

AI Risk Report

Project Title

Team: Srimukha Sarma, Project Name: Bias Detection in Loan Approval Dataset

1. Problem Overview

- **What was the task? (e.g., Predict loan approval)**

The task was to detect bias in a Loan approval dataset.

- **Why does it matter in a real-world or ethical context?**

This helps us identify which groups in our society are being unfairly favoured. It also helps us identify where the problem lies and come to an effective solution.

- **What dataset were you given, and what were the known sensitive attributes?**

I was given a dataset which had loan approval reports, The sensitive attributes are:

'Gender', 'Race', 'Employment Type', 'Education Level', 'Citizenship Status', 'Disability Status', 'Criminal Record', 'Age', 'Age Group', 'Income', 'Credit Score', 'Zip Code Group', 'Language Proficiency', 'Loan Amount'

2. Model Summary

- **What model(s) did you use and why?**

After training and evaluating multiple models (Logistic Regression, XGBClassifier, Random Forest Classifier and Custom neural network) I've decided to use the logistic regression model as it had the highest accuracy.

- **Key preprocessing, feature engineering, or hyperparameter choices**

FEATURE ENGINEERING: I added 3 more input features (Income Group, Credit score group, loan amount group) to better evaluate the fairness metrics and I label encoded all the columns for training.

Hyperparameters: I trained the logistic regression model for a maximum iterations of 1000.

- **Performance on internal validation data (accuracy, precision, recall, etc.)**

The accuracy, precision, recall of the model are as follows: Accuracy: 0.626, Precision: 0.5833333333333334, Recall: 0.4652241112828439

3. Bias Detection Process

- **Methods used to detect bias (e.g., SHAP, Disparity Impact Ratio, False Positive Rate comparisons, Grouped AUC, Fairlearn audit, etc.)**

I considered metrics like Equal Opportunities Difference(EOD), Average Odds Difference(AOD), Theil Index(TI), Statistical Parity Difference(SPD) and Disparate Impact(DI) to detect bias. I also used SHAP and Lime to check which of the features had the highest contributions to the model's output. Then I checked the intersection fairness of two groups at a time to check if there was a certain combination which was being favoured.

- **Did you audit raw data, model output, or both?**

I audited both raw data and model output.

- **Were audits performed at the individual or group level?**

Audits were performed at a group level and individual level.

