# Resuminer – A Résumé Recommender System for Job Seekers using Cluster Analysis

Orlando Ferrer Hernández
Master of Engineering in Computer Engineering
Dr. Jeffrey Duffany
Electrical & Computer Engineering and Computer Science Department
Polytechnic University of Puerto Rico

***Abstract*** — *Résumés or curriculum vitae are an essential requirement for job seekers in any industry and are usually the first interactions with third parties that are actively hiring. For a single person, there is an infinite combination of résumés possible. They can vary not only in style, but in the content and the overall presentation. Resuminer is a system for comparing and tracking résumés and their performance against jobseekers. By tracking how a third party interact online with a résumé, and by gathering a large quantity of résumés, the system can analyze new résumés as input and make recommendations as to what elements can improve the chances of getting a callback. The system will assign "interaction points" with different values to the various elements in a résumé.*

***Key Terms*** — *Bag of Words, Clustering, Corpus, Natural Language Processing (NLP).*

## INTRODUCTION

In order to measure what are the qualities of a desirable résumé, we need to be able to collect data. In order to collect the desired data on what is usually a physical or electronical document we cannot track, we need to be able to move all documents to the Web. In every résumé created, physical or electronic, there will be a link to an electronic version of the résumé using a short URL and instructions inviting the reader to visit said URL for additional information and examples of work. Every different document will have a unique URL. Documents can also incorporate a QR code for quick access. By logging all visits each URL, we can identify how many times any version of a specific document is been looked at, and create a scoring system that will track all interactions. All third party interactions with a résumé will be referred to as "interactions", and the scoring system will be referred to as "interaction points". Third party interactions with the résumé will be translated into "interaction points".

Actions that will increase the "interaction points" of a particular résumé are:

- Visiting a résumé via short URL
- Visiting a résumé via QR code scan
- Clicking on links in the résumé website
- Receiving a callback and/or interview

Callback data will be entered manually into the system.

## BACKGROUND

Machine learning can be described as a "Field of study that gives computers the ability to learn without being explicitly programmed" [1]. It is currently at the forefront of many recent and upcoming new technologies. Data can be considered the new currency of technology [2], and Big Data holds the answer to many interesting questions, some of them yet unknown. Machine learning is not a new topic, but with growing processing power and storage capabilities of modern computing, new applications are being found. It can be used from a simple classification, to complex computer vision used in self driving cars. At its core, machine learning consists of a set of inputs that generate a set of outputs using a model [3]. Once we have a model, we can make predictions. The model can evolve as more data is gathered.

Natural language processing is a way for a computer to analyze, understand and manipulate text to do useful things. At its core it extracts semantic information from human language, and can be used for simple tasks such as word frequency analysis, but can be used for complex tasks such as sentiment analysis, language translation, automatic summarization, and many others. In this work, some simple natural language processing will be used.

## DATABASE SCHEMA

The back end of the online web application will need a database to save and display data. Below is the database schema.

**Table 1**
**Page_visits**

| Field | Description |
|---|---|
| Resume_id | Unique id |
| Incoming_URL | Referrer URL |
| timestamp | Date and time |

**Table 2**
**Interactions**

| Field | Description |
|---|---|
| Resume_id | Primary id |
| timestamp | Date and time |
| Action_performed | Clicks |

**Table 3**
**Interaction_points_value**

| Field | Description |
|---|---|
| Visit_resume | Points for visiting |
| Scan_QR | Points for scanning QR |
| Click on link | Points for clicking |
| Click on Github | Points for clicking |
| Points | Points |

**Table 4**
**Job_requirements**

| Field | Description |
|---|---|
| Resume_id | Unique id |
| job_title | Name of the position |
| Company_name | Company name |
| Position_description | Description |

**Table 5**
**Resume**

| Field | Description |
|---|---|
| Resume_id | Unique id |
| Short_url | Short URL for resume |
| Resume_text | Resume text |
| Company_name | Company name |
| Job_title | Job title |

