

# EFFECTS OF DIFFERENT PARAMETERS ON THE RESULTS OF TWITTER COMMUNITY DETECTION ALGORITHMS

A Partial Thesis Document  
Presented to  
the Faculty of the College of Computer Studies  
De La Salle University Manila

In Partial Fulfillment  
of the Requirements for the Degree of  
Bachelor of Science in Computer Science

by

FERNANDEZ, Ryan Austin  
POBLETE, Clarisse Felicia M.  
SAN PEDRO, Marc Dominic  
TAN, Johansson E.

Charibeth K. CHENG  
Adviser

November 2, 2016

## **Abstract**

In the past few years, social media has become an avenue for users to build relations with friends and share thoughts and opinions. These thoughts, opinions, and relations eventually give rise to communities of users, which have more in common within the community relative to outside the community. These communities are not readily visible in social media; community detection is necessary to extract them. Community detection in social media has become a popular field of study because of the large volume and variety of data that could be extracted. This study aims to determine which combination of community detection algorithms and similarity parameters can produce appropriate communities on Twitter and to produce a visualization of the communities. Multiple iterations of community detection will be performed for each selected combination of algorithm and parameter. This also means that, using the generated visualization, meaningful data can already be extracted, which can be used by businesses or politicians to determine patterns in their target audience.

**Keywords:** Social networks, community detection, sentiment analysis

# Contents

<b>1</b>	<b>Research Description</b>	<b>1</b>
1.1	Overview of the Current State of Technology . . . . .	1
1.2	Research Objectives . . . . .	2
1.2.1	General Objective . . . . .	2
1.2.2	Specific Objectives . . . . .	2
1.3	Scope and Limitations of the Research . . . . .	3
1.4	Significance of the Research . . . . .	4
1.5	Research Methodology . . . . .	4
1.5.1	Preparation . . . . .	4
1.5.2	Iterative Experimentation . . . . .	5
1.5.3	Analysis and Finalization . . . . .	7
1.6	Calendar of Activities . . . . .	8
<b>2</b>	<b>Review of Related Literature</b>	<b>9</b>
2.1	Community Detection . . . . .	9
2.2	Similarity Parameters . . . . .	15
2.2.1	Sentiment Analysis . . . . .	15
2.2.2	Other Parameters . . . . .	16

2.3	Community Evaluation Metrics . . . . .	18
<b>3</b>	<b>Theoretical Framework</b>	<b>22</b>
3.1	Data Collection . . . . .	22
3.2	Community Detection . . . . .	23
3.3	Similarity Parameters . . . . .	27
3.3.1	Sentiment Analysis . . . . .	27
3.3.2	Other Parameters . . . . .	28
3.4	Community Evaluation Metrics . . . . .	29
<b>4</b>	<b>The System Model, Algorithm, and Design</b>	<b>30</b>
4.1	System Overview . . . . .	30
4.2	System Objectives . . . . .	30
4.3	System Scope and Limitations . . . . .	31
4.4	Architectural Design . . . . .	31
4.5	System Functions . . . . .	34
4.5.1	Select Algorithm . . . . .	34
4.5.2	Select Parameter . . . . .	35
4.5.3	Generate Communities . . . . .	37
4.6	Physical Environment and Resources . . . . .	38
<b>A</b>	<b>Research Ethics Documents</b>	<b>39</b>
<b>B</b>	<b>Resource Persons</b>	<b>40</b>
	<b>References</b>	<b>41</b>

# List of Figures

4.1	UML Diagram of the Model Component of the Proposed System .	32
4.2	Sample Screen of User Story 4.5.1 . . . . .	34
4.3	Sample Screen of User Story 4.5.2 . . . . .	36
4.4	Sample Screen of User Story 4.5.3 . . . . .	37

# List of Tables

1.1	Timetable of Activities . . . . .	8
2.1	Summary of Review of Related Literature . . . . .	19

# Chapter 1

## Research Description

This chapter is an overview of the research undertaken in the field of community detection in social networks. This chapter is divided into four sections which are the current state of the technology, research objectives, scope and limitations, and significance of the research.

### 1.1 Overview of the Current State of Technology

Social media has become much more prevalent in recent years. People can now participate in what is called microblogging, a way for people to share their thoughts, status, and opinions in short posts, like Twitter, where posts are limited to one-hundred and forty characters (Java, Song, Finin, & Tseng, 2007). As such, these social media platforms are a prime opportunity to mine sentiments and to detect communities in the social network. Social media can also be used to build user networks through users following other users, for example, in Twitter, creating a network of accounts.

Community detection is clustering multiple users into groups where users within the group are more similar than users from outside the group (Tang & Liu, 2010). Community detection is necessary because it observes the interaction of multiple users in relation to specific similarity parameters such as followed users, post topics, and mentioned users.

Numerous studies on community detection have already been done. Zhang, Wu, and Yang (2012) defined features such as text content similarity, URL similarity, hashtag similarity, following similarity, and retweeting similarity that can

be used to identify similarity between two nodes and aggregated these similarities to detect communities. Their software provided the listings of users in each community. K. Lim and Datta (2012) performed an inverted version in which they first defined interests and based on these interests, sought to extract communities from the network by identifying users that follow the top six celebrities, users with more than 10,000 followers, relating to the given interest. Their software generated the size and clustering coefficient of each community

Individual opinions of one specific user towards another were also studied by West, Paskov, Leskovec, and Potts (2014) combining sentiment analysis, using an L2-regularized logistic-regression classifier, with network analysis, inferring one user's opinion of another by analyzing their common links with other users. Their output was the area under the curve of the receiver operating characteristic (AUC/ROC) and the precision-recall curves (AUC/negPR).

