

Master Thesis

Recommendation System for Galaxy Tools and Workflows

(Find similar tools and predict next tools in workflows)

Anup Kumar

Examiners: Prof. Dr. Rolf Backofen
Prof. Dr. Wolfgang Hess
Adviser: Dr. Björn Grüning

University of Freiburg
Faculty of Engineering
Department of Computer Science
Chair for Bioinformatics

July, 2018

Thesis period

10.01.2018 – 09.07.2018

Examiners

Prof. Dr. Rolf Backofen and Prof. Dr. Wolfgang Hess

Adviser

Dr. Björn Grüning

Declaration

I hereby declare, that I am the sole author and composer of my thesis and that no other sources or learning aids, other than those listed, have been used. Furthermore, I declare that I have acknowledged the work of others by providing detailed references of said work.

I hereby also declare, that my Thesis has not been prepared for another examination or assignment, either wholly or excerpts thereof.

Place, Date

Signature

Acknowledgment

I would like to extend my sincere gratitude to all the people who encouraged and supported me to accomplish this work. I am grateful to my mentor Dr. Björn Grüning who entrusted me with the task of building a recommendation system for Galaxy. He facilitated this work by providing me with all the indispensable means. Being precise, his pragmatic suggestions concerning the Galaxy tools and workflows helped me discern them better and improve the overall quality of the work. His advice to create a visualizer for showing the similar tools worked wonders as it enabled me to find and rectify a few bugs which were tough to establish. For the next task, creating a separate visualizer for looking through the next predicted tools was conducive in all merits. I appreciate and thank Eric Rasche who extracted the workflows for me from the Galaxy Freiburg server. I offer thanks to Dr. Mehmet Tekman and Joachim Wolff for their expert feedback, insights and general advice. At length, I wish to thank all the other members of Freiburg Galaxy team for their continuous support and help.

Abstract

The study explores two concepts to devise a recommendation system for Galaxy. A recommendation system can apprise a Galaxy user of the latent relations that exist among the tools in terms of their functions and types. Exhibiting an array of next possible tools at each step of picking a tool while creating workflows can also be a meaningful addition to this system.

To find similarities among tools, we need to extract information about each tool from its attributes like name, description, input and output file types and help text. We take into account these attributes one by one and compute similarity matrices. We compute three similarity matrices, one each for input and output file, name and description and help text attributes. Each row in a similarity matrix holds similarity scores of one tool against all the other tools. These similarity scores depend on the similarity measures (jaccard index and cosine similarity) used to compute the score between a pair of tools. To combine these matrices, one simple solution is to compute an average. But, assigning equal importance weight to each matrix might be sub-optimal. To find an optimal combination, we use optimization to learn importance weights for the corresponding rows for each tool in the similarity matrices. To define a loss function for optimization, we use a true similarity value based on the similarity measures. The similarity scores are positive real numbers between 0 and 1. We take an array of 1.0 as the true value.

Next task analyzes workflows to predict a set of next tools at each stage of creating workflows. While creating workflows, it would be convenient to leaf through a set of next possible tools as a guide. It can assist the less experienced (Galaxy) users in creating workflows when they are unsure about which tools can further be connected. In addition, it can curtail the time taken in creating a workflow. To achieve that, we need to learn the connections among tools in order to be able to predict the next possible ones based on the previously connected tools. To preprocess the workflows to make them usable by downstream machine learning algorithms, we compute all the paths bridging the starting and end tools in all workflows. We follow a classification approach to predict the next tools and use LSTM (long short-term memory), a variant

of recurrent neural networks. It performs well for learning long range, sequential and time-dependent data (tools connections) [1, 2]. We report the accuracy as precision.

Zusammenfassung

Contents

1	Introduction	2
1.1	Galaxy	2
1.2	Galaxy tools	3
1.3	Motivation	4
2	Approach	6
2.1	Extract tools data	6
2.1.1	Select tools attributes	7
2.1.2	Clean tools data	7
2.2	Learn dense vector for a document	12
2.2.1	Latent semantic indexing	12
2.2.2	Paragraph vectors	15
2.3	Similarity measures	18
2.3.1	Cosine similarity	18
2.3.2	Jaccard index	19
2.4	Optimization	20
3	Experiments	25
3.1	Amount of help text	25
3.2	One and two sources instead of three	25
3.3	Similarity measures	25
3.4	Latent Semantic Analysis	25
3.5	Paragraph Vector	26
3.5.1	Distributed bag-of-words	26
3.6	Gradient Descent	26
3.6.1	Learning rates	27
4	Results and Analysis	28
4.1	Latent Semantic Analysis	28
4.1.1	Full-rank matrices	28

4.1.2	70% of full-rank	28
4.1.3	30% of full-rank	29
4.1.4	5% of full-rank	30
4.1.5	Improvement verification	30
4.2	Paragraph vectors	32
4.2.1	Optimization	32
4.2.2	Visualizer for paragraph vectors approach	32
5	Conclusion	42
5.1	Tools data	42
5.2	Approaches	42
5.3	Optimization	43
6	Future Work	44
6.1	Get true value	44
6.2	Correlation	44
	Bibliography	45

Part 1: Find similar Galaxy tools

1 Introduction

1.1 Galaxy

Galaxy¹ is an open-source biological data processing and research platform. It supports numerous types of extensively used biological data formats like FASTA, FASTAQ, GFF, PDB and many more. To process these datasets, Galaxy offers tools and workflows which either transform these datasets from one type to another or manipulate them. A simple example of data processing is to merge two compatible datasets to make one. Another example can be to reverse complement a sequence of nucleotides².

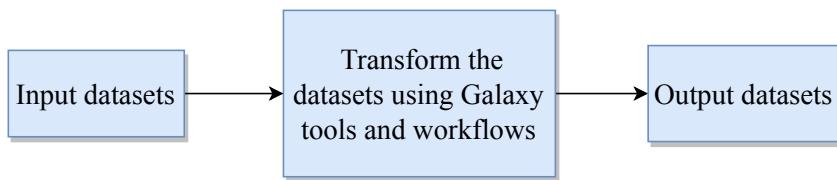


Figure 1: Dataset transformation: The image shows a general flow of dataset transformation using Galaxy tools and workflows.

A tool is a data-transforming entity which allows one or more types of datasets, transforms these datasets and produces output datasets. The tools are classified into multiple categories based on their functions. For example, the tools which manipulate text like replacing texts and selecting first lines of a dataset are grouped together under "Text Manipulation" group.

These tools form the building blocks of workflows. The workflows are data processing pipelines where a set of tools are joined one after another. The connected tools need to be compatible with each other which means the output types of one tool should be present in the input types of the next tool. A workflow can have one or more starting and end tools.

¹<https://usegalaxy.eu/>

²https://usegalaxy.eu/?tool_id=MAF_Reverse_Complement_1&version=1.0.1&__identifier=zmk9dx9ivbk

1.2 Galaxy tools

A tool entails a specific function. It consumes datasets, brings about some transformations and produces output datasets which can be fed to other tools. A tool has multiple attributes which include its input and output file types, name, description, help text and so on. They carry more information about a tool. When we look at the collective information about all these attributes for multiple tools, we find that some tools have similar functionalities based on their similarities in their corresponding attributes. For example, there are tools which share similarities in their respective functions and the input and output dataset types they are glued to. For example, a tool "hicexplorer hicpca"³ has an output type named "bigwig". Hence, if there is a tool or a set of tools which also has "bigwig" as their input and/or output type, we consider there could be some similarity between "hicexplorer hicpca" and the other tools as they do transformations on similar types of datasets. In addition, we can find similar functions of tools by analyzing their "name" and "description" attributes. Let's take an example of two tools (Figure 2):

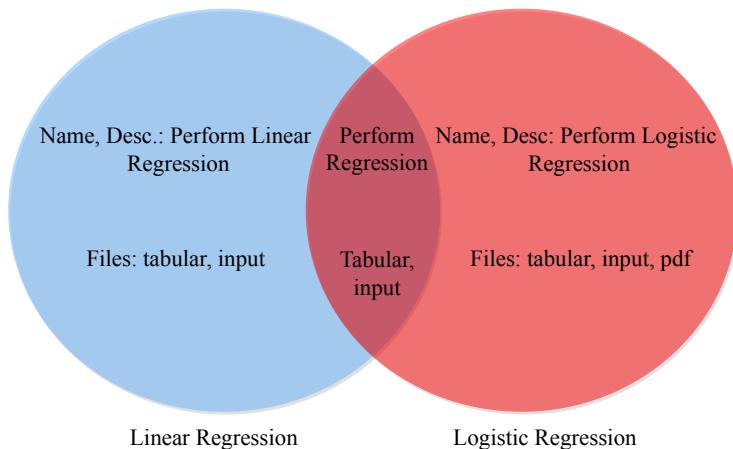


Figure 2: Common features: The venn diagram shows the common features for two tools (linear regression and logistic regression). Based on these common features, the extent of similarity between two tools are computed.

In figure 2, we take two tools - "Linear Regression" and "Logistic Regression" and collect their respective information from their input, output file types, name and description attributes. We see that these tools share features in the venn diagram. They both do regression and few file types are also common. In the same way, if we

³https://usegalaxy.eu/?tool_id=toolshed.g2.bx.psu.edu/repos/bgruening/hicexplorer_hicpca/hicexplorer_hicpca/2.1.0&version=2.1.0&__identifier=5kcqmvb71gx

extrapolate this venn diagram and match one tool against all other tools, we hope to find a set of tools similar in nature to the former tool. While searching for the related tools for a tool, it is possible that we end up with an empty set.

1.3 Motivation

From figure 2, we see that there can be tools which share characteristics. Galaxy has thousands of tools having a diverse set of functions. Moreover, new tools keep getting added to the older set of tools. From a user's perspective, it is hard to keep knowledge about so many tools. It is important to make a user aware of the presence of new tools added. If there is a model which dispenses a clue that there is a set of say n tools which are similar to a tool, it would give more options to a user for her/his data processing.

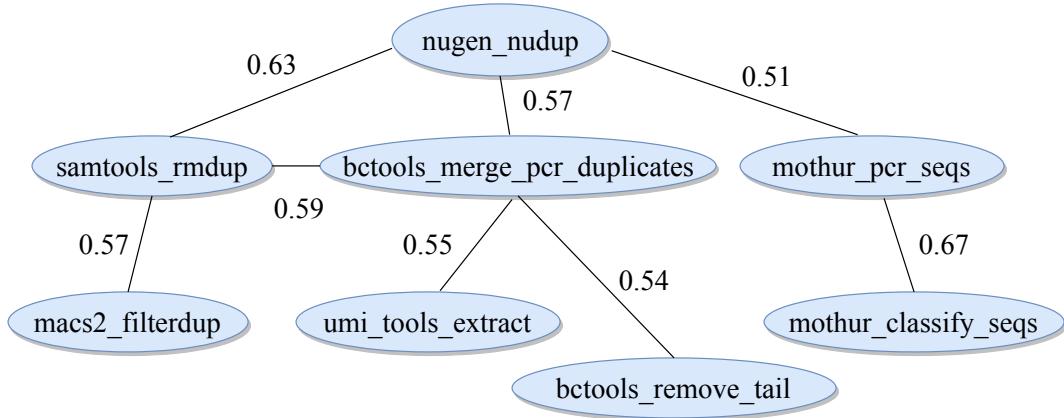


Figure 3: Similarity knowledge graph: The graph defines a network of relationship among the tools using the similarity scores between a pair of nodes (tools). The real numbers on an edge is the relation strength between two nodes.

