



From Bias to Balance: A Data-Driven Approach to Fair Recruitment Practice

JSM 2025

Hairu Fan



COLLEGE OF
**SCIENCE &
ENGINEERING**
CENTRAL MICHIGAN
UNIVERSITY

Presenter



- Hairu Fan, Ph.D. student, Department of Statistics, Actuarial and Data Sciences, Central Michigan University
- Double Major: Applied Statistics & Business Information Systems
- Email: fan2h@cmich.edu

From Bias to Balance: A Data-Driven Approach to Fair Recruitment Practices



COLLEGE OF
**SCIENCE &
ENGINEERING**
CENTRAL MICHIGAN
UNIVERSITY

Agenda

- 1. Introduction**
- 2. Datasets & Data Preprocessing**
- 3. Methods**
- 4. Key Findings**
- 5. Bias Analysis**
- 6. Conclusion**



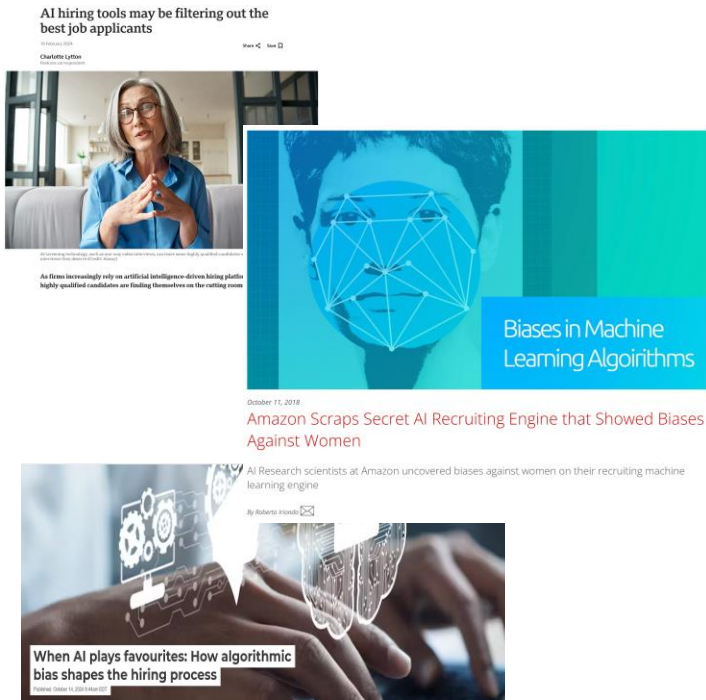
Introduction



COLLEGE OF
**SCIENCE &
ENGINEERING**
CENTRAL MICHIGAN UNIVERSITY

Introduction

Why Fairness Matters in Recruitment?



In recent years, with the widespread use of data-driven recruiting systems, there has been a growing interest in recruiting fairness in both academia and industry.

- Machine learning has been widely used in recruitment, such as resume screening & salary prediction
- Gender, salary bias in historical data may be amplified by algorithms
- Unfairness in recruitment systems may affect business diversity



Introduction

Recently Existing Studies on Algorithmic Bias

- Chandra et al. (2023) found that algorithms may inherit **biases from historical data**, leading to unfairness in the recruitment process.
- Pessach & Shmueli (2022) noted that imbalance in **training data and algorithm goal setting** may affect the fairness of hiring predictions.
- Quer et al. (2024) found that hiring scores may be **biased against candidates younger** than 25 years old. However, most studies still focus on overall bias and lack specific analysis on salary fairness and hiring scores.
- There is a **trade-off between accuracy and fairness**, and as we strive for greater fairness, accuracy may be impacted (Pessach & Shmueli, 2022) .
- Fabris et al. (2023) study notes that while many companies try to reduce bias in hiring algorithms, current technology is still **unable to eliminate** it completely.



Introduction

Recently Existing Studies on Algorithmic Bias

- Machine Learning Recruitment Systems May Have Gender and Salary Bias
- Data bias and model objectives may lead to unfair hiring decisions
- Fairness constraints may reduce bias but impact prediction accuracy.
- Existing research focuses on algorithmic bias and lacks in-depth analysis of salary and hiring scores.



Introduction

Research Questions

Bias in Data

- Are gender, education, and job category distributions imbalanced in the dataset?
- Do these factors significantly impact salary and hiring scores?

Bias in Algorithms

- Do machine learning models amplify or inherit historical biases in the data?

Fairness Constraints (Future Research)

- Can fairness modeling strategies be used to reduce gender bias in hiring systems?
- What is the trade-off between fairness and prediction accuracy?



Introduction

Research Objectives

1. Quantifying Bias in Salary & Hiring Decisions

- Analyze gender, education, and job category influences on salary and hiring scores.
- Examine salary distribution across industries and education levels.

2. Investigating Salary Growth & Occupational Segregation

- Assess how salary growth trends vary across genders and industries.
- Explore the role of occupational segregation in salary discrepancies.

3. Provide data-driven recommendations for fair recruitment

- Identify bias patterns in recruitment systems.
- Propose algorithmic fairness methods for mitigating bias (future research).



Datasets & Data Preprocessing



COLLEGE OF
**SCIENCE &
ENGINEERING**
CENTRAL MICHIGAN UNIVERSITY

Datasets Overview

Data Sources

- This study uses three public datasets provided by UCI, GitHub, and ProPublica for analyze bias in the hiring process.
- The data source and variable descriptions are as follows:

Datasets	Description	Data Sources
Adult Income	There are a total of 15 variables, including 48,842 records of personal income, age, gender, occupation, education, and other factors. Missing value shows '?'. 	UCI Machine Learning Repository
Job Salary	6,704 salary forecast data with 6 key variables, 17 records have missing values. 	GitHub Repository
COMPAS	60,843 recruitment scoring/ risk prediction data, 28 variables. 45,240 records have missing values (mostly MiddleName)	ProPublica GitHub



Datasets Overview

Key Variable Example: Adult Income

age	workclass	education	educational-num	occupation	gender	hours-per-week	native-country	income
25	Private	11th	7	Machine-op-inspct	Male	40	United-States	<=50K
38	Private	HS-grad	9	Farming-fishing	Male	50	United-States	<=50K
28	Local-gov	Assoc-acdm	12	Protective-serv	Male	40	United-States	>50K
44	Private	Some-college	10	Machine-op-inspct	Male	40	United-States	>50K
18	?	Some-college	10	?	Female	30	United-States	<=50K

Key Variable Example: Job Salary

Age	Gender	Education Level	Job Title	Years of Experience	Salary
32.0	Male	Bachelor's	Software Engineer	5.0	90000.0
28.0	Female	Master's	Data Analyst	3.0	65000.0
45.0	Male	PhD	Senior Manager	15.0	150000.0
36.0	Female	Bachelor's	Sales Associate	7.0	60000.0
52.0	Male	Master's	Director	20.0	200000.0

Key Variable Data Example: COMPAS

Sex_Code_Text	DateOfBirth	DecileScore	RawScore	ScoreText	RecSupervisionLevelText	RecSupervisionLevelText
Male	12/05/92	4	-2.08	Low	Low	Low
Male	12/05/92	2	-1.06	Low	Low	Low
Male	12/05/92	1	15.0	Low	Low	Low
Male	09/16/84	2	-2.84	Low	Low	Low
Male	09/16/84	1	-1.5	Low	Low	Low

