

# Salary Prediction for Data Science and Related Jobs

Alex Ojemann Department of Computer Science University of Colorado Boulder Boulder CO, USA Alexander.Ojemann@colorado.edu	Tyler Kimbell Department of Computer Science University of Colorado Boulder Boulder CO, USA Tyler.Kimbell@colorado.edu	Adrian Reghitto Department of Computer Science University of Colorado Boulder Boulder CO, USA adre5577@colorado.edu
---	--	---

## Introduction

When thinking about applying for a job or starting a career, it is important to understand how different positions, skill sets, and geographical location will affect someone's salary. With many different job postings it would be beneficial to give an estimate of an average salary given some rudimentary data. Stakeholders include students looking at majors they want to go into, people planning out careers, and people looking for a job. The significance of this problem is that these stakeholders could gain useful insight into what skills to invest time and effort into developing based on their correlation with a higher salary and what jobs may pay the best based on the skills they have.

## Literature Survey

Our data set contains 742 entries of data science and related jobs from glassdoor. One publication that explores the same data set explores a number of techniques including Ridge and Lasso regression, random forest, support vector regression, gradient boosting, and adaboosting, along with a dual adaboosting system of their design<sup>1</sup>. The authors use min-max scaling to preprocess their numeric variables and transform all categorical variables into binary dummy variables for each category using one hot encoding. They use adaboosting feature importance in order to limit the excessive dummy variables, filtering out any with a feature importance of less than 0.001. They ultimately evaluate the root mean squared error on their

---

<sup>1</sup> "Salary Prediction Based on the Dual-Adaboosting System | Semantic Scholar."

training set, the root mean squared error on their test set, and the mean absolute error on their test set for each of their models. The adaboost and dual adaboost models perform the best but exhibit a high degree of overfitting given the large split between their training root mean squared error and testing root mean squared error. Their accuracies can be used as a reference point for the accuracies of our model because they use the same data set.

One of the challenges of our data set is that the categorical features have many distinct values. For example, one of the features in our data set is the job title, which has 264 unique values. If we were to one hot encode this variable, we would get 263 binary dummy variables. This subjects us to the curse of dimensionality, which is the effect that large numbers of features leave you more vulnerable to finding trends in your data that aren't generalizable<sup>2</sup>. Two frequently used tactics for addressing this are Ridge<sup>3</sup> and Lasso<sup>4</sup> regression, which penalize feature weights in a regression model to limit the number of features for which the optimal model will have a nonzero coefficient. While these methods maintain the interpretability of the model, they may limit accuracy. Principal component analysis, or PCA, compares favorably to other methods for encoding variables in terms of maintaining accuracy in some applications with high dimensional categorical data<sup>5</sup>. However, this method sacrifices interpretability as PCA creates features using orthogonal linear components of all the variables that explain the most variation possible.

One of the other challenges that occurs when creating a model to predict salaries is hyperparameter optimization. Many machine learning models require hyperparameters that are computationally expensive to estimate with grid search and can lead to a suboptimal model if manual trial and error doesn't arrive at an optimal value. Randomized parameter optimization avoids the computational expense of grid search but is less prone to missing optimal combinations than manual exploration<sup>6</sup>.

## **Proposed Work**

As previously mentioned, the dataset we will use for our project is located on [Kaggle](#). To preprocess the data, we first analyze each of the features and determine if there are significant outliers that need to be removed. Next we will fill NA data with the mean for numeric features or the mode for categorical features. Then we will perform min max scaling on all of the numeric features. For the categorical features we plan to take a few different approaches to encoding them and comparing the results. One is to one hot encode them which will result in thousands of dummy variables and put them into ridge and lasso regression models to regularize the features. This approach is similar to that of the first paper mentioned in the literature survey but

---

<sup>2</sup> Debie and Shafi, "Implications of the Curse of Dimensionality for Supervised Learning Classifier Systems."

<sup>3</sup> Marquardt and Snee, "Ridge Regression in Practice."

<sup>4</sup> Tibshirani, "Regression Shrinkage and Selection via the Lasso."

<sup>5</sup> Farkhari et al., "New PCA-Based Category Encoder for Cybersecurity and Processing Data in IoT Devices."

<sup>6</sup> Bergstra and Bengio, "Random Search for Hyper-Parameter Optimization."

we will let the ridge and lasso regression models choose from all of the features rather than using adaboost beforehand. This approach could still be subject to the curse of dimensionality in that one of the many dummy variables could appear to be significant in the training data so the model includes it but it may not generalize. For example, the top feature in ADABOOST feature importance was “GRM Actuarial”, a dummy variable only found in three entries. Another approach will be to generate principal components for the dummy variables and use these as features in our models along with the numeric variables. The number of principal components will be decided later. The final approach will be to attempt to group these categorical variables into a manageable number of dummy variables representing characteristics of the variable. For example, the headquarters variable which contains the location of the company headquarters could be grouped into urban or rural and a dummy variable could be made for whether the description contains a defining word about an industry. We plan to put these PCA features and manually engineered features into a random forest, an xgBoost model, and an adaboost model using randomized search to optimize hyperparameters, then compare the performance of each of these models along with the ridge and lasso regression models with all of the one hot encoded variables.

## **Evaluation**

To evaluate each of our models, we plan to use root mean squared error and mean absolute error, each of which are defined below.

$$\begin{aligned}\text{Root Mean Squared Error} &= (\text{sum}(y - y\_pred)/n)^{1/2} \\ \text{Mean Absolute Error} &= (\text{sum}(\text{abs}(y - y\_pred))/n)\end{aligned}$$

We chose these evaluation metrics because they allow for direct comparison to the models in the first citation of the related work section. We plan to train each of these metrics on each model using 10-fold cross validation to avoid the arbitrary nature of using one single train test split. We are using k=10 so that the sample size isn't too small in the test sets for each fold as they will contain approximately 80 entries.

## **Milestones**

- Our first milestone will be to perform the outlier analysis, NA filling, and numeric scaling described in the proposed work section. This will be done by October 15th.
- Our second milestone will be to complete the encodings of categorical features described in the proposed work section. This will be done by November 9th.
- Our third milestone will be to train our models, optimize hyperparameters, and compare using our evaluation metrics. This will be done by December 12th, along with the final report.

## **Changes to Project Proposal**

The main change to our project proposal is that we now intend to limit the features we use based on further exploration of the data set. Some redundant columns were discovered, such as 'Founded' and 'age', where age is just the number of years since the founding year as of 2020 when this data set was scraped, so the 'Founded' column was removed. In addition, there are multiple columns for salary, our response variable, one of which is text and three of which are numeric representing the minimum, maximum, and mean. The text based variable corresponds to the numeric variables in all cases except for those in which the text based one specifies that the salary is hourly rather than yearly, in which the numeric variables appear to be scaled to yearly assuming 2000 hours worked per year. We decided to use the mean salary as our response so the text based salary, the minimum numeric salary, the maximum numeric salary, and the binary variable representing whether the salary was hourly or not were removed. Some additional features such as 'Industry' and 'Sector' may capture similar information. We didn't remove these, however, because one of our feature encoding methods, PCA, will separate the variation into orthogonal components and in our other feature encoding method, manual encoding, we can encode them as we see fit.

## **NA Filling**

NA values were filled using the mean for numeric features and the mode for categorical features as described in the proposed work section. One challenge faced here was that the seniority feature was populated with NAs to indicate that the job was not a senior position and 'Senior' if it was. Replacing the NA values with the mode would not make sense here, so the feature was converted into a binary encoding, `is_senior`, that was a 1 if the job was a senior position and a 0 otherwise.

## **Outlier Analysis**

Histograms and boxplots were generated for all numeric variables to analyze outliers. The 'age' and 'rating' features contained values of -1, which didn't make sense in context and caused a right skew when scaled in the next step. These variables were filled with the mean value similar to the na values so this was no longer the case. Some other variables had upward outliers, but none that didn't make sense in context so they weren't removed.

## **Numeric Scaling**

As described in the proposed work section, all continuous numeric features were scaled using a min-max scaler. Feature scaling is necessary for PCA, one of our proposed dimensionality reduction techniques. This included just the avg\_salary, age, desc\_len, num\_comp, and Rating features after some were removed as described in the 'Changes to Proposed Work' section.

## **Feature Encoding and Dimensionality Reduction**

We planned three ways of encoding the categorical features in the proposed work section. One of these is using PCA. Generating principal components was accomplished using the Scikitlearn library in Python. All of the dummy variables generated from categorical features were used to generate these principal components, of which there are several thousand, but none of these components explained more than 6% of the variation in the data on their own, so 20 principal components were used to retain a sufficient percentage of the variation while still having each principal component explain at least 1%. The other two methods for encoding and dimensionality reduction, one hot encoding with ridge and lasso regression and manual encoding, have not yet been completed.

## **Challenges**

One of the most significant challenges we've faced so far is that using PCA with all categorical features one hot encoded appears to overfit the data. When I initially tried using PCA, the first principal component explained 99.7% of the variation in the data. Since it can take into account all of the thousands of dummy variables, it's likely creating a principal component tailored to each data point which will not generalize to new data. However, when only using the encoded dummy variables from the categorical columns, the variation explained by the first principal component is just 5.42%, a much more reasonable figure.

## **Timeline**

We are on track with our project milestones as we completed milestone 1 which included na filling, outlier analysis, and feature scaling which we planned to complete on October 15th and we're progressing on our feature encoding and dimensionality reduction techniques which we plan to complete on November 9th.

## **Feedback Incorporation**

The introduction and evaluation sections were expanded based on feedback from our project proposal.

## References

- Bergstra, James, and Yoshua Bengio. "Random Search for Hyper-Parameter Optimization," n.d.
- Debie, Essam, and Kamran Shafi. "Implications of the Curse of Dimensionality for Supervised Learning Classifier Systems: Theoretical and Empirical Analyses." *Pattern Analysis and Applications* 22 (May 1, 2019). <https://doi.org/10.1007/s10044-017-0649-0>.
- Farkhari, Hamed, Joseanne Viana, Luis Miguel Campos, Pedro Sebastiao, and Luis Bernardo. "New PCA-Based Category Encoder for Cybersecurity and Processing Data in IoT Devices." arXiv, May 23, 2022. <http://arxiv.org/abs/2111.14839>.
- Marquardt, Donald W., and Ronald D. Snee. "Ridge Regression in Practice." *The American Statistician* 29, no. 1 (1975): 3–20. <https://doi.org/10.2307/2683673>.
- "Salary Prediction Based on the Dual-Adaboosting System | Semantic Scholar." Accessed September 26, 2023. <https://www.semanticscholar.org/paper/Salary-Prediction-Based-on-the-Dual-Adaboosting-Chen-Zhan/641c323bce27ea7f1791f176991a3c10df5d0f70>.
- Tibshirani, Robert. "Regression Shrinkage and Selection via the Lasso." *Journal of the Royal Statistical Society. Series B (Methodological)* 58, no. 1 (1996): 267–88.