## ALGORITHMS & MODELS

The first level of text analysis will be a simple "bag of words" with frequency analysis. The "bag of words" is a model used in natural language processing where a document is represented as a bag (multiset) of unordered words, disregarding the grammatical structure but keeping duplicate words. After converting a document into this "bag of words" representation, it is possible to calculate various characteristics of the text, the most common being term frequency- the number of times each word appears in a document. After this is done, there are some techniques to normalize the text. Before normalization, we can remove some of the most common words such as prepositions and articles (the, on, at, in, which, on, and) and any other words that will not be relevant for the analysis- these are commonly referred to as stop words. Each individual file is called a "document", and the collection of all documents is called the "corpus".
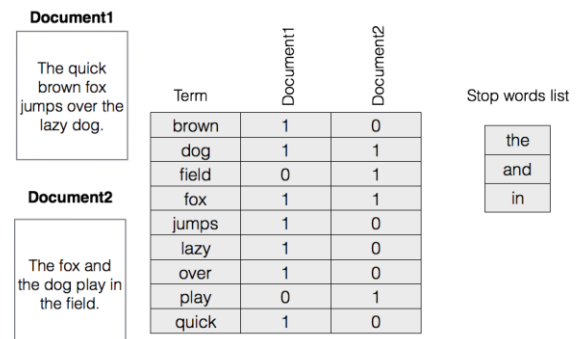


**Figure 1**
**Bag of Words Representation**

This information is saved internally in two separate matrixes, one for the terms and one for the term frequency.

**Figure 2**
**Runtime Terms and Term Frequency**

After removing stop words, the next step towards normalization calculating the term frequency (tf). As the name suggests, the frequency of every word is calculated in every individual document. This information is relevant, but not enough when we have multiple document we want to compare. In order to analyze multiple documents, the inverse document frequency (idf) will be needed as well. The inverse document frequency will measure how important a term is. When computing the term frequency, all the term are considered equally important. The inverse document frequency searches for the terms through all the corpus, and measures the relevance of the word though all the documents. By calculating both of these, we can now calculate the term frequency-inverse document frequency (tf-idf), which is the product of both. This is a numerical statistic that reflects the relevance of a word to a document in a corpus. This value increases proportionally to the number of times it appears in a document, but is also offset by the frequency of the word in the corpus. Mathematically speaking, the term frequency tf (t, d) can be represented as "(1)"

$$tf(t,d) = 0.5 + 0.5 \frac{f_{t,d}}{\max\{f_{t',d}:t'\in d\}} \qquad (1)$$

The inverse document frequency can be represented as "(2)".

$$idf(t,D) = \log \frac{N}{|\{d\in D:t\in d\}|} \qquad (2)$$

The second level of analysis will involve clustering. After we have multiple résumés that have gone through the text analysis, and have been viewed online multiple times, we can do a k-means cluster analysis. The top key terms are selected for each cluster.

K-means clustering aims to partition n observations into k clusters in which each observation belongs to the cluster with the nearest mean [4]. The way this is computed is by generating k initial means randomly. k clusters are created by associating observations to the nearest mean. The centroid of each of the k clusters now becomes the new mean. These steps are repeated until convergence. It is now possible to generate a cosine similarity matrix using the tf-idf matrix, then generate a distance matrix (1 − similarity matrix). After this we generate a plot. In the plot we should see the clusters of the different documents.
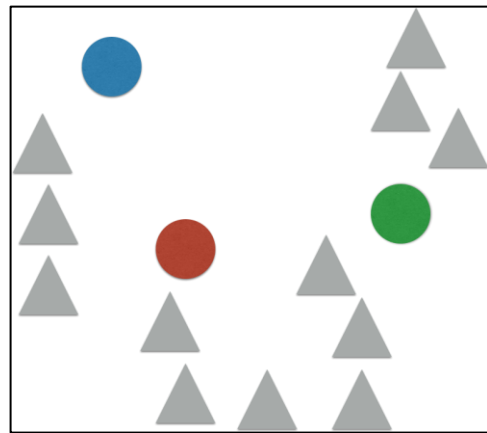


**Figure 3**
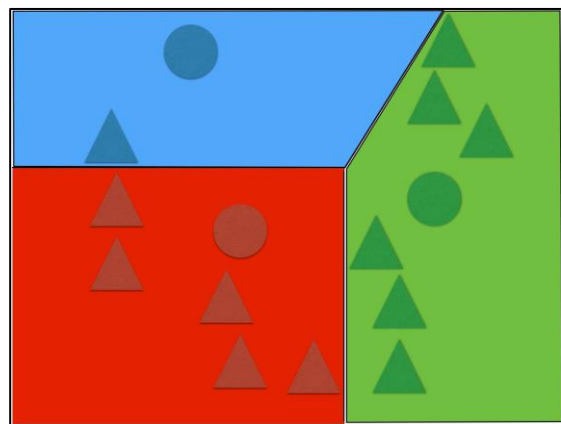**k Means Step 1- Random Generated k (shown in color)**



**Figure 4**
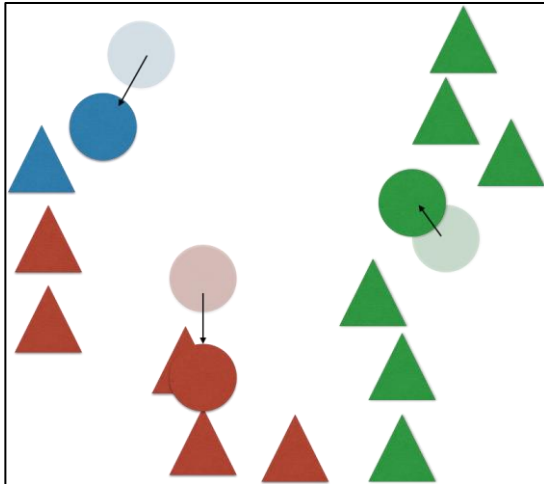**k Means Step 2 - Clusters Created by Associating with Nearest Mean**
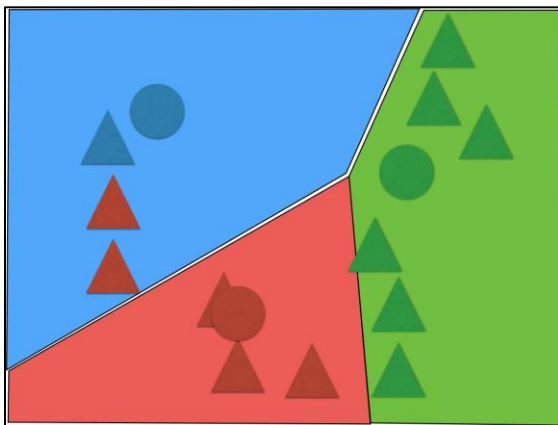
**Figure 5**
**k Means Step 3 - Centroids become New Mean**



**Figure 6**
**k Means Step 4 - Previous Steps Repeated until Convergence**



**Figure 7**
**Resuminer Process**

The complete data analysis process, from reading the data to obtaining the cluster data can be summarized with figure 7. It can be summarized in six steps: reading the résumé data from text files, tokenizing each document into individual words, generating the term frequency-inverse document frequency (tf-idf) matrix, generating the clusters (using k-means), calculating the similarity, and plotting the results. These same steps can also be used when analyzing document style (such as fonts, font size, font style, document format). These data points can be entered into the system, to get a similar analysis about which document styling options make a document stand out. By comparing both of the resulting plots, it is possible to find what keywords attract more interest as well as which styling options seem to be preferred.

By feeding an untested résumé (without any "interaction points") to this system, and running it thought the same analysis, we'll be able to determine which cluster it is closest to and we can make an educated prediction on how it will perform. We can determine if this résumé is closer to the cluster with the desirable scores, or with a low score. As the dataset grows, the approximations to the correct cluster will improve.

## TOOLS

The data analysis will be done using scikit-learn, a library for machine learning for Python. The scikit-learn libraries are well documented and available at http://scikit-learn.org. It is currently used extensively in the data science community for

both academic research, as well as industry projects. Scikit-learn provides a range of supervised and unsupervised learning algorithms. In addition to scikit-learn, Jupyter Notebook (previously known as IPython) is being used as an interactive programming environment and documentation engine. It is available from http://jupyter.org. Jupyter Notebook runs as a server-client application in the web browser, it can be installed and executed locally or in a remote machine. Any notebooks created in this environment is easily portable to a different machine, and the code that is being executed can be documented.
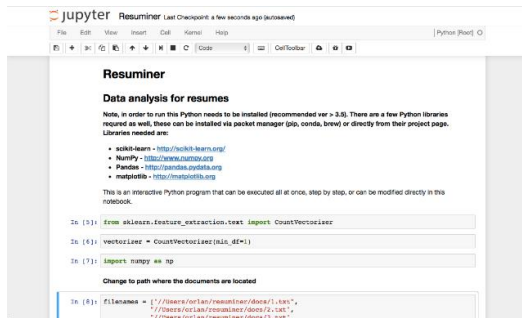