In addition to these works, some visualizations have already been created such as SocialHelix by Cao, Lu, Lin, Wang, and Wen (2015), which uses the temporal extent of social communities, topics or events discussed, and user responses to topics and events to classify users into two sides of the argument. They then depict the two sides of an argument as strands in a double helix and their intersections defines events.

However, most of these studies already had preselected features and algorithms to work with or had specific features in mind before implementing their models. It would be interesting to compare the quality of communities produced by different combinations of algorithms and similarity parameters.

## **1.2 Research Objectives**

### **1.2.1 General Objective**

To determine which combination of community detection algorithm and similarity parameter produces the best communities

### **1.2.2 Specific Objectives**

1. To build a corpus of social media data;
2. To determine the various techniques and algorithms in detecting communi-



ties;

3. To determine the parameters/features to be used in detecting the communities;
4. To determine how to evaluate the correctness of the detected communities given a particular combination of algorithm and similarity parameter;
5. To implement a tool for the visualization of detected communities using the gathered information

### 1.3 Scope and Limitations of the Research

In order to perform a study on community detection, it is necessary to build a corpus of social media data. This research will cover searching for Application Programming Interfaces (API) that will allow extraction of data from Twitter and then using these APIs to build a body of data where community detection can be performed on. Data will include posts, profile information, and network information such as following list, follower list, and group membership.

Different techniques have been used in community detection. Among these techniques are the Infomap algorithm and the speaker-listener label propagation algorithm (Deitrick & Hu, 2013). This research will consider algorithms used in other community detection works, including the Markov stability model, clique percolation method, k-means clustering, and divisive hierarchical clustering.

Before the selected algorithms are implemented, it is necessary to identify which parameters/features/attributes indicate one user's similarity to another. Inquiry will be done to identify these parameters, specifically how to extract them based on the raw data. The research will be limited to sentiment analysis and elements which can be extracted from a user's post, which may include follow networks, hashtags, mentions, and retweets, which were mostly inspired from literature which focused on Twitter (Deitrick & Hu, 2013; Zhang et al., 2012; K. Lim & Datta, 2012).

After community detection, it is necessary to determine whether the detected communities are sensible. Inquiry will be done to find appropriate metrics in determining the accuracy of detected communities. These algorithms will include average mutual following links per user per community or FPUPC (Zhang et al., 2012), modularity (Deitrick & Hu, 2013), and clustering coefficient (K. Lim & Datta, 2012).

After inquiring about multiple community detection algorithms and similarity parameters, it will be necessary to implement the selected algorithms in a working system. The community detection model will then be implemented and augmented by a visualization tool which will be created that will only consider communities detected from data gathered from Twitter.

## **1.4 Significance of the Research**

Community detection is already a widely researched topic in the field of computer science. Our study will contribute to that field by exploring a domain that is not a frequently explored which is a comparative analysis of combinations of algorithms and parameters.

Our study can also be a very useful tool in the domains of viral marketing and political endorsement. This means that companies and governments may benefit from this research. Interested companies may use the result of our study to improve their sales and marketing. The government may use our study to gauge public opinion on certain issues and to see analytics about which geographical areas have a particular opinion.

## **1.5 Research Methodology**

This chapter details the research activities to be done for the duration of this thesis. Our study will be in three main parts: the preparation phase, the iterative experimentation phase, and the analysis and finalization phase.

### **1.5.1 Preparation**

This phase constitutes the gathering of information pertinent to the study. This includes reviewing related literature and building a theoretical framework. The review of related literature step will only take two weeks for the initial bibliography. The theoretical framework step will overlap with the review of related literature step and possibly take two more weeks. This step will include gathering implementation details for the variables in the experimentation phase: community detection algorithms found in the review of related literature; computing similarity parameters, including sentiment analysis and features in Twitter; and

evaluation metrics.

This phase includes finding the necessary API's to use in collecting data from Twitter. Selection of the platform to host the data and a programming language to implement the model is also part of this phase's activities.

This phase will take up the first month of the study. It is necessary because it provides the theoretical framework around which the entire study will be based on, as well as deciding the platform and programming language to be used throughout the study. All design and implementation performed in the project will be based on the information gathered in this phase.

### **1.5.2 Iterative Experimentation**

This phase deals with design, implementation, and testing of multiple community detection models. Each iteration would differ based on two variables: a different similarity parameter and a different community detection algorithm. Each iteration will take two to three weeks to perform, depending on how many of the aforementioned variables differ and how different the given variable's implementation details are from the previous iteration's.

#### **Similarity Parameter Selection**

Since this study aims to produce an accurate visualization of communities in social media, it is necessary to find out which parameter would produce the best communities based on the evaluation metric found in the preparation phase. This is the reason why there will be multiple iterations.

Based on the RRL and the Theoretical Framework, a similarity parameter will be selected from the features identified in the preparation phase.

#### **Community Detection Algorithm Selection**

In producing the visualizations of communities, it is also necessary to select the best community detection algorithm for the data, which is determined by using the evaluation metric on the produced communities. Each iteration would then deal with a different combination of similarity parameter and community detection algorithm.

