



**THE GENDER
PAY GAP IS REAL**



Chris Pratt was paid around \$10 million more than this female extra for doing Jurassic World because he's a man.



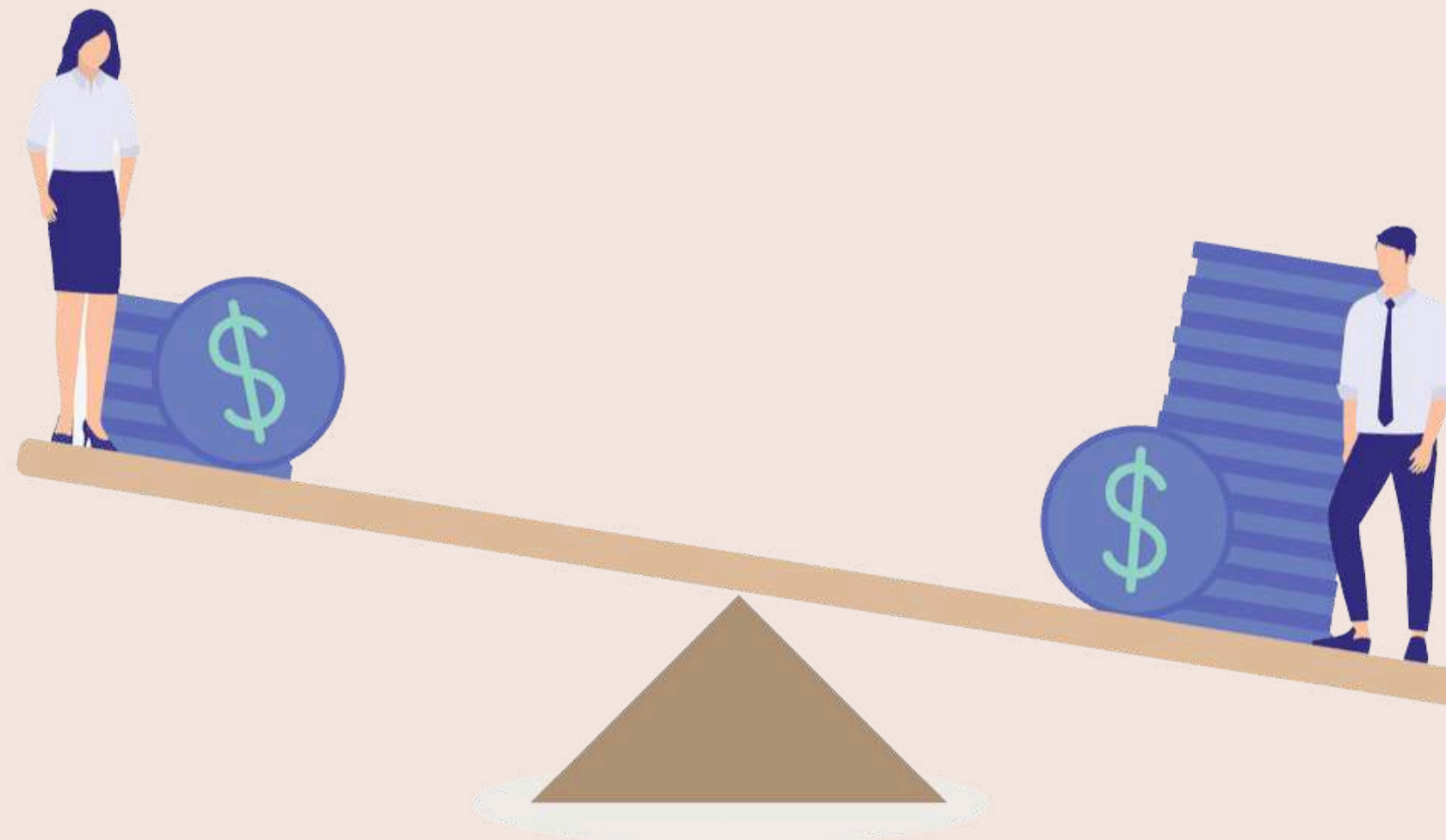
Pay Parity Across Multiple Demographics



Kavyaa Agrawal, Sadhika Anand and Sahil Gada

Problem Statement

We examine whether **income** differences between different demographics, especially men and women persist within comparable groups by analyzing factors such as **education, occupation, workclass, race, marital status,** and **hours worked**, using exploratory data analysis, **XGBoost** for income prediction, and **SHAP** to explain the impact of demographic features on predictions.



LITERATURE REVIEW: KEY FINDINGS

83 cents

Women earn for every 1 \$
earned by a man

38%

of the wage gap remains
unexplained

74 cents

Working mothers earn for every \$1
earned by a man

57 cents

earned by Hispanic women per 1\$ earned by a
white man

64 cents

earned by Black women per 1\$ earned by
a white man

\$964,000

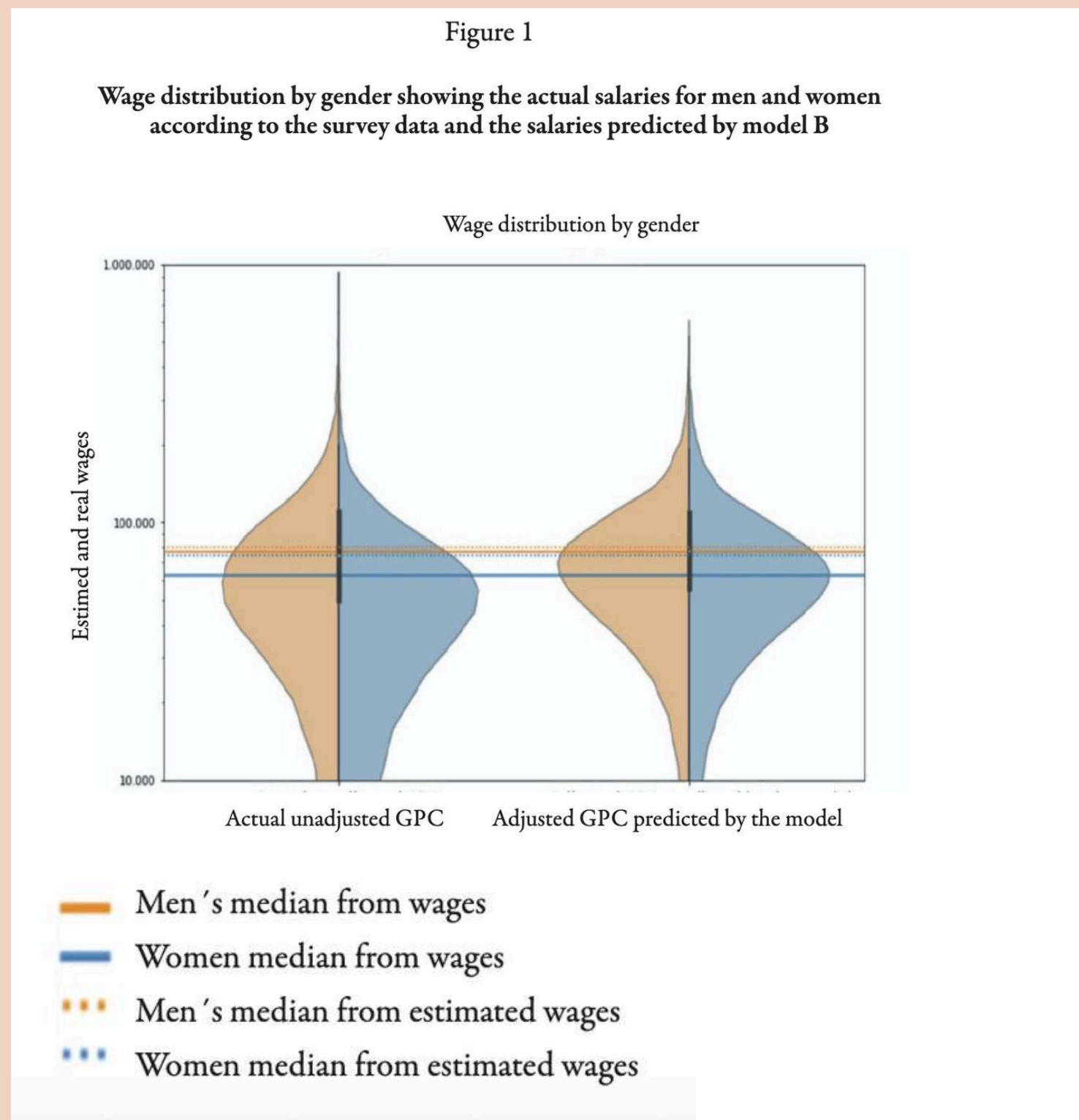
estimated income losses over a 40 year
career for Black women

\$1.16 million

estimated income losses over a 40 year
career for Hispanic women

Literature review 2 - Argentinean Model

Goal: Estimate gender pay gap using salary prediction.



DATASET

5,700+ IT WORKERS

MODEL

XGBoost Regression

OBJECTIVE

Explained vs unexplained gap

LIMITATION

Limited intersectional analysis

20%

Total gender pay gap

Literature review 2 - Chilean Model

Goal: Estimate the causal effect of gender on wages..

$$\ln(w_{ij}) = \beta_0 + \beta_1 Female_{ij} + \beta X_{ij} + \varepsilon_{ij}$$

$$\log(w_m) - \log(w_f) = (\bar{X}_m - \bar{X}_f) \beta_f + \bar{X}_m(\beta_m - \beta_f)$$

Field of study	2017			2022			Change in gender wage gap
	Coef.	St.Err.	N	Coef.	St.Err.	N	
(0) All fields	-0.171***	0.012	15,340	-0.151***	0.009	22,691	[-]***
(1) Education	-0.088***	0.028	3162	-0.041*	0.023	3891	[-]***
(2) Arts & Humanities	-0.015	0.073	381	-0.047	0.046	693	
(3) Social Sciences & Communications	0.025	0.064	634	-0.163***	0.048	781	[+]**
(4) Business & Laws	-0.157***	0.025	3529	-0.159***	0.019	5122	[+]***
(5) Natural Sciences, Math & Statistics	-0.031	0.122	178	-0.220***	0.077	290	[+]***
(6) Information & Communication Technology	-0.195***	0.069	584	-0.112**	0.056	935	[-]***
(7) Engineering, Industry & Construction	-0.090***	0.035	2901	-0.129***	0.024	5068	[+]***
(8) Agriculture, Forestry, Fishing & Veterinary	-0.091	0.071	435	-0.114*	0.060	447	[+]***
(9) Health	-0.235***	0.036	2741	-0.144***	0.024	4203	[-]**
(10) Services	-0.156***	0.052	795	-0.119***	0.034	1261	[-]***

DATASET

National Labour Survey from Chile

MODEL

ML Regressors + SHAP

INTERPRETATION

Shap values identifying influential factors

LIMITATION

Limited intersectional analysis

15%
Estimated pay gap

Literature review 3 - Research GAPS

LIMITATION 1: Existing Model Limitations

Studies mainly focus on occupation and education differences.

Intersectional factors such as race, marital status, and motherhood effects are ignored.

Average salary gaps are measured but hidden opportunity and leadership gaps remain unexplored

Hidden inequality patterns remain unexplored.

LIMITATION 2: Policy & Interpretation Gaps

Existing models explain wage disparities but provide few actionable policy recommendations.

