

DS 5230 Unsupervised Machine Learning

Project Report

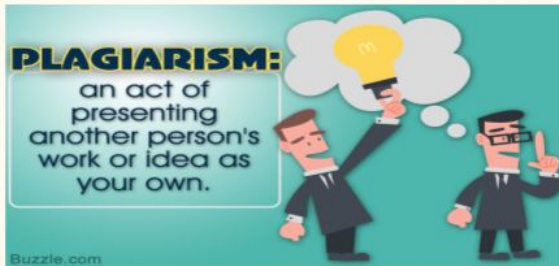
Plagiarism Detection

Authors



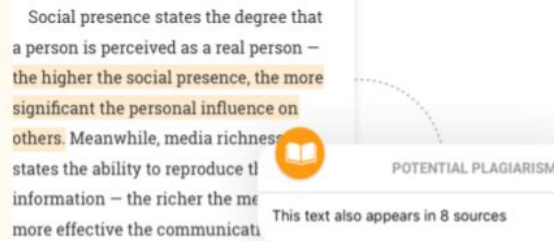
Aishwarya Vantipuli, Spatika Krishnan, Nandini Jampala

Introduction



Business Applications

- Verifying Originality
- Tracking Content Misuse
- Finding plagiarized websites
- In-Depth Plagiarism Analysis



Introduction:

Plagiarism is the “wrongful appropriation” and “stealing and publication” of another author’s “language, thoughts, ideas, or expressions” and representing them as one’s own original work. Checking assignments for plagiarism is essential in academics as assignments are used to evaluate students in their courses. If a course consists of a large number of students, it is impractical to check each assignment by a human inspector. Therefore it is essential to have automated tools in order to assist detection of plagiarism in assignments.

Business Applications:

- Students may plagiarize by copying ideas from friends, web or private tutors. Most courses in universities evaluate students based on the marks of their assignments.
- A quality plagiarism detector has a strong impact on law suit prosecution.
- To create websites quickly with less effort and to follow the quality content of a reputed website people commit plagiarism.
- Some level plagiarism checkers also provide the accuracy of grammar and paragraphing.

Dataset

Corpus of Plagiarised Short Answers.

	File	Task	Category	Text
0	g0pA_taska.txt	a	non	inheritance is a basic concept of object orien...
1	g0pA_taskb.txt	b	cut	pagerank is a link analysis algorithm used by ...
2	g0pA_taskc.txt	c	light	the vector space model also called term vector...
3	g0pA_taskd.txt	d	heavy	bayes theorem was names after rev thomas bayes...
4	g0pA_taske.txt	e	non	dynamic programming is an algorithm design tec...
5	g0pB_taska.txt	a	non	inheritance is a basic concept in object orien...
6	g0pB_taskb.txt	b	non	pagerank pr refers to both the concept and the...
7	g0pB_taskc.txt	c	cut	vector space model is an algebraic model for r...
8	g0pB_taskd.txt	d	light	bayes theorem relates the conditional and marg...
9	g0pB_taske.txt	e	heavy	dynamic programming is a method for solving ma...
10	g0pC_taska.txt	a	heavy	inheritance in object oriented programming is...

Student id_task id Question name Level of plagiarism Uncleaned Text

Dataset Statistics

Number of questions	5
Number of answers per question	20
Number of non plagiarised answers	38
Number of lightly plagiarized answers	19
Number of heavily plagiarized answers	19
Number of cut level plagiarized answers	19
Mean number of words per answer	216

Each question has 19 students answers and 1 wikipedia answer

Link

https://ir.shef.ac.uk/cloughie/resources/plagiarism_corpus.html

Dataset:

The dataset is taken from a freely available Clough-Stevenson corpus . The corpus consists of answers to five short questions on a variety of topics in the Computer Science field. The five short questions are:

1. What is inheritance in object oriented programming?
2. Explain the PageRank algorithm that is used by the Google search engine.
3. Explain the Vector Space Model that is used for Information Retrieval.
4. Explain Bayes Theorem from probability theory.
5. What is dynamic programming?

These questions are represented as tasks(a-e) in the dataset where each task maps to each of the five questions respectively.

Each question has 19 students answers and 1 corresponding wikipedia answer.