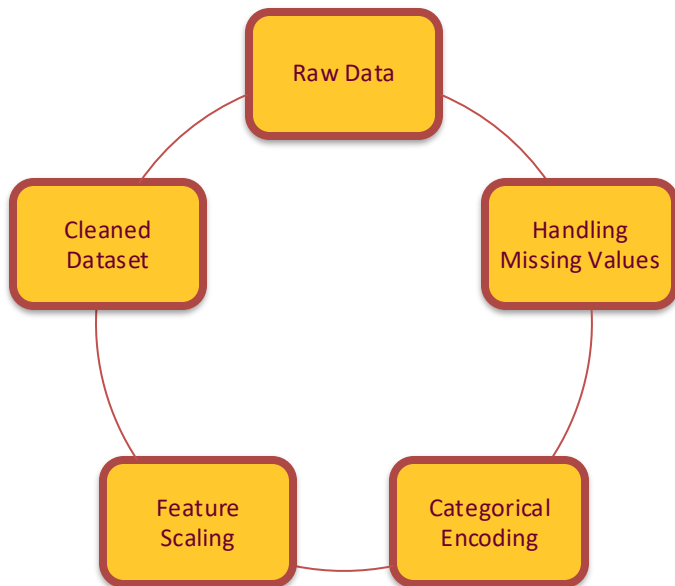
From Bias to Balance: A Data-Driven Approach to Fair Recruitment Practices



COLLEGE OF
**SCIENCE &
ENGINEERING**
CENTRAL MICHIGAN UNIVERSITY

Data Preprocessing

Data Cleaning



Step 1: Handle Missing Values

Remove incomplete records less than 5%
Fill missing categorical data with mode imputation

Step 2: Convert Categorical Variables

Gender, Income → Binary encoding: Male (0), Female (1)

Job Categories, Workclass → Reclassification reducing category sparsity & One-Hot Encoding

Education → Ordinal

Step 3: Normalize Data

Winsorization handles abnormal data
Hiring scores normalization
Standardized salary & experience for consistency

Step 4: Unify all the variable name



Data Cleaning Results

Adult Income

- 42166 records
- 19 columns
- Data type: int
- No missing values

Job Salary

- 6684 records
- 15 columns
- Data type: int
- No missing values

COMPAS

- 60842 records
- 6 columns
- Data type: int
- No missing values



Data Preprocessing

Summary of Data Preparation

- The final dataset has balanced categorical distributions.
- Data cleaned, missing values addressed, categorical variables encoded, such as gender and job category variables.
- Final dataset missing values handled, data standardized, and ready for bias analysis.



Methods



COLLEGE OF
**SCIENCE &
ENGINEERING**
CENTRAL MICHIGAN UNIVERSITY

Statistical Analysis Methods

Bias in Data Analysis

- Descriptive Statistics
- Correlation Analysis
- T-Test
- Chi-Square
- ANOVA

Regression Analysis

- Ordinary Least Squares
- Logistic Regression
- Ordinal Logistic Regression



Methods

Fairness- Aware Modeling(Future Research)

Fairness Metrics

Future research will focus on the fairness evaluation.

Data level

- Disparate Impact
- Kolmogorov- Smirnov (KS) Test

Model level

- Equal Opportunity
- Equalized Odds

Fairness Optimization Methods

Future research will focus on the practical effects of fairness optimization. Three types of fairness optimization methods research program:

Data Preprocessing

- Reweighting
- Feature De-biasing

In – Processing

- Fairness Regularization
- Equalized Odds Constraints

Post- Processing

- Output Calibration
- Result Re-ranking



Key Findings



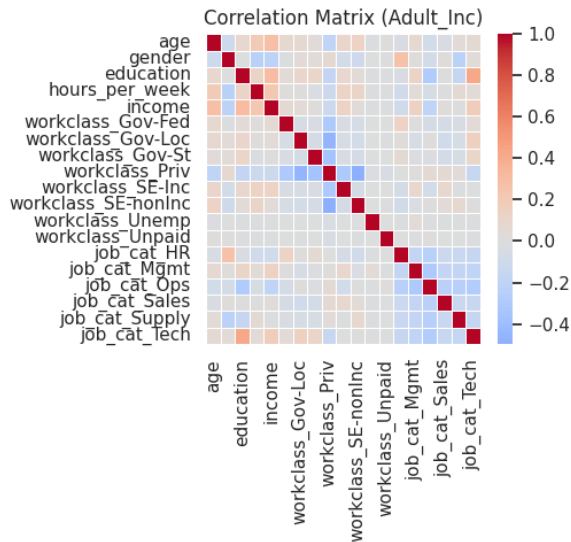
COLLEGE OF
**SCIENCE &
ENGINEERING**
CENTRAL MICHIGAN UNIVERSITY

Key Findings

Bias in Data Analysis : Correlation Analysis

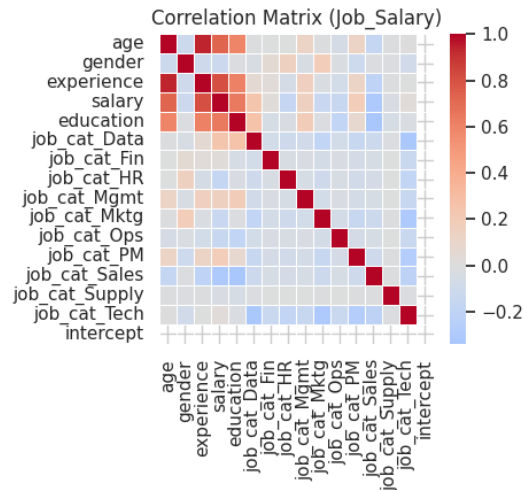
Adult_Inc:

- Education appears to be a major factor influencing income, so salary models should account for it.
- Gender does not show a strong direct effect, but further analysis is needed.



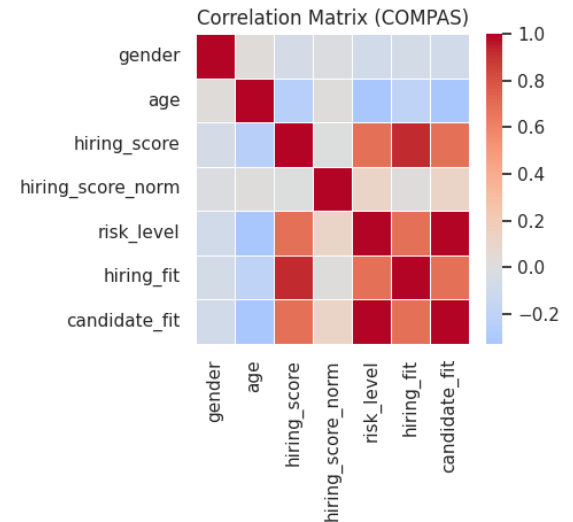
Job_Salary:

- Experience is a stronger predictor of salary than education.
- Job category has a significant influence on salary



COMPAS:

- Potential bias in hiring scores against older candidates.
- Hiring fit is heavily influenced by risk level.



Key Findings

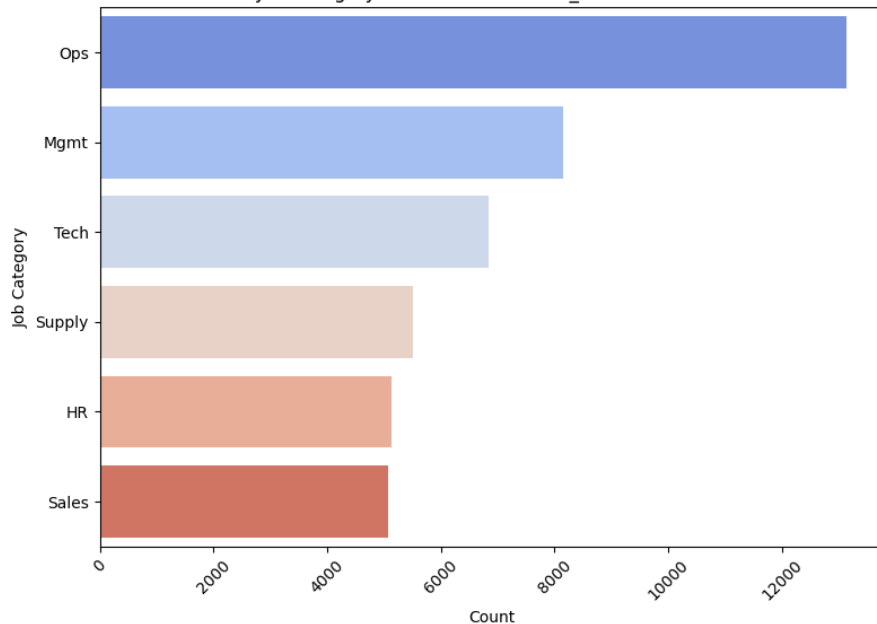
Bias in Data Analysis : Descriptive Statistics Analysis

