

# **Lab 2: Who Pays More?**

**Examining Rate Spread Variation by Race in Mortgage Loans**  
**[https://github.com/mids-w203/lab\\_2\\_mean\\_machine/tree/main](https://github.com/mids-w203/lab_2_mean_machine/tree/main)**

Ender Ricart | Kelvin Yang | Alex Kim

April 17, 2025

## 1. Introduction

In 2023, CNN published an investigative report alleging racial disparities in the Navy Federal Credit Union’s (NFCU) mortgage lending practices. (Tolan and Marsh (2023)) Using NFCU’s loan-level data, CNN’s analysis showed that “Black/African American” borrowers experienced significantly lower approval rates than “White” borrowers. As a fair lending analytics team at the Consumer Financial Protection Bureau (CFPB), we have conducted a descriptive analysis that builds on CNN’s findings to explore the relationships between a successful borrower’s race or ethnicity and mortgage pricing outcomes for NFCU. Even slight differences in the mortgage rate has a long-term financial consequences for borrowers.

## 2. Description of the Data Source

We obtained a prefiltered subset of the Home Mortgage Disclosure Act (HMDA) Modified Loan Application Register (LAR) data directly from the CFPB’s website. It consists of loan applications for home purchases submitted to the NFCU in 2023 and contains loan-level information for mortgages reported by financial institutions and individual-level mortgage application records, including applicant demographic information, loan details, and loan outcomes.

We looked at individual mortgage loan applications. We only included applications that were approved, were conventional conforming loans, and were for homes where people live (primary residences). By focusing only on these types of loans, we could better compare how different borrowers were treated under the same lending rules. By excluding denied applications, selection bias has been introduced that prevents us from analysing racial or ethnic disparities in loan successes as CNN has done. In the full dataset, approximately 51% of applicants identified as White and 23% as Black or African American. However, after applying the sample restrictions, White applicants represent over 65% of the sample, while Black applicants account for only 17% (see Figure 2). This shift suggests that approval rates vary by race, potentially contributing to the increased disparity observed in the filtered data. This limitation means our findings may understate disparities that are more evident in the full applicant pool. Therefore, our results should be interpreted as describing rate differences conditional on approval, not disparities in access to credit overall. This limitation is relevant because our focus on previous reporting from CNN already shows racial approval disparities at NFCU.

## 3. Data Wrangling and Operationalization

### Data Wrangling

Our data-wrangling process prioritized analytic clarity and supported independent and identically distributed (i.i.d.) assumptions. After filtering our dataset for approved loans, we then narrowed our set more to include only conventional and conforming loans. These loan types follow standardized underwriting guidelines established by Fannie Mae and Freddie Mac. In doing so we control for variability in how lenders determine a borrower’s creditworthiness

and that we are comparing loans that are evaluated similarly. These standards provide consistent treatment of borrowers and identically distributed observations. We also excluded introductory-rate loans, as these types of loans have different pricing and may not follow standardized underwriting guidelines. We eliminated records with missing values for critical predictor variables such as rate spread, loan-to-value ratio, income, and debt-to-income ratio.

As a result, we removed 22,934 applications that did not meet criteria for standardized pricing analysis. This included 16,322 non-conventional loans (e.g., Federal Housing Administration), 1,068 non-conforming loans, 1,514 loans for non-primary residences, 3,718 with introductory rates, 116 second trust deed loans, and 196 records missing key variables. Our final dataset contained 6,527 observations (22.2% of the original data, home purchase loans for primary residence approved by NFCU in 2023), providing a foundation for examining racial disparities in mortgage pricing.

### **Operationalization: Rate Spread as the Outcome Variable**

Rate spread is defined as the difference between annual percentage rate (APR) and average prime offer rate (APOR), a market benchmark rate. We selected rate spread as our outcome variable because it allows us to evenly assess whether similarly situated applicants receive interest rates that systematically differ from the benchmark. This accommodates daily fluctuations in the lending market, as opposed to measuring APR outright.

