

Towards partial fulfillment for Undergraduate Degree Level Programme
Bachelor of Technology in Computer Engineering

A Final Stage Evaluation Report on:

Job Recommendation System

Prepared by :

Admission No.

Student Name

U18CO009

Tarsariya Keyur

U18CO010

Jadhav Ganesh

U18CO018

Shekhaliya Shubham

U18CO039

Nainuji Jigar

Class : B.TECH. IV (Computer Engineering) 8th Semester

Year : 2021-2022

Guided by : Dr. Sankita J. Patel



**DEPARTMENT OF COMPUTER SCIENCE AND ENGINEERING
SARDAR VALLABHBHAI NATIONAL INSTITUTE OF TECHNOLOGY,
SURAT - 395 007 (GUJARAT, INDIA)**

Student Declaration

This is to certify that the work described in this project report has been actually carried out and implemented by our project team consisting of

Sr.	Admission No.	Student Name
1	U18CO009	Tarsariya Keyur
2	U18CO010	Jadhav Ganesh
3	U18CO018	Shekhaliya Shubham
4	U18CO039	Nainuji Jigar

Neither the source code there in, nor the content of the project report have been copied or downloaded from any other source. We understand that our result grades would be revoked if later it is found to be so.

Signature of the Students:

Sr.	Student Name	Signature of the Student
1	Tarsariya Keyur	
2	Jadhav Ganesh	
3	Shekhaliya Shubham	
4	Nainuji Jigar	

Certificate

This is to certify that the project report entitled Job Recommendation System is prepared and presented by

Sr.	Admission No.	Student Name
1.	U18CO009	Tarsariya Keyur
2.	U18CO010	Jadhav Ganesh
3.	U18CO018	Shekhaliya Shubham
4.	U18CO039	Nainuji Jigar

Final Year of Computer Engineering and their work is satisfactory.

SIGNATURE:

GUIDE

JURY

HEAD OF DEPT.

Abstract

Talent acquisition is the most important task for the success of the company. In the current situation for a given job, thousands of job seeker apply which make it hard for the hiring team to go through each and every resume manually and check for the credibility of the applicant. Similar to the job seeker in the large market thousands of jobs are available which makes finding a suitable job difficult for the user. This project aims to solve this problem by making automation of the resume matching process by using the various technologies for the data extraction from the given text or description and finding similarities between the job seeker's profile and job description. The Doc2Vec embedding is used to learn the features from the highly unstructured text data. An entity ruler is used to extract the skills from the resume and the job description data set. The method of tuning the hyperparameter is also discussed. Later similar job descriptions were found to collect the explicit skills. Furthermore, a combination of the cosine similarity and skill matching is used to find the similarity that gives better results than using them individually.

Keywords: Recommendation System - Hyperparameters tuning - POS tagging - Job recommendation System - Implicit skills - Explicit skills - Doc2Vec - Cosine similarity - Skills Matching - Classification55

Contents

List of Figures	vi
List Of Acronyms	vii
List of Symbols	viii
1 Introduction	1
1.1 Applications	1
1.2 Motivation	2
1.3 Objectives	2
1.4 Contribution	3
1.5 Organization of project report	3
2 Theoretical Background	4
2.1 Resume	4
2.2 Job Description	6
2.3 Overview of Recommendation System	6
2.3.1 Collaborative Filtering Recommenders	7
2.3.2 Content Based Filtering	7
2.3.3 Knowledge Based Recommender	8
2.3.4 Hybrid Recommender System	9
2.4 Doc2Vec	9
2.4.1 Distributed Memory Model (DMM)	9
2.4.2 Distributed Bag-of-Words (DBOW)	10
2.5 Classification :	10
2.5.1 Support Vector Machine Algorithm	10
2.5.2 Naive Bayes	10
2.5.3 Logistic regression	11
2.6 Performance measures for classification problem	11
2.7 Finding similarity	12
2.7.1 Cosine similarity	12
2.7.2 Euclidean distance and similarity	12
2.7.3 Supremum distance	13
2.8 EntityRuler	13
3 Literature Survey	14
3.1 Overview of the recommendation system	14
3.2 Boolean matching technique	15

3.3	Context to recommend	15
3.4	Hybrid job recommender systems	16
3.4.1	A probabilistic hybrid approach	16
3.4.2	A proactive job recommender system	16
3.5	Content-based job recommender systems	17
3.5.1	Machine learned recommender system	17
4	Proposed Work	18
4.1	Logical Development	18
4.2	Proposed Methodology	19
4.2.1	Data Extraction and Data Preprocessing	20
4.2.2	Feature Extraction	21
4.2.3	Implicit and Explicit Skill Extraction	21
4.2.4	Doc2Vec and It's Hyperparameters	22
4.2.5	Computing similarity	22
4.3	Summary	23
5	Simulation And Results	24
5.1	Data Extraction and Pre-Processing	24
5.2	Extracting data from Document	25
5.3	Experiment setup	26
5.4	Skills Extraction Module	26
5.4.1	Skills Extraction	26
5.5	Tuning of Hyperparameters of Doc2Vec model	28
5.5.1	Hyperparameters for training Doc2Vec	28
5.5.2	Hyperparameters tuning for the resume dataset.	28
5.5.3	Hyperparameters tuning for the job description dataset	29
5.6	Finding similar Job descriptions.	30
5.6.1	Need of Finding Similar Job descriptions	30
5.6.2	Merging skills	30
5.7	Recommendation Engine	31
5.7.1	finding similarities	31
5.7.2	Recommending the Resumes to the Job Description	32
5.7.3	Recommending the job descriptions to the resume	32
5.8	Classification of the job description	33
5.8.1	justification	33
5.9	Conclusion	34
6	Conclusion and Future Work	36
6.1	Conclusion	36
6.2	Future works	36
	References	37
	Acknowledgement	39

List of Figures

2.1	Resume Block	5
2.2	Content Based filtering System	8
2.3	Confusion Matrix	11
4.1	Block diagram describing the proposed methodology	19
4.2	flow diagram of proposed system	22
5.1	Chunk	25
5.2	Token Pattern	27
5.3	Phrase Pattern	27
5.4	Accuracy vs. Vector Size	33
5.5	F1 - Score vs. Vector Size	33
5.6	Accuracy vs. Window	33
5.7	F1 - Score vs. Window	33
5.8	Accuracy vs. Epochs	34
5.9	F1 - Score vs. Epochs	34

List Of Acronyms

CF Collaborative filtering

CBF Content-based filtering

RS Recommendation System

JRS Job Recommender Systems

SVM Support Vector Machine

JD Job Description

NLP Natural Language Processing

List of Symbols

$\sqrt{}$ Square root

θ Angle between vector

\sum Summation

\cdot Dot product

α weighted constant for cosine similarity

β weighted constant for skill matching

Chapter 1

Introduction

More and more applications have been broadly developed, and new techniques have emerged to support human decisions suggesting services, products, and various types of information to customers. One field of research in this direction is that of Recommender Systems. Recommender Systems use various techniques and algorithms to isolate irrelevant information from a vast amount of data and generate personalized suggestions of a small subset of them that a user can examine in a reasonable amount of time. The increasing usage of the Internet has heightened the need for online job hunting. The critical problem is that most job-hunting websites display recruitment information to website viewers. Many websites on the Internet give employment opportunities, but the task is tedious as students need to go through a large amount of information, taking lots of time and energy and suffering from unwanted or less helpful information. Jobseekers have to retrieve all the information to find jobs they want to apply for. The whole procedure is tedious and inefficient. One field of research in this direction is that of Recommender Systems.

1.1 Applications

Dealing with the tremendous amount of recruiting information over the Internet, a job seeker always spends hours finding useful ones. With a huge number of different job roles existing today along with the typically large number of applications received, short-listing poses a challenge for the human resource department. This is only further worsened by the lack of diverse skill and domain knowledge within the HR department, required for effective screening. Being able to weed out non-relevant profiles as early as possible in the pipeline results in cost savings, both in terms of time as well as money.[5]

1.2 Motivation

Talent acquisition is an important, crucial, complex, essential task in industries that requires a significant amount of time. Talent acquisition has the most challenging part. The lack of a standard structure and format for a resume makes a short listing of desired profiles for required roles very tedious and time-consuming. Effective screening of resumes requires domain knowledge to understand the relevance and applicability of a profile for the job role. With a massive number of different job roles existing today and the typically large number of applications received, short-listing challenges the human resource department.