Job Category Distribution

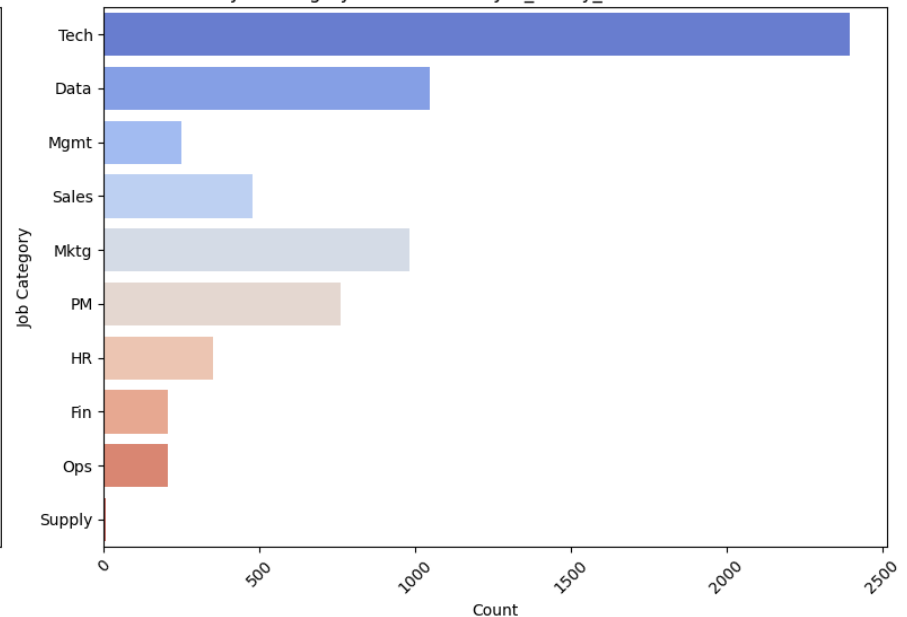
Adult_Income dataset: Operations & Administration (Ops) is the most common job category, followed by Marketing & Advertising (Mktg) and Tech & Engineering (Tech).

Job_Salary dataset: Tech is the most common category, followed by Data & Analytics (Data) and Mktg.

Job Category Distribution in Adult_income dataset



Job Category Distribution in Job_Salary_Cleaned Dataset



From Bias to Balance: A Data-Driven Approach to Fair Recruitment Practices



COLLEGE OF
SCIENCE & ENGINEERING
CENTRAL MICHIGAN UNIVERSITY

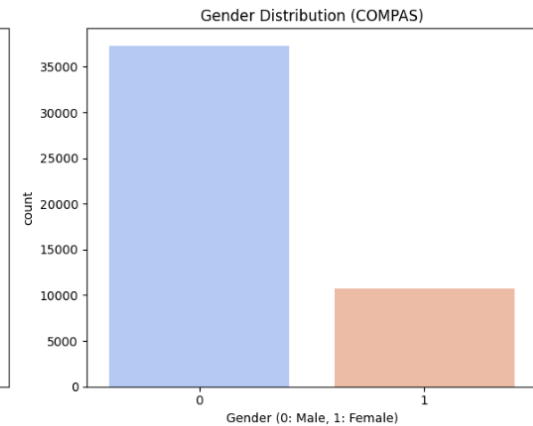
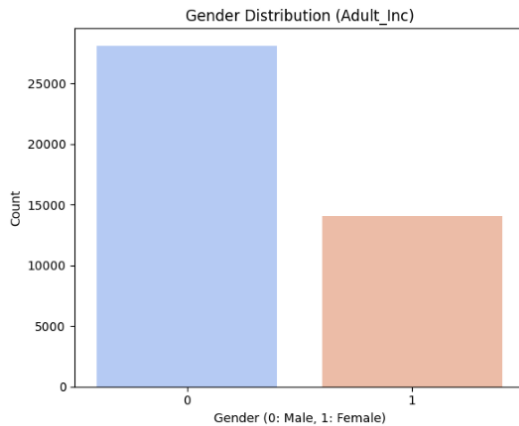
Key Findings

Bias in Data Analysis : Descriptive Statistics Analysis

Gender Distribution

Adult_Income & COMPAS dataset: The gender (0=Male, 1=Female) ratio is more unbalanced (higher for males).

Job_Salary dataset: The gender ratio is relatively balanced.



From Bias to Balance: A Data-Driven Approach to Fair Recruitment Practices



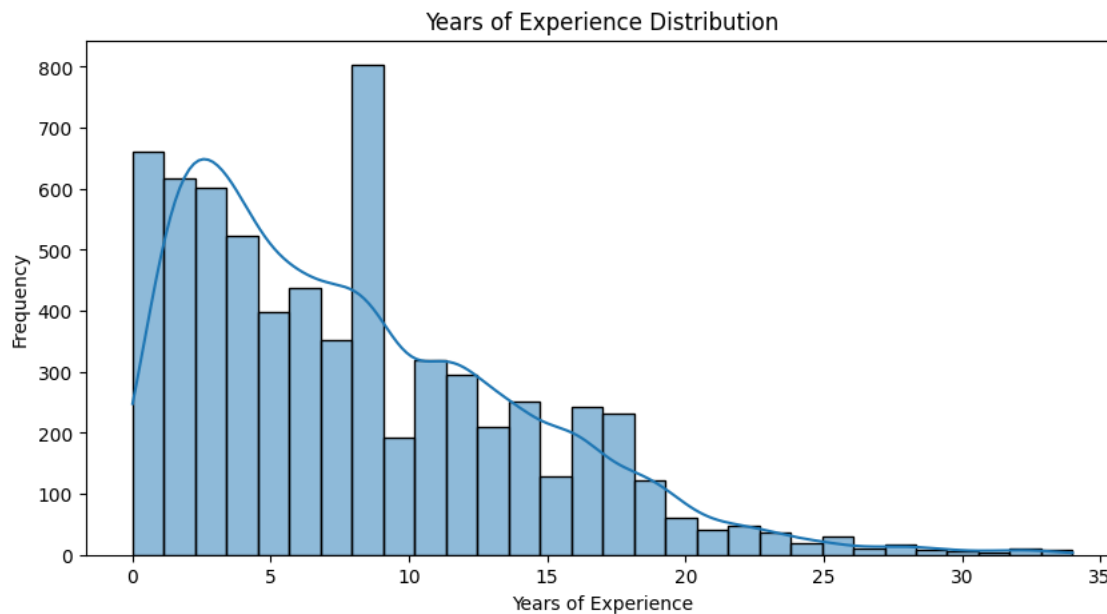
COLLEGE OF
**SCIENCE &
ENGINEERING**
CENTRAL MICHIGAN UNIVERSITY

Key Findings

Bias in Data Analysis : Descriptive Statistics Analysis

Year of Experience Distribution

Job_Salary dataset: Right-skewed (positively skewed) distribution.



