

Unbiased Prediction of Loan Repayment

Hsiao-Yu Chiang, Yara Mubarak, Jordan Fan,
Yang Shi, Bizuayehu Whitney, Florin Langer





Introduction

- How do biased vs unbiased predictors perform when predicting loan default or late repayment?
- How it can be useful :
 - Companies can use this as their main screening technique of borrowing candidates
 - Companies can use this baselines model to compare their current screening process' biases.

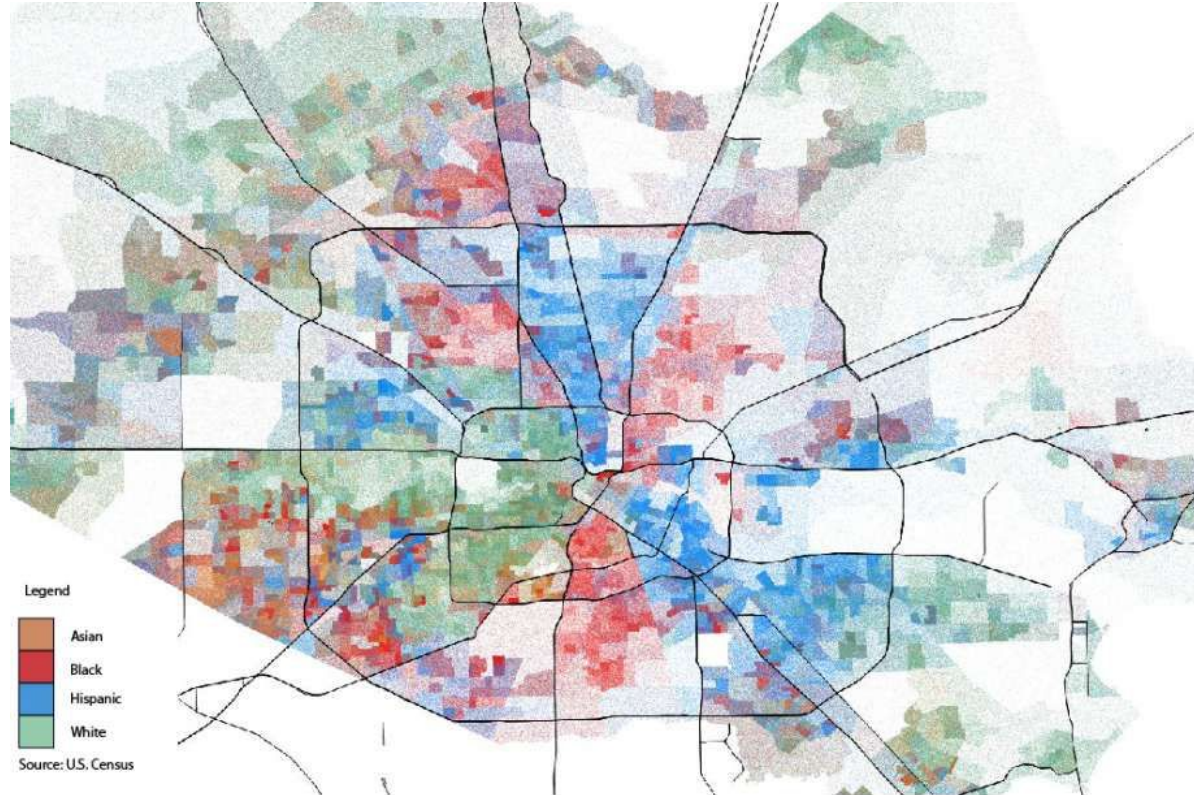


Dataset

- <https://www.kaggle.com/wendykan/lending-club-loan-data>
- 2.26 million rows, 145 columns from LendingClub's Data
- Features include loan amount requested by borrower, amount funded by investors, interest rate, occupation of borrower
- Loan status include whether loan is fully paid, whether the loan was defaulted, whether it's still currently being paid, whether loan is late by x amount of day.