In addition, in most Recommendation System (RS)s, the most general application of recommendation algorithms uses Collaborative filtering (CF) algorithms without considering the user's resume and job description. That means candidates' resumes and details of recruiting information. So we proposed an improved algorithm based on Content-based filtering (CBF). Our aim is to give an effective method of recommendation for online job hunting and talent hunting. We hope to offer candidates a personalized service that can help them find ideal jobs quickly and conveniently.

1.3 Objectives

The e-recruiting platforms are usually based on Boolean search and filtering techniques that cannot sufficiently capture the complexity of a person-job fit as selection decisions. Much literature has applied the recommender system concept to the job problem. Recommendation between entities of the domain: users and opportunities

The job recommendation problem is a bidirectional recommendation between job-seeker and job.[3] Two viewpoints are distinguished: from recruiters and job seekers. The recruiters generate the job description by determining the set of requirements and constraints on skills, expertise levels, and degrees. The job-seeker, on the other hand, generates the candidate's resume by specifying the academic background, previous work experience and skills[3].

Based on the requirement that a good match between jobs and persons needs to take into account both the preferences of the candidate and the preferences of the recruiter to recommend the job.

1.4 Contribution

This project aims to recommend relevant resumes and jobs for particular job titles and resumes respectively from large data. The project uses Doc2Vec embedding technique for vectorising the job description and the cosine similarity measure is used to find similarities between two vectorised judgements. And For Feature Extraction Rule-based entity Recognition technique used.

This project performs the following in order to achieve these objectives.

1. Resumes and job description's Feature extraction: To Extract the various skills efficiently from resumes and job description spacy predefined entity-ruler used.
2. Training Doc2Vec embedding on each job description and resume: Document Embedding and a bag of the word is trained on each Job description and resume using various features like skills and job description words or tokens.
3. Finding similarity and skill matching: To find document similarity cosine used and also extracted skill used as a combination of score of document similarity to recommend top N job for each queried resume and top n resumes for each queried job description.

1.5 Organization of project report

This chapter covers the introduction to the project along with its application, motivation, objective and overview. Chapter 2 presents a theoretical background of the terminologies in the job recommendation as well as other important concepts needed to understand the project better. Chapter 3 is the Literature Survey which summarises the work done in the job recommendation to recommend the job. Later in Chapter 3, a review of various job recommendation algorithms is discussed. Our proposed methodology and logic development of the same is covered in Chapter 4. Chapter 5 discusses the result and simulation of our experiment. Chapter 6 ends the report with conclusion and future work proposed.

Chapter 2

Theoretical Background

This chapter discusses various recommendation techniques and how this technique works and the advantages and limitations of the various recommendation techniques. Finally, it contains an overview of different RS.

2.1 Resume

A resume is a formal document created by a job seeker to list their qualifications for a particular position. A customized cover letter is typically sent with a resume, in which the candidate expresses interest in a specific job or organization and highlights critical information on the CV.

Skills :

knowledge of different technologies in which job seekers have experience or he/she learns that

Language: java, python, c++, c, HTML

Technologies: Spring boot, Django, node.js.

Past Experience, Internships and Certification :

Job seeker's experience in the previous companies as a full-time worker or an intern. Job Seeker does certificates of Internships and different courses.

Includes the List of the companies that Job Seeker worked for, employment/internship dates, their positions, and brief descriptions of their work responsibilities, enriched with keywords and enhanced with bulleted lists of quantifiable achievements done in Job / Internship. It also includes a List Of certifications that a Jobseeker has.

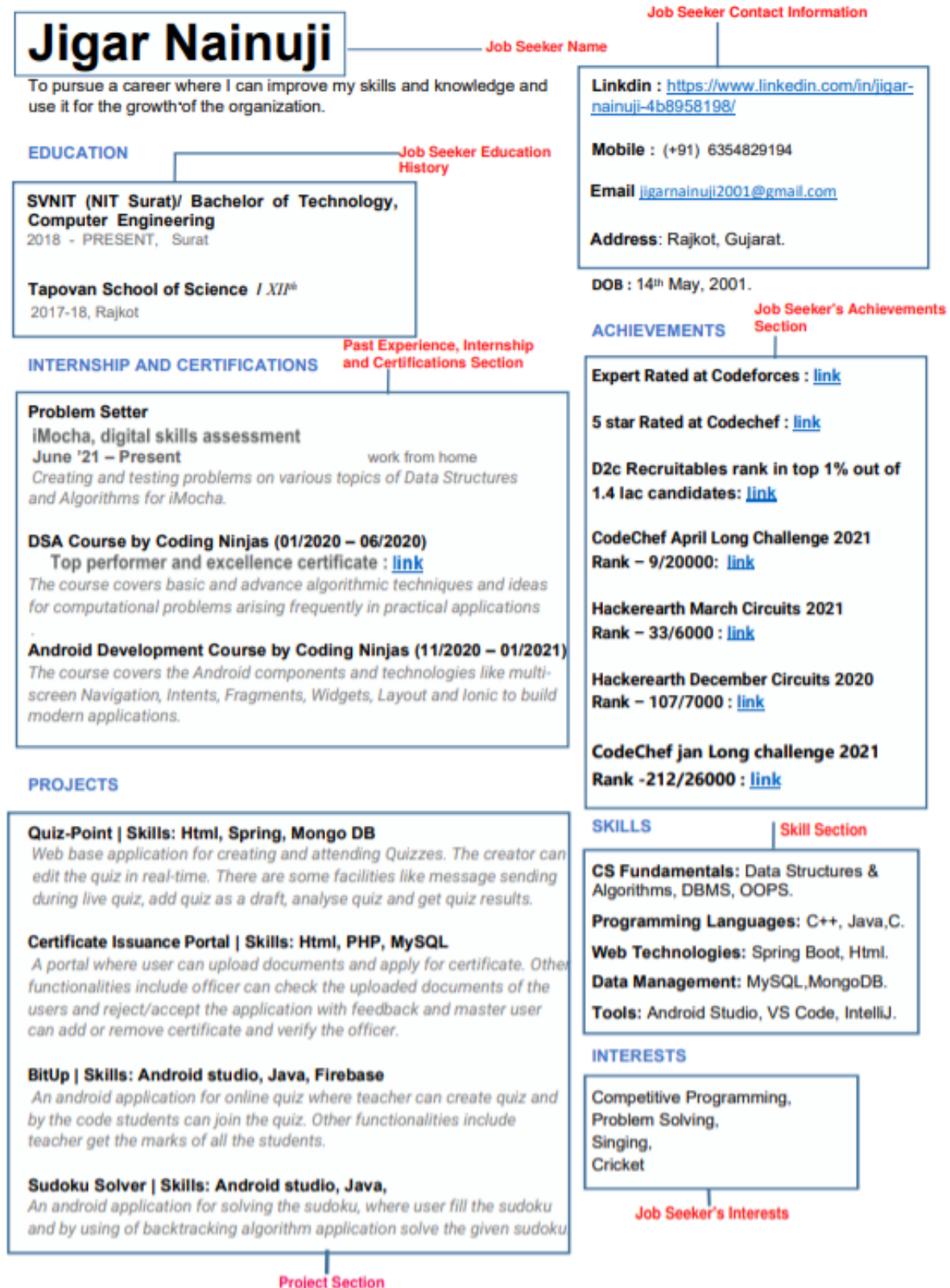


Figure 2.1: Resume Block

Projects :

Current students or recent graduates can use university projects to highlight their relevant skills in a more practical setting. For others, it can include freelance projects and personal projects also.

Education :

Includes the job seeker's educational background. Employers look for a few essential aspects of resumes, including education. This information will give interviewers a better idea of their background, which might help them figure out if they are a good fit for the job.

Interest :

Includes a job seeker's shared interests. The themes and broad ideas that you appreciate in your daily life are considered interests. They are usually more all-encompassing concepts that you are enthusiastic about. Interests are broad notions that guide your real-life decisions and activities.

Tools :

Includes Common tools That a job seeker used in the past or wants to work with them.

Personal profile :

Includes Job seeker personal information like Contact Number, Email id, URLs of portfolios, and LinkedIn profiles.

2.2 Job Description

The Job description provides a high-level overview of the role, level, and scope of responsibility. A job description is a document containing information about the job title, job purpose, required qualification, required minimum experience, knowledge, skills, and abilities where experience Identifies the minimum number of full-time experience required in terms of years and the type of work experience that an employee needs to be qualified for the job.