Based on the RRL and the Theoretical Framework, a community detection algorithm will be selected. This algorithm must be compatible with the selected similarity parameter.

## **Data Collection**

User data will be collected from Twitter. The API's will be used to gather the data. This will be done in order to have a corpus of data to perform the algorithms on. Afterwards, each user will be anonymized. This includes the username and the real name. This anonymization is done to preserve the terms and conditions of using public data extracted from social media.

Any other data transformation necessary for the parameter or the algorithm will then be performed.

## **Model Design**

The proponents will design a model for the selected algorithm and similarity parameters. For the first iteration, this step should also include designing the model for the evaluation module i.e., the module which will evaluate the detected communities. This should take three to four days depending on whether the algorithm or the parameter was used in a previous iteration. This is done to organize the given algorithm and feature in a model that is ready for implementation.

## **Model Implementation**

The majority of the iteration will involve implementing the given model in the language. This should take one to two weeks depending if the algorithm or parameter has already been implemented in a previous iteration. For the first iteration, this step should also include the implementation of the evaluation metrics. This is done in order to have a working model that can be run on the collected data and tested for accuracy.

## **Model Evaluation**

For the next two to three days of the iteration, the model will be run on the collected data for Twitter and the evaluation module will be run on the detected

communities in order to measure the communities' accuracy. This is done in order to evaluate how the selected parameter and algorithm performs on the collected data, which will be comparable at the end of the study to the results from other iterations, allowing the proponents to select which parameter-algorithm pair produces the most accurate communities in which particular social network.

## **Documentation**

For the last one to two days of the iteration, the proponents will finalize the documentation of the iteration. Note that documentation should have been done regularly throughout the iteration, but this step is to ensure the quality and correctness of the documentation. This step will also include a retrospective on what worked in the previous iteration, what did not work, and how the development process can be improved. This step is done in order to ensure integrity in the data and documentation as well as to constantly improve the development process during the duration of the study.

### **1.5.3 Analysis and Finalization**

This phase will involve revisiting the data collected from the multiple iterations and selecting which combination of parameters and algorithm resulted in the most accurate communities. This phase may include supplementary research in an attempt to see why some combinations produced better communities than others to have a more thorough understanding of the results. Finally, the proponents will produce a visualization using the best parameter-algorithm combination, satisfying the objective of the study. This phase should take two to three weeks, mirroring the steps of one iteration in the previous phase. This step is necessary in order for the information gathered in this study to be presentable and to have a tangible output based on the results of the study.

## 1.6 Calendar of Activities

Table 1.1 shows a Gantt chart of the activities. Each bullet represents approximately one week worth of activity.

Table 1.1: Timetable of Activities

Activities (2016-2017)	Jul	Aug	Sept	Oct	Nov	Dec	Jan	Feb	Mar	Apr	May	Jun	Jul	Aug
Preparation (RRL)	••• _•••	••• •••_											•• --	• ----
Preparation (TF)	••• _•••	••• •••_	--••											
Iterative Experimentation				••••	••••	•• _•••	•••	••••	••••	••••	••••	••••	•• --	
Analysis and Finalization													•• --	• ----

# Chapter 2

## Review of Related Literature

This chapter discusses the features, capabilities, and limitations of existing research, algorithms, or software that are related or similar to the proposed research. This chapter is divided into three sections. The first section discusses algorithms for community detection. The second section discusses algorithms for sentiment analysis and other similarity parameters used for community detection. The third section discusses community evaluation metrics.

### 2.1 Community Detection

Clauset, Newman, and Moore (2004) presents an improvement to community detection algorithms. Other existing algorithms presented in this paper are correct logically, but are computationally expensive and cannot run over extremely large datasets in a reasonable amount of time. This paper presents an algorithm that is identical in terms of output but is significantly faster than existing algorithms in terms of runtime.

The algorithm is based on the greedy optimization of modularity, which is a property that indicates the strength of division of a network into communities. The higher the modularity, the better the community structure. The algorithm finds the largest modularity  $Q$  that would result from merging two arbitrary communities, and merging those two communities. The algorithm's main offering is that it skips calculating  $Q$  when the two communities have no edges between them. offering an increase in performance.

The algorithm was run against purchase data from amazon.com. The data

graph worked on was quite big, with 409,687 items and 2,464,630 edges. The algorithm was able to successfully structure the data into communities based mainly on purchasing information. The proponents were successful in discovering clear communities that correspond to specific topics or genres of books or music, indicating that the purchasing tendencies of Amazon customers are strongly correlated with subject matter.

The proponents hope that this algorithm will allow datasets with millions of vertices and tens of millions of edges to be processed using current computing resources in an efficient manner (Clauset et al., 2004).

Tang and Liu (2010) wrote a textbook on community detection where they introduced characteristics of social media, reviewed representative tasks of computing with social media, and illustrated associated challenges because with the emergence of social media websites, they felt it was an avenue to study human interaction and collaboration on an unparalleled scale. Multiple community detection approaches were discussed such as node-centric, network-centric, and hierarchy-centric.

Node-centric algorithms involve the maximum clique detection problem which involves searching for a maximum complete subgraph of nodes in a network graph that are all adjacent to each other. The clique percolation method can find overlapping communities by finding cliques of size  $k$ , and then producing a clique graph, wherein two cliques are connected if they share  $k - 1$  nodes. Each connected element in this clique graph is then a community.

