

See discussions, stats, and author profiles for this publication at: <https://www.researchgate.net/publication/344646493>

Real-Time Resume Classification System Using LinkedIn Profile Descriptions

Conference Paper · October 2020

DOI: 10.1109/CISPSSE49931.2020.9212209

CITATIONS

2

READS

1,584

2 authors:



Sivakumar V

Manipal Academy of Higher Education

18 PUBLICATIONS 237 CITATIONS

[SEE PROFILE](#)



Ramraj Santhanam

Sri Ramachandra Institute of Higher Education and Research

15 PUBLICATIONS 109 CITATIONS

[SEE PROFILE](#)

Some of the authors of this publication are also working on these related projects:



Routing in underwater sensor networks [View project](#)



De Novo Algorithm Analysis when pipelined with Read Stack algorithms [View project](#)

Real-Time Resume Classification System Using LinkedIn Profile Descriptions

1st Mr. Ramraj S

*Assistant Professor, Department of
Software Engineering
SRM Institute of Science and Technology
Chennai, India
ramrajs@srmist.edu.in*

2nd Dr.V. Sivakumar

*Assistant Professor, Department of
Software Engineering
SRM Institute of Science and Technology
Chennai, India
sivakumv1@srmist.edu.in*

3rd Kaushik Ramnath G

*Student, Department of Computer
Science
SRM Institute of Science and Technology
Chennai, India
kaushikram1999@gmail.com*

Abstract - In the domain of online job recruitment, accurate job and resume classification is vital for both the seeker and the recruiter. We have built an automatic text classification system that utilizes various techniques like Term frequency-inverse document frequency with Machine Learning and Convolution Neural network for training the model with texts and classifying them into labels and finally to compare their results. Using resume data of applicants, we have categorized them into different categories. Due to the sensitive nature of resume data, we have used domain adaptation. A classifier is trained on a large dataset of job description snippet, which is then used to classify resume data. Despite having a small dataset, consistent classification performance is seen. The primary filter for this type of work is the efficiency the system can provide. We aim to compare the results obtained by various algorithms that are generated using the same data so that the efficiency of each algorithm can be evaluated. From the result, it is evident that character-level CNN gives a better F1 score compared to other models.

Keywords—Domain Adaptation, CNN, SVM

I. INTRODUCTION

In today's fast-moving corporate industry, recruiters often need to go through vast amounts of resume to analyze the applicants reliably to decide upon the deserving candidates. But it is not possible to keep up with the pace today. So, automated classification of resumes is needed to ease out the process. To do the same, a bulk of labeled resume data is required, and job openings are divided into a certain number of predefined job categories.

The labeled datasets have job titles and job descriptions, along with several other parameters. Resumes of the same job category are semantically similar. This similarity is looked for to match the applicants to their respective job categories. The data set is derived by extracting real-time resume data of applicants from LinkedIn using a web scrapping method. Individual LinkedIn URLs were extracted using a URL extracting platform. The URLs were then processed upon to finally scrape off the resume data in the form of a CSV file with all the features well distinguished. Convolution Neural network (CNN) has been used in several image and text classification applications over the years. We propose a model where a CNN based algorithm is used to train a classifier with a dataset of over 1000 applicants. After the training process, it is

applied to classify unlabelled resumes. As it is a comparison study, various algorithms are used with multiple combinations of features on the same data set, and the results were noted. Algorithms like Support Vector Machine (SVM), Naïve Bayes (NB), Term frequency-inverse document frequency (TFIDF) were initially applied individually and then later in an ensemble with another algorithm. Finally, all the different results are compared to find the most effective method of them all.

The paper is organized in the following manner: In sections 2 & 3, related works and our contributions are explained. In section 4, the system architecture is explained along with the process to generate the dataset. In sections 5 & 6, we have described the efficiency of each model and have analyzed the results of all the models.

II. LITERATURE SURVEY

In [1], Ohma B. Hashemi, Amir Asiae, and Reiner Kraft enhanced the search results by understanding the user query's intent. CNN methodology used for feature extraction. Word2vec used that is pre-trained with filters and feature maps. Over 10000 queries were derived from the logs of a search engine. The rule-based model was taken for comparison. For 14 high-level intent, CNN gives 0.47 as an average F1 score. For 125 class of low intent, CNN gave 0.50 as an F1 score.