2.3 Overview of Recommendation System

We often seek suggestions from friends, colleagues or known ones whenever we want to buy something like a refrigerator, TV, mobile phone or washing machine or even when planning for the Trip or which book to refer to or which movie or song for entertainment. Even with their best intentions, these friendly suggestions sometimes do not fit us or are effective in our case. Not just in decision making plays an imperative part to settle on choices which help to pick up benefits by connecting the best alternative as a suggestion. The point is that it is very arduous

to highlight a precise suggestion on the items on which we might be interested.

One field of research in this direction is that of Recommender Systems. RSs are tools that use various techniques and algorithms to isolate irrelevant information from a huge amount of data and generate personalized suggestions of a small subset of them that a user can examine in a reasonable amount of time. An RS is an intelligent computer-based technique that predicts on the basis of users' adoption and usage and helps them to pick items from a vast pool of online stuff[17], Or it identifies the users' needs automatically by inferring the needs from the user's item interactions. Alternatively, the recommender system asks users to specify their needs by providing a list of keywords or through some other method[4].

Following are the approaches of the job RS.

1. Collaborative Filtering recommenders
2. Content Based Filtering recommenders
3. Knowledge-based recommenders
4. Hybrid Recommenders

2.3.1 Collaborative Filtering Recommenders

CF uses similarity between users and items simultaneously to provide recommendations. CF RS finds users with similar interests as the target user and suggests recommendations to him/her based on their liked items. The key function in CF RS is the computation of similarities among users. [6]

2.3.2 Content Based Filtering

CBF is based on a description of the item and a profile of the user's preferences. Items are recommended having similar content information to those a user has. CBF analyses the similar characteristics of the item and target users based on that build the profile for the user. In this system keywords are extracted from the item and user's description to find similarity between them. only the most descriptive features are used to model an item and users and these features are typically weighted.

As show in figure 2.2 it extracts attributes from the user and also from the job description later find similarity between them using various known technique.

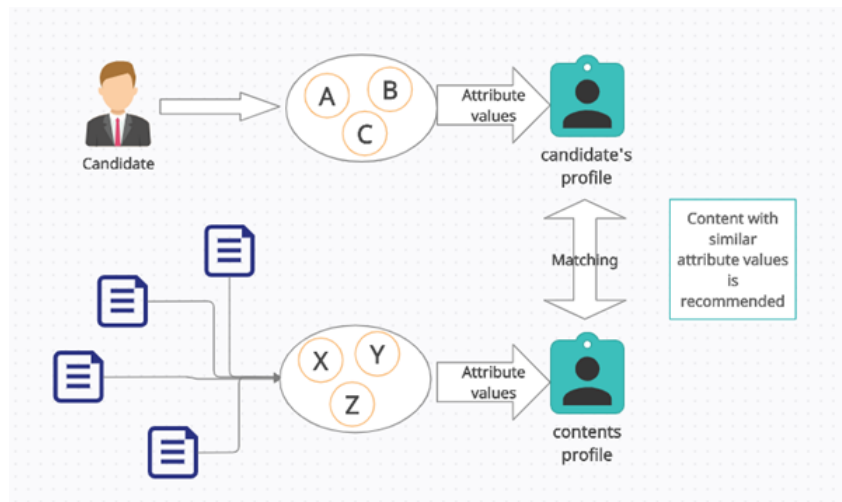


Figure 2.2: Content Based filtering System

Advantages

- The model doesn't need any data about other users, since the recommendations are specific to this user. This makes it easier to scale to a large number of users.
- The model can capture the specific interests of a user, and can recommend niche items that very few other users are interested in.
- New items can be recommended; just data for that item is required

Limitations

- The model can only make recommendations based on existing interests of the user. In other words, the model has limited ability to expand on the users' existing interests.
- Business cannot be expanded as the user does not try a different type of product.

2.3.3 Knowledge Based Recommender

To recommend the items which are less frequently used. In this technique, the relationship between user and item can be explicitly modelled. By using the knowledge of an item based on rules and patterns, we can recommend how a particular item is suitable for the user.[3]

Advantage

- It can recommend the new item to the user even when item is new in the system as it solves the problem of cold start

2.3.4 Hybrid Recommender System

Hybrid recommender technique is a mix of other techniques to override the drawback of the existing techniques.

As show in figure 2.4 Hybrid Recommender System takes input from the all suitable recommendation technique for the job recommendation problem.

2.4 Doc2Vec

Doc2Vec provides a technique for extracting word embedding from corpus paragraphs. Instead of vector representations of the entire corpus, we can think of these word vectors as paragraph vectors. Doc2Vec modelling is a technique for learning text vector representations.[9] Vectorization is also applied to short portions of text documents, ranging from phrases or sentences to huge documents. This technique can be used to anticipate words in paragraphs. These models can predict a word in a particular context by concatenating the section with numerous word vectors from a paragraph.

Using the word vectors as a starting point for learning is a good concept. Even though we can predict the next word in the sentence using a random factorization word vector, we can also expect the next word using the paragraph vector, which contains information about the semantic relationship and can help us get better results.

Doc2Vec models have two variants :

- Distributed Memory Model
- Distributed Bag Of Words

The continuous Bag Of Word Models is similar to the Distributed Memory Model, while distributed bag of words is alike to skip gram model.

2.4.1 Distributed Memory Model (DMM)

Distributed Machine Learning is a multi-node Machine Learning system that improves performance, increases accuracy, and scales to larger input data sizes. It reduces errors made by the machine and assists individuals in making informed decisions and analyses from large amounts of data. Distributed machine learning algorithms have evolved to handle enormous data sets.

Distributed ML algorithms are integral to large-scale learning because of their ability to allocate learning processes onto several workstations to enable faster learning algorithms.

2.4.2 Distributed Bag-of-Words (DBOW)

The Distributed Bag-Of-Words (DBOW) model aids in determining the context words associated with a target word. The allocated memory model approximates the word using the context of surrounding words, whereas the dispersed bag of words model approximates the word using the context of surrounding words. The target word is used in the dispersed bag of words model to match the word's context. When comparing a distributed bag of words with skip-gram models, remember that the skip-gram model takes the target word as input. We find that the distributed bag-of-words models outperform the distributed memory models in several comparisons.

2.5 Classification :

2.5.1 Support Vector Machine Algorithm

SVM stands for Support Vector Machine and is one of the most widely used Supervised Learning algorithms for classification and regression tasks [12]. The purpose of the SVM method is to find the best line or decision boundary for categorising n-dimensional space into classes so that fresh data points can be placed in the correct category rapidly in the future. This best decision boundary is called a hyperplane. The extreme points(vectors) that assist create the hyperplane are chosen via the Support Vector Machine. Support vectors describe these severe circumstances, and the process is known as a Support Vector Machine.

2.5.2 Naive Bayes

Naive Bayes is a classification technique established on Bayes' Theorem and the presumption of predictor independence [13]. A Naive Bayes classifier, in simple terms, asserts that the existence of one part in a class is irrelevant to the presence of any other feature. The Naive Bayes model is simple and useful for large data sets. Naive Bayes outperforms even the most advanced classification systems due to its simplicity.

2.5.3 Logistic regression

The method of modeling the probability of a discrete result given an input variable is known as logistic regression. The most frequent logistic regression models have a binary outcome, which might be true or false, yes or no, and so forth. Multinomial logistic regression is a type of logistic regression that can describe events with more than two distinct outcomes. Logistic regression is a useful analysis tool for determining if a fresh sample fits best into a category in classification tasks. Because our project requires the classification of documents into more than two categories, we will use multinomial logistic regression.

2.6 Performance measures for classification problem

Confusion Matrix is an evaluation matrix used to evaluate a prediction model's performance by comparing the actual values to the predicted values. It has 4 parameters that are used to calculate different measures: True Positive (TP), True Negative (TN), False Positive (FP), False Negative (FN).

		Actual Values	
		Positive (1)	Negative (0)
Predicted Values	Positive (1)	TP	FP
	Negative (0)	FN	TN

Figure 2.3: Confusion Matrix

Accuracy:

Accuracy is the ratio of correct predictions to the total number of predictions. The accuracy evaluation matrix is observed to be a good evaluation metric when a balanced data set was considered.

$$Accuracy = \frac{TP + TN}{TP + FP + FN + TN} \quad (2.1)$$

Precision:

Precision is the ratio of correct positive predictions to the total positive predictions. Higher

precision indicates that false positives are fewer.

$$Precision = \frac{TP}{TP + FP} \quad (2.2)$$

Recall:

Recall is the ratio of correct positive predictions to the total actual true values. Higher recall indicates that false negatives are fewer.

$$Recall = \frac{TP}{TP + FN} \quad (2.3)$$

F1 Score:

F1 score is harmonic mean. F1 score is the mean of precision and recall, which aims to generalise precision and recall metrics. This evaluation measure has been taken into consideration since the testing dataset is imbalanced.

$$F1-Score = \frac{2 * Precision * Recall}{Precision + Recall} \quad (2.4)$$

2.7 Finding similarity

2.7.1 Cosine similarity

The similarity of two vectors in an inner product space is measured by cosine similarity. It detects if two vectors are pointing in the same general direction by measuring the cosine of the angle between them. In text analysis, it's frequently used to determine document similarity.