### **Operationalization: Race**

In our dataset, Race and Ethnicity are self-selected by applicants or visually determined and reported by Navy Federal Credit Union personnel. We operationalized race based on the following logic:

1. White race group: If either applicant or co-applicant identified as White, the application was classified as “White.” This approach follows fair lending guidance and reflects how NFCU may view joint applications.
2. Primary applicant’s self-reported race: If neither the applicant nor co-applicant is White, then we default to the primary applicant’s self-reported race. We consolidated disaggregated Asian and Pacific Islander subgroups into broader categories (e.g., classifying Asian Indian, Chinese, etc. as Asian).
3. Hispanic/Latino: We created a category for Hispanic/Latino in our operationalization of race, even though US Census defines Hispanic/Latino as an ethnic category. If self-reported race is missing but either applicant identified as Hispanic/Latino, we counted the row counted the application as Hispanic/Latino.
4. Observed race: If no self-reported race or ethnicity information was available, we used visually observed race as reported by NFCU.

We considered alternative classification approaches, including using only the primary applicant's self-reported race, to maintain consistency across all observations and avoid mixed-race categorization. We also considered assigning Hispanic ethnicity first or using only primary applicant race, but our approach better reflects lender perception and attempts to avoid misclassification as much as possible.

## Control Variables

We included several variables commonly used in underwriting and pricing models to isolate racial disparities that applicant's credit risk cannot explain:

- Debt-to-income (DTI) ratio: Reflects the borrower's ability to repay and may influence pricing tiers. Since the LAR reports exact values (within the 36%-49% range) and categorical ranges, we converted ranges to numeric midpoints to retain those observations.
- Loan-to-value (LTV) ratio: Reflects borrower equity; higher LTVs typically receive higher pricing.
- Income: Serves as a proxy for repayment capacity.
- Loan amount: Affects pricing through risk-based adjustments and eligibility thresholds.

## Splitting the data into an exploration set and a confirmation set

To minimize overfitting from sequential model decisions, we randomly split the original dataset into a 30% exploration set ( $n = 1,958$ ) and a 70% confirmation set ( $n = 4,569$ ). All transformation and modeling choices were made using only the exploration set.

## 4. Model Specification

To explain racial disparities in mortgage pricing, we estimated two linear regression models using rate spread as the dependent variable. The simple model uses race as the sole predictor variable, and the expanded model uses control variables, including income, debt-to-income ratio (DTI), loan-to-value ratio (LTV), and loan amount. These variables were selected because they are all common credit risk indicators that are publicly available.

```
model_simple <- lm(rate_spread ~ race_group, data=nfcu_data)
model_expanded <- lm(rate_spread ~ race_group + loan_to_value_ratio +
                     dti_clean + income + loan_amount, data=nfcu_data)
```

The race group is a categorical variable, with White as the reference category. This approach allows us to assess average differences in pricing between White applicants and other race groups. All other variables are metric. Full model specs, estimates, and robust standards errors are shown in Table 1.

## 5. Model Assumptions

### i.i.d

While we have taken steps to support the i.i.d. of data, we still expect some clustering. We limited our sample to a single year (2023) and primary residence home loans, minimizing the likelihood of multiple applications by a single borrower. However, a small number of borrowers may still appear multiple times in the data set; the dataset is anonymized, so we cannot directly identify duplicate applications from the same person. Furthermore, there may be some level of geographic clustering in the sample of approved loans. Borrowers from the same geographic areas may face similar lending conditions, leading to correlated outcomes within regions.

### Linearity and Zero Conditional Mean

When plotting the expanded model's residuals against fitted values, it shows a funnel-shaped pattern and curvature, suggesting nonlinearity and heteroscedasticity. To address this, we first applied a signed log transformation to the dependent variable, rate spread, which improved the overall linear relationship. We then examined each predictor individually. The relationships between residuals and the predictors, income, loan-to-value ratio, and loan amount, showed evidence of nonlinearity. Applying a log transformation to income and loan amount improved the linear fit and transforming the loan-to-value ratio using a second-degree polynomial addressed the curvature in its residual plot.

