

# Project 2: Final

Econ 1680: Machine Learning, Text Analysis, and Economics

Lena Kim

May 10, 2024

github link: <https://github.com/wkim31/Econ1680-TA-Project>

## 1 Introduction

The job market is overcrowded. Google receives millions of applications a year, and by that point it is simply impossible to read each and every resume to its fullest. As such, many companies now rely on algorithms to pick the resumes of best fit, motivating questions of algorithmic fairness in the hiring process- are the right people being hired, or are they simply being chosen based on certain words in their resume?

In fact, a New York Times coverage of the ever-increasing reliance on resume algorithms has even led to identification of certain keywords and skills to include in your resume (sectioned off by position type) to make it past this resume sifter (Weed (2021)). As a Post article states, including even one of these key words can mean the difference between getting hired and being discarded (Abril (2023)).

Taking inspiration from this trend, this project aims to approach it from the other side of the table and accurately predict the salary levels of a job based on keywords used in its job description. With some companies remaining tight-lipped in regards to salary transparency, this text analysis project may be useful for job seekers who want to know more about compensation before taking any steps toward subjecting themselves to a cold, algorithmic hiring process.

**Research Question:** Using keywords and skills found in the job description, can we accurately predict the salary of the job? Which words contribute the most weight to a higher/lower salary?

## 2 Data Sources and Descriptions

APIFY, a cloud platform for web scraping, is used to scrape recent jobs data from the

employment website Indeed.com. The dataset originally included 700 jobs (white collar, blue collar, service jobs, and skilled trades) with locations in New York; however, in the exploratory analysis it was discovered that it had also scraped some Spanish-language jobs. After filtering out jobs not in the English language and additional data preprocessing to filter out potential duplicate posts, the final set consisted of 639 unique jobs posted on Indeed in 2024 ready for text analysis.

The main input features are: 'company', 'position', and 'description'. The main target variable is 'yearlysalary', which is a manually created column based on the 'salary' column that could take the form of hourly salaries, daily salaries, and yearly salaries. Descriptive statistics- in this case, only applicable for our quantitative target variable hourly salary- is shown in section 7.

In addition to these numerical descriptive statistics above, the text data in the job descriptions were also explored in the exploratory data analysis step. Word clouds and more are shown below, in Section 7.

Certain jobs require certain skills, and as such words more common to that job type appear more often. For example, "service" and "customer" appear more often in client-facing roles, and "data" and "team" appear more in analyst roles. Across the board, "experience" appears in all of the wordclouds, meaning experience is desired no matter the role.

### 3 Method

The main methods used for the analyses are: Ridge Regression and Neural Networks.

Before the implementation of these methods, pre-trained sentiment analysis models were implemented to see the sentiments of different job descriptions, yielding unsurprising results: most job advertisements were labelled positive, as they are essentially advertisements encouraging people to apply. However, a significant finding by the pretrained model is: words associated with service jobs ("dishwashing", "ironing") are predictive of lower salaries. See Section 7 for details.

For the main regression task of predicting salaries from job descriptions, a Ridge Regression is used of the form defined in equation 1. In this example,

$$(\text{predicted salary})_i = \beta_0 + \beta_1 (\text{TFIDF vector})_i + \varepsilon_i + \lambda \sum_{j=1}^p \beta_j^2 \quad (1)$$

where

$$(\text{TF-IDF vector})_i = \frac{\text{count of word } w \text{ in job description}}{\text{total words in job description}} \cdot \left( \frac{\text{count of job descriptions } d \text{ containing } w}{\text{total job descriptions in corpus}} \right)^{-1} \quad (2)$$

$(predicted\ salary)_i$  is the target variable of interest, the yearly predicted salary for job description  $i$ .  $X_i$  is our main independent variable, the vector of coefficients (weights) corresponding to the TF-IDF vector constructed for job description  $i$ . The TF-IDF vector itself consists of measurements evaluating how relevant a job descriptor word is to that document/job description  $i$ , which is itself in the larger corpus.  $\varepsilon_i$  is the error term for each job description regression, and  $\lambda \sum_{j=1}^p \beta_j^2$  represents the regularization term added to prevent overfitting.

The *conditions under which* this analysis is valid rests on common assumptions of linearity, as ridge regressions are regularized forms of linear regressions. In the current example, we would require some form of conditional independence assumption to hold. We can express this as

$$E[\varepsilon_i | X_i] = 0 \quad (\text{A1})$$

This means that once we know the values of our independent variable (the TF-IDF vector constructed for job description  $i$ ), the expected value of the error term should be 0. That is, the error term is not systematically related to the values of the TF-IDF vector; there are no systematic biases left in the residuals between the observed and predicted salaries. Given that each document of job descriptions seems to be independent of each other and our large enough sample size, we will proceed with our Ridge regression.