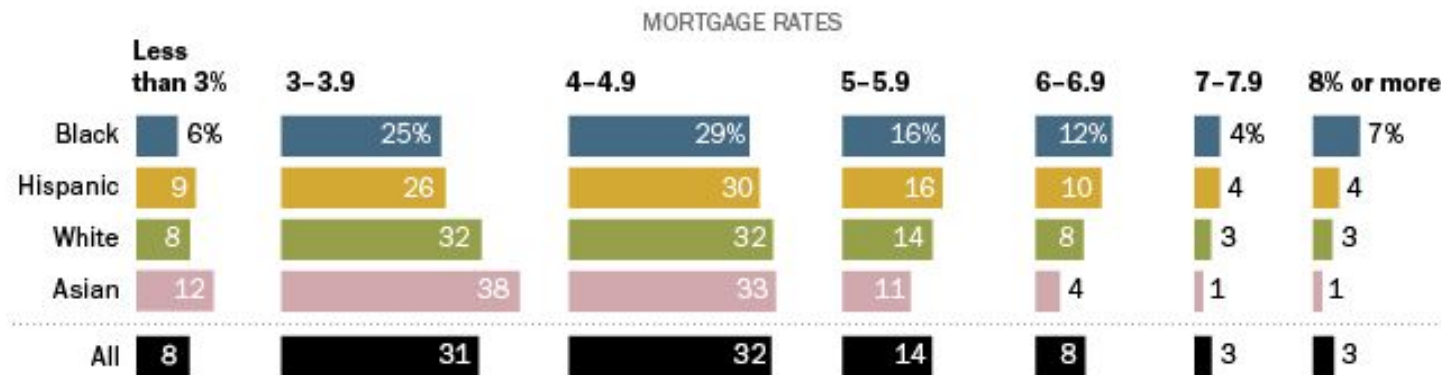
Houston Segregation By Race Example



Mortgage Rate Bias Based on Race

Blacks, Hispanics more likely to pay higher mortgage rates

Among households in 2015 with at least one regular mortgage, % of each group paying these rates