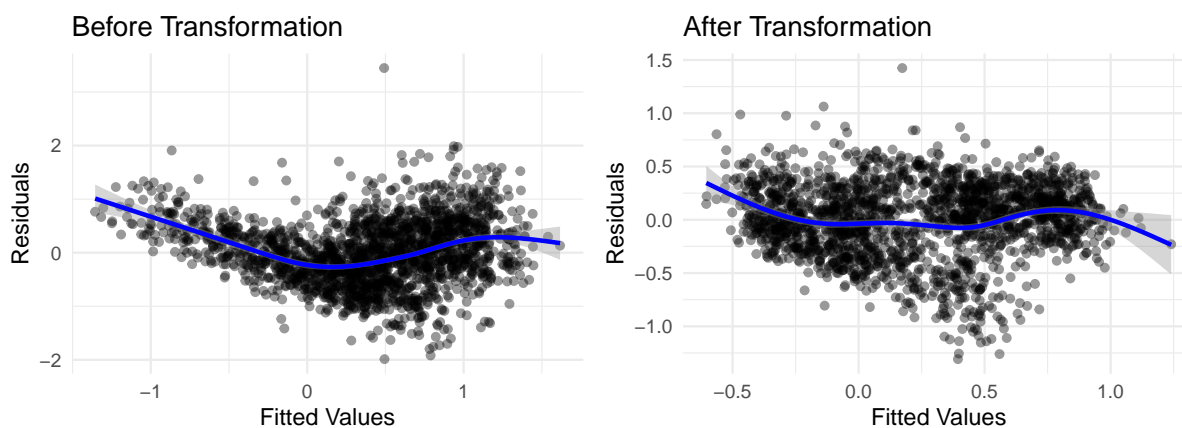


Figure 1: Residual Plots of Expanded Model: Before and Ater Transformation

After these transformations, the model's performance improved substantially: the R-squared increased from 0.45 to 0.57, and the adjusted R-squared increased accordingly, indicating better model fit and enhanced adherence to OLS assumptions. Although some mild non-linearity remains at the extremes, the transformation substantially improved model behavior, bringing it closer to satisfying the linearity and zero conditional mean assumptions.

## No Perfect Collinearity

The assumption of no perfect collinearity checks whether a predictor variable can be described completely by a combination of the others. The variance inflation factor (VIF) measures how much the variance of each regression coefficient is inflated due to multicollinearity. Our continuous variables have VIF values ranging between 1.17 and 2.04, indicating minimal concern with multicollinearity. Our categorical variable, `race_group`, has an adjusted generalized VIF value of 1.01, again suggesting minimal concern. After transformation, VIF values increased to range between 1.51 and 2.65 for continuous variables but is still not a cause for concern with multicollinearity. Overall, both models meet the assumption of no perfect collinearity.

## 6. Model Results and Interpretation

When we applied our simple model to our confirmatory data set we found that Black (+0.56%), American Indian (+0.30%), and Hispanic (+0.20%) applicants had higher average rate spreads than White applicants. In the expanded model, which includes credit risk indicators (LTV, DTI, income, and loan amount), we found a reduction in the coefficients. As shown in Figure 3, based on our expanded model the average rate spread varies by race when controlling for LTV, DTI, income, and loan amount. Even after accounting for these credit risk factors, the predicted rate spread remains meaningfully higher for most minority race groups relative to White applicants. Statistically significant gaps remain for Black (+0.19%) and American Indian (+0.21%) applicants relative to White applicants, indicating residual differences in pricing not accounted for by observable underwriting criteria.