Network-centric algorithms involves vertex similarity, which is the similarity of the node's social circles or how many common friends the two nodes have. This is what structural equivalence deals with. Nodes that are structurally equivalent belong to a community.

Hierarchy-centric algorithms come in two forms: divisive and agglomerative. In divisive clustering, the entire set of nodes starts out in one set and each time, each set is divided into two until each community only has one member. The division is done by finding the node with the lowest edge betweenness and removing it, since that node is most likely the node connecting two communities. Agglomerative clustering starts with each node in their own community and communities are joined if they increase the overall modularity of the set of communities.

These algorithms may be considered in the final community detection phase of the proposed research (Tang & Liu, 2010).

Lancichinetti, Radicchi, Ramasco, and Fortunato (2011) details an approach to



performing community detection across a network. The authors of this paper argued that while there already exists a large variety of techniques for achieving this, there is still a need for more in-depth techniques that can handle different types of datasets, and the “subtleties” of community structure. This paper presents a technique called OSLOM (Order Statistics Local Optimization Method), which they claim to be the first method capable to detect clusters in networks accounting for edge directions, edge weights, overlapping communities, hierarchies and community dynamics. They claim that the algorithm performs just as well as other existing ones, and have been applied on several real networks. It is also freely available for anyone who wants to use it (Lancichinetti et al., 2011).

K. Lim and Datta (2012) aimed to detect communities that share common interests on Twitter, based on linkages among followers of celebrities representing an interest category because they wish to help markets identify target groups with common interests. However, their approach differs from the typical paradigm of “identify communities then, for each, identify interests”, for they first identified interests they wish to extract communities from and from these interests, they then extracted the communities.

Given this set of fans common to the most popular celebrities in the specific interest,  $P$ , they used the Infomap Algorithm and the Clique Percolation Method to detect communities in  $P$ . For each interest they wish to extract communities for, they chose the top 6 most popular celebrities based on follower count. Google and Wikipedia were used to identify which interests a celebrity represents. Afterwards, all users that follow the 6 celebrities were selected. They then selected 200,858 random users as a control group. Their algorithm produced more communities and larger communities than the control group, as well as more consistent communities, having a higher clustering coefficient.

K. Lim and Datta (2012) provides an interesting alternative to detect communities by first specifying the interest of the community before detecting it, which may be used in the proposed research on top of the usual algorithms.

Papadopoulos, Kompatsiaris, Vakali, and Spyridonos (2012) discussed methods of community detection in social media, comparing different aspects of specific methods, and discussing possible incremental applications of these methods. This survey aims to address two main elements left unaddressed in existing related survey articles, namely performance, in terms of aspects such as computational complexity and memory requirements, as well as the interpretation of results of community detection by social media applications.

Classes of community detection methods discussed include cohesive subgraph discovery, vertex clustering, community quality optimization, divisive, and model-

based methods. Cohesive subgraph discovery comprises of methods that require the specification of certain structural properties that must be satisfied in order for a subgraph of the network to be considered a community. Vertex clustering makes use of traditional data clustering methods. Community quality optimization methods focus on the optimization of a graph-based measure of community quality. Divisive methods make use of identified edges and vertices in a network. Model-based methods consider dynamic processes or statistical models in the process of detecting communities.

The survey led to conclusions about the concept and structure of communities in the context of social media, to a rough classification of existing community detection methods, and to determining which methods are most appropriate for social media mining applications.

Papadopoulos et al. (2012) provides an overview of existing methods of community detection as well as a comparison of these methods, which may aid the proponents in deciding the appropriate methods to use for the proposed research.

Xie (2012) details a study into two topics of social network analysis, namely opinion dynamics and community detection. In terms of community detection, one of the main difficulties presented is particularly challenging in large-scale networks. It presents an algorithm called Speaker-listener Label Propagation Algorithm (SLPA) for fast overlapping community detection. Another challenge presented was detecting communities in dynamic networks where changes happen often and in real-time, as with most real world applications. An algorithm for incrementally updating communities instead of profiling each snapshot is presented as well, called LabelRankT. This algorithm is claimed to drastically outperform existing detection algorithms, with similar results (Xie, 2012).

Zhang et al. (2012) sought to identify communities in Twitter based on common interests. This is because Twitter has become very popular recently but little is known of it in the user level. This study would help in user recommendation and tweet recommendation as well as viral marketing to specific target groups. To identify the communities, they first compute specific feature similarities, then aggregate these features to compute for the final user similarity, and then they used classical clustering algorithms to detect the communities. To identify the communities, they first compute specific feature similarities, then aggregate these features to compute for the final user similarity, and then they used classical clustering algorithms to detect the communities.

The specific features they used were textual contents. Each data point was the entirety of a users tweets. Latent Dirichlet Allocation was used to identify latent topics from the users tweets. URL similarity was also detected, finding which users

share similar links. Hashtag similarity was also analyzed. The social structure of users was also analyzed, which includes following similarity and retweeting similarity.

In aggregating these similarities, the weighted sum of the previous similarities was computed to get the final similarity. Finally, k-means clustering was used to detect the communities based on their computed similarities.

Zhang et al. (2012) presents a possible framework for detecting the communities. The k-means clustering algorithm will be one of the proposed algorithms for use in the proposed research.