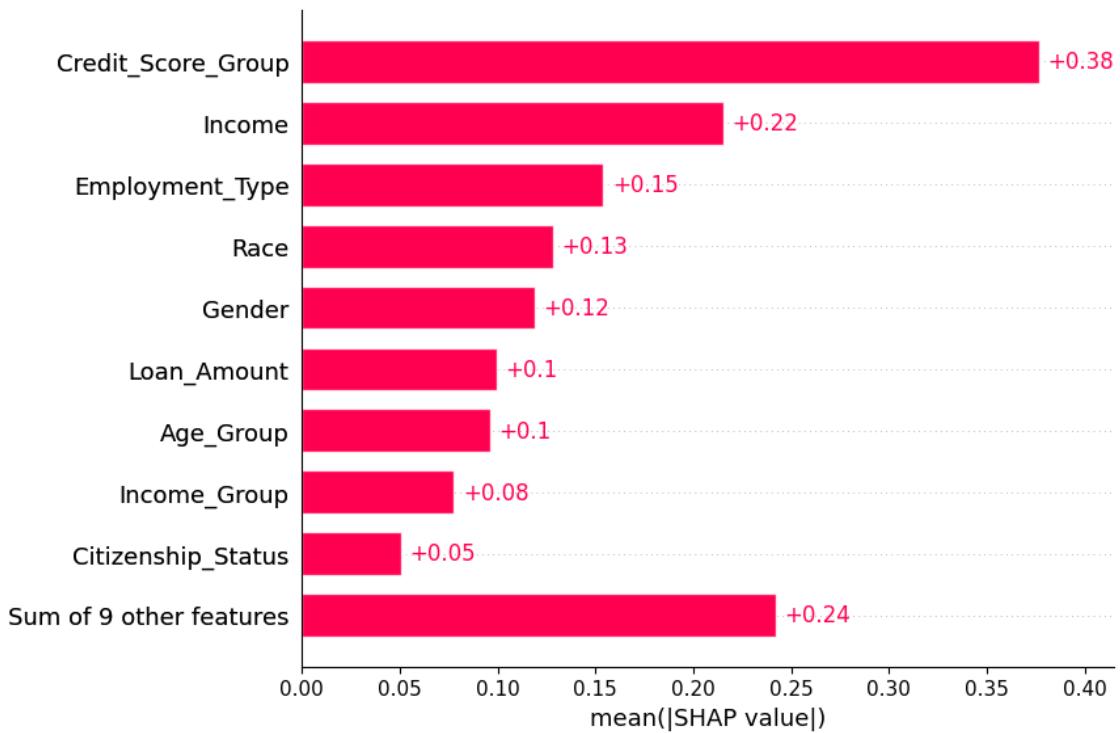
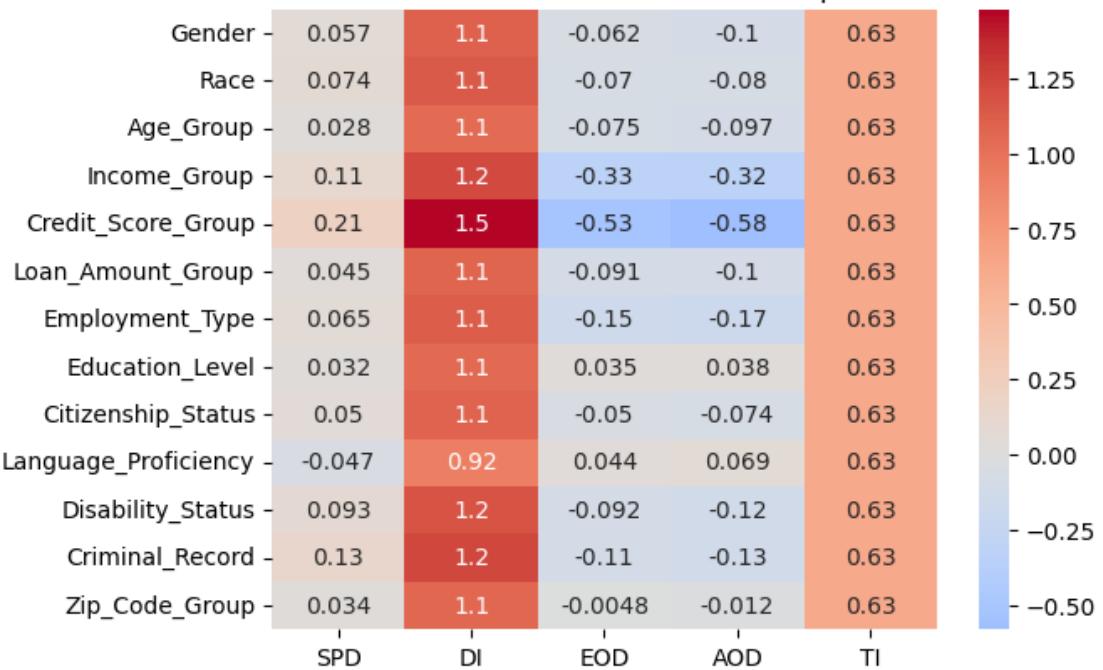
4. Identified Bias Patterns

Bias Type	Affected Group	Evidence	Metric	Comment
Negative	Low Income Group	Aif360 heatmap	EOD,AOD, TI,SPD ,DI	Low income groups have much lower approval rates
Negative	Low Credit Score Group	Aif360 heatmap	EOD,AOD, TI,SPD ,DI	Theres a strong bias against people with a low credit score.