$\beta_1$  is our main metric of interest- a vector of coefficients corresponding to the TF-IDF vector. Each element within  $\beta_1$  corresponds to the contribution of a particular word in the TF-IDF towards the predicted yearly salary. These estimates are presented in Section 4.

Finally, a neural network architecture is used to predict yearly salaries from job descriptions. Neural networks complement the regression task well, as they may be more effective at capturing the complex relationship between the input features of the TF-IDF vector and yearly-salary. The architecture consists of multiple layers of neurons interconnected by weighted edges. A feedforward neural architecture is employed with three layers: an input layer, a hidden layer, and an output layer of interest.

The input layer is fed the constructed TF-IDF vector for job description  $i$ . The hidden layer is composed of 64 neurons with a nonlinear ReLU activation function, defined as:  $f(x) = \max(0, x)$ . This layer is the primary processing unit where the network learns the complex patterns between the TF-IDF and the yearly salary. The output layer consists of one neuron representing the predicted yearly salary that is outputted, which is necessarily activated by a linear activation function for our regression task.

All of this information is encapsulated in the equation:

$$(predicted\ salary)_i = \text{ReLU}(\beta_0 + \beta_1(\text{TFIDF vector})_i) + \varepsilon_i$$

where the ReLU is applied element-wise to the linear combination of the input features and their corresponding weights.

The metric of interest is again  $\beta_1$ , and is interpreted in much the same way as the ridge regression: the impact of each word in the TF-IDF vector on the predicted salary. A one unit change in X represents a higher relevance of that word in the job description, and an increase in this is associated with a  $\beta_1$  increase/decrease (based on if  $\beta_1$  is positive/negative) in the associated salary. Job descriptions containing more relevant terms or words are more likely to be associated with higher predicted salaries.

## 4 Results

The Ridge Regression and Neural Network here yielded unsurprising results: words associated with skills in tech and finance predict higher salaries.

First, after constructing TF-IDF vectors for each job description, a Ridge Regression is trained on an evaluation metric of Mean Squared Error, which is minimized. Looping over the alpha parameter found an optimal alpha of around 0.3, and this graph is shown below in Section 7.

The coefficient table (Section 7) specified the highest magnitude positive and negative words. Positive words are associated with higher salaries according to our ridge specification; of note is the pattern of tech and finance-related words, such as "ai", "risk", and "business". This suggests that most positions that contain higher salaries in our dataset are likely to be finance, business, or technology-related.

By contrast, the negative words are much more service-related. For example, "shift", "hour", "week", and "customer" seem to be related more to service or part-time jobs, with their emphasis on customer-oriented, scheduled work.

The neural network, with a ReLU activation function, has been trained with 10 epochs and an Adam optimizer. The top positively weighted words are outputted in Section 7, and include words such as "duties", "teams", and "programs". The negative words include words such as "stakeholders", "solutions" and "competitive". The neural network results seem to highlight the importance of a team and project-oriented worker, while under-emphasizing the value of individual-oriented, competitive workers.

## 5 Conclusion

The Ridge Regression results hold the most power in answering my original research question of predicting salaries from job descriptions. It has been hypertuned with the optimal alpha, and produced concrete results in identifying the highest magnitude positive and negative results in determining salaries. The limitation of the research rests on the fact that the dataset is data focused on jobs in New York, and a similar project may not be generalizable to the rest of the nation. To extend this project, there will be a need to consider addressing this limitation by reproducing the project on data outside the field of New York, and out into all corners of America.

## 6 References

### References

Abril, Danielle (2023) “Your résumé isn’t the only thing popular job sites evaluate,” *Washington Post*.

Weed, Julie (2021) “Résumé-Writing Tips to Help You Get Past the A.I. Gatekeepers,” *New York Times*.

## 7 Graphs and Tables

### Summary Statistics and EDA

Figure 1: Sample job description

'Every day, Customer Service Representatives (CSRs) at Maximus are entrusted to serve some of the most vulnerable communities by providing customer care to millions of New Yorkers who need to maneuver through complex healthcare plans. During these uncertain times we ensure that we are delivering the best outcomes possible for our clients and customers – ensuring every action is thoughtful, open, transparent, and done with integrity. To prepare you for this role, Maximus provides paid, comprehensive training which ensures our customer service representatives care for each caller with the highest levels of knowledge and professionalism.\nPlease note this job posting is for upcoming classes in Albany, NY. This position requires 10 days of onsite training after successful completion of training, there is the opportunity to work remotely. Training is a total of 6 weeks.\n\nEssential Duties and Responsibilities:\n\nProcess new applications for health care coverage via telephone i

Note: Sample job description.