In [2], Kevin and Neha classified comments into toxic and nontoxic as well as specify the type of toxicity. Representation of word and character are fed to CNN. Sparse category cross-entropy was used with SoftMax function, whereas for multi-label, a sigmoid function was used with binary cross-entropy. Dataset was gathered from Kaggle of over 159571 labeled samples. Subsampling technique used for binary classification and multi labeled, the type of toxicity was included. A multi-layer perceptron was taken for comparison with LSTM and CNN. LSTM gives better results in word-level analysis with 0.886 F1 scores, and at a character level, CNN proved to be better with 0.8 F1 score.

As demonstrated in [3], Multi-label classification of text documents is used in Czech. The input vector has 1s and 0s. A bag of words is used. Dictionary is built using 20000 words. Forty different kernels used for different sizes. Data was taken from a corpus of 2974040 words from 11955 documents. They concluded that the Proposed topology with output thresholding is efficient in multi-label classification. SoftMax works better

with FDNN, whereas sigmoid works better with CNN. Neural nets improve maximum baseline entropy with pre-processing and without parametrization.

LSTM [4] automation tool used in feature extraction instead of statistical or context properties for the detection of DGA domains. The dataset contains Top 1 million domains, Alexa - NOT DGA, OSINT DGA (about 30 classes) with 750,000 examples of DGA. They achieved a detection rate of 90% with a 1:10000 false-positive (FP) rate—an improvement of twenty times FP improvement over the next best method. The F1 score of 0.9906 was micro-averaged.

In [5], Andrej Karpathy et al. analyzed the semantic content of various applications like search and the summarization of video classification. CNN model was used to get results of image recognition, segmentation as well as detection and retrieval. The model trained by processing approximately five clips a second in case of full-frame networks and 20 clips a second for the multiresolution system. The performance was improved from 43.9% to 63.3%. They concluded that CNN architecture is capable of learning features. Low-resolution architecture speeds up CNN without letting go of the accuracy.

Alex Krizhevsky et al. [6] classified 1.2 million images by training a deep CNN. GPU implementation used to avoid overfitting. The regularization method is known as dropout used. The dataset contains 15 million images from 22000 categories collected from the web and trained the network on raw values of RGB. THE seven CNN model gave good results. The top one is 36.7%, and for the top five, 15.4%.

Using data from HiQ Lab, Eric Boucher, and Clement Renault [7] classified job titles based on Linked Summaries. They grouped the job title in small chunks of categories. Recurrent Neural Network was used and achieved an accuracy of 31.7 percent on 133 classes with a baseline at 16 percent.

In [8], Richard Girshick et al. presented a simple solution for object detection. They achieved a 30% improvement using the Convolution Neural Network algorithm and PASCAL VOC dataset. Ilya Sutskever [9] elaborately describes the nature of the RNN learning problem and also addresses the difficulty of the RNN Learning problem using second-order optimization. According to [10], when the neural network is trained on the small training dataset, it usually performs poorly. This can be rectified by reducing half of the features on the training set. Random “dropout” gives significant improvements in many benchmark tasks and sets new records for speech and object recognition.

In [11] validated the performance, efficiency of the XGBoost 4j package in a distributed programming environment called Apache Spark for predicting the speed of the wind. XGBoost is used in the training dataset to predict a target variable. In [12], a new method was proposed by combining CNN with Naive Bayes for classifying the generated domains of DGA. Dataset provided by DMD 2018 Shared Task was used for this classification problem. In [13], V. Sivakumar and D.Rekha used the evolutionary Genetic Algorithm (GA) to maximize the usage of the underwater acoustic bandwidth. The results showed that by reducing the time slots and maximizing the throughput using parallel

transmission, UWASN could transmit in average minimal turnaround time. A deep learning method [14] called CNN is used to classify the health-relevant web pages because of its learning power. Also, Character level embedding was used to extract appropriate features using CNN.

III. CONTRIBUTIONS

The main contributions of this work includes the following:

1. The dataset is generated by scraping more than 1000 linkedin URLs with attributes like name,job description,skill set etc.
2. Convolutional neural network is applied for the first time to learn the character level features from the given attributes of each LinkedIn profiles.

IV. METHODOLOGY

A. System Architecture

The given system architecture depicts the basic flow of operations in the system of resume classification. The first is of an applicant uploading his/her resume. This document is raw and unstructured that needs to be worked upon. Various modules are depicted as a data module, model training, and finally, testing. In the data module, the process involved is to access LinkedIn profile URLs from the net and process them with Selenium that will locate the necessary fields and write the data. In model training, Feature extraction is performed by passing data to embedded layers and convolutions.

