

Algorithmic Bias in AI Resume Screening

AUTHOR

Aidan Bugayong

AFFILIATION

Department of Mathematics, Applied Mathematics, and
Statistics, Case Western Reserve University

PUBLISHED

November 14, 2025

1 Introduction

1.1 Context and Background

Resume screeners were developed in order to screen candidates more efficiently and reduce the human bias in the screening process. However, there are concerns with whether or not these automated resume screening systems are truly unbiased in their decision making process. As more and more companies use some form of AI automation in their hiring process, the question of whether or not these systems are unbiased has become more important.

According to Naveen Kumar in his article on AI recruitment statistics, roughly 87% of companies are using AI in the hiring and recruitment process¹. As such, the University of Washington decided to conduct their own study to determine if there is any bias in the decision making process of resume screeners and what sorts of factors contribute to the scoring process of a resume. Their research found “significant racial, gender and intersectional bias in how three state-of-the-art large language models, or LLMs, ranked resumes. The researchers varied names associated with white and Black men and women across over 550 real-world resumes and found the LLMs favored white-associated names 85% of the time, female-associated names only 11% of the time, and never favored Black male-associated names over white male-associated names”². Their research came to these conclusions by handing the AI a list of identical resumes with different names and then having the AI give them scores. The article ends with the research team noting that more research should be done in this area by looking at different attributes and more LLMs in order to better align these AI systems with the real world policies to reduce bias and harm.

1.2 Project Purpose

This project aims to determine if there is any bias in the decision making process of resume screeners and what sorts of factors contribute to the scoring process of a resume. More specifically, looking at a wide variety of applicant attributes to determine what factors have the highest contribution to the bias in the decision making process. Ideally, professional experience and education will be the best predictors of the resume scores but we will also look into other factors such as gender, ethnicity, and institutional prestige.

2 Methods

2.1 Data Collection and Purpose

The original list of resumes is from a dataset in huggingface³, which is comprised of both real and synthetic resume data in JSON formatting. The purpose of this dataset was for training natural language processing (NLP) models for resume parsing. Specifically, the resumes are oriented around technical roles and is designed for NLP models to be trained on and used for candidate matching / screening in this field. This dataset was posted on huggingface February 21st, 2025 and has not been updated since (excluding the minor changes to the readme).

The sources used for this dataset are sourced from anonymized CV submissions as well as synthetic resumes generated using “Faker Library” and filled with realistic and role appropriate information. All resumes are anonymized by removing PII (Personally Identifiable Information) but many fields (such as names) contain realistic placeholders. The makers of the dataset note that the data is oriented around technical roles and the synthetic resumes may not capture the same nuances of a real resume. As such, the makers note that this dataset should only be used for NLP, data augmentation, or exploratory data analysis and should not be used for non-technical roles or personalized hiring decisions.

The dataset contains over 4500 resumes in a JSON format. Each resume entry contains personal information, work experience, education information, skills and projects. Since these are technical resumes (oriented around the computer science / information technology field), the skills and projects fields contained a candidate’s coding projects and/or coding languages. For the scope of this project, all fields were used for analysis or scoring of the resume.

2.2 Data Processing

The collection of resumes was processed in order to create the score datasets used for this project. The score datasets has the following columns: Name (acting as the primary key), Resume Score, Gender, Ethnicity, Institutional Prestige, Years of Experience, Skill Relevance, Experience Relevance and Project Relevance. The score datasets are specific to each of the language models used in this project. This is because we are trying to understand the relation between how the model scores a resume and the model’s own perception of the candidate (gender, ethnicity, and institutional prestige).

For creating the score datasets themselves, the resumes were scored by the model and then the model was asked to guess the gender, ethnicity, and institutional prestige of the candidate. For determining the skill/project/experience relevance, that was also a call to the language model. It is important to note that each column of data is a separate instance of the language model to reduce the likelihood of previous responses affecting the current responses. All of this information was then used to create the score dataset for each of the models used in this project.

The models used for this project include IBM’s *granite3.3:8b*, Microsoft’s *Phi4:14b* and Meta’s *llama3.1:8b*. These models were chosen because they are highly rated models despite the lower parameter size and are all open source.

2.3 Statistical Tools and Approach

After the resumes have been processed by the language models, there was some data cleaning to ensure the models followed the instructions. The initial approach was to use linear regression to predict the resume score from all of the other variables in the dataset. The formula used for this model was: $\text{score} \sim \text{gender} + \text{ethnicity} + \text{prestige} + \text{skill_score} + \text{project_score} + \text{experience_score} + \text{years_experience}$. Here, the response variable is the resume score which are integers in the range 0-100. The categorical variables are: gender, ethnicity, and prestige. The numeric variables are: years of experience, skill score, project score, and experience score. All categorical variables were setup such that “Unknown” was the reference category. The numeric variables, like the resume score, are in range 0-100, with the exception of years experience which is a positive decimal.

Next there was analysis of interaction effects and determining if any interactions were considered statistically significant. If the interaction was considered significant, then the interaction was added to the model. Additionally, the model diagnostics were conducted to ensure that the model’s assumptions were met. Specifically, the model was checked for normality in the residuals and multicollinearity by using the Shapiro-Wilk test and variance inflation factor, respectively. Interaction effects were added to the model if they were statistically significant. Due to multicollinearity issues, the model used a ridge regression.

3 Results and Discussion

Initial Modeling

Table 1: Initial Model Coefficients

Predictor	Meta β (SE)		Microsoft β (SE)		IBM β (SE)	
	p-value		p-value		p-value	
Intercept	93.81 (4.934)	$< 2 \times 10^{-16}$	32.87 (4.896)	3.07×10^{-11}	77.70 (6.487)	$< 2 \times 10^{-16}$
Gender						
Female	-6.32 (4.681)	0.177	-2.409 (2.718)	0.376	-2.25 (2.923)	0.441
Male	-6.08 (4.681)	0.194	-2.078 (2.725)	0.446	-6.902 (2.91)	0.018
Ethnicity						
African American	-0.945 (1.022)	0.355	-0.1997 (3.15)	0.950	8.545 (5.715)	0.135
Asian	-0.343 (1.252)	0.784	2.761 (3.536)	0.435	6.521 (6.675)	0.329
Caucasian	-0.867 (0.960)	0.367	-0.098 (3.058)	0.974	11.31 (5.456)	0.039
Hispanic	0.3980 (0.997)	0.690	-0.741 (3.103)	0.811	10.61 (5.618)	0.059
Prestige						
High	-0.415 (0.650)	0.524	1.214 (3.122)	0.697	-3.686 (2.145)	0.086
Medium	0.1902 (0.541)	0.725	-0.1486 (2.023)	0.941	0.687 (1.786)	0.701
Low	-0.021 (0.516)	0.968	3.422 (3.481)	0.326	-3.462 (1.698)	0.042
Scores						
Skills	0.0007 (0.009)	0.933	0.1147 (0.045)	0.011	-0.1036 (0.05)	0.024

Projects	0.0081 (0.009)	0.345	0.1480 (0.021)	6.75×10^{-12}	-0.0763 (0.04)	0.058
Experience	-0.001 (0.011)	0.932	0.3693 (0.028)	$< 2 \times 10^{-16}$	-0.0804 (0.04)	0.028
Years Experience	0.0168 (0.056)	0.764	-0.2541 (0.083)	0.002	-0.7438 (0.18)	4.73×10^{-5}

Let's interpret all significant predictors at a value of $\alpha = 0.05$ according to the values in Table 1. For the Meta model, there were not any significant predictors, instead only the intercept was significant. Specifically for the Meta model, when all predictors are at 0 and the categorical variables are at their reference group (Unknown), then the average expected score is 93.81. This suggests either the response variable (resume score) is either independent of the predictors or some transformations need to be applied to the response variable and/or the predictors.

For the Microsoft model, the only statistically significant predictors were the numerical ones. In the context of this problem, this is most likely a good thing as it suggests that the Microsoft language model is unbiased towards gender, ethnicity and institutional prestige. For the experience score predictor, the linear regression resulted in a coefficient of 0.3693, suggesting that the resume score increases by that value on average per unit increase of the experience score while holding all other predictors constant. Similarly, the skills score increased the resume score by 0.1147 on average per unit increase of skills score while holding the other predictors constant. Project score also increased the resume score by 0.1480 on average per unit increase of project score while holding the other predictors constant. Most surprising to me was the years experience had a negative coefficient, suggesting that the resume score decreased by 0.2541 per year increase of experience while all other predictors remained constant. This could suggest that the Microsoft model is discriminating against older individuals who are closer to retirement but more data analysis would need to be done.

The IBM model was quite sporadic with its significant predictors. All the categorical variables (gender, prestige, and ethnicity) were significant as well as the skills and experience scores and years experience. For ethnicity, knowing the candidate was white (compared to unknown ethnicity) increased the resume score by an average of 11.31 while holding all other predictors constant. This was also the highest increase for any categorical variable across any model. For gender, knowing the candidate was male (compared to unknown gender) decreased the resume score by an average of 6.902 while holding all other predictors constant. This was the greatest decrease for any categorical predictor across all models. For prestige, knowing the candidate was from a school with less prestige (compared to unknown prestige) decreased the resume score by an average of 3.462 while holding all other variables constant.

Model Improvements and Diagnostics

First I applied the Shapiro-Wilk test to all the models to determine if the normality assumption is satisfied. All models resulted in p-values less than $\alpha = 0.05$ which suggests that none of the residuals are normally distributed. More specifically, the Meta and IBM models had p-values less than 2.2×10^{-16} and the Microsoft model had a p-value of 5.042×10^{-13} . Thus, I will apply a box-cox transformation in an attempt to make the residuals normally distributed.

Next according to the VIF test, the IBM model had VIF's greater than 10 for both the Caucasian (16.24) and Hispanic (11.22), suggesting that multicollinearity is present in the IBM model. For the Meta model, the Gender variable had VIFs greater than 10 (both Male and Female had VIFs greater than 134) suggesting the Meta model may also be subject to multicollinearity. Lastly the microsoft model had VIFs of over 20 for the gender column and many ethnicity categories also had VIFs greater than 10. Thus all models were subject to multicollinearity and will use Ridge Regression to help mitigate the effects of multicollinearity.

To test the presence and significance of interaction effects, I fitted new models on each of the datasets with the same formula mentioned in section 2.3 with the addition of the interaction term. Specifically the interaction terms added were Ethnicity×Gender, Gender×Prestige, and Prestige×Gender. Each of these interaction terms were added one at a time and compared to the original formula using an anova test to determine if the interaction term was significant. All anova tests had p-values greater than $\alpha = 0.05$ for each interaction term in all of the models so none of the interaction effects were significant and will not be added to the final modeling process.

Improved Model

Ridge Regression Cross-Validation Results

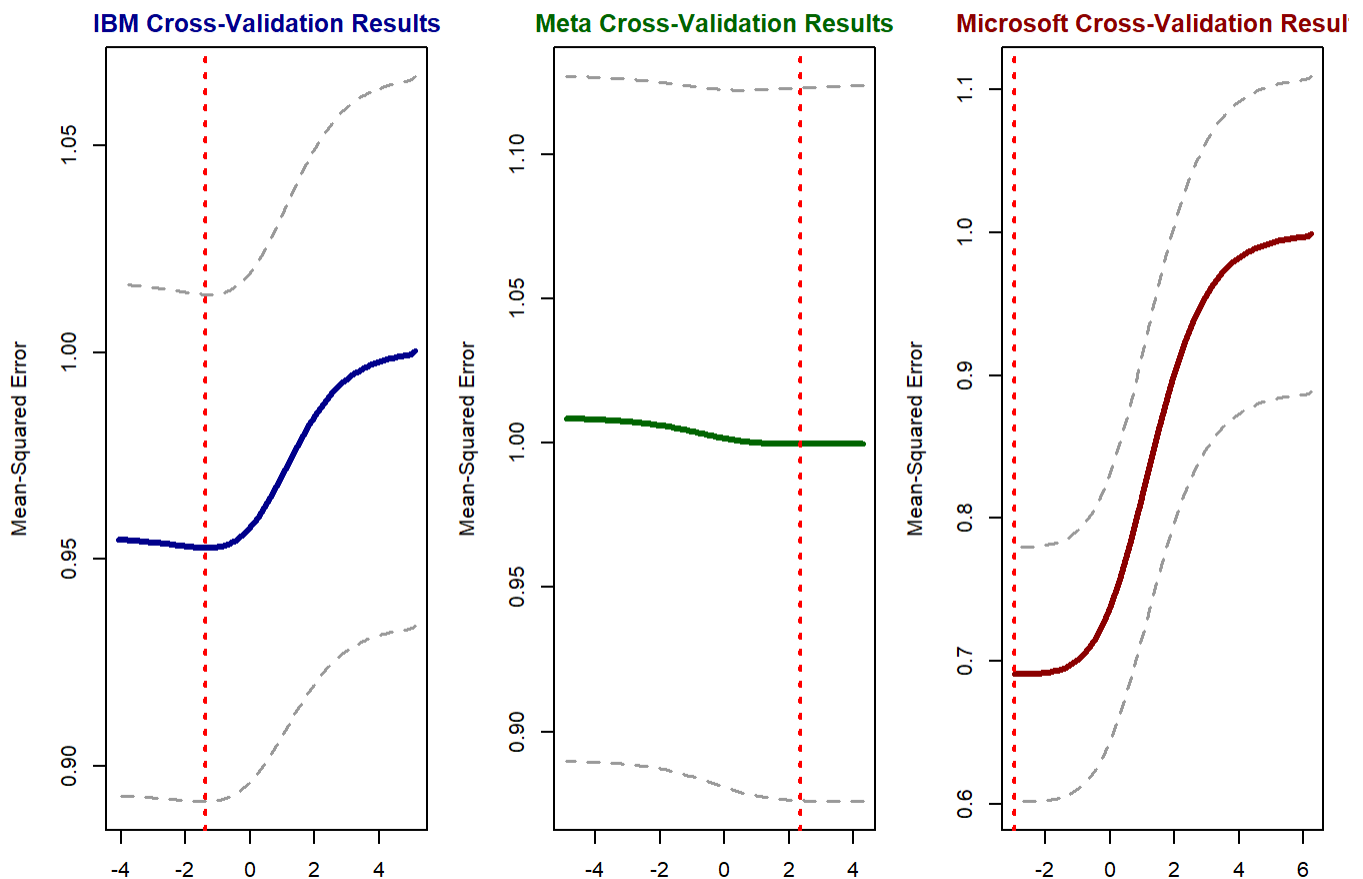


Figure 1: A line plot

When performing the box-cox transformation, all models resulted in $\lambda = 2$ which will then be used to transform the response for the ridge regression. Thus, all the ridge regressions will be performed on Score^2 as the response variable. Note this $\lambda = 2$ refers to the exponent for our response variable and not the weight of the penalty term in the ridge regression. It is also important to center and scale all numerical predictors during ridge regression

because ridge regression is sensitive to the scale of inputs.

Next was performing Cross-Validation for the ridge regression to find the optimal parameter for the penalty term. According to Figure 1, the IBM model had an optimal lambda of 0.2498. For Meta, the optimal lambda was 55.98. For Microsoft, the lambda was 0.0514.

After finding the optimal lambda for each model, we can report the 4 largest coefficients, which represent the variables most important for explaining the variation of the response variable. For the IBM model, the 4 largest coefficients were low institutional prestige (-0.1376), high institutional prestige (-0.1219), years experience (-0.1230) and asian ethnicity (-0.1129). For the meta model, all coefficients were extremely small suggesting a rather poor fit overall. Specifically, the largest were hispanic ethnicity (3.018×10^{-3}), caucasian ethnicity (-1.748×10^{-3}), african american ethnicity (-1.312×10^{-3}) and high institutional prestige (-1.183×10^{-3}). While these models are poor fits ($R^2_{ibm} = 0.0689$, $R^2_{meta} = 0.00042$) this still demonstrates subtle biases within these models which should aim to be removed. This demonstrates this because the most important factors for the ridge regression were the categorical variables (gender, prestige, ethnicity) that should not influence a candidates' fit for the job.

The microsoft model had by far the best results for minimizing bias and also the highest performing model of the bunch ($R^2_{microsoft} = 0.3221$). The top 4 largest predictors for this model were experience score (0.3817), low prestige (0.2705), asian ethnicity (\$ 0.2625\$), and project score (0.2203), with most other predictors being below 0.1. It is important to note when interpreting the coefficient for a ridge regression, we are looking at an average increase in the response (which is a scaled and centered Score²) per unit increase in the predictor (for numerical variables a unit would be one standard deviation since the variables are scaled, for categorical interpretation is the same as before) while holding the other variables constant.

Conclusion

Thus, this project found that there is likely bias in the decision making process of resume screeners. More specifically, for most of the models, gender, ethnicity, and institutional prestige seemed to play more of an important role in determining a candidate's fit for the job than other professional factors like work experience and projects. Additionally, the data was extremely variable with each model having completely different distributions of scores, ethnicities and prestiges, with certain models not appearing to have a very consistent scoring process or method. It is important to note however that the model that utilized the professional experience the most in determining the resume's overall score was the Microsoft model (*phi4_14b*) which also had the largest parameter size of the group. More studies should be done to analyze if model size plays a factor in presence of bias as well as potentially look at a temporal aspect by seeing if newer models tend to have less biases than older models.

Footnotes

1. <https://www.demandsage.com/ai-recruitment-statistics/> ↩
2. <https://www.washington.edu/news/2024/10/31/ai-bias-resume-screening-race-gender/> ↩
3. <https://huggingface.co/datasets/datasetmaster/resumes> ↩