Deitrick and Hu (2013) sought to use sentiment classification to analyze communities in Twitter because harvesting information from these online social networks (OSN) would aid in the fields of politics and marketing.

Their process is as follows: The follower network was represented as a weighted directed graph, each with initial weight of 1. To augment this, replies, mentions, retweets, hashtags, and sentiment classification of tweets were also harvested. These factors adjusted the weights in the graph. For community detection, the Infomap algorithm and Speaker-listener Label Propagation Algorithm(SLPA) were run.

Generally, the network with updated weights produced communities with greater modularity. Of the two algorithms, the Infomap algorithm performed better. Recurring sentiment analysis was also helped by performing the aforementioned algorithms on the accounts that have already been placed in detected communities, which permits more in-depth analysis into the users sentiment since it could be analyzed within the context of the detected community.

Deitrick and Hu (2013) provides a new way to represent the network, with their updated weights, as well as more possible algorithms to consider in detecting communities.

Bakillah, Li, and Liang (2015) sought to contribute to the field of extracting relevant information from social media by detecting geo-located communities in Twitter in disaster situations. The main disaster they focused on is the occurrence of typhoon Haiyan in the Philippines.

Social graphs of Twitter users related to the focus are created by comparing Twitter's different interaction nodes like follow relations, mentions and tweet content. The fast-greedy optimization of modularity (FGM) clustering algorithm enhanced with semantic similarity is used in order to handle the complex social graphs created. Modularity measures the quality of divisions of a network into

communities. By maximizing the modularity between the generated graph structure and a random graph structure, the optimal clustering results can be obtained.

Together with FGM, the varied density-based spatial clustering of applications with noise spatial (VDBSCAN) clustering algorithm is used to get spatial communities at different time periods. This is done to divide thematic communities discussing same topics formed by the FGM algorithm into more meaningful sub-clusters. The discovery of geo-located communities could potentially help in identifying and locating incidents occurring during emergency situations.

Bakillah et al. (2015) provides algorithms that could prove useful in getting the optimal clustering for detecting communities. It also gives an insight in considering the spatial and thematic properties of these communities.

Amor et al. (2015) sought to detect communities and to identify roles in the Twitter network regarding on the subject of the care.data debate using graph-theoretic methods, one of them being the Markov Stability method. There are two networks constructed from the obtained data relating to the care.data debate: follower network and retweet network. The flow-based community detection method Markov Stability was used for identifying interest communities in the follower network which resulted into a 13-way partition composed of 4 large communities and 9 minor ones. It is also used for the retweet network in order to find conversation communities which resulted into 8 communities.

The Markov Stability method works on the behaviour of dynamical processes on a network. This potentially reveals meaningful structure about the graph. It can extract different coarse-grained descriptions of the graph at different time scales. In addition, this method can find non-clique communities.

Amor et al. (2015) gives another community detection method that can be used for the proposed research. Its main advantage is its scalability especially over time scales.

Cao et al. (2015) proposed a visual analysis system, SocialHelix, because social media is a grand avenue for people to express their opinions and the researchers believed that an intuitive visualization that unfolds the process of sentiment divergence would have a far-reaching impact on multiple domains.

They first identified the key domain problems of social divergence before employing a data abstraction design to convert the raw data into a form that captures all the key factors of the aforementioned domain problems. This abstracted data is then represented in a visualization based on a visual DNA metaphor. In identifying the key domain problems, it is determined when divergences start and end,

how they evolve, who is involved, what roles do they play, and why does divergence occur. In the abstraction phase, the raw data is decomposed into temporal extent of social communities, topics or events, and user responses to these topics or events. In the visualization phase, the opposite sides of the helix represent the two sides of a divergence. The helix curves represents the changes in the communities sentiment. Nucleobase pairs represents events that connect the two communities.

In implementation, the data was first filtered, removing unrelated posts and people. Statistical Linguistic Sentiment Analysis was used to determine the users sentiment. Finally, clustering was done using Hadoop, producing a cluster with 30 nodes.

In the end, all test users were impressed by the visualization and agreed with the researchers model for the visualization. All test users felt that divergence identification was made easy due to the visualization.

Cao et al. (2015) gives a sample visualization that can inspire the proponents own visualization, albeit the goal is not to model single divergences but an entire community. A possible tool, Hadoop, is also mentioned, which may be used to cluster the data.

## **2.2 Similarity Parameters**

This section outlines the basis/features/parameters used in detecting communities. It is divided into two subsections. The first subsection deals solely with sentiment analysis. The second subsection deals with other network and node parameters not related to sentiment analysis.

### **2.2.1 Sentiment Analysis**

Zhang et al. (2012) provided a formula to determine similarity in terms of text. This research provides a metric to determine the similarity of two users in terms of post content, which can be used in the proposed research (Zhang et al., 2012).

Bryden, Funk, and Jansen (2013) aimed to determine whether or not members of identified communities had similar word usage and language features on their social media posts. This was done through the analysis of 75 million mutual tweets among 189 thousand Twitter users. This study focused on the connec-

tion of language has to network structure, in order to explore the potential of understanding society through analysis of communication on social media. Communities were characterized through the words used in messages sent by members of the community; the most representative words from each community were identified through the Z-scores of each words usage. The Euclidean distances between word usage frequencies for each pair of communities was the basis for determining how significant the differences between these communities word usages were. The research determined that there were many similarities in words, word fragments and word lengths among tweets from users in identified groups, including word usage that was not related to subject matter. Through language structure alone, the researchers were also able to determine a users' network communities. This research focused on the detection of communities through language used on social media. As it involves on community detection on social media as well, the proposed research may make use of the approach presented in this research (Bryden et al., 2013).