B. Data Source

Since the pre-existing dataset was not according to the standards which were required, we scrapped the data on our own from LinkedIn. It contained multiple fields such as job title, job description, skills of a person, location, past experiences, images, and many more fields. We used an online API to scrap all the data from LinkedIn using skill-specific keywords to attain a high level of accuracy in the result section.

Further, we try to be more accurate with our result, so we will propose a first in class subcategory result generation through our system since there has been extensive development in every category, which is complicated on their own. For this out of these 27 categories, we have confined our search for a single category IT (Java developer, Cloud computing, Web developer, Testing field, Python developer, etc.). We have further broken down this field into five more subcategories to generate a sense of preciseness. In this way, a more efficient result will be generated, which would allow the recruiters to generate precise data.

LinkedIn does not allow scraping through code without permission, but there are specific tools through which it is possible to scrap real-time user data but in a limited way. We used phantom buster API to generate data from scratch. There is a restriction of 100 entries per user per day, so we used multiple accounts and combined the data generated. The data generated was stored in the form of an excel sheet, which

allows various cleaning operations on the raw data generated. This helped us to refine the data and make it something on which we can work on later.

V. EXPERIMENT

We trained our model using 95% of our resume data and used the remaining 5% to test the trained model. On testing with different algorithms and combinations, results were obtained and were compared with one another. On feature extraction initially, two features, namely, job title and job description, were used, and later, it was done only using the job description to avoid the hints that are by default conceded by the title.

The novelty of our research is that we have used the same data for different algorithms as opposed to various subsets of the dataset. This allows us to achieve better accuracy with the algorithms, and that way, it's easier to study the behavior of the different algorithms under a certain scenario. The different methods used were CNN and TF-IDF used in combination with other algorithms like TF-IDF+NB(Naïve Bayes), TF-IDF+SVM(Support Vector Machine) and TF-IDF+XGB. The results entail precision, recall, F1 score values. Each row in the result table represents each job category(e.g., Python, java developer, etc.).

In Fig. 1, the Feature Extraction module of the CNN algorithm is shown. This module extracts all the essential features using the CNN algorithm. This module is then fed to the training module, where the model trains for the given data and finally sent for the prediction to the testing module. In Fig. 2, the model accuracy is little over 60% when Resume Description and Job Title variables are taken into account, whereas in Fig. 3, the model accuracy is about 50% when only the Resume Description variable is considered. This tells us that the accuracy shown in Fig. 3 is better because of the inclusion of the Job Title variable with Resume Description.

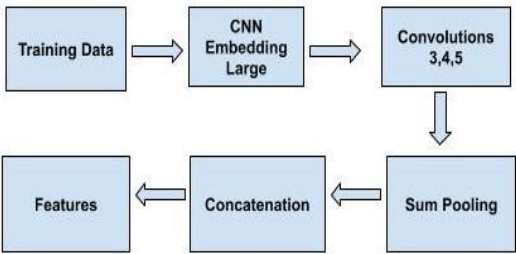


Fig. 1. . Feature Extraction using CNN

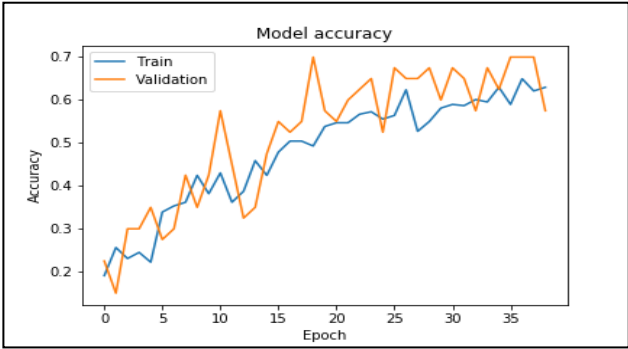


Fig. 2. . Resume Description + Job Title Accuracy

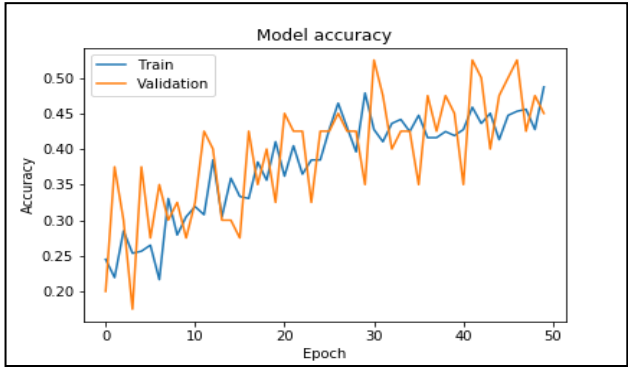


Fig. 3. Resume Description Accuracy

VI. RESULTS AND ANALYSIS

The combination of TF-IDF with other algorithms generated far better results. The most prominent result is that of CNN. It is an algorithm that is generally opted for image classification. But, our results are a testament to its ability to perform well in applications of text classification as well. Here, CNN has generated the best results of them all.