From Bias to Balance: A Data-Driven Approach to Fair Recruitment Practices



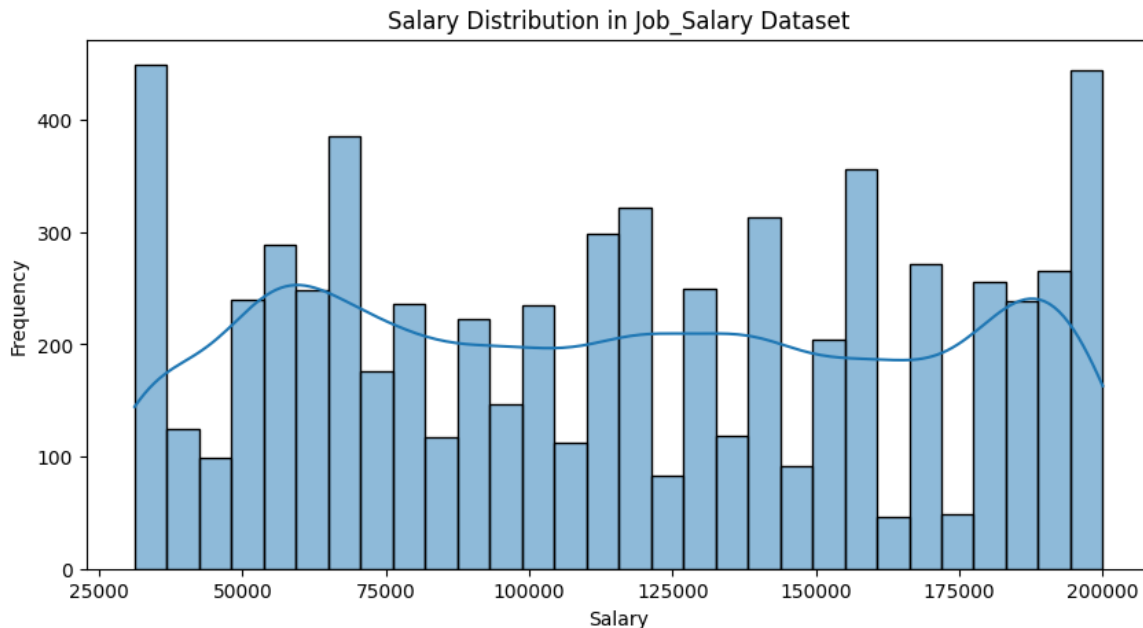
COLLEGE OF
**SCIENCE &
ENGINEERING**
CENTRAL MICHIGAN UNIVERSITY

Key Findings

Bias in Data Analysis : Descriptive Statistics Analysis

Salary Distribution

Job_Salary dataset: Relatively scattered, with multiple peaks and large differences in salary levels.



Key Findings

Bias in Data Analysis : Descriptive Statistics Analysis/ ANOVA

Education Level Distribution

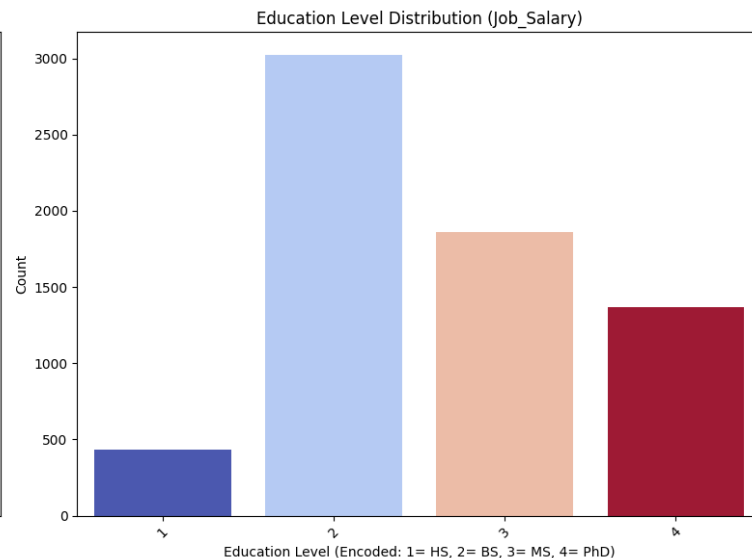
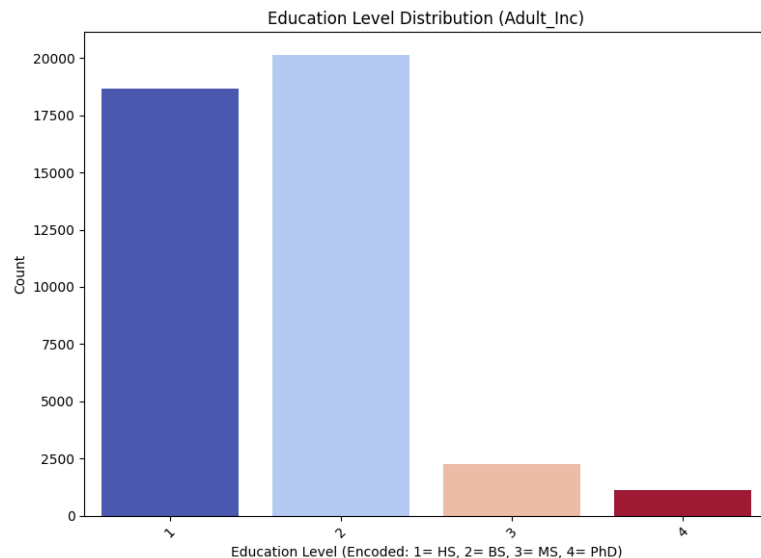
Adult_Income dataset: The educational background is mainly high school (HS) and bachelor's degree (BS).

Job_Salary dataset:

BS is the largest, while the proportion of master's degree (MS) and doctorate (PhD) is relatively high.

F-statistic > 1,000 → Large F-value, indicating strong between-group variance.

p-value < 0.01, education level significantly affects salary.



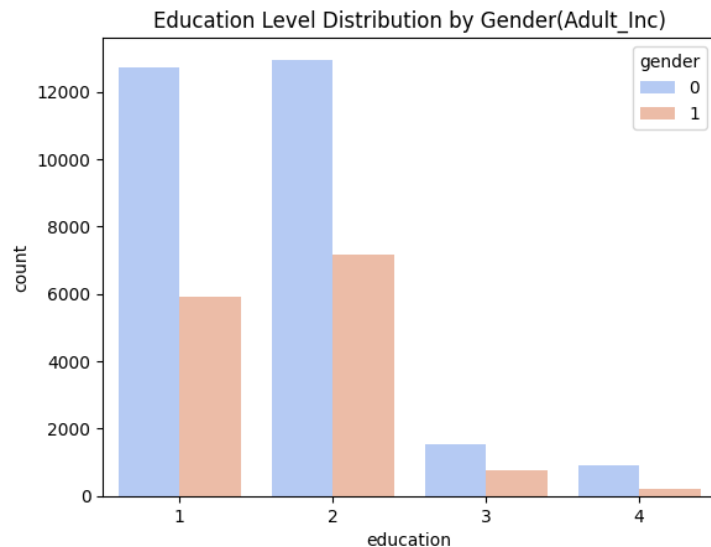
Key Findings

Bias in Data Analysis : Descriptive Statistics Analysis

Education Distribution by Gender

Adult_Income dataset: Male (0) dominate all education levels, especially at lower levels. Few individuals, regardless of gender, have the highest education levels (PhD).

Job_Salary dataset: Female have higher representation in BS and MS In contrast to Adult_Inc, the gender gap in Master's and PhD degrees is smaller, possibly due to dataset selection.



