

Application of Markov Chain in Linguistic Analysis

Cloris Li, Irene Li, Hongda Li, Khoa Vo

March 2020

1 Introduction

In this paper, our objective is to model novels written by different authors through Markov chain with letter to letter transition matrix. We will focus on two authors: Charles Dickens and Mark Twain, which are chosen based on group members' knowledge and preference on English writers. We consider the above two authors are well-known to English readers, and they have quite different writing styles. Charles Dickens was an British writer from 1812 to 1870 with style more poetic and focuses on Victorian romanticism[1], while Mark Twain was an American writer from 1835 to 1910 known for keen satire and careful choice of diction[28].

Our purpose is to investigate whether the two authors can be represented by two separate clusters of transition matrices generated from their novels or there is an overlap between their transition matrices under various assumptions and conditions. We will also examine the closeness of author's writing to his own profile. This is analyzed through calculating distances of each novel's transition matrix to his own author's profile matrix, as well as the cross distances of each novel's transition matrix to the other author's matrix. Transitional matrix is first be defined by letter by letter base and later modified to two letter by two letter sequence. The effect of changing assumption on syntactic elements such as handling proper noun and emphasizing rarely used words is studied. We also investigated the effect of using different distance measures for comparison.

2 Background

Our group project was interested in the real-world application of Markov Chain. We first started with the idea of using Markov Chain to model the texts of a few authors as trained data to identify the author of a text. After some discussions, we decided to narrow our project down to investigate the novels of two authors, Charles Dickens and Mark Twain, to see whether transition matrices generated by the novels of each author has a "cluster" that is distinctly theirs or there is a large overlap between their works.

Markov chains had been a commonly technique applied in literature analysis. Since each author has their own writing style, the literature written by different authors have specific syntactic

elements that could be detected and traced back to its author. Using the transition matrix of the Markov Chain, we could construct a profile for the author's writing and compare it to other author's writings. The application of authorship attribution was commonly used in plagiarism detection, deducing the writer of threatening or harassing e-mails and other problems with unclear authorship.[16]

A study on this topic was done and published in 2014 by Osipovs, Rinkevics, Kuleshova and Borisov. Their research is focused on using Markov chains in the task of text authorship identification. In the research, they created a model using Markov chain transition matrix to represent the author's writing style. This model is constructed based on the texts for which the authorship of the targeted author is known. The result produced from the model should indicate whether a new text, whose author is unknown, is written by the target author or not. They used texts with more than 26000 words to make sure it's large enough to contain information about specifics of writing styles. They also implement various ways to pre-process the texts such as dealing with punctuation and short words to see if yields different results. Results of this study indicate that most useful information about the author's style is contained in long words, and removal of short words and punctuation marks does not affect the model.[15] With this reference, we will ignore the effect of punctuation by considering them as space that separates words, and detail will be explained later in model assumption.

A paper published in 2006 by Sanderson, Conrad, and Simon Guenter used similar techniques in relatively short texts (300 to 5000 words) authorship attribution. They worked on a large dataset which includes 50 authors with each author covering several topics. Because they wanted to avoid the probability that the model is based on topics rather than the author's writing styles. The completeness of their model was verified by matching testing texts to the correct authors. In conclusion, they claimed that analyzing texts using character and word Markov chains is less useful when dealing with short texts. A good performance of the model should be constructed base on texts with more than 5000 words.[16] This is consistent with the later research mentioned above, where long texts are used as data for text analysis. With these reference in mind, we decide to focus on analyzing novels from selected authors since the long texts will provide enough information to build models that represents unique writing style for the two chosen author. We

also use different ways of pre-processing our texts in adjustment section to see how results changes with or without proper nouns.

Authorship identification is not only applicable to book literature, but also social media text post. A research by Daniel Freeman aim to predict the author of an unknown post on the social media network Twitter character Markov chains. Freeman chooses to analyze Twitter posts since Twitter posts are different with formal writings in both length and writing standard, such as the use of hashtags. Markov chain transition matrix for this model is constructed based on tokens of characters and also tokens of words. A collection of Markov chain transition matrices from different authors are generated. Freeman calculates the probability that a new text could be represented by the each existing Markov chains, and the highest probability corresponds to the authorship predicted by this model. The author found that it is easier to identify authorship with Twitter posts using character as tokens, since there are more unique transitions between characters in informal short texts that captures writer's individuality.[12] Our project is intended to perform data modeling by tokens of characters on long novels, on which the technique mentioned in this paper did not work best. We consider authors can variate their writing style greatly in novels, which capture more individuality compare to formal literature like academic writings. This is interesting because with different types of writing, we will need to model our data differently.

A paper written by Jagtap, Avinash, and Limbore has many similarities with our project. We create model by using Markov chains that based on the probabilities of subsequent letters, and all punctuation and formatting are removed. Also, the capital letters in the texts are converted into lower cases. Thus, the transitional matrices in both papers have 27 state spaces which corresponds to 26 letters and space. However, the difference between their model and ours is that they calculate the average of the elements of the matrices from all the texts by each author and then consider the probability of new text being generated by each of the transition matrices. Whereas we have different matrices for each text by the same author.[14] Instead of attributing testing texts to the correct authors, we verify our model by illustrating the clouds of target authors, which is constructed by the matrices from their novels, are not overlapped. In order to show the relationship of the text clouds, we calculate and analyze the distances between matrices. The

procedure of creating our model and analyzing the distances are described in more details below.

3 Model

In this model, we will generate transition matrices for novels written by Charles Dickens and Mark Twain that capture their writing styles. In addition, a centroid of each author will also be generated using the aggregated texts, which will be specified below. We will compare the distances of each author’s literature to their own centroid as well as distance between the two author’s writing.

3.1 Data

We will be using 10 novels written each author of our interest. Names of the literature are specified in the table below.

Charles Dickens	Mark Twain
A Tale of Two Cities	A Connecticut Yankee in King Arthur’s Court
Bleak House	Adventures of Huckleberry Finn
David Copperfield	Personal Recollections of Joan of Arc
Dombey and Son	Life on the Mississippi
Great Expectations	The Adventures of Tom Sawyer
Little Dorrit	The American Claimant
Olive Twist	The Gilded Age: A Tale of Today
Our mutual friend	The Mysterious Stranger
The Life And Adventures Of Nicholas Nickleby	The Prince and the Pauper
The Pickwick papers	The Tragedy of Pudd’nhead Wilson

The text files of these novels can be found and downloaded from Project Gutenberg. These data files are pre-processed by deleting all the text related to Project Gutenberg so that only writing by the author is considered when generating the transition matrix.

3.2 Model Assumptions

3.2.1 Transition matrix

When constructing the transition matrix, we assumed that an author’s writing style can be characterized by the sequence of individual letter and space used in his/her writing. Thus, we will be considering the probability of one character following another character in the transition matrix of our model.

3.2.2 Syntactic elements in text pre-processing

Regarding syntactic elements, we assume the following when pre-processing the text files.

- Difference between upper and lower case letters are ignored. Upper case letters are generally used as start of the sentence or proper nouns, and we consider changing the upper cases to lower cases has no effect on individual author’s style.
- Apostrophes are removed. We believe it usually does not separate two meaningful word. For example, it is not meaningful to separate the word “don’t” into “don” and ”t”.
- Punctuation is consider to have the same function as space, which is to separate words.
- A space followed by a space has no useful information on apprehension of author’s writing style. So consecutive spaces are treated as a single space in the texts.

By employing the above assumptions, our model will only be focused on authors’ writing in terms the way they use English words. More specifically, we believe author’s writing feature can be captured by the sequences of letters and spaces used in their literature. We do realize that by omitting the difference in capitalized and lower case letters, we might have a frequently occurred sequence of the book character’s name. We will be caution with this and will later address in result interpretation and alter in adjustment.

To implement the assumptions above, all alphabetical letters are converted to lower case after the above process. Characters other than the alphabetical letters(i.e. punctuation) are replaced with space character, with the exception of apostrophe, which is removed. All consecutive spaces are collapsed into one single space. We will pre-process all texts used by this method before generating our transition matrices.

3.2.3 Measure of distance

There are many different ways to compare how similar two authors' writings are by determining the "difference" between two transition matrices. One of those is to measure the distance between two transition matrices created from the writings of each author. In our model, we will calculate the Euclidean distance between two transition matrices (i.e. Frobenius or 2 norm of the difference of matrices). The formula to calculate the distance is:

$$\text{dis}(A, B) := \|A - B\|_F = \sqrt{\sum_{i=1}^n \sum_{j=1}^n |a_{ij} - b_{ij}|^2}$$

where A, B are 2 transition matrices with n states under comparisons, and a, b are entries of matrices. Using this distance measure, two identical matrices will have a distance of 0; matrices with different entries will have a distance greater than 0, and the more different the entries, the greater the distance.

3.2.4 Measure of author metric space(centroid and cluster)

To compare the similarity of two authors, we will define a centroid as the profile of the author. The centroid is a $\mathbb{R}^{27 \times 27}$ transition matrix generated from aggregating all 10 texts from the author(i.e. appending all 10 texts together). The cloud of an author would be interpreted by calculating the distances of 10 transition matrices of individual novels to the author's centroid. The distance between two authors will be the Euclidean distance of two centroids of the authors.

3.3 Generate Transition Matrices

By our assumptions, the transition matrix of our model will be in $\mathbb{R}^{27 \times 27}$, with the first 26 states corresponding to the 26 letters in English alphabet, and the last state representing space that separates the characters. We generate the transition matrix using Python. We iterate through the entire text to count the frequencies of letter α following letter β , where α and β are the 27 English letters and space character. The probability of letter α following letter β is then calculated by dividing the frequency of letter α following letter β by the total frequency of letter β occurs.

The matrix \mathbf{P}_{tm27} can be represented as:

$$\mathbf{P}_{\text{tm27}} = \begin{bmatrix} \mathbb{P}(a|a) & \mathbb{P}(b|a) & \dots & \mathbb{P}(z|a) & \mathbb{P}(*|a) \\ \mathbb{P}(a|b) & \mathbb{P}(b|b) & \dots & \mathbb{P}(z|b) & \mathbb{P}(*|b) \\ \vdots & & & & \\ \mathbb{P}(a|z) & \mathbb{P}(b|z) & \dots & \mathbb{P}(z|z) & \mathbb{P}(*|z) \\ \mathbb{P}(a|*) & \mathbb{P}(b|*) & \dots & \mathbb{P}(z|*) & \mathbb{P}(*|*) \end{bmatrix}$$

Where $\mathbb{P}(a|b)$ denotes the probability of letter “a” following by a letter “b” in the given text, etc. The character “*” in the matrix represents the space character.

4 Results

Using the method described above, we generate 10 transition matrices for the novels of each author, and 1 transition matrix as the profile matrix (i.e. centroid) of each author. Python 3.8 is used to construct the model, and code can be found in the appendix section. It takes approximately 1 minute to calculate all distances on macOS system with processor 2.4 GHz Quad-Core Intel Core i5. By calculating the distances of each transition matrices to the authors’ centroids, we can analyze the “cloud” information of each author indicated in the table below.

For Charles Dickens and his 10 novels, the distances from his novels to his centroid and cross distances of his novels to Mark Twain’s centroid are shown in the table below.

Charles Dickens' Novels	Distance to Centroid	Distance to Twain's Centroid
A Tale of Two Cities	0.3646	0.4526
Bleak House	0.2511	0.4212
David Copperfield	0.2416	0.3723
Dombey and Son	0.3304	0.4780
Great Expectations	0.3721	0.4796
Little Dorrit	0.2710	0.4357
Oliver Twist	0.3789	0.5340
Our mutual friend	0.4407	0.6524
The Life And Adventures Of Nicholas Nickleby	0.2519	0.4071
The Pickwick papers	0.3772	0.4724

From the table above, we find that all 10 books from Charles Dickens are closer to his own profile than to Mark Twain's profile. The average distance of the 10 novels to Charles Dickens' centroid is approximately 0.3274 with standard deviation of 0.0655. To visualize how the distances are distributed, we plot a histogram of the distances of Charles Dickens' novels to his centroid, with a blue vertical line indicating the average distances.