The final model, which incorporates log and polynomial transformations to better meet OLS assumptions, improves model fit (adjusted  $R^2$  increases from 0.46 to 0.59). The estimated gap for Black applicants is +2.5% in log rate spread, but is not statistically significant. The gap for American Indian applicants remains statistically significant, with a +7.6% in log rate spread. Asian applicants are associated with a +4.6% in log rate spread compared to White applicants, a statistically significant result only in the transformed model. These patterns suggest that differences in observable applicants' credit risk indicators and nonlinear relationships account for much of the variation seen in the simpler models. However, some group-level differences remain after adjusting for these factors.

## 7. Overall Effect

Many credit risk indicators like income, credit-score, and location are correlated with race. Our descriptive analysis finds that even after controlling for some of these indicators, there continue to be small but systemic rate differentials for Black and Native American applicants. This raises concerns regarding fair lending practices, warranting further investigation. Since the publicly available LAR data does not capture all factors that influence rate pricing, additional nonpublic data can drive a deeper understanding of these disparities. We hope future analysis can build on our findings to reduce systemic housing inequities.

## Appendix

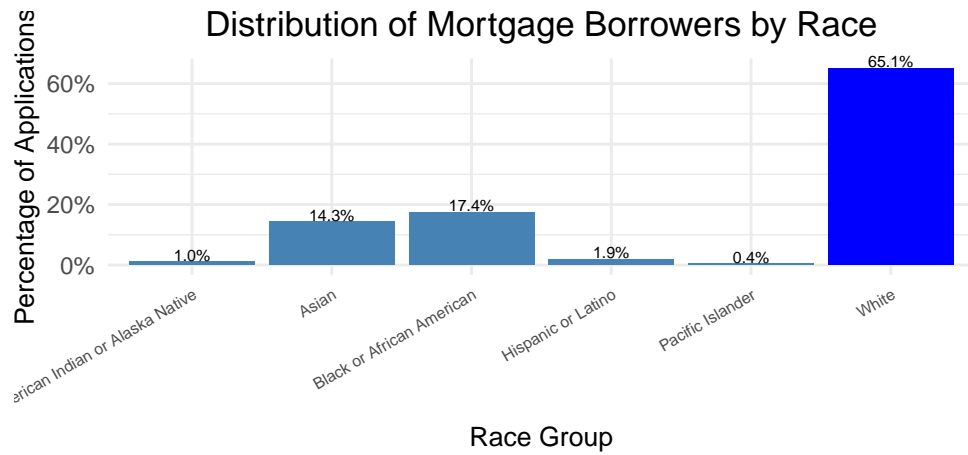


Figure 2: Distribution of Applicants by Race Group

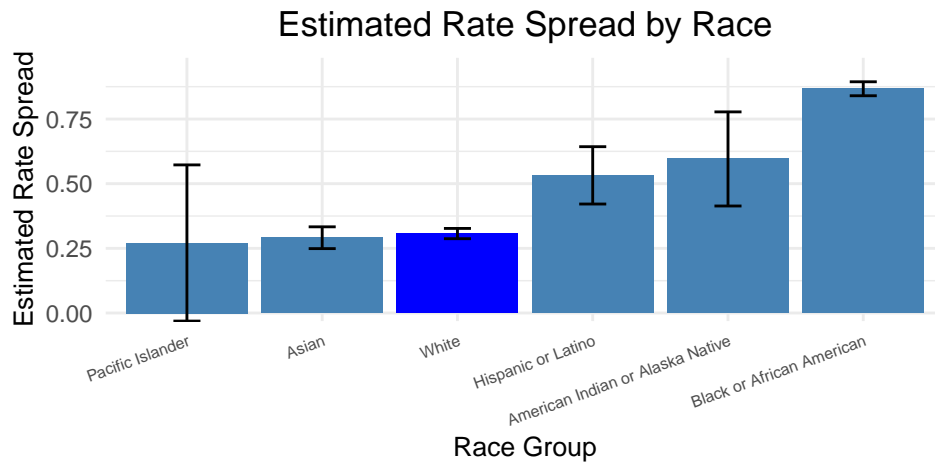


Figure 3: Average Estimated Rate Spread by Race Group with 95% CI (Adjusted for income, debt-to-loan ratio, loan-to-value ratio, and loan amount)

