# Smart Employment Navigator - A Job Recommendation System

Aboli Wankhade, Dharma Teja  Kolluri, Saranya Gondeli

Department of Applied Data Science, San Jose State University

DATA 240 Sec 12: Data Mining/Analytic

Dr. Shayan Shams

December 7, 2023

## Background

In the constantly evolving and competitive landscape of today's job market, individuals often face the formidable challenge of identifying job opportunities that align seamlessly with their unique skill sets and career aspirations. This challenge is further amplified by the overwhelming volume of job listings available online, each with its distinct set of requirements and expectations. Our study embarks on the development of an advanced Job Recommendation System, aimed at bridging this gap. This system will leverage a comprehensive dataset, meticulously compiled through web scraping from various job portals, encompassing a broad spectrum of candidate qualifications and job openings.
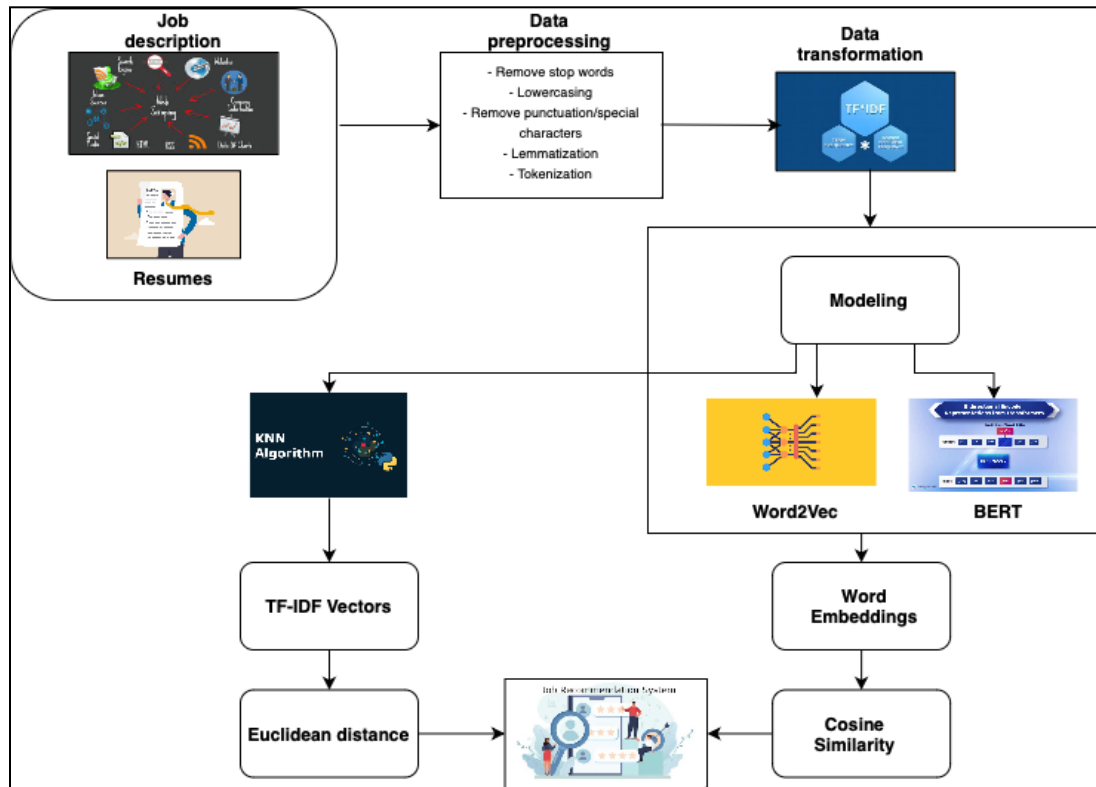
## Motivation

The primary motivation behind our project stems from the urgent need to address the multifaceted challenges confronted by job seekers in the modern digital era. One of the main challenges include navigating through the sheer abundance of online job postings to find relevant opportunities; and counteracting the inefficiencies inherent in current recruitment methodologies, which often lead to protracted hiring processes and missed career opportunities. Our initiative aspires to revolutionize the job-matching process by constructing an effective and intelligent recommendation system. Through meticulous analysis of both candidate profiles and job descriptions, our system is designed to deliver bespoke job suggestions. These suggestions go beyond the conventional criteria of matching skills and instead provide a holistic view by also considering factors such as work experience and technical skills. This comprehensive approach aims to facilitate a more streamlined and personalized job search experience, ultimately leading to more successful and fulfilling career placements.