Limited insights on which sectors or demographic groups require intervention.

Most studies focus on salary prediction rather than identifying structural inequality patterns.

Low Policy Interpretability.

SOLUTION: Our Proposed Approach

Combines gender with race, marital status, occupation, and education factors.

Uses explainable ML techniques (SHAP) to identify drivers of inequality.

Focuses on both wage disparities and leadership parity as well.

Targets both wage and leadership parity.

Dataset Overview

US Census American Community Service data, 2022

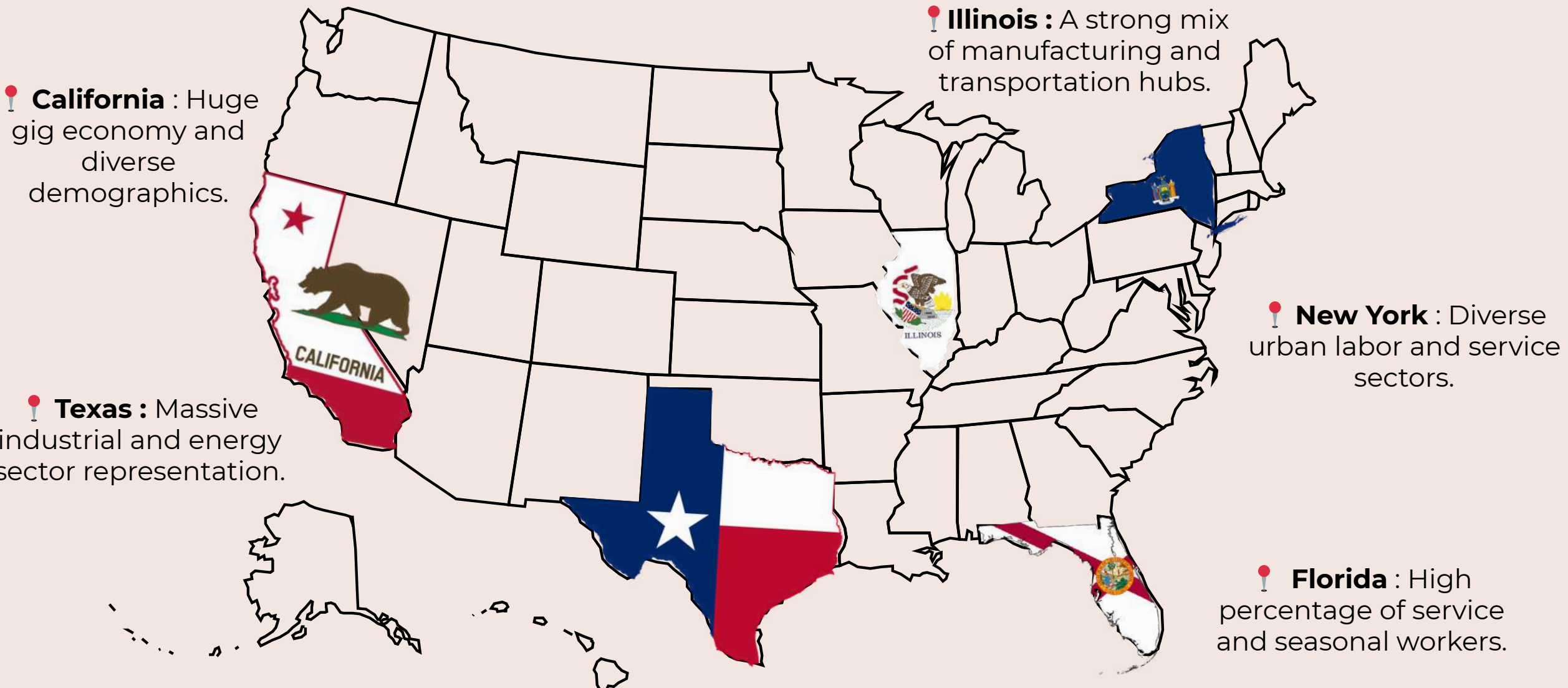
AGEP	COW	MAR	SCHL	SEX	WAGP	WKHP	INDP	OCCP	PINCP	RAC1P	STATE
28	5.0	1	19.0	1	28500.0	55.0	9670.0	9825.0	28500.0	1	CA
20	1.0	5	19.0	1	7800.0	30.0	7680.0	3930.0	7800.0	7	CA
27	1.0	5	16.0	2	30000.0	40.0	7680.0	3930.0	30000.0	9	CA
24	4.0	5	21.0	2	31700.0	40.0	7870.0	2205.0	31720.0	6	CA
20	5.0	5	16.0	1	25000.0	40.0	9770.0	9825.0	25000.0	8	CA

567,921

Raw observations from 5 US states

14 columns

Which include our main features



WHY WE CHOSE THIS DATASET:

- Exact Feature Alignment
- Absence of Viable Alternatives

Dataset Overview

US Census American Community Service data, 2022

Data Collection

- Administered by the U.S. Census Bureau (2022)
- Randomly sampled ~3.5 million households nationwide.
- **Validity:** Legally mandatory participation ensures high statistical representation.
- **Methodology:** Mixed-mode gathering (Internet, Mail, Phone, In-Person) to minimize non-response bias.



Ethical Vulnerabilities

High Sensitivity:

Contains granular personal data including exact wages, racial identity, and disability status.

Re-identification:

Highly specific feature combinations risk unmasking individuals.

Representation Bias:

Historic difficulty in accurately counting unhoused or undocumented populations.

Mitigations

Census Bureau removed personally identifying information before releasing the data, and we did not use columns with personal data which were not relevant.

Used data from multiple states and demographics to ensure fairness and inclusion.

Feature Engineering & Final Output

Transforming raw census variables into interpretable ML-ready features

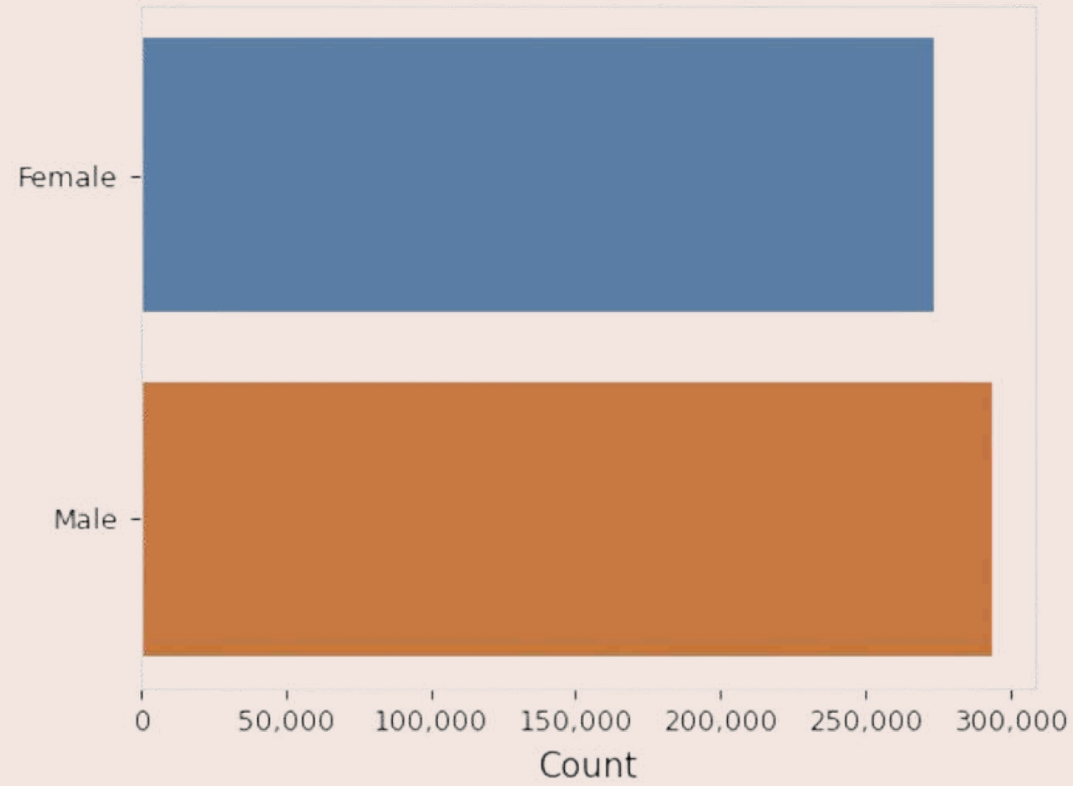
	AGEP	WKHP	OCCP_ENC	SCHL_ENC	MAR_ENC	RAC1P_ENC	COW_ENC	INDP_ENC
count	340752.000000	340752.000000	340752.000000	340752.000000	340752.000000	340752.000000	340752.000000	340752.000000
mean	0.344857	0.561871	10.631680	10.632267	10.632101	10.632116	10.632134	10.631935
std	0.198147	0.310824	0.593217	0.398380	0.283486	0.152679	0.122905	0.464363
min	0.000000	0.000000	9.324738	9.730481	10.262305	10.345734	9.240772	9.285984
25%	0.185430	0.425000	10.192891	10.230744	10.262305	10.488080	10.624326	10.328339
50%	0.331126	0.625000	10.678888	10.552298	10.876694	10.713765	10.624326	10.698234
75%	0.503311	0.625000	11.134760	10.938305	10.876694	10.713765	10.624326	11.040457
max	1.000000	1.000000	11.878112	11.390496	10.876694	10.850973	10.906005	11.840542

DECODING PROCESS:

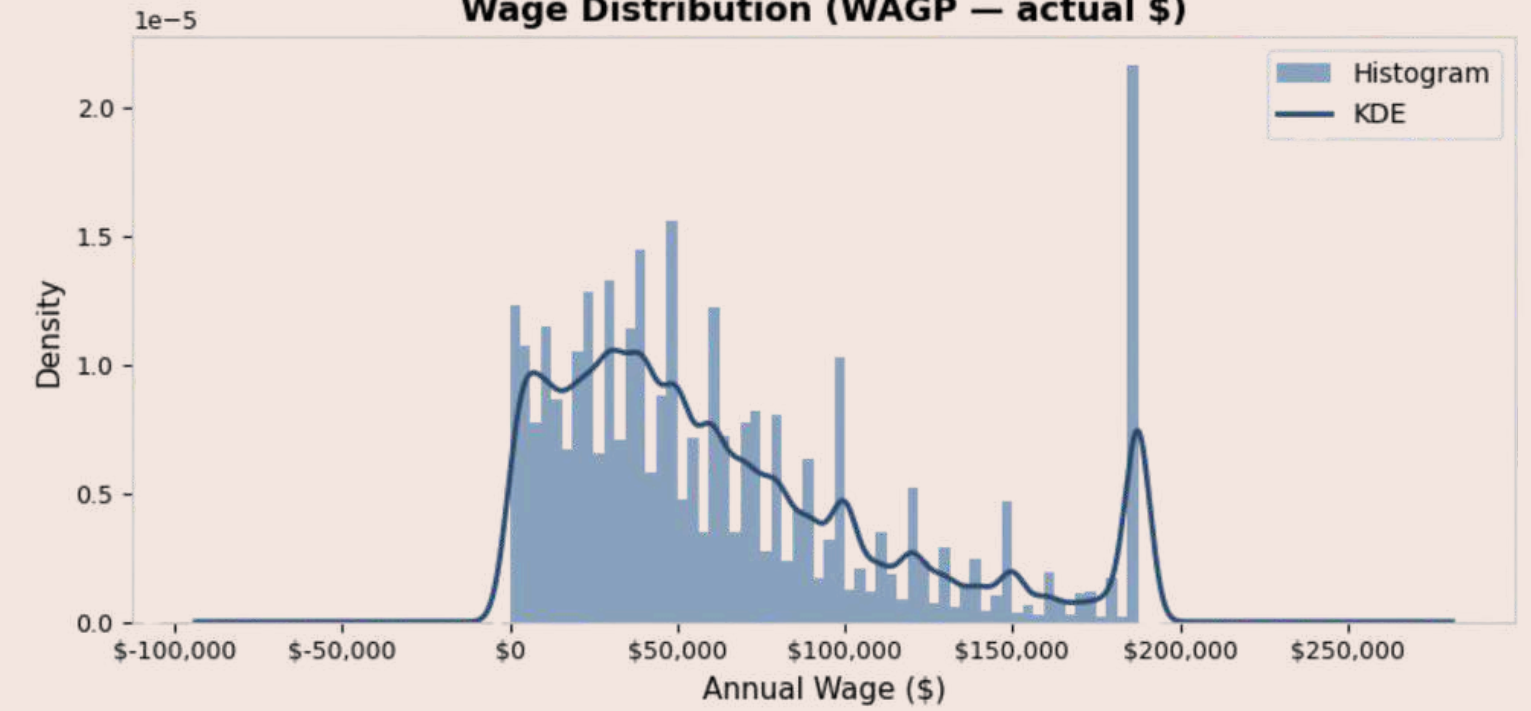


HOW THE DATA IS SPREAD

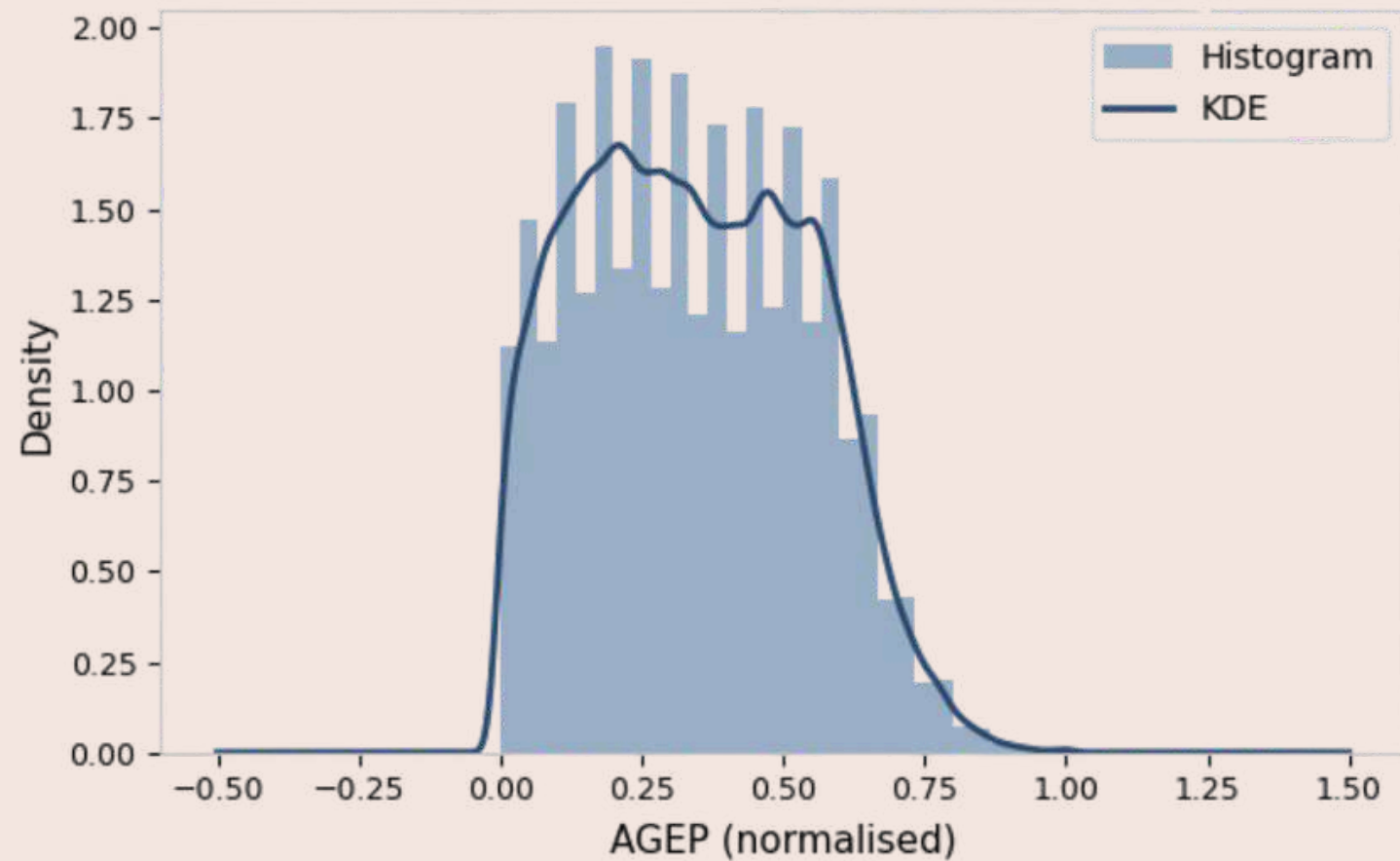
SEX Distribution



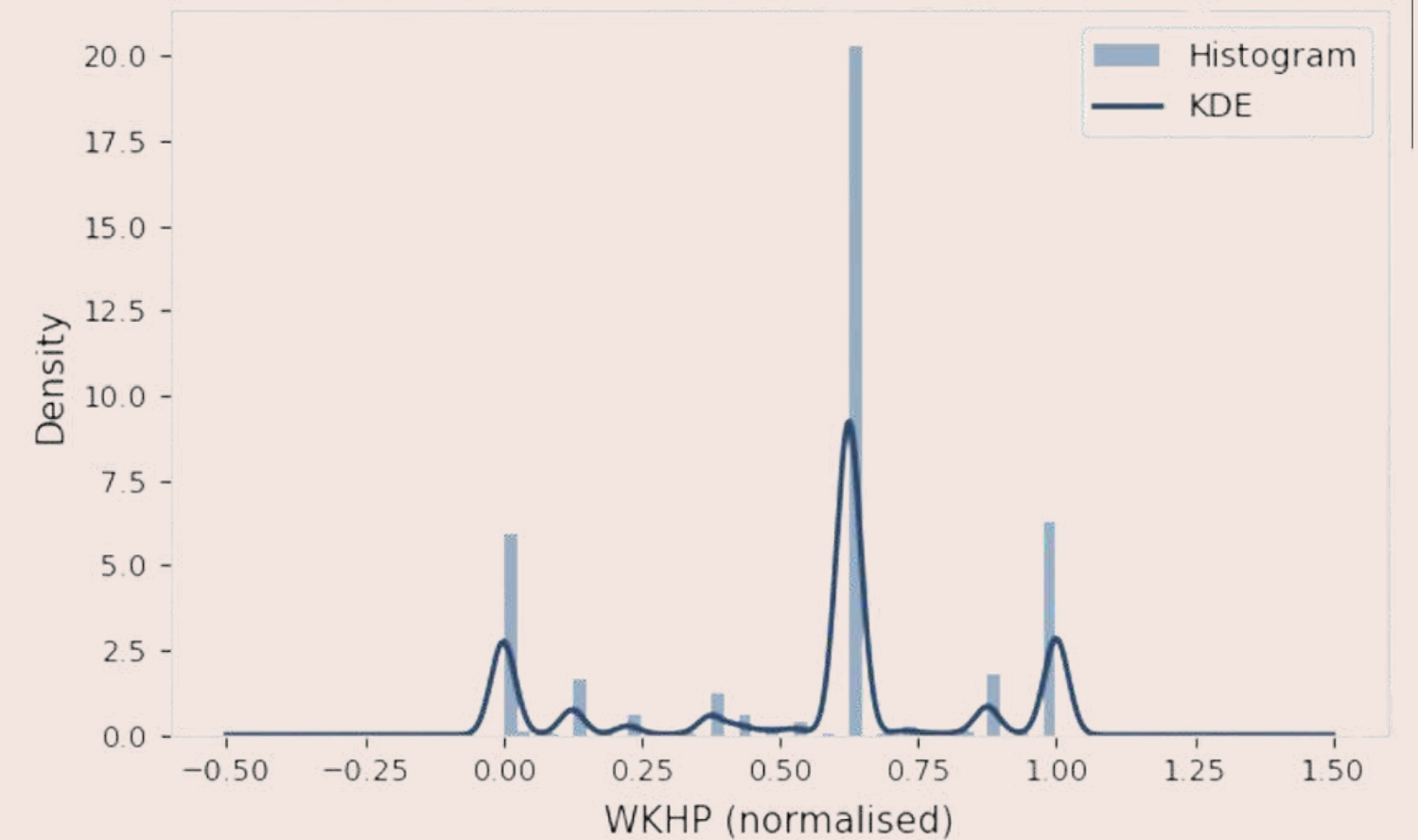
Wage Distribution (WAGP — actual \$)