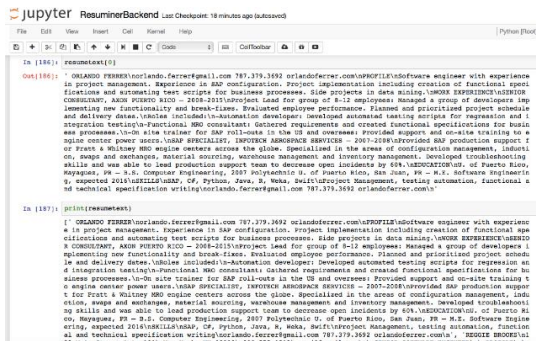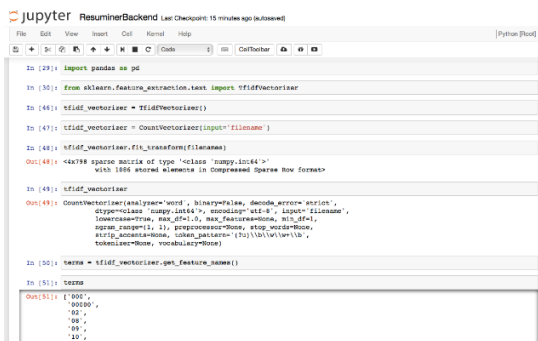**Figure 8**
**Resuminer Running in Jupyter Notebook**



**Figure 9**
**Full Document Loaded in Jupyter Notebook**



**Figure 10**
**Performing tf-idf Analysis in Jupyter Notebook**

## USER INTERFACE

Users will be able to interface with the system through a regular web browser. The interface will be simple, it will display the database entries and will allow entering new data.



**Figure 11**
**Display of the Database Entries in the Browser**

## CONCLUSION

All else being equal, there are infinite ways in which the same information can be conveyed. There are some qualities or styles of writing that make a document (a résumé, for the purpose of this work) stand out from the others. By tracking the interest in each unique document and applying data science techniques it is possible to have an understanding of which characteristics make a résumé stand out over other ones. Impact keywords could be a factor, but there are many others waiting to be discovered by applying other algorithms and gathering more training data.

Below is the result of plotting and analyzing all the points with the current test data in the Resuminer engine.



**Figure 12**
**Resulting Clusters**

## FUTURE WORK

In its current state, only the backend engine has been implemented, so it is not easy to create new data, it needs to be loaded manually. In order to be able to easily add new résumés and track their performance, a web based user interfaced needs to be developed. It will interface with the existing database, and will allow for a simpler data input process. It will need to interface with the current Python backend to process the data. Another feature to be added is to not only analyze the text in résumés, but the styling information as well- such as font, font size, font style, and any others styling properties. Analyzing the style data in conjunction with the text, may yield interesting results as to which styles of résumés grab the attention of recruiters. It is also possible that the current implemented algorithm is not the most efficient one, so other algorithms and data processing techniques can be explored. As more data is collected, there will be more training data available to use and test the predictions made by the algorithm.

## REFERENCES

[1]   P. Tan, M. Steinbach and V. Kumar, *Introduction to data mining*, Boston: Pearson Addison Wesley, 2005.

[2]   W. Eggers, R. Hamill and A. Ali, "Data as the new currency", *Deloitte Review*, no. 13, 2013, pp. 19-23.

[3]   M. Fierro, "A Gentle Introduction to the Basics of Machine Learning", *SciBlog*, 2016.

[4]   I. Witten, E. Frank and M. Hall, *Data mining*. Burlington, MA: Morgan Kaufmann, 2011.