



# **GRAPH KERNELS AS PREPROCESSING FOR UNSUPERVISED CLUSTERING: TWO CASE STUDIES**

by

LEVI C. NICKLAS

A Thesis Submitted to the Faculty of the  
DEPARTMENT OF COMPUTER SCIENCE

In Partial Fulfillment of the Requirements

For the Degree of

MASTER OF SCIENCE

In the Graduate College

Florida Polytechnic University

2021

# GRAPH KERNELS AS PREPROCESSING FOR UNSUPERVISED CLUSTERING: TWO CASE STUDIES

by

LEVI C. NICKLAS

A Thesis Submitted to the Faculty of the  
DEPARTMENT OF COMPUTER SCIENCE

In Partial Fulfillment of the Requirements

For the Degree of  
MASTER OF SCIENCE

In the Graduate College

Florida Polytechnic University

2021

Approved by:

	<u>Signature</u>	<u>Date</u>
Dr. Reinaldo Sanchez-Arias (Committee Chair, Advisor)	_____	_____
Dr. Grisselle Centeno (Committee Member)	_____	_____
Dr. Harish Chintakunta (Committee Member)	_____	_____
Dr. Tom Dvorske Vice Provost of Academic Affairs (Graduate Division)	_____	_____

To my wife Julianne, whose love, support, and encouragement kept me  
on the right track. I look forward to this life we are making.

Ecclesiastes 9:9

And to my dear baby girl, I hope this work is something you can be  
proud of, and that this be a reminder that hard work and perseverance  
are always rewarded.

1 Corinthians 9:24-25

## **ACKNOWLEDGMENTS**

This work would not have been possible without the guidance and support of my mentor, advisor, and friend, Dr. Reinaldo Sanchez-Arias. His example of excellence at what he does, is one to chase; I hope to catch you one day.

I would like to thank Florida Polytechnic University for their support, both in the form of scholarships and in the community that they foster on campus. I look forward to returning to campus and seeing how our small school is making a big impact.

# CONTENTS

---

<b>List of Figures</b>	<b>iii</b>
<b>List of Tables</b>	<b>iv</b>
<b>List of Algorithms</b>	<b>v</b>
<b>Abstract</b>	<b>1</b>
<b>1 Introduction</b>	<b>1</b>
<b>2 Literature Review</b>	<b>5</b>
<b>3 Methods</b>	<b>8</b>
3.1 Introduction . . . . .	8
3.1.1 Skip-grams . . . . .	8
3.1.2 Graph Kernels . . . . .	9
3.1.3 Using Kernel for Clustering . . . . .	10
3.1.4 Kernel Density Estimation Clustering for Linear Kernel . . . . .	10
3.2 Processing Text into Graphs . . . . .	10
3.3 Graph Kernels for Similarity . . . . .	12
3.4 Graph Kernels for Similarity . . . . .	13
3.5 Development of Methods . . . . .	14
<b>4 Results and Analysis</b>	<b>18</b>
4.1 NHTSA Special Crash Investigations . . . . .	18
4.2 reddit Threads . . . . .	24

4.3 Comparative Analysis . . . . .	27
<b>5 Conclusion</b>	<b>28</b>
<b>6 Future Work</b>	<b>30</b>
<b>Appendices</b>	<b>31</b>
<b>A My first appendix</b>	<b>32</b>

# LIST OF FIGURES

---

4.1	Hierarchical Clustering Variation (NHTSA) for differing hyperparameters.	19
4.2	Five hyper-parameter sets and their corresponding dendrograms. . . . .	20
4.3	Co-membership of documents across hyper-parameter sets A,B,C,D, and E.	21
4.4	Term-Frequency/Inverse Document-Frequency for NHTSA reports, using hyper parameter set A. . . . .	22
4.5	Skip-grams which occurred more than 50 times in hierarchical clustering with hyper-parameter set A, in cluster 3. . . . .	23
4.6	Within sum of squares as a function of number of clusters . . . . .	25
4.7	Left - Proportion of cluster's threads coming from each query word. Right - Dendrogram for reddit thread analysis, with 4 clusters. . . . .	26

## LIST OF TABLES

---

2.1	Summary of Literature Review . . . . .	6
4.1	Hyperparameter for Variation Study in Figure X . . . . .	19



## LIST OF ALGORITHMS

---

## Abstract

Abstract goes here

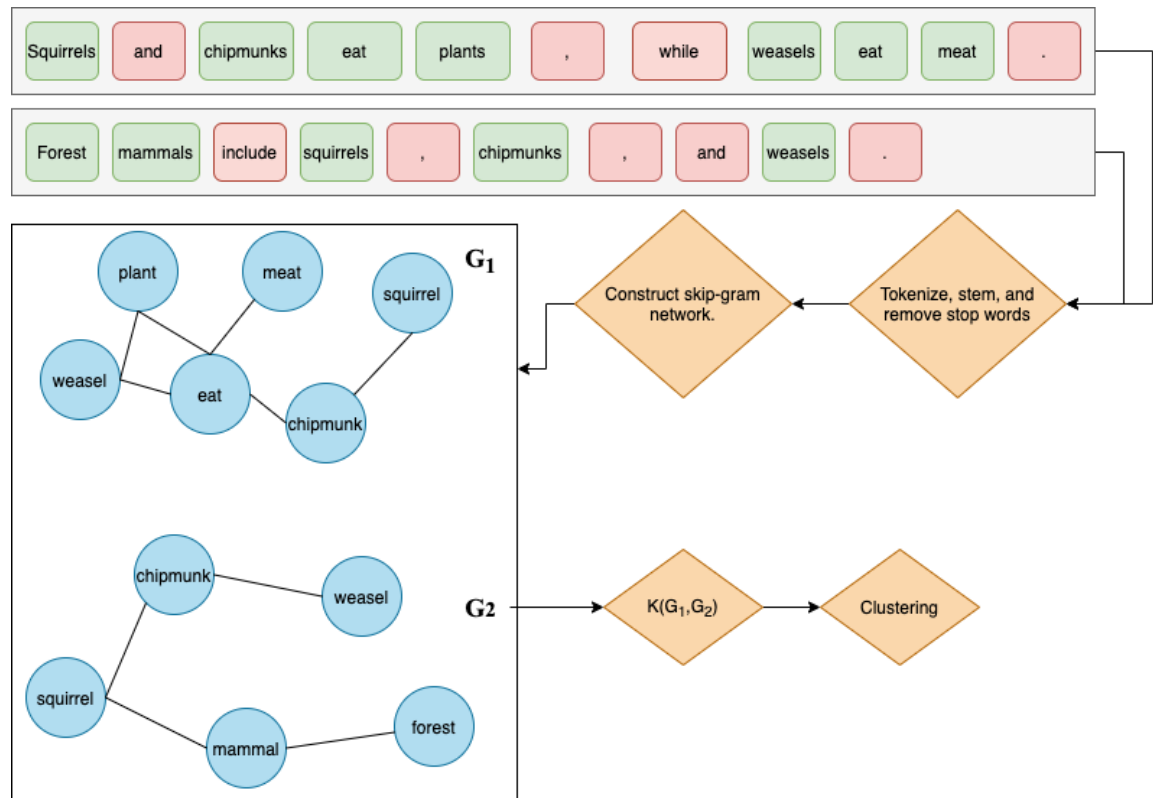
# CHAPTER 1: INTRODUCTION

---

## Introduction Chapter Here

Text mining is the process of extracting insight from text documents using computational and statistical methods. A common task in text mining is document clustering, where the natural language of the documents is reformatted and used to group documents into clusters of similar topics or text content. Popular methods to cluster text include k-means, hierarchical methods, and topic models. Most of these methods utilize term frequency-inverse document frequency measures or bag-of-words methods to do clustering, though these methods often divide words and thus may lose meaning, or context, surrounding the word of concern. In an effort to cluster documents while still preserving the relationship between words, text can be modeled with graphs which better preserve the context surrounding a word. Text can be represented with bigram graphs, where the edge between two vertices is the result of two words appearing adjacent in the text. This idea can be extended to more general n-grams, and further extended to skip-n-grams, resulting in even richer representations of the text. The similarity of two of graphical representations can then be assessed using modern methods called graph kernels. Graph kernels produce a measure of similarity between graphs, thus allowing for further machine learning to take place. With a measure of similarity between two graphs, we can do an assortment of clustering methods.