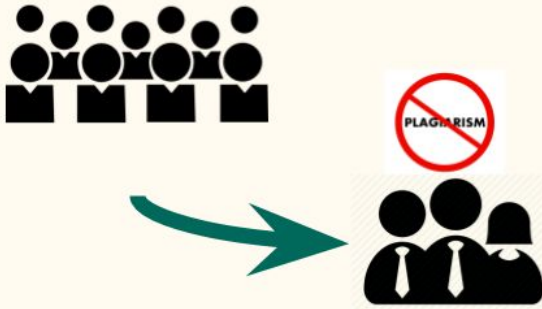
We are given the extent of plagiarism for every answer in a separate file_information file. We use this file to check accuracy and precision of the results achieved by our algorithm.

There are four different levels of plagiarisms. They are cut, light, non and heavy. Cut refers to answers that are cut, copy, pasted from Wikipedia answers. Heavy refers to answers that are heavily copied from Wikipedia answers, light refers to answers which are slightly copied from Wikipedia answers with extreme changes to the structure of answers. Non refers to non-plagiarised work.

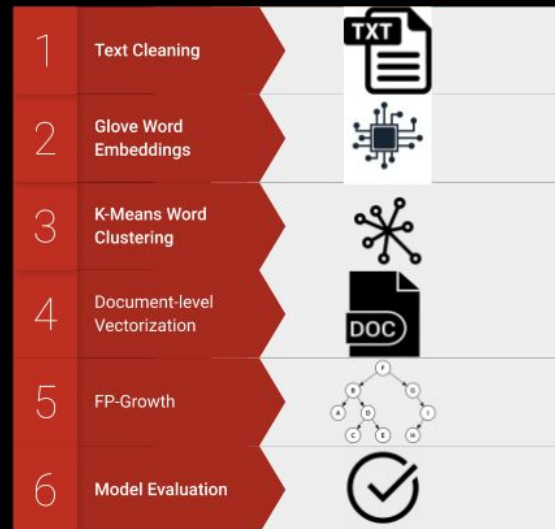
Goal

- Finding the list of students with plagiarised work.

For this we proposed two methods.



METHOD 1



Project Goal:

The goal of this project is to find the extent of plagiarism in each student's submission for every question using unsupervised machine learning methods.

For this, we have Proposed two methods.

Method 1 Overview:

- 1. Text Cleaning:** Initial step for building text models. Accuracy of machine learning models on textual data depends heavily on text pre-processing.
- 2. Glove Embeddings:** GloVe is a log-bilinear model that is trained on the non-zero entries of a global word-word co-occurrence matrix, which tabulates how frequently words co-occur with one another in a given corpus. Output of this step is vector representation for each unique word in the corpus.
- 3. K-Means Clustering:** Cluster these word vectors into K clusters based on their contextual meaning. The output of this step is K word clusters.
- 4. Document Level Clustering:** From the obtained clusters and its respective labels from the K-means clustering, we now use these clusters to convert our data in each document to vectors.
- 5. FP-Growth Algorithm:** In this step, we mine frequent patterns in documents based on sentence similarities which outputs a list of items which were plagiarised.
- 6. Model Evaluation:** We evaluated the results manually calculating accuracy, precision, F1- score .

METHOD 1

Text Pre-Processing

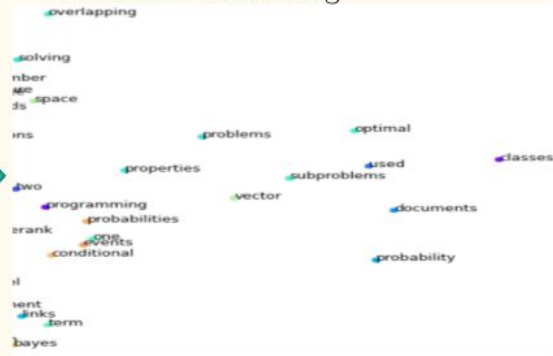
- Stopword Removal
- Punctuation Removal



Glove Embeddings



K-Means Clustering



Text Pre-Processing:

For pre-processing we have to make sure that all the words are lowercase with all the punctuations removed. We use the NLTK sentence tokenizer to get a list of sentences so that we know where each sentence begins and ends. There are various latin words in our answers as well so we use "latin1" encoding for reading the text. Implemented EDA on text data to get insights like top most used words in each question, Bigram representation in word clouds etc.,

Glove Embeddings:

For creating vector word embeddings, we used a method called GloVe. GloVe is an unsupervised learning algorithm for obtaining vector representations for words. List of words in sentences of the documents is fed to the Algorithm. After performing GloVe on our preprocessed data we get vectors of each unique word in our corpus. We can use this representation for clustering and further processing.