From Bias to Balance: A Data-Driven Approach to Fair Recruitment Practices



COLLEGE OF
SCIENCE & ENGINEERING
CENTRAL MICHIGAN UNIVERSITY

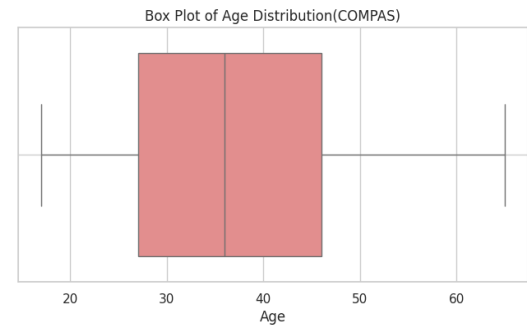
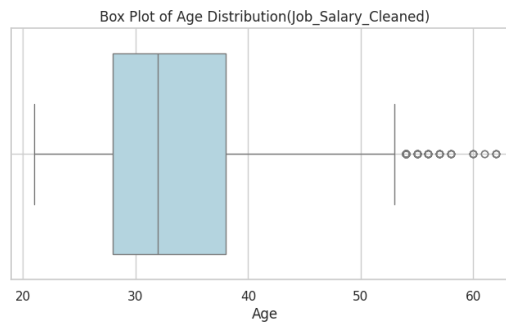
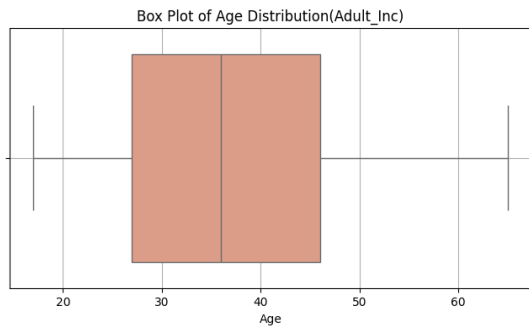
Key Findings

Bias in Data Analysis : Descriptive Statistics Analysis

Age Distribution

Adult_Inc & COMPAS dataset : Similar age distributions.

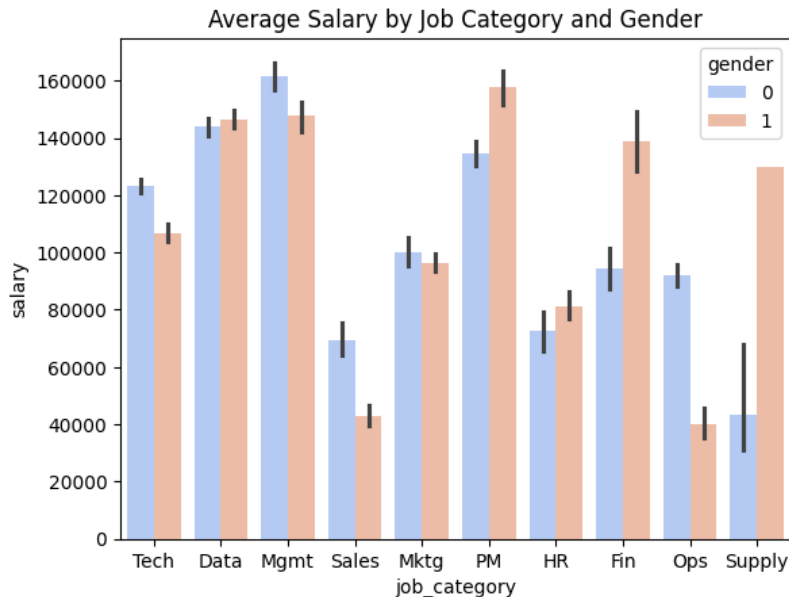
Job_Salary dataset: younger workforce with some outlier.



Key Findings

Bias in Data Analysis : Descriptive Statistics Analysis/T - Test

Salary Distribution by Gender



T-Test (p-value < 0.001)

1. Gender Pay Gap Varies by Job Category

Tech, Data, Mgmt, and Project & Product Management (PM) :
Female salaries closely match or exceed male salaries.
PM is the only category where women earn more on average than men.

2. Male Salaries Are More Dispersed, With Higher Peaks

Male have a wider salary distribution, more high-earning.
Female salaries are more concentrated, fewer reach top salaries .

3. Significant Male Salary Advantage in Certain Fields

Sales & Business Development (Sales), Ops, and Supply Chain & Logistics (Supply) :
Male salaries are significantly higher than female salaries.
This could indicate occupational segregation, where male dominate leadership roles.

4. Female Salary Advantage in Finance & HR

Fin and HR :
Female earn more on average than male.
These fields tend to have higher female representation, which may contribute to this trend.



Key Findings

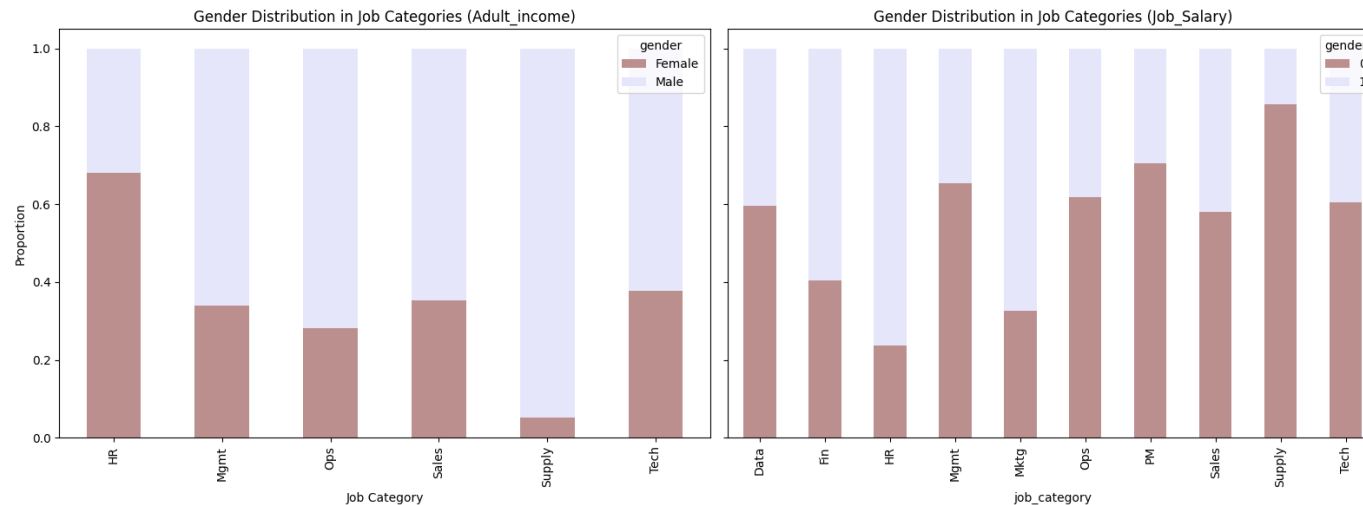
Bias in Data Analysis : Descriptive Statistics Analysis

Job Category by Gender

Adult_Income Dataset: Human Resources(HR) has the highest proportion of females (~65%). Supply Chain & Logistics (Supply) is male-dominated (~90% male).

Job_Salary Dataset: The gender distribution is more even.

- Tech and Supply Are Strongly Male-Dominated
- HR and Finance & Accounting (Fin) Have More Female Representation



Key Findings

Bias in Data Analysis : Descriptive Statistics Analysis/ANOVA

Job Category by Gender

Job_Salary Dataset:

F-statistic > 250 → Large F-value, at least one job category has a significantly different salary distribution
p-value < 0.01, job categories significantly impact salary levels.



From Bias to Balance: A Data-Driven Approach to Fair Recruitment Practices



COLLEGE OF
SCIENCE & ENGINEERING
CENTRAL MICHIGAN UNIVERSITY

Key Findings

Bias in Data Analysis : Descriptive Statistics Analysis/ ANOVA

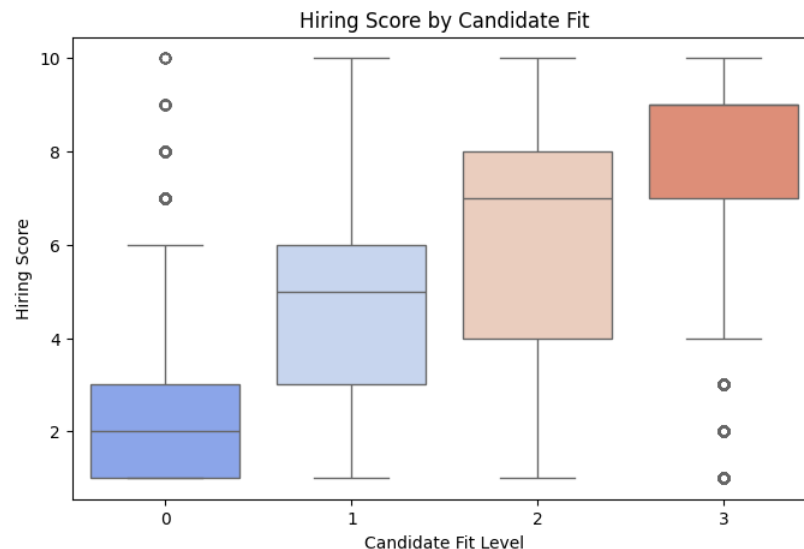
Hiring Score Differences Across Risk Levels

COMPAS dataset:

F-statistic > 1,4000 → Large F-value, indicating strong between-group variance.

p-value < 0.01, hiring scores differ significantly across different risk levels.

This suggests that higher-risk candidates may receive systematically lower hiring scores.



Key Findings

Bias in Data Analysis : Descriptive Statistics Analysis/T - Test

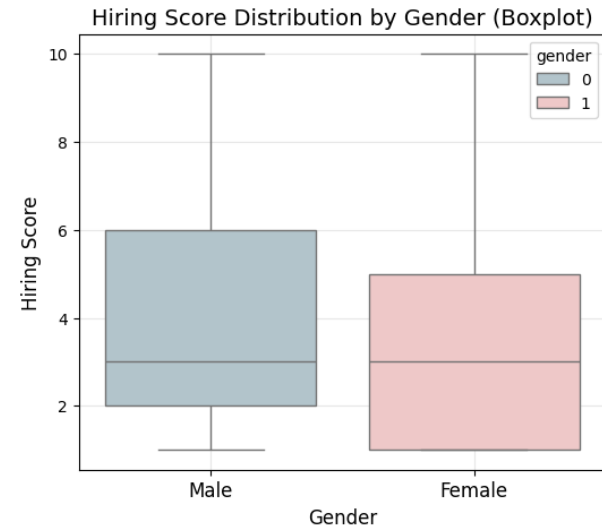
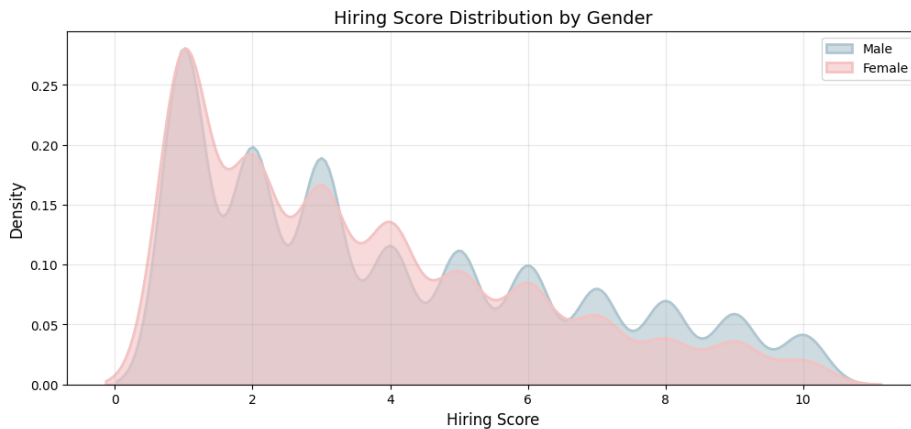
Hiring Scores vs. Gender Analysis

COMPAS dataset:

Blue (male) and pink (female) represent the probability density distribution of recruitment scores for both sexes. This shows that males have a higher overall performance in recruitment scores.

Although the median score for female is slightly higher, the proportion of high scores is smaller, while the proportion of high scores for male is larger.

T- Test: P-value < 0.001, there are significant differences in recruitment scores by gender.



From Bias to Balance: A Data-Driven Approach to Fair Recruitment Practices



COLLEGE OF
**SCIENCE &
ENGINEERING**
CENTRAL MICHIGAN UNIVERSITY

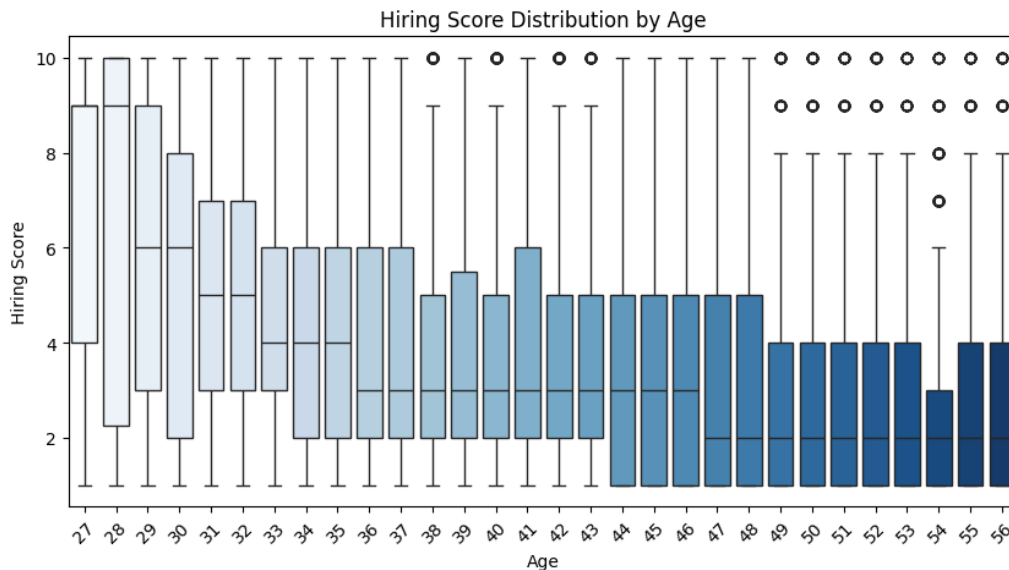
Key Findings

Bias in Data Analysis : Descriptive Statistics Analysis

Age vs. Hiring Score

COMPAS dataset:

Younger candidates have a wider distribution of hiring scores, while older candidates have more concentrated scores.
Lower scoring age groups (over 30) are more stable than younger people



Bias Analysis



COLLEGE OF
**SCIENCE &
ENGINEERING**
CENTRAL MICHIGAN UNIVERSITY

Key Findings

Analysis of Sources of Bias: Ordinary Least Squares(Ols) Model

Gender vs. Salary bias

Job_category dataset:

Model 1: Use of gender only to predict salary

Salary ~ gender, $P < 0.001$, $R^2 < 0.02$.

Gender has a significant effect on salary, but R^2 is extremely low ability to explain salary changes by gender variables

Model 2: Control of education and experience

Salary ~ gender + education + experience, $P < 0.001$, $R^2 > 0.06$.

Gender still has a significant impact on salaries, but education and experience can partially reduce the gap. Education and experience are the main salary determinants.

Model 3: Controlling for education, experience and job_category

salary ~ gender + education + experience + job_category, $P > 0.05$, $R^2 > 0.07$.

The effect of gender on salary is no longer significant after controlling for job category.