The chosen graph representation of the observational unit of text is an extension of bigram graphs. The use of skip-grams here is an attempt to connect ideas between words that would have normally attained meaning through the "context" of the sentence. By connecting words that appear close to one another while not being immediately adjacent, we are creating "skip-grams"—an extension of the bigram where the window for which words are considered adjacent is widened.



In the above figure, a demonstration of bigram construction is available (**Levi**: This graphic needs to be a skipgram for which  $k > 1$ , and comments thereafter need to be about the skip gram set up...this will be a place holder for now).

Many methods used in text mining use bag-of-words techniques, but these methods do not always preserve the context and relationships between words. The graph representations for text, used here, aim to preserve the relationship and context between words. The idea of bigrams (pairs of two words) can be extended to “skip-grams”; skip-grams are words which appear within some width  $w$  of each other. This leads to more connections, and a richer graph representation of the text. This skip-gram method is similar to that used in the popular machine learning package Word2Vec, which uses the “continuous bag of words” or “continuous skip-grams”. In

this paper, a graph representation of the text is constructed, using skip-gram methods, for each comment in the dataset.

The graph representations are compared with a graph kernel, which is a measure of the similarity of two graphs. In this study, the graph kernel of choice is a “edge histogram kernel”, because the kernel utilizes the labels of the vertices to assess the similarity of two graphs– not just the topology. For each graph, the kernel produces a measure of similarity to the other  $n-1$  graphs in the dataset. We can use this similarity measure between two graphs, to assess how similar other graphs in the dataset may be. That is, we can compare the similarity of graphs B and C through their similarity to A. To do this, a KDE (Kernel Density Estimation) to estimate a kernel density curve, the curve is then used to partition the rest of the dataset into clusters. The KDE bandwidth is modified to produce more or less local minima/maxima that can be then used to identify potential clusters. Basic calculus can be used to locate to inflection points that will serve as intervals for which a cluster will be defined.

Alternatively, we can use machine learning methods which work nicely with kernels, i.e. support vector machine (supervised) or hierarchical clustering (unsupervised). These methods are more tried and true, and can serve as a comparison to the success of the KDE methods.

To test these methods, two different text data sets will be considered. First, reddit comments from subreddits pertaining to mental health. Using the redditExtractoR package by (NAME), the comments are collected if they match criteria set by the query. From the subreddit r/mentalhealth, posts are collected if they have chosen keywords and at least 10 comments. The keywords considered are “anxious”, “anxiety”, “depressed”, “depression”, “mental”, “illness”, “scared”, “afraid”, “sad”, “emotion”, “anger”, “angry”, “upset”, “suicide”, “abuse”, “emotional”, “help”, and “addiction”. These words do not represent all words which define or indicate mental health topics

may be present; these words were chosen by the researcher arbitrarily. The result of the queries is thousands of comments. Mental health continues to be a growing focus at the national level, and the discussion on reddit is often candid and open, due to reddit's reputation as an anonymous website. The second data set being considered in the study, is NHTSA report data. These reports are from the Special Crash Investigation (SCI) reports from the NHTSA (National Highway Traffic Safety Administration). The pdf documents that are considered in this study are those that involved ambulance(s) in the incidents highlighted in the reports. These events are considered "edge cases" and could prove useful to those working in autonomous vehicles.

These datasets were chosen because they are two very different styles of writing. The reddit posts are often in contemporary and casual use of language, possibly including "netspeak", while the NHTSA reports feature a technical writing style with plain language, absent of hyperbole or sarcasm. These datasets should show two very different sides of text data, and comparisons will be made to how each dataset responds to the methods outlined above.

## CHAPTER 2: LITERATURE REVIEW

---

The introduction of graph kernels first arose in 2002, from the paper “*Diffusion kernels on graphs and other discrete structures*” [?]. Since then, graph kernels have been used in a variety of fields, and modified to suit different problems. These developments have resulted in applications in biology, chemistry, social media analytics, and machine learning. Graph kernels have since formed into 3 large sub-groups: random walk kernels, sub-graph (also called graphlet) kernels, and tree-based methods. Within these three categories there are a multitude of modifications which include things such as edge or vertex attributes into the calculations [?]. The graph kernels of particular interest to this study are those which incorporate vertex attributes into the calculation of the kernel. This inclusion of vertex attributes will allow the assignment of a word from a text document to a vertex.

The use of graph kernels in text processing has been largely focused on applying them to a neural networks for NLP. However, a 2017 paper by Nikolentzos et.al, covered how their use first arrived at the idea to apply graph kernels to what they called “graph-of-words” [?]. They utilized a window width to construct a connections between words, effectively “skipping” some words in between; this concept has also been called a “skip-gram” network in other contexts [?]. Members of this research group, and others close to them, have produced software packages to perform some basic graph kernel methods. In addition, there have been a small number of papers from the group about the application of these methods to text mining [?]. Their work centered around use of neural networks for classification tasks.

Table 2.1: Summary of Literature Review

Topic	Author	Title	Year
Graph Kernels	Vishwanathan	Graph Kernels	2010
Graph Kernels	Kriege et al.	A survey on graph kernels	2020
Graph Kernels	Nikolentzos et al.	[graph kernels for doc. similarity]	2017
Graph Kernels	Kondor and Lafferty	Diffusion Kernels on Graphs [...]	2002
Skip-Grams	Cheng et al.	From n-gram to skipgram [...]	2006
Text Mining	Vazirgiannis et al.	GraphRep: [...]	2018
Application	Rosenfeld et al.	Kernel of Truth: [...]	2020
Software	Silge and Robinson	tidytext	2016
Software	Sugiyama et al.	graphkernels	2018
Software	Casardi et al.	igraph	2013

Using a similar framework, the performance of the graph kernels as a pre-processing step for unsupervised learning will be assessed with data sources other than what the originators of the idea used. This assessment will result in validation, or disagreement, with their methods. The data sources will come from social media websites, which will serve as small document size examples, and technical reports from the National Highway and Transportation Safety Administration (NHTSA), which will serve as a larger document size example. These datasets, while very different in both writing style and length, provide an excellent opportunity to compare how performance, computation, and accuracy vary across differing types of text data. The NHTSA data was curated within a research group at Florida Polytechnic University this summer, and the process of data cleaning and transformation is in its final stages for initial exploratory data analysis and tests. This application would appeal to those within the Advanced Mobility Institute (AMI) as a way to text mine these technical documents and identify “edge cases” within the reports, further supporting the efforts to provide closer to real life scenarios for testing and validation. The social media data will be collected off of reddit, and will remain focused on social health forums on the platform. This will be a continuation of work done at Florida Polytechnic University as well. By text mining these forums, there is an opportunity to better compare comments, threads, or communities at large, thus allowing for additional machine learning processes to take place.



The work proposed here will center not only on unsupervised methods but also empirical studies on the effectiveness of graph kernels on these types of text data. The effects of varying window width, edge weighting, other graph structures, and document types will be considered. These contributions to the topic will aid in demonstrating the effectiveness on differing document types, and more importantly an exploration into using graph kernels as a preprocessing method for unsupervised learning methods. Using data sources I am familiar with, I will be building upon previous work in unsupervised text mining [?]. Data and text mining have been a focus of my education at Florida Polytechnic and this project will be utilizing various skills I have learned through the program. In addition, this work will be relevant to my professional career as well.