Let say There is two Vector A and B then we can find similarity between them using cosine

$$similarity = \cos(\theta) = \frac{AB}{|A||B|} = \frac{\sum_{i=1}^n A_i B_i}{\sqrt{\sum_{i=1}^n A_i^2} \sqrt{\sum_{i=1}^n B_i^2}} \quad (2.5)$$

2.7.2 Euclidean distance and similarity

Let say There is two Vector A and B then we can find similarity between them using Euclidean distance

$$d(p, q) = \sqrt{\sum_{i=1}^n (p_i - q_i)^2} \quad (2.6)$$

Similarity can also be find using Euclidean distance using below formula

$$\frac{1}{1 + d(p, q)} \quad (2.7)$$

2.7.3 Supremum distance

The weighted euclidean distance can be calculated by assigning a weight to each property based on its perceived value.

$$d(i, j) = \sqrt{w_1|x_{i1} - x_{j1}|^2 + w_2|x_{i2} - x_{j2}|^2 + \dots + w_m|x_{ip} - x_{jp}|^2} \quad (2.8)$$

Weighting can also be applied to ther distance measures as well.

2.8 EntityRuler

The EntityRuler is a factory that allows one to create a set of patterns with corresponding labels. A factory is a set of preloaded classes and functions that perform set tasks. In EntityRuler, the factory at hand allows the user to create an EntityRuler, first we have to give it a set of instructions, and then use these instructions to find and label entities.

Chapter 3

Literature Survey

This chapter briefly discusses the existing literature in the field of Job recommendation (JRD) system. Extracting the data from the resume of the user and job profile and relate them to show recommendations.

3.1 Overview of the recommendation system

Because of the Internet, companies have changed their hiring process by using the online platform. Companies choose to use online platforms because recruiting the appropriate person is a challenge faced by most companies. The unavailability of specific candidates in some skill areas has long been identified as a significant obstacle to the company's success. RSs are a useful alternative to search algorithms since they help users to. These are the systems that help us to select similar things whenever we select something online. The concept of understanding a user's preference by their online behaviour, previous purchases, or history in the system is called a recommender system [14]

The recommender systems techniques can be used to address the problem of information overload by prioritizing the delivery of information for individual users based on user preferences. Recommender Systems are tools that use various techniques and algorithms to isolate irrelevant information from a huge amount of data and generate personalized suggestions of a small subset of them that a user can examine in a reasonable amount of time. So The task of the RS is to help the user to concentrate on the area of interest.

3.2 Boolean matching technique

Online channels like Internet job portals, social media applications, or a firm's career website have driven this development. While the companies established job positions on these portals, job-seekers use them to publish their profiles. For each posted job, thousands of resumes are received by companies. Consequently, a huge volume of job descriptions and candidate resumes are becoming available online. This vast volume of information gives a great opportunity for enhancing the matching quality; this potential is unused since search functionality in recruiting applications is mainly restricted to Boolean search methods. The need increases for applying the recommender system technologies that can help recruiters to handle this information efficiently.[7]

3.3 Context to recommend

We must consider unary attributes such as individual skills, mental abilities, and personality that control the fit between the individual and the tasks to be accomplished [16], as well as the relational attributes that determine the fit between the individual and the upcoming team members.

In this context, literature usually distinguishes between

1. person-job
2. person-team
3. person-organization fits

Many types of research have been conducted to discuss different issues related to the recruiting problem as well as the application of recommender system technologies. However, job recommendation is still a challenging domain and a growing area of research.[3]

Some of the followings are existing systems for a job recommendation.

- Hybrid job recommender System
 - A probabilistic hybrid approach
 - A proactive job recommender system
- Content-based job recommender systems
 - Machine learned recommender system

3.4 Hybrid job recommender systems

3.4.1 A probabilistic hybrid approach

The recommendation approach used both concepts: CBF and CF simultaneously. It understands the individual preferences as a combination of preference factors. In a basic approach for CF, we look at each value of user/ object pairs (x, y) , where x is a set of users and y is a set of objects. The model can then be represented as a variable z which is associated with each value of (x, y) , assuming that x and y are independent conditioned on z . The model parameters are then estimated using the Expectation Maximization (EM) algorithm.[8]

This model produced a rating matrix that assigns assessed values to candidate's profile containing the probability that recruiter x rates candidate y with value v . Later, they defined $v =$ "qualified", "not qualified". Then, they transformed the rating matrix by replacing variable y with a variable a to represent the attributes that were extracted from the candidate resumes. As many attributes are assigned to several profiles, we will see the attribute several times with different values v .

To improve the match between people and jobs: a CV-recommender and a job recommender, separately. In the first step, they built a system recommending CVs that are similar to resumes previously selected by the recruiter for a specific job profile. In the second step, they developed a second RS that recommends jobs to candidates based on their preference profiles which are in turn based on previous preference ratings.[8]

Limitation

- It answers in binary only either 0 or 1 cannot answer in rank wise to give recommendations.

3.4.2 A proactive job recommender system

The proactive recommender system is an adaptive system that attempts to integrate the idea of recommender systems.[15] This system contains five components: web spider, ontology checker, profile analyzer, preference analyzer, and user interface generator. Web spider is a parser that periodically acquires job information from an exterior source. The ontology checker matches information with ontologies and performs the classification. Then, the job data is stored in a pre-designated form. The profile analyzer makes the recommendations whenever the users modify the group of favorites by comparing the weight differences with current open jobs. Then, a list of recommended jobs is generated. Finally, the preference analyzer deduces the explicitly

defined user's preferences and gives a recommendation for preferred jobs after calculating the similarity of jobs to the user's preference.[5]

Limitations

- One way recommendation only recommend to the job seeker
- Cold start problem as user change profile

3.5 Content-based job recommender systems

3.5.1 Machine learned recommender system

The recommendation problem is treated as a supervised machine learning problem. They build an automated system that can recommend jobs to applicants based on their past job histories, in order to facilitate the process of choosing a new job. An item in this learning model represents a person who is hired in an organization. Each item is characterized by a set of features extracted from the candidate's resumes. Given a person who is currently working in an organization, they want to predict the next organization. If the accuracy of such predictions is sufficiently high, the model can be used to recommend organizations to employees who are seeking jobs. This approach uses all past job transitions as well as the data of both employees and organizations to predict an employee's next job transition. They train a machine learning model using a large number of job transitions extracted from person profiles available on the web.[11]

Limitations

- As it takes previous or historic data into the consideration, the problem of sparsity and cold start could occur.

Chapter 4

Proposed Work

The Project aims to analyse the information from the resume of the candidate and the job description posted by the employer/organiser and extract the useful information from the resume and job description to find the similarity. For that, job descriptions are classified into several categories to identify the role of the job. Then Based on the various similarity methods and using extracted implicit and explicit skills, the model can recommend the job to the candidate or it helps the organiser/employer for shorting the candidate in the initial phase.

4.1 Logical Development

Finding the relevant job based on the candidate profile from a large amount of information present on the internet is time-consuming and cumbersome. Also for particular jobs posted on the internet, a massive number of applications are received to the employer which makes it hard to go through each resume manually. Thus, an effective technique is required to recommend an appropriate job for the candidate and to shortlist the most suitable resumes for an employer. For that purpose, From the resumes and job descriptions features are extracted like tools/technologies, education background, skills and expertise. And to handle large data and to identify job description title or class the job descriptions are classified into several categories using the classifier model such as Logistic Regression, Naïve Bayes and Support vector machine. Using the above three models one having better accuracy will be chosen as a model. And this features and implicit and explicit skills or job description is used to find or recommend a suitable and appropriate job to the candidate over the pool of jobs. Jobs are recommended by finding the similarity of the documents which can be calculated using the cosine similarity to recommend the top-n matching items.