In TABLE I., the results of the Convolution Neural Network are shown. Convolution Neural Network gives the maximum F1-score because of its power to extract features and learn from them quickly and effectively. The results of all algorithms used are displayed in TABLE II.

TABLE I. PERFORMANCE OF CONVOLUTION NEURAL NETWORK

S.no	Precision	Recall	F1-score
1.	0.63	0.60	0.62
2.	0.83	0.55	0.48
3.	1.00	0.60	0.70
4.	0.64	0.50	0.55
5.	0.24	0.78	0.37
Average/ Total	0.68	0.67	0.65

TABLE II. PERFORMANCE OF ALL MODELS

Algorithm	Precision	Recall	F1-score
CNN	0.68	0.67	0.65
TF-IDF + NB	0.58	0.61	0.57
TF-IDF + SVM	0.63	0.63	0.61
TF-IDF + XGB	0.61	0.60	0.59

VII. CONCLUSION AND FUTURE WORK

We have compared the various state of the art algorithms using similar self-generated data which allowed us to compare them in an unbiased way. The data generation was the most challenging task as the data we generated is very sensitive data and is not accessible for everyone easily. Further, the data was raw and it had more deformities, so we had to clean it up various times. After all the efforts, the results generated were promising. We attained an F1 score of 0.57 for naïve Bayes algorithm, 0.59 for extreme gradient boost (XGB), 0.61 for support vector machine (SVM), and 0.65 in case of convolution neural network. As a result, CNN acquires the maximum result when compared with similar state of the art algorithm; we can very proudly conclude that CNN can also be used to classify text rather than just images and show exceptional results when compared to other text classification algorithms. In the future, Generative Pre-Trained Transformer (GPT -2) algorithm can be used to increase the size of the data set, which helps the training model to learn better with more information, and as a result, the accuracy increases. Also, other Deep Learning

algorithms like LSTM (Long Short-Term Memory) can be used to increase the performance on this data

REFERENCES

- [1] Homa B. Hashemi, Amir Asiaee, Reiner Kraft, Query Intent Deduction using CNN, ACM, 2016
- [2] Kevin, Neha, Detecting and classifying toxic comments, Stanford University, 2016
- [3] Ladislav Lenc, Pavel Kr, Deep Neural Networks for Czech Multi-label Document Classification, CoRR, 2017
- [4] Jonathan Woodbridge, Hyrum S. Anderson, Anjum Ahuja, and Daniel Grant, Predicting Domain Generation Algorithms with Long Short-Term Memory Networks, CoRR, 2016
- [5] Andrej Karpathy, George Toderici, Sanketh Shetty, Thomas Leung, Rahul Sukthankar, and Li Fei-Fei, Large-scale Video Classification using CNN, Stanford University, 2014
- [6] Alex Krizhevsky, Ilya Sutskever, and Geoffrey E. Hinton, Large-scale Video Classification using CNN, Stanford University, 2014
- [7] Eric Boucher, Clement Renault, Job Classification Based on LinkedIn Summaries, Stanford University, 2015
- [8] R. Girshick, J. Donahue, T. Darrell, and J. Malik, Rich feature hierarchies for accurate object detection and semantic segmentation, CVPR, 2014
- [9] Ilya Sutskever, Training Recurrent Neural Networks, University of Toronto, 2013
- [10] G. E. Hinton, N. Srivastava, A. Krizhevsky, I. Sutskever, and R. R. Salakhutdinov, Improving neural networks by preventing, University of Toronto, 2012
- [11] Ramraj, S; Saini, Aeshita; Kaur, Gurleen, Comparative Study of XGBoost4j and Gradient Boosting for Linear Regression, International Journal of Control Theory and Applications
- [12] Rajalakshmi, R; Ramraj, S; Kannan, R Ramesh, Transfer learning approach for identification of malicious domain names, International Symposium on Security in Computing and Communication
- [13] V. Sivakumar, D. Rekha, Node scheduling problem in underwater acoustic sensor network using genetic algorithm, Personal and Ubiquitous Computing, 2018
- [14] R. Rajalakshmi and S. Ramraj, A Deep Learning Approach for URL based Health Information Search, International Journal of Innovative Technology and Exploring Engineering, 2019