K-means clustering:

K-means clustering is a type of unsupervised learning, which is used when we have unlabelled data. The goal of this algorithm is to find groups in the data, with the number of groups represented by the variable K. We determined K value using Grid Search, a hyperparameter selection method. The output of this step is K (=25) word clusters.

Document Level Vectorization

Sentence vector: each entry is a word representing cluster label it belongs to.

```
taska ('g0p0': [0, 1, 2, 0, 1, 1, 2, 4, 4, 5, 2, 5, 5, 0], [2, 2, 4, 4, 5, 4, 2, 5, 4, 4, 4, 2, 4, 6], [0, 7, 7, 2, 6, 5, 2, 2, 2,
#####
taskb ('g0p0': [10, 5, 5, 9, 0, 9, 2, 1, 2, 2, 2, 1, 1, 12, 7, 7, 10, 2, 2, 2, 2, 1], [9, 2, 2, 1, 2, 2, 10], [2, 2, 2, 9], [10], [
#####
taskc ('g0p0': [0, 0, 16, 0, 16, 0, 0, 5, 5, 17, 18, 18, 17], [2, 2, 16, 2, 4], [16, 0, 0, 0, 0, 5, 6, 16, 0], [6, 6, 2, 5, 0, 2, 5, 16,
#####
taskd ('g0p0': [22, 22, 0, 7, 9, 0, 22, 0, 12, 0, 0, 1, 7, 12, 1, 6, 1, 7, 1], [0, 24, 24, 1, 0, 1, 24, 24], [22, 0], [0, 12, 22], [2, 7, 7,
#####
taskd ('g0p0': [0, 0], [5, 1, 0, 9, 0, 0, 21], [2], [21], [0, 0], [6, 2, 12, 7, 7, 0], [0, 0, 0], [24], [0, 0], [0], [0, 9], [1]], 'g0p0': [
```

Question type

Student id

Document Vector: each entry is sentence numbered according to the rule.

```
taska ('g0p0': [0, 1, 2, 3, 4, 5, 6, 7], 'g0p0': [8, 9, 10, 11, 920, 13, 1075, 15, 16, 17, 18, 19], 'g0p0': [20, 21, 22, 1032, 24, 1067, 26,
####
taskb ('g0p0': [390, 227, 228, 424, 230, 424, 232, 233, 391, 235, 938, 409, 410, 239, 1052], 'g0p0': [241, 242, 243, 244, 245, 246, 247], 'g0
####
taskc ('g0p0': [428, 429, 430, 431, 432, 433, 434, 435, 436, 437, 920, 439, 920, 441, 442], 'g0p0': [443, 444, 445, 446, 447, 920, 567, 920,
####
taskc ('g0p0': [673, 793, 794, 879, 677, 678, 1036, 884, 681, 732, 683, 888, 890], 'g0p0': [686, 876, 875, 794, 879, 880, 1075, 883, 884, 695
####
taskd ('g0p0': [1074, 892, 920, 1076, 1074, 896, 1070, 1071, 1074, 1075, 1013, 1059], 'g0p0': [903, 904, 905, 920, 950, 948, 1066, 920],
####
```



Vector Representation for document-level embeddings:

Sentence Vector: Each entry is a word representing the cluster label it belongs to.

Document Vector: Each entry represents sentence number according to the following rule.

From the obtained word clusters and its respective labels from the K-means clustering, we now use these clusters to convert data in each document to vectors. For this, we do the following:

- Determine the cluster label for each of the words in sentences and substitute the word with its respective cluster label.
- From the above step, we merge all the words and form sentence vectors.
- From these sentence vectors, we form an answer(document) vector by assigning a unique number if two sentence vectors are different. If two sentence vectors are the same, we assign each of them the same number.

Thus we get a vector for each answer(document). This vector can be viewed as an item-set and each answer can be viewed as a transaction. Thus we get a list of transactions which contain item-sets. We can feed this into the FP-Growth Algorithm to get frequent item sets, i.e., plagiarised answers.

