# What Determines Employee Salaries: A Regression Analysis Study ?

Neeraj Namani, Srisailam Gitte

University at Albany, State University of New York

## INTRODUCTION

This study utilizes multiple regression analysis conducted in R to explore the determinants of employee salaries, focusing on variables such as experience, education, job role and industry. Through meticulous data preprocessing and model selection, we aim to develop a robust predictive model. The objective is to offer actionable insights into a salary structuring and strategies for enhancing employee retention.

How do factors such as experience, education, job role and industry influence employee salaries and how can this information be utilized to inform salary structuring and retention strategies in organizations ?

## EDA

- EDA revealed strong positive correlations between salary and key variables such as age and years of experience. The matrix plots and correlation coefficients indicate that years of experience have the highest correlation with salary, suggesting that experience is a significant predictor of salary levels.
- Scatter plots display a distinct trend where salary increases with both age and experience, with density plots and histograms showing the distribution of these variables. Additionally, box plots reveal the spread and central tendency of age, years of experience and salary, suggesting variability in the data that could affect salary predictions.
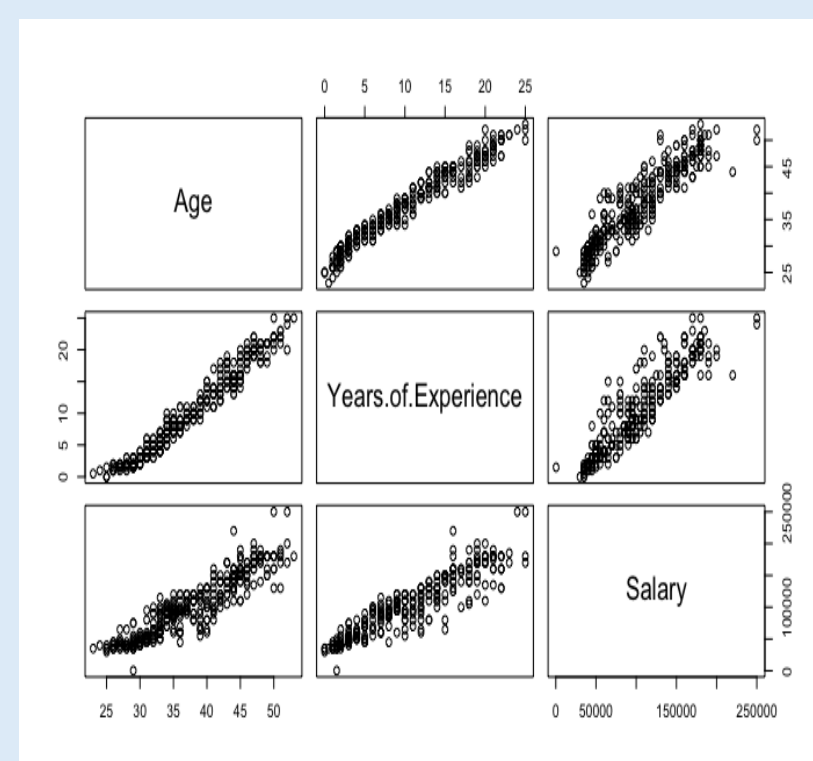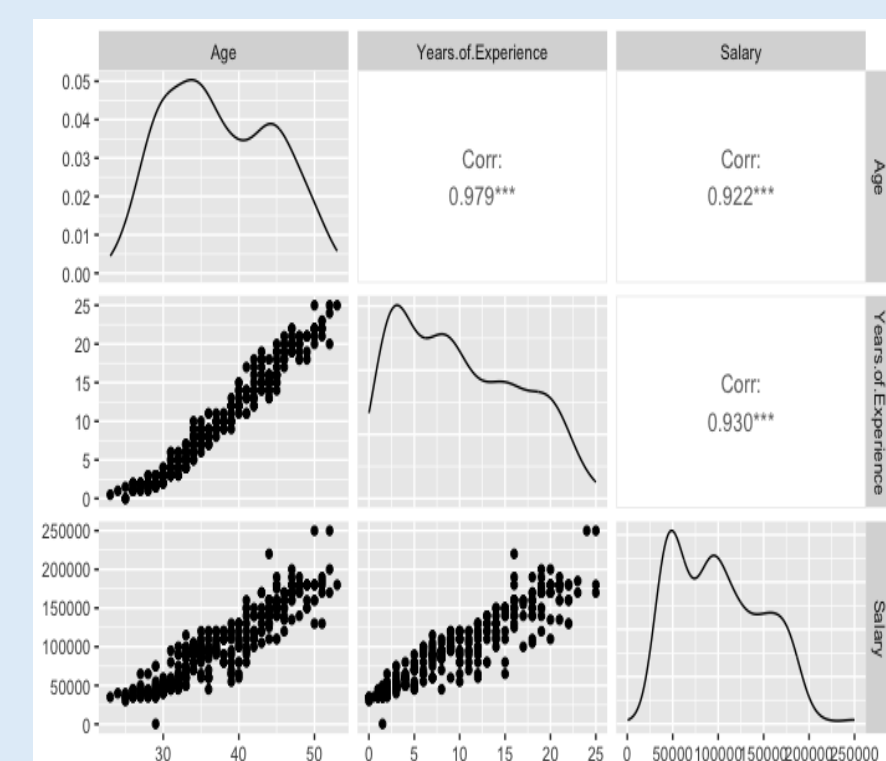