To elaborate it more, let's take an example of a tool "nugen nudup"⁴. It is used to find and remove PCR duplicates. The similar tools for it can be "samtools rmdup" and "bctools merge pcr duplicates" which work on related concepts. These similar tools would further have their respective set of similar tools thereby making a network of related entities (tools). This "knowledge network" can help a user find multiple ways to process her/his data and exhibits "connectedness" among tools. The strength of this relation may vary from being small to large. To ascertain that, this study

⁴https://toolshed.g2.bx.psu.edu/repository?repository_id=4f614394b93677e3

learns a continuous representation of the relation strength. Figure 2 shows how this knowledge graph can evolve. First, we find similar tools for "nugen nudup" and connect them to their source tool specifying the similarity values as real numbers at the edges. These similar tools further have their own sets of similar tools and so on.

2 Approach

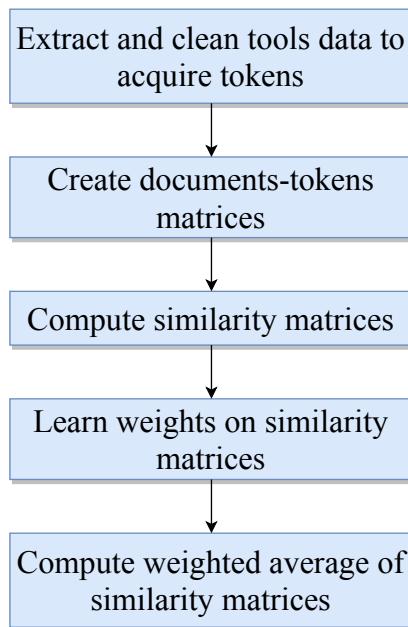


Figure 4: Sequence of steps to find similar tools: The image shows a series of steps used to establish similarity among tools using approaches from natural language processing to compute text similarity scores and optimization to combine these similarity scores.

2.1 Extract tools data

There are multiple repositories of Galaxy tools stored at GitHub¹. In each of the tool repository, there are *xml* files starting with a *< tool >* tag. We read all of these *xml* files, extract information from a few of the attributes and collect them in a *tabular* file.

¹One example:<https://github.com/galaxyproject/tools-iuc/tree/master/tools>

2.1.1 Select tools attributes

A tool has multiple attributes like input and output file types, help text, name, description, citations and some more. But all of these attributes are not important and do not generally identify a tool exclusively. We consider some of these attributes:

- Input and output file types
- Name and description
- Help text

Moreover, we combine the input and output file types and name and description respectively as they are of similar nature. These two combined attributes give complete information about a tool file types and its functionality. We also consider help text attribute which is larger in size compared to the previous two. At the same time, they are empty for few tools. Apart from being large in size, this attribute is noisy as well. It provides more information about the usage of a tool. Generally, in the first few lines, it gives a detailed explanation of the tool functions. Further, it explains how the input data should be supplied to a tool or how an input data looks like. Hence, much of the information present in this attribute is not important. Because of noise present in this attribute, we decide to use only upto first 4 lines which illustrates the core functionality of the tool. The decision to select only first 4 lines is empirical. The rest of the information in help text is discarded.

2.1.2 Clean tools data

Remove duplicates and stopwords

The collected data for tools is raw containing lots of commonplace and duplicate items which do not add value. These items should be removed to get *tokens* which are unique and useful. For example, a tool *bamleftalign* has input files as *bam*, *fasta* and output file as *bam*. While combining the file types, we discard the repeated file types and in this case, we consider file types as *bam*, *fasta*. The other attributes we deal with are different from the file types. The files types are discrete items but in attributes like name and description and help text, the account is in a human language. The explanation contains complete or partially complete sentences in *English*. Hence, to process this information, we need startegies that are prevalent in natural language processing². The sentences we write in *English* contain many words

²<https://www.ncbi.nlm.nih.gov/pmc/articles/PMC3168328/>

and has different parts. These parts include subject, object, preposition, interjection, verbs, adjectives, adverbs, articles and many others. For our processing, we need only those tokens (words) which categorize a tool uniquely and do away with multiple parts of speech present in the statements. For example, a tool named *tophat* has name and description as "TopHat for Illumina Find splice junctions using RNA-seq data". The words like *for*, *using* and *data* do not give much value as they must be present for many tools. These words are called as "stop words"³ and we selectively discard them. In addition, we remove numbers and convert all the tokens to lower case.

Use stemming

After removing duplicates and stop words, our data is clean and contain tokens which uniquely identify corresponding tools. When we frame sentences, we follow grammar which constrains us to use different forms of the same word in varying contexts. For example, a word *regress* can be used in multiple forms as *regresses* or *regression* or *regressed*. They share the same root and point towards the same concept. If many tools use this word in varying forms, it is beneficial to converge all the different forms of a word to one basic form. This is called stemming⁴. We use nltk⁵ package for stemming.

Learn relevance for words

We use a term "token" for each word. For example, a tool's name contains "regress, perform" as a set of tokens (words). After discarding duplicate tokens and stopwords and using root words, we have a set of good tokens for all the three attributes - input and output file types, name and description and help text. Let's call these sets as *documents*. The tokens present in these documents do not carry equal importance. Some tokens are more relevant to the document and some not so relevant . We need to find out importance factor for all tokens in a document. Using these factors, we can arrange them in big, sparse documents-tokens matrix. In this matrix, each row represents a document and each column belongs to one token. To compute these relevance scores, we use bestmatch25. Let's associate some variables to be used in explaining this algorithm.

³<https://www.ranks.nl/stopwords>

⁴<https://nlp.stanford.edu/IR-book/html/htmledition/stemming-and-lemmatization-1.html>

⁵<http://www.nltk.org/>

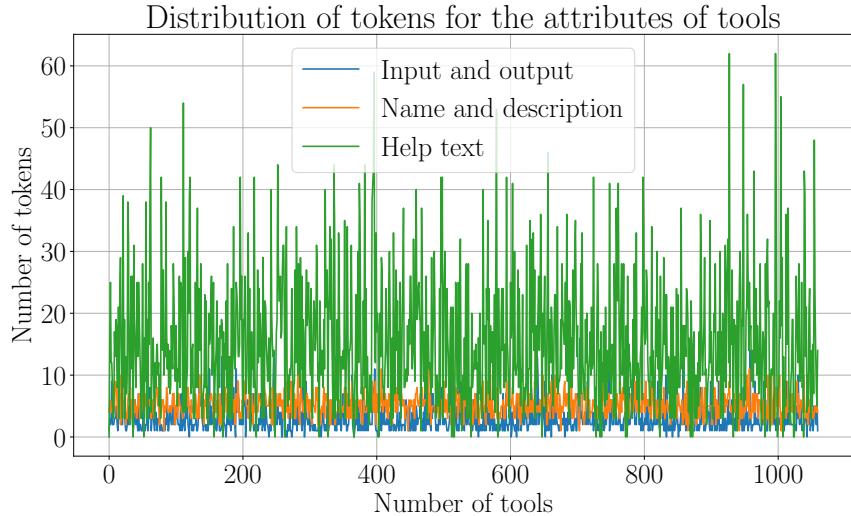


Figure 5: Distribution of tokens (words): The plot shows a distribution of tokens for input and output file types, name and description and help text attributes of tools. The help text attribute contains more number of tokens compared to the other two. The input and output file types attribute contains lower number of tokens compared to the other two attributes.

- Token frequency ⁶ tf
- Inverted document frequency idf
- Average document length $|D|_{avg}$
- Number of documents N
- Size of a document $|D|$

First of all, token frequency (tf) specifies the count of a token's occurrence in a document. If a token *regress* appears twice in a document, its tf is 2. This can also be understood as a weight given to this term. Inverted document frequency for a token is defined as:

$$idf = \log \frac{N}{df} \quad (1)$$

where df is the count of the documents in which this token is present and N is the total number of documents. If we randomly sample a document, then the probability

⁶<https://nlp.stanford.edu/IR-book/pdf/06vect.pdf>

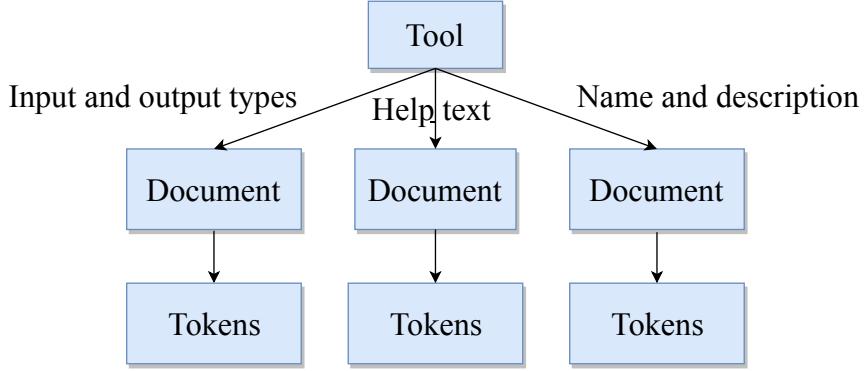


Figure 6: Relationship between a tool, its documents and their tokens:

The image shows that a tool has three documents corresponding to each attribute and each document contains tokens. For all the tools, we have documents equal to the number of tools for each attribute. The number of tokens in each document varies. Minimum number of tokens for any document can be 0.

of this token to be present in this document is $p_i = \frac{df}{N}$. From information theory, we can say that the information contained by this event is $-\log p_i$. The entity idf is higher when a token appears less number of documents which means that this token is a good candidate for representing the document and possesses higher power to distinguish between documents. The tokens which appear in many documents are not good representatives. Average document length is the average number tokens for all the documents. Size of a document is the count of all the tokens for that document [3].

$$\alpha = (1 - b) + \frac{b \cdot |D|}{|D|_{avg}} \quad (2)$$

$$tf^* = tf \cdot \frac{k + 1}{k \cdot \alpha + tf} \quad (3)$$

$$BM25_{score} = tf^* \cdot idf \quad (4)$$

where k and b are hyperparameters. Using the equation 4, we compute the relevance score for each token in all the documents. Table 1 shows some sample scores for a few documents where the tokens are present with their respective relevance scores. In this way, we arrange document-tokens matrix for all the attributes of tools. For input and output file types, these matrix entries will have only two value, 1 if a token

is present for a document and 0 if not. For other attributes, relevance scores are positive real numbers. This strategy of representing documents with their tokens is called vector space model as each document represents a vector of tokens.

Documents/tokens	regress	linear	gap	mapper	perform
LinearRegression	5.22	4.1	0.0	0.0	3.84
LogisticRegression	3.54	0.0	0.0	0.0	2.61
Tophat2	0.0	0.0	1.2	1.47	0.0
Hisat	0.0	0.0	0.0	0.0	0.0

Table 1: Documents-tokens matrix: This table shows a matrix of tools (documents) arranged along the rows and tokens along the columns. Each value in the matrix is a weight (relevance-factor) assigned to a token for a document. This matrix is sparse containing mostly zeros as the number of tokens is significantly large compared to the number of tokens present in a document. This table shows a sample of how actual documents-tokens matrix would look like.

Figure 4 shows the heatmaps for documents-tokens matrices that belong to input and output file types and name and description. We can see that these plots are sparse. Each entry in these matrices contain BM25 score for each token in every document. The representation shows how to find tokens which are good representatives of documents with a weighted by their relevance factors. But, they do not tell us anything about the co-occurrence of a few tokens in a document. It tells us that a token is important for a document if the BM25 score is higher but it does not tell us anything about its relation to other tokens. Due to this shortcoming, it does not acknowledge the presence of "concepts" or "context" hidden in a document. A concept in document can be realised when we see the relation among a few words. To illustrate this idea, let's take an example of three words - "New York City". These three words mean little or point to different things if looked at separately. But, if we see them together, it points towards a concept. This vector space model lacks the ability to find the correlation among tokens. To enable the vector space model to learn this hidden concepts and find correlation among multiple tokens, we explore two ideas

- Latent Semantic Indexing/Analysis ⁷
- Paragraph Vectors

⁷<http://lsa.colorado.edu/papers/dp1.LSAintro.pdf>

Using these approaches, we learn dense, n dimensional vector for each document instead of using sparse vectors as shown in figure 5.