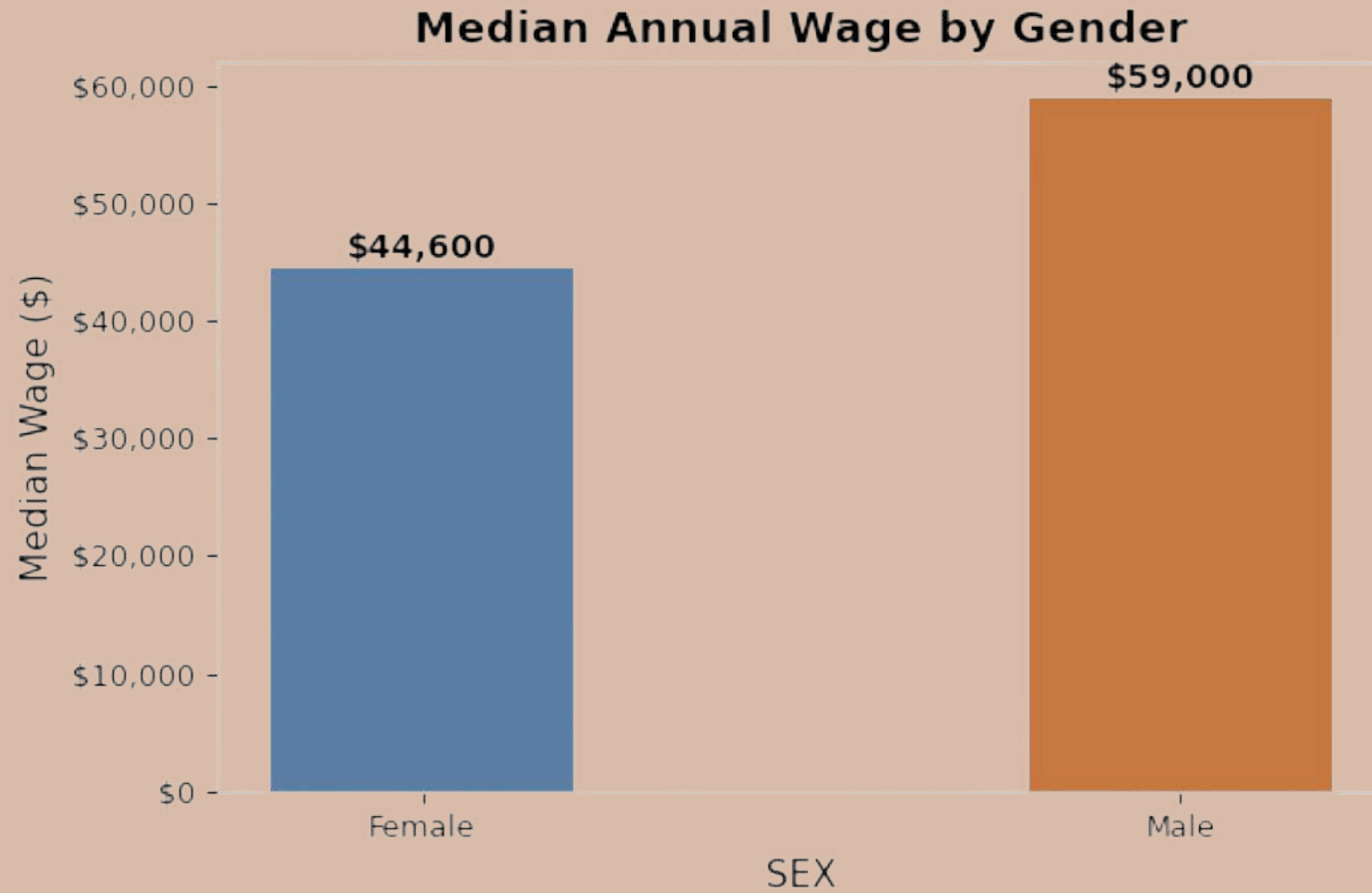
AGEP Distribution (Normalised 0-1)



WKHP Distribution (Normalised 0-1)



Initial Insights from Census Data

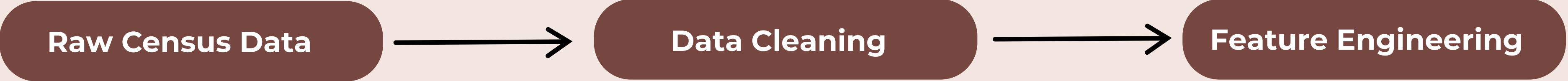


Initial visualizations reveal clear gender disparities in wages across states and occupations.

- Male median wage is consistently higher across all states.
- California and New York show the highest wage levels.
- Alaskan natives and American Indians show the lowest wage levels.

Methodology - XGBoost

STAGE 1: DATA PREPARATION



STAGE 2: MODELING



STAGE 3: INTERPRETATION



Methodology - XGBoost

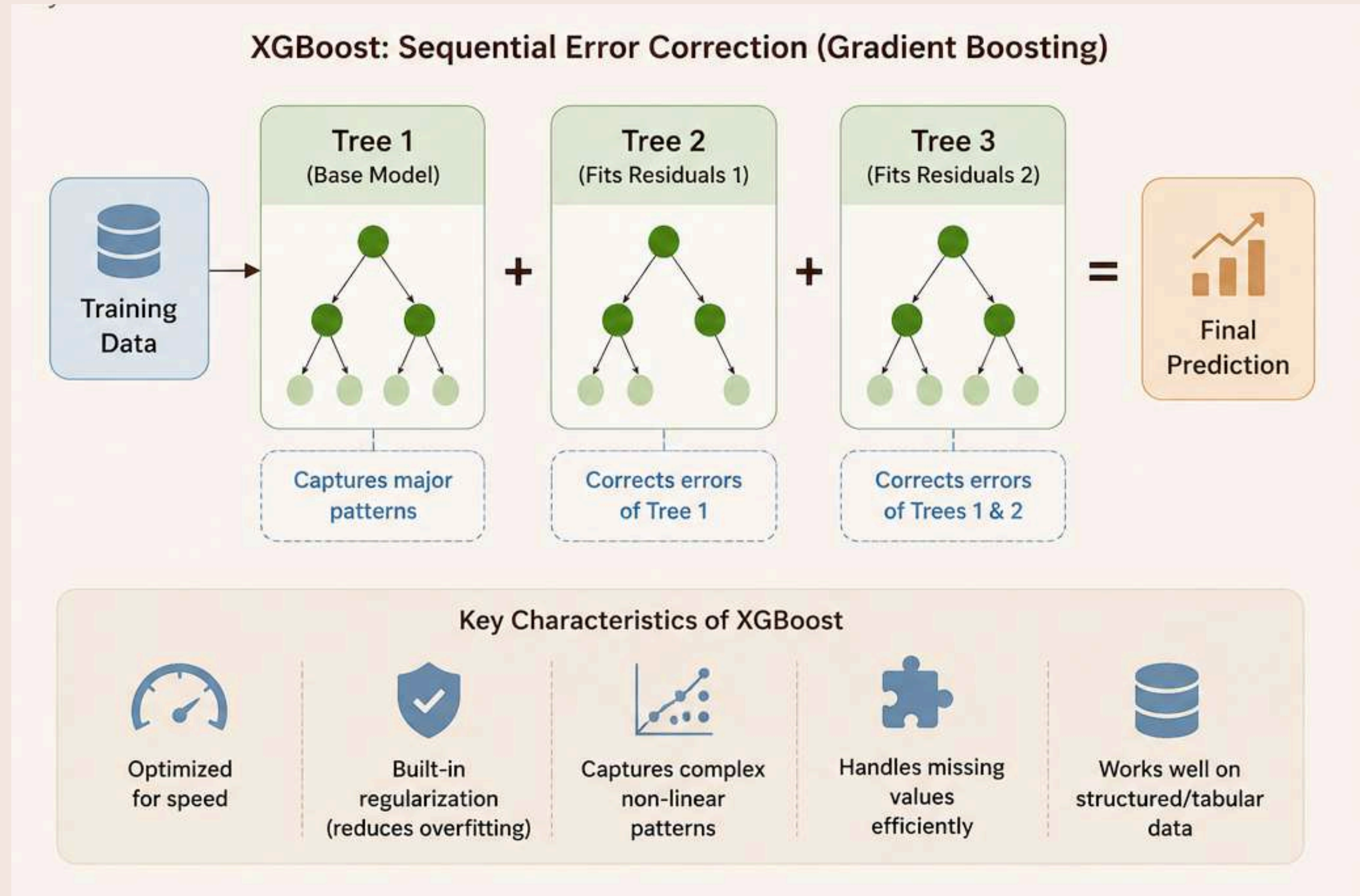
Objective: Accurately map complex, non-linear relationships between demographics, employment factors, and income.

How It Works

- Builds decision trees sequentially
- Each new tree corrects previous errors
- Final prediction combines all trees

Why XGBoost?

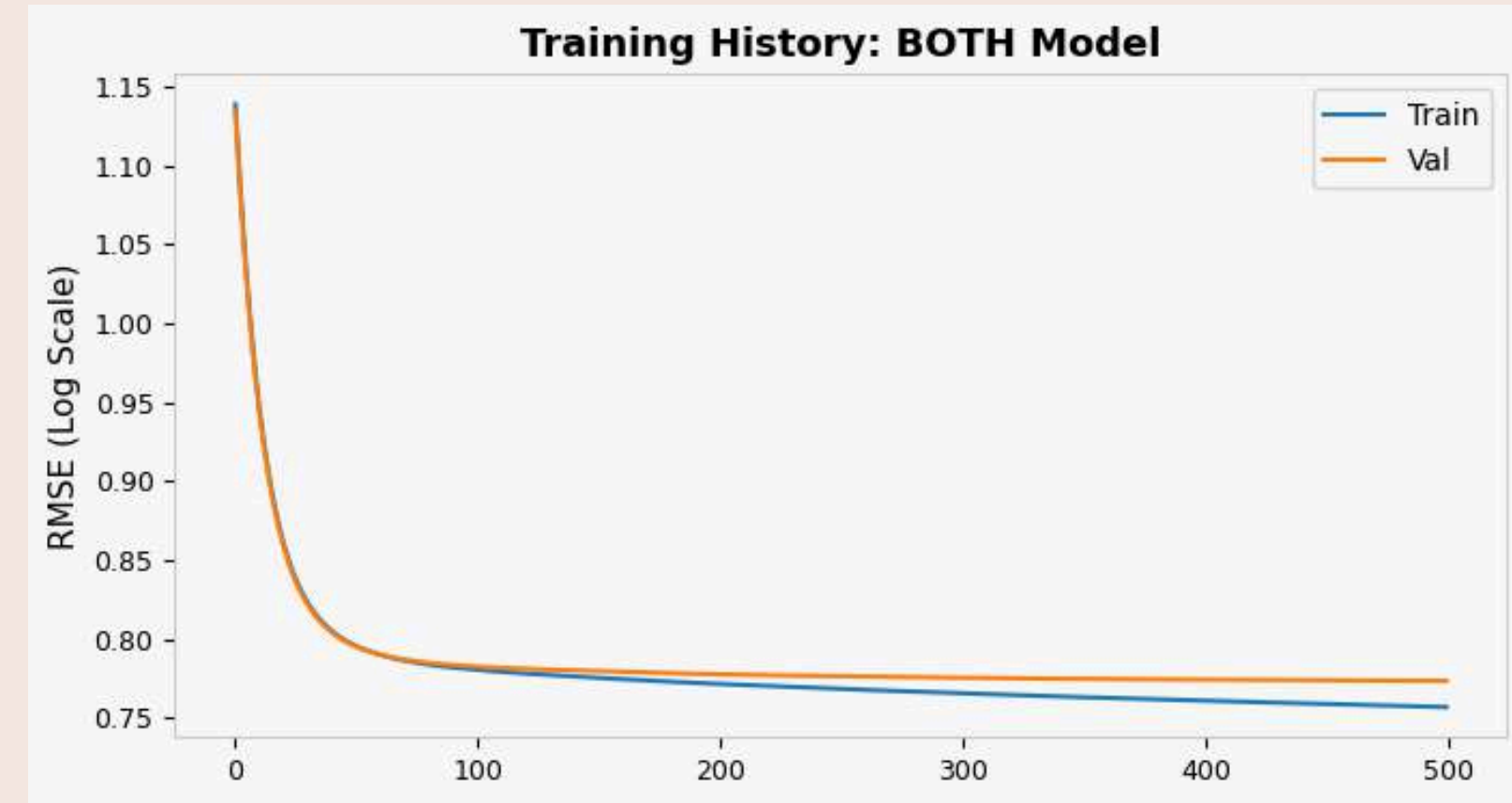
- Handles non-linear patterns
- Captures interaction effects
- Strong performance on tabular census data



Training Strategy & Optimization

Objective: Train a robust model that generalizes well to unseen data without memorizing noise (overfitting).

Hyperparameter Tuning: Checking training vs validation error to get lowest root mean square error (RMSE) while ensuring no overfitting or underfitting.