# CHAPTER 3: METHODS

---

## 3.1 Introduction

### 3.1.1 Skip-grams

As an alternative to natural language processing (NLP) methods, which are reliant on “bag-of-words” methods, the methods used here utilize a graph representation of the text. Consider a bigram, a pair of two words—like “hot dog” or “peanut butter”, these bigrams can be constructed for a text document where every pair of adjacent words is a bigram. The bigrams can then be used to make a graph, where each word is a vertex, and each bigram is an edge. This graph representation holds more context than the bag-of-words methods; for example seeing the words “cake” and “carrot” in a bag of words may not show that “carrot cake” was the real intent of the text. This is an important concept for modeling text, as we should strive to achieve a representation of the text that makes for effective modeling that will capture the true meaning of the text in question. Keeping this in mind, with the example of “carrot cake”, what about the idiom “beating a dead horse”? Each word individually may mean something other than the idiom. Even the bigrams “beating dead” and “dead horse” do not capture what the idiom means. We can expand the number of words in the  $n$ -gram to be 3 or 4 words, or alternatively, we can make more “edges” or connect more words. We can connect words that are not immediately adjacent but perhaps within  $k$  words away. These bigrams that appear within  $k$  words of each other are called “skip-grams”. The skip-gram allows to capture context of larger sequences of words since the graph representation will show how the  $k$  wide neighborhood of words was connected. In the idiom example, using skip-grams with window width  $k = 2$ , and removing common words (e.g. “a”, “at”, “the”), will produce a graph like:

$$E(G) = \{\text{beat} \longleftrightarrow \text{dead}, \text{dead} \longleftrightarrow \text{horse}, \text{horse} \longleftrightarrow \text{beat}\}$$

This graph representation contains a cycle, of length 3, where most native english speakers will identify the meaning behind the graph representation. As ideas, idioms, figures of speech, and other concepts (that may be explained in a non-literal fashion) grow in size as they include more words, it becomes more difficult to capture the meaning behind the text. However, leveraging the concept of a skip gram can produce such a rich graph representation of the text that the original meaning is more likely to be preserved.

### 3.1.2 Graph Kernels

The next natural question is, “how can we compare these graph representations?”, and we address this with graph kernel methods. These methods are generally used to compare the similarity of graphs. These use of a graph kernel to compare graphs was first published in 2003, and since then various applications and adaptations have been made to the methods. In the case of text mining, the graph kernel must assess vertex labels —if one intends to map words to vertices, otherwise they will be assessing the topology alone. In this study, the Edge-Histogram kernel is the kernel used to compute similarity. This kernel was chosen as it uses labels on the graph structure, and is not as computationally intensive as other methods (CITE). In the specific implementation used for these studies, the computation time was the shortest when compared with other kernel methods like: graphlet, random walk, and Weisfeiler-Lehman kernel. (CITE) Since the data sets of concern in the studies feature either large graphs or a large number of graphs, the kernel had to be cheap computationally.

To compute and edge histogram kernel on two graphs,  $G_1$  and  $G_2$ , first define the set of edges  $E_i = \{(u_1, v_1), (u_2, v_2), \dots, (u_n, v_n)\}$  where  $(u_n, v_n)$  is the  $n$ -th edge connecting  $u_n$  to  $v_n$ . Then the edge label histogram is defined to be  $\vec{g} = \{g_1, g_2, \dots, g_s\}$  where (FINISH

THIS, link in comment)

### 3.1.3 Using Kernel for Clustering

The output of the kernel is useful for a variety of tasks. Some other popular applications have included classification with support vector machines, which are popular with other kernel methods. In this case, the kernel is used for unsupervised clustering. Within the kernel matrix,  $K$ , the entry  $k_{i,j}$  represents the similarity between graphs  $i$  and  $j$ . This matrix which contains measures of similarity between points can be used as a distance matrix for hierarchical clustering. Before using the graph kernel as a distance matrix, normalization or standardization takes place, and principal component analysis may be used. The end result is each row is a single graph-document being described by its similarity to all the other graphs, which are the column values. Once the values are transformed or rotated by preprocessing methods, the points are just represented by their similarity to one another, but in a transformed space. Various hyper parameters can be tuned for successful clustering; the graph kernel has a parameter that can be tuned, and the hierarchical clustering can be tried with differing types of linkage.

### 3.1.4 Kernel Density Estimation Clustering for Linear Kernel

In addition to hierarchical clustering, a method was developed to find potential clusters based on the kernel similarity measure, but while measuring similarity to a single graph. For example, we can compare how similar graphs  $B$  and  $C$  without computing their similarity, by comparing how similar they are to graph  $A$ ; this extension of transitive property logic allows for focusing on the similarity of the graph list as it relates to just one graph.

## 3.2 Processing Text into Graphs

For both datasets, the reddit comments and the SCI papers, the text data is collected and stored in its raw form for reproducibility. The text is then cleaned and prepared for analysis using a suite of text processing tools from the *Tidyttext* (CITE) package for R.

Using this package, the text undergoes tokenization into skip-grams, stop word removal, and filtering to remove punctuation, numbers, short words, etc.

Skip-grams are produced for  $k = 2, 3, 4$ , where  $k$  is the window width of the skip-gram for which bigrams are formed. Stop words are gathered from popular lexicons: "snowball", "SMART", and "onix". As stated above, all numeric values, and lingering punctuation were also removed through use of the *stringR* package (CITE). This processing is an essential preprocessing step for using graph kernels to gauge similarity of documents, because words and strings (like punctuation or numbers) that get repeated frequently across many documents in the set do not actually mean the documents are similar—they just have common reoccurring words. Removing these types of words forces the text data to be more unique across the document set.

The result of the skip-gram tokenizer is a data frame where each observation is a bigram, a pair of words, that the skip-gram window captured. The data frame is then cleaned of stop words through using anti-joins on each bi-gram; any row with the appearance of a stop word in either position, the first or second word, was removed. The same method was applied for punctuation and numbers. The result was a data frame of bigrams which appeared in a fixed window width,  $k$ , of one another and where both words are of a length greater than 3 letters, not a stop word, and do not contain numbers or punctuation. This cleaned dataset is now ready to be converted into a graph object.

Now that the skip-grams are cleaned, the data frame can be converted to a graph object, through use of the *igraph* package. To do this, each word pair in the data frame is converted to an edge and vertex pair in the graph object. For example, the word pair "data frame" will become two vertices, labeled "data" and "frame", with an undirected edge connecting them. In this study, directed graph edges are not used, but could be considered in future iterations of this work. When this completed for all the skip-gram pairs that were generated, it produces a singular connected graph, however a great deal of filtering occurred and there may be disconnected portions. When words are removed from the graph, a vertex will disappear and can potentially split the graph. This is not too common in the NHTSA dataset for two reasons. First, the text data is quite long,

and so if a word is reused at a later time in the text there will be an additional connection to keep it included in the main graph. Secondly, the benefit of the skip-gram, as opposed to plain bigrams, is that the larger window width means words get connected and can "skip" over words that will get removed through the data cleaning process. So at this point, if there is a group of words that are isolated and not connected to the main graph, they are often quite small and are only several words. To keep computation simple, these lingering small isolated graphs are removed. The vertices that are not members of the main graph are removed from the graph object and that leaves a single connected graph for each text document. These graphs can be compared with graph kernels easily now.

### 3.3 Graph Kernels for Similarity

Graph kernels can be used to compare the similarity of two graphs. The development of these methods arose out of the need to determine if graphs were isometric in a faster way. The solution was a graph kernel which produces a scalar value for how similar two, or more, graphs are. The result of a graph kernel is a matrix where the similarity of graphs  $i$  and  $j$  is in the kernel's  $i$ -th row and  $j$ -th column entry. This resulting matrix can be used in a variety of ways, but it will be used as a distance matrix in applications here.

Amongst graph kernels which consider edge or vertex labels in their computation, various surveys and studies have found edge label histogram kernels or vertex label histograms to be the most efficient. They may not out perform the other methods, like a Weisfeiler-Lehman or other subgraph methods, but they are computationally cheap. The datasets being studied here are both large: one contains many smaller graphs, and one contains 48 very large graphs. So for this study, computation efficiency was prioritized.