Negative	People with Criminal Records	Aif360 heatmap	EOD,AOD,DI,SPD	Theres a strong bias against people with criminal records.
Negative	People with Disabilities	Aif360 heatmap	EOD,AOD,DI,SPD	Theres a moderate bias against people with no disabilities
Positive	Full-Time worker	Aif360 heatmap	EOD,AOD,DI,SPD	Theres a moderate bias toward people with full time jobs.
Negative	High loan amount	Aif360 heatmap	EOD,AOD,DI,SPD	Theres a mild bias toward people taking a large loan.
Positive	Male	Aif360 heatmap	EOD,AOD,DI,SPD	Theres a mild bias toward males.
Negative	Multi-Racial people with low income	Intersection Fairness heatmap	Loan Approval rate	Theres a moderate bias towards them
Positive	25-60 year olds with high credit score	Intersection Fairness heatmap	Loan Approval rate	Theres a strong bias towards them.

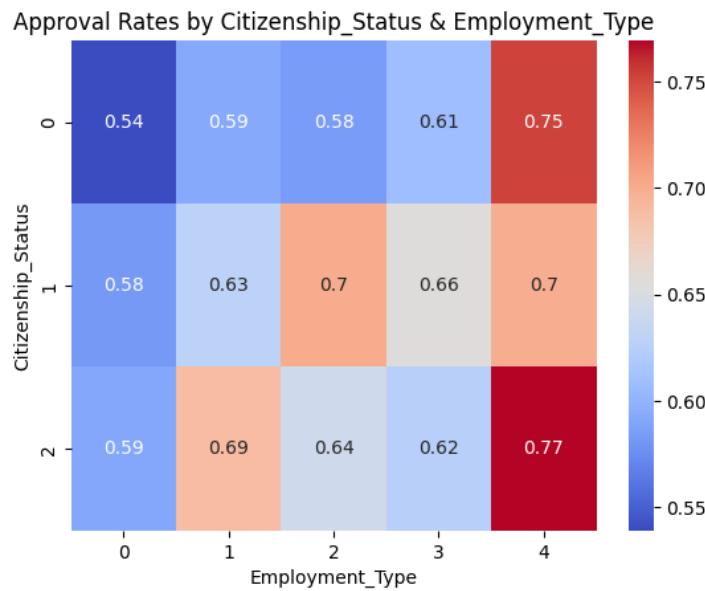
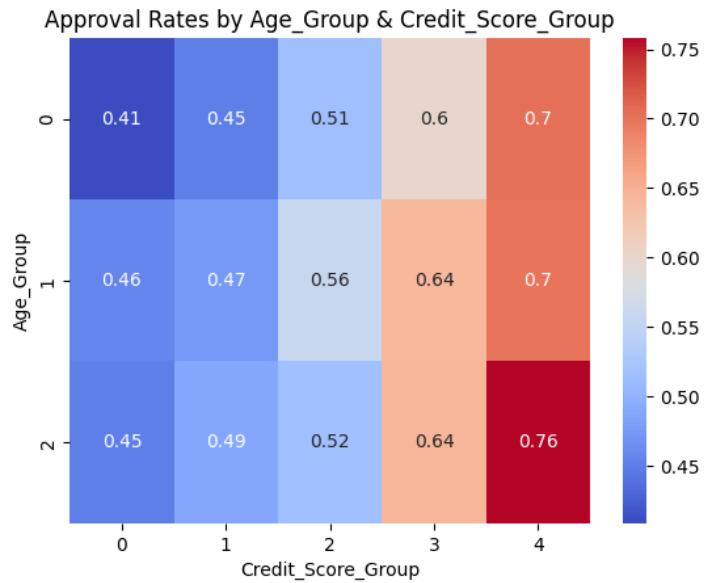
5. Visual Evidence