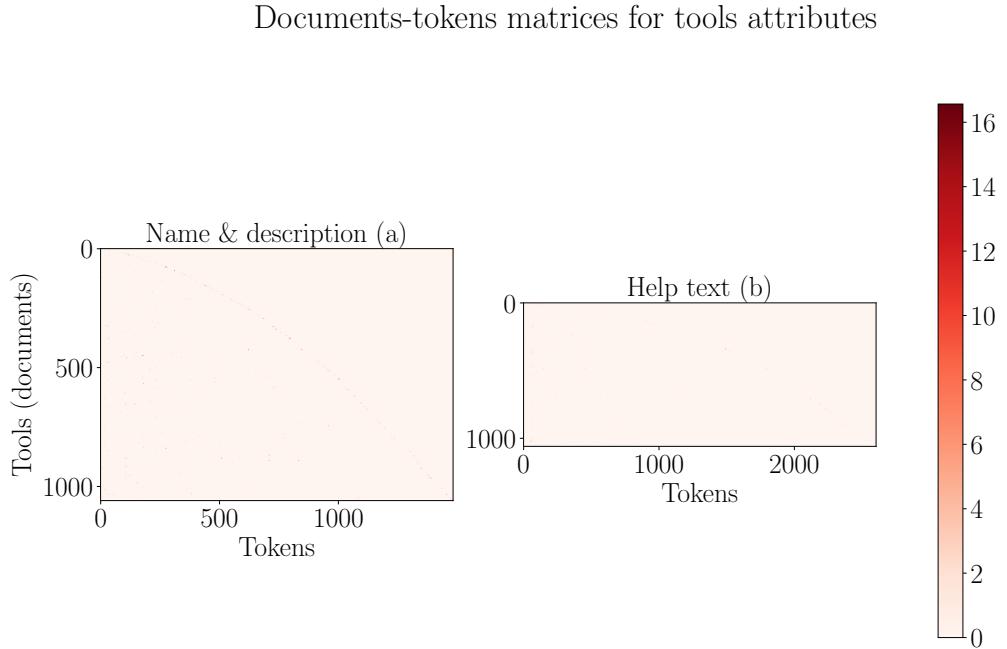


Figure 7: Heatmap for documents-tokens matrices: The plot shows a heatmap of documents-tokens matrices for name and description and help text attributes. We see that the matrices are sparse containing only few darker spots. The help text matrix is more sparse than name and description matrix as the former contains more number of tokens. We exclude the documents-tokens matrix of input and output file types because we use it as such. For the other two matrices, we would estimate their dense, lower-rank approximations.

2.2 Learn dense vector for a document

2.2.1 Latent semantic indexing

It is a mathematical way to learn the hidden (latent) concepts in documents by computing a low-rank representation of a documents-tokens matrix [4, 5, 6]. This low-rank matrix is dense (figure 6). We use singular value decomposition (*SVD*) to decompose the full-rank matrix into a significantly lower rank matrix. The optimal

rank to which a matrix to be decomposed is empirical in nature. We would choose ranks from higher to lower value and consequently the sum of singular values would also decrease with the rank. This decomposition follows the equation:

$$X_{n \times m} = U_{n \times n} \cdot S_{n \times m} \cdot V^T_{m \times m} \quad (5)$$

where n is the number of documents and m is the number of tokens. S is a diagonal matrix containing the singular values in descending order. It contains the weights of the concepts present in the matrix. The matrices U and V are orthogonal matrices which satisfy:

$$U^T \cdot U = I_{n \times n} \quad (6)$$

$$V^T \cdot V = I_{m \times m} \quad (7)$$

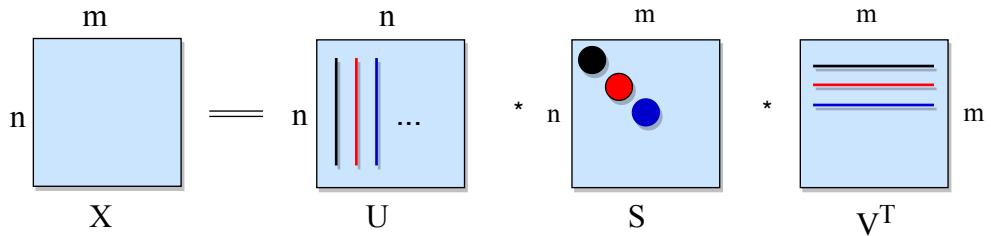


Figure 8: Singular value decomposition: The image shows that how a matrix is decomposed using singular value decomposition. A matrix X is decomposed into three factor matrices, U , S and V . The k most important dimensions are kept and the rest are discarded. The value of k is empirical and can be estimated using optimization taking frobenius norm as an error function.

Figure 6 explains ⁸ how SVD of a matrix is carried out. The matrix U contains information about how the tokens, arranged along the columns, are mapped to concepts and matrix V stores information about how the concepts are mapped to documents arranged along the rows.

Low-rank approximation

The low-rank approximation of a matrix is important to discard the features which are non-repeating or noise. Using this, we can collect the latent relations present in

⁸<http://theory.stanford.edu/~tim/s15/l19.pdf>

the documents-tokens matrices. We saw that our documents-tokens matrices suffer from sparseness and exhibit no relation among tokens. The approximation deals with these issues as well. The resulting matrices are dense and contain most of the singular values. The singular values which are small (the last entries of the S matrix along the diagonal) are discarded [7]. The low-rank approximated matrix X_k is computed as:

$$X_{n \times m} = U_k \cdot S_k \cdot V_k^T \quad (8)$$

where U_k is the first k columns of U , V_k is the first k rows and S is the first k singular values. k is an empirical parameter. X_k is called as the rank- k approximation of the full rank matrix X . Figure 7 shows the variation of the sum of singular values with the percentage rank of matrices. The percentage rank is $k \div K$ where $1 \leq k \leq K$ and K is the original (full) rank of a matrix. We take it as the original ranks of the three matrices are not comparable. By doing this, we can show them on one plot. Similar idea we use for the sum of the singular values of each matrix. From figure 7, we can say that if we reduce the ranks of matrices to 70% of the full-rank, we can still capture $\approx 90\%$ of the singular values. The reduction to half of the full-rank achieves $\approx 80\%$ of the sum of singular values.

We reduce the rank of the original documents-tokens matrices and compute the dense and approximated low-rank matrices. Figure 8 shows the low-rank matrices for input and output file types and name and description attributes. For computing this, we use only 40% of the full-rank. We can compare it with figure 5 and deduce that figure 8 is denser than figure 5. In these low-rank matrices, we have dense vector representations for documents along the rows. In each matrix, each row contains a vector for one document. Using these documents vectors, we can compute the correlation or similarity using any similarity measure. There are multiple similarity measures that can be used like euclidean distance, cosine similarity, manhattan distance. We use cosine angle similarity to compute the correlation between vectors. By this, we get a positive real number between 0.0 and 1.0 specifying how similar a pair of documents are. The higher the score, higher is the similarity. Computing this similarity for all the documents give us a similarity matrix $S_{n \times n}$ where n is the number of documents. This square, symmetric matrix is called as similarity or correlation matrix. We compute three such matrices each corresponding to one attribute.

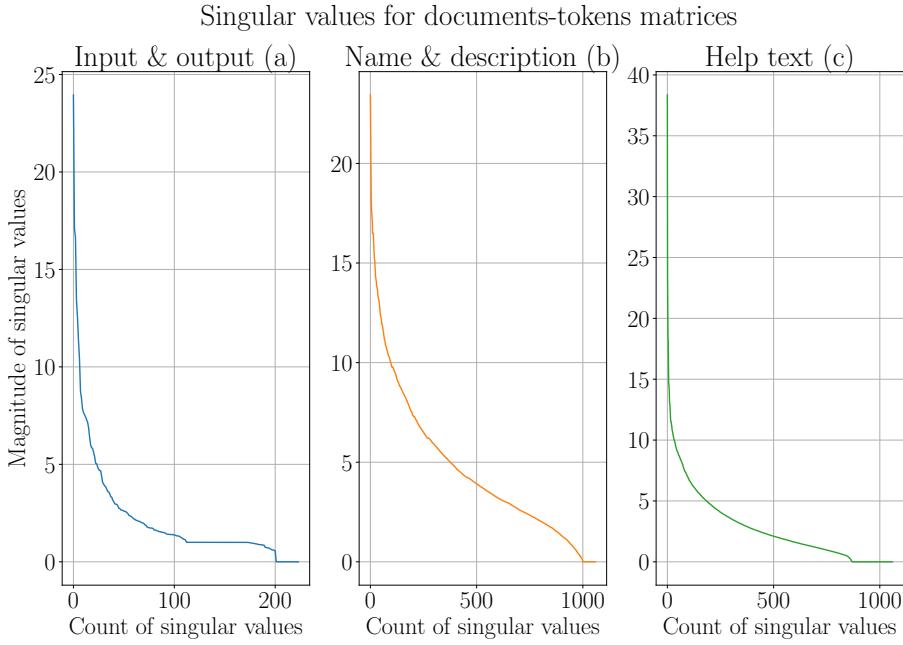


Figure 9: Singular values of the documents-tokens matrices: The plot shows singular values computed using singular value decomposition (equation 5). The diagonal matrix S contains these singular values sorted in descending order. We can see that in (a), (b) and (c) that very few singular values have higher magnitude and most of the singular values are smaller.

2.2.2 Paragraph vectors

Using latent semantic indexing (LSI), we learnt dense vectors to represent each document. It learns better vector representations for documents compared to using full-rank documents-tokens matrices. One main limitation is to assess the quantity by which we need to lower the rank of a matrix in order to find the optimal results. There are ways to find the optimal reduced rank optimizing the frobenius norm but it is not simple. We would see in the analysis section that similar tools are more dominated by the scores shared by the input and output file types due to which the tools which similar in their functions do not come up. In order to avoid these limitations, we use an approach known as *doc2vec* (document to vector). It learns a dense, fixed-size vectors for each document using neural networks. These vectors are unique in a way that captures the semantics present in the documents. The documents which share similar context are represented by similar vectors. When we compute cosine distance between these vectors sharing similar context, we get a

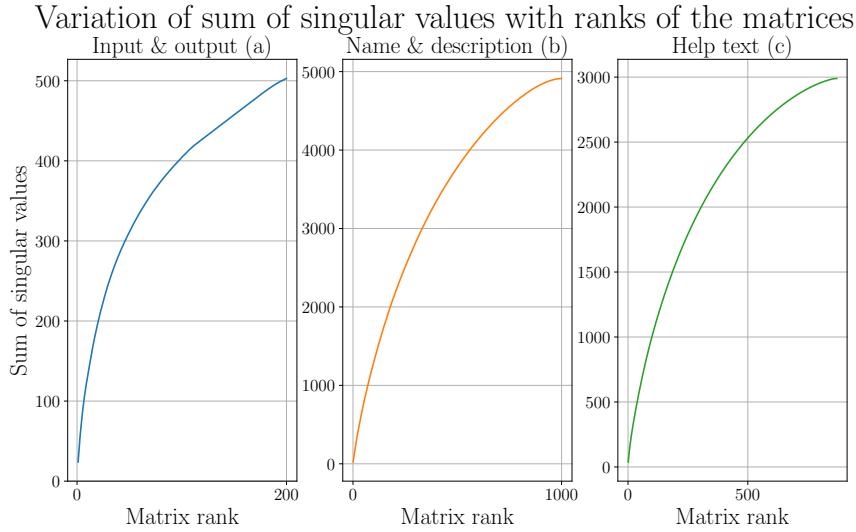


Figure 10: Sum of singular values with matrix rank: The plot shows an easier way to see how the sum of singular values varies with a documents-tokens matrix rank for the three attributes. Here the (a), (b) and (c) show separately this variation as the ranks of these matrices and sum of singular values differ.

higher number (close to 1.0). It allows the documents to have variable lengths.