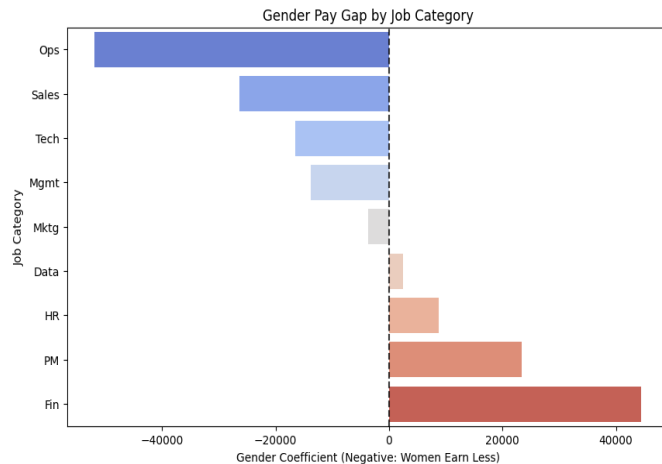
This suggests that the gender salary gap arises mainly from differences in job categories rather than gender.



Key Findings

Analysis of Sources of Bias: Ordinary Least Squares(OLS) Regression Model

Gender salary differentials in Job Category



job_category	gender_coef	p_value
Tech	-	<0.001
Data	+	>0.001
Mgmt	-	<0.001
Sales	-	<0.001
Mktg	-	>0.001
PM	+	<0.001
HR	+	>0.001
Fin	+	<0.001
Ops	-	<0.001

Job_category dataset:

Bar Plot : The salary difference associated with gender.

Negative values (blue bars) indicate female earn less than male. Positive values (red bars) indicate female earn more than male.

OSL Model: salary ~ gender for each industry, grouped by job_category.

Females face the largest salary gaps in Ops, Sales, and Tech.

Mktg and Data show no significant pay gap.

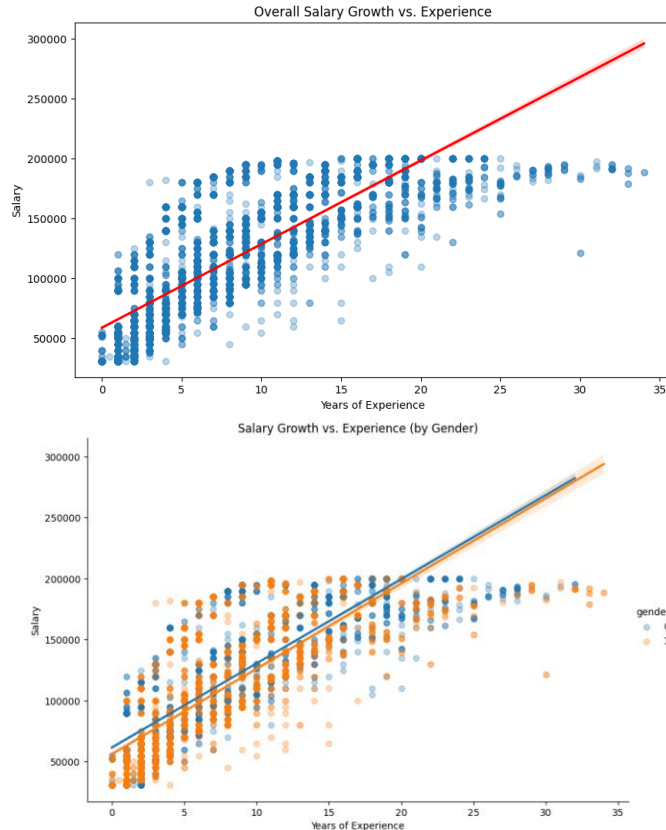
Fin and PM favor women, with significantly higher salaries for them.



Key Findings

Analysis of Sources of Bias: Ordinary Least Squares(OLS) Regression Model

Salary Growth vs. Year of Experience



Job_category dataset:

Overall Salary vs. Experience : Salary growth is linear, with higher experience leading to higher salaries, but salary growth slows down after 20 years.

Salary growth by gender: Experience strongly impacts salary for both of them. Female (orange) and Male (blue) salary growth curves are similar but start at different points (lower for females)

OSL Regression: The salary gap is persistent and could be due to other factors, such as job category or promotions.

Model 1: salary ~ experience, $P < 0.001$, $R^2 > 0.06$. Suggests that the effect of experience on salary is highly significant.

Model 2: salary(gender = male) ~ experience, $P < 0.001$, $R^2 > 0.06$.

Model 3: salary(gender = female) ~ experience, $P < 0.001$, $R^2 > 0.06$.

From Bias to Balance: A Data-Driven Approach to Fair Recruitment Practices



COLLEGE OF
**SCIENCE &
ENGINEERING**
CENTRAL MICHIGAN UNIVERSITY

Key Findings

Analysis of Sources of Bias: Polynomial Regression Model

Non-linear Trends in Salary Increases

Job_category dataset:

Polynomial Regression to analyze the effect of experience² on salary growth.

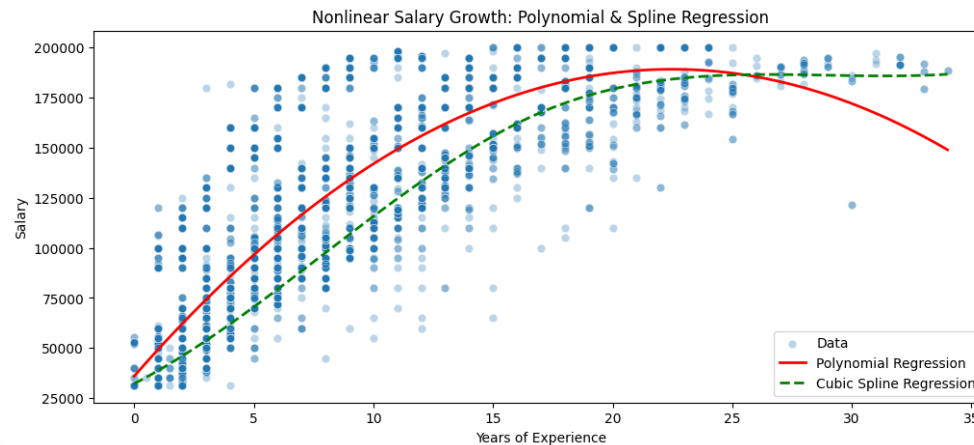
Cubic Spline to get a smoother nonlinear trend.

Polynomial Regression Model : Salary ~ Experience + Experience², $P < 0.001$, $R^2 > 0.07$. Suggests that the effect of experience on salary is highly significant.

Experience Early stage (0-10 years): Salary growth is rapid and the curve is steep.

Experience Mid-term (10-20 years): Salary growth slows down and the growth rate decreases.

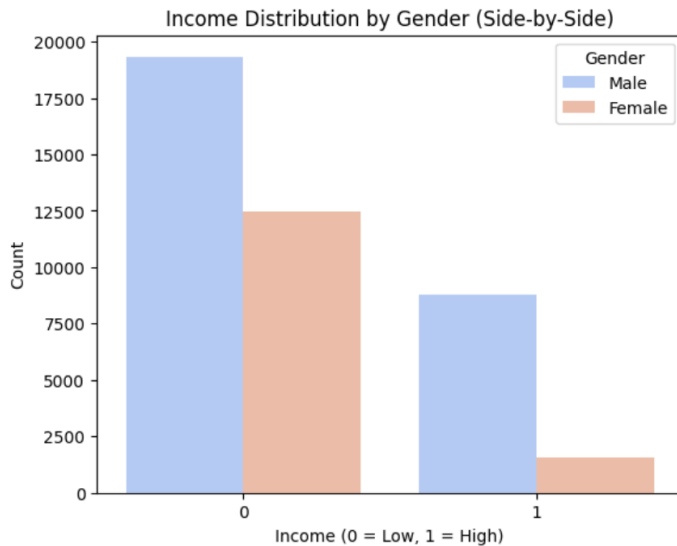
Experience Late stage (20+ years): Salary growth tends to be stable or even slightly declines.



Key Findings

Analysis of sources of bias: Logistic Regression Model

Gender vs. Income



Adult_ income dataset:

Logistic Regression Model :

income ~ gender

The gender coefficient is negative and highly significant ($p < 0.001$), indicating that women are less likely to earn a high income compared to men.