The edge label histogram kernel is the graph kernel that was chosen for these datasets, and it can be computed using either a linear kernel or a Gaussian radial basis

function (RBF) kernel between the edge label histograms.

An edge label histogram is defined as  $\vec{g} = (g_1, g_2, \dots, g_i)$  such that  $g_i = |\{(u, v) \in E | \phi(u, v) = i\}|$  for each  $i$  (CITE SUGIYAMA HALTING KERNELS). Where  $g_i$  is a histogram bin for a unique edge label's magnitude,  $E$  is the set of edges, and  $\phi$  is a function that maps each label to a scalar value in the range of unique values. The edge label histograms are then passed through a kernel, either a linear kernel or a Gaussian RBF kernel.

Computation using a linear kernel takes two graphs,  $G$  and  $G'$ , and uses their edge label histograms  $\vec{g}$  and  $\vec{g}'$ . The kernel is computed as:

$$K(\vec{g}, \vec{g}') = \vec{g}^T \vec{g}'$$

The resulting value is stored in the graph kernel matrix as the measure of similarity between the two graphs in the corresponding row and column for the pair. (CITE SUGIYAMA)

Alternatively, the Gaussian RBF kernel takes the edge label histograms of  $G$  and  $G'$ , that we call  $\vec{g}$  and  $\vec{g}'$ , and the kernel is computed as:

$$K(\vec{g}, \vec{g}') = e^{-\left(\frac{\|\vec{g} - \vec{g}'\|^2}{2\sigma^2}\right)}$$

The resulting value is stored in the graph kernel matrix as the measure of similarity between the two graphs in the corresponding row and column for the pair. (CITE SUGIYAMA)

Through either of these kernels, we obtain a kernel of dimensions  $n \times n$  for a list of  $n$  graphs. This kernel can then be used for clustering methods.

### 3.4 Graph Kernels for Similarity

The resulting graph kernel matrix, with dimensions  $n \times n$ , is high dimensional data. We can take this kernel and perform a principal component analysis (PCA) to reduce the

high dimensionality of this data. Each row, can be considered as an observation, and each column as it's similarity value or in the  $j$ -th graph's dimension. PCA then allows us to express the data in a lower dimensionality, hopefully in 2 or 3 dimensions which can be visualized much better.

Then,

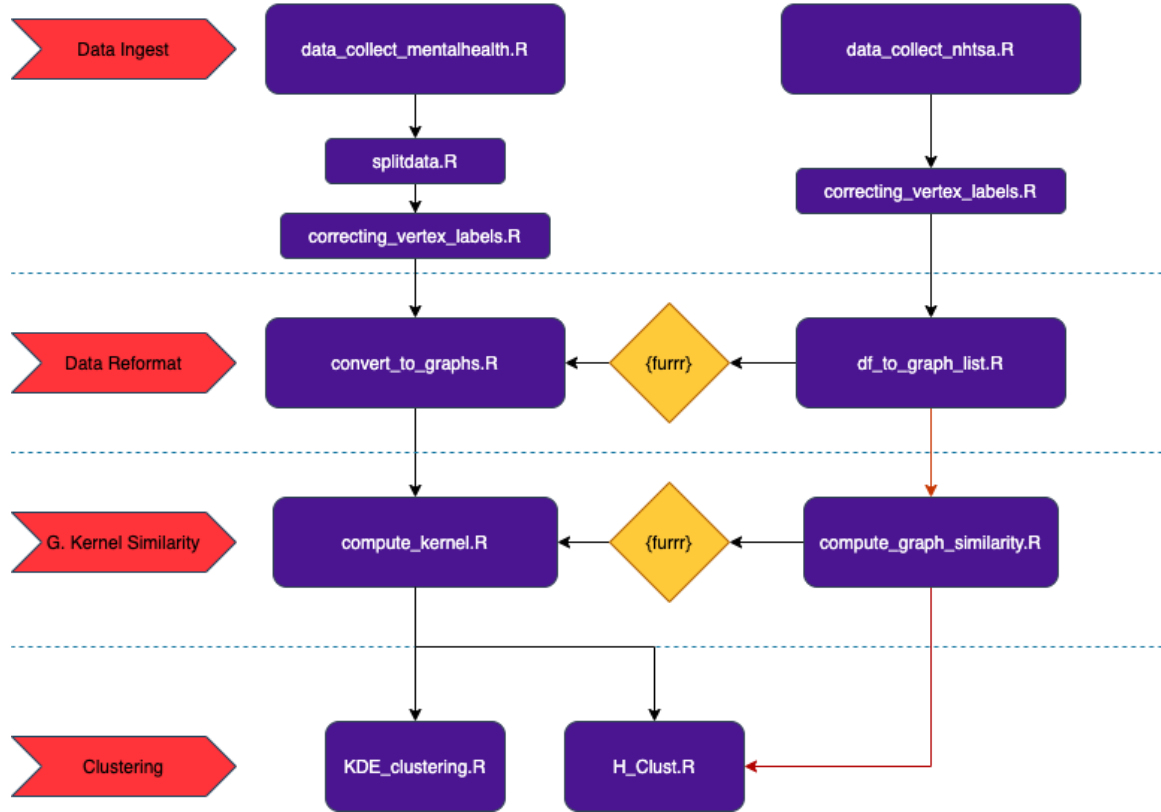
### 3.5 Development of Methods

All of the methods outlined in preceding chapters were developed in *RStudio* and are stored in a public code repository on *GitHub*. Initial development and testing of all the packages and libraries was completed using *R Markdown* notebooks, which allow for coding in chunks with explanation in markdown text. These netbooks are kept in the code repo to serve as more explanatory documents demonstrating how all the scripts and datasets work in conjunction with one another.

The resulting and final scripts are kept in the code repository as well. The final scripts were developed to use the output from one another and to produce intermediate outputs so that small studies could be conducted. How the scripts work with one another is mapped out in the "Script Map" in figure X.