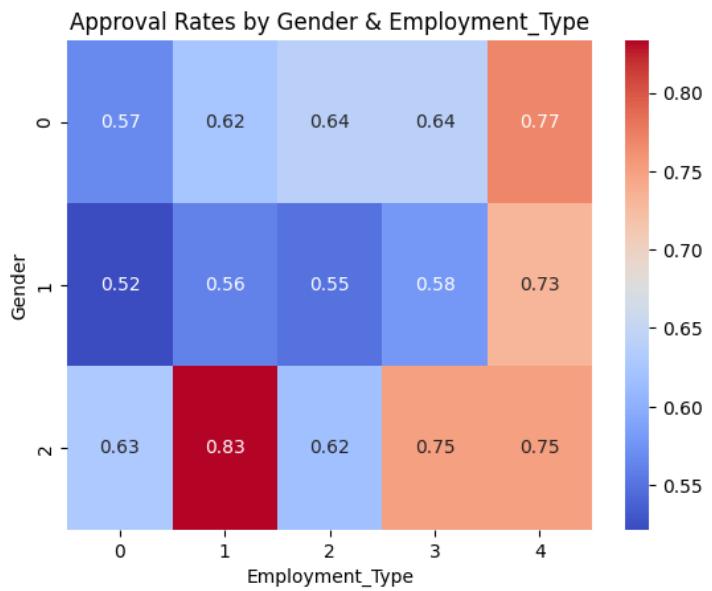
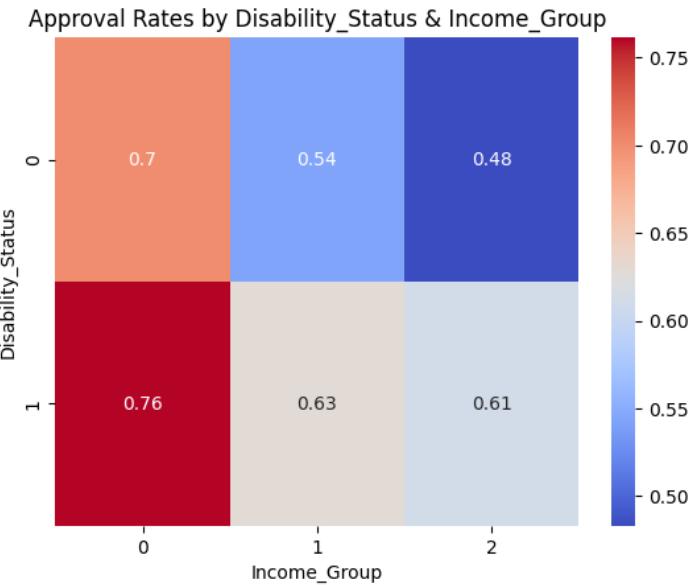
AIF360 Fairness Metrics Heatmap

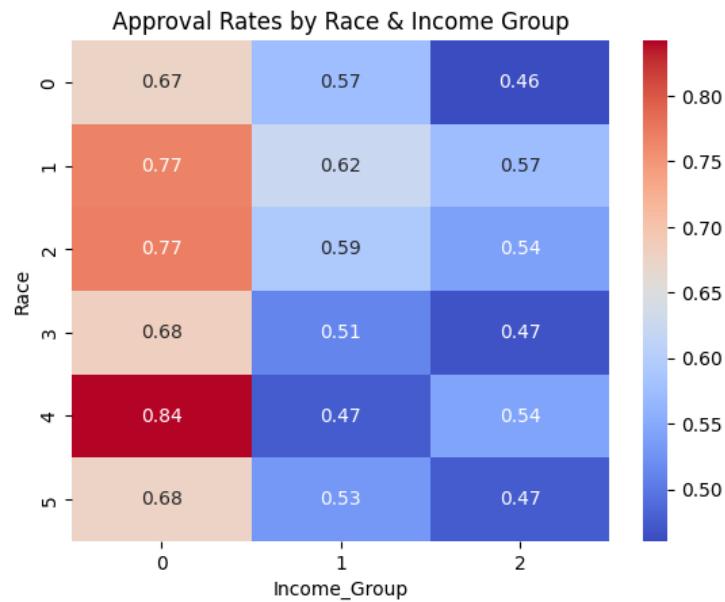
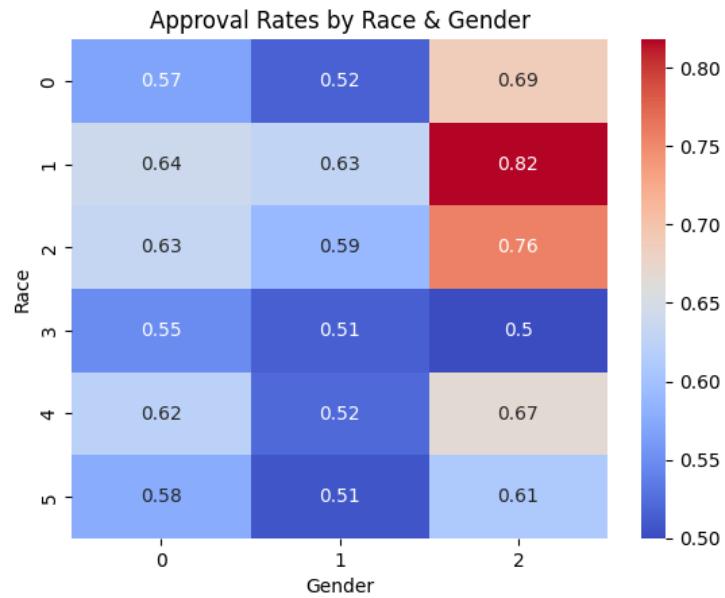