Odds Ratio around 0.3, female are about 70% less likely to earn a high income compared to male.



Key Findings

Analysis of Sources of Bias: Logistic Regression Model

The Impact of Education vs. Gender on Income

Adult_ income dataset:

Logistic Regression Model :

income ~ gender + education

The negative effect is still significant ($p < 0.001$), indicating that even after controlling for education level, it is still more difficult for women to earn high income than men.

Higher education significantly increases the probability of earning a high income ($OR > 2.5$). Education explains part of the income gap, but not all of it.



Key Findings

Analysis of Sources of Bias: Logistic Regression Model

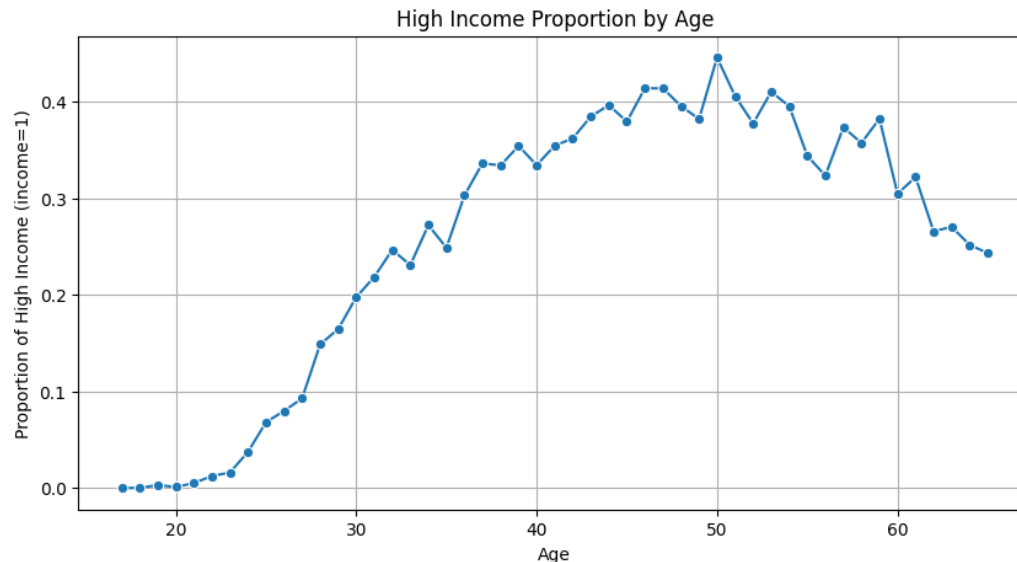
Age vs. Income Analysis

Adult_ income dataset:

Logistic Regression Model: Analyze the relationship between age and high income (income=1)

Model 1: $\text{Income}(\text{income}=1) \sim \text{age}$, $p < 0.001$, Age is significantly positively correlated with high income

Model 2: $\text{Income}(\text{income}=1) \sim \text{age} + \text{gender}$, $p < 0.001$, After controlling for age, gender still significantly affects income, and women are significantly less likely to enter the high-income group than men.



Key Findings

Analysis of Sources of Bias: Ordinal Logistic Regression

Hiring Scores vs. Gender, Age, Hiring_fit, Candidate_fit

COMPAS dataset:

Ordinal Logistic Regression Model: Hiring Scores ~ Gender + Age + Hiring_fit + Candidate_fit

- gender, $p < 0.05$
- age, $p < 0.001$, the older the age, the lower the recruitment score.
- hiring_fit, $p < 0.001$, indicating that candidates with high recruitment fit are more likely to receive higher scores
- candidate_fit, $p < 0.001$, also has a significant impact on the recruitment score.

The recruitment score is most affected by hiring_fit and candidate_fit.

Older candidates have significantly lower scores than younger ones



Key Findings

Analysis of Sources of Bias: Ordinal Logistic Regression

Hiring Scores vs. Gender, Age, Hiring_fit, Candidate_fit

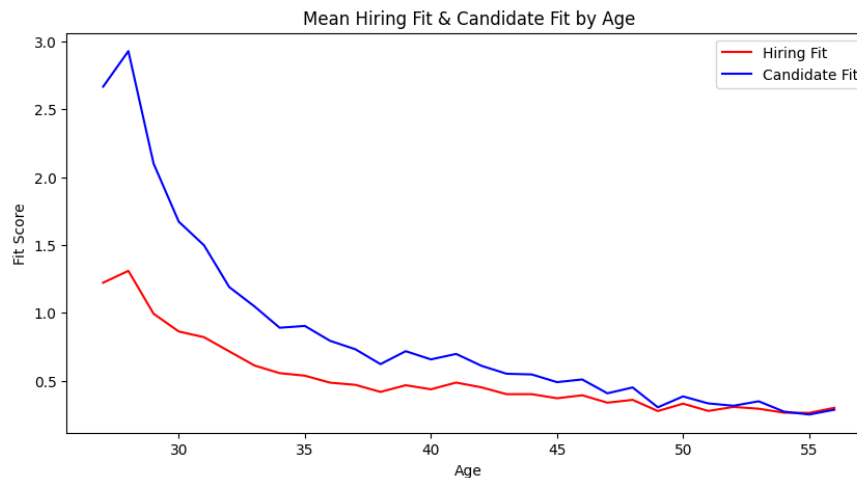
COMPAS dataset:

Model1 : hiring_fit ~ age

Age was significantly and negatively related to fitness for recruitment, $p < 0.001$

Model2 : candidate_fit ~ age

$p < 0.001$, age has a stronger negative effect on candidate fit than recruitment fit



Conclusion



COLLEGE OF
**SCIENCE &
ENGINEERING**
CENTRAL MICHIGAN UNIVERSITY

Contributions

- Empirical Bias Analysis, study quantifying salary bias across multiple datasets, and analyzing gender salary gaps
- Hiring Score Evaluation evaluated hiring score biases, identifying potential inconsistencies in AI-driven recruitment models
- Career Growth Modeling, understanding long-term salary trends, offering to track salary progression disparities over time.
- Beyond academic contributions, our findings hold practical implications for industry professionals, policymakers, and AI developers.

From Bias to Balance: A Data-Driven Approach to Fair Recruitment Practices



COLLEGE OF
**SCIENCE &
ENGINEERING**
CENTRAL MICHIGAN
UNIVERSITY

Further Analysis

- Current research mainly measures gender bias in salary and recruitment scores, but has not yet achieved a comprehensive fairness assessment.
- Use fairness measurement methods (Equalized Odds and Disparate Impact) to further verify whether the model has discriminatory predictions.
- Use Fairness optimization strategies (Reweighting)
- Optimizing Algorithmic Fairness (Data Augmentation)
- Comparison of the effectiveness of different fairness optimization methods

From Bias to Balance: A Data-Driven Approach to Fair Recruitment Practices



COLLEGE OF
**SCIENCE &
ENGINEERING**
CENTRAL MICHIGAN
UNIVERSITY

Limitations

This study reveals gender bias in salary and hiring scores, but the following limitations remain.

- Limitations at the data level: the publicly available dataset (Adult_Inc, Job_Salary, COMPAS) used in this study may not fully reflect the fairness of the actual hiring market.
- Limitations of the statistical analysis: it was not possible to fully control for all variables affecting salary.
- Model fairness modeling is not yet complete, and this study has not yet implemented fairness optimization on a machine learning model

From Bias to Balance: A Data-Driven Approach to Fair Recruitment Practices



COLLEGE OF
**SCIENCE &
ENGINEERING**
CENTRAL MICHIGAN
UNIVERSITY

Thank You For Listening!

From Bias to Balance: A Data-Driven Approach to Fair Recruitment Practices
Hairu Fan, Department of Statistics, Actuarial and Data Sciences, CMU
Email: fan2h@cmich.edu



COLLEGE OF
**SCIENCE &
ENGINEERING**
CENTRAL MICHIGAN UNIVERSITY