4.2 Proposed Methodology

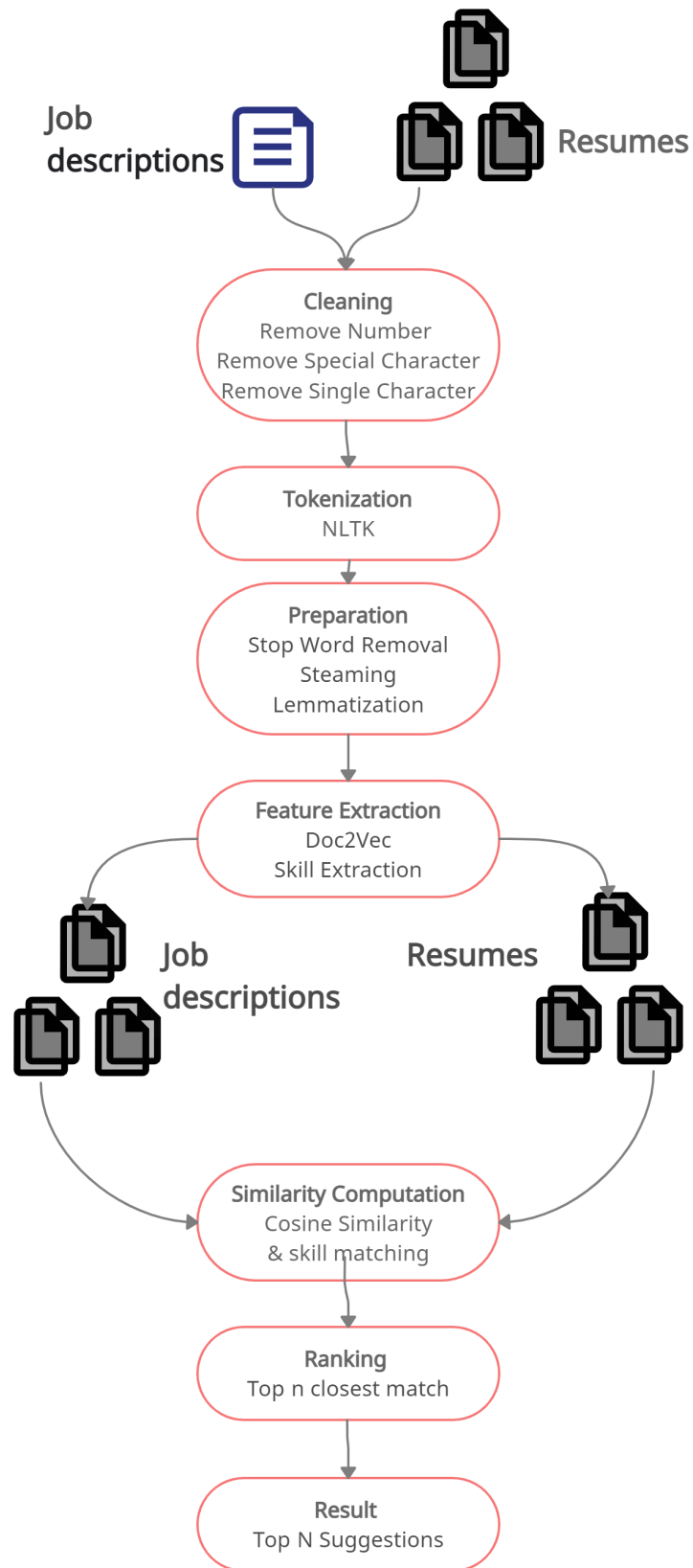


Figure 4.1: Block diagram describing the proposed methodology

A Systematic approach is proposed in order to achieve the objectives mentioned in the above Figure 4.1. The proposed work consists of the four mother phases. The first phase is a Data extraction and Data preprocessing. The next phase is Feature Extraction where various features will be extracted from the resume and the job description which are useful to finding or recommending a job. The third phase is the finding explicit skills for job descriptions where for each job description top n similar job description's skills are used as explicit skills and finding these explicit skills can be done efficiently using classification of the job description. The final step is about computing the similarity between the documents to recommend using cosine similarity and skill matching.

4.2.1 Data Extraction and Data Preprocessing

The resumes are collected using web scraping and this scrapped resumes can contain the noise or undesired character or spacing, blocked structure not as the simple line so to extract the text which is related to its context the data preprocessing happens. Data preprocessing is needed for both resumes and job descriptions.

Data preprocessing[2] is one of the most required steps in data analysis in order to achieve maximum accuracy and throughput . It includes techniques to remove incomplete data, making data consistent and ready to use for experiments. Mostly, a library called pandas[10] is used for such preprocessing. Preprocessing text involves converting it into a form that is predictable and analyzable for the task. Data cleaning, data transformation, and data reduction are procedures involved in data preprocessing.

Data Cleaning involves handling of missing data, noisy data. Strategies to handle missing data involve removing the tuples, filling the missing values. In noisy data Lowercasing, Stemming, Lemmatization, Stop words removal such as 'a', '.', ',', 'an', 'the', removing extra spaces, new line and the unknown character from the description such as ' ', 'a', 'o', 'u' outlier analysis can be done to clean irregular and inconsistent data such as experience of the candidate. Lowercasing is one of the simplest and most effective forms of text preprocessing.

Tokenization is essentially splitting a phrase, sentence, paragraph, or an entire text document into smaller units, such as individual words or terms. Each of these smaller units are called tokens. The Resume usually doesn't consist of the whole english sentence, it only consists of the words so the resume can be tokenized as words using NLTK. While the Job description consists of full english sentences which can be done using tool gensim.

Stemming is the process of reducing inflection in words (e.g. experienced ,experience) to their

22 root form (e.g. experience). The "root" in this case may not be a real root word, but just a canonical form of the original word. NLTK provides the implementation of stemming. Stemming is desirable as it may reduce redundancy as most of the time the word stem and their derived words mean the same. Lemmatization is very similar to stemming, where the main aim is to remove inflections and map a word to its root form. The only difference is that lemmatization tries to do it properly. It doesn't just remove things off, it links words with similar meaning to one word.

4.2.2 Feature Extraction

To extract the features from the resume such as tools, languages, work experience, projects model can use a bag of words which match with the heading of each block described in block structure of the resume.

- **Tools:** tools, IDE, editors
- **Work experience:** work experience, past history, job role, achievement
- **Projects:** projects, publication, certificate

From the job description useful features to recommend jobs such as tools, programming languages and required skills, any brownie or plus point can be useful which can be extracted using

- **Required skills:** required skills, Technical skills, soft skills
- **Plus point:** plus point, brownie point, good to have, advantageous

4.2.3 Implicit and Explicit Skill Extraction

Implicit Skills : Implicit Skills are the skills which are defined in the Resume and Job Description.

Explicit Skills: Explicit Skills are the skills which are derived by comparing one Job Description to other Job descriptions.

Implicit skills of resumes and job descriptions are extracted using rule-based entity recognition in which skills rules are predefined and this was implemented using spacy python library. Spacy library is used for advanced natural language processing.

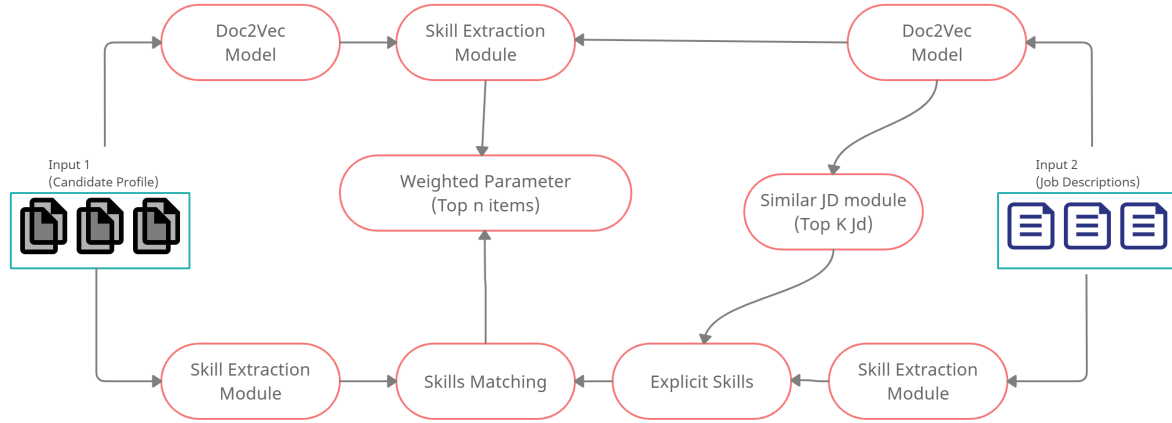


Figure 4.2: flow diagram of proposed system

For Resume only implicit skills are required because adding explicit skill will improve the candidate skills which is not appropriate for job recommendation and it can cause insufficient and unnecessary recommendation.