Deitrick and Hu (2013) used a subjective/objective and positive/negative Naive Bayes classifier. To do this, all tweets were converted to lowercase; hashtags, usernames, urls were replaced with twitterhashtag, twitterusername, and twitterurl respectively; the tweet text was tokenized; repeated punctuation was replaced with the + sign e.g. “!!!” would become “!+”; sentence punctuation was split into separate tokens; non-sentence punctuation was removed. Ten-fold cross validation was used in training the classifier. Weights in the graph mentioned in section 2.1 were then updated if two users posted something with a similar sentiment and similar hashtag. This research shows a clearly defined process in performing sentiment analysis, particularly the data cleaning step. This process could be adapted for the proposed research (Deitrick & Hu, 2013).

Bakillah et al. (2015) enhanced the FGM algorithm with a similarity measure. A threshold  $T$  for text similarity is used to determine whether two communities are similar enough to increase the priority of merging them. 0.2-0.3 was used as the value of  $T$ . The cosine similarity measure is used to compute similarity between the communities' set of terms. This measure can be used as a means for getting the similarity between different communities' set of words when merging similar communities will be relevant to the proposed research (Bakillah et al., 2015).

## 2.2.2 Other Parameters

Zhang et al. (2012) provided a few formulas to determine similarity in terms of URL, hashtag, following, and retweeting similarity. This research gives formulas

that can be used in the proposed study to measure similarity as well as a means to aggregate similarities from multiple parameters. (Zhang et al., 2012).

Bakillah et al. (2015) created graphs with weighted edges with similarities based on Twitter's various interaction modes. This study presents alternative formulae that may be used in the proposed research (Bakillah et al., 2015).

Amor et al. (2015) includes some sentiment analysis on the tweets they handled, particularly on negative tweets as these comprised most from sample they took. They divided the concerns of these negative tweets into 3:

1. Implementation - concerns regarding information provision, the opt-out process, and communication with the public
2. Scheme concept - concerns about privacy, sharing of personal data, and the use or sale of the data
3. Execution - Concerns around security, effectiveness of pseudonymisation, and cyber attacks

No formula or representation was given as to how tweets were categorized between these 3 concern categories. However, this opens up the idea of having specific parameters related to the focus of the community detection, in this case with regards to the care.data debate, instead of general parameters concerning the social site's interaction modes (Amor et al., 2015).

Darmon, Omodei, and Garland (2015) aimed to present an approach to community detection that is multifaceted, focusing not only on structure-based communities, but on other types as well, namely activity-based, topic-based, and interaction based communities. Communities can be defined similarly or differently according to these types, so in order to come up with a more accurate and dynamic picture of a community, all types of communities, as well as the overlaps among these communities, should be taken into account. This study was done through the analysis of a Twitter dataset in order to assign representative weights for each community type. Activity-based communities were derived through the timing of users'tweets, topic-based communities were derived from hashtag similarities, and interaction-based communities were derived from retweets and mentions.

For topic-based communities, edges on the network of users and followers are weighted depending on the number of common hashtags between each user and follower pair. Interaction-based communities are defined by three weighting schemes. The first scheme considers the number of tweets follower  $f$  retweeted from user

$u$ . The second scheme considers the number of tweets wherein user  $u$  mentions follower  $f$ . The third and final scheme takes the arithmetic mean of the mentions and retweets.

Darmon et al. (2015) determined that the multifaceted approach to community detection could aid in better understanding the structure of online communities and in finding communities in social media that would otherwise be hidden. His study provides an approach that may be considered as well as algorithms that may be used in the detection of communities in the proposed research (Darmon et al., 2015).

## 2.3 Community Evaluation Metrics

Zhang et al. (2012) used the average number of mutual following links per user per community (FPUPC) to evaluate their communities. Based on this, appropriate weights for the aggregation were found by first performing their k-means clustering algorithm using only one feature similarity for each of the similarities and extracting the FPUPC. Afterwards, they gave each feature similarity a weight based on a formula. The number of clusters,  $k$ , used in the k-means clustering algorithm was also tweaked to get the maximum FPUPC. They concluded that they were successful in generating relatively accurate communities due to the incrementally increasing FPUPC after adjusting the weights. This provides one possible evaluation metric for the proposed research, as well as a method to provide weights for the feature similarities that the proponents will eventually be using for community detection (Zhang et al., 2012).



Table 2.1 shows a summary of our review of related literature with respect to community detection, similarity parameters, and evaluation metrics for each paper.