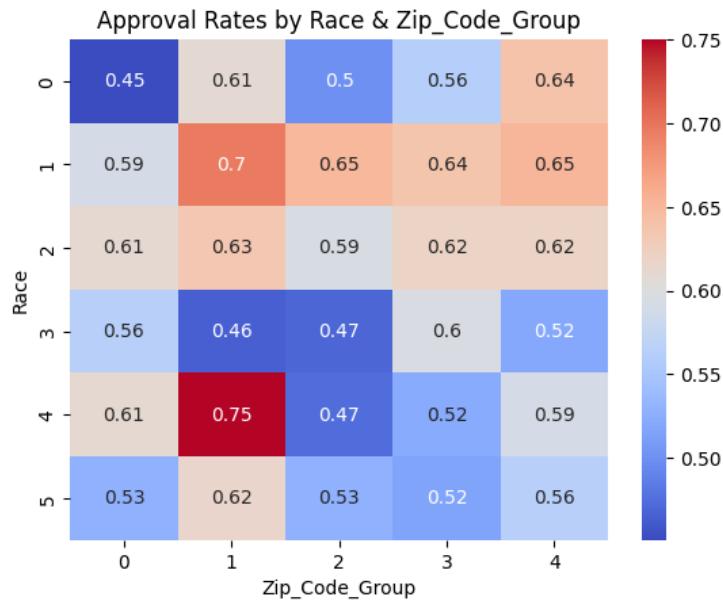


LIME RESULT









6. Real-World Implications

- **Who is most at risk if your model were deployed as-is?**

All the groups which have been under a negative bias would suffer as the model is trained on a biased dataset and all the groups which have a positive bias will continue to have an unfair advantage.

- **What are the ethical or social consequences?**

Biased loan models can unfairly deny access to credit for certain groups, deepening social and economic inequalities.

They risk violating anti-discrimination laws and AI principles, leading to legal consequences.

Such bias also destroys public trust and causes injustice.

- **Would your model pass a fairness audit in a regulated setting?**

No, the current model would likely not pass a fairness audit.

To pass an audit, the model would need:

- Bias mitigation applied in training or decision thresholds.
- Regular audits with fairness metrics.
- Transparent documentation explaining design decisions and their social implications.

7. Limitations & Reflections

- **What didn't work?**

Purely basing the detection on approval rates and confusion matrices was very ineffective due to the data imbalances. Metrics like EOD,AOD,TI,SPD,DI worked better in detecting bias.

- **What would you try next time with more time or data?**

I would try intersection fairness with many more combinations to detect specific biases. I would also try out more metrics available in the aif360 tool.

- **Any lessons learned on fairness or auditing?**

I've learnt the different parameters which need to be considered for bias detection and have gotten very comfortable with fairness detection tools like IDM aif360.