For job description, explicit skills are merged from other top n similar job description skills which can be found using cosine similarity. In the job description, explicit skills can be combined because all required skills are not mentioned in each job description and combining skills will help to improve similarity finding method accuracy.

4.2.4 Doc2Vec and It's Hyperparameters

Doc2Vec model is used to create a vector representation of a group of words taken collectively as a single unit. It doesn't only give the simple average of the words in the sentence. Doc2Vec model used to represent words of resume and job descriptions. This vector will be further used to find similarity between documents like job descriptions and resume.

The hyperparameters of the training Doc2Vec model for which optimal values need to be found are Vector size, Window and Epochs. This hyperparameter is chosen optimally from experiments on resumes and job descriptions to embed doc2vec model.

4.2.5 Computing similarity

In case of bag of words or vector of words, In First cosine similarity between the vector of the job description and the resume can be computed to find likeness and Second with explicit/im-

plicit skill score was computed and this both similarity finding method combined using various combination of weighted parameter to finally find top n resumes for particular job description and top n job descriptions for particular resume. [1]

4.3 Summary

This project, Thus systematically presented an overview of the methodology proposed. By dividing our problem into parts such as skill extraction and the similarity computation using weighted parameter, we intend to modularize our problem solving approach

Chapter 5

Simulation And Results

This chapter describes our experimental simulations and results. It consists of a description of our initial data preprocessing techniques, skill extractions techniques, doc2vec model training on the dataset, and computing the similarity between the documents to fetch top n results. It also discusses the selection of optimal hyperparameters for the doc2vec model. Later, the classification of the job descriptions is done using the various algorithms using Logistic regression and SVM with linear kernel, and their corresponding results are discussed. Finally, this chapter presents the performance obtained by various combinations of weighted constants.

5.1 Data Extraction and Pre-Processing

The experiments have been conducted using the dataset of the dice website, which is the leading American website for IT jobs. This dataset contains the job descriptions of the 22000 and the title of the jobs. The dataset of job descriptions contains lots of variants in the job title name, so to overcome this problem, various proper job titles are listed, and as a part of data preparation, the similarity between job titles is calculated using the cosine similarity one having result more than the threshold value of 90 percentage is selected for the model training. The resume dataset is prepared using the web scraping of the livecareer website, which is a resume builder tool. The dataset contains resumes of 9 categories. Each has around 50 resumes.

In the working prototype of this experiment, the user is supposed to upload a job description and resume in the pdf document file. The tika tool and pdfminer are used to extract the data from the pdf document. The characteristic of the tika tool is that it extracts the information block-wise from the pdf so that valuable and related information remains together. At the same time, pdfminer helps extract the information line by line. Extracting the data from the pdf files may contain noise such as non-ASCII keywords. For that, Traditional text preprocessing techniques

are used to remove the inconsistency and non-uniformity in the text. The following traditional preprocessing steps are applied to highly unstructured data.

The following traditional preprocessing steps are applied to highly unstructured data.

- Transforming the text content into the lower case
- Individual words are extracted and tokenized with the help of white spaces and line breaks
- stemming and lemmatization are performed to convert the given word into a root word
- Elimination of white spaces, newline and punctuation marks.
- Removal of stop-words such as "a", "an", "the" etc.
- Removal of non-ASCII keywords

5.2 Extracting data from Document

This section defines the various techniques utilised for extracting the needed information such as the name, email, phone numbers, or experience from the document file. The regex string matching method extracts the emails and phone numbers from the document file by scanning the entire Document after preprocessing the data.

The POS tagging method extracts the candidate's name and experience. After POS tagging on the text, each token is categorised in the verb, adverb, noun Etc., based on the corpus context. In the resume document, generally, the candidate's name comes as two consecutive nouns. This is implemented using the nltk library. The regex grammar for it prepared, and it passed to the chunking by traversing all possible chunked tree names can be found.



Figure 5.1: Chunk

To extract the experience from the text is tokenized by the new line character and tokenized by line. An important observation is that lines containing the experience keyword have a cardinal token. This way can find experience from the document text. In the same way, one can find a brownie point and a minimum qualification.

5.3 Experiment setup

The project has four significant motives. The first is to extract the skills from each resume and job description using the entity ruler. The second is tuning the hyperparameters to find the optimal value to embed the doc2vec model on the job description and resume datasets. The third is to find similar job descriptions to merge the skills of the most relevant job descriptions. The fourth is to perform the 2-way recommendation using the cosine similarity. The experiments are performed on a machine with an Intel Core i5-8250U processor running at 1.60 GHz using 8 Gb of Ram with 240 Gb of SSD, running on windows 10.

5.4 Skills Extraction Module

This section discusses the method for extracting the skills from the resume text and the job description text. The skills are extracted using the spacy's entity ruler pipeline. Furthermore, show the result of the extracted skills.

5.4.1 Skills Extraction

A candidate acquires skills through formal education, internships, or previous job experience. The candidate may start identifying relevant jobs based on these acquired skills. The key function of a job search engine is to help the candidate by recommending those jobs which are the closest match to the candidate's existing skillset. This recommendation can be provided by matching the candidate's skills with the skills mentioned in the available job descriptions.

To extract the skills Entity ruler of the spacy is used. In which the Entity patterns are as the dictionary with two keys. The first 'label' is the label assigned to the 'SKILL' entity. Furthermore, the second 'pattern' is a match pattern. For which algorithm is going to search over the text. Entity ruler accepts two types of patterns. The first is Phrase patterns to search for the skills having a single word, for example, c, java, and python. The second is Token patterns to search for the skills having multiple words, for example, data structure and computer network.

Figure 5.2: Token Pattern

```
{
  "label": "SKILL",
  "pattern": [
    {"LOWER": "aws"},
    {"LOWER": "lambda"}
  ]
}
```

Figure 5.3: Phrase Pattern

```
{
  "label": "SKILL",
  "pattern": [
    {"LOWER": "java"}
  ]
}
```

Then the spacy en-core-web-lg-model loaded where more than 2000 patterns are added in the pipeline. After which, the text is passed through the NLP model, and the whole text is searched against Token having the 'SKILL' Label. Below is the approach for the same.

Algorithm 1: Skill Extraction Algorithm

```
procedure getSkill( text )
  doc ← nlp( text )                                ▷ Apply nlp model
  skills ← []                                         ▷ Initialize list
  for ent in documents do

    if ent.label = "SKILL" then
      skill.append(enttext)
    end if

  end for
  return uniqueSkills( subset )                    ▷ return unique skill set
```

Below is an example of the extracted skills from the text.

```
{
  skills = [
    "java", "data structure", "machine learning"
  ]
}
```

Following this method, the resume dataset and the job description dataset are passed from this algorithm, and extracted skills are stored along with the id of the resume and stored in the CSV file for the further processing of the data.

5.5 Tuning of Hyperparameters of Doc2Vec model

This section describes the training of the doc2vec model on the job description and resume dataset and also discusses the finding of the optimal value of the hyperparameters to train the model and the Evaluation function to find the optimal value of the hyperparameters for the job description and resume dataset.

5.5.1 Hyperparameters for training Doc2Vec

The hyperparameters of training the Doc2Vec model for which optimal values need to find out are:

Vector size (vSize): dimensionality of the feature vectors

Window (w): the maximum distance between the current and predicted word within a sentence.

Epochs (iter): Number of iterations (epochs) over the corpus.

To tune the model, selected are the ranging values of the parameters.

```
dm = [0, 1]
vector_size = [200, 300, 400, 500, 600, 700]
window = [3, 5, 10]
epochs = [10, 20, 30]
```

Above are the values of the hyperparameters. All the 108 combinations of the parameters are used to train the model and tested against the evaluation functions.

Above is the pseudocode for the evaluation function in which for one document top n similar documents are found. If the categories are the same as a document and top n fetched documents, then the ith ranked document is given (n-i) score.

5.5.2 Hyperparameters tuning for the resume dataset.

The resume dataset is passed through all the various combinations of the hyperparameters and, given the score below, is the result of the experiment.