Table 2.1: Summary of Review of Related Literature

Reference	Community Detection Algorithms	Sentiment Model	Analysis	Other parameters	Community Evaluation
(Clauset et al., 2004)	Greedy Optimization of Modularity				
(Tang & Liu, 2010)	Clique percolation method, similarity divisive and agglomerative clustering				
(Lancichinetti et al., 2011)	Order Statistics Local Optimization Method				
(K. H. Lim & Datta, 2012)	Topic driven community detection, Infomap method, Clique percolation method				
(K. Lim & Datta, 2012)	Topic driven community detection, Infomap method, Clique percolation method				
(Papadopoulos et al., 2012)	Comparison of Existing Methods				

(Xie, 2012)	Speaker-listener Label Propagation Algorithm, Label- RankT			Correlations between different snapshots of the network over time	
(Zhang et al., 2012)	k-means clustering	Similarity Formula for Text	Similarity Formula for URL, Hashtag, Fol- lower, and Retweeting		FPUPC metric
(Bryden et al., 2013)		Characterization of communities through word usage			
(Deitrick & Hu, 2013)	Weighted directed graph, Infomap Algorithm, SLPA	Subjective / Objec- tive, Positive / Nega- tive Naive Bayes Clas- sifier	replies, mentions, retweets, hashtags		
(Bakillah et al., 2015)	enhanced fast-greedy optimization of mod- ularity (FGM) algo- rithm with similar- ity measure, varied density-based spatial clustering of applica- tions with noise spa- tial (VDBSCAN) al- gorithm	cosine similarity mea- sure	mentions, follow rela- tions, shared URLs, Tweet content		
(Amor et al., 2015)	Markov Stability			care.data debate - im- plementation, scheme concept and execution	

(Cao et al., 2015)	Data abstraction design, Hadoop tool	Temporal extent of posts, topics and events, user responses, Statistical Linguistic Sentiment Analysis		
(Darmon et al., 2015)			Activity-based communities, Topic-based communities, Interaction-based communities	

## Chapter 3

# Theoretical Framework

This chapter discusses the algorithms, formulas, and existing procedures with regards to our study. This chapter is divided into four sections. The first section discusses data collection libraries. The second section discusses algorithms for community detection. The third section discusses algorithms for sentiment analysis and other similarity parameters used for community detection. The fourth section discusses community evaluation metrics.

### 3.1 Data Collection

To construct a graphical representation of the users and their social connections (friends lists), a custom crawler was designed by (Ferrara, 2012). The crawler, written in Java, made use of the Apache HTTP library (<http://httpd.apache.org/apreq>) to acquire data about Facebook users and their friends. One simple acquisition method used by this paper was the Breadth-First Search (BFS) algorithm. Starting from a root node (user), the nodes public friend list was obtained, and then for each friend, the algorithm is repeated. The algorithm terminated if a.) the third depth-level of friendship was completely covered, or b.) the algorithm exceeded a runtime of 10 days without termination. Data was then anonymized to hide references to the actual users.

In a paper on Community Detection in Twitter, Java et al. (2007) says that through the Twitter Developer API, it is possible to obtain time-bound data of the social network of users. Through the API, a directed graph  $G(V,E)$  was constructed, representing users ( $V$ ), and their following networks ( $E$ ). Public information per user was obtained, like profile information and location (using the

Yahoo! Geocoding API).

The proponents were also able to obtain a crawler for Twitter written by a fellow researcher conducting a similar study (Lam, 2016). The crawler, written in Python, utilizes the Tweepy framework (<http://www.tweepy.org/>), which is a wrapper for easily using the Twitter Developer API with Python. The crawler can be fed various search terms, and will return tweets matching those search terms.

## 3.2 Community Detection

Clauset et al. (2004) presents an improvement to community detection algorithms. Other existing algorithms presented in this paper are correct logically, but are computationally expensive and cannot run over extremely large datasets in a reasonable amount of time. This paper presents an algorithm that is identical in terms of output but is significantly faster than existing algorithms in terms of runtime.

The algorithm is based on the greedy optimization of modularity, which is a property that indicates the strength of division of a network into communities. The higher the modularity, the better the community structure. Given a graph with  $n$  vertices and  $m$  edges, the algorithm defines two quantities, namely:

$$e_{ij} = \frac{1}{2m} \sum_{vw} A_{vw} \delta(c_v, i) \delta(c_w, i) \quad (3.1)$$

which is the fraction of edges that join vertices in community  $i$  to vertices in community  $j$ , and

$$a_i = \frac{1}{2m} \sum_v k_v \delta(c_v, i) \quad (3.2)$$

which is the fraction of ends of edges that are attached to vertices in community  $i$ . The algorithm then defines the modularity  $Q$  as:

$$Q = \sum_i (e_{ii} - a_i^2) \quad (3.3)$$

The algorithm finds the largest  $Q$  that would result from merging two arbitrary communities, and merging those two communities. The algorithm's main offering is that it skips calculating  $Q$  when the two communities have no edges between

them, offering an increase in performance. The algorithm maintains a matrix  $\Delta Q_{ij}$  for each pair  $i, j$  of communities with at least one edge between them, a max-heap  $H$  containing the largest  $Q$  for each row of the matrix, and a vector array with elements  $a_i$ .  $\Delta Q_{ij}$  is defined as follows:

$$\Delta Q_{ij} = \begin{cases} \frac{1}{2m} - \frac{k_i k_j}{(2m)^2} & \text{if } i, j \text{ are connected,} \\ 0 & \text{otherwise} \end{cases} \quad (3.4)$$

and  $a_i$  as:

$$a_i = \frac{k_i}{2m} \quad (3.5)$$

The algorithm runs as follows: (1) Calculate the initial values of  $\Delta Q_{ij}$  and  $a_i$ , and populate the max-heap with the largest element of each row of the matrix  $\Delta Q$ . (2) Select the largest  $\Delta Q_{ij}$  from  $H$ , join the corresponding communities, update the matrix  $\Delta Q$ , the heap  $H$ , and  $a_i$ , and increment  $Q$  by  $\Delta Q_{ij}$ . (3) Repeat step 2 until only one community remains.