Approach

Paragraph vectors approach learns vector representations for all the words present in a set of documents. The words which are used in a similar context have similar vectors. For example, words like "data" and "dataset" which are used, in general, in similar context have are represented by close vectors. The vector representations of words in a corpus is learnt by finding the maximum of the following equation:

$$\frac{1}{T} \cdot \sum_{t=k}^{T-k} \log p(w_t | w_{t-k}, \dots, w_{t+k}) \quad (9)$$

where T is the total number of words in a corpus, k is the window size. We take a few words which make a context and using this context we try to predict each word. The probability p is computed using a softmax classifier and backpropagation is used to compute the gradient and the vectors are optimized using stochastic gradient descent. To learn paragraph vectors, in addition to using words vectors, paragraph vectors are also used to learn probability of the next words in a context.

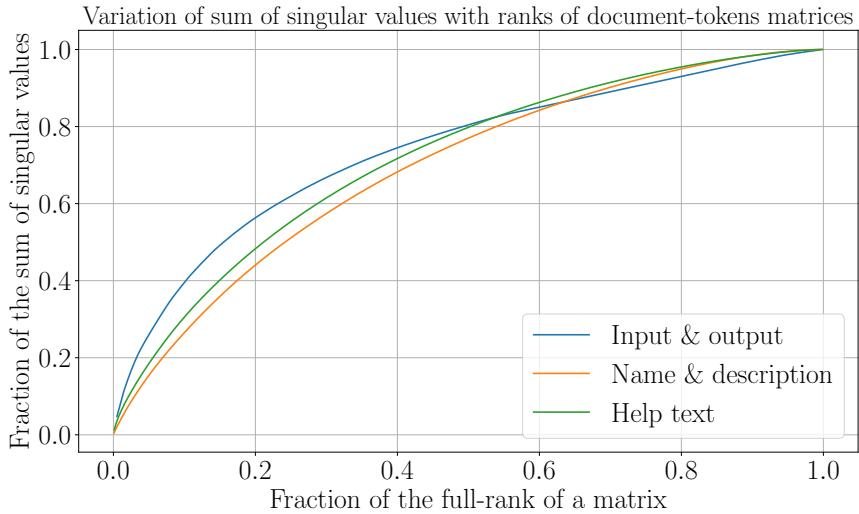


Figure 11: Matrix rank and singular values: This plot merges the results of the figure 10 into one plot. As the ranks of document-term matrices and sum of singular values vary, we convert them to respective percentages. $rank_{fraction} = \frac{k}{N}$ where k is the reduced rank and N is the full-rank of a matrix. For example 0.2 on the rank axis (x-axis) means 20% of the original rank of a matrix. Similarly, y-axis shows the fraction of the sum of all singular values $sum_{fraction} = \frac{\sum_{i=1}^k}{\sum_{i=1}^K}$ where K is the number of all singular values.

The paragraph and word vectors are averaged to make the classifier which predict next words in a context. There are two way how to choose the context.

- Distributed memory: In this approach of learning paragraph vectors, fixed length window of words are chosen and paragraph vectors and word vectors are used to predict the words in this context. The word vectors are shared across all paragraphs and paragraph vector is unique to each paragraph.
- Distributed bag of words: The words are randomly extracted from paragraphs and in this set of words, a word is chosen randomly and predicted using the paragraph vectors. The order of the randomly chosen words is ignored.

The figures 9 and 10 are inspired from the original work - Distributed Representations of Sentences and Documents ⁹. The second form of learning paragraph vectors is simple and we use it to learn documents (paragraphs) vectors for name

⁹<https://arxiv.org/abs/1405.4053>

Low rank representation of documents-tokens matrices

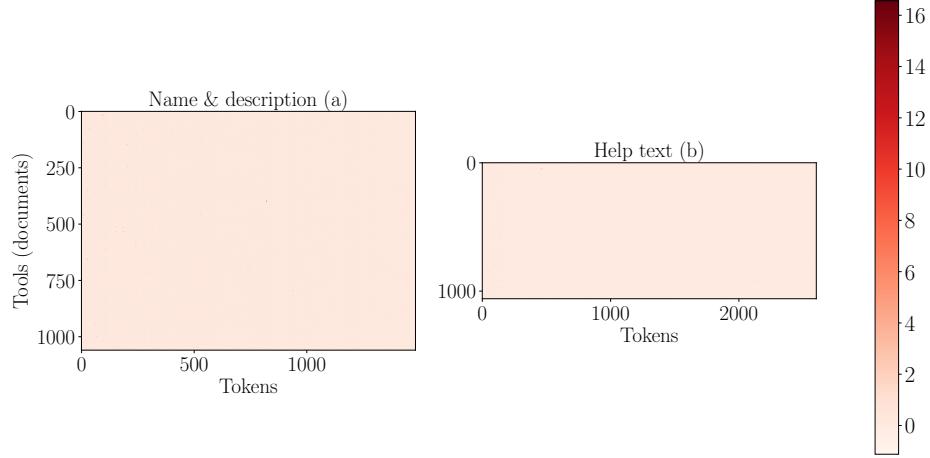


Figure 12: Documents-tokens matrices: The heatmap shows the low-rank (5% of full-rank) representations of documents-tokens matrices for name and description and help text. These matrices are more dense compared to figure 7 which shows these matrices corresponding to the full-rank representations.

and description and help text attributes. We learn only the paragraph vectors which makes it less computationally expensive [8].

2.3 Similarity measures

We have vectors representing the documents for all the three attributes we consider to compute similar tools. To find the similarity in a pair of vectors, we need to apply a similarity measure to get a similarity score. This score quantifies how much similar a pair of documents are which means how close or distant these documents are.

2.3.1 Cosine similarity

It gives the cosine angle value between a pair of documents vectors. Let's say we have two vectors, x and y . We can write:

$$x \cdot y = |x| \cdot |y| \cdot \cos \theta \quad (10)$$

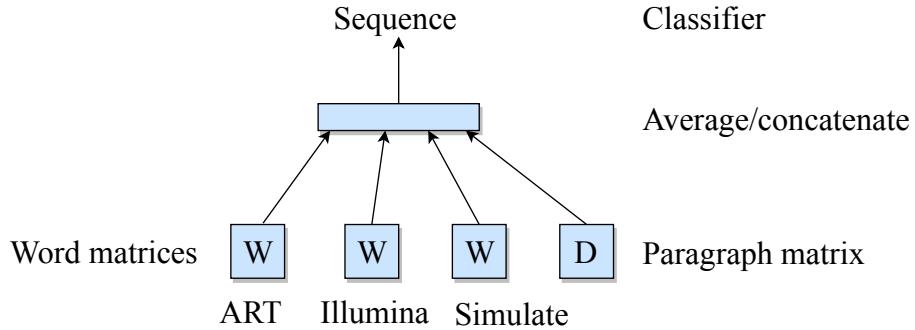


Figure 13: Distributed memory model for paragraph vectors: This image shows the mechanism for learning a vector for a paragraph. W is a word matrix where each word is represented by a vector. D is a paragraph matrix where each paragraph (document) is represented by a vector. The word vectors are shared across all paragraphs (documents) but not the paragraph vectors. The three words "Art", "Illumina" and "Simulate" represent a context. The average or concatenated words and paragraph vectors are used to predict the "Sequence" word.

$$\cos \theta = \frac{x \cdot y}{|x| \cdot |y|} \quad (11)$$

where $|.|$ is the norm of the vector x . If the norm is 0, we use 0 for the value of $\cos \theta$. The values emitted by this similarity follows a natural progression that if the documents are dissimilar, then it is 0 and if completely similar, it is 1.0. Otherwise it lies between 0.0 and 1.0. This score can also be understood as a kind of probability of similarity between a pair of documents ¹⁰.

2.3.2 Jaccard index

Jaccard index is a measure of similarity between two sets of entities given by the equation:

$$j = \frac{A \cap B}{A \cup B} \quad (12)$$

where A and B are two sets. \cap is the number of entities present in both the sets and \cup is the sum of unique entities present in sets A and B [9]. We use this measure to compute the similarity between two tools based on their file types. For example, *LinearRegression* has 3 file types: *tabular*, *input* and *pdf*. Another tool

¹⁰<https://nlp.stanford.edu/IR-book/html/htmledition/dot-products-1.html>

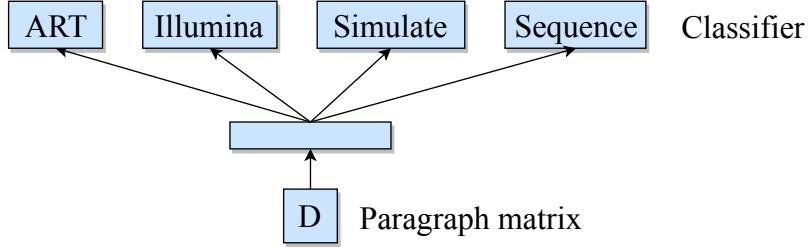


Figure 14: Distributed bag of words model for paragraph vectors: This image shows how only paragraph vectors are learned by predicting a random word chosen from a randomly selected set of words. In this approach, the order of the words does not matter.

LDAAnalysis has *tabular* and *txt* as file types. The jaccard index for this pair of tools would be:

$$j = \frac{\text{Length}[(\text{tabular}, \text{input}, \text{pdf}) \cap (\text{tabular}, \text{txt})]}{\text{Length}[(\text{tabular}, \text{input}, \text{pdf}) \cup (\text{tabular}, \text{txt})]} \quad (13)$$

$$j = \frac{1}{4} = 0.25 \quad (14)$$

2.4 Optimization

Using the paragraph vectors and applying cosine similarity, we get similarity or correlation matrices one each for input and output file types, name and description and help text attributes. The matrices have same dimensions ($n \times n$, n is the number of tools). To combine these matrices, one simple idea would be to take an average of the corresponding rows of scores from the matrices for a document. This would generate a matrix of the same dimension where each row would correspond to one document. This row would contain similarity scores of a document against all other documents. The diagonal entries of this matrix would be 1.0. All other entries would be a positive real number between 0.0 and 1.0. Another way to find the combination is to learn the weights on the rows from three matrices and then combine them to obtain optimal similarity scores for a tool. The weight would be a positive real number between 0.0 and 1.0. Instead of using a fixed importance factor (weights) of $1/3$ (3 is the number of matrices), we use optimization to find these real numbers and then combine the matrices by multiplying with these weights.

$$S_k^{optimal} = w_{io}^k \cdot S_{io}^k + w_{nd}^k \cdot S_{nd}^k + w_{ht}^k \cdot S_{ht}^k \quad (15)$$

where w is a positive, real number and satisfy $w_{io}^k + w_{nd}^k + w_{ht}^k = 1$. S_{io}^k , S_{nd}^k and S_{ht}^k are the similarity scores (corresponding matrix rows) for k^{th} tool corresponding to input and output file types, name and description and help text attributes respectively.

We have these similarity scores. But we need to learn these importance weights. We use gradient descent optimizer to learn these weights against a loss function. If we look closely the similarity measure, we see that the maximum similarity between a pair of documents can be 1.0. Hence, the ideal similarity scores for one document against all other documents:

$$S_{ideal} = [1.0, 1.0, \dots, 1.0] \quad (16)$$

where S_{ideal} is the ideal similarity scores for one document against all other document. Using this ideal scores, we can compute the error we accrue for all the similarity scores from three attributes. After computing the error, we can verify which attribute is closer to the ideal score and which are not. The attributes which exhibit lower error get higher weights and those which score higher error get lower weights. The next section elaborates the online way to do the optimization.

Gradient descent

Gradient descent is a popular algorithm for optimizing an objective function with respect to its parameters. The parameters are the entities which we want to learn. In our case, these are the weights. The algorithm minimizes an error function. Let's frame the error function:

$$Error_{io}^k(w_{io}) = \sum_{t=1}^{N-1} \cdot (w_{io} \cdot S_{io}^t - S_{ideal})^2 \quad (17)$$

$$Error_{nd}^k(w_{nd}) = \sum_{t=1}^{N-1} \cdot (w_{nd} \cdot S_{nd}^t - S_{ideal})^2 \quad (18)$$

$$Error_{ht}^k(w_{ht}) = \sum_{t=1}^{N-1} \cdot (w_{ht} \cdot S_{ht}^t - S_{ideal})^2 \quad (19)$$

$$Error^k(w) = \sum_{t=1}^{N-1} \cdot (w \cdot S^t - S_{ideal})^2 \quad (20)$$

$$\arg \min_w Error^k(w) \quad (21)$$

where $Error$ is a vector of $\langle Error_{io}, Error_{nd}, Error_{ht} \rangle$, w is $\langle w_{io}, w_{nd}, w_{ht} \rangle$ and S is $\langle S_{io}, S_{nd}, S_{ht} \rangle$.

To minimize the equation 17, we need to compute the gradient. The gradient specifies the rate of change of error with respect to the weights.

$$Gradient^k(w) = \frac{\partial Error^k}{\partial w} = 2 \cdot \sum_{t=1}^{N-1} (w \cdot S^t - S_{ideal}) \cdot S^t \quad (22)$$

Using the gradient vector, we update the weights. For this we need to set the learning rate.

$$w^k = w^k - \eta \cdot Gradient(w^k) \quad (23)$$

To find the right learning rate is important as a high value can diverge the gradient learning and a low value can slow down the learning. For each iteration of gradient descent, we employ time-based decay of learning rate.

Learning rate decay

To find an optimal learning rate forms an important part of gradient descent optimization. If the learning rate is high, it poses a risk of optimizer divergence. On the other hand, if is small, the optimizer can take a very long time to converge. Both these situations are undesirable and we can avoid them by starting out with a small value and gradually decrease the learning rate over iterations. It is called as a time-based decay of the learning rate [10].

$$lr^{t+1} = \frac{lr^t}{(1 + (decay * iteration))} \quad (24)$$

where lr^{t+1} and lr^t are the learning rates for $t+1$ and t iterations, $decay$ controls how steep or flat the learning rate curve is and $iteration$ is the gradient descent iteration number. A higher value of decay makes the learning rate curve steep as the learning rate drops quickly. A lower value can make the curve flat which can slow down the learning.

Momentum

To reach the minimum point of our convex error function (equation 20), we need to go down continuously without being blocked at the saddle points. These saddle points are where the derivative of a function is zero. Adding a momentum term to

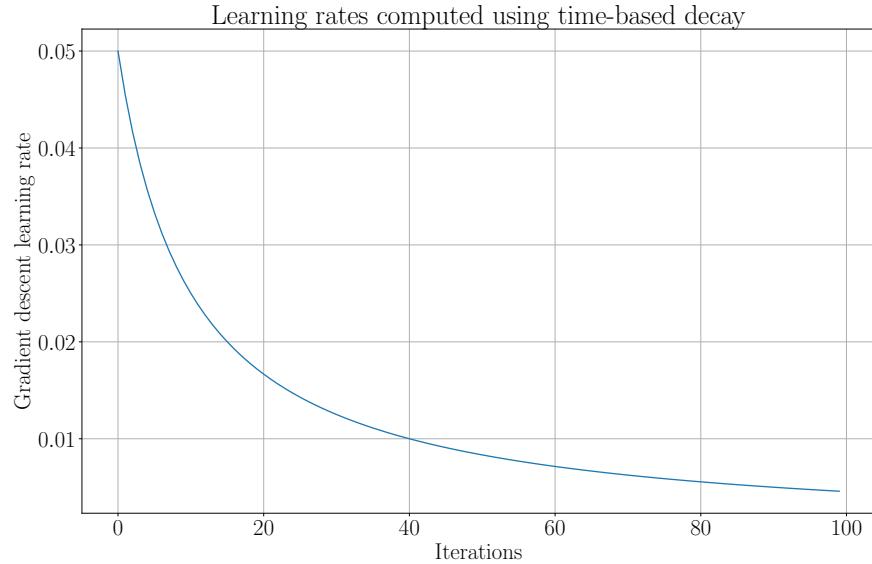


Figure 15: Learning rate for gradient descent optimizer: The plot shows how the learning rate for gradient descent evolves with iterations. It starts with a small value and decreases gradually over time. It is essential to have learning rates which neither drops too quickly or too slowly. Both of these ways can lead to divergence or slow convergence of the optimizer.

the weight parameter, we expect to avoid the local minima and should be able to converge to the lowest point quickly. It gives the necessary push to keep going down the convex error function [10, 11]. We compute the weight parameter update for each iterations using:

$$update_{t+1} = \gamma \cdot update_t - \eta \cdot Gradient(w^t) \quad (25)$$

$$w_{t+1} = w_t + update_{t+1} \quad (26)$$

where $update_{t+1}$ is the update for changing the weight parameter for the current iteration $t + 1$. $update_{t-1}$ is previous iteration update. η is the learning rate and $Gradient$ is with respect to the weight parameter w_t .

Gradient verification

We compute gradient using equation 22 and use it to update our weight parameters. To verify that the computed gradient is correct, we can approximate this gradient using the error function we formulated in equation 20.

$$\text{Gradient}(w) = \frac{\partial \text{Error}}{\partial w} \approx \frac{\text{Error}(w + \epsilon) - \text{Error}(w - \epsilon)}{2 \cdot \epsilon} \quad (27)$$

where ϵ is a very small number.

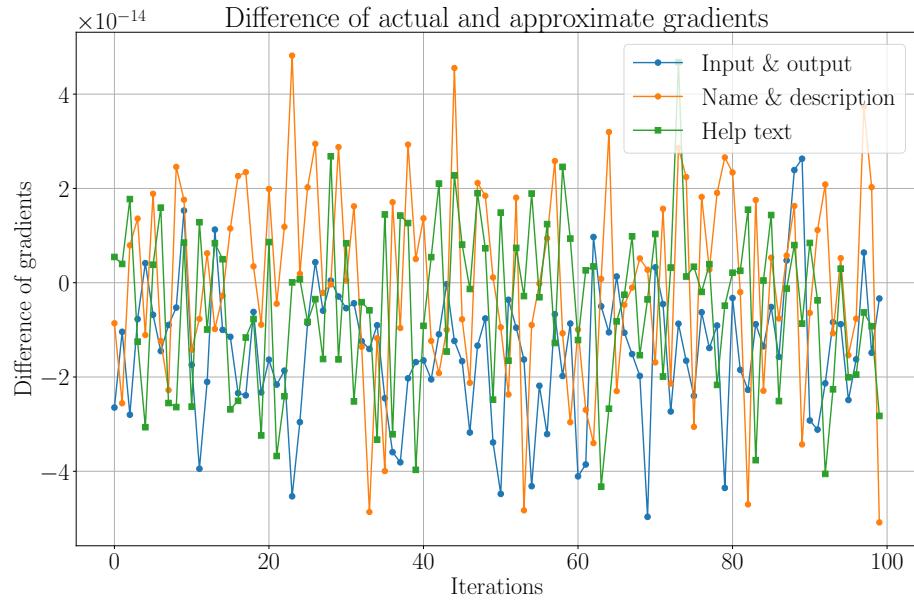


Figure 16: Gradient verification: The plot checks that the difference between the actual and approximated gradients for all the attributes across all tools computed over 100 iterations is close to 0. Gradient plays an important role in learning and it should be estimated correctly.

3 Experiments

3.1 Amount of help text

We use a maximum of 4 lines of text from the help text attribute while reading the xml files of tools. This attribute is noisy and contains text which is not useful to be set up as a basis for finding similarity. We need to be careful of the amount of text which we extract from this attribute. Experimentally, we verify that using 4 lines of text from this source works better in most of the cases.

3.2 One and two sources instead of three

We have used three different attributes for compiling the collection of tokens. These attributes are different from each other. Using them together to make one set of tokens for each tool would not be beneficial. We compute similarities for these different sources and combine them using optimization by learning optimal weights on each attribute.

3.3 Similarity measures

For finding similarity using input and output file types, we use jaccard index because we do not need to learn any "concept" hidden in a set of file types for a tool. For similarities computed for name and description and help text, we use cosine similarity. Both these similarity measures give a real number between 0.0 and 1.0 as a similarity score.

3.4 Latent Semantic Analysis

Using latent semantic analysis, we learn dense vector representation for a document. It reduces the rank of the document-tokens matrix. We need to find out this reduction factor which can improve the similarity scores among tools compared to using full-rank

document-tokens matrix. We follow an approach of lowering the rank by certain factors and look at how the values in these matrices are spread. We reduce the rank to 70%, 50%, 30%, 10% and 5% of the full-rank value. Moreover, we also verify the similarity matrices corresponding to these low-rank matrices representations. We expect the documents-tokens and similarity matrices to be more dense compared to no rank reduction. To reduce the ranks, we consider documents-tokens matrices of name and description and help text and leave the input and output matrix in its full-rank state. We singular value decomposition routine from *numpy* linear algebra package.

3.5 Paragraph Vector

In this approach, we learn a fixed-length dense vector representations for each document. We set the dimensions of these vectors. When the size of a document is lower, we use a lower number for computing the fixed-length vector. For example, in figure 5 we see that the number of tokens is higher for help text compared to name and description. To learn fixed-length vectors, we set the vector's length to be 50 for name and description and 300 for help text. Here as well, we learn paragraph vectors only for name and description and help text and not for input and output file types. We use *gensim* model to learn these paragraph vectors.

3.5.1 Distributed bag-of-words

We use this approach to learn vectors for each document. It does not use word-vectors and rely on paragraph vectors to predict the randomly chosen words from a sampled text window. Learning only the paragraph vectors is faster and computationally less expensive.

3.6 Gradient Descent

To optimize the similarity scores from coming from multiple attributes, we use gradient descent optimizer to learn weights on the scores of these attributes. Based on the similarity measures, we construct a error function to the optimizer which minimizes it.

3.6.1 Learning rates

We use backtracking line search to set learning rate for each iteration. It expedites the learning by computing a learning rate that sufficiently diminishes the error function in each iteration of gradient descent. We set a maximum iteration to 100 within which the learning saturates when the learning rate start with a value of 0.1. When we start with 0.01 or smaller, the optimization does not converge within 100 iterations.

4 Results and Analysis

4.1 Latent Semantic Analysis

We experiment with multiple values of matrix rank reduction and we find that as we reduce the rank of documents-tokens matrices, they become more dense. As we reduce the ranks, the distribution of the learnt weights also change.

4.1.1 Full-rank matrices

Figure 17 shows the similarity matrices computed using full-rank documents-tokens matrices for input and output (17a), name and description (17b) and help text (17c) attributes. We do not apply rank reduction on documents-tokens matrix of input and output file types. Using these similarity matrices, we learn each row's respective importance factor using gradient descent and then combine to get a weighted average similarity matrix (17d). Figure 18 shows the distribution of these importance factors (weights) for multiple tool attributes and we see that the magnitude of weights estimated for input and output file types is higher than the other two attributes. The high magnitude is associated with the higher values captured for the similarity matrix of input and output file types (17a).