Algorithm 2: Evaluation Function Algorithm

```
procedure doc2vec( resume , resumeDataset )  
  resumeVec  $\leftarrow$  model.inferVector(resume['text'])  
  similaritiesList  $\leftarrow$  []  
  score  $\leftarrow$  0  
  for res in resumeDataset do  
    similarity  $\leftarrow$  cosine(resumeVec, resVec)  
    similaritiesList.append(similarity, res)  
  
  end for  
  for (index, res) in similaritiesList.sort(n, ascending =  
    False, By = similarity) do  
  
    if res['category'] = resume['category'] then  
      score  $\leftarrow$  score + (n - index)  
    end if  
  end for  
  return score
```

dm	vector size	window	epochs	score
0	500	5	10	13681
0	700	5	10	13575
1	600	3	30	13543
1	700	5	30	13304
1	200	5	20	13008
0	300	10	30	12727
1	200	10	10	12164

As shown in the above table, the evaluation function gives the highest score when the parameters are vector size = 500, dm = 0, window = 5, and epochs = 10, which gives the highest score of 13681.

5.5.3 Hyperparameters tuning for the job description dataset

The job description dataset is passed through all the various combinations of the hyperparameters and, given the score below, is the result of the experiment.

dm	vector size	window	epochs	score
0	500	5	10	18769
0	200	3	10	18617
1	700	3	30	16959
0	600	10	30	16765
1	200	5	20	16362
1	300	10	10	15411
1	200	10	30	14838

As shown in the above table, the evaluation function gives the highest score when the parameters are vector size = 500, dm = 0, window = 5, and epochs = 10, which gives the highest score of 18769.

5.6 Finding similar Job descriptions.

This section describes the method for finding similar job descriptions among the dataset of the Job description for each Job description. Moreover, it discusses the need to find similar job descriptions.

5.6.1 Need of Finding Similar Job descriptions

As discussed in the previous chapter the Explicit skills which are useful in recommending the job description to the candidate's resume. As Explicit skills not only contains the skills mentioned in the particular job description but also skills in similar job descriptions. Explicit skills are the more relevant skills required for the given job role of the job description for evaluating our model. We are using the Job role as a parameter. So while recommending the Job descriptions to the candidate's resume, this is a potential factor that affects the recommendation. The same is not true when recommending the candidate's resumes for a Job description.

5.6.2 Merging skills

To find the similar Job descriptions of a particular job description the, all the datasets trained on with Doc2Vec model with the optimal value of hyperparameters found in the section 5.5 value of vector size = 500, dm = 0, window = 5, and epochs = 10 and cosine similarity is used to find the similarity between the numerical vectors of the job descriptions then top 10 most similar job descriptions are used to merge the skills.

5.7 Recommendation Engine

This chapter discusses two-recommendation of the job description and candidate's resume. That is recommending the Job descriptions to the candidate's resume and recommending the candidate's resume to the job descriptions. Also describe about the computing the similarity score between the documents.

5.7.1 finding similarities

For computing the similarities between the job descriptions and candidates' resumes. To recommend and find similarity using the cosine similarity function is used on the trained doc2vec model. In another way to find the similarity between the skills extracted from the documents is used. And in another way, the combination of both methods is used.

The first is to use the Cosine similarity in this method to recommend the resume to the job description. The resume dataset is trained on the doc2vec model with the value of hyperparameters found in section 5.6. Later for an input job description, the numeric vector of the job description is inferred from the doc2vec, and similarities with all resumes are found. Later, top 10 or top 20 found resumes are recommended to the job descriptions. To recommend the job descriptions to the candidate's resume, the job description dataset is trained on the doc2vec model with the value of hyperparameters found in the section 5.6 later for an input resume numeric vector of the resume is inferred from the doc2vec and similarity with the all job descriptions are found. Later, top 10 or top 20 found job descriptions are recommended to the candidate's resume.

The second is to use the skill matching the acquired skills from the candidate's resume to the required skills of the job descriptions to recommend the job descriptions to the candidate's resume. In the same way we recommend the candidate's resume to the job description.

The third approach is to use both the cosine similarity and the skills matching. That is, rather than using the cosine similarity and skills matching independently, we combined both approach to get better results. For combining both, the scoring method below the formula is used.

$$similarityscore = \frac{\alpha * (cosine \ similarity) + \beta * (skill \ Matching)}{\alpha + \beta} \quad (5.1)$$

For the above formula both different value of the alpha and beta is used. To find the optimal

values of both. For that for each combinations alpha and beta for each documents in dataset it is calculated against the every other documents in dataset and for top 20 or top 10 category is matched with the particular document's category is the both are same the score is given.

5.7.2 Recommending the Resumes to the Job Description

For the various combinations of the alpha and beta in formula number the approach described in sections 5.7.1 is used.

	cosine simi- larity	skills match- ing score	$\alpha=9, \beta=3$	$\alpha=7, \beta=1$	$\alpha=3, \beta=1$	$\alpha=5, \beta=1$	$\alpha=9, \beta=5$	$\alpha=5, \beta=3$
top 20	3707	1112	3846	3767	3846	3544	2886	2808
top 10	984	118	873	987	873	1028	658	622

Above is the table for the same. It gives better result on the weighted value of $\alpha=3$ and $\beta=1$.

5.7.3 Recommending the job descriptions to the resume

For the various combinations of the α and β in formula number the approach described in sections 5.7.1 is used.

	cosine simi- larity	skills match- ing score	$\alpha=9, \beta=3$	$\alpha=3, \beta=1$	$\alpha=5, \beta=3$	$\alpha=7, \beta=5$	$\alpha=1, \beta=3$	$\alpha=1, \beta=5$
top 20	3437	1778	3571	3571	3490	3422	2401	2161
top 10	1117	441	1259	1259	1148	1052	637	573

Above is the table for the same. It gives better result on the weighted value of $\alpha=3$ and $\beta=1$.

5.8 Classification of the job description

In the above discussed the methods we calculated similarities with the all the available documents in the dataset. Which is time inefficient. To overcome this we used classification techniques on the job description for that various classifier such as logistic regression, svm linear kernel and Naive bayes.

5.8.1 justification

To use the classification in the classified data rather than computing the similarity all over the dataset we only find similarity with documents of the related class. For that we trained the doc2vec model on the job description dataset. To embed the model optimal Hyperparameters are found using the tuning method. With the naive bayes classifier.

For evaluating the performance of Doc2Vec models on our dataset Accuracy and F1 Score have been utilised as the evaluation metrics First, Vector Size is varied, keeping the window fixed at 10 and epochs set to 30. As shown in Figure 5.3 and Figure 5.4, the highest accuracy of 0.7386 and an F1 - Score of 0.7246 is obtained for Vector Size = 325; thus, it is selected.