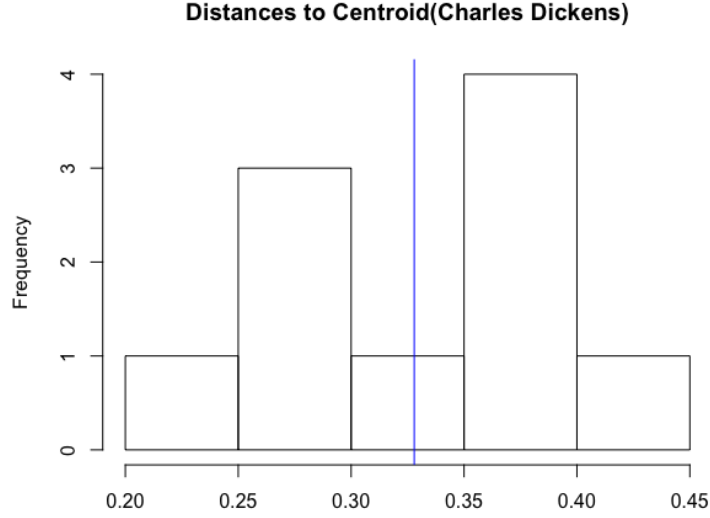


Figure 1: Histogram of Euclidean distances of Dickens' novels to own centroid in original model.

We can see that distribution of distances is quite symmetry around the the average distance from centroid, with two peaks close to the mean and low frequencies at two tails. The thin tails of histogram imply that we have 1 novel really close to the author's profile and 1 novel really different from Dickens' profile. If we imagine the transition matrices of the author lying in a higher than 3 dimensional space, it would be a cloud of matrices where most of the matrices would lie with a distance around 0.25 to 0.4 (i.e. roughly the average distance to centroid) around the centroid, with just a couple matrices much closer or much further away from the centroid.

Similarly, the distances from Mark Twain's novels to his own profile and cross distances of his novels to Charles Dickens' profile are shown in following table.

Mark Twain's Novels	Distance to Centroid	Distance to Dickens' Centroid
A Connecticut Yankee in King Arthur's Court	0.2466	0.4514
Adventures of Huckleberry Finn	0.7091	0.9068
Personal Recollections of Joan of Arc	0.4697	0.5084
Life on the Mississippi	0.2173	0.4241
The Adventures of Tom Sawyer	0.2820	0.5176
The American Claimant	0.2895	0.4600
The Gilded Age: A Tale of Today	0.2648	0.4208
The Mysterious Stranger	0.3999	0.5800
The Prince and the Pauper	0.4282	0.4207
The Tragedy of Pudd'nhead Wilson	0.4677	0.6441

Unlike Charles Dickens, most of Mark Twain's novels have a smaller distance to own centroid than cross distance to Twain's centroid with the exception of *The Prince and the Pauper*. The average distance of the 10 novels to Mark Twain's centroid is approximately 0.3775 with a standard deviation of 0.1421, and we can visualize the distribution of distances by plotting a histogram.

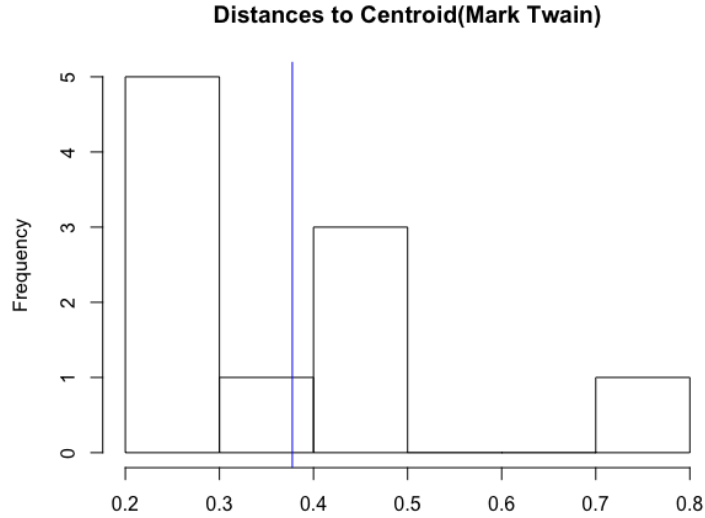


Figure 2: Histogram of Euclidean distances of Twain's novels to own centroid in original model.

Unlike Charles Dickens' distance distribution, the distances of Mark Twain's texts to his centroid seems to be skewed left, with more novels similar to his profile and 1 novel really different from the profile. This can be represented by cloud of matrices with most of the matrices lied with a distance of 0.2 - 0.5 to the centroid, and 1 matrix almost twice the distance away from the centroid. We also notice that the standard deviation of the distances is greater than that of Charles Dickens, and this could imply a greater variation in Mark Twain's writing styles. The majority of that variation can possibly explained by the novel with furthest distance to centroid.

From the table of distance above, we observe that the novel with furthest distance(however not neccesarily an outlier) to Mark Twain's profile is *Adventures of Huckleberry Finn*, which is surprising because many people consider it as his most iconic work. We then take a closer look into the transition matrix of *Adventures of Huckleberry Finn* and compare it to the Mark Twain's centroid. A possible explanation is the effect of special nouns in this book, such as character names. In the adjustment section, we remove the proper nouns to see if this explanation is reasonable. To further investigate this, we generate the difference matrix by taking their difference: $A - B$, where A is the transition matrix for *Adventures of Huckleberry Finn* and B is the transition matrix for Twain's centroid. A color coded matrix is shown in the following figure.

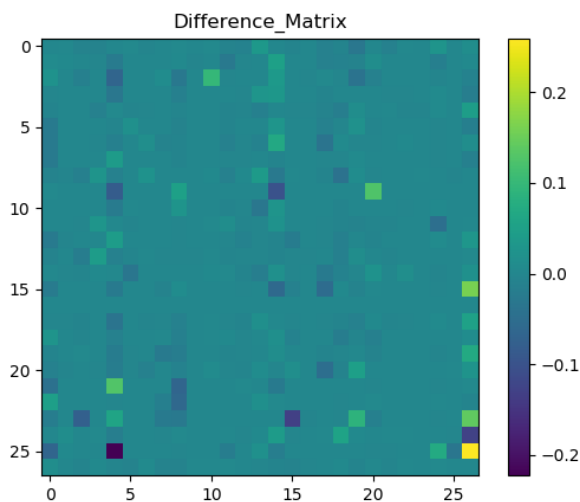


Figure 3: Difference matrix of *Adventures of Huckleberry Finn* and Mark Twain's centroid.

As we can observe that the largest difference occurs in row 26 and column 27 of the difference matrix(index in figured is shifted by 1 due to programming technicality). This corresponds to the sequence "z[space]", which means a significant difference in the use of words ending with letter "z" in *Adventures of Huckleberry Finn*. We look up main characters name in *Adventures of Huckleberry Finn*, and none of those ends with letter "z". Thus, proper noun does not seem to be a factor that cause the large distance from centroid. This imply that *Adventures of Huckleberry Finn* has relatively distinct writing style compare to Twain's other literature.

We also calculate the Euclidean distance between the two centroids of Charles Dickens and Mark Twain, which is 0.3595. This distance between two authors is greater than 0.3274, which is Charles Dickens' cloud "radius", but smaller than Mark Twain's cloud "radius" of 0.3775, where "radius" is the average distance of novels to own author's centroid. Since the distance between two authors' centroids does not seem significantly larger than the authors' clouds "radius", we cannot make a conclusion based on this information. However, we can plot a comparative box-plot of novel's distances to own author's centroid and their cross distances to the other author's centroid and observe if there exist any significant difference.

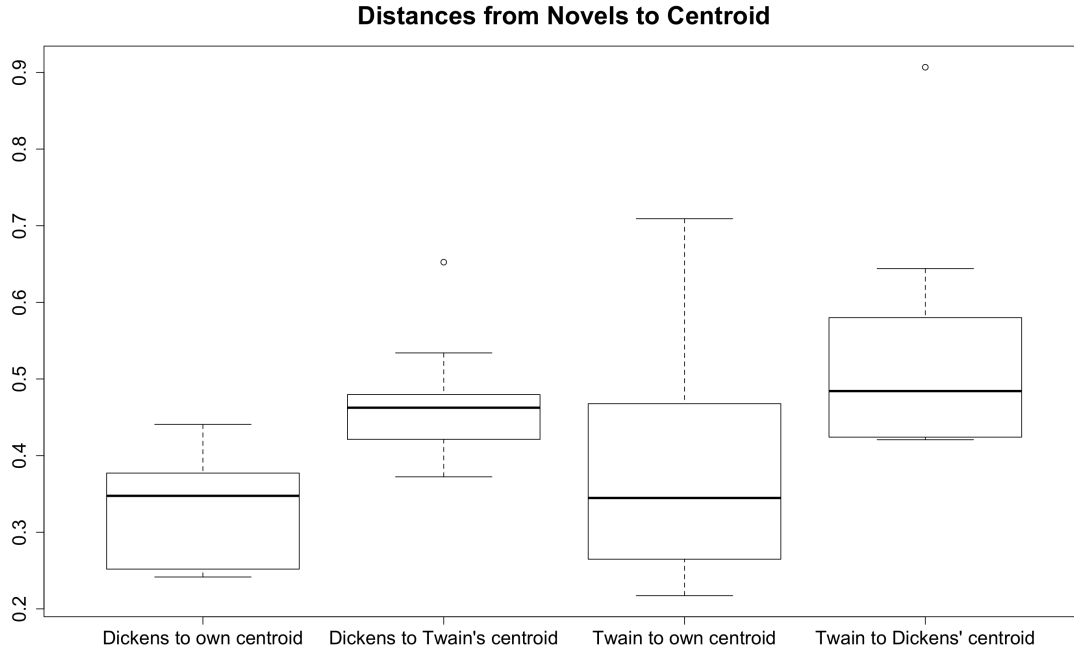


Figure 4: Comparative box-plot of novel's distances to own author's centroid and their cross distances to the other author's centroid in original model.

The left two boxes show the distances from Charles Dickens' 10 novels to his centroid and to Mark Twain's centroid. We conclude that these 10 texts are closer to Charles Dickens's centroid since the second box is higher than the first one without overlapping. The spreads of these two groups of data have big difference. The second box is much smaller than the first box means the variation of the cross distances of Charles Dickens' novels is small. We could make a guess that the ten books are spread around Charles Dickens's centroid but all of them are far from Mark Twain's centroid. Thus, we say most of Charles Dickens's ten texts were written in his own writing style compared to Mark Twain's styles. We can see that there is overlap between cross distances and distances of Mark Twain's texts, meaning it is hard for us to make a clear conclusion about his cloud, even though the median distance from Mark Twain's texts to Charles Dickens' centroid is greater than these texts to their own author's centroid. Thus, because of the overlap of the box areas, we do not have enough evidence to conclude that the two authors have their own clusters

without overlapping with each other. One thing to notice is that both cross distances box-plots have extreme high values. Twain’s *Out Mutual Friend* is much further away from Dickens’ profile, and Dickens’ *Adventures of Huckleberry Finn* is exceptionally further from Twain’s profile. This could potentially indicate that those two novels have very distinct style to another author’s writings. However, we do aware that it could potentially be the effect of repeating proper noun, and we will adjust the model later to see if these outliers still present.

5 Adjustments and Extensions

In the previous section, we create transition matrices to represent texts and analyze authors’ writing styles. However since the matrices live in a rather high dimensional space, we cannot easily conclude that based on the distances we have calculated so far. The distances used in the previous model is calculated based on one of the ways to compare matrices. Some interesting results occurs in original model results due to possible effects of special nouns or just difference in writing styles. To study how model results changes as the above condition changes, we make several adjustment and extensions in this section.

