this to kaggle three years ago, Freddie Mac is a government sponsored liquidity provider for the US residential mortgage market. Founded in 1970, it was a part of a push by the US government to turn as many Americans into homeowners as possible.

The csv file has 27 variables and are listed below, although the loan sequence number won't serve any analytical purpose. A potential source of error in the data may be the debt to income ratio, because in the first half of this time interval, income verification by the originator was notoriously lax.

The common perception today is that lending standards have tightened since the 2008 financial crash, but is this true? Did delinquent loans in the early 2000's have similar traits to loans in the 2010s? The data stops at 2018, but are we due for another housing crash?

- [1] "CREDIT\_SCORE"
- [2] "FIRST\_PAYMENT\_DATE"
- [3] "FIRST\_TIME\_HOMEBUYER\_FLAG"
- [4] "MATURITY\_DATE"
- [5] "METROPOLITAN\_STATISTICAL\_AREA"
- [6] "MORTGAGE\_INSURANCE\_PERCENTAGE"
- [7] "NUMBER\_OF\_UNITS"
- [8] "OCCUPANCY\_STATUS"
- [9] "ORIGINAL\_COMBINED\_LOAN\_TO\_VALUE"
- [10] "ORIGINAL\_DEBT\_TO\_INCOME\_RATIO"
- [11] "ORIGINAL\_UPB"
- [12] "ORIGINAL\_LOAN\_TO\_VALUE"
- [13] "ORIGINAL\_INTEREST\_RATE"
- [14] "CHANNEL"
- [15] "PREPAYMENT\_PENALTY\_MORTGAGE\_FLAG"
- [16] "PRODUCT\_TYPE"
- [17] "PROPERTY\_STATE"
- [18] "PROPERTY\_TYPE"
- [19] "POSTAL\_CODE"
- [20] "LOAN\_SEQUENCE\_NUMBER"
- [21] "LOAN\_PURPOSE"
- [22] "ORIGINAL\_LOAN\_TERM"
- [23] "NUMBER\_OF\_BORROWERS"
- [24] "SELLER\_NAME"
- [25] "SERVICER\_NAME"
- [26] "PREPAID"
- [27] "DELINQUENT"