Figure 2: Top 10 Most Common Words and Summary Stats

```
[('work', 2245),  
 ('experience', 1896),  
 ('''', 1620),  
 ('team', 1269),  
 ('skills', 1186),  
 ('ability', 1092),  
 ('job', 1047),  
 ('required', 1019),  
 ('new', 969),  
 ('including', 959),  
 ('benefits', 955),  
 ('position', 929),  
 ('time', 839),  
 ('must', 754),  
 ('company', 750)]
```

Table 1: Summary Statistics

yearly_salary	
Count	632
Mean	69343.35
Stdv	45391.71
Min	325.00
Median	57510.00
Max	440000

Note: Descriptive statistics for yearly\_salary.

Methods:

Figure 3: Wordclouds by Job Type



Note: Top words for client-facing job descriptions.

Top words for analyst job descriptions

Figure 4: Wordcloud for Overall Job Descriptions



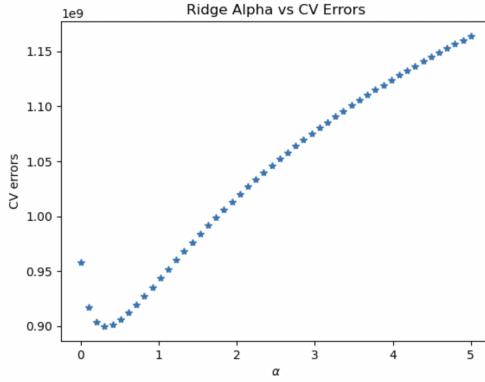
Note: Top words overall, for all positions.

Figure 5: Top 5 Negative Words (Sentiment Analysis)

	word	coef	abs coef
<b>8148</b>	osi	-30.963059	30.963059
<b>11411</b>	toiletry	-23.447377	23.447377
<b>3861</b>	dishwashing	-21.384278	21.384278
<b>6903</b>	linens	-19.405892	19.405892
<b>6435</b>	ironing	-7.661715	7.661715

Note: Negative Words are those associated with service work, which are associated with lower salaries.

Figure 6: Ridge Alpha Plot and Coefficients



Note: Optimal Alpha of 0.3.

Figure 7: Ridge Coefficients

				<b>var</b>	<b>var_ridge</b>
<b>6867</b>	management	19025.781729	<b>9952</b>	shift	-16359.918868
<b>954</b>	ai	17078.131791	<b>5638</b>	hour	-13543.910505
<b>5947</b>	info	15680.692006	<b>7367</b>	must	-13073.772226
<b>1145</b>	anthropic	15285.377074	<b>3248</b>	customer	-12897.677279
<b>6527</b>	lead	14754.615668	<b>5346</b>	guests	-10931.255592
<b>9498</b>	risk	14516.375279	<b>3949</b>	duties	-10929.621474
<b>3556</b>	develop	13526.747383	<b>9302</b>	required	-10511.365903
<b>2697</b>	compensation	13008.824596	<b>9876</b>	service	-10131.617852
<b>2011</b>	business	12824.403413	<b>11816</b>	week	-10062.748386
<b>9314</b>	research	12158.195567	<b>643</b>	able	-9919.959307

Top magnitude positive words. Top magnitude negative words

Figure 8: Weights for Neural Net

Top positive-weighted features:		Top negative-weighted features:	
work	<b>0.193608</b>	programs	<b>-0.021118</b>
duties	<b>0.191875</b>	office	<b>-0.021008</b>
equal	<b>0.190836</b>	stakeholders	<b>-0.020624</b>
requirements	<b>0.189184</b>	high	<b>-0.020516</b>
schedule	<b>0.188626</b>	solutions	<b>-0.020515</b>
of	<b>0.187863</b>	location	<b>-0.020352</b>
in	<b>0.185911</b>	range	<b>-0.020281</b>
teams	<b>0.184459</b>	state	<b>-0.020023</b>
with	<b>0.183957</b>	competitive	<b>-0.020023</b>
program	<b>0.183939</b>	quality	<b>-0.019939</b>

Top positive weight words. Top negative weight words