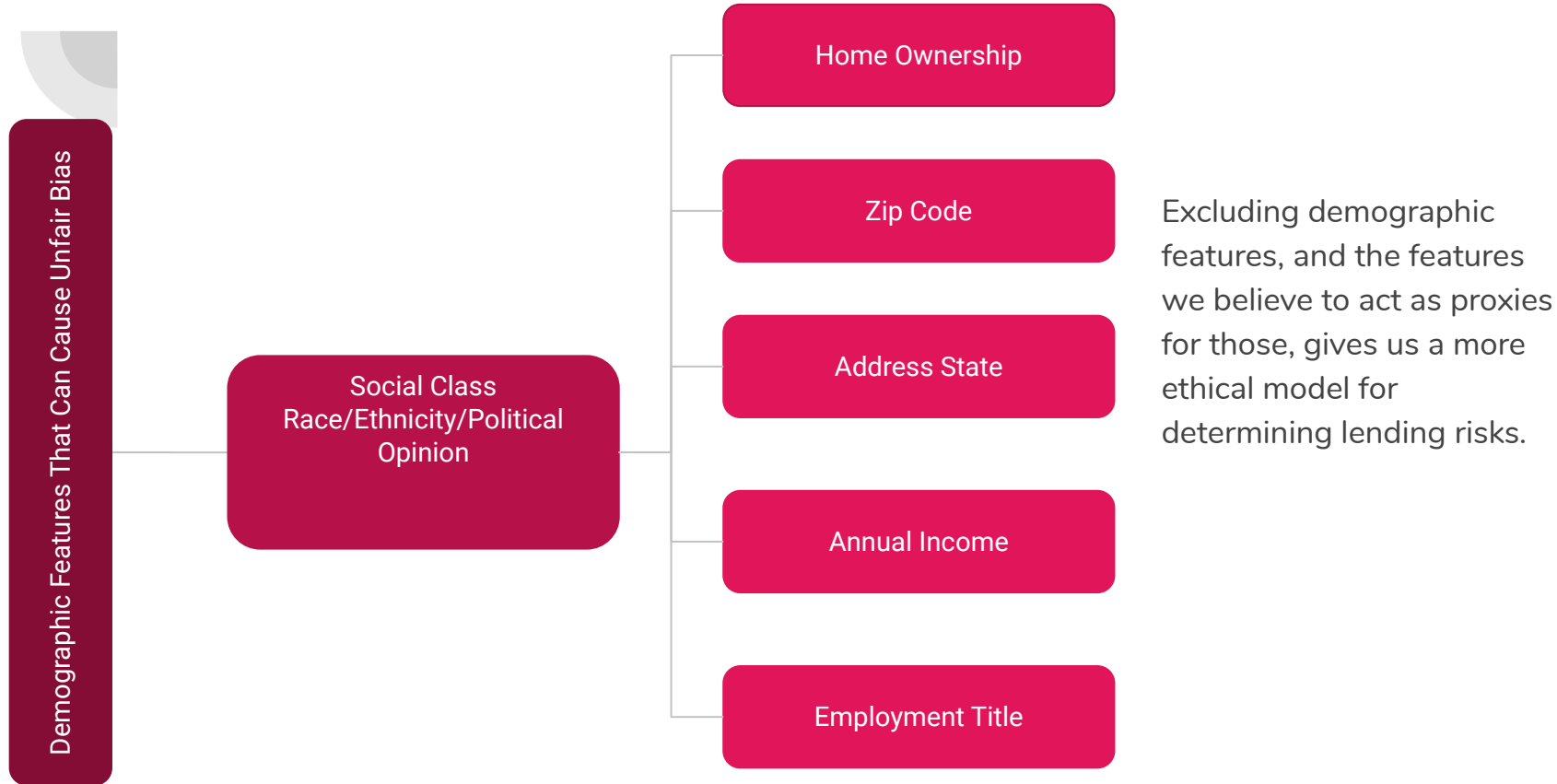


Note: Hispanics may be of any race. *Not reported* categories not shown. Data on whites, blacks and Asians refer to single-race groups.

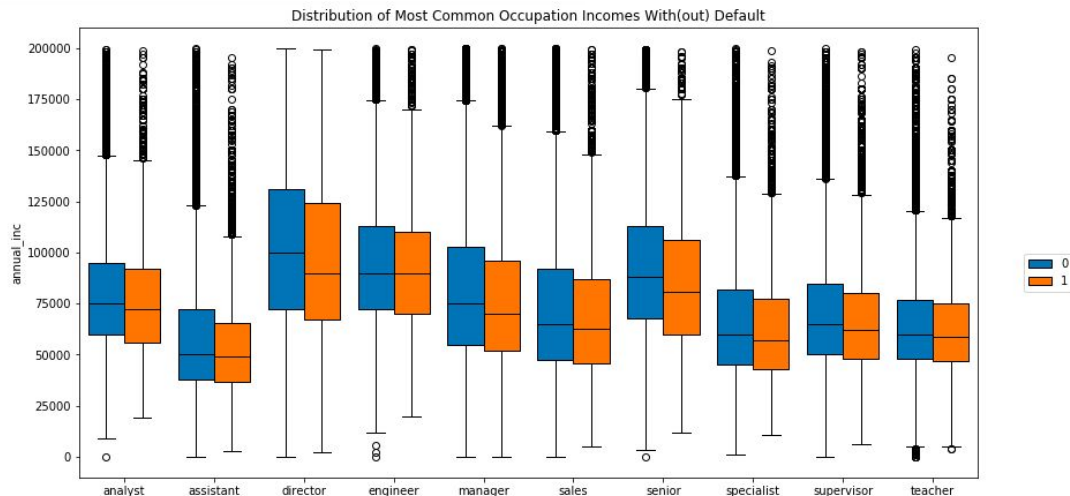
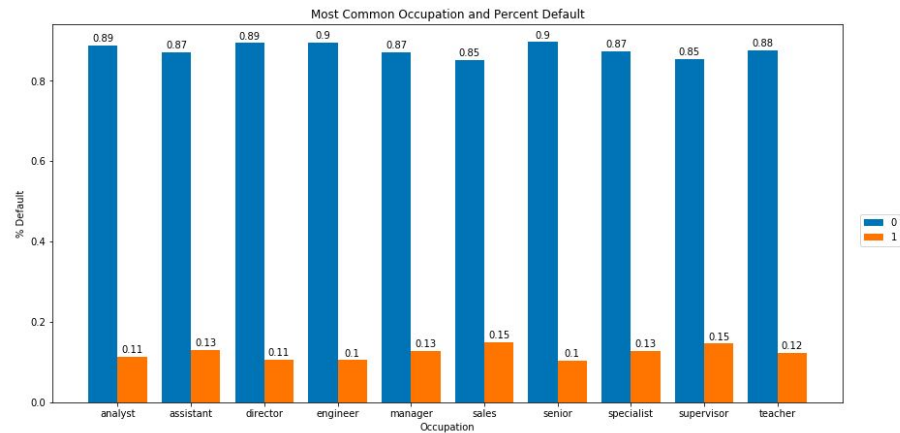
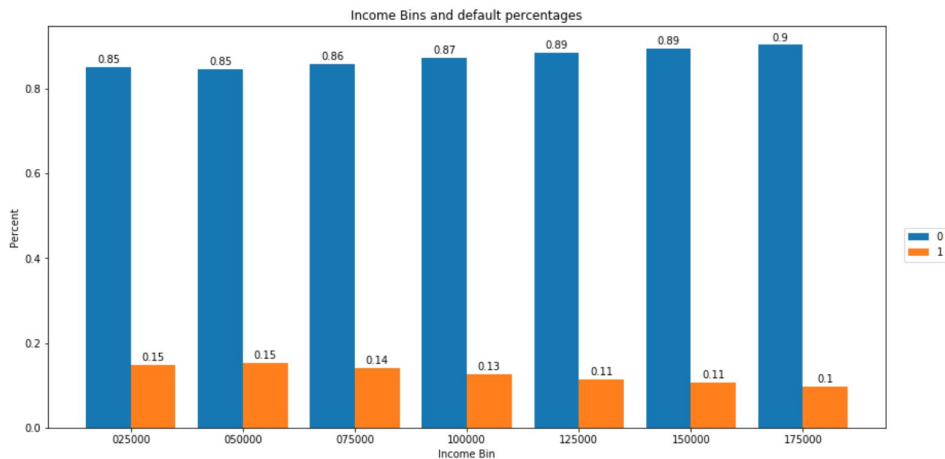
Source: Pew Research Center analysis of American Housing Survey data