5.1 Omitting Proper Nouns

In our original model, we convert the capitalized letters to lowercase so that we don’t lose any written information(i.e. all letters and space) from the author. However, characters’ name or names of the locations in a particular novel usually have uncommonly used letter sequences. This could have noticeable effect when generating character by character transition matrix, which could potentially bias the results. For example, if we change the main character’s name from *Adventures of Huckleberry Finn* from “Huck” to “Jack”, the distance between these two texts will not be zero, while syntactically they are identical.

To prevent this situation, we can remove all repeated proper nouns designated to a particular novel. However, the implementation to remove just those proper noun is complicated. For the sake of simplicity, we will remove all the words with capitalized first letter when pre-processing the texts files. We are aware of the fact that words at the start of sentences would also be removed. Since the novels we choose are long texts, the rest of the texts will still capture enough

syntactic information for us to analyze with capitalize words removed. We then analyze the two clouds of authors with the new adjustment of ignoring proper noun.

The distances of Charles Dickens to each author's centroid with new assumption is shown in the table below.

Charles Dickens' Novels	Distance to Centroid	Distance to Twain's Centroid
A Tale of Two Cities	0.2030	0.3670
Bleak House	0.1352	0.3939
David Copperfield	0.1737	0.4002
Dombey and Son	0.1330	0.4567
Great Expectations	0.1783	0.3609
Little Dorrit	0.1465	0.4125
Oliver Twist	0.2098	0.4063
Our mutual friend	0.1673	0.4556
The Life And Adventures Of Nicholas Nickleby	0.1768	0.4233
The Pickwick papers	0.2234	0.3430

We can see all novels are closer to Dickens' own centroid than to Twain's, which is consistent with original model result. We also visualize the distribution of distances in a histogram.

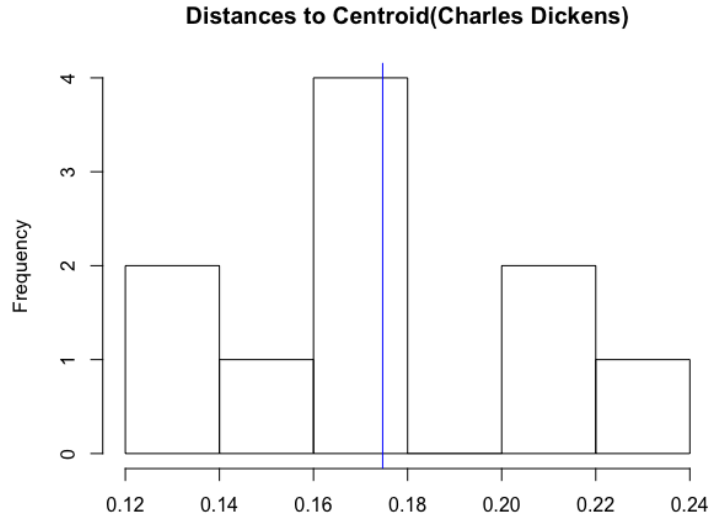


Figure 5: Histogram of Euclidean distances of Dickens' novels to own centroid in Adjustment 1(omit special noun).

The distribution of distances also seems symmetric and roughly normal. The new average distance of Dickens' novels to the centroid is 0.1746, which is much smaller than 0.3274, the average distance when capitalized words are not removed. This indicates that the clouds of Charles Dickens' novels to his profile is less spread when the proper nouns are not considered as his writing style. The new standard deviation of distances is 0.0293, which roughly reduces the original standard deviation 0.0655 by half. This implies that when we exclude the affect of proper nouns in Charles Dickens' writing, the variation on writing is much lower.

For Mark Twain, the table of distances and histogram are shown below.

Mark Twain's Novels	Distance to Centroid	Distance to Dickens' Centroid
A Connecticut Yankee in King Arthur's Court	0.1476	0.3637
Adventures of Huckleberry Finn	0.6298	0.8827
Personal Recollections of Joan of Arc	0.1735	0.4013
Life on the Mississippi	0.1829	0.3701
The Adventures of Tom Sawyer	0.2152	0.4945
The American Claimant	0.2338	0.3833
The Gilded Age: A Tale of Today	0.1613	0.3137
The Mysterious Stranger	0.2440	0.4858
The Prince and the Pauper	0.3972	0.3719
The Tragedy of Pudd'nhead Wilson	0.2188	0.3471

After we ignoring the special nouns, most novels from Mark Twain are still closer to his own profile except for *The Prince and the Pauper*. Since special noun is no longer a factor for the difference, this imply that the writing style of *The Prince and the Pauper* is quite ambiguous between Dickens and Twain.

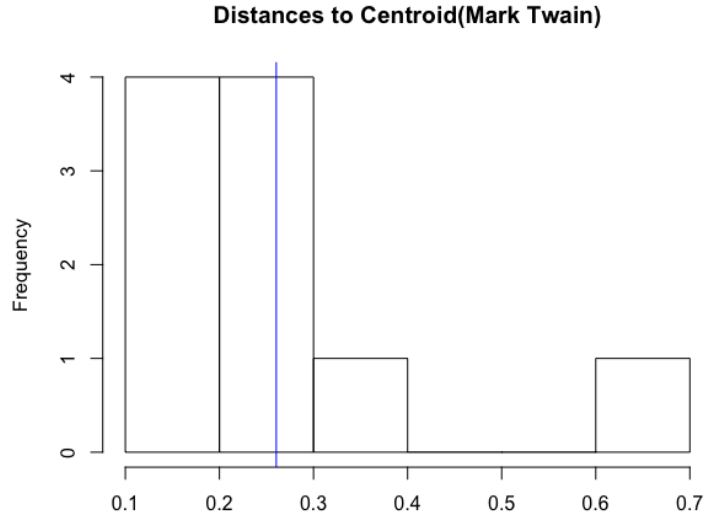


Figure 6: Histogram of Euclidean distances of Twain's novels to own centroid in Adjustment 1(omit special noun).

We can see that distribution of Mark Twain's novels to his centroid is still skewed left as in original model, but with a mean distance 0.2604 smaller than the original mean of 0.3775. This also indicates that removing the special nouns makes Twain's cloud of novels closer to the centroid. The new standard deviation is 0.1400, which is slightly less than 0.1421 in the original model. This is a potential indicator that special nouns in Twain's literature is not a significant source of his variations in writing style, which is consistent with original model result.

To better compare the two authors clusters, a comparative box-plot of the distances and cross distances is shown below.

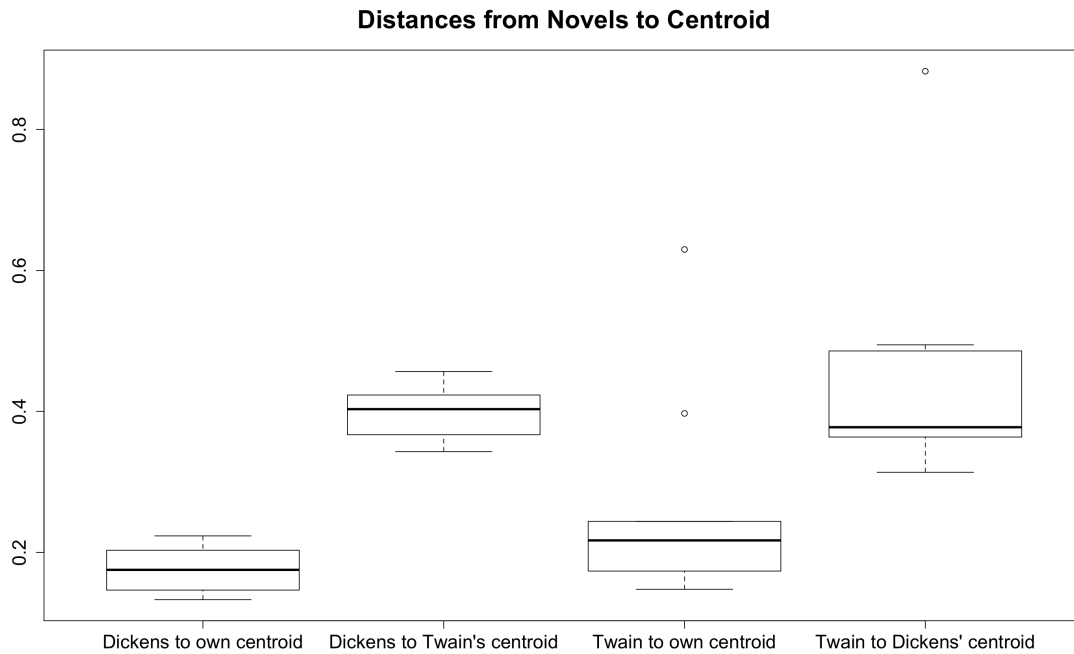


Figure 7: Comparative box-plot of novel's distances to own author's centroid and their cross distances to the other author's centroid in Adjustment 1(omit special noun).

After ommiting special nouns, the box-plots of each author's books are much less overlapped,

which indicates that ignoring proper nouns better separates the two authors works and now all books are closer to the center of their writers. The two boxes of cross distances are much higher than the two boxes of distances to author’s own centroid, with a couple exceptions. This suggest that the two authors novels are very likely to form their own cluster that are separated from the other author. Moreover, all four boxes are smaller than those in the original model, which implies that removing the special nouns create less variations and make the texts closer to each other. The changes allow us to be able to avoid meaningless effects caused by special nouns on transition matrices, and thus reveal two dense clouds that are relatively away from each other.

Two outlier distances in Mark Twain novels to own centroid, which in decreasing distance order are *Adventures of Huckleberry Finn* and *The Prince and the Pauper*. The existence of *Adventures of Huckleberry Finn* as an outlier is consistent with our conclusion in original model results, where we found out the entry with largest difference is not one of the repeated character names. We look up the novel *Adventures of Huckleberry Finn* and find out that one special characteristic in this novel is Mark Twain’s heavy use of vernacular. This novel is important in the history of American literature since it was the first time introducing dialect in public reading. In order to show the background of a character, Huck, Mark Twain used a number of dialects in his words, which include: “the Missouri “Negro” dialect, the backwoods Southwestern dialect, the Pike County dialect, and a few others”. Moreover, contractions are also frequently used in this novel, such as ain’t, they’d, you’ve, wouldn’t, they’ll, it’, I’ll, etc. At that time, using contractions in books was not a popular writing style.[13] We also find many use of the phrase “wuz” and “uz” in the writing, which is rarely seem in English writing. Thus, the large distance between *Adventures of Huckleberry Finn* and the centroid can be explained by the specially use of vernacular. *The Prince and the Pauper* was also a outlier but less far away from the centroid might be due to the fact that it is written in 1881, which is recent to 1884 when *Adventures of Huckleberry Finn* first published, so they might have similar use of language.

In sum, by omitting proper noun better the two authors novel form a denser cluster with less overlap than when proper noun is include. *Adventures of Huckleberry Finn* has distance significantly far away from Twain’s profile, and it reflectse distinct writing style of using special dialect for this novel.

5.2 Two-letter by Two-letter Transition Matrix

We suspect that one letter to one letter sequence might not carry enough syntactic information of individual author, so the difference between two matrices is relatively small. With a goal to capture more syntactic elements of the author, we now build our transition matrix based on the probability of two-letter sequence following by another two-letter sequence. In this case, we will have $27^2 = 729$ states indicating the 2-character permutations from the 26 English letters and a space character. Matrix generated by this method is referred as the second order transition matrix.