Table 1

	<i>Dependent variable:</i>		
	Rate Spread		Log Rate Spread
	Simple	Expanded	Transformed
	(1)	(2)	(3)
American Indian	0.29** (0.12)	0.21*** (0.08)	0.08** (0.04)
Asian	-0.02 (0.03)	-0.04 (0.03)	-0.05*** (0.02)
Black	0.56*** (0.03)	0.19*** (0.03)	0.02* (0.01)
Hispanic	0.23** (0.10)	0.06 (0.07)	-0.001 (0.04)
Pacific Islander	-0.04 (0.18)	0.12 (0.12)	0.01 (0.06)
LTV		0.03*** (0.001)	
Centered LTV			0.03*** (0.0004)
Centered LTV Squared			0.0004*** (0.0000)
DTI		0.01*** (0.001)	0.01*** (0.001)
Income		0.002*** (0.0002)	
Loan Amount		-0.0000*** (0.0000)	
Log Income			0.15*** (0.02)
Log Loan Amount			-0.25*** (0.02)
Constant	0.31*** (0.01)	-1.98*** (0.06)	2.31*** (0.14)
Observations	4,569	4,569	4,569
R <sup>2</sup>	0.07	0.46	0.60
Adjusted R <sup>2</sup>	0.07	0.46	0.60

*Note:*

\*p<0.1; \*\*p<0.05; \*\*\*p<0.01  
Robust standard errors in parentheses. White applicants are the reference group.



## References

Tolan, Audrey Ash, Casey, and Rene Marsh. 2023. “The Nation’s Largest Credit Union Rejected More Than Half Its Black Conventional Mortgage Applicants.” *CNN*.

Bartlett, Robert, Adair Morse, Richard Stanton, and Nancy Wallace. “Consumer-Lending Discrimination in the FinTech Era.” Working Paper, Berkeley Haas School of Business, 2019. <https://faculty.haas.berkeley.edu/morse/research/papers/discrim.pdf>.

Ambrose, Brent W., James Conklin, and Luis A. Lopez. “Does Borrower and Broker Race Affect the Cost of Mortgage Credit?” *The Review of Financial Studies* 34, no. 2 (2021): 790-826. <https://pure.psu.edu/en/publications/does-borrower-and-broker-race-affect-the-cost-of-mortgage-credit>

Bhutta, Neil, and Aurel Hizmo. “Do Minorities Pay More for Mortgages?” *The Review of Financial Studies* 34, no. 2 (2021): 763-789. <https://academic.oup.com/rfs/article/34/2/763/5827007>

Giacoletti, Marco, Rawley Heimer, and Edison G. Yu. “Using High-Frequency Evaluations to Estimate Discrimination: Evidence from Mortgage Loan Officers.” Federal Reserve Bank of Philadelphia Working Paper (2021). <https://www.philadelphiafed.org/-/media/frbp/assets/working-papers/2021/wp21-04.pdf>

## Reference Materials:

Consumer Financial Protection Bureau. “A Beginner’s Guide to Accessing and Using Home Mortgage Disclosure Act Data.” Washington, DC: Consumer Financial Protection Bureau, June 2022. [https://files.consumerfinance.gov/f/documents/cfpb\\_beginners-guide-accessing-using-hmda-data\\_guide\\_2022-06.pdf](https://files.consumerfinance.gov/f/documents/cfpb_beginners-guide-accessing-using-hmda-data_guide_2022-06.pdf).

CFPB’s HMDA documentation page: <https://ffiec.cfpb.gov/documentation/>

Modified LAR Schema documentation: <https://ffiec.cfpb.gov/documentation/publications/modified-lar/modified-lar-schema>