Data Splitting Strategy

Full Dataset (567,921 Rows)



Training Set (80%)
Model learns patterns here

Eval
(10%)

Test
(10%)

Used for Early Stopping
(Prevents Overfitting)

Unseen Data
(Final Model Evaluation)

Oaxaca - Blinder

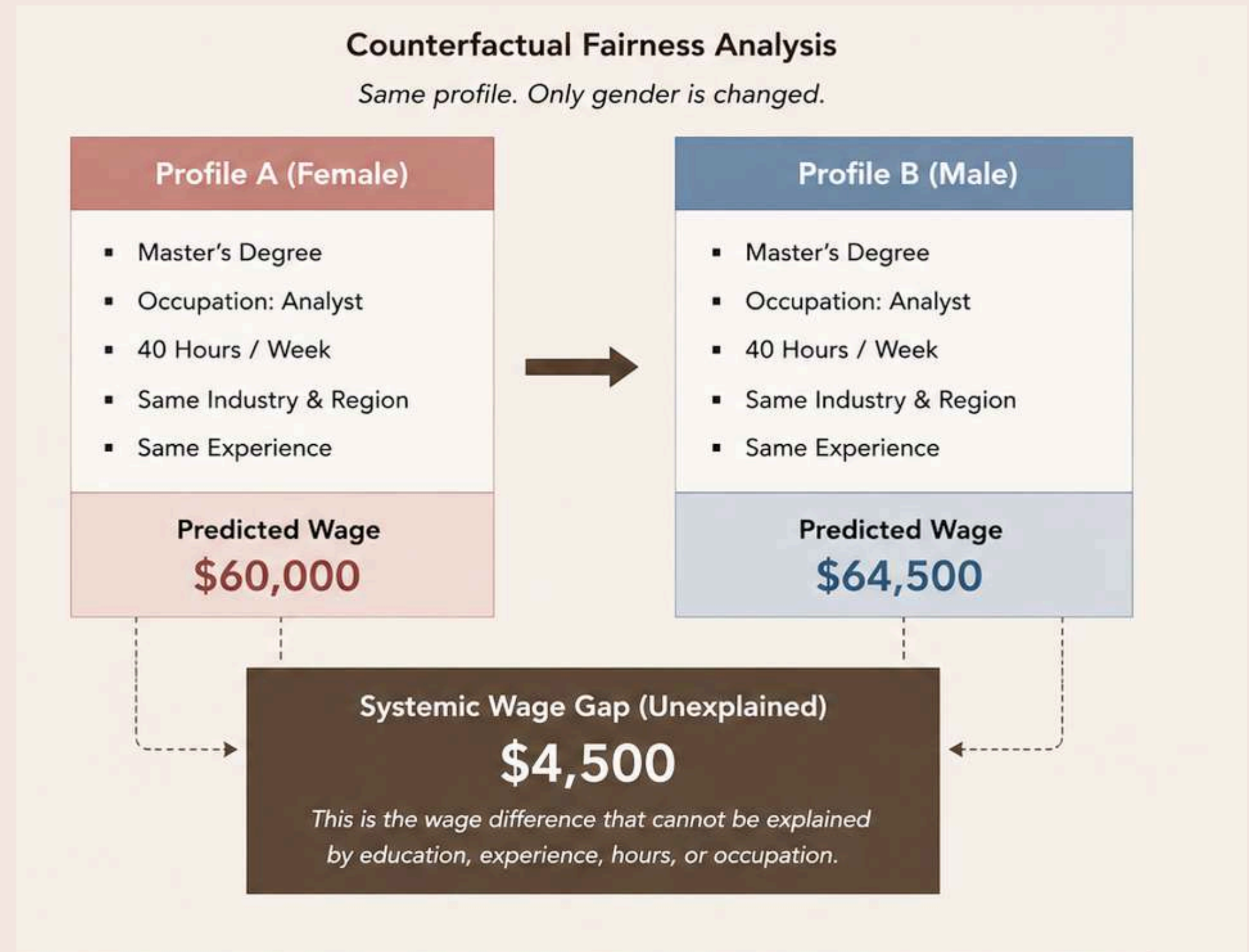
Objective: Mathematically separate justifiable wage differences (productivity, hours) from systemic disparities (the unexplained gap).

The Baseline: The overall expected average wage of the dataset.

Method: Train 2 models for each gender, and put the other gender's datapoints in the model to see the differences in pay.

UNEXPLAINED WAGE GAP:
Calculated by isolating the SHAP contribution of Gender while holding all other variables constant.

COUNTERFACTUAL DIFFERENCE



Interpretability (SHAP)

Objective: Open the "black box" of the XGBoost model to quantify exactly what drives wage predictions.

The Methodology: SHapley Additive exPlanations (SHAP).

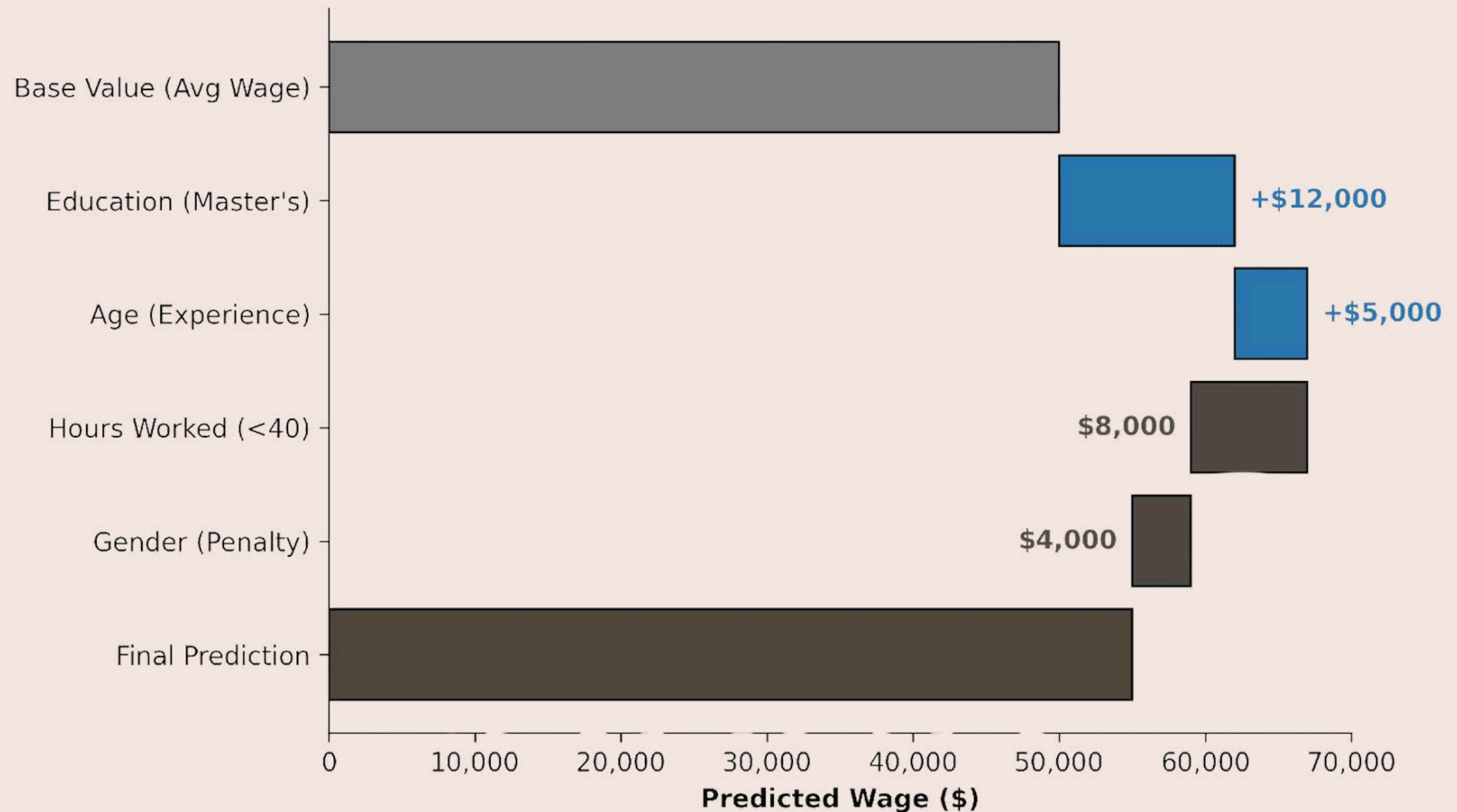
How It Works:

- Every feature contributes to wage prediction
- SHAP measures exact feature contribution

Methodological Justification:

SHAP is mathematically proven to fairly distribute the prediction value among features, ensuring consistent and exact local explanations.

SHAP Explainability: Additive Feature Contributions



Results

Male

When we put male data in the female model

\$71,878.10

Y-Test average male salary

\$55,520.38

Men earn 22.76% MORE than women with identical features.

21.47%

Raw gender gap

Female

When we put female data in the male model

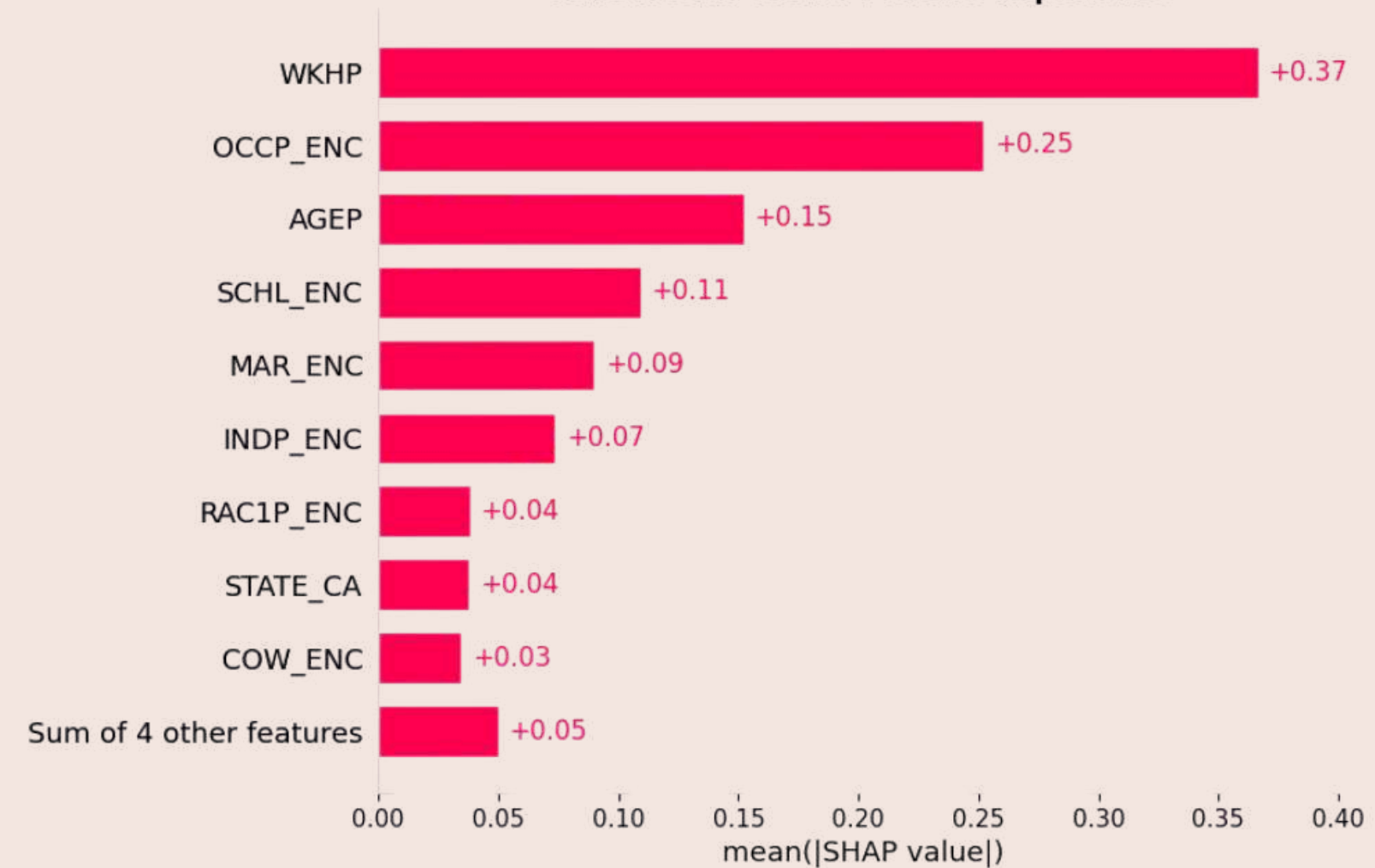
\$56,446.99

Y-Test average female salary

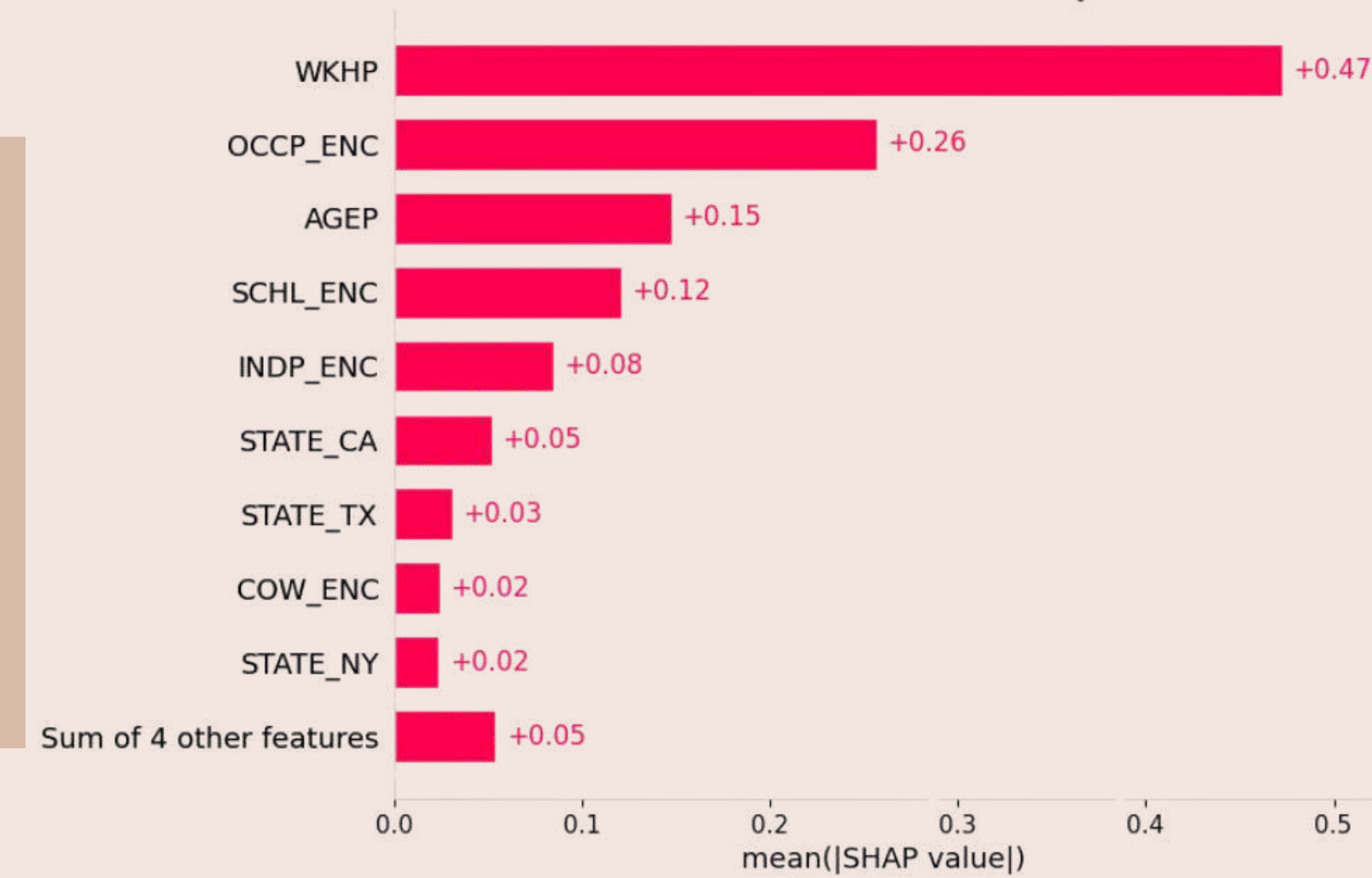
\$52,566.60

Women earn 7.38% MORE than men with identical features.

Male Model: Global Feature Importance



Female Model: Global Feature Importance



Additional Insights

HYPOTHESIS TEST: MODEL PREDICTION DIVERGENCE

T-Statistic: **150.4156**

P-Value: **0.0000e+00**

\$62,029.98

Predicted average male salary

\$55,520.38

Predicted average female salary

\$6,509.60

Systematic Gap

HYPOTHESIS TEST: THE MARRIAGE PREMIUM

Male Married vs Others:

t = **160.74**

p = **0.0000e+00**

Female Married vs Others:

t = **79.71**

p = **0.0000e+00**

0.6568

Mean Log Advantage (Male)

0.3553

Mean Log Advantage (Female)

0.3015

Difference in 'Premium'

Mid-Career Age Bonus (Scaled Age 0.2-0.4)

0.1016

Male Avg Age Bonus

0.0840

Female Avg Age Bonus

-0.0176

Difference

Comparing accuracy

XGBoost slightly outperforms Neural Networks while remaining significantly more interpretable.

XGBoost

78.85%
Accuracy

80.95%
Precision

78.83%
Recall

79.88%
F1 - score

Neural Network

77.29%
Accuracy

81.50%
Precision

75.78%
Recall

78.54%
F1 - score

Confusion Matrix: 50000 Threshold



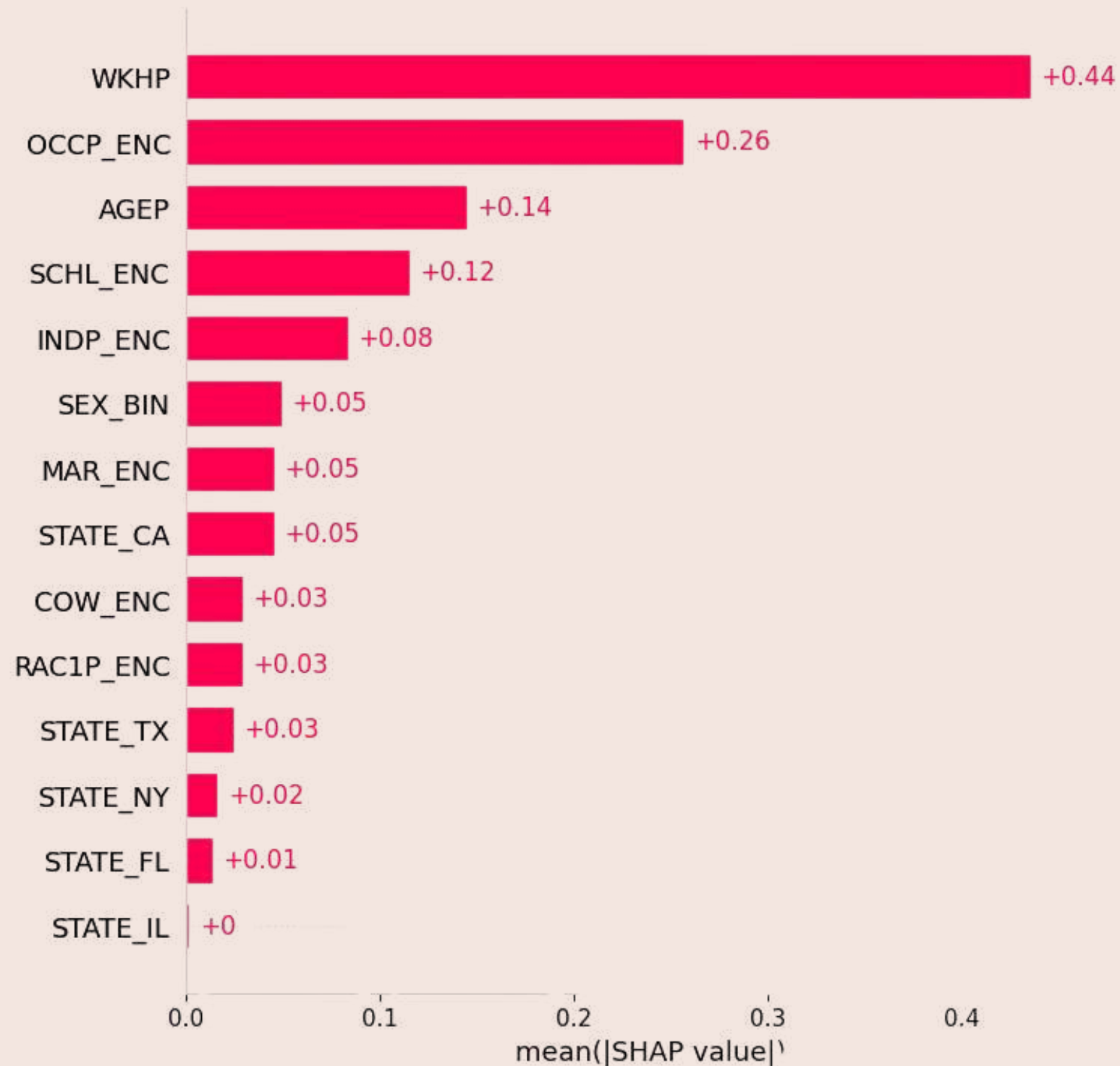
Confusion Matrix: Predicting Salary > \$50k