The second order transition matrix is a $\mathbb{R}^{729 \times 729}$ matrix with states of the matrix being a list of ordered tuple produce by the cross product from the states of the \mathbf{P}_{tm27} matrix. For example, the first 4 states in the markov chain is: $aa, ab, ac, ad...$ Here is how the \mathbf{P}_{2ndtm} looks like when written out:

$$\mathbf{P}_{2ndtm} = \begin{bmatrix} \mathbb{P}(aa|aa) & \mathbb{P}(ab|aa) & \dots & \mathbb{P}(az|aa) & \mathbb{P}(**|aa) \\ \mathbb{P}(aa|ab) & \mathbb{P}(ab|ab) & \dots & \mathbb{P}(az|ab) & \mathbb{P}(**|ab) \\ \vdots & & & & \\ \mathbb{P}(aa|*z) & \mathbb{P}(ab|*z) & \dots & \mathbb{P}(az|*z) & \mathbb{P}(**|*z) \\ \mathbb{P}(aa|**) & \mathbb{P}(ab|**) & \dots & \mathbb{P}(az|**) & \mathbb{P}(**|**) \end{bmatrix}$$

For this type of transition matrix, there could be some rows that are entirely zeros if a given 2 characters sequence is never observed in the written text from the author. This is absolutely true for the last row of the \mathbf{P}_{2ndtm} due to the way text are pre-processed. Two consecutive spaces never appear in the text after pre-processing. While having a row consists of entirely zeros breaks the definition of a transition matrix, it keeps the size of transition matrices consistent for all the written works from different authors, making the comparison using matrix norm possible for all authors.

Then, we combine the Adjustment 5.1, which is to ignore the proper nouns in the novels with this new method of generating transition matrices to analyze two clouds of the authors. The following tables shows the distances of two authors' novels to their own centroid and cross distances to

other author's centroid with the new adjustment.

Charles Dickens' Novels	Distance to Centoid	Distance to Twain's Centoid
A Tale of Two Cities	8.8439	9.9892
Bleak House	7.8239	10.0659
David Copperfield	7.5429	9.8883
Dombey and Son	7.5891	9.9286
Great Expectations	8.3164	9.9106
Little Dorrit	7.4535	9.6967
Oliver Twist	8.4592	9.9044
Our mutual friend	7.8250	9.6872
The Life And Adventures Of Nicholas Nickleby	7.3731	10.0668
The Pickwick papers	7.4124	10.3332

Mark Twain's Novels	Distance to Centoid	Distance to Dickens' Centoid
A Connecticut Yankee in King Arthur's Court	8.1742	10.6459
Adventures of Huckleberry Finn	10.5916	12.7403
Personal Recollections of Joan of Arc	8.9146	10.3789
Life on the Mississippi	7.8095	10.4310
The Adventures of Tom Sawyer	8.8183	10.7295
The American Claimant	9.0189	10.7002
The Gilded Age: A Tale of Today	8.2194	9.8532
The Mysterious Stranger	9.9974	11.0395
The Prince and the Pauper	9.5853	10.4050
The Tragedy of Pudd'nhead Wilson	9.4333	10.8971

The comparative box-plot of the distances is shown below.

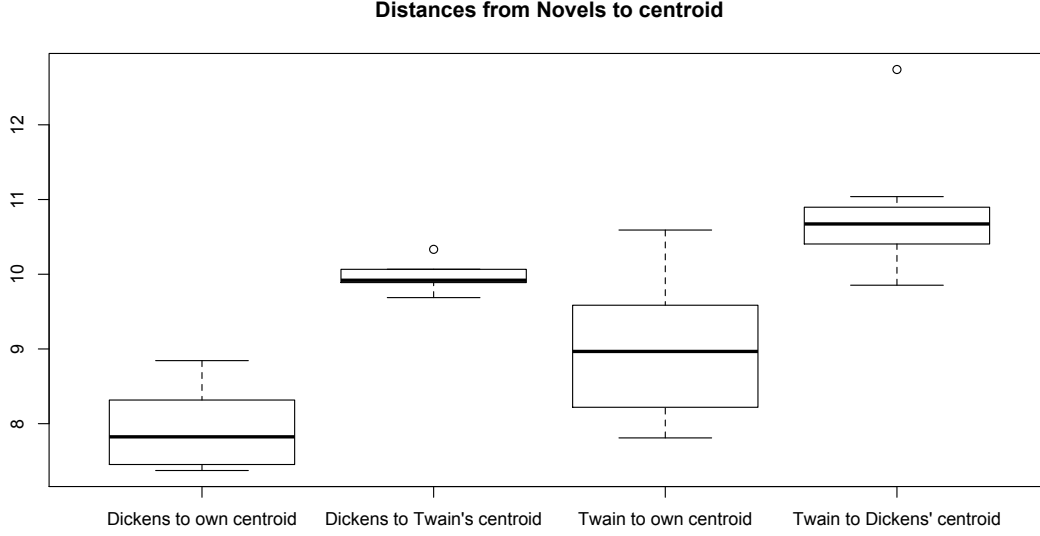


Figure 8: Comparative box-plot of novel's distances to own author's centroid and their cross distances to the other author's centroid in Adjustment 2(second order matrix omitting special noun).

The first noticeable point is that the distances from each author's novels to his centroid is significantly greater than the distances in the previous models, the Charles Dickens' and Mark Twain's mean distances to their own centroids are 7.8639 and 9.0562 respectively. This is because we summing differences between two $\mathbb{R}^{729 \times 729}$ matrices, which is has much more entries than the $\mathbb{R}^{27 \times 27}$ matrices in the previous models. Two letter by two letter sequence is also more restrictive on the letter sequence and thus captures more syntactic information from writings, which could potentially reflect the different writing style better.

Secondly, Figure 8 shows that the distances between each author's novels to his own centroid is smaller than the distances between his novels to the other author's centroid. The distance between two centroids in this model is 9.8827, which is larger than the average distances from Charles Dickens' and Mark Twain's novels to their own centroids, which are 7.8639 and 9.0562. This indicates that each author are likely to have his own cluster that are separated from the

other author. The 2 box-plots of cross distances being higher than distances to own centroid with little overlap also support this conclusion.

A minor fact to notice is that all Twain’s novels are now closer to his own centroid than to Dickens’, even *The Prince and the Pauper* which was ambiguous in previous models. The outlier of *Adventures of Huckleberry Finn* is still present, and it is shown to be significantly far away from Dickens’ centroid. This is also consistent with Twain’s use of dialect and idioms, which is not an characteristic of Dickens’ writing.

5.3 Other Distance Measures

In all of the models above, Euclidean distance(i.e. 2-Norm of matrix) are used as the measure of distance between two matrices. In this section, we want to investigate how author’s clusters will change under different distance measure: matrix one norm and entry-wise one norm.

5.3.1 One Norm

To use matrix one norm, we first calculate the difference matrix of probability of a two letter sequence following by another two letter sequence. We then sum up the probability difference in absolute value of each column, and use the maximum column sum as the distance measure. The formula to calculate one norm is

$$\text{dis}(A, B) := \|A - B\|_1 = \max_{1 \leq j \leq n} \sum_{i=1}^n |a_{ij} - b_{ij}|$$

where A, B are two second-order transition matrices with $n = 729$ states under comparisons, and a, b are entries of matrices.

The intuition behind using one norm is that if there is a significant difference in a specific sequence of words used, it will be reflected in distance without being weighted out by other entries, especially when differences are less than 1. For example, if one novel use excessive amount of the word "quad" while the other use none, we should expect transitioning entry of "qu" to "ad" to be significantly large in the difference matrix, which will very likely to make the column with that entry to have maximum sum. Following adjustment 2, we calculate the one norm distances

using the two letter by two letter transition matrix for two authors, and results are shown in tables below.

Charles Dickens' Novels	Distance to Centoid	Distance to Twain's Centoid
A Tale of Two Cities	7.1986	9.0822
Bleak House	6.7416	9.1862
David Copperfield	6.3498	10.2326
Dombey and Son	6.4366	8.8756
Great Expectations	7.4247	8.8523
Little Dorrit	5.9755	8.3802
Oliver Twist	7.6053	10.1967
Our mutual friend	5.8113	8.5028
The Life And Adventures Of Nicholas Nickleby	6.5210	10.5184
The Pickwick papers	6.1336	10.7013

Mark Twain's Novels	Distance to Centoid	Distance to Dickens' Centoid
A Connecticut Yankee in King Arthur's Court	6.1060	10.3811
Adventures of Huckleberry Finn	9.9373	13.4310
Personal Recollections of Joan of Arc	7.7319	9.0927
Life on the Mississippi	5.2345	9.5424
The Adventures of Tom Sawyer	8.2151	10.4789
The American Claimant	7.9103	11.1265
The Gilded Age: A Tale of Today	6.3769	9.5817
The Mysterious Stranger	8.1830	10.6664
The Prince and the Pauper	8.8797	11.0327
The Tragedy of Pudd'nhead Wilson	8.5975	11.3223

All novels are closer to their author's centroid and further away from another author's centroid, and this is consistent throughout all the models so far. This again is an evidence that Charles Dickens and Mark Twain's novels form their own clusters. The comparative box-plots of one norm distances is shown below.

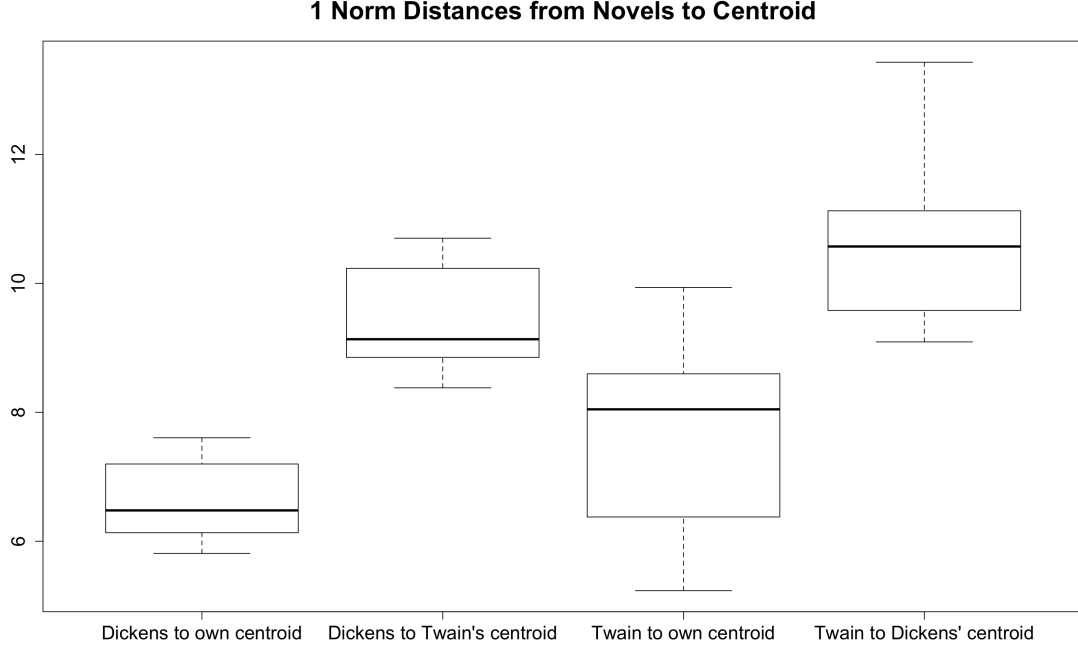


Figure 9: Comparative box-plot of novel's one norm distances and cross distances in Adjustment 3.1(second order matrix omitting special noun).

Similar to previous model, both author have their writings closer to their own centroid. The cross distances box-plots being higher than distances to own centroid allow us to make the conclusion that two clusters constructed by the texts from different authors are separate. Moreover, we observe the distances in this model have bigger spreads than previous models. This is because of the feature of one norm, which uses the absolute value while computing the results. As a result of increased spread, although we still have large distances for novels at the upper tail, the extreme values are no longer considered as outliers.

5.3.2 Entry-wise Matrix One Norm

To calculate entry-wise matrix one norm, we use the formula

$$\text{dis}(A, B) := ||\text{vec}(A - B)||_1 = \sum_{i=1}^n \sum_{j=1}^n |a_{ij} - b_{ij}|$$

The intuition of using this "entry-wise" one norm is that most of the probabilities are less than 1 and so are their differences in absolute value. Using two norm in previous models, we square individual differences which makes them even smaller. Distances calculated by entry-wise one norm and the comparative box-plot are shown below.

