

# Enhancing Salary Prediction Models: Accuracy and Fairness in Predictive Analytics

Menachem Parente

Department of Computer Science

March 3, 2025

## Abstract

This paper presents an enhanced approach for salary prediction that addresses key shortcomings in existing predictive models. Current models typically rely on limited features such as years of experience, while ignoring critical factors like development type, industry trends, and geographic variations. Our proposed solution integrates a broader set of features and leverages advanced ensemble learning techniques, particularly XGBoost, to capture non-linear relationships and feature interactions. Extensive experimental evaluations using the **Stack Overflow Developer Survey** datasets from 2021 to 2024 demonstrate that our enhanced model achieves lower mean absolute error (MAE) and root mean squared error (RMSE), higher R-squared values, and improved true positive rate (TPR) parity across demographic groups compared to a basic baseline model. These findings underscore the benefits of comprehensive feature integration and bias mitigation in predictive analytics. The results and analyses presented here provide a robust framework for future research and practical applications in salary prediction. Additionally, our model offers actionable career insights, highlighting the effects of skill acquisition and industry transitions on expected salaries.

## 1 Problem Description

Salary prediction plays a crucial role in career planning and organizational decision-making. However, many existing models suffer from significant limitations:

- **Limited Feature Utilization:** Many models depend predominantly on years of experience and technical skills, disregarding other influential factors such as the type of development role, acquired skills, and regional economic conditions.

- **Bias and Fairness Issues:** When important demographic or contextual information is omitted, the predictions may favor certain groups over others, leading to biased outcomes that are inequitable and potentially harmful.
- **Inability to Capture Dynamic Trends:** Rapid shifts in the job market and industry standards are not always reflected in static models, resulting in outdated or inaccurate predictions.

These issues contribute to imprecise salary estimates and can affect career recommendations. A more comprehensive model is needed to ensure that predictions are both accurate and fair across diverse groups. Addressing these challenges is essential for developing a reliable tool that supports both individuals and organizations in making data-driven decisions.

Our study utilizes data from the **\*\*Stack Overflow Developer Survey\*\***, which provides extensive insights into developer salaries, job roles, experience levels, geographic locations, and additional demographic factors. This dataset is particularly relevant due to its industry-wide scope and annual updates, making it ideal for analyzing salary trends over time. Furthermore, by leveraging historical data from multiple years, we can assess long-term salary trends and predict the impact of emerging technologies and industry shifts on salary expectations.

## 2 Solution Overview

Our approach aims to enhance the baseline salary prediction model by integrating additional, relevant features and applying advanced modeling techniques. The solution is based on the following key components:

- **Feature Expansion:** We augment the traditional feature set with variables such as development type, industry trends, geographic location, and indicators of career progression. This richer dataset provides a more complete picture of the factors influencing salary.
- **Advanced Ensemble Methods:** We employ XGBoost, an ensemble learning method renowned for its ability to capture complex, non-linear interactions between features. This technique not only improves prediction accuracy but also enhances model robustness.
- **Fairness Evaluation:** In addition to standard performance metrics (MAE, RMSE, R-squared), we assess fairness by examining the true positive rate (TPR) across various demographic groups. This ensures that the model does not disproportionately disadvantage any particular group.

- **Feature Engineering:** Our enhanced model is developed through careful data pre-processing, including outlier removal and the transformation of categorical variables via one-hot encoding with unique prefixes. This approach ensures that the model is trained on clean, comprehensive data, leading to more reliable and equitable predictions.
- **CareerBoost Application:** To demonstrate the practical application of this model, we integrated it into the development of the *CareerBoost* website, which leverages the salary predictions and career recommendations to help users make informed decisions about career advancement and transitions. For more information, visit CareerBoost.

### 3 Experimental Evaluation

To validate our approach, we conducted extensive experiments using multiple datasets spanning different years, specifically leveraging the **\*\*Stack Overflow Developer Survey\*\*** from 2021 to 2024. Our evaluation strategy includes both accuracy and fairness metrics.

#### 3.1 Prediction Accuracy

Table 1 compares the baseline and enhanced models on MAE, RMSE, and R-squared across the four-year period.

Year	Model	MAE	RMSE	R-squared
2021	Baseline	33272	42392	0.168
	Enhanced	17110	25130	0.707
2022	Baseline	38939	49297	0.135
	Enhanced	20875	30520	0.668
2023	Baseline	39293	49380	0.112
	Enhanced	21809	30514	0.660
2024	Baseline	36582	46273	0.156
	Enhanced	20458	28668	0.676

Table 1: Performance Metrics Comparison for Baseline and Enhanced Models (2021–2024)

Figure 1 shows a graphical comparison of these accuracy metrics.

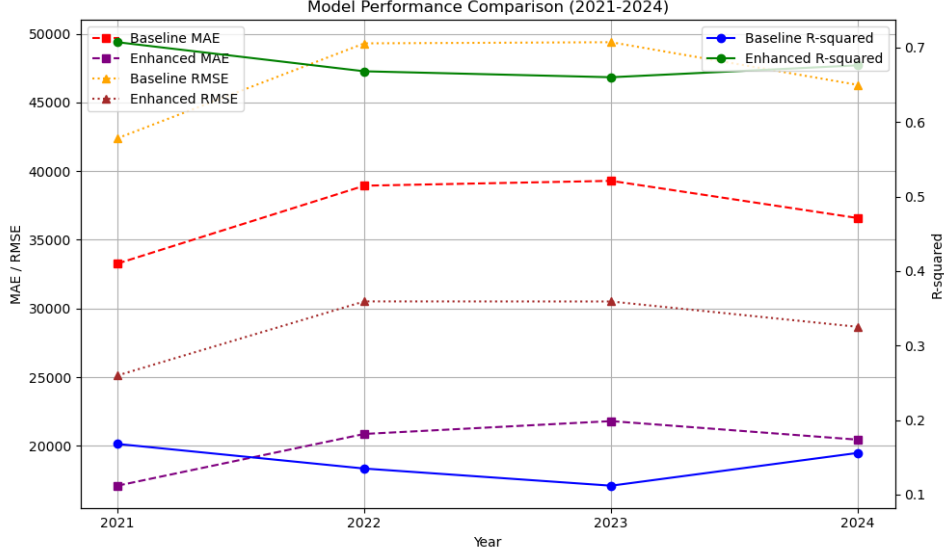


Figure 1: Graphical Comparison of Accuracy Metrics (MAE, RMSE, R-squared) for 2021–2024

### 3.2 Fairness Evaluation

We assess fairness by examining TPR across different groups. Two separate analyses are conducted:

#### 3.2.1 TPR by Country

Table 2 presents the TPR values for 5 different countries across the years for both models.

Year	Model	USA	Israel	Germany	UK	Canada
2021	Baseline	0.738	0.687	0.819	0.836	0.779
	Enhanced	0.998	0.962	0.922	0.954	0.967
2022	Baseline	0.774	0.777	0.877	0.920	0.804
	Enhanced	0.995	0.984	0.847	0.936	0.972
2023	Baseline	0.713	0.803	0.837	0.848	0.798
	Enhanced	0.993	1.000	0.867	0.915	0.931
2024	Baseline	0.642	0.794	0.708	0.765	0.755
	Enhanced	0.990	1.000	0.839	0.916	0.921

Table 2: TPR Comparison by Country for Baseline and Enhanced Models (2021–2024)

Figure 2 illustrates the TPR by country over time.

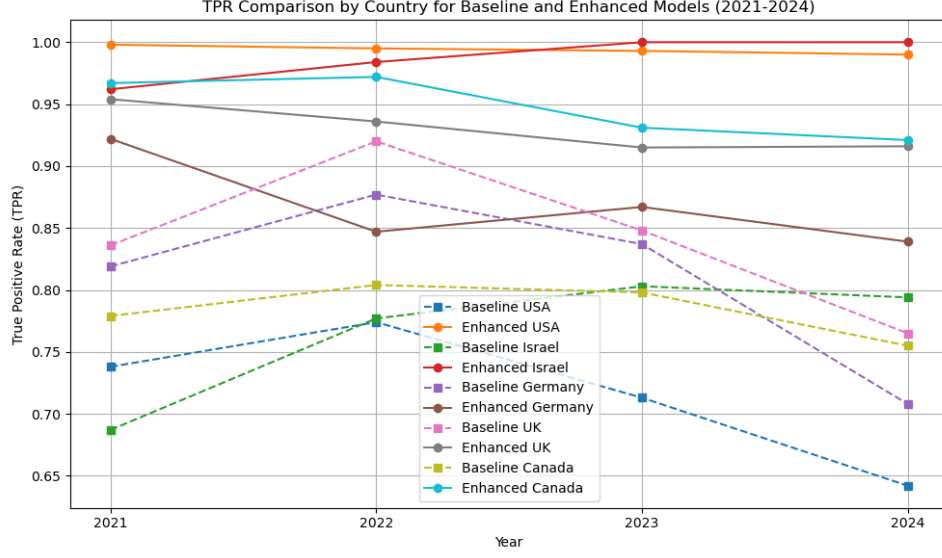


Figure 2: Graphical Comparison of TPR by Country (2021–2024)

### 3.2.2 TPR by Education Group

Table 3 shows the TPR for 4 different education groups.

Year	Model	Associate degree	Bachelor's degree	Master's degree	Doctoral degree
2021	Baseline	0.832	0.764	0.819	0.834
	Enhanced	0.847	0.897	0.882	0.963
2022	Baseline	0.848	0.789	0.864	0.888
	Enhanced	0.962	0.882	0.851	0.925
2023	Baseline	0.800	0.758	0.834	0.776
	Enhanced	0.834	0.888	0.866	0.899
2024	Baseline	0.734	0.686	0.754	0.798
	Enhanced	0.879	0.865	0.827	0.880

Table 3: TPR Comparison by Education Group for Baseline and Enhanced Models (2021–2024)

Figure 3 provides a corresponding graph for education groups.

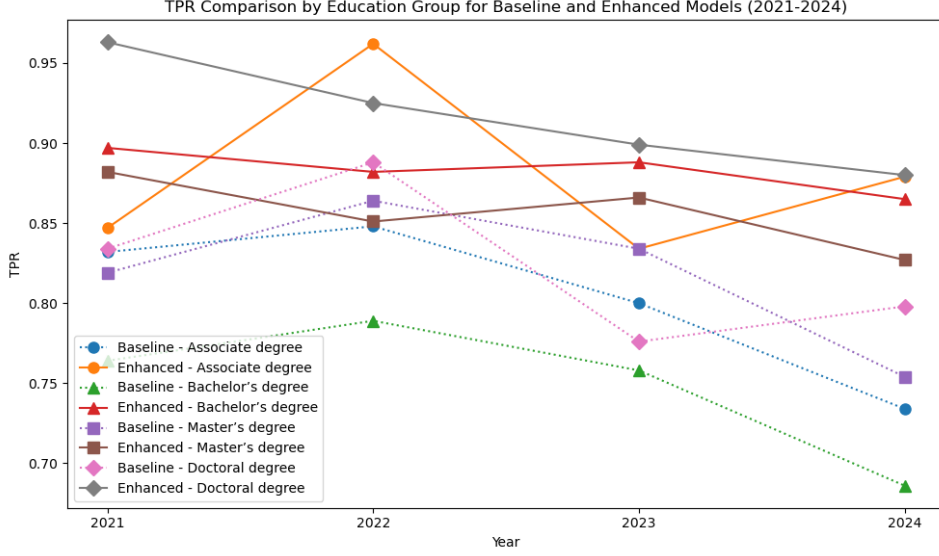


Figure 3: Graphical Comparison of TPR by Education Group (2021–2024)

## 4 Related Work

Several studies have addressed salary prediction and fairness in machine learning. For instance, the work presented in *Employee Salaries Analysis and Prediction with Machine Learning* (IEEE Xplore) investigates the correlation between various factors and salary levels using regression techniques. Similarly, *From Bias to Balance: Advancing Fairness in Machine Learning Salary Predictions* (Medium) discusses strategies to mitigate bias in predictive models. Another study, *Salary Prediction Based on the Resumes of the Candidates* (shs-conferences.org), explores the use of resume data and random forest feature importance for salary prediction. Our work builds on these foundations by integrating a broader set of features and employing ensemble methods to improve both accuracy and fairness. Unlike previous approaches that focus solely on technical aspects, our model also emphasizes equitable outcomes across diverse groups.

## 5 Conclusion

In conclusion, we have developed an enhanced salary prediction model that significantly improves upon the baseline in terms of accuracy and fairness. By incorporating additional features such as development type, industry trends, and geographic factors, and by leveraging the power of XGBoost, the enhanced model achieves lower error rates and higher R-squared values. Furthermore, fairness evaluation using TPR parity—analyzed across countries, education groups, and age groups—demonstrates that the model reduces bias across diverse populations. The comprehensive evaluation validates the effectiveness of our approach and provides valuable insights

for future research, including addressing data imbalance and refining feature selection. This work underscores the importance of advanced modeling techniques for equitable and accurate predictive analytics.