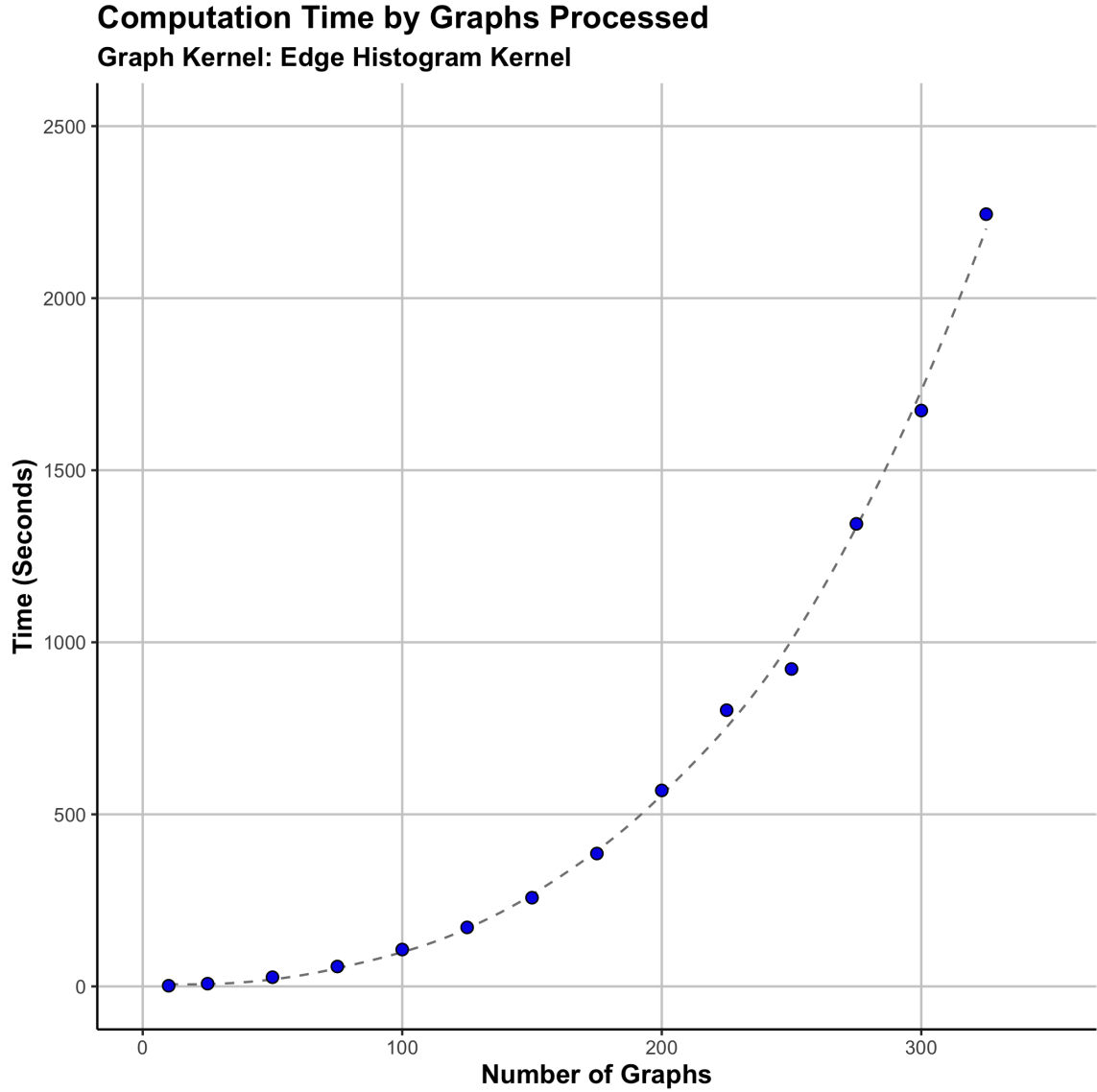


## Script Map



The Script Map shows how these R scripts can be used for either standard computation, in the case of the NHTSA dataset, or parallel computation, as was used in the reddit dataset. Parallel computation was completed through use of the `{furr}`, an R package that makes parallization of functions similar to mapping functions from the popular `{purrr}` package. Implementation of a parallel option was necessary because computational time for a list of graph objects grows rapidly as the number of graphs in the list increases. A small study was conducted on the reddit dataset to plot this

relationship.



Above, in Figure X, there is strong evidence that the computation time follows  $O(n^2)$ . For this reason, the parallel workflow was developed—to handle datasets with a large number of documents to compare. Even though the amount of text per document in the NHTSA dataset is much larger than the amount of text per document in the reddit dataset, the reddit dataset takes much longer to compute.

Referring back to the Script Map, either workflow is decomposed into four steps: Data Ingest, Data Reform, Graph Kernel Similarity, and Clustering. This is the basic work flow developed here. There are some things to note in the Script Map.

First, vertex labels need corrected to work with the `{graphkernels}` package. This is done through reassignment of all text vertex labels to be a scalar value. The labels are mapped from the union of two graphs' labels to values ranging from 1-n, where n is the size of the text label set from the union. This script is run within the `convert_to_graphs.R` script, so it operates on two graphs at a time as it iterates through the whole graph set of concern.

Second, the scripts routed through `{furrr}` are functions which are parallelized in the `convert_to_graphs.R` and `compute_kernel.R` scripts. These parallelized scripts are done so through a function being executed with no argument, so the dataset which they consider is hard coded into the script. This is not good from a reproducibility aspect, but it is well documented in the code and in the repository, and has reproducible results.

Third, the `splitdata.R` function is there to simply split up the dataset from reddit, which would be too large to store in the repository on GitHub. This segmentation of the data still allows for workflow and reproducibility, since the datasets can all be read in to the environment all the same.

Lastly, a KDE clustering analysis was not performed on the NHTSA dataset, as that dataset has few observations ( $N = 48$ ) and will not perform well with the density estimation methods developed here; there would be too few observations to generate a number of local maxima/minima needed to apply the clustering methods developed.

# CHAPTER 4: RESULTS AND ANALYSIS

---

Hierarchical Clustering was performed on both of the reddit and NHTSA datasets. These clustering methods lend themselves to be used with similarity/distance matrices, and the results shown here can be displayed in a variety of formats. One can learn more about the clustering results through dendrograms, n-gram graphs, term-frequency/inverse document-frequency.

## 4.1 NHTSA Special Crash Investigations

For the NHTSA’s Special Crash Investigation dataset, a hierarchical clustering was performed on the resulting graph kernel computed from the skip-gram graphs. Since this dataset was pretty small at  $N = 48$ , it served as a good subject to study the interactions between modifying the skip window width, kernel hyper parameters, and number of clusters. This study on these interactions informs decisions made on which hyper parameters to use in final analysis of the NHTSA data as well as the reddit thread data.

For the study on these interactions, hyper-parameters were varied and the value of cluster within sum of squares was tracked. Specifically, the values in Figure X were tested.

In Figure X, we see some promising values occurring at more prominent “elbow” points in the plots. These points will be further analyzed and compared. The hyper parameters for these points are:

These hyper-parameter sets, and their corresponding graph kernel matrix, are used in a hierarchical clustering analysis. The results from the analysis perform well, and

## Hierarchical Clustering Variation

Varying Skip-gram window (column) and  
Gaussian RBF parameter  $\sigma$  (row)