In summary, our project not only seeks to alleviate the challenges faced by job seekers but also aims to transform the dynamics of the job market by introducing a more refined and user-centric approach to job matching.

## Literature review

Chou and Yu (2020) outlined the deployment of a job recommendation system for job seekers and employers at a career fair, utilizing resumes as input and concentrating on 179 available job listings. They applied cosine similarity for similarity assessment and employed content-based, collaborative filtering, and hybrid recommendation techniques for data training. The system generated online reports for employers, presenting the top 20 candidates derived from the term frequency matrix and document vectors. Jain and Kakkar (2019) constructed a job recommendation system that has a textual approach for mapping the skills of the professionals with jobs in the industries, for this they employed the vector space model also known as Tf-Idf. They calculated the Tf, idf, and tf-idf for every skill and the weighted factor by assigning equal probabilities to each skill. They assigned the most suitable job for a candidate with the one that has the highest weighted factor. Chaudary et al. (2020) employed count vectorization and Tf-Idf vectorization for feature extraction and applied the Gradient Boosting algorithm using a one-vs-rest classification approach to extract skills, education, and experience from job descriptions. They achieved an ROC score of 0.80 for count-based and TF-IDF-based classifiers. Wang and Shi (2022) implemented an article recommendation system using Word2Vec and TfIdf methods by extracting the vectors of the articles and then recommending valid articles based on LSH methodology. Word2Vec was implemented beforehand to get the vectorized values and then Tfidf was used to get the most important vectors from the obtained vectors but hasn't performed any validation or evaluation of the results.

# Methodology



## Data Collection:

In this step, the Job Details dataset was extracted from the Indeed website using Python with libraries such as BeautifulSoup and Selenium to scrape job listings. The key steps in this data extraction process included configuring Selenium WebDriver with Chrome options to navigate and interact with web pages. This setup included handling user-agent strings to prevent detection as a bot. Then each job listing's URL was saved for further detailed scraping, providing a more granular level of data. Using the URLs collected, fetched detailed job descriptions, including job titles, company names, locations, salaries, job ratings, and other relevant details.

Then another job was implemented to retrieve the Skills and Work Experience from the resume which is in PDF format with which the job recommendations will be done

## Data Preprocessing:

Data preprocessing involves cleaning and transforming the job description data from the job listing dataset. This cleaning process involves removing the Null value job descriptions, eliminating duplicate job descriptions, lowercasing the text, removing HTML tags, eliminating punctuations, special characters, and stop words. After performing these data cleaning tasks, cleaned data was utilized for tokenization and then lemmatizing techniques were applied to lemmatize the words.

Similarly, the data that was extracted from the resume should be cleaned by lowercasing the text, removing HTML tags, eliminating punctuations, special characters, and stop words. Similar to descriptions, the dataset, tokenization and lemmatization techniques were implemented. These two cleaned datasets will be used for data transformations in the next

steps. Since the transformation process is different for each model that we are using, during that modeling of that data, transformations will be performed accordingly.

### Modeling and Training:

After preprocessing the data, we employed several models, namely K-Nearest Neighbors (KNN), Word2Vec, and BERT, for our job recommendation system. The specifics of this phase included:

**Model Selection:** Based on the nature of our dataset, we selected three models for implementation: KNN for its simplicity and effectiveness in recommendation systems, Word2Vec for capturing semantic meanings in job descriptions, and BERT for its advanced capabilities in understanding contextual relationships in text.

### KNN Model

After preprocessing the data, we employed TFIDF Vectorizer to transform the job descriptions into numerical vectors based on the TF-IDF scores of words. We used an n-gram approach to capture more context from the text, allowing for a nuanced understanding of job requirements and candidate capabilities.

Building and training the model involved initializing the NearestNeighbors model and we opted for Euclidean distance as the metric for several reasons. Firstly, Euclidean distance is intuitively appealing and easy to understand, as it simply calculates the straight-line distance between two points in the feature space. This characteristic makes it highly suitable for our application, where we need a straightforward yet effective way to measure the similarity between job descriptions and candidate profiles. Secondly, when dealing with high-dimensional data, Euclidean distance provides a reliable measure of closeness, ensuring that our recommendations are both relevant and accurate. The recommendation process involved matching candidate resumes to job descriptions using the trained KNN model. The model identified job descriptions closest to the candidate's resume in the feature space, emphasizing skills and experience compatibility.