4.1.2 70% of full-rank

In figure 17, we can see that the similarity matrices for name and description and help text are sparse. To reduce the sparsity, we attempt to reduce the ranks of the respective documents-tokens matrices of these two attributes to 70% of full-rank. For example, if the rank of a matrix is 100, we reduce the rank to 70 using singular value decomposition. Along with reducing the ranks, we reduce the singular values of these matrices as well. Comparing figures 17 and 19, we see that the name and description and help text similarity matrices start becoming more dense. The distribution of the weights also change (figures 18 and 20). At this stage, it is hard to see the effect of rank reduction. Further we reduce the ranks drastically to see the effect.

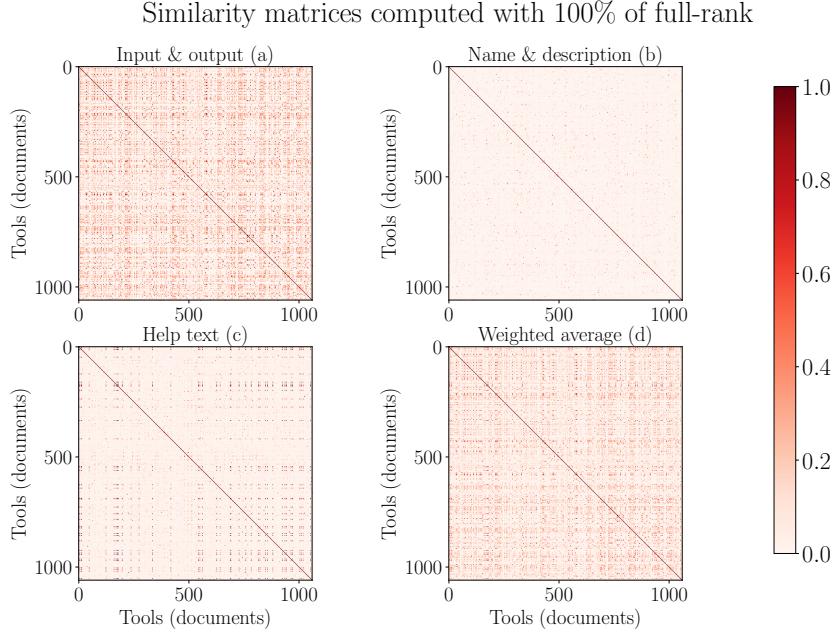


Figure 17: Similarity matrices using full-rank: The heatmap shows documents-documents (tools-tools) correlation matrices for input and output (a), name and description (b) and help text (c) attributes. The (d) shows a documents-documents (tools-tools) correlation matrix which is the weighted average computed using (a), (b) and (c) and weights (figure 18) given by the gradient descent optimizer (equation 15). The corresponding documents-tokens matrices contain their full-ranks.

4.1.3 30% of full-rank

To see the effect of rank reduction, we reduce the ranks of two documents-tokens matrices to 30% of full-rank. This is a large reduction and would amount to keeping $\approx 60\%$ (removing $\approx 40\%$) (figure 11) of sum of singular values for all the documents-tokens matrices. In figures 21 and 22, the effect of rank reduction is more visible compared to 70%. We see that the magnitudes of weights learned for input and output files start to decrease and the magnitudes of weights for name and description and help text start to increase because the corresponding similarity matrices score become more dense (figure 21).

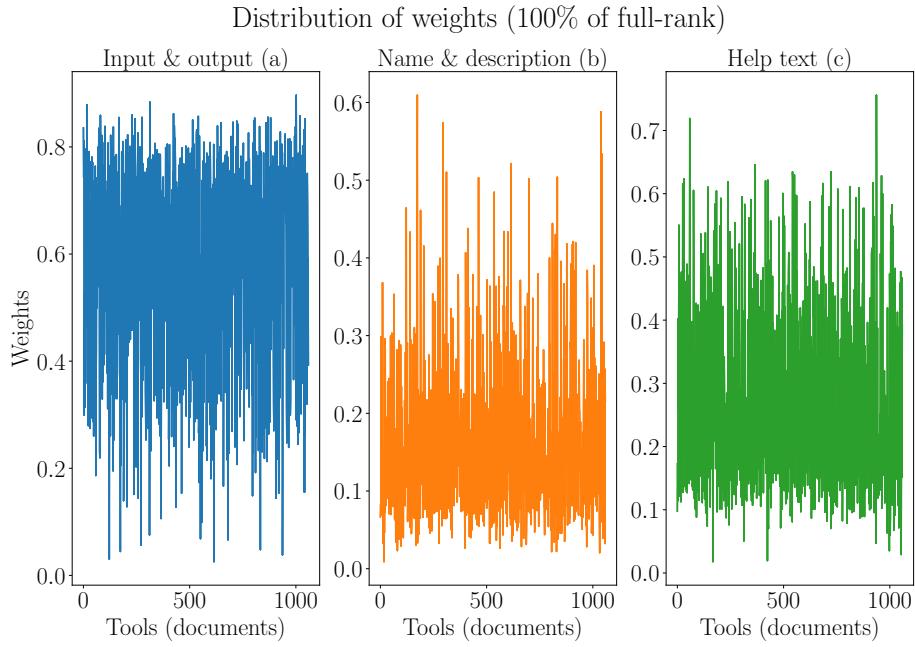


Figure 18: Weights distribution full-rank: The plot shows the distribution of weights learned by gradient descent optimizer on the similarity matrices for the input and output, name and description and help text attributes. The corresponding documents-tokens matrices contain their full-ranks.

4.1.4 5% of full-rank

To verify the reduction in sparsity, we reduce the ranks of two documents-tokens matrices to 5% of the full-rank. By choosing this low value, we consider only top $\approx 20\%$ of the sum of singular values for all the attributes (figure 11). From figure 23, we can see that all the similarity matrices corresponding to the attributes become more dense compared to figure 17, 91 and 21. Due to this, the weights distribution also change (figure 24). We learn higher weights for name and description (24b) and help text (24c). Along side, the weights on the input and output file types decrease (24a).

4.1.5 Improvement verification

To verify that the matrix rank reduction actually works and learns better similarity scores for tools (documents) which are similar in actual case, we showcase two ways.

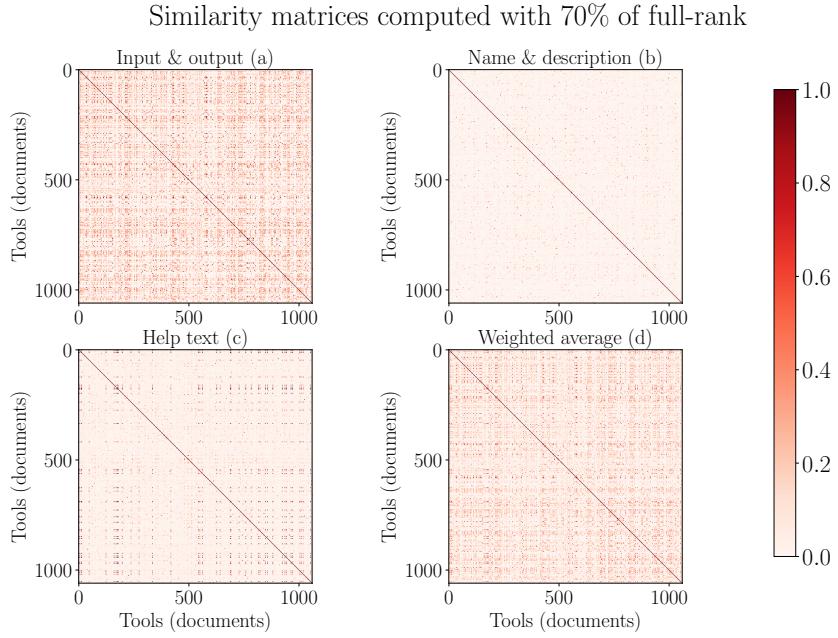


Figure 19: Similarity matrices using 70% of full-rank: The heatmap shows documents-documents (tools-tools) correlation matrices for input and output (a), name and description (b) and help text (c) attributes. The (d) shows a documents-documents (tools-tools) correlation matrix which is the weighted average computed using (a), (b) and (c) and weights (figure 20) given by the gradient descent optimizer (equation 15). The corresponding documents-tokens matrices are reduced to 70% of their respective full-ranks.

Reduction in error

We observe drop in mean squared error when we decrease the ranks of the matrices. When we reduce the ranks, the similarity scores for the name and description and help text increase. This increase accounts for their larger weights learned by the optimizer. The weights on an average become more balanced and together with higher similarity scores account for the decrease in the mean squared error. In the absence of true similarity values, this viewpoint might not be completely correct way to establish that we actually improve the performance. To be more certain, we created a visualizer using javascript and html to look through the similar tools and their respective scores and weights for all tools. The next section explains it.

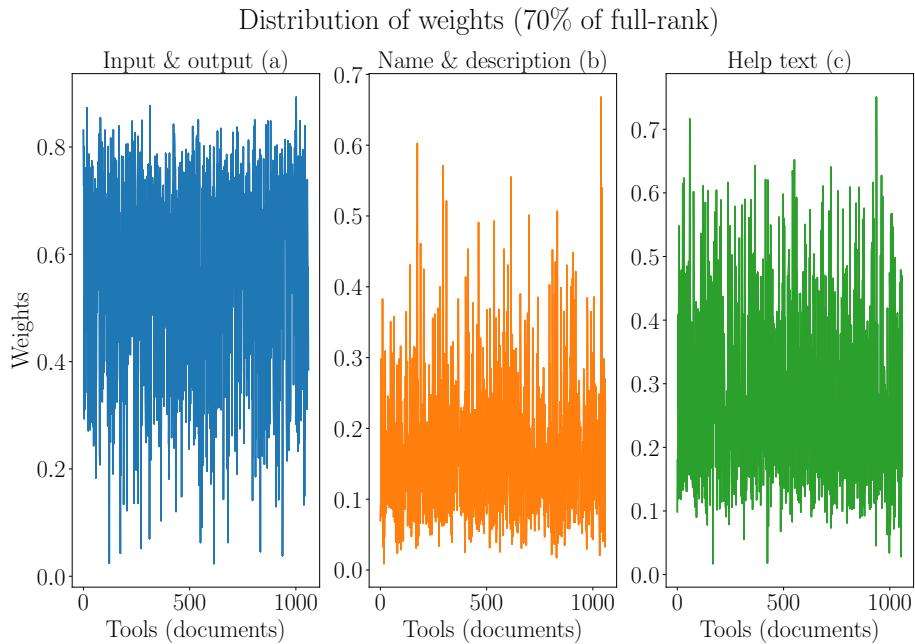


Figure 20: Weights distribution using 70% of full-rank: The plot shows the distribution of weights learned by gradient descent optimizer on the similarity matrices for the input and output, name and description and help text attributes. The corresponding documents-tokens matrices contain 70% of their full-ranks.

Visualizer for latent semantic analysis approach

We see the similar tools for a few selected tools for the different stage of rank reduction and verify if we actually fetch the similar tools.