Charles Dickens' Novels	Distance to Centroid	Distance to Twain's Centroid
A Tale of Two Cities	272.7793	341.2822
Bleak House	213.3771	334.3947
David Copperfield	202.8842	335.7517
Dombey and Son	200.9683	339.2411
Great Expectations	247.7411	335.9497
Little Dorrit	198.0833	323.7918
Oliver Twist	258.7954	348.5776
Our mutual friend	209.2619	326.3303
The Life And Adventures Of Nicholas Nickleby	199.5594	343.6294
The Pickwick papers	223.1247	363.6674

Mark Twain's Novels	Distance to Centroid	Distance to Dickens' Centroid
A Connecticut Yankee in King Arthur's Court	250.8516	373.8739
Adventures of Huckleberry Finn	402.8773	511.7869
Personal Recollections of Joan of Arc	276.2573	366.0022
Life on the Mississippi	251.0494	366.0901
The Adventures of Tom Sawyer	304.5456	390.6068
The American Claimant	306.2761	381.4775
The Gilded Age: A Tale of Today	253.7137	324.3947
The Mysterious Stranger	346.5795	406.7008
The Prince and the Pauper	329.9942	375.6150
The Tragedy of Pudd'nhead Wilson	319.8175	396.9573

Using entry-wise one norm, the distance measures are significantly larger than any of the distance measure in previous models. All novels are closer to own author's centroid than others, which

is consistent with previous model. A comparative box-plot of distances is shown in the following figure.

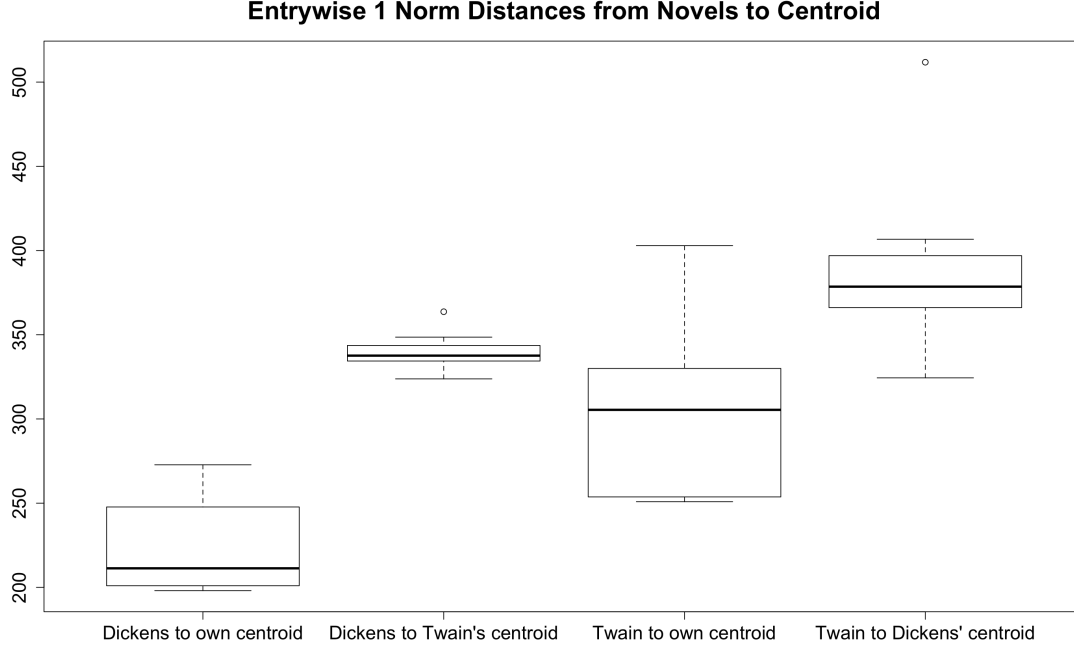


Figure 10: Comparative box-plot of novel's entry-wise one norm distances and cross distances in Adjustment 3.2(second order matrix omitting special noun).

Similar with the previous models, the box-plots supports that there are two separate clouds, since the distances of texts to their authors' profile is much less than the cross distances. We also observe that the cross distances are generally higher than distances from novels to own centroid while having a much smaller spread, especially Dickens' novels to Twain's centroid. The smaller spread of cross distances is a result of increasing magnitude of distance value. Under entry-wise one norm formula, all distances to centroid are in larger value and so is the distance between two centroids. Imagine in a high dimensional space($\mathbb{R}^{729 \times 729}$), it would be two clouds of matrices closely located around a centroid but clouds are far away from each other. The distances within each cloud are relatively small while the two clouds are far away, so the cross distances from one cloud to the other's centroid be large but with small spread.

5.4 Log Scaled Transition Matrix Accounting for Uncommon Sequence

We are aware that some sequence of English letters are very rarely used, such as “az” or “xe”, or even the case of “wuz” mentioned earlier in *Adventures of Huckleberry Finn*. We think the frequencies of using these rare words might be an important factor of an author’s writing style. In previous model, we fail to capture this factor since the probability of these transitions are very small. As a result, it does not carry much weight when we calculated the distance between two matrices(i.e. norm of difference matrix).

In order to make our matrices having an emphasis on the use of rare words, we will take the \log of the count of one character following another character before calculating the probability. Rarely used sequence with smaller counts would now weight a little more than when without taking \log . We are aware that there are sequences that might not be used at all by the authors or do not form a word in English, which would have a count of 0. In that case, we keep the count of those sequences as 0 instead of taking \log , since $\log(0)$ is undefined. We generate the log scaled second order transition matrices for each author(proper noun omitted), and the Euclidean distances are shown in the following table.

Charles Dickens’ Novels	Distance to Centoid	Distance to Twain’s Centoid
A Tale of Two Cities	9.5910	9.5923
Bleak House	8.1064	9.3101
David Copperfield	8.4421	9.3244
Dombey and Son	8.0977	8.7139
Great Expectations	9.2536	9.5502
Little Dorrit	8.5623	8.8847
Oliver Twist	9.3095	9.4043
Our mutual friend	8.2005	9.0116
The Life And Adventures Of Nicholas Nickleby	8.1721	9.0656
The Pickwick papers	7.9420	9.2946

Mark Twain's Novels	Distance to Centroid	Distance to Dickens' Centroid
A Connecticut Yankee in King Arthur's Court	8.7710	10.1071
Adventures of Huckleberry Finn	10.1208	11.5048
Personal Recollections of Joan of Arc	8.7629	10.0970
Life on the Mississippi	8.0998	9.7903
The Adventures of Tom Sawyer	9.5983	10.8898
The American Claimant	9.8571	11.0426
The Gilded Age: A Tale of Today	8.6106	9.8611
The Mysterious Stranger	10.6961	11.2969
The Prince and the Pauper	9.7923	10.6779
The Tragedy of Pudd'nhead Wilson	9.9615	10.8796

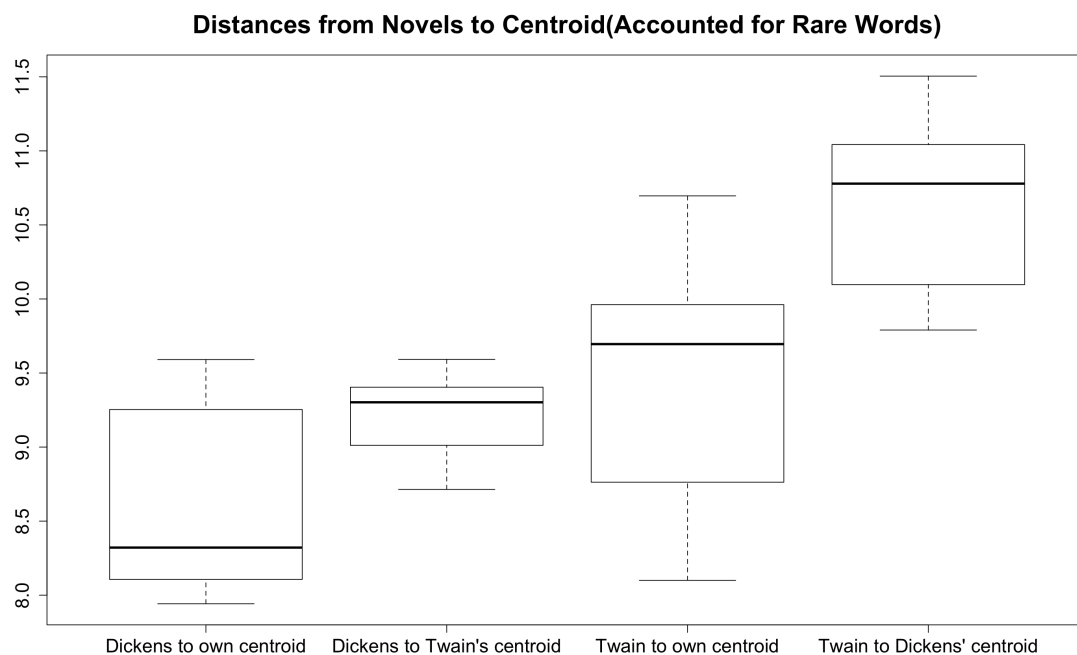


Figure 11: Comparative box-plot of novel's Euclidean distances and cross distances in Adjustment 4 (log scaled second order matrix omitting special noun).

From figure 11, we can see that the pair box-plot of Dickens novels to itself and the cross distance to Twain’s centroid is very overlapped, meaning we cannot conclude that Dickens have a cluster for his profile even though the average cross distance is much higher than average distance to his own centroid. For Mark Twain, there is little overlap on the two box-plot on the right. Cross distances of Twain’s novels to Dickens’ centroid are shown to be greater than distances to his own centroid, meaning Twain is very likely to have his small cluster of literature relatively far from Dickens’ profile.

Comparing the box-plot to those from previous adjusted models, log scaling the counts to emphasize rarely used sequence does not seem to be very effective in defining clouds of authors that are separated, especially for Dickens novels in this case. A possible reason could be that those rarely used words are just not an important factor of the author’s writing style. For example, if the word "za" is used a couple times in one of Dickens’ writing, it does not necessarily imply using "za" should be consider as a characteristic of his writing. Thus, it is not always reasonable to emphasize the use of rarely used sequence. It might work for author who are known for the use of uncommon diction to generate better clusters with smaller "radius" and far from other cluster. In our case, rarely used sequence does not seem to play an important role in Dickens and Twain’s writing.

6 Conclusion

After the above analysis on transition matrices generated from two authors novels, we can conclude that using Markov transition matrix can captures syntactic elements or writing style of authors. We do have evidence showing that transition matrices generated from novels of Charles Dickens and Mark Twain form two distinct clusters as we use different distance measure and impose more restrictive assumption such as omitting proper noun or using second order transition matrix. This is consistent with real world phenomenon, because the two authors we have chosen are generally known to have different writing style.

When defining transition matrix using probability of letter by letter sequence and with proper noun intact, we cannot conclude two authors have distinct cluster separated from each other since

there is overlap in Twain's cross distances to other centroid with distances to own centroid. After omitting proper noun, overlap is much more reduced which implies two denser cluster of authors' novel relatively away from each other.

As we move to second order transition matrix define by two letter by two letter probabilities, distance values are significantly greater. No overlap between cross distances and distances to own centroid is observed and the spread of distances seem to be reduced, which further supports the conclusion on two separate clouds.

When we use one norm and entry-wise one norm measure used for distance comparison, same conclusion on two separate cluster still holds. However, we observe the spread of one norm cross distances being slighter larger than Euclidean distance, and the spread of entry-wise one norm being smaller than Euclidean distance. This imply that using entry-wise one norm seem to define two distinct cluster better for the two chosen authors.

Using log scaled transition matrices for the intention to emphasize use of rare words does not support previous conclusion of separate clusters, since there is overlapped in Dickens' cross distances to other centroid and distances to own centroid. We believe that emphasizing rare words might not be reasonable to use for the two chosen author, since it is not one of the characteristic of their writing style.

An observation consistent though out all models is that the texts from Charles Dickens are always closer to his centroid with smaller spread than those of Mark Twain. This imply that Charles Dickens had a more stable writing style for his literature than Mark Twain. This can also be supported by the fact that *Adventures of Huckleberry Finn* occur as an outlier for several models. As mentioned previously, Twain's frequent use of special dialect and idioms is very distinct from both Dickens' and his own writings, which contribute to the large variation in Twain's novels. This specialty is not captured the the models using one norm distance measure and log scaled matrix.

There are many aspect of the model we can consider for future modifications. If the author have a

large profile of literature, aggregating all writings could be challenging. We can change the centroid measure to the average of individual transition matrices of author's writing. When comparing two authors distance to each other using just one value, we can use the distance between two closest matrices from two author instead of distance between centroid to give us a sense of how close can the two potential clouds be.

References

- [1] *Charles Dickens*. Charles Dickens - The Greatest Literature of All Time,
www.editoreric.com/greatlit/authors/Dickens.html
- [2] Dickens, Charles. *A Tale of Two Cities*. Project Gutenberg. N.p. N.d. EPUB. Downloaded Feb 25, 2020.
<https://www.gutenberg.org/ebooks/98>
- [3] Dickens, Charles. *Bleak House*. Project Gutenberg. N.p. N.d. EPUB. Downloaded Feb 25, 2020.
<http://www.gutenberg.org/ebooks/1023>
- [4] Dickens, Charles. *David Copperfield*. Project Gutenberg. N.p. N.d. EPUB. Downloaded Feb 25, 2020.
<https://www.gutenberg.org/ebooks/766>
- [5] Dickens, Charles. *Dombey and Son*. Project Gutenberg. N.p. N.d. EPUB. Downloaded Feb 25, 2020.
<https://www.gutenberg.org/ebooks/821>
- [6] Dickens, Charles. *Great Expectations*. Project Gutenberg. N.p. N.d. EPUB. Downloaded Feb 25, 2020.
<https://www.gutenberg.org/ebooks/1400>
- [7] Dickens, Charles. *Little Dorrit*. Project Gutenberg. N.p. N.d. EPUB. Downloaded Feb 25, 2020.
<http://www.gutenberg.org/ebooks/963>
- [8] Dickens, Charles. *Oliver Twist*. Project Gutenberg. N.p. N.d. EPUB. Downloaded Feb 25, 2020.
<https://www.gutenberg.org/ebooks/730>
- [9] Dickens, Charles. *Our mutual friend*. Project Gutenberg. N.p. N.d. EPUB. Downloaded Feb 25, 2020.
<http://www.gutenberg.org/ebooks/883>
- [10] Dickens, Charles. *The Life And Adventures Of Nicholas Nickleby*. Project Gutenberg. N.p. N.d. EPUB. Downloaded Feb 25, 2020.
<http://www.gutenberg.org/ebooks/967>

- [11] Dickens, Charles. *The Pickwick papers*. Project Gutenberg. N.p. N.d. EPUB. Downloaded Feb 25, 2020.
<http://www.gutenberg.org/ebooks/580>
- [12] Freeman Daniel, *Predicting the author of Twitter posts with Markov chain analysis..* (2017). Honors Theses.
<https://scholar.utc.edu/cgi/viewcontent.cgi?article=1126&context=honors-theses>
- [13] *Huck's Vernacular throughout The Adventures of Huckleberry Finn..* Living Language, 29 Apr. 2009,
livingwithlanguage.wordpress.com/2009/04/03/huck's-vernacular-throughout-the-adventures-of-huckleberry-finn/
- [14] Jagtap, Avinash S, and Jaya L Limbore. *Authorship Attribution Using Markov Chain..* International Journal of Latest Trends in Engineering and Technology. Vol.(9) Issue(3), pp.119-125.
<http://dx.doi.org/10.21172/1.93.21>
- [15] Pavels Osipovs, Andrejs Rinkevics, Galina Kuleshova and Arkady Borisov. *Markov Chains in the Task of Author's Writing Style Profile Construction..* Information Technology and Management Science Volume 17 Issue 1 (2014). Information Technology and Management Science. Sciendo, March 11, 2015.
<https://doi.org/10.1515/itms-2014-0018>
- [16] Sanderson, Conrad, and Simon Guenter. *Short Text Authorship Attribution via Sequence Kernels, Markov Chains and Author Unmasking: An Investigation..* ACL Anthology. Empirical Methods in Natural Language Processing 2006, pages 482–491, July 2006.
<https://www.aclweb.org/anthology/W06-1657>
- [17] Twain, Mark. *A Connecticut Yankee in King Arthur's Court*. Project Gutenberg. N.p. N.d. EPUB. Downloaded Feb 25, 2020.
<http://www.gutenberg.org/ebooks/86>
- [18] Twain, Mark. *Adventures of Huckleberry Finn*. Project Gutenberg. N.p. N.d. EPUB. Downloaded Feb 25, 2020.
<http://www.gutenberg.org/ebooks/76>

- [19] Twain, Mark. *Life on the Mississippi*. Project Gutenberg. N.p. N.d. EPUB. Downloaded Feb 25, 2020.
<https://www.gutenberg.org/ebooks/245>
- [20] Twain, Mark. *Personal Recollections of Joan of Arc Vol.1*. Project Gutenberg. N.p. N.d. EPUB. Downloaded Feb 25, 2020.
<https://www.gutenberg.org/ebooks/2874>
- [21] Twain, Mark. *Personal Recollections of Joan of Arc Vol.2*. Project Gutenberg. N.p. N.d. EPUB. Downloaded Feb 25, 2020.
<https://www.gutenberg.org/ebooks/2875>
- [22] Twain, Mark. *The Adventures of Tom Sawyer*. Project Gutenberg. N.p. N.d. EPUB. Downloaded Feb 25, 2020.
<https://www.gutenberg.org/ebooks/74>
- [23] Twain, Mark. *The American Claimant*. Project Gutenberg. N.p. N.d. EPUB. Downloaded Feb 25, 2020.
<https://www.gutenberg.org/ebooks/3179>
- [24] Twain, Mark. *The Gilded Age: A Tale of Today*. Project Gutenberg. N.p. N.d. EPUB. Downloaded Feb 25, 2020.
<https://www.gutenberg.org/ebooks/3178>
- [25] Twain, Mark. *The Mysterious Stranger*. Project Gutenberg. N.p. N.d. EPUB. Downloaded Feb 25, 2020.
<https://www.gutenberg.org/ebooks/3186>
- [26] Twain, Mark. *The Prince and the Pauper*. Project Gutenberg. N.p. N.d. EPUB. Downloaded Feb 25, 2020.
<https://www.gutenberg.org/ebooks/1837>
- [27] Twain, Mark. *The Tragedy of Pudd'nhead Wilson*. Project Gutenberg. N.p. N.d. EPUB. Downloaded Feb 25, 2020.
<https://www.gutenberg.org/ebooks/102>

[28] *Writing Style - Mark Twain Author Study..* Google Sites,
sites.google.com/site/marktwainauthorstudy3/home/writing-style

Appendix

Code to print the result

```
1 from core import *
2
3
4 AUTHORS_FOLDER = "data"
5 from os import listdir
6 from os.path import isfile, isdir
7 from typing import Callable, Type
8
9 class Authors:
10     """
11     This class models a bunch of authors.
12     self.__AuthorList:
13         This is an array of authors created from the root folder, it contains
14         all instance of the author, with generating function and a
15         ignore case option specified by the self.prepare() function.
16     """
17     def __init__(self, MatrixGeneratingFunction: Callable,
18                 IgnoreSpecialNoun: bool=False):
19         self.__AuthorList = None
20         self.prepare(MatrixGeneratingFunction, IgnoreSpecialNoun)
21         self.__IgnoreSpecialNoun = IgnoreSpecialNoun
22
23     def prepare(self, MatrixGeneratingFxn: Callable, IgnoreSpecialNoun: bool):
24         """
25         This function set up all the authors in the folders.
26         * It buffers all the text for each of the authors.
27         * Specifies Matrix generationfunction and Ignore special nouns option.
28         * Please repeatedly call the function cause that will me it slower than
29         necessary.
30         * This function must be called in the __init__ to establish the field
31         of the class, such as the transition matrix generating function, ignore
32         case option, and the
33         author list.
34
35         :param MatrixGeneratingFxn:
```

```

34         This is the matrix function you want to use
35         to generate the transition matrices of the authors.
36     :param IgnoreSpecialNoun:
37         This is set to true if you want to remove
38         all special nouns in the text.
39     :return:
40         None
41     """
42     AuthorsFolders = []
43     for d in listdir(AUTHORS_FOLDER):
44         d = AUTHORS_FOLDER + "/" + d
45         if isdir(d):
46             AuthorsFolders.append(d)
47     assert len(AuthorsFolders) > 0, f"Author's Folder is empty: AUTHORS_FOLDER:{
AUTHORS_FOLDER}"
48
49     self.__AuthorList = []
50     self.__IgnoreSpecialNoun = IgnoreSpecialNoun
51     self.__GeneratingFunction = MatrixGeneratingFxn
52     for AuthorDir in AuthorsFolders:
53         self.__AuthorList.append(Author(AuthorDir, MatrixGeneratingFxn,
IgnoreSpecialNoun=IgnoreSpecialNoun))
54
55     def produce_results(self,
56                         CentroidType:Type[CentroidOption]=CentroidOption.AggregateMatrix,
57                         Mmetric:Type[MatrixMetric]=MatrixMetric.TwoNorm,
58                         Ametric:Type[AuthorMetric]=AuthorMetric.CentroidDis,):
59         """
60         This function takes a series of options on metric of the authors and the
61         matrices and
62         centroid type.
63         * I will return a string, which is a verbal description of all the
64         author's works relative to each other.
65     :param CentroidType:
66         An instance of the class: CentroidType
67     :param Mmetric:
68         An instance of the class: MatrixMetric
69     :param Ametric:
70         An instance of the class Ametric.

```

```

70     :return:
71         A string that is a report of all the works of the authors.
72     """
73     # Specify all the metrics for all the authors
74     Author.CentroidType = CentroidType
75     Author.MMetricType = Mmetric
76     Author.AMetricType = Ametric
77     NameList = [Author.name() for Author in self.__AuthorList]
78     MaxNameLen = max(len(n) for n in NameList)
79     NameList = [Name.ljust(MaxNameLen) for Name in NameList]
80     Res = f"||{'|'}.join(NameList)||\n"
81
82     DistanceTriangle = []
83     for I, EachAuthor in enumerate(self.__AuthorList):
84         Info1, Info2 = EachAuthor.author_cloud()
85         Row = [Info1[0]]
86         for J in range(I + 1, len(self.__AuthorList)):
87             Row.append(dis_between_authors(EachAuthor, self.__AuthorList[J]))
88         DistanceTriangle.append(Row)
89
90     for I, EachRow in enumerate(DistanceTriangle):
91         EachRow = [f"{{: <{MaxNameLen + 1}.4f}}".format(Value) for Value in
EachRow]
92         Res += " " + " " * I * (MaxNameLen + 1) + " ".join(EachRow) + "\n"
93
94     Res += " ".join(str(Author) for Author in self.__AuthorList)
95     Res += "\n" * 50 + "\n"
96     return Res
97
98
99     def __repr__(self):
100         s = ""
101         for author in self.__AuthorList:
102             s += str(author)
103         N = len(self.__AuthorList)
104         A = self.__AuthorList
105         s += "-----ALL AUTHORS TOGETHER-----\n"
106         s += "Author's Distance matrix with triangular part printed:\n"
107         s += " | " + " | ".join(self.__AuthorNameList) + " |\n"