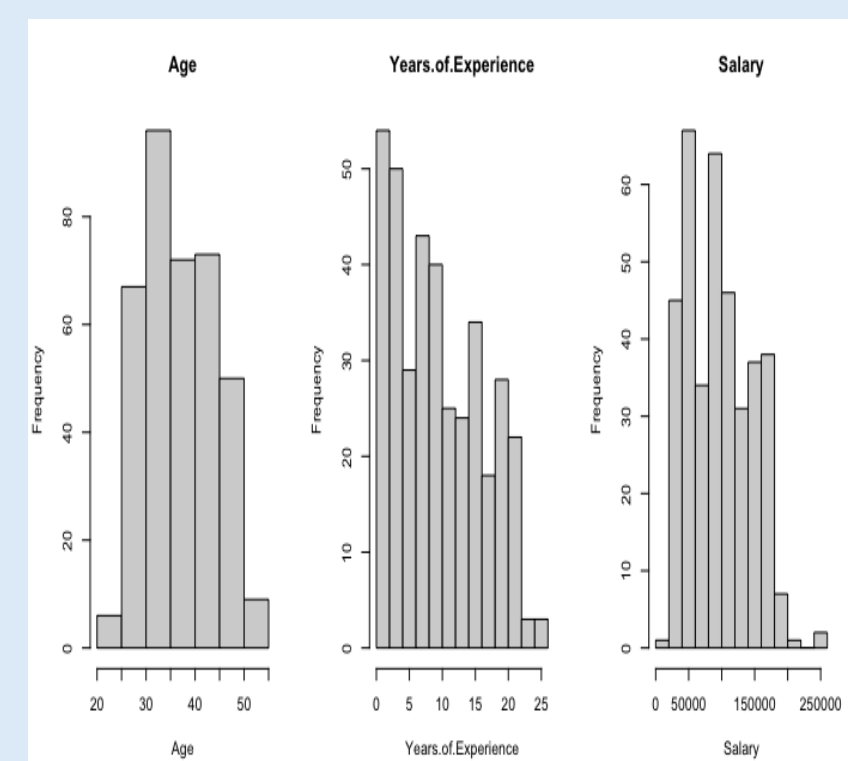Figure 1: Matrix Plot 1
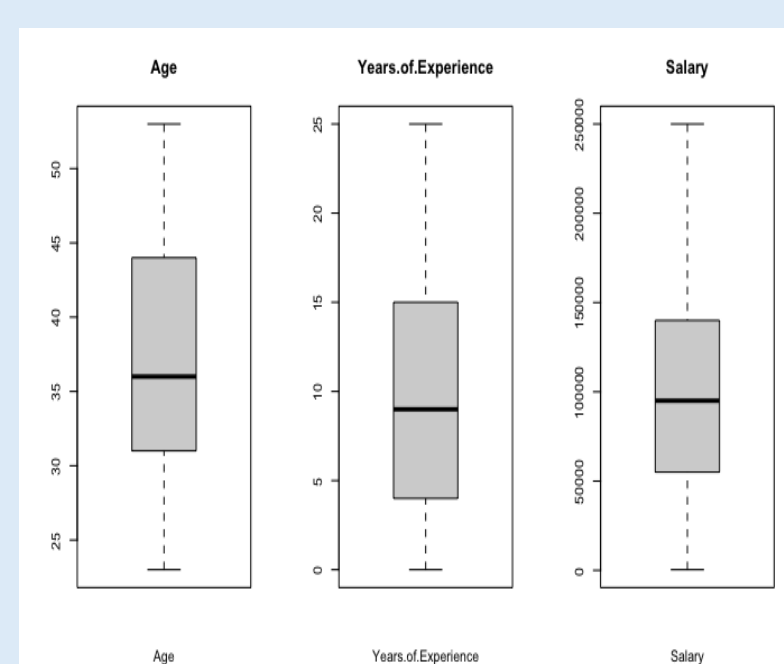

Figure 2: Matrix Plot 2


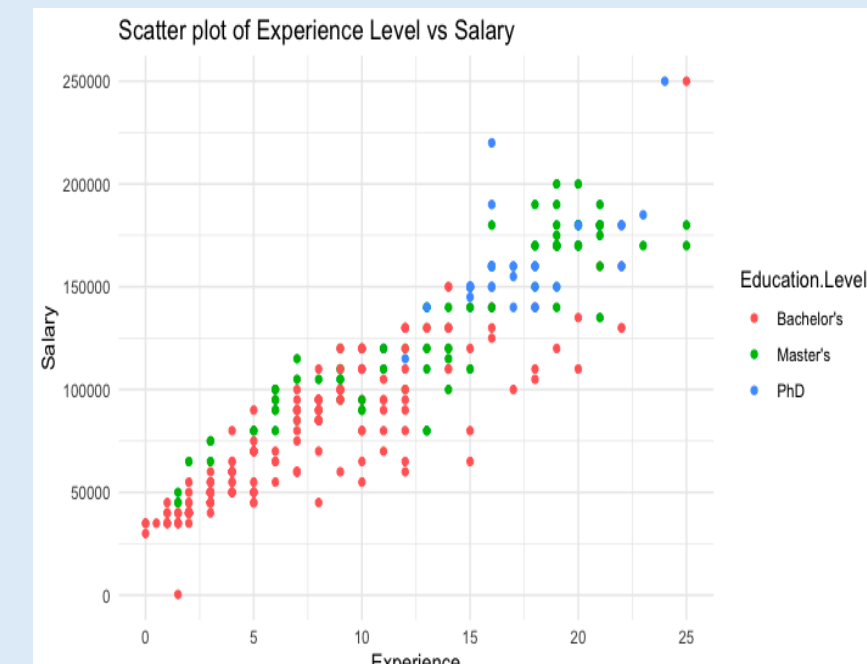Figure 3: Histogram Plots


Figure 4: Box Plots


Figure 5: Scatter Plot (Salary vs Experience)


Figure 6: Scatter Plot (Salary vs Age)

## MODELS


Figure 7: Linear Regression Summary

Age and Years of Experience positively impact salary coefficients of $2880.30 and $2873.50 respectively, indicating higher wages with increased age and experience.


Figure 8: Random Forest Summary

Years of Experience and advanced degrees (PhD, Master's) are the most influential in salary prediction. These variables show higher importance scores based on IncNodePurity, which is a measure of the decrease in node impurity from splitting on the variable, averaged over all trees.


Figure 9: Variable Importance Plot (Random Forest)


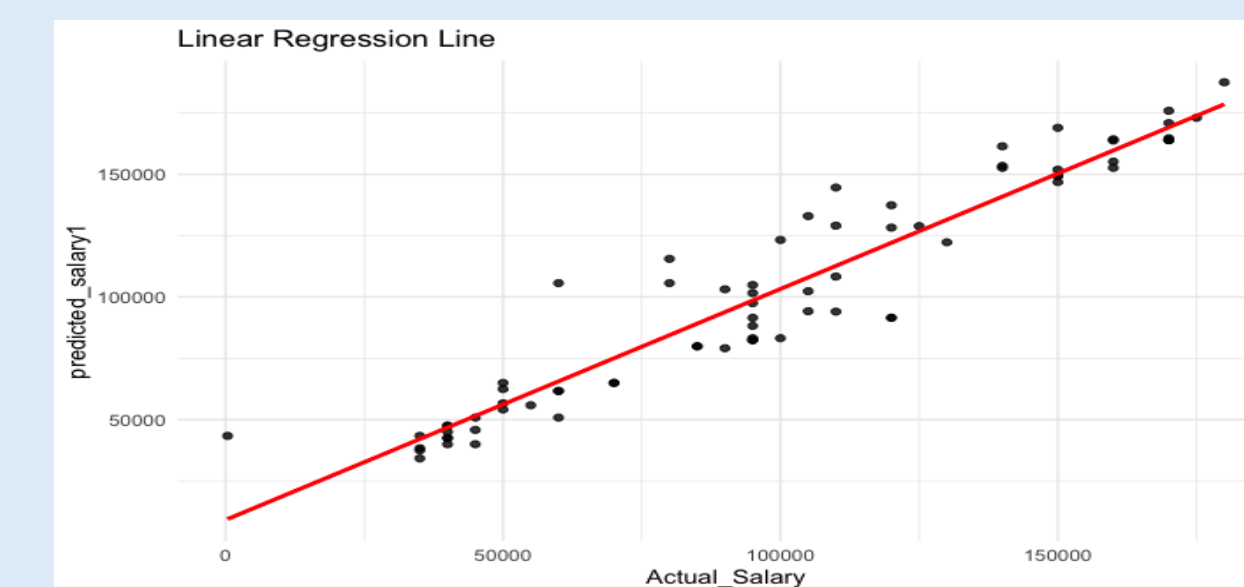Figure 10: Scatter plot (Salary vs Age Fitted)


Figure 11: Scatter plot (Salary vs Experience Fitted)


Figure 12: Linear Regression (Cross Validation)


Figure 13: Random Forest (Cross Validation)

| Model | RMSE <dbl> | Rsquared <dbl> | MAE <dbl> |
|---|---|---|---|
| Linear Regression | 15658.61 | 0.8973575 | 11291.03 |
| Random Forest | 15065.59 | 0.9071382 | 10307.93 |

Figure 14: Statistical Comparisions for two models

## MODEL VISUALIZATIONS


Figure 15: Linear Regression Plot


Figure 16: Random Forest Plot

| | Actual_Salary <int> | Predicted_Salary <dbl> | Error <dbl> |
|---|---|---|---|
| 7 | 120000 | 128287 | 8287 |
| 16 | 125000 | 128905 | 3905 |
| 26 | 45000 | 45806 | 806 |
| 35 | 170000 | 170879 | 879 |
| 36 | 45000 | 39993 | -5007 |
| 43 | 60000 | 50832 | -9168 |

Figure 17: Predicted vs Actual Salary

| | Actual_Salary <int> | Predicted_Salary <dbl> | Error <dbl> |
|---|---|---|---|
| 7 | 120000 | 118748 | -1252 |
| 16 | 125000 | 126267 | 1267 |
| 26 | 45000 | 46730 | 1730 |
| 35 | 170000 | 163362 | -6638 |
| 36 | 45000 | 42036 | -2964 |
| 43 | 60000 | 49672 | -10328 |

Figure 18: Predicted vs Actual Salary

- Random Forest performed better than Linear Regression

## CONCLUSION

After examining the results of our study alongside the referenced papers. It is evident that both our analysis and the literature concur on the significance of experience and education in predicting employee salaries. Where our study extends these findings is in the deployment of a Random Forest model that outshines Linear Regression in 10-fold cross-validation, offering a more intricate understanding of the determinants of salary.

Our study concludes that experience and advanced education significantly influence salary outcomes, with the Random Forest model providing a more nuanced understanding of these relationships than Linear Regression.

## REFERENCES

1. G. Wang, "Employee Salaries Analysis and Prediction with Machine Learning," 2022
Link: https://ieeexplore.ieee.org/document/9943146

2. "Employee Salary Prediction", 2022 by Tiasa Mukherjee, MS. B. Satyasaivani
Link: https://www.ijariit.com/manuscript/employees-salary-prediction/