4.2 Paragraph vectors

4.2.1 Optimization

4.2.2 Visualizer for paragraph vectors approach

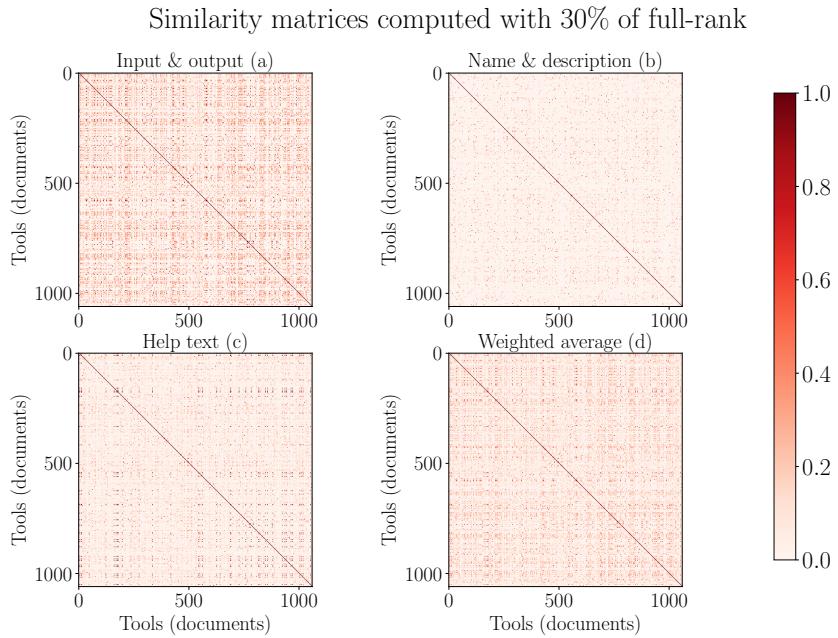


Figure 21: Similarity matrices using 30% of full-rank: The heatmap shows documents-documents (tools-tools) correlation matrices for input and output (a), name and description (b) and help text (c) attributes. The (d) shows a documents-documents (tools-tools) correlation matrix which is the weighted average computed using (a), (b) and (c) and weights (figure 22) given by the gradient descent optimizer (equation 15). The corresponding documents-tokens matrices are reduced to 30% of their respective full-ranks.

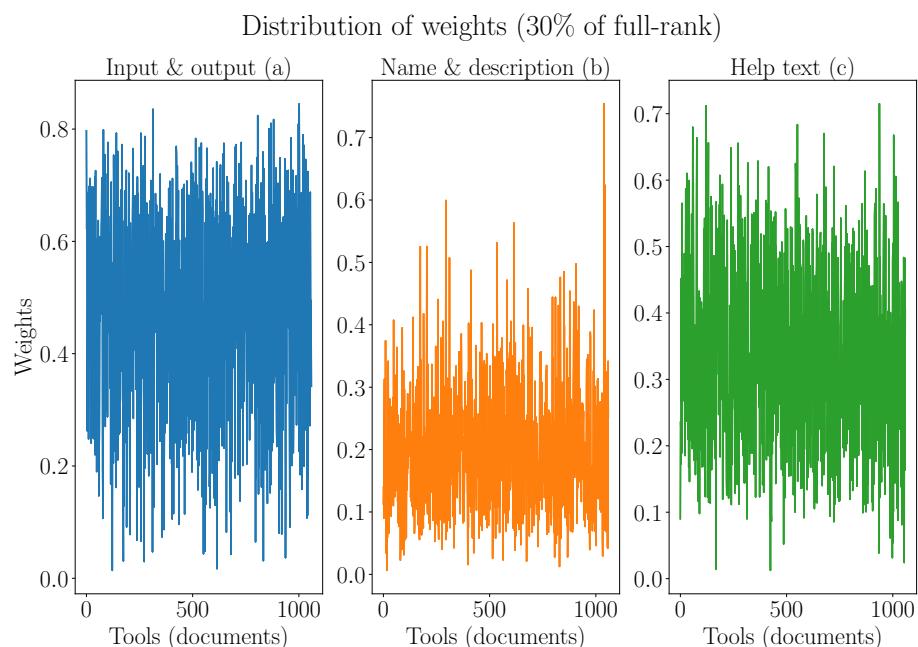


Figure 22: Weights distribution using 30% of full-rank: The plot shows the distribution of weights learned by gradient descent optimizer on the similarity matrices for the input and output, name and description and help text attributes. The corresponding documents-tokens matrices contain 30% of their full-ranks.

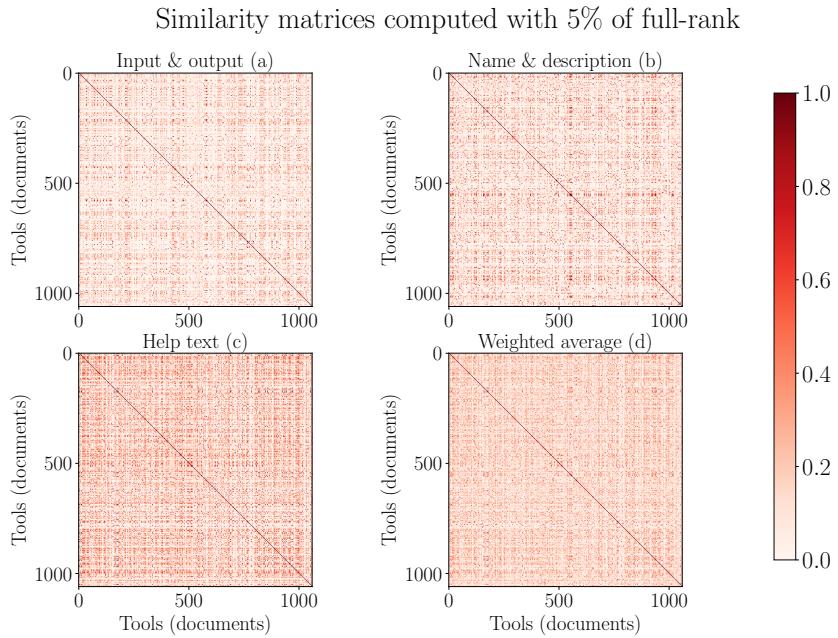


Figure 23: Similarity matrices using 5% of full-rank: The heatmap shows documents-documents (tools-tools) correlation matrices for input and output (a), name and description (b) and help text (c) attributes. The (d) shows a documents-documents (tools-tools) correlation matrix which is the weighted average computed using (a), (b) and (c) and weights (figure 24) given by the gradient descent optimizer (equation 15). The corresponding documents-tokens matrices are reduced to 5% of their respective full-ranks.

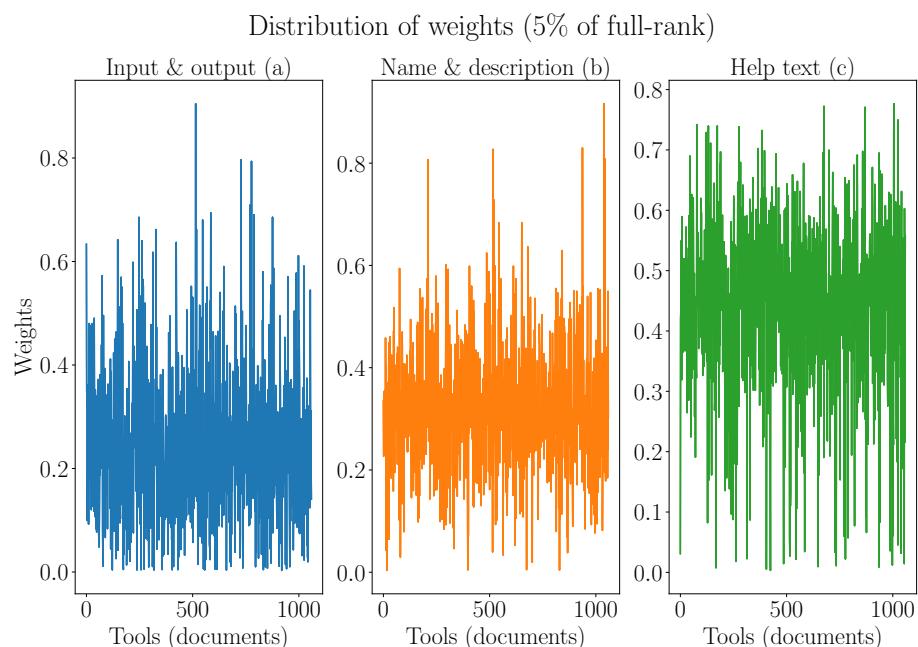


Figure 24: Weights distribution using 5% of full-rank: The plot shows the distribution of weights learned by gradient descent optimizer on the similarity matrices for the input and output, name and description and help text attributes. The corresponding documents-tokens matrices contain 5% of their full-ranks.

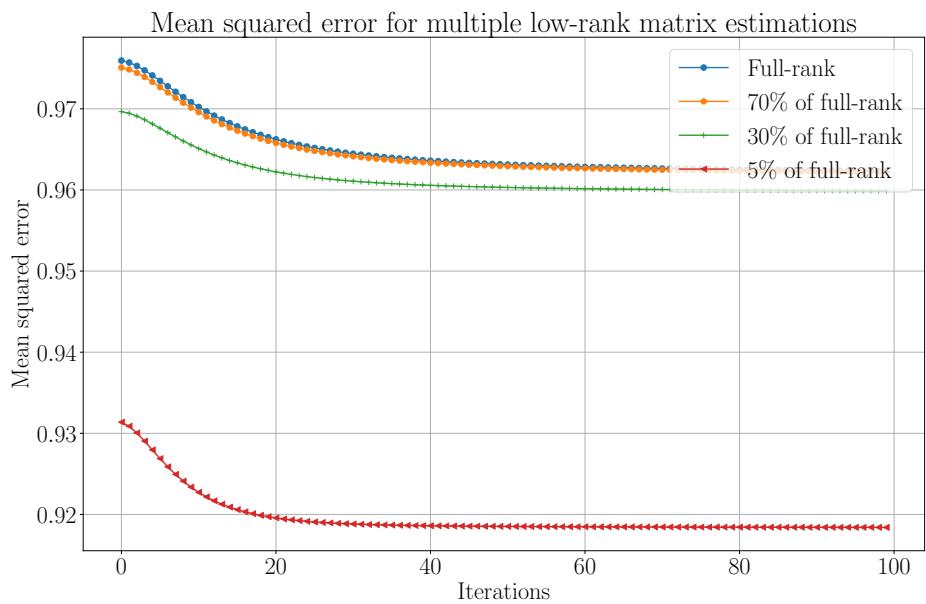


Figure 25: Mean squared error using full-rank and multiple estimations of low-rank documents-tokens matrices: This shows an mean squared error comparison computed using full-rank and various estimations of low-rank documents-tokens matrices. Each line plot shows an average error over all the tools and attributes which drops as we move along the iterations of gradient descent.

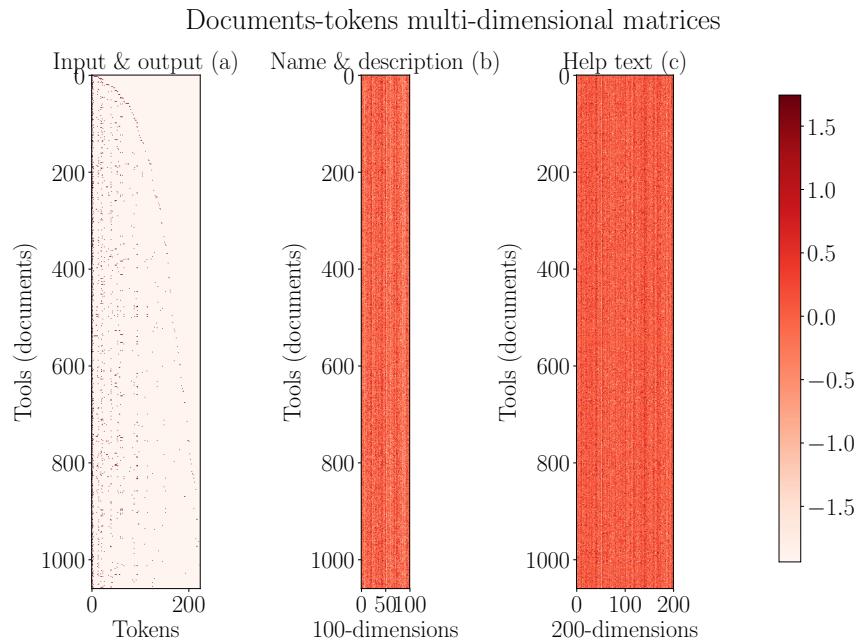


Figure 26: Documents-tokens matrices for paragraph vectors approach:

The heatmap shows documents-tokens matrices for input and output file types (a), name and description (b) and help text (c) attributes. The matrix in (a) is a document-token matrix where each entry shows a relative frequency of the token's occurrence. The matrices in (b) and (c) are 100 and 200 dimensional (each row is a document vector) respectively which means that each row of a matrix belongs to a document (paragraph) and is fixed-length and dense in nature.