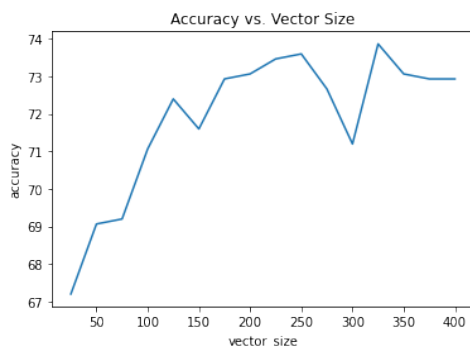


Figure 5.4: Accuracy vs. Vector Size

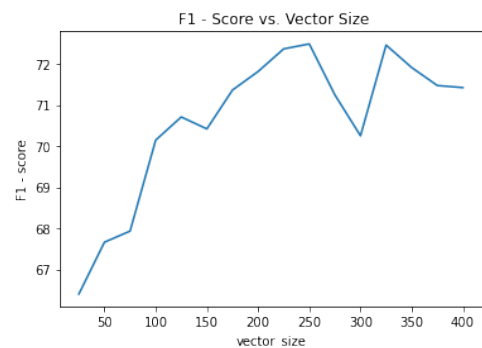


Figure 5.5: F1 - Score vs. Vector Size

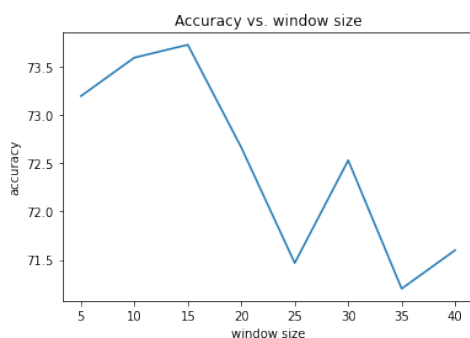


Figure 5.6: Accuracy vs. Window

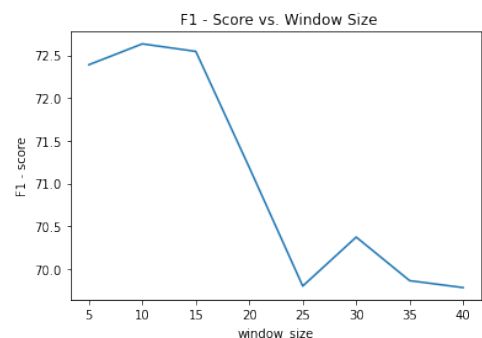


Figure 5.7: F1 - Score vs. Window

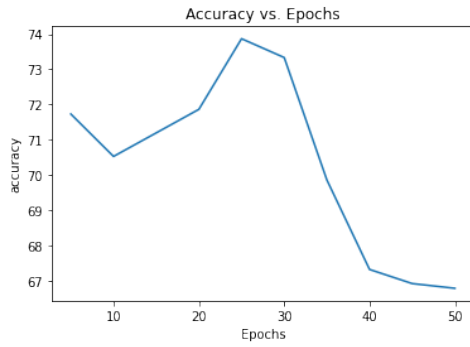


Figure 5.8: Accuracy vs. Epochs

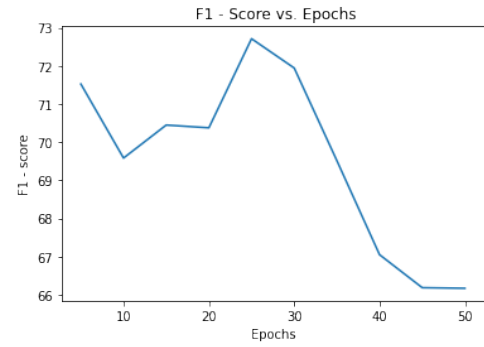


Figure 5.9: F1 - Score vs. Epochs

Now, the window values are varied, keeping the Vector Size fixed at 325 and epochs set to 30. As it can be seen from Figure 5.5 and Figure 5.6, a Window = 15 gives the highest accuracy of 0.7373 and an F1 - score of 0.7254; thus, it is selected.

Finally, epochs were varied, keeping the Vector Size fixed at 325 and Window fixed at 15. It is evident from Figure 5.7 and Figure 5.8 that epochs = 25 gives the highest accuracy of 0.7386 and an F1 score of 0.7271; thus, it is selected. Hence, from the various experiments performed, Vector Size = 325, Window = 15 and epochs = 25 are obtained as the best optimal hyperparameters for training the Doc2Vec model.

For these hyperparameters value classification is done on the svm linear kernel and logistic regression.

Classification Method	Accuracy	F1 score
logistic regression	0.7266	0.7166
svm linear kernel	0.6626	0.6590
Naïve Bayes	0.7373	0.7386

Above is the table showing the accuracy of the model. As show in the above table Naïve bayes gives the better results.

As show in the classification models are having lower accuracy on the dataset so it result in the lower accuracy in the recommendation.

5.9 Conclusion

This chapter presents the implementation of the proposed methodology in chapter 4. According to that, the skills are extracted from the resume and job description dataset. Later the skills of the most relevant job descriptions are merged to create the explicit skills for each job description.

Later, the Doc2vec model is trained for the resume dataset and job description dataset, and tuning of the hyperparameters is done. Different 108 combinations of the Hyperparameters are used against the evaluated function. Value of the vector size = 500, dm = 0, window = 5, and epochs = 10. These parameters are the same for both the training model.

Later, the SimilaritySimilarity finding methods are discussed, one only using the cosine similarity function and the other using the skills matching. The proposed combination method is used to find the SimilaritySimilarity, which uses both cosine similarity and skills matching. The weight consists of 3 for cosine similarity and 1 for skills matching, which gives better results than the traditional one.

Chapter 6

Conclusion and Future Work

6.1 Conclusion

Talent acquisition is the most important task for the success of the company. In the current situation for a given job, thousands of job seeker apply which make it hard for the hiring team to go through each and every resume manually and check for the credibility of the applicant. Similar to the job seeker in the large market thousands of jobs are available which makes finding a suitable job difficult for the user. This project aims to solve this problem by making automation of the resume matching process by using the various techniques for the data extraction from the given text or description and finding similarities and experimenting with the similarity calculated as the combination of the cosine function and skills matching. That gives better results than the individual one.

6.2 Future works

The work proposed in this project can be extended by considering all the job types. Later Hadoop distributed files running each on a separate node to train for the better optimal parameters. Here the skills dictionary used is the static set. Later feedback mechanisms can be used to make the skill dataset dynamic. The system can be hosted on the cloud to provide the service.

References

- [1] Hemant Misra Akshay Gugnani. In: *Implicit Skills Extraction Using Document Embedding and Its Use in Job Recommendation*. 2007, pp. 169–169.
- [2] Suad A Alasadi and Wesam S Bhaya. “Review of data preprocessing techniques in data mining”. In: *Journal of Engineering and Applied Sciences* 12.16 (2017), pp. 4102–4107.
- [3] Shaha Alotaibi. “A survey of job recommender systems”. In: *International Journal of the Physical Sciences* 7 (July 2012). doi: 10.5897/IJPS12.482.
- [4] Joeran Beel et al. “Research-Paper Recommender Systems: A Literature Survey”. In: *Int. J. Digit. Libr.* 17.4 (Nov. 2016), pp. 305–338. issn: 1432-5012. doi: 10.1007/s00799-015-0156-0. url: <https://doi.org/10.1007/s00799-015-0156-0>.
- [5] P. Brusilovsky and D. H. Lee. “Fighting Information Overflow with Personalized Comprehensive Information Access: A Proactive Job Recommender”. In: *Autonomic and Autonomous Systems, International Conference on*. Los Alamitos, CA, USA: IEEE Computer Society, June 2007, p. 21. doi: 10.1109/CONIELECOMP.2007.76.
- [6] Google Developer. *Collaborative Filtering and Matrix Factorization*. url: <https://developers.google.com/machine-learning/recommendation>. (accessed: 23.09.2021).
- [7] Tim; Faerber Frank; Weitzel and Tobias Keim. “”An Automated Recommendation Approach to Selection in Personnel Recruitment”. In: 2003.
- [8] Thomas Hofmann and Jan Puzicha. “Latent Class Models for Collaborative Filtering”. In: *IJCAI’99*. Stockholm, Sweden: Morgan Kaufmann Publishers Inc., 1999, pp. 688–693.
- [9] Donghwa Kim et al. “Multi-co-training for document classification using various document representations: TF-IDF, LDA, and Doc2Vec”. In: *Information Sciences* 477 (2019), pp. 15–29. issn: 0020-0255. doi: <https://doi.org/10.1016/j.ins.2018.10.006>. url: <https://www.sciencedirect.com/science/article/pii/S0020025518308028>.
- [10] Wes McKinney et al. “pandas: a foundational Python library for data analysis and statistics”. In: *Python for high performance and scientific computing* 14.9 (2011), pp. 1–9.

- [11] Ioannis K. Paparrizos, Berkant Barla Cambazoglu, and Aristides Gionis. “Machine learned job recommendation”. In: *RecSys '11*. 2011.
- [12] Ashis Pradhan. “SUPPORT VECTOR MACHINE-A Survey”. In.
- [13] Irina Rish. “An Empirical Study of the Naïve Bayes Classifier”. In: *IJCAI 2001 Work Empir Methods Artif Intell* 3 (Jan. 2001).
- [14] Baptiste Rocca. *Introduction to recommender System*. url: <https://towardsdatascience.com/introduction-to-recommender-systems-6c66cf15ada>. (accessed: 23.09.2021).
- [15] J. Ben Schafer, Joseph Konstan, and John Riedl. “Recommender Systems in E-Commerce”. In: *Proceedings of the 1st ACM Conference on Electronic Commerce*. EC '99. Denver, Colorado, USA: Association for Computing Machinery, 1999, pp. 158–166. isbn: 1581131763. doi: 10.1145/336992.337035.
- [16] Tomoki Sekiguchi. “Person-Organization Fit and Person-Job Fit in Employee Selection: A Review of the Literature”. In: *Osaka Keidai Ronshu* 54 (Jan. 2004).
- [17] Pradeep Singh et al. “Recommender Systems: An Overview, Research Trends, and Future Directions”. In: *International Journal of Business and Systems Research* 15 (Jan. 2021), pp. 14–52.

Acknowledgement

We would like to express our deep gratitude to our project guide, Dr. Sankita J Patel, Associate Professor, Computer science Engineering Department, SVNIT Surat, for their valuable guidance, helpful feedback, and co-operation with a kind and encouraging at the initial stage. We would also like to thank Dr. Mukesh A. Zaveri, Professor, Computer science Engineering Department. We are also thankful to SVNIT Surat and its staff for providing this opportunity which helps us to gain sufficient knowledge to make our work successful. Special thanks and appreciation to our colleagues in development and people who have willingly helped us out with their abilities.