PEW RESEARCH CENTER

Biased VS Unbiased



Income vs. Loan Amounts/Jobs & Defaults





Approaches Key Qualities

Ability To Deal With Class Biased Data

Only 13% of our dataset is labelled as positive (late or default). Therefore, our model must be able to deal with a biased dataset.



High Accuracy

In order for companies to use this as their main model, we would need higher accuracies. Otherwise, that would pose the companies too much of a trade off between risk and being socially unbiased.

Interpretability

Since companies might need to use this to compare their decision factors vs the unbiased decision factors, features that cause a borrower to be late need to be analyzable. Moreover, we need to be able to analyze the difference between the biased and unbiased model.



Models That Reflect Key Qualities

Interpretability

We get Interpretability with the following models :

- Linear Models
- Decision Trees
- Clustering Algorithms

Class Biased
Data

Using Ensemble learning we deal with biased datasets using the bootstrapping, weighted sampling and random sampling techniques. The following exist in random forests, boosted trees, and bagged predictors.

Optimal
Complexity

Since we choose the number of simple predictors we want to use in ensemble methods, we are free to fit our dataset with as high or low of an accuracy as we want.



Success Indicators Criteria

- Cross Validation was not needed as our Dataset was sufficiently large
- 75/145 features were used
- Low False Negative Rates
- High Accuracy

Model Results

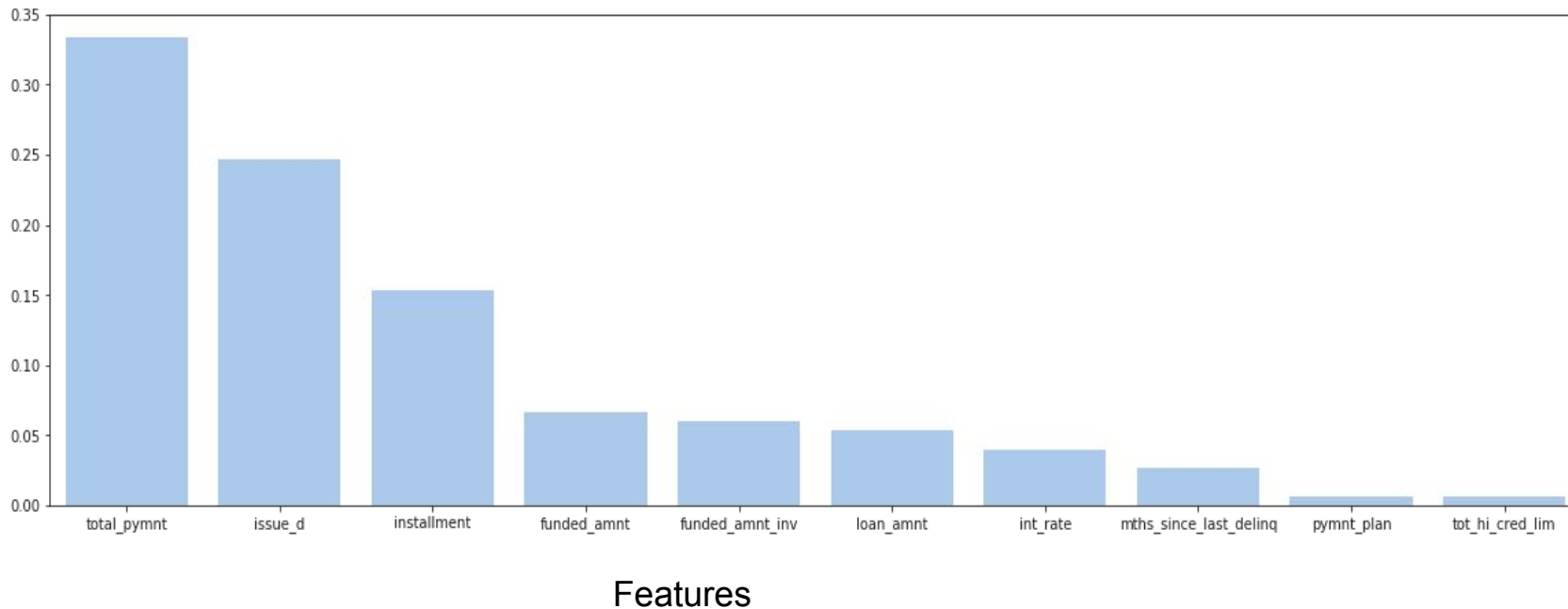
		Biased Model Accuracy	Unbiased Model Accuracy	False Negative Rates
1	K-Nearest Neighbors Clustering	86.4 %	88.5 %	<ul style="list-style-type: none">• 10.9%• 9.5%
2	Linear Regression	86.5 %	86.3 %	<ul style="list-style-type: none">• 13.4 %• 13.4 %
3	Random Forest	97.11 %	97.23 %	<ul style="list-style-type: none">• 2.73%• 2.6%
4	Boosted Gradient Trees	95.5 %	96 %	<ul style="list-style-type: none">• 1.1%• 1.1%