### Word2Vec Model

After preprocessing, TfidfVectorizer was used to filter out words with lesser significance, retaining those terms in job descriptions that carry more weight and relevance. Subsequently, we train two distinct Word2Vec models on our job descriptions dataset: one employing the Continuous Bag of Words (CBOW) method and the other using the Skip-gram (SG) approach. Each job description is then transformed into a vector by weighted averaging the vectors of the words using TF-IDF scores it comprises, thereby enabling a semantic and context-based representation of the job descriptions. Similar embeddings were obtained for resume data and employed cosine similarity to gauge the similarity between vector representations of resumes and job descriptions. For each resume, the system identifies the top job descriptions with the highest cosine similarity, suggesting the most relevant job opportunities to the candidate. This approach ensures that the recommendations are not only accurate but also highly personalized.

### BERT Model

BERT (Bidirectional Encoder Representations from Transformers) was employed for its superior ability to capture context and semantics in text. After retaining the most significant terms using TF-IDF, a pre-trained BERT model was implemented on the Jobs description dataset to get the embeddings that serve as a high-dimensional representation of the text, encapsulating complex semantic relationships and similar to the Word2Vec model, job descriptions is then transformed into a vector by weighted averaging the vectors of the words. Similarly embeddings were extracted from Resume data using the above pretrained model. For matching resumes with job descriptions, we use cosine similarity to measure the degree of

similarity between the BERT embeddings of the resumes and the job descriptions. This method quantitatively evaluates the alignment between a candidate's profile and the job requirements. The system then identifies the job descriptions that exhibit the highest cosine similarity with each resume, pinpointing the most fitting job opportunities for the candidates.

## Experiments and Evaluation

### Experiment Design

Initially, we conducted a test using a Data Engineer's resume as input for our model. The process began with extracting the Skills and Experience sections through a specific function. This extracted text was then pre processed and modified to suit the model's input format. Subsequently, we utilized the KNN algorithm along with TF-IDF to generate vectors. Based on these vectors, similarity metrics were applied to recommend the top 5 job matches. Additionally, we employed Word2Vec and BERT models, training them with job descriptions, to derive embeddings and recommend the top 5 job options. Upon analyzing the results, we observed that the majority of the recommended jobs pertained to Data Engineering, confirming the effectiveness of each model in proposing job matches aligned with the resume's content.

To theoretically assess the effectiveness of our models, we undertook a thorough evaluation of our job recommendation system, which incorporates the three detailed models previously described. This evaluation aimed to determine the precision with which each model recommends pertinent jobs when provided with different resumes.

**Data Collection and Preparation:** Our dataset comprised a diverse range of job descriptions, covering various job scopes. We meticulously preprocessed this data and categorized all the job descriptions into 4 categories namely Data Scientist, Data Analyst, Data Engineer and Software Engineer and created a new column in the job description dataset. This column will further be used in evaluating the results of the model.

***Resume Processing and Recommendation Generation:*** Each model processed resumes from the different job categories and subsequently generated job recommendations based on the analysis of the resume content. While uploading the resume to the system for getting recommendations, it will first check the folder name from which the resume was uploaded and compare it to the category of the job description that was given in the recommendation. If all of the recommended categories are matching with the resume category we get 100% accuracy, if not we get varied accuracy according to the recommended jobs.

### Evaluation and Results

The key metric for our evaluation was accuracy defined as the ratio of relevant job recommendations to the total number of recommendations made for each resume and MAP@K defined as a measure of precision across multiple queries at top K results respectively. A recommendation was deemed relevant if it matched the job scope of the resume.

| Models | Accuracy | MAP@K |
|--------|----------|-------|
| KNN | 76% | 69.9% |
| Word2Vec | 82% | 78.4% |
| BERT | 80% | 76.6% |

**KNN Model Results:**

The KNN model showed varying accuracy across different job categories, reflecting its capability to align resumes with suitable job recommendations based on syntactic similarity. An overall accuracy of 76% was achieved by averaging the individual accuracies for each resume. MAP@K score of 69.9% was achieved.

**Word2Vec Model Results:**

The Word2Vec model, focusing on semantic analysis, also displayed differential accuracy across job categories. The overall system accuracy of 82% was achieved through the average accuracy across all processed resumes. MAP@K score of 78.4% was achieved.