```



```

108         for I in range(N - 1):
109             for J in range(I + 1, N):
110                 s += '{:10.4f}'.format(dis_between_authors(A[I], A[J]))
111                 s += "\n"
112             s += f"The matrix metric is: {Author.MMetricType}\n"
113             s += f"The Author metric is: {Author.AMetricType}\n"
114             s += "All special nouns are IGNORED\n" if self.__IgnoreSpecialNoun else ""
115             s += "===== "
116         return s
117
118
119 def print_experiment(
120     MatrixGeneratingFxn: Callable,
121     IgnoreSpecialNouns: bool,
122     CentroidType: CentroidOption,
123     MatrixMetric: MatrixMetric,
124     AuthorMetric: AuthorMetric
125 ):
126     """
127     Given all the options, this function will print out everything and run the
128     experiment based on that option.
129     :param MatrixGeneratingFxn:
130     :param IgnoreSpecialNouns:
131     :param CentroidType:
132     :param MatrixMetric:
133     :param AuthorMetric:
134     :return:
135     """
136     TheInstance = Authors(MatrixGeneratingFunction=MatrixGeneratingFxn,
137                           IgnoreSpecialNoun = IgnoreSpecialNouns)
138     print(TheInstance.produce_results(CentroidType=CentroidType,
139                                     Mmtric=MatrixMetric,
140                                     Ametric=AuthorMetric))
141
142     return
143
144 def print_cross_compare_experiemment(Author1Directory:str,
145                                     Author2Directory:str,
146                                     MatrixGeneratingFxn: Callable = get_tm27,

```

```

146             IgnoreSpecialNouns: bool = True,
147             MatrixMetric: MatrixMetric = MatrixMetric.

TwoNorm):
148     """
149         This function takes in 2 instances of author and cross compare then and then
150         print all
151         the report on each of the individual works, average distances to the
152         opposite centroid.
153     :param Author1Directory:
154         A string representing the folder that contains all the text files of the
155         author
156     :param Author2Directory:
157         A string, representing another author.
158     :return:
159         None
160     """
161     assert Author2Directory != Author1Directory, "You should not cross compare the
162     same author. "
163
164     Author.MMetricType = MatrixMetric
165
166     Author_A, Author_B= Author(Author1Directory, MatrixGeneratingFxn,
167     IgnoreSpecialNouns), \
168     Author(Author2Directory, MatrixGeneratingFxn,
169     IgnoreSpecialNouns)
170     Author1, Author2 = Author_A.cross_distance_stats(Author_B)
171
172     def print_author(A, B):
173         for K, V in A[0].items():
174             print(f"{K}: {V}")
175         print("-"*40)
176         print(f"The average distance is: {A[1]}, The Standard Deviation is: {A[2]}")
177
178     print("Here is a list of distances when we compare each individuals works of "
179           f"{Author_A.name()} to {Author_B.name()}'s centroid")
180     print("-"*40)
181     print_author(Author1, Author2)
182     print("Here is a list of distances when we compare each individuals works of "
183           f"{Author_B.name()} to {Author_A.name()}'s centroid")
184     print("-"*40)

```

```

178     print_author(Author2, Author1)
179     print("="*40)
180     print("Here is a list the cloud info for both author: ")
181     print("="*40)
182     print(Author_A)
183     print(Author_B)
184     return
185
186
187 if __name__ == "__main__":
188     print_cross_compare_experiemment("data/Charles Dickens", "data/Mark Twain")

```

Code to define functions needed to generate matrices

```

1  """
2  Group 4, this is for the Project 2.
3
4
5  """
6  import numpy as np
7  import matplotlib.pyplot as plt
8
9  from string import ascii_letters
10 from typing import List, Callable, Type
11 import re
12 import enum
13 import math
14 from os import listdir
15 from os.path import isfile
16
17 __all__ = ["Author", "dis_between_authors", "get_tm27", "get_2ndtm", "get_2ndlogarithm",
18           "CHARLES_DICKENS",
19           "MARK_TWAIN", "CentroidOption", "MatrixMetric", "AuthorMetric", "dis"]
20
21 # A list of authors' directory:
22 CHARLES_DICKENS = "data/Charles Dickens"
23 MARK_TWAIN = "data/Mark Twain"
24 ENCODING = "utf8"

```

```

25 PLOT = "plot"
26
27 np.set_printoptions(threshold=10, precision=2, linewidth=1000)
28
29
30 def process_text(filepath: str):
31     """
32     Function reads lines from the file, with the new line character stripped off
33     from
34     the line.
35     :param filepath:
36         The file path.
37     :return:
38         A array of string. Each string is a line in the file.
39     """
40     with open(filepath, 'r', encoding=ENCODING) as f:
41         return [l.strip() for l in f.readlines()]
42
43 def trim_line(s:str, IgnoreCapitalizedWord=False):
44     if IgnoreCapitalizedWord:
45         s = ' '.join([word for word in s.split() if word[0].islower()])
46     Res = ""
47     for c in s:
48         if c == "\'":
49             continue
50         Res += c if c in ascii_letters else (" " if len(Res) >= 1 and Res[-1] != ' '
51         else "")
52     return Res.lower()
53
54 def get_tm27(lines: List[str], ignoreSpecialNoun=False):
55     """
56     Function takes the path of a file and returns the transition
57     matrix based on the 26 letters in the alphabet,
58     The last states of the matrix is space.
59     * All spaces are collapsed.
60     * All punctuations are ignored.
61     * The apostrophe is stripped off from the text.

```

```

62     :param lines:
63         An array of lines that is in the file.
64     :return:
65         """
66     matrix = np.zeros((27, 27))
67     characters = ascii_letters[0:26] + " "
68     for line in lines:
69         trimmed_line = trim_line(line, IgnoreCapitalizedWord=ignoreSpecialNoun)
70         line = list(trimmed_line)
71         for i, c2 in enumerate(line[1:]):
72             c1 = line[i]
73             indx1 = characters.find(c1)
74             indx2 = characters.find(c2)
75             if (indx1 == -1 or indx2 == -1):
76                 continue # Skip non alphabetical characters.
77             matrix[indx1, indx2] += 1
78     for i in range(27):
79         s = np.sum(matrix[i])
80         if s > 0:
81             matrix[i] /= s
82     return matrix
83
84
85 def get_2ndtm(lines: List[str], skipSpecialNoun=False):
86     """
87         Given the content of the file separated by lines, this function will return
88         the
89         26^2 by 26^2 transition matrix.
90         * It's a second order transition matrix based on the letters of the alphabet
91         .
92         * Spaces will be included as the last states of the matrix.
93     :param lines:
94         The content of the file represented in the an array of lines.
95     :return:
96         The np matrix.
97     """
98     Alphabet = ascii_letters[0:26] + " "
99     l = len(Alphabet)
100    n = l ** 2

```

```

99     npmatrix = np.zeros((n, n))
100
101     def s(letter):
102         return Alphabet.find(letter)
103
104     for Line in lines:
105         Line = trim_line(Line, IgnoreCapitalizedWord=skipSpecialNoun)
106         for I in range(len(Line) - 3):
107             i = s(Line[I]) * 1 + s(Line[I + 1])
108             j = s(Line[I + 2]) * 1 + s(Line[I + 3])
109             npmatrix[i, j] += 1
110
111     for i in range(npmatrix.shape[0]):
112         s = np.sum(npmatrix[i])
113         if s > 0:
114             npmatrix[i] /= s
115     return npmatrix
116
117
118 def get_2ndlogarithm(lines: List[str], skipSpecialNoun=False):
119     """
120     Takes the logarithm after counting the frequency,
121     This is a transition matrix that will amplify the
122     occurrences of rare sequence compare to traditional sequences.
123     :param lines:
124         All the lines in the file
125     :param skipSpecialNoun:
126         All special nouns are ignore if this is set to true.
127     :return:
128         A numpy matrix
129     """
130     Alphabet = ascii_letters[0:26] + " "
131     l = len(Alphabet)
132     n = l ** 2
133     npmatrix = np.zeros((n, n))
134
135     def s(letter):
136         return Alphabet.find(letter)
137

```

```

138     for Line in lines:
139         Line = trim_line(Line, IgnoreCapitalizedWord=skipSpecialNoun)
140         for I in range(len(Line) - 3):
141             i = s(Line[I]) * 1 + s(Line[I + 1])
142             j = s(Line[I + 2]) * 1 + s(Line[I + 3])
143             npmatrix[i, j] += 1
144         for i in range(npmatrix.shape[0]):
145             for j in range(npmatrix.shape[1]):
146                 if npmatrix[i][j] > 0:
147                     npmatrix[i][j] = np.log(npmatrix[i][j])
148         for i in range(npmatrix.shape[0]):
149             s = np.sum(npmatrix[i])
150             if s > 0:
151                 npmatrix[i] /= s
152         return npmatrix
153
154
155 class CentroidOption(enum.Enum):
156     """
157     An enum class to represent the options for center of the author cloud.
158     """
159     AggregateMatrix = 1 # Taking the average among all works of the author.
160     AverageMatrix = 2 # Treating all works as one single block of text.
161
162
163 class MatrixMetric(enum.Enum):
164     """
165     An enum class to represent the options of measuring distance between matrices.
166     """
167     OneNorm = 1 # Matrix 1 norm
168     TwoNorm = 2 # Matrix Euclidean distance
169     # WeightedNorm = 3 # Matrix weighted by PD matrix.
170     # HighPower2Norm = 4 # Raising matrix to high power and take the 2 norm.
171     Vectorized1Norm = 5
172
173
174 class AuthorMetric(enum.Enum):
175     """
176     This is a enum class consists of method that can be used to determine the

```

```

177 distance
178 between the cloud of an author to a transition matrix.
179 ! This option only specifies how to determine the distance of an author with a
180 transition
181 matrix!
182 """
183
184 MinimumDis = 1 # The distance between the author and a given transition matrix
185 is the minimum distance
186 # of any work of the author to that transition matrix.
187
188 AverageDis = 2 # Taking the average distance of the given transition matrix with
189 respect to
190 # All the matrices of the author.
191
192 CentroidDis = 3 # This metric take the matrix norm on the difference of 2
193 centroids of the author.
194
195 MM = Type[MatrixMetric]
196 CO = Type[CentroidOption]
197 def dis(Matrix1, Matrix2, Metric:MM):
198     """
199     This function returns the distance between 2 matrices, given
200     the type of Metric space and the weights.
201     :param Matrix1:
202         A numpy matrix.
203     :param Matrix2:
204         A numpy matrix.
205     :param WeightVec1:
206         A numpy vector
207     :param WeightVec2:
208         A numpy vector.
209     :return:
210         A float.
211     """
212     if Metric == MatrixMetric.OneNorm:
213         return np.linalg.norm(Matrix1 - Matrix2, 1)
214     elif Metric == MatrixMetric.TwoNorm:

```



```

211         return np.linalg.norm(Matrix1 - Matrix2)
212     # elif Metric == MatrixMetric.WeightedNorm:
213     #     raise RuntimeError("WeightedNorm not yet implemented. ")
214     # elif Metric == MatrixMetric.HighPower2Norm:
215     #     return np.linalg.norm(Matrix1**10 - Matrix2**10)
216     elif Metric == MatrixMetric.Vectorized1Norm:
217         return np.linalg.norm(np.matrix.ravel(Matrix1) - np.matrix.ravel(Matrix2),
218                                1)
219     else:
220         raise RuntimeError("Invalid Matrix metric space. ")
221
222 class Author:
223     """
224     Files for an author and transitional matrix for the author.
225     * Transitional Matrix classified by each files in the folder
226     * For each file, there will be several transitional matrices for parts of the
227       files.
228     * All the matrices will be in the same order as the list of works.
229     """
230     CentroidType = CentroidOption.AggregateMatrix
231     MMetricType = MatrixMetric.TwoNorm
232     AMetricType = AuthorMetric.CentroidDis
233
234     def __init__(self, dir: str, matrixfunction: Callable = get_tm27,
235                  IgnoreSpecialNoun=False):
236         """
237         Create an instance of an author by specifying:
238         * A directory containing all text files written by the author.
239         :param dir:
240             The directory of the folder.
241         :param matrixfunction:
242             A function you want to use to generate the transition matrices for the
243             authors.
244         """
245         FilePathList = []
246         for filename in listdir(dir):
247             filepath = dir + "/" + filename
248             if isfile(filepath):