FP-GROWTH

Transactions	Itemset (Min sup > 3)
g0pB_taske	[686, 876, 875, 794, 879, 880]
g0pB_taskc	[241, 242, 243, 244, 245, 246, 247]
g0pA_taska	[0, 1, 2, 3, 4, 5, 6, 7]
g0pB_taska	[8, 9, 10, 11, 920, 13, 1075, 15, 16, 17, 18, 19]
g1pB_taskc	[489, 631, 632, 634, 635, 636, 1052]

100 Documents

Model Failures: Accuracy is low with Multi-Class Classification because of imbalances in dataset.

Question	List of students with plagiarism
Task e	{'g1pB', 'orig', 'g3pB', 'g0pE', 'g4pB'}
Task a	{'g3pC', 'orig', 'g2pC', 'g0pE', 'g4pC'}
Task d	{'g2pA', 'g2pB', 'g0pB', 'orig', 'g3pA', 'g0pC', 'g1pA'}

Model Evaluation

Binary Classification

Accuracy	79.03%
Precision	90%
Recall	53.55%
F1-Score	68.81%

Multi-Class Classification

Accuracy	61%
Precision	75%
Recall	45.5%
F1-Score	66.23%

FP-Growth Algorithm:

The FP-growth algorithm is currently one of the fastest approaches to frequent itemset mining. We can find which documents are plagiarised by just looking at the transactions which in our case are student answers. Here each student answer is one transaction and their vector representation of sentences are itemsets.

Model Evaluation:

Based on Support counts output is classified into 4 classes. For heavy, minsup is 8, for cut its 6, for light its 4, for non its 3.

Multi-Class Classification: Evaluated final results with actual labels given in dataset i.e., light, heavy, cut, non. Accuracy is as low as 60% but high precision is to be noted. This is usually due to class imbalances.

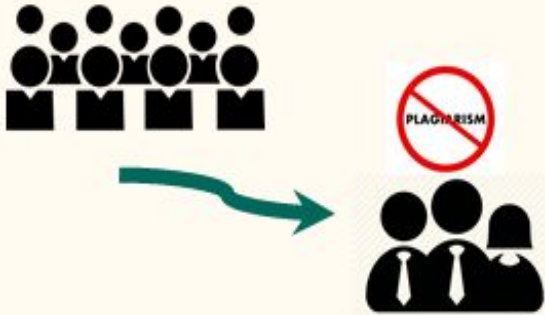
Binary Classification: To Tackle Class Imbalance problems, Binary Classification is applied. Documents with light,cut,heavy labels are categorized as Plagiarized and with non as Non-Plagiarized. Accuracy is as high as 80% and precision is 90%

Method Failures:

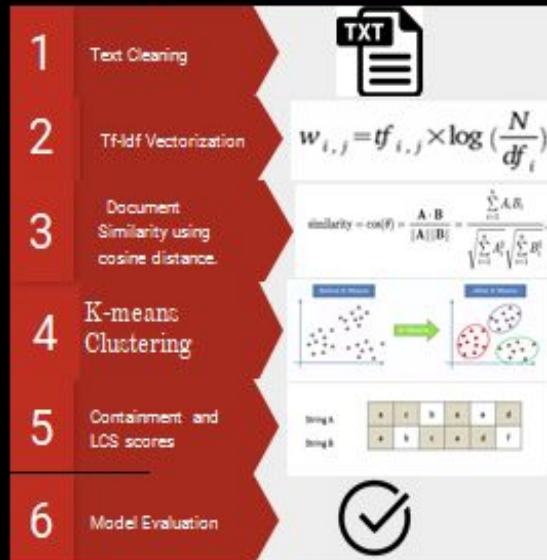
This approach of using FP-Growth for Detecting Plagiarism got pretty decent results. Although, not absolutely great, the results were better than the traditional methods used for this task. The main drawback is insufficient data and class imbalances. One of the possible ways to tackle this problem is collecting more data as we can't regenerate data in these types of applications and experimenting with different support counts for class thresholds.

Goal

- Finding the list of students with plagiarised work.
- Using a different approach



METHOD 2



Pre-Processing Techniques

We have experimented with other unsupervised machine learning algorithms, to get better results for multi-class classification. So, we have implemented a whole new approach to achieve our objective.

We started our method-2 by cleaning and preprocessing our documents as the optimal quality and performance of the clustering methods depend largely on the pre-processing techniques. For each document in the corpus we followed the 3 main steps:

1) Normalization: This means to have each word in the uniform format. We achieved this by converting all text to the lower case, removing punctuation and converting numbers to their word equivalents.