**BERT Model Results:**

Leveraging its deep contextual understanding, the BERT model exhibited its proficiency in matching resumes to job descriptions, with varying degrees of accuracy. An overall accuracy of 80% was achieved for the BERT model and calculated as an average of individual resume accuracies. MAP@K score of 76.6% was achieved.

**User Interface**

The job recommendation system is made accessible through a Streamlit-based application, providing a user-friendly interface for candidates to upload their resumes in PDF format. Once a resume is uploaded and processed, the system dynamically displays the top job recommendations. These recommendations are comprehensive, including vital details such as job title, location, and direct links to the job listings.

**Results**

Below are the screenshots of the UI where we upload the resume in the system and recommendations of all the models are displayed.

Word2Vec Recommendations:

| Job_Title | Job_Detail_Link | Company_A_Loc... |
|---|---|---|
| Senior AWS Data Engineer | https://www.indeed.com/... | Cognitive Medical Systems\nRemote Little Rock, AR |
| Big Data Engineer | https://www.indeed.com/... | Uber\nSan Francisco CA |
| Data Analyst | https://www.indeed.com/... | MANDO TECHNOLOGIES INC\nPlano, TX |
| Data Engineer - Python, Snowflake | https://www.indeed.com/... | eTeam Inc.\nManhattan, |
| Senior Cloud Data Engineer | https://www.indeed.com/... | Ursus, Inc.\nHybrid remote in Westlak Village, CA 91362 |

BERT Recommendations:

| Job_Title | Job_Detail_Link | Company_A_Locatio... |
|---|---|---|
| Azure DevOps Solution Architect | https://www.indeed.com/... | MIST GLOBAL LLC\nBellevue, WA |
| AWS Cloud Developer | https://www.indeed.com/... | CTIS, Inc.\nRockville, MD 20850 |
| Full-Stack Cloud Webpage Architect | https://www.indeed.com/... | Realign LLC\nCalifornia |
| Senior Data Engineer Machine Learning (ML) | https://www.indeed.com/... | Halvik\nRemote |
| Senior AWS Database Engineer | https://www.indeed.com/... | PRO IT\nSt. Louis, MO |

## Discussion

Our experiments highlighted each model's unique strengths: KNN and Word2Vec in capturing syntactic elements and BERT in understanding deeper contextual meanings. The variance in accuracy across different job scopes suggests that the effectiveness of each model may depend on specific job characteristics and data quality. These results illuminate the potential of these models in automated job recommendation while also indicating areas for improvement, especially in handling diverse job descriptions and complex resume data. In summary, the evaluation of our job recommendation system using KNN, Word2Vec, and BERT models has provided valuable insights into their effectiveness. These findings will be instrumental in further refining our system, aiming to enhance accuracy and personalization in our job matchmaking process.

## Future Improvement

This recommendation system works better when considering the exact experience of the candidate which is one of the most important and complex characteristics to implement in job recommendations. Integration of Education details for the recommendations could enhance in recommending more user specific jobs. Collaborative filtering can be implemented for the current work which will become a hybrid model to recommend the job using Single Value Decomposition (SVD) modeling techniques.

# Bibliography

Jain, H., & Kakkar, M. (2019). Job Recommendation System based on Machine Learning and Data Mining Techniques using RESTful API and Android IDE. *Job Recommendation System Based on Machine Learning and Data Mining Techniques Using RESTful API and Android IDE*. https://doi.org/10.1109/confluence.2019.8776964

Chou, Y., & Yu, H. (2020). Based on the application of AI technology in resume analysis and job recommendation. *Based on the Application of AI Technology in Resume Analysis and Job Recommendation*. https://doi.org/10.1109/iccem47450.2020.9219491

Chaudary, A., Nasar, Z., Mubasher, M. M., & Qounain, S. W. U. (2020). Extraction of Useful Information from Crude Job Descriptions. *Extraction of Useful Information From Crude Job Descriptions*. https://doi.org/10.1109/inmic50486.2020.9318132

R. Wang and Y. Shi, "Research on application of article recommendation algorithm based on Word2Vec and Tfidf," 2022 IEEE International Conference on Electrical Engineering, Big Data and Algorithms (EEBDA), Changchun, China, 2022, pp. 454-457, doi: 10.1109/EEBDA53927.2022.9744824.