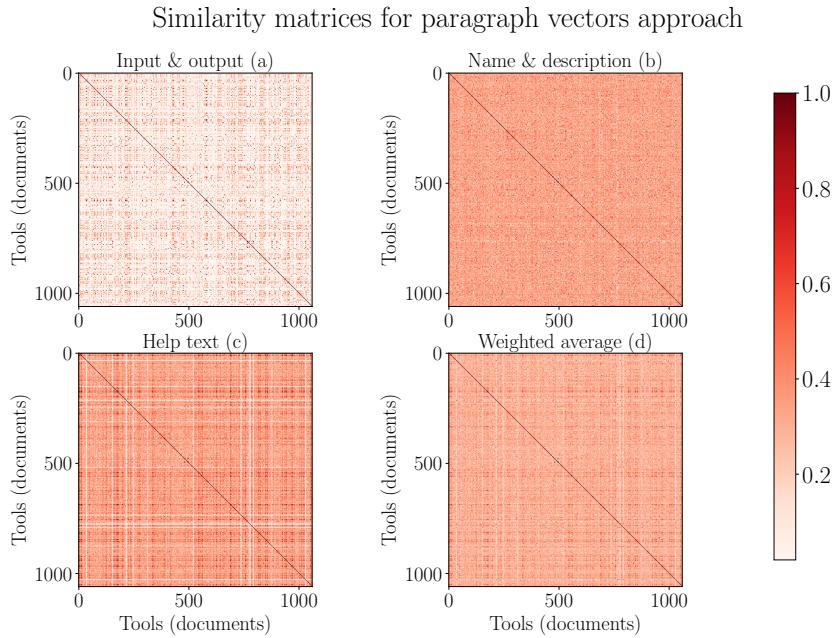


Figure 27: Similarity matrices using paragraph vectors approach: The heatmap shows documents-documents (tools-tools) correlation matrices for input and output (a), name and description (b) and help text (c) attributes. The (d) shows a documents-documents (tools-tools) correlation matrix which is the weighted average computed using (a), (b) and (c) and weights (figure 22) given by the gradient descent optimizer (equation 15). The corresponding documents-tokens matrices are computed as shown in figure 26.

Distribution of weights for paragraph vectors approach

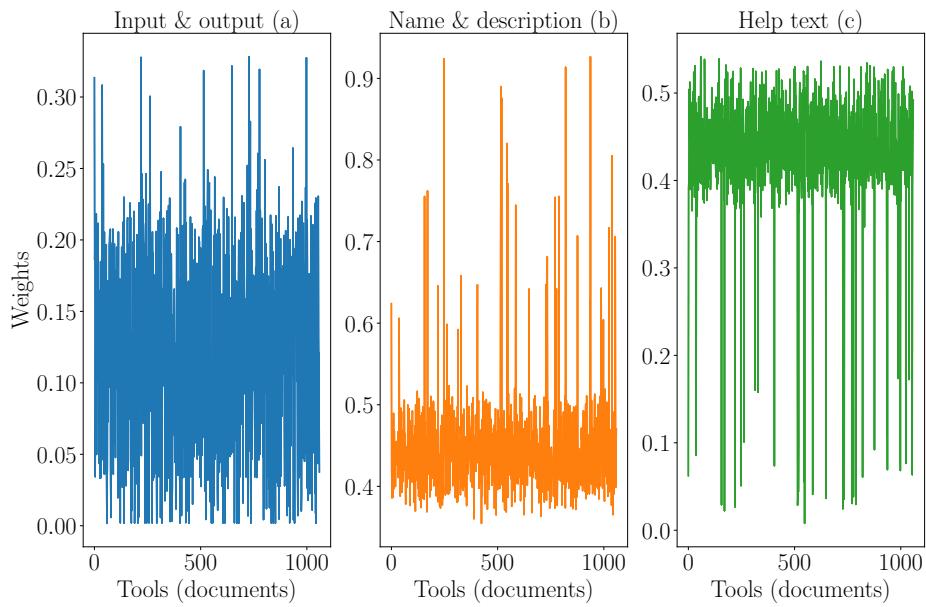


Figure 28: Weight distribution learnt using paragraph vectors approach:

The plot shows the distribution of weights learned by gradient descent optimizer on the similarity matrices for the input and output (a), name and description (b) and help text (c) attributes. The corresponding documents-tokens matrices are computed as shown in figure 26.

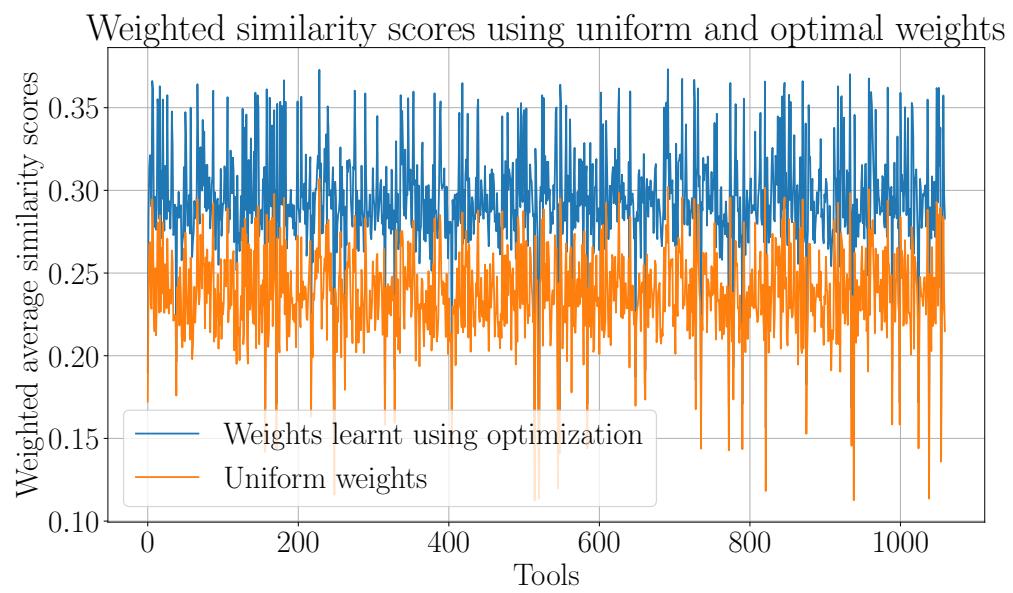


Figure 29: Similarity scores learned using optimal and uniform weights for paragraph vectors approach: The plot shows a comparison of average similarity scores computed across all tools using weights learned using optimization and uniform weights.

5 Conclusion

5.1 Tools data

The text for help text attribute was noisy containing lots of words which were generic and provided little information to categorize tools. However it was necessary as it supplied more information and due to this, it needed more filtering. The data for other attributes, input and output file types and name and description was helpful. The extraction of tools data from GitHub's multiple repositories was slow¹ and due to which we read information from the xml tools files into our local tabular file and executed analysis using that tabular file.

5.2 Approaches

We investigated two approaches to find similarity in Galaxy tools. The latent semantic analysis approach which relied on matrix rank reduction in order to remove unimportant concepts or dimensions worked better than using full-rank documents-tokens matrices. However, we were not sure about the quantity of reduction of noise dimensions and because of this we ran the risk of losing important dimensions while reducing the noise. During optimization, it gave more importance to input and output file types which were many times undesirable because we saw that But, in general, this approach is simple and take less time to execute. We can easily compute the low-rank estimations of sparse matrices using singular value decomposition.

The paragraph vectors approach worked better than the previous one in terms of finding relevant and similar tools. The documents-tokens matrices it learned for name and description and help text were dense and encoded the semantic similarities among the documents. The documents which were similar in context leaned similar vectors as proved by the dense similarity matrices for name and description and help text. However, this approach is slow. We trained to learn vectors for each document for 10 epochs and each epoch had 800 iterations. Running for 200 iterations over 10

¹The extraction of ≈ 1050 tools took ≈ 30 minutes

epochs did not fetch relevant results. Moreover, to specify the number of dimensions of paragraph vectors was an important concern. From figure 5, we saw that the average of tokens for name and description was 2-3 times higher than that of help text. So, the number of dimensions for name and description should be set to a smaller number as higher dimensions for smaller dataset would overfit.

5.3 Optimization

Learning the weights on the similarity scores from multiple attributes worked in a good way and we reached the saturation already $\approx 50^{th}$ iteration. We used standard gradient descent optimizer with mean squared error as loss function as we expected to be as close to the true distribution (based on similarity measures) as possible. We decreased the risk of getting stuck at saddle points or local minima by using momentum with accelerated gradient to update the weight parameters.

6 Future Work

6.1 Get true value

To quantify the improvements, we need to set absolute true values among tools. It means that we should create a dictionary of say 10 similar tools for each tool. Using this constraint, we can verify our findings in a reliable way. Otherwise, we need to have a visualizer to look through the results to find whether we improve. On the other hand, it is not an easy task to create logical categories of thousands of tools.

6.2 Correlation

Our similarity matrices are dense. We can nullify some low scores to retain only higher similarity scores. We can do it by finding the median value of similarity scores for a tool and set the scores to 0 which are lower than median. After this, we can optimize the scores and see the distribution and results if we improve.

Part 2: Predict next tools in Galaxy workflows

Bibliography

- [1] Z. C. Lipton, D. C. Kale, C. Elkan, and R. C. Wetzel, “Learning to diagnose with LSTM recurrent neural networks,” *CoRR*, vol. abs/1511.03677, 2015.
- [2] H. Sak, A. W. Senior, and F. Beaufays, “Long short-term memory based recurrent neural network architectures for large vocabulary speech recognition,” *CoRR*, vol. abs/1402.1128, 2014.
- [3] S. Robertson and H. Zaragoza, “The probabilistic relevance framework: Bm25 and beyond,” *Found. Trends Inf. Retr.*, vol. 3, pp. 333–389, Apr. 2009.
- [4] P. W. Foltz, “Latent semantic analysis for text-based research,” *Behavior Research Methods, Instruments, & Computers*, vol. 28, pp. 197–202, Jun 1996.
- [5] A. M. Shapiro and D. S. McNamara, “The use of latent semantic analysis as a tool for the quantitative assessment of understanding and knowledge,” *Journal of Educational Computing Research*, vol. 22, no. 1, pp. 1–36, 2000.
- [6] T. K. Landauer, “Learning and representing verbal meaning: The latent semantic analysis theory,” *Current Directions in Psychological Science*, vol. 7, no. 5, pp. 161–164, 1998.
- [7] J. Yang, “Notes on low-rank matrix factorization,” *CoRR*, vol. abs/1507.00333, 2015.
- [8] Q. V. Le and T. Mikolov, “Distributed representations of sentences and documents,” *CoRR*, vol. abs/1405.4053, 2014.
- [9] G. I. Ivchenko and S. A. Honov, “On the jaccard similarity test,” *Journal of Mathematical Sciences*, vol. 88, pp. 789–794, Mar 1998.
- [10] S. Ruder, “An overview of gradient descent optimization algorithms,” 09 2016.
- [11] I. Sutskever, J. Martens, G. Dahl, and G. Hinton, “On the importance of initialization and momentum in deep learning,” pp. III–1139–III–1147, 2013.