Gradient Boosted Trees

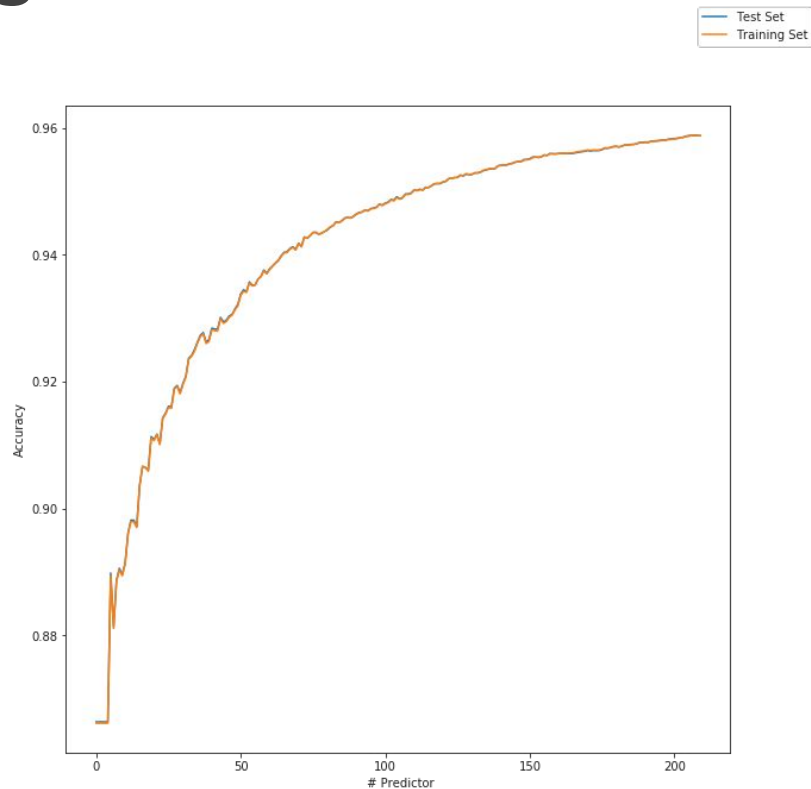
Biased Model Feature Importances

Importance by Percentage





Training And Validation Accuracies





Conclusion

- Neither Income nor Zipcode played a role in our unbiased Gradient Boosted Model with 96% accuracy, which shows that companies do not need such information for their screening process.
- Types, amounts and installments qualities of the loan matter the most in whether a borrower is going to be late. This would help companies design loans that would be more in the bracket where borrowers can pay on time.
- If a company is screening based on income, gender, race or location, it was shown that the data does not support such decisions.



Biased VS Unbiased

Even though there is a pattern observed between loan repayment and area as well as loan repayment and income bracket, we don't use these biased features.