```

```

246         FilePathList.append(filepath)
247
248     assert len(FilePathList) > 0, f"There is no file under the directory: {dir}"
249     self.__IgnoreSpecialNoun=IgnoreSpecialNoun
250     FilePath2Lines = {}
251     for f in FilePathList:
252         FilePath2Lines[f] = process_text(f)
253     # A map that maps the file path to array of lines containing the content of
the file.
254     self.__FilePathToLines = FilePath2Lines
255     self.__AuthorName = dir.split("/")[-1]
256     self.__TMFunction = matrixfunction
257     self.__NpMatrices = None # a list of np matrices for each works of the
author
258     self.__AggregateMatrix = None # Instance of transition matrix.
259     self.__AuthorItems = list(self.__FilePathToLines.items())
260
261     def get_fp2lines(self):
262         """
263         A map where the key is the path of the file of an author's work, and the
264         value is a list of string, representing the raw content of the
265         work written by the author.
266         * The text in the line is un-processed.
267         :return:
268         """
269         return self.__FilePathToLines
270
271     def list_of_works(self):
272         return [work.split("/")[-1] for work in list(self.__FilePathToLines.keys())]
273
274     def list_of_works_content(self):
275         return list(self.__FilePathToLines.values())
276
277     def work_matrix_dict(self):
278         """
279         Give a dictionary that maps the name of the works to the
280         transition matrices.
281         :return:
282         """

```

```

283         return dict(zip(self.list_of_works(), self.get_matrices()))
284
285     def name(self):
286         return self.__AuthorName
287
288     def get_matrices(self):
289         """
290         This function returns a transition matrix for each work of the author.
291         each work of the author is a file in the author's folder.
292         :return:
293             A list of np matrix.
294         """
295         res = None
296         if self.__NpMatrices is not None:
297             res = self.__NpMatrices
298         else:
299             self.__NpMatrices = [self.__TMFunction(lines, self.__IgnoreSpecialNoun)\
300                                 for lines in self.list_of_works_content()]
301             res = self.__NpMatrices
302         return res
303
304     def __aggregate_matrix(self):
305         """
306         combine all the lines in the folder into one single work.
307         Then create the transitional matrix for this author,
308         treating all his works as one text.
309         * Results are stored after first time computing it.
310         :return:
311             An npmatrix
312         """
313         if self.__AggregateMatrix is not None:
314             return self.__AggregateMatrix
315         alllines = []
316         for writing in self.list_of_works_content():
317             alllines += writing
318         self.__AggregateMatrix = self.__TMFunction(alllines, self.
__IgnoreSpecialNoun)
319         return self.__AggregateMatrix
320

```

```

321     def __average_matrix(self):
322         """
323             This function returns the centroid matrix.
324             * A centroid matrix is the average for each of all the matrices
325             from this authors.
326             * It's not necessarily a transition matrix anymore.
327         :return:
328             A numpy matrix.
329         """
330         TransitionMatrices = self.get_matrices()
331         N = len(TransitionMatrices)
332         CentroidMatrix = np.zeros(TransitionMatrices[0].shape)
333         for TM in TransitionMatrices:
334             CentroidMatrix += TM
335         CentroidMatrix /= N
336         return CentroidMatrix
337
338     def get_center(self):
339         """
340             Returns a matrix representing the center of the author,
341             The type of matrix depends on the global options for the authors.
342         :return:
343             An Np matrix.
344         """
345         Center = None
346         if Author.CentroidType == CentroidOption.AggregateMatrix:
347             Center = self.__aggregate_matrix()
348         elif Author.CentroidType == CentroidOption.AverageMatrix:
349             Center = self.__average_matrix()
350         else:
351             raise RuntimeError("Unspecified Centroid Option.")
352         return Center
353
354     def work_distances(self):
355         """
356             The function will calculate the distance for all the transition matrix
357             from each file.
358         :param mode:
359             mode == 1:

```

```

360         Using the centroid as the center of this author.
361         mode != 1:
362             Using the aggregate matrix as the center of this author.
363     :return:
364         A map; the key is the name of the file and the
365         value is the distance from teh centroid matrix.
366     """
367     DistanceMap = {}
368     Center = self.get_center()
369     for Writing, Matrix in zip(self.list_of_works(), self.get_matrices()):
370         DistanceMap[Writing] = dis(Matrix, Center, Metric=Author.MMetricType)
371     return DistanceMap
372
373 def author_cloud(self):
374     """
375         * The average distance.
376         * the standard deviations of the distance.
377         * A map describing the distance from the center, mapping
378         author's works to a distance represented in a float value.
379     :return:
380         2 items:
381         1. [<average distance>, <standard deviation>], distancelist
382         2. A map, string to float.
383     """
384     DistanceList = self.work_distances()
385     Sum = 0
386     Squaresum = 0
387     L = len(DistanceList)
388     for Distance in DistanceList.values():
389         Sum += Distance
390         Squaresum += Distance**2
391     Sum/=L
392     Squaresum/=L
393     return [Sum, math.sqrt(Squaresum - Sum**2)], DistanceList
394
395 def distance_to(self, m2):
396     """
397         The function return the distance of this author to a given transition
398         matrix of the same size.

```

```

399         * How the distance is calculated depends on the input
400         parameters.
401     :param norm:
402         The matrix to norm to calculate the distances.
403     :param mdoe:
404         mode==1:
405             Distance from the centroid of the author.
406         mode != 1:
407             Distance from the aggregate matrix of the author.
408     :return:
409         A float.
410     """
411     centroid = self.get_center()
412     if Author.AMetricType == AuthorMetric.CentroidDis:
413         return dis(centroid, m2, Metric=Author.MMetricType)
414     temp = [dis(m1, m2, Metric=Author.MMetricType) for m1 in self.get_matrices()]
415 ]
416
417     if Author.AMetricType == AuthorMetric.MinimumDis:
418         return min(temp)
419     if Author.AMetricType == AuthorMetric.AverageDis: # Not sure if symmetry
420         property is preserved.
421         return sum(temp)/len(temp)
422     raise("Invalid Author Metric. ")
423
424 def cross_distance_stats(self, AnotherAuthor):
425     """
426     * This function will compare each of the works of THIS author
427     to another author using the matrix metric.
428     * This function will return detailed statistics about the
429     distance for all the works from this author to that author, and THAT
430     author to THIS author.
431
432     :param AnotherAuthor:
433         An in stance of another author that is not this author.
434     :return:
435         - w lists are returned, each is a packed information about all information.
436         [<Distance Dict>, Avg_distance, SD],
437         [<Distance Dict>, Avg_distance, Sd],
438     """

```

```

436     def sd_avg(Arr):
437         ArrSquared = list(map(lambda x: x**2, Arr))
438         return sum(Arr)/len(Arr), math.sqrt(sum(ArrSquared)/len(ArrSquared) - (
sum(Arr)/len(Arr))**2)
439
440     def cross_compare(Matrices, WorkList, Centroid):
441         Distances = [dis(M, Centroid, Metric=Author.MMetricType) for M in
Matrices]
442         DistancesDict = dict(zip(WorkList, Distances))
443         Avg, SD = sd_avg(Distances)
444         return [DistancesDict, Avg, SD]
445
446         ThisCompareToThat = cross_compare(self.get_matrices(), self.list_of_works(),
AnotherAuthor.get_center())
447         ThatCompareToThis = cross_compare(AnotherAuthor.get_matrices(),
AnotherAuthor.list_of_works(), \
448                                         self.get_center())
449         return ThisCompareToThat, ThatCompareToThis
450
451
452     def __repr__(self):
453         s = "-----AUTHOR INFO-----\n"
454         s += f"Author's Name: {self.__AuthorName} \n"
455         s += "Distances of his works from the center:\n"
456         Cloud, DistanceList = self.author_cloud()
457         s += f"Average Distance from the center: {Cloud[0]}\n"
458         s += f"Standard deviation of the distances: {Cloud[1]}\n"
459         TitleMaxLength = max(len(W) for W in DistanceList.keys())
460         for Work, dis in DistanceList.items():
461             s += f"{{(Work+':')}.ljust(TitleMaxLength)} '{{:10.4f}}'.format(dis)} \n"
462         s += f"Matrix Norm used: {Author.MMetricType}\n"
463         s += f"Centroid Matrix is: {Author.CentroidType}\n"
464         s += f"Function used to generate transition matrix: {self.__TMFunction.
__name__}\n"
465         return s
466
467
468     def dis_between_authors(author1, author2):
469         """

```

```

470         This function returns 1 number to represent the distance between 2 author's
471         works.
472     :param author1:
473         An instance of an Author class.
474     :param author2:
475         An instance of an Author class.
476     :param metric
477         A type of author metric.
478     :return:
479         a float.
480     """
481     metric = Author.AMetricType
482     if metric == AuthorMetric.CentroidDis:
483         return author2.distance_to(author1.get_center())
484     temp = [author1.distance_to(Author2Works) for Author2Works in author2.
485             get_matrices()]
486     if metric == AuthorMetric.AverageDis:
487         return sum(temp)/len(temp)
488     if metric == AuthorMetric.MinimumDis:
489         return min(temp)
490     raise RuntimeError("Invalid AuthorMetricOption")
491
492 def main():
493     """
494         Adventures of Huckleberry Finn: tm27, 2ndtm, and their centroid.
495     """
496     print("Producing and print...")
497
498     def plot_matrix_savefig(m, filedir: str):
499         plt.imshow(m)
500         plt.colorbar()
501         plt.savefig(filedir)
502         plt.clf()
503
504     def save_nparray(m, name):
505         np.savetxt(X=m, fname=name)
506     global Author
507

```



```

508 Mark27 = Author(MARK_TWAIN, get_tm27)
509 Mark2nd = Author(MARK_TWAIN, get_2ndtm)
510 Mark27_ignore = Author(MARK_TWAIN, get_tm27, IgnoreSpecialNoun=True)
511 Mark2nd_ignore = Author(MARK_TWAIN, get_2ndtm, IgnoreSpecialNoun=True)
512 filenames = ["Mark27", "Mark2nd", "Mark27_ignore", "Mark2nd_ignore"]
513 authors = [Mark27, Mark2nd, Mark27_ignore, Mark2nd_ignore]
514 for Author, Filename in zip(authors, filenames):
515     for WorkName, Matrix in zip(Author.list_of_works(), Author.get_matrices()):
516         if "Huckle" in WorkName:
517             save_nparray(Matrix, Filename + " Huckle.txt")
518             save_nparray(Author.get_center(), f"{Filename} Centroid.txt")
519 print("Ok, we are going to save some centroid for both of the authors now: ")
520
521
522 def save_matrices_forall_data():
523     """
524     This function will save all the works of the authors's transition
525     matrices as a text files, this is for the concerns of correctness of
526     the codes, so that people can just take all the text file and confirm the
527     correctness by themselves.
528     :return:
529     None
530     """
531     global Author
532     def save(NpMatrix, dir:str, filename:str):
533
534         np.savetxt(fname=filename, X=NpMatrix)
535         return
536
537     def generate_all_authors():
538         AllMatrices = [get_tm27, get_2ndtm, get_2ndlogarithm]
539         FileLocations = [CHARLES_DICKENS, MARK_TWAIN]
540         ListofAuthors = []
541         for L in FileLocations:
542             for G in AllMatrices:
543                 ListofAuthors.append(
544                     Author(dir=L,
545                           matrixfunction=G,
546                           IgnoreSpecialNoun=True))

```

```

547         ListofAuthors.append(
548             Author(dir=L,
549                   matrixfunction=G,
550                   IgnoreSpecialNoun=True))
551
552     return ListofAuthors
553
554     # TODO: FINISH THIS SHIT.
555     for Aut in generate_all_authors():
556         AuthorName = Aut.name()
557         for Work, Matrix in zip(Aut.list_of_works(), Aut.get_matrices()):
558             save(Matrix, "")
559         pass
560
561     pass
562
563
564 if __name__ == '__main__':
565     main()
566     pass

```

Github Repository for the codes:

<https://github.com/iluvjava/Math-381-Project-2>