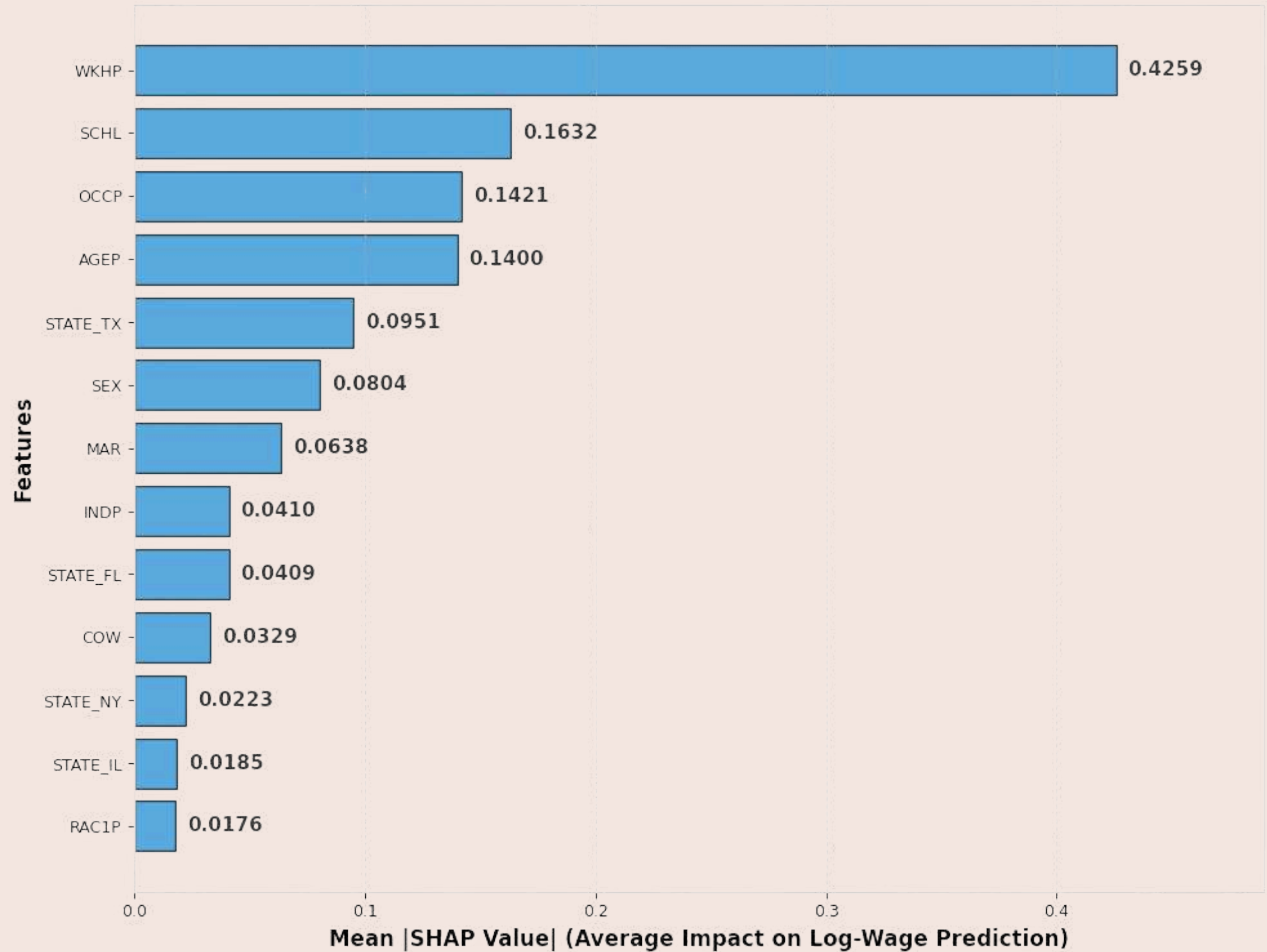
Comparing explainability

XGBoost produces more concentrated and interpretable feature importance than Neural Networks.

XGBoost

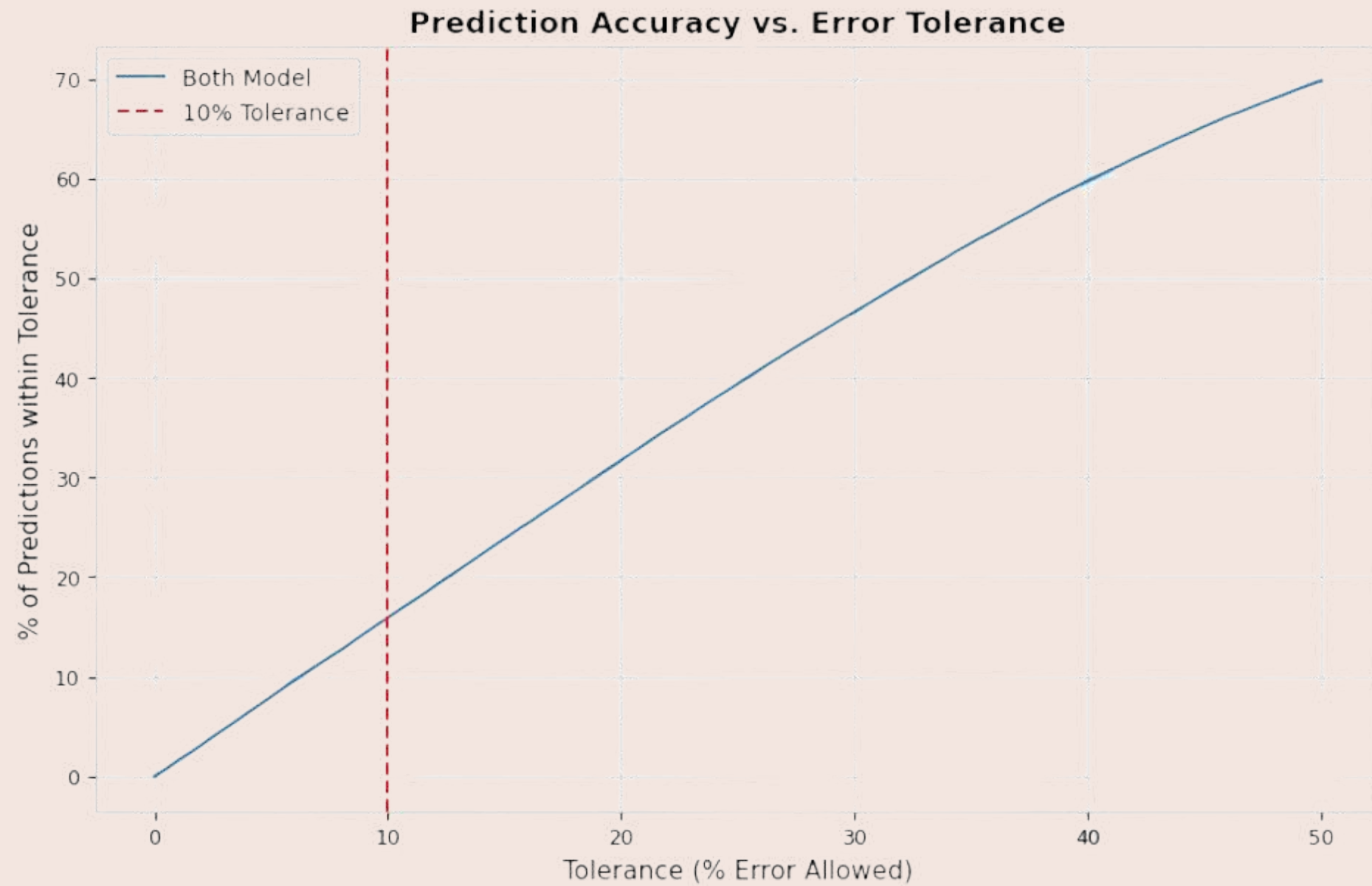


Neural Network

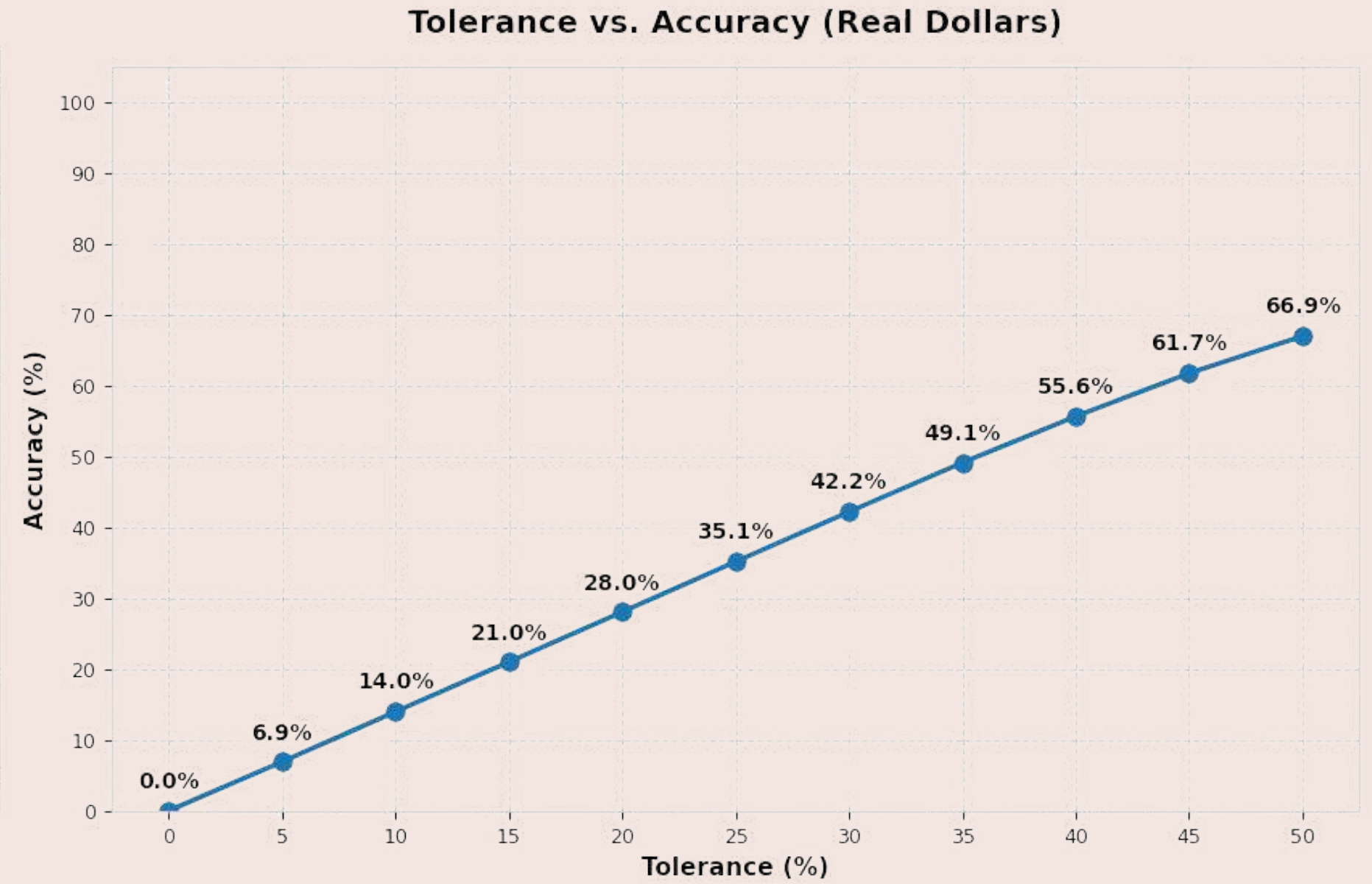


Comparing explainability

XGBoost



Neural Network



Literature review for PROMOTION

Goal: ML framework which accurately predicts employee promotions using performance and HR-related indicators.

DATASET

54, 808 TRAINING RECORDS

MODEL

XGBoost Regression, Random Forest

OBJECTIVE

Predict employee promotions using XGBoost, HR performance metrics, training scores, ratings, tenure, department, and employee attributes.

LIMITATION

Does not examine how gender parity exists at the highest levels of leadership that too across different industries. Predicts promotions but does not analyze gender disparities across wages, industries, class-of-worker groups, or elite leadership positions.

Z. Xiong and W. Lu, "Employee Promotion Prediction: Machine Learning in HR Management," ResearchGate, Dec. 2025. [Online].

Y. A. Ibrir and M. Çavur, "Forecasting Employees' Promotion Based on Personal Indicators by Using a Machine Learning Algorithm," International Journal of Management Information Systems and Computer Science, vol. 8, no. 2, pp. 75-98, Dec. 2024, doi: 10.33461/uybisbbd.1471499.

Methodology

Classify “elite” leaders

- Using occupation titles and industry codes along with searching for terms like “executive”, “manager” to define someone as a leader.
- Using salary above 80th percentile in an industry group to be leader

XGBoost + Metrics

- Features - WKHP, AGE, MAR, SEX, SCHL, COW
- Accuracy : 0.9331
- Precision: 0.1162
- Recall : 0.2872
- F1 Score: 0.3932
- AUC: 0.9447
- Tolerance: 0.65

Limitations

- SMOTE, undersampling
- Models trained for different industries
- Dataset columns dont explain promotion well

So how can we use this information for Plaksha?

Internal Pay Equity (Faculty & Staff Auditing)



The Problem: Traditional payroll analysis struggles to separate valid salary drivers from implicit biases.

The SHAP Solution: Plaksha can run this exact model on its own payroll. By holding variables like tenure, research output, courses taught, and department mathematically constant, the university can definitively isolate and eliminate any unexplained wage penalties tied to gender, race, or background.

Placement Fairness & Career Services

OFFICE OF
CORPORATE
PARTNERSHIPS
AND CAREERS

The Problem: Are graduating students receiving fair market value, or are certain demographics being lowballed by recruiters?

The SHAP Solution: By analyzing campus placement data, the model can separate valid predictors (GPA, internship experience, technical stack) from invalid ones.

Actionable Outcome: If the model detects that recruiters are offering lower starting packages to specific groups for the exact same roles, CPC can use this data to renegotiate minimum baseline packages with corporate partners.

POTENTIAL APPLICATIONS OF THE MODEL

Policy & Government Applications

EQUAL PAY RECOGNITION

Detect unexplained wage gaps after controlling for occupation, education, and hours worked.

Supports evidence-based pay legislation

LEADERSHIP DIVERSITY MONITORING

Track executive and board-level gender representation across sectors.

Measures leadership parity over time

INDUSTRY LEVEL INTERVENTION

Identify industries with the highest wage and leadership disparities.

Enables targeted workforce reforms

HIGH-RISK INDUSTRY DETECTION

Identify industries with the largest hidden gender disparities in pay and leadership.

Enables targeted equity interventions

POTENTIAL APPLICATIONS OF THE MODEL

Corporate Implications & Applications

PAY EQUITY AUDITS

Detect hidden wage gaps between employees with similar qualifications and roles.

Supports fair pay decisions.

TRANSPARENT COMPENSATION

Use SHAP to explain which factors drive salary and promotion outcomes.

Improves HR compensation.

LEADERSHIP PIPELINE ANALYSIS

Identify barriers limiting women's progression into senior leadership roles.

Strengthen leadership diversity.

BIAS DETECTION

Detect demographic patterns influencing hiring, promotion, or compensation.

Support inclusive workforce hiring.



Occupation - Agentic AI Intern - Zomato

Age - 19

Marital Status - Not married

Working hours per work - 40

Course - DSEB

Wage - Rs 50,000



Occupation - Finding (hopefully)

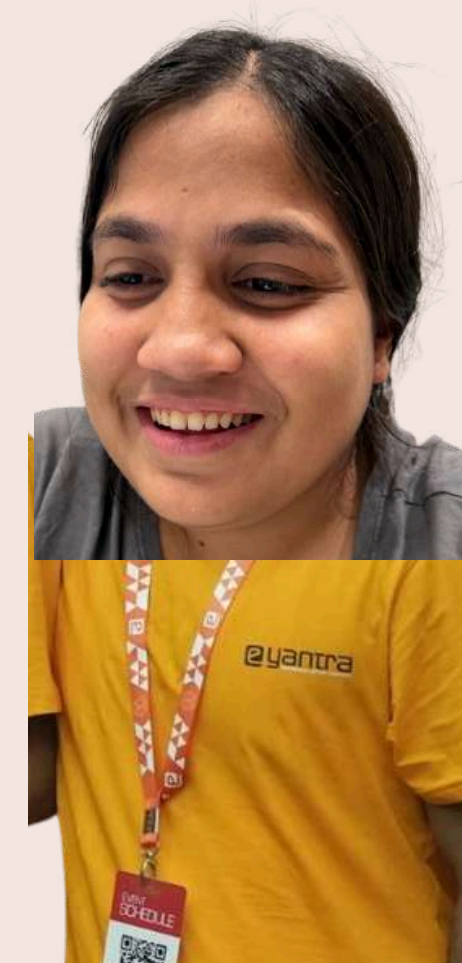
Age - 19

Marital Status - Forever single

Working hours per work - 400

Course - DSEB

Wage - Rs 0



Occupation - Research Internship IIT Bombay

Age - 19

Marital Status - Not married

Working hours per work - 54

Course - CSAI

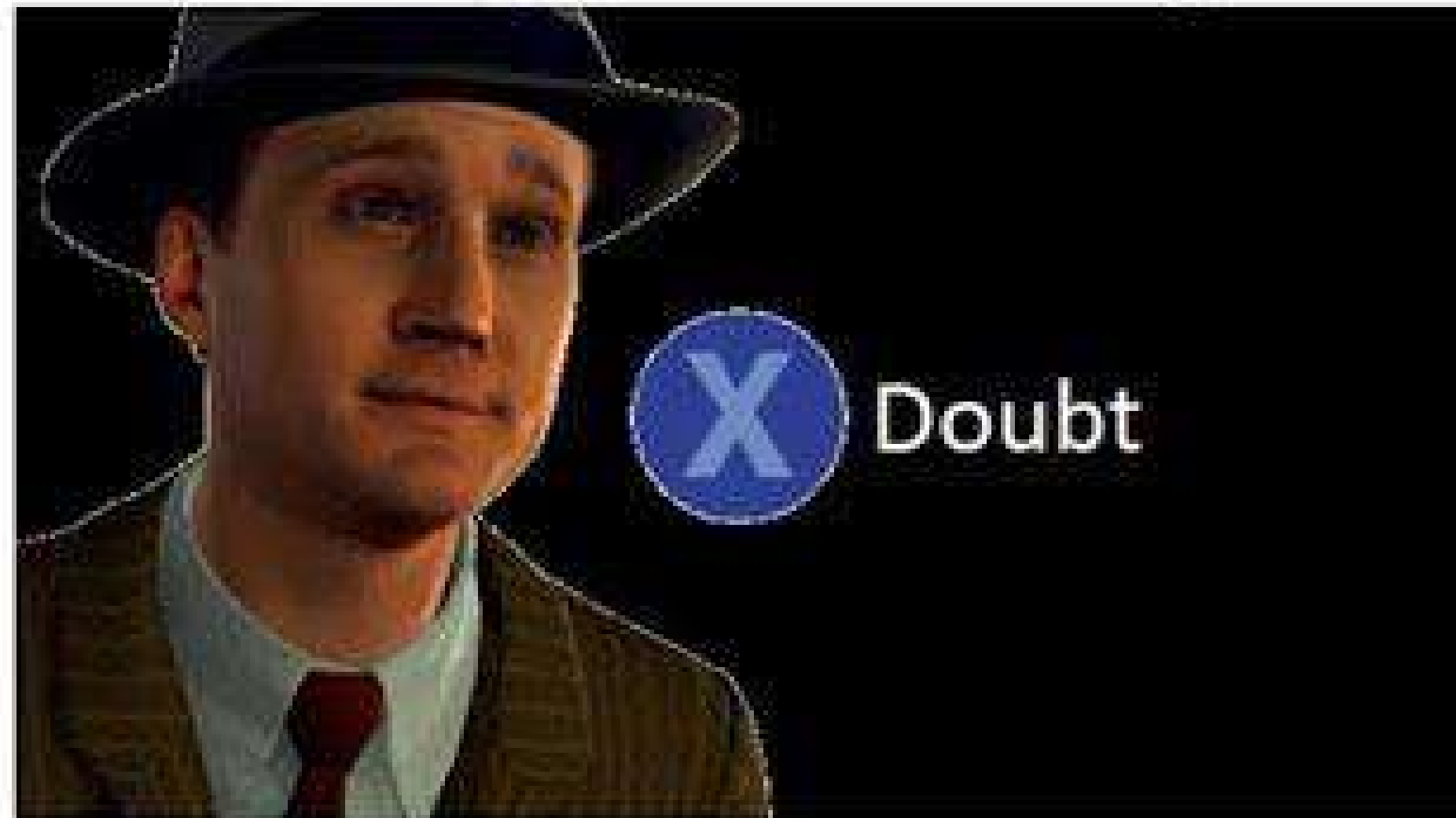
Wage - Rs 12,500

Sahil be like :

Feminist: A woman working
your job would earn 20% less

NO

Me working ~~minimum~~ wage:



REFERENCES:

1. S. Gupta and A. Sharma, "Gender Pay Gap and Workforce Equality Analysis Using Machine Learning," All Commerce Journal, vol. 5, no. 1, pp. 42–48, 2024.
2. S. Barocas, M. Hardt, and A. Narayanan, "Fairness and Machine Learning," arXiv preprint arXiv:2004.13332, 2020.
3. F. Blau and L. Kahn, "The Gender Wage Gap: Extent, Trends, and Explanations," IZA Journal of Labor Economics, vol. 7, no. 1, pp. 1–26, 2018.
4. OECD Strategic Foresight Programme, OECD, accessed May 16, 2026.
5. A. Rahman and S. Ahmed, "Machine Learning Approaches for Wage Prediction and Economic Inequality," Journal of Economics and Political Economy, vol. 10, no. 2, pp. 88–101, 2023.
6. M. Bussmann, "Machine Learning in Human Resource Analytics," in Springer Reference Work Encyclopedia, Springer, 2019, pp. 1–15.
7. Z. Xiong and W. Lu, "Employee Promotion Prediction: Machine Learning in HR Management," UYBIS Bilisim Sistemleri ve Teknolojileri Dergisi, vol. 6, no. 2, pp. 45–55, 2024.
8. Z. Xiong and W. Lu, "Employee Promotion Prediction: Machine Learning in HR Management," ResearchGate, Dec. 2025.
9. A. Glynn, "Women of Color and the Wage Gap," Center for American Progress, Mar. 2023.
10. S. Benard and I. Paik, "The Motherhood Penalty," National Institutes of Health / Annals of the American Academy of Political and Social Science.
11. World Economic Forum, "How to Reduce the Motherhood Penalty and Close the Gender Pay Gap," May 2022
12. M. A. Budig, J. Misra, and I. Boeckmann, "The Motherhood Penalty in Cross-National Perspective," Journal of Marriage and Family, vol. 80, no. 5, pp. 1287–1313, 2018.