2) Stemming: It is a process of reducing words to its root form. We have used Porter stemming algorithm for this. This is used to catch the similar words even if they are in its non-root forms.

3) Tokenization and removal of stop-words: This is the process of splitting the given text into smaller pieces called tokens. So this represents every document as a list of tokens without having the stop-words.

Tf-idf Vectorization And Cosine Similarity

$$\text{tfidf}(t, d, D) = \text{tf}(t, d) \times \text{idf}(t, D)$$

	Doc 1	Doc 2	...	Doc m
Term(s) 1	12	2	...	1
Term(s) 2	0	1	...	0
...
Term(s) n	0	6	...	3

Output of Tf-Idf :

Vector form of g0pA_taska document

```
[ [ 0, 0, 0, 0, 0, 0, 0, 0,
  0, 0, 0.57501852, 0.36113499, 0, 0, 0,
  0, 0, 0, 0.47152341, 0.57501852, 0, 0, 1,
  0, 0, 0, 0.40824829, 0.40824829, 0, 0, 0,
  0.40824829, 0, 0, 0, 0, 0, 0, 0,
  0.40824829, 0.40824829, 0.40824829, 0, 0, 0, 0, 0,
  0, 0.490779, 0.490779, 0, 0, 0, 0, 0,
  0, 0, 0, 0.490779, 0.490779, 0, 0, 0,
  0.26394385, 0, 0, 0, 0, 0, 0.3397124, 0,
  0, 0.79189114, 0, 0, 0, 0.15659897, 0, 0,
  0, 0, 0, 0, 0, 0, 0, 0,
  0, 0, 0, 0, 0, 0, 0, 0,
  0, 0.59430676, 0, 0.35253792, 0.59430676, 0, 0,
  0, 0, 0, 0, 0, 0, 0.41144595, 0,
  0, 0, 0, 0, 0, 0, 0, 0,
  0, 0, 0, 0.26999725, 0, 0.9099135,
  0.31497221 ] ]
```

Output of Cosine Similarity :

	0	1	2	3	4	5
0	1.000000	1.000000	-0.001518	0.029564	0.029697	0.317527
1	1.000000	1.000000	-0.001518	0.029564	0.029697	0.317527
2	-0.001518	-0.001518	1.000000	0.040489	0.031902	0.042664
3	0.029564	0.029564	0.040489	1.000000	0.999642	0.081541
4	0.029697	0.029697	0.031902	0.999642	1.000000	0.077855
5	0.317527	0.317527	0.042664	0.081541	0.077855	1.000000

After preprocessing to apply data mining techniques we have converted these words to a machine-readable format. So, to achieve a relationship between the important words that are not just placed adjacently in a document but also placed throughout the corpus we used Tf-idf vectorization.

Tf-idf Vectorization:

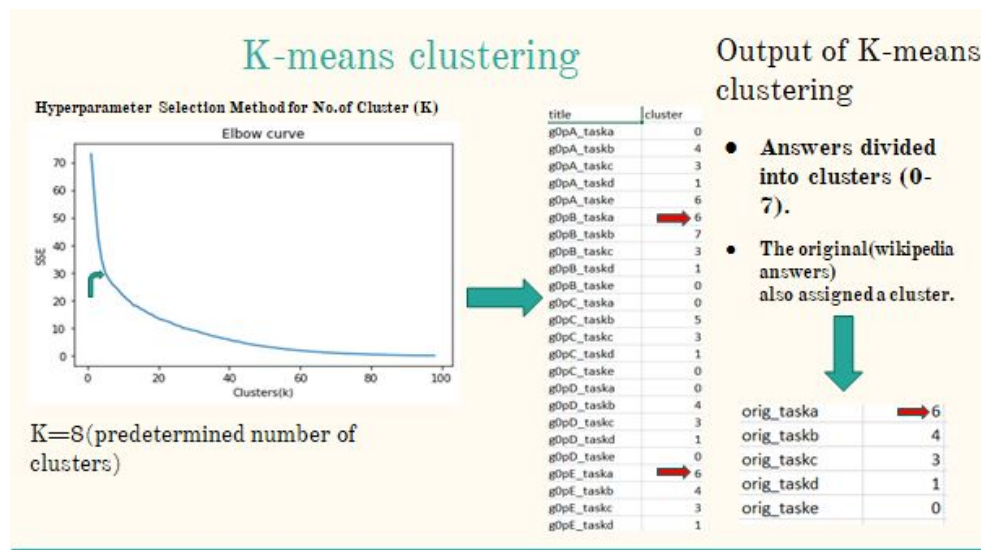
Term frequency (Tf): Count of how frequently that expression (term, word) occurs in the document.

Inverse document frequency (idf): It is used to calculate the weight of rare words across all documents in the corpus. This value is obtained by dividing the total number of documents by the number of documents containing the term, and then taking the logarithm of that quotient.

Words that occur frequently within a document but not frequently within the corpus receive a higher weighting as they are assumed to contain more meaning in the document. Now each document in the corpus is represented in a vector form that contains the product of Tf and idf values for each word in the document.

Cosine Similarity:

This is measured against the tf-idf matrix and can be used to generate a measure of similarity between each student's answers and the other answers in the corpus. High values indicate a high level of similar words between two documents.



K-means clustering

The K-means clustering was iterated multiple times with different values of k(number of clusters) to compare the SSE(Sum of Squared Errors) and find the optimal value of k.

K-Means clustering was performed in two ways:

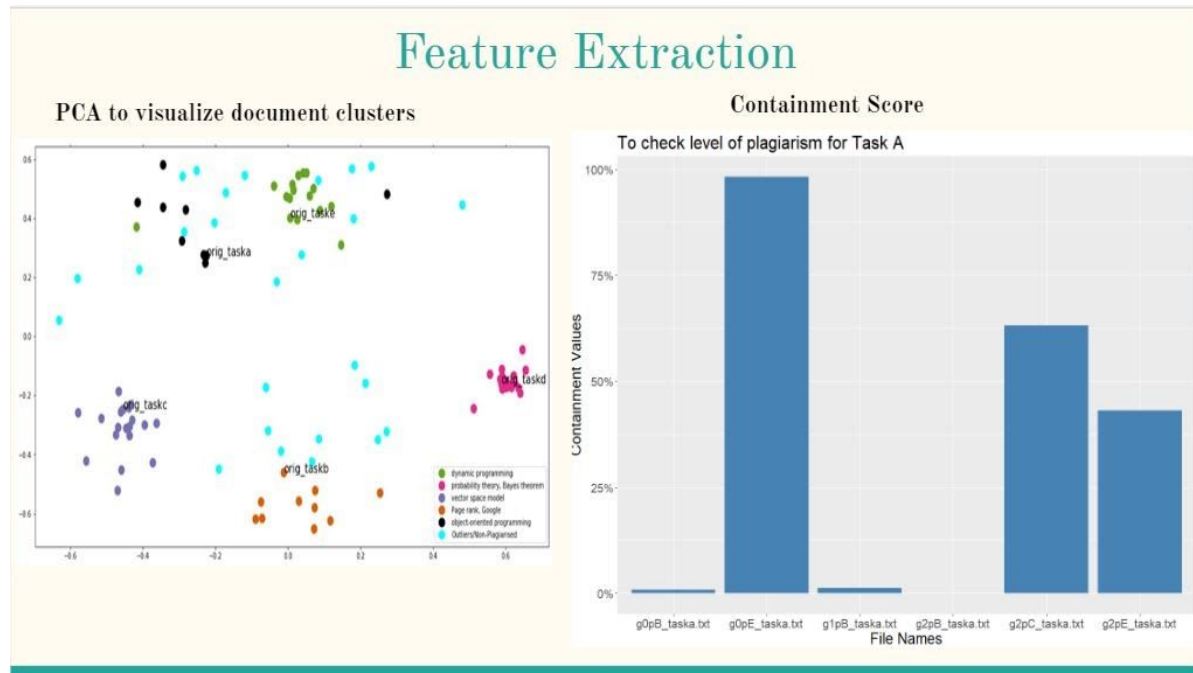
- Only Using the Tf-idf matrix.
- Reducing the dimensionality using PCA i.e by calculating Eigenvectors and Eigenvalues from the Cosine Similarity matrix and using k Eigenvectors corresponding to the topmost k Eigenvalues to cluster into k-dimensional document vectors.

The Elbow Method is an effective way to find the k parameter. If the line chart looks like an arm, then the "elbow" on the arm is the value of k that is the best.(In our case,k=8)

Failure of method using Tf-idf:

We have implemented method 'b' because we observed that method 'a' gave a rather smooth SSE curve with all high values of SSE whereas Method 'b' gave significantly low SSEs for any comparable value of k because of the reduced dimensions. For method 'b' we see that the SSE curve is shown above with an elbow near k=8.

Now our student answers and our reference(wikipedia) answers are assigned a cluster label from(0-7). So if the same cluster label is assigned to the reference and to the student answers that would imply there is some level of plagiarised work.



As shown in the above cluster figure we can see that there are student answers belonging to the same cluster as the reference answer, and there are some outliers not belonging to any of the clusters from the reference answers, thus the outliers are falling under non-plagiarised work.

Now to classify, upto which extent the students have plagiarised the work we calculate the Containment and the LCS scores in each cluster between the student answer document and the reference(wikipedia) document that is falling under the same cluster.

Containment Score:

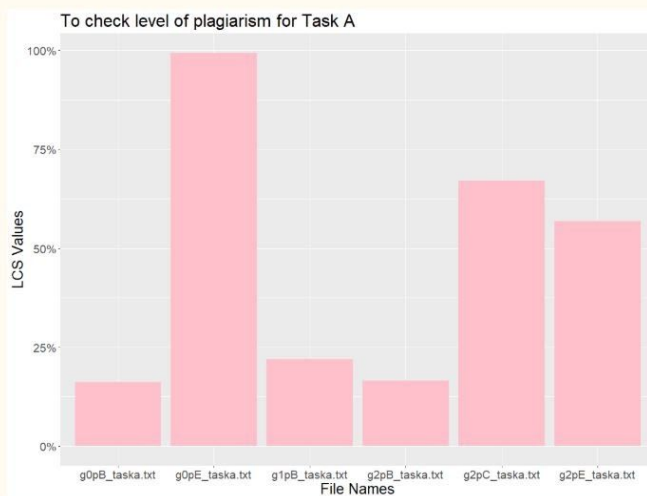
Containment is a good way to find overlap in word usage between two documents. it may help identify cases of cut-and-paste as well as paraphrased levels of plagiarism.

We are calculating the containment values between an answer and source text belonging to the same cluster according to the following equation.

$$\frac{\sum count(ngram_A) \cap count(ngram_S)}{\sum count(ngram_A)}$$

High containment scores refers to high levels of plagiarised work. Since plagiarism is a complex task with varying levels of plagiarism, it's useful to include other measures of similarity too like LCS.

Longest Common Subsequence Score



Model Evaluation

Classification Metrics

Accuracy	92%
Precision	93%
Recall	92%
F1-Score	92%

Longest Common Subsequence:

The longest common subsequence is the longest string of words that are the same between the Wikipedia Source Text and the Student Answer Text. This value is then normalized by dividing by the total number of words in the Student Answer Text. High LCS score refers to high levels of plagiarised work.

We can classify the student answers into different levels of plagiarism by looking at their LCS and containment scores by comparing them with the wikipedia answers. For eg : If the containment and the LCS scores are in the range of 80-100 % then we can say that it falls under the “cut” classification. If it is in the range of 50-80% then we can say that it falls under the “heavy” classification. If it is in the range of 15-50% then we can say that it falls under the “light” classification.

Model Evaluation:

Multi-Class Classification: Evaluated final results with actual labels given in dataset i.e., light, heavy, cut, non. We saw a high rise in accuracy as it is 92% and also a high rise in precision is to be noted.

Conclusion

- Found Plagiarism using FP-Growth and Containment Scores
- 2nd method results in better classification as it eliminates the use of imbalanced classes.
- Hierarchical Clustering gives inaccurate results hence K- Means is chosen.
- Proposed methods can be used for detecting Plagiarism in Programming languages as well.

Conclusion:

In our first method we have tried to find the extent of plagiarism using FP-growth and in the second method we have used Containment and LCS scores to the clusters. We have developed a system based on principles of vector representations. We also compared the results between Hierarchical Clustering and K-Means. K-means give better results than hierarchical clustering as the latter is not suitable for large datasets and Initial seeds have a strong impact on the final results thus giving us inaccurate cluster labels.

We have seen a rise in the accuracy for multi-class classification in the second method because we have clustered the abundant class and when we apply containment and LCS scores between documents in each cluster, we do not face the problem of imbalanced classes as it itself finds the similarity score and classifies according to the score.

As an extension to this work, we can include vector representations using other techniques such as word2vec or fast text. Another future work is to change the domain from documents to programming. In programming, we would need stricter noise removal as well as the support count of FP-growth algorithms will increase.