The algorithm was run against purchase data from amazon.com. The data graph worked on was quite big, with 409,687 items and 2,464,630 edges. The algorithm was able to successfully structure the data into communities based mainly on purchasing information. The proponents were successful in discovering clear communities that correspond to specific topics or genres of books or music, indicating that the purchasing tendencies of Amazon customers are strongly correlated with subject matter.

The proponents hope that this algorithm will allow datasets with millions of vertices and tens of millions of edges to be processed using current computing resources in an efficient manner (Clauset et al., 2004).

Node-centric algorithms involve the maximum clique detection problem which involves searching for a maximum complete subgraph of nodes in a network graph that are all adjacent to each other. The clique percolation method can find overlapping communities by finding cliques of size  $k$ , and then producing a clique graph, wherein two cliques are connected if they share  $k - 1$  nodes. Each connected element in this clique graph is then a community (Tang & Liu, 2010).

Hierarchy-centric algorithms come in two forms: divisive and agglomerative. In divisive clustering, the entire set of nodes starts out in one set and each time, each set is divided into two until each community only has one member. The di-

vision is done by finding the node with the lowest edge betweenness and removing it, since that node is most likely the node connecting two communities. Agglomerative clustering starts with each node in their own community and communities are joined if they increase the overall modularity of the set of communities. Modularity is given by

$$Q = \frac{1}{2m} \sum_{l=1}^k \sum_{i \in C_l, j \in C_l} (A_{ij} - \frac{d_i d_j}{2m}) \quad (3.6)$$

where  $m$  is the number of edges,  $d_i$  is the degree of node  $v_i$ ,  $C_l$  being the  $l$ th community, and  $A_{ij}$  being the value in the adjacency matrix for node  $v_i$  and  $v_j$  (Tang & Liu, 2010) .

Bakillah et al. (2015) sought to contribute to the field of extracting relevant information from social media by detecting geo-located communities in Twitter in disaster situations. The main disaster they focused on is the occurrence of typhoon Haiyan in the Philippines.

Social graphs of Twitter users related to the focus are created by comparing Twitter's different interaction nodes like follow relations, mentions and tweet content. The fast-greedy optimization of modularity (FGM) clustering algorithm enhanced with semantic similarity is used in order to handle the complex social graphs created. Modularity measures the quality of divisions of a network into communities. By maximizing the modularity between the generated graph structure and a random graph structure, the optimal clustering results can be obtained. This modularity is expressed through the quality function  $Q$ :

$$Q = \sum_{c=1}^n \left[ \frac{l_c}{m} - \left( \frac{d_c}{2m} \right)^2 \right] \quad (3.7)$$

where  $n$  is the number of clusters,  $m$  is the total number of edges,  $l_c$  is the total number of edges joining the vertices of cluster  $c$  and  $d_c$  is the sum of the expected random graph degree of the vertices of  $c$ . To achieve the largest quality change ( $\Delta Q$ ), communities are progressively merged. The usage of quality function  $Q$  and the merging of communities are the core steps of the FGM algorithm. Now, the FGM algorithm was enhanced by integrating graph based and text-based (text similarity measure) models to balance the importance of shared content versus graph structure. This will be discussed more in detail in section 3.3.1.

Together with FGM, the varied density-based spatial clustering of applications with noise spatial (VDBSCAN) clustering algorithm is used to get spatial communities at different time periods. This is done to divide thematic communi-

ties discussing same topics formed by the FGM algorithm into more meaningful sub-clusters. The discovery of geo-located communities could potentially help in identifying and locating incidents occurring during emergency situations.

The aim of the VDBSCAN spatial clustering algorithm is to find spatial clusters based on regions with higher density. The algorithm is based on two parameters:

1.  $r$ : the value of the radius that will be used to select members of a cluster
2. MinPts: minimal density to form a cluster.

Let  $D$  be the set of points corresponding to the geo-located tweets found in a given thematic community and that were issued during time period  $T$ .

The neighbour of point  $p$  is expressed as:

$$N(p) = \{q \in D | dist(p, q) \leq r\} \quad (3.8)$$

A cluster is generated based on the following properties, called density-reachable and density-connected:

1. A point  $q$  is directly density-reachable from a point  $p$  if

$$q \in N(p, r) \text{ and } |N(p, r)| \geq MinPts \quad (3.9)$$

2. A point  $q$  is density-reachable from a point  $p$  if there exists a sequence of points  $p_1, \dots, p_n$  where  $p_1 = p$  and  $p_n = q$  such that  $p_{i+1}$  is directly reachable from  $p_i$ , for all  $i$ .
3. A point  $q$  is density-connected to a point  $p$  if there is a point  $o$  such that  $p$  and  $q$  are density-reachable from  $o$ .

A cluster is a non-empty subset of  $D$  that satisfies the following properties:

1.  $\forall p, q$ , if  $p \in C$  and  $q$  is density - reachable from  $p$ , then  $q \in C$  (maximality).
2.  $\forall p, q \in C$ ,  $p$  is density - connected to  $q$  (connectivity).



The parameters  $r$  and  $\text{MinPts}$  are optimized automatically based on the variation in density of the data set.

This research provides algorithms that could prove useful in getting the optimal clustering for detecting communities. It also gives an insight in considