Data Mining Final Project

Ander Eguiluz

Bellevue University

Predicting Salaries for Data Science Roles: A Data-driven Approach

**Introduction**

In today's data-driven world, the field of data science has seen exponential growth, leading to increased demand for skilled professionals. One of the critical aspects of a data scientist's career is negotiating a competitive salary based on various factors such as experience, job title, and company size. However, the salary range for data science roles can vary widely, making it challenging to determine an appropriate compensation package. This project aims to develop a predictive model to estimate salaries for data science roles, providing job seekers and employers with valuable insights for informed decision-making.

**Problem Justification**

Accurate salary predictions are essential for both job seekers and employers. For job seekers, having a clear understanding of the salary range for different roles enables them to make informed decisions about their career path and negotiate better offers. On the other hand, employers benefit from setting competitive salaries to attract top talent and retain valuable employees. By developing a reliable salary prediction model, we can improve transparency and fairness in the hiring process, ultimately benefiting both job seekers and employers.

**Stakeholder Pitch**

To pitch this problem to stakeholders, I would highlight the importance of leveraging data-driven approaches in salary negotiations and hiring practices. By developing a predictive model for salary estimation, we can provide job seekers with valuable insights into the market value of their skills and help employers make informed decisions about compensation packages. This project has the potential to revolutionize the way salaries are determined in the data science field, leading to better outcomes for both job seekers and employers.

**Data Source**

The dataset used for this project was obtained from a survey of data science roles, including information on work years, salary, job title, and company details. The dataset was carefully curated and preprocessed to ensure data quality and relevance to the problem. Missing values were handled using appropriate imputation techniques, and categorical variables were encoded for modeling purposes.

**Organized Summary of Milestones 1-3**

**Exploratory Data Analysis (EDA):**

Conducted a comprehensive analysis of the dataset to identify key features and relationships between variables. Visualizations such as histograms, box plots, and scatter plots were used to gain insights into the data distribution and identify potential patterns.

**Data Preparation:**

Prepared the data for modeling by handling missing values, encoding categorical variables, and splitting the data into training and test sets. This ensured that the data was in a suitable format for model training and evaluation.

**Model Building and Evaluation:**

Built several models, including linear regression and ridge regression, to predict salaries based on the available features. The models were evaluated using metrics such as Mean Absolute Error (MAE) and R-squared to assess their performance.

**Model Optimization and Performance Improvement:**

To address the issue of overfitting observed in the initial model, several steps were taken

to optimize the model and improve its performance. These steps included feature selection,

hyperparameter tuning, and model evaluation using cross-validation.

**Results and Findings**

The analysis revealed several key findings:

- Experience level, job title, and company size are significant factors influencing salary

  levels in the data science field.

- The models developed in this project showed reasonable performance in predicting

  salaries, with R-squared values indicating a good fit to the data.

**Numerical Analysis of Models:**

After adjusting the models to mitigate overfitting, the performance metrics improved as

follows:

| Metric | Initial Model | Adjusted Model |
|---|---|---|
| Training MAE | 25811.54 | 26761.58 |
| Training RMSE | 38504.50 | 39005.29 |
| Training R2 | 0.72078 | 0.71347 |
| Test MAE | 36394.72 | 30565.20 |
| Test RMSE | 59174.69 | 45101.22 |
| Test R2 | 0.08635 | 0.46925 |

**Recommendations**

Further optimization of the model is recommended, including exploring advanced techniques such as ensemble methods and feature engineering to improve prediction accuracy. Incorporating additional features or datasets, such as geographic location and industry sector, could enhance the model's performance and applicability. Regular updates to the model are essential to ensure it remains relevant and effective in predicting salary trends in the data science job market.

**Challenges and Opportunities**

Challenges include addressing potential biases in the dataset, ensuring the model's scalability, and adapting to new data sources and trends. Opportunities for further research include exploring additional factors that may influence salary levels, such as educational background, certifications, and specific skills.

**Conclusion**

In conclusion, this project demonstrates the potential of data-driven approaches in predicting salaries for data science roles. By leveraging data science techniques, we can gain valuable insights that benefit both job seekers and employers, ultimately contributing to a more efficient and transparent job market.