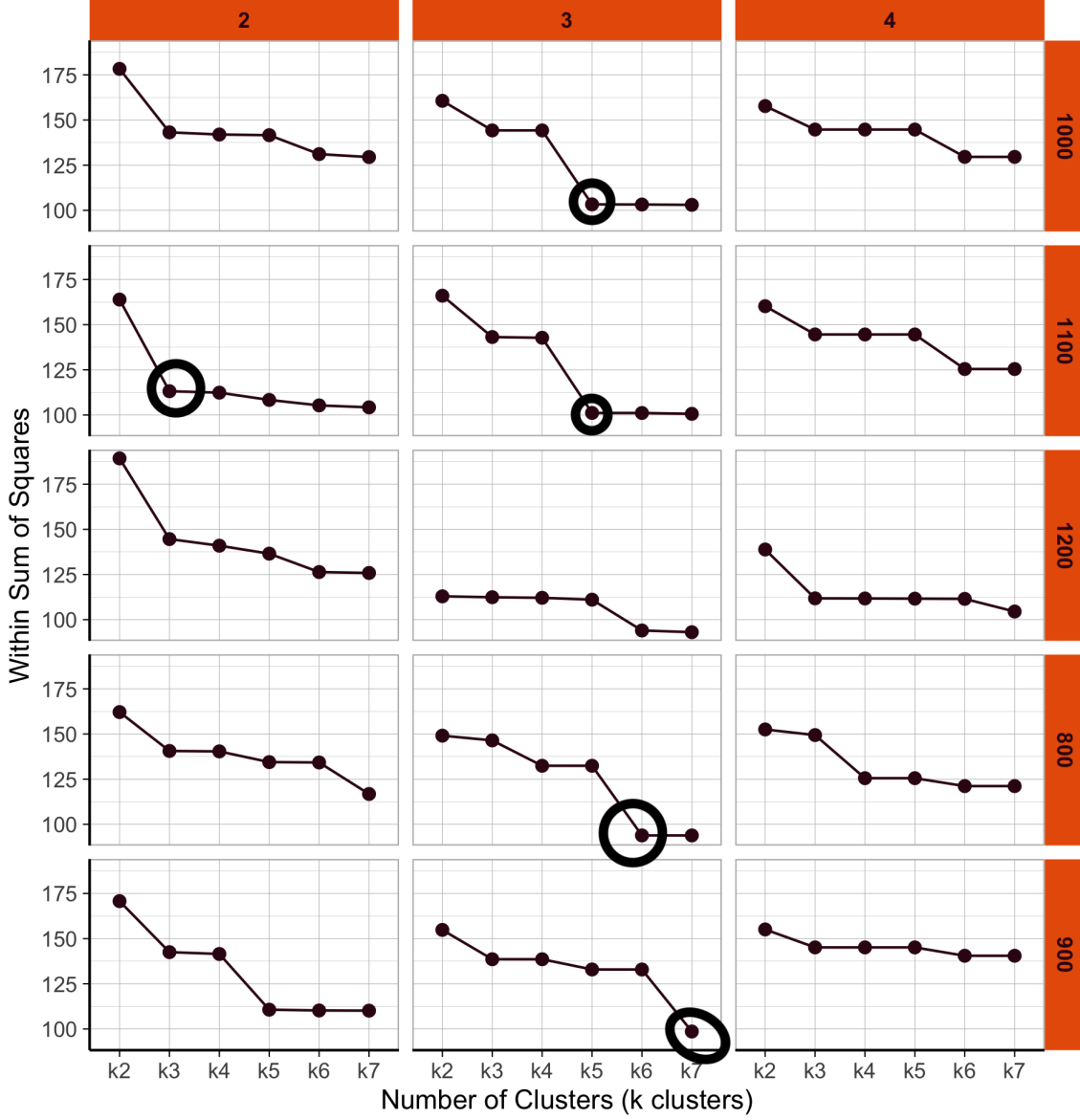


Figure 4.1: Hierarchical Clustering Variation (NHTSA) for differing hyperparameters.

Point	Skip-gram k	Graph Kernel sigma	No. of Clusters
A	3	1000	5
B	2	1100	3
C	3	1100	5
D	3	800	6
E	3	900	7

Table 4.1: Hyperparameter for Variation Study in Figure X

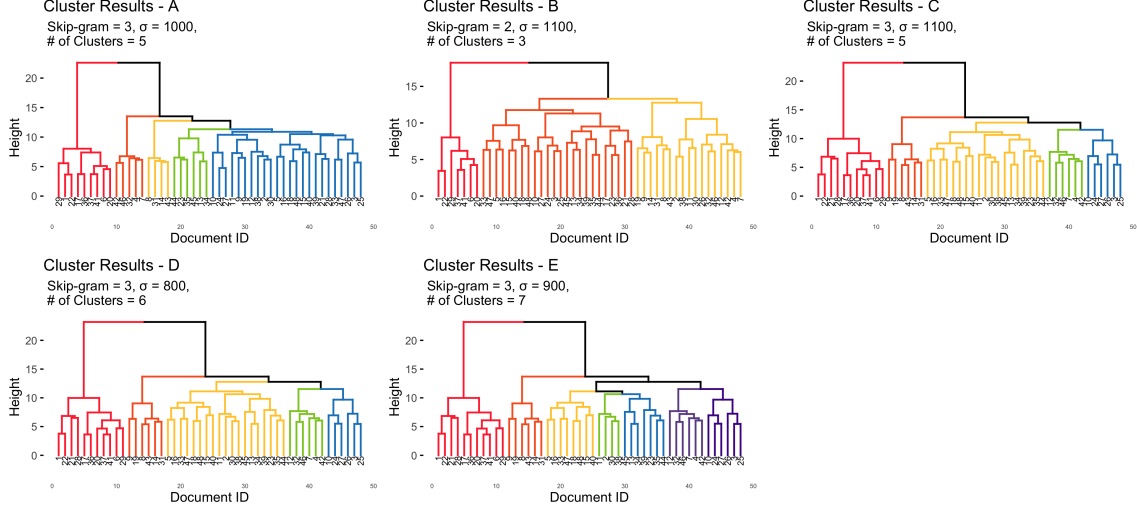


Figure 4.2: Five hyper-parameter sets and their corresponding dendrograms.

display prominent clusters, often of comparable sizes. In Figure X, we see the 5 dendrograms produced from hyper-parameter sets A, B, C, D, and E. Upon some examination, one will notice that documents often appear in the same clusters together, regardless of the hyper-parameter sets. This is most encouraging, as it indicates that these clusters are not a result of user-chosen parameters, but that the documents are similar.

To compare how often this co-membership in groups was appearing, a matrix of co-membership was created. In Figure X, we see that each document has a set of other documents with which it is always clustered with, regardless of hyper-parameter configurations. These small document sets are the document sets which we can expect to have the most in common with one another.

We can inspect the results by checking some of the co-members' text. For example, document two (which was an ambulance crash in Angola, Delaware), has only 3 other co-members. Document two's co-members are: document 11, document 30, and document 38. Document 11 was an ambulance crash in Sheridan, Indiana. Document 30 was an ambulance crash in Pasadena, California. Document 38 was an ambulance crash in Ocilla, Georgia.

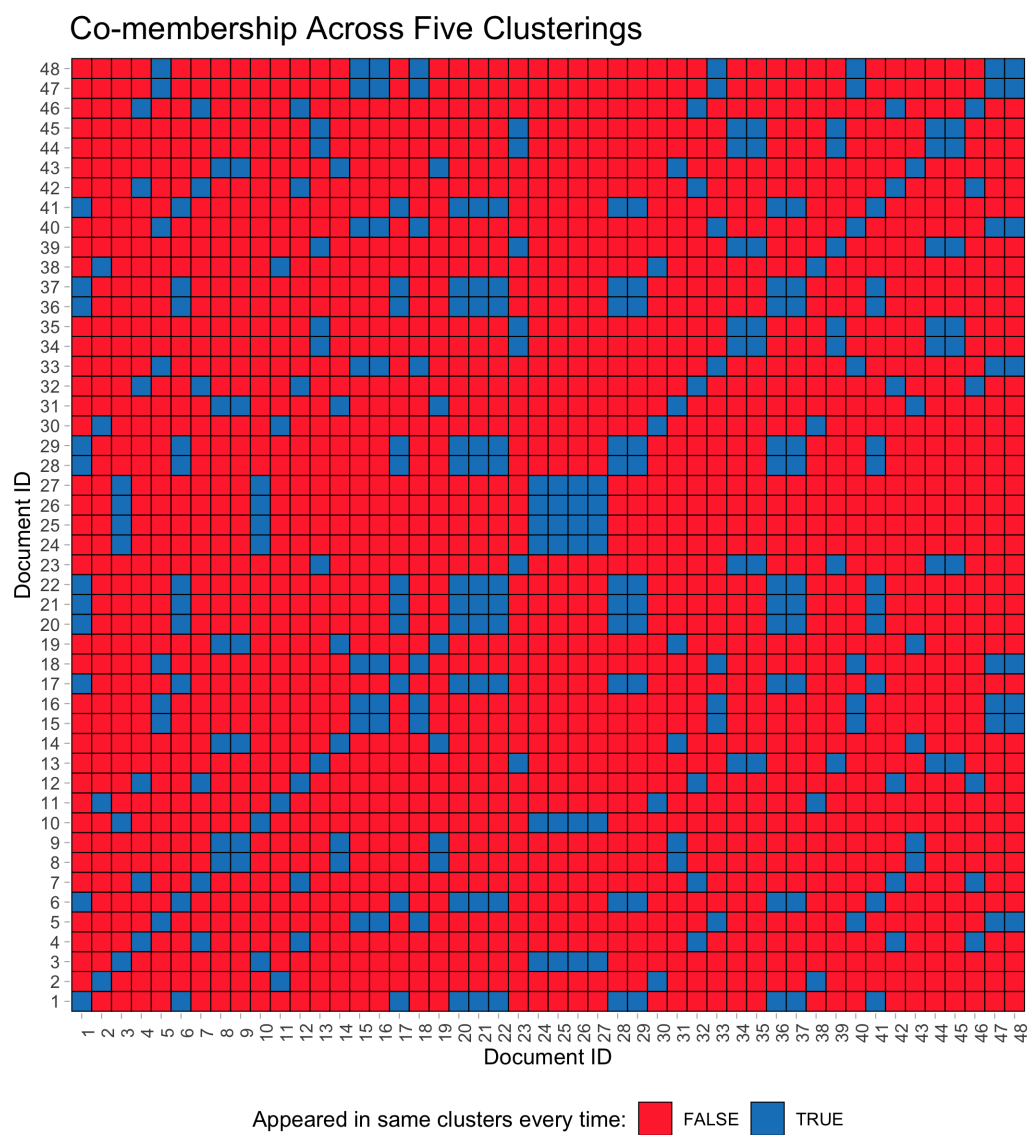


Figure 4.3: Co-membership of documents across hyper-parameter sets A,B,C,D, and E.

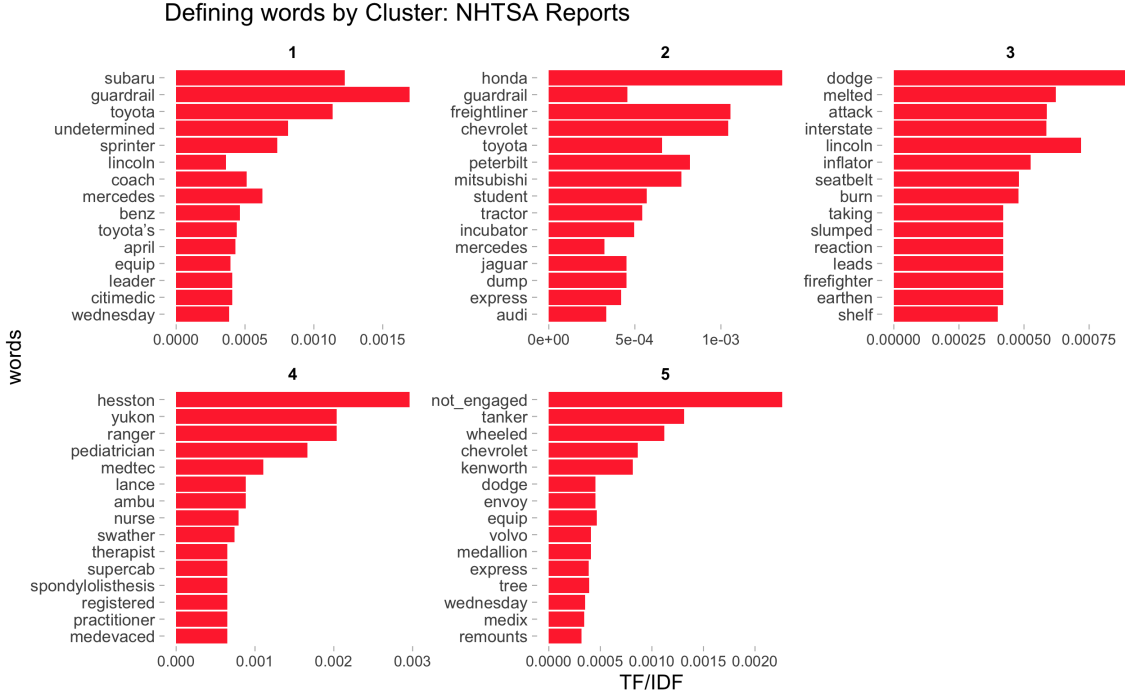


Figure 4.4: Term-Frequency/Inverse Document-Frequency for NHTSA reports, using hyper parameter set A.

Now, we can examine some information about these four incidents that may explain their clustering. For example, all of these incidents had fatalities, three of them had roll over events, and these four were all front end collisions to the ambulance where they struck another vehicle, or object, head on. While this is just one example of the clustering methods identifying similarities in the text, we can examine similarities across the dataset through use of term-frequency/inverse document frequency.

In Figure 4.4, we see that words that appear frequently, and are unique to that cluster appear to be vehicle manufacturers, medical terms, and words that describe a crash situation (e.g. guardrail, tree, interstate). We can use these to get an idea of the crash situation. Alternatively, we can also examine portions of the skip-gram graph produced, which was used for the graph kernel calculation.





We can put together some of the ideas and topics discussed in the cluster. For example, in Figure 4.5 we see that "Front Left Passenger", "Rear Facing Seat", and "Patient Compartment" are all skip-grams which appeared frequently in this cluster.

## 4.2 reddit Threads

For the reddit threads dataset, analysis was focused on the ability of the clusterings to organize threads based on their query words. The words used to query the MentalHealth subreddit all returned threads which were then clustered into new groups. Some of these clusters expressed strong preference for posts that corresponded with specific key words.

To assess how many clusters to use, an "elbow blot" was formed similar to those in Figure 4.1. For the reddit threads, the hyper parameters from the NHTSA analysis was used, and the consistently lowest metric parameter set was chosen; skip-grams of  $k = 3$  and kernel parameter  $\sigma = 1200$  were used in the reddit clustering analysis. In Figure X, we see that 4 clusters is the optimal value, where additional clusters provides diminishing results.

After computing the kernel and using these parameters, four cluster are produced. Examining the dendrogram, in Figure X, we see that we have 2 very distinct groups, which then break into two smaller groups. These groups are very clear and provide confidence in the clustering results. Along side the dendrogram, there is a heat map of what proportion of documents (or threads) in the cluster came from a query word. In Figure X, we see that clusters 3 and 4 collected posts from the "angry" query word. Additionally cluster 1 displays a strong amount of its posts coming from the "sad" query word.

The prominence of some of these query words in the clusters indicates that the tone or language used by the thread authors was picked out by the clustering methods. This

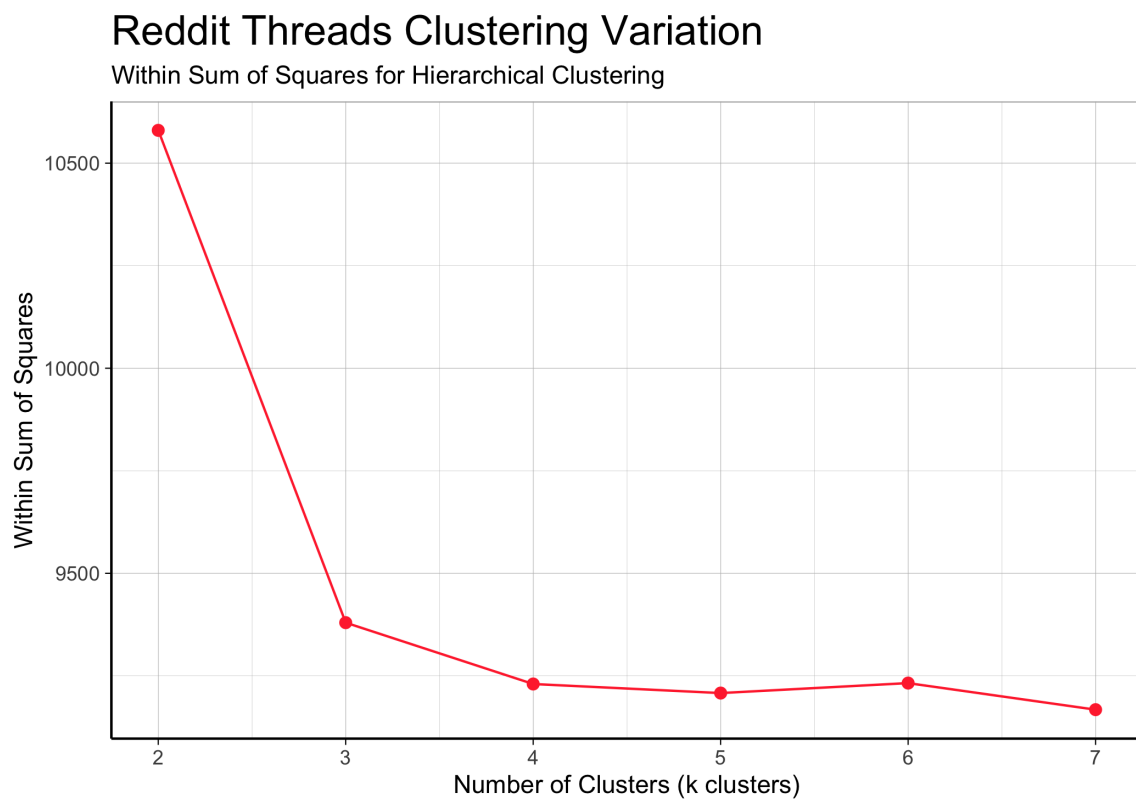


Figure 4.6: Within sum of squares as a function of number of clusters

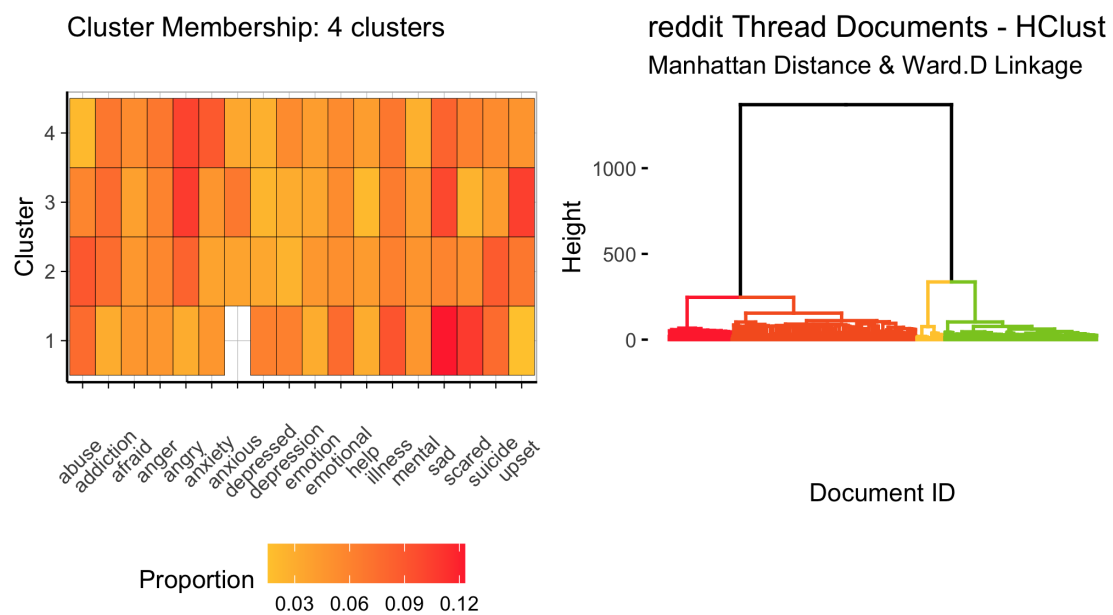


Figure 4.7: Left - Proportion of cluster's threads coming from each query word. Right - Dendrogram for reddit thread analysis, with 4 clusters.

is of particular interest, as the goal of these methods was to perform unsupervised text clustering which utilized more of the rich intricacies of language.

### 4.3 Comparative Analysis

Looking at the results from both hierarchical clusterings on the graph kernels it is apparent that these methods work much better on larger documents. The performance of the graph kernels being applied to the skip-gram graph and then clustering was much better in the NHTSA dataset than the reddit threads dataset. By measuring within sum of squares, the NHTSA data set performed much better, and the computation time was much smaller despite the documents being substantially longer than the reddit comments. However, the reddit dataset displayed more prominent clusters, which was also indicative of success.

Comparing the elbow plots for both the NHTSA and the reddit threads datasets, it is clear that the optimal number of clusters for these methods will be dependent upon the dataset being analyzed. The results from the NHTSA study that indicated that skip-grams with  $k = 3$  did seem to apply to the reddit dataset as well. These hyper parameters will need tuned for any new application.

Overall, the largest difference between the two analyses was that although the reddit threads displayed strong, prominent clusters, the NHTSA dataset was being computed in reasonable amount of time. Since the application of these methods are scaling better for increasing document size instead of increasing document lists, there is a case to be made that aggregating reddit threads into larger documents may be more effective, but to keep the observational unit tied to a reddit thread, this was not completed.

These methods could be applied in any smaller text mining application. Since the methods were able to be parallelized during their most computationally expensive portions of code, the methods could scale with sufficient hardware. In a high performance computing environment, with a much larger dataset, it may become clear that these methods are either more effective or less effective than they were found to be

here.

# CHAPTER 5: CONCLUSION

---

Though applying several different methods in conjunction with one another (skip-grams, graph kernels, and clustering), unsupervised clustering can be done on text while preserving as much of the rich information from the original text as possible.

The NHTSA results, and their clusters, can be used to inform other research initiatives of similar scenarios for ambulance crashes. Identifying these clusters, providing summarized results, like n-gram graphs and term frequency/inverse document frequency lists, allows for information about similar crash reports to be compiled. This compiled information can then be used to generate more "edge-case" scenarios for Florida Polytechnic University's Autonomous Mobility Institute and their work on autonomous vehicles; rare cases are not always available in the datasets used to train autonomous driving models, and this kind of research could help inform their decision making moving forward.

The reddit thread dataset, was just too large for these methods to be applied to all observations in the dataset. This problem will need to be addressed in future work. The large samples ( $N = 1000$ ) used in the study here were accelerated through use of parallel programming and made this analysis possible. The results of the analysis however, were prominent clusters were formed. The dendrogram in Figure 4.7 displays very encouraging results, being that the clusters not only stood out from one another, but that the clusters showed that they were collecting some threads according to their query word that was used to gather the original posts/threads.

In summary, the methods outlined here will continue to develop into more mature methods that could be applied to a wide range of natural language processing tasks, being that the graph kernel matrix could be used for a variety of machine learning tasks—not just hierarchical clustering. One these methods are rebuilt into a clean package for use in high performance computing environments, they will see use in

industries that process large amounts of text data.

Lastly, building upon the published work from 2020, these methods could be used for intelligent navigation of comments and posts on social media. In a previous paper, "A Framework for Intelligent Navigation Using Latent Dirichlet Allocation on Reddit Posts about Opiates", smart ways to cluster text and how that could be used to filter harmful, or otherwise unwanted, content from a web user were outlined. These same principles could be applied with the methods applied to the reddit threads here. Additionally, the NHTSA analysis will benefit Florida Polytechnic University's AMI, which study autonomous vehicles. The fringe cases which are often absent from training datasets, such as ambulance crashes, can be used to better inform their efforts to improve autonomous vehicles. The many ambulance accidents can be clustered so that instead of recreating 50 ambulance incidents, they can select one or two that represents each cluster.

These methods are all available, publicly, on a GitHub repository at <https://github.com/Levi-Nicklas/GraphDocNLP>. Updates may be made and the software written can be used, adapted, recreated, scrutinized, and more by the software development and data science communities that use GitHub. This open sourcing of the software leave the work conducted here open to examination, and will lead to better science for the community.



# CHAPTER 6: FUTURE WORK

---

Future developments on this project would be focusing on adapting the methods to scale on high performance computing (HPC )environments. The methods currently do not do well as the number of graphs for comparison grows. While the methods do fairly well as the size of the graph scales up, the issue of many graphs needs to be addressed for application at scale.

Additionally, supervised learning methods could be explored on these datasets. Use of support vector machines in conjunction with these methods is well documented. Once the datasets are labeled, it would be an easy study to complete. The NHTSA data set, with only 48 observations, could be labeled with a number of different dependent variables that could be gleaned from the crash reports. The reddit thread dataset, as large as it is, may be a good candidate for a semi-supervised application. The clustering results from this study could be applied in a supervised application on reddit threads from the same subreddit, and efficacy could be assessed.

Lastly, these methods we used with edge histogram kernels exclusively, because they are computationally cheaper than other methods, and exploration of the other kernel methods which utilize edge/vertex labels would be an important study to conduct.

# Appendices

# APPENDIX A: MY FIRST APPENDIX

---

```
reddit_thread_kernel_matrix <- matrix(data = rep(0, 1000*1000), nrow = 1000)
```

```
for(i in 1:length(reddit_thread_kernel)){  
  reddit_thread_kernel_matrix[i,] <- reddit_thread_kernel[[i]]  
}
```

```
rm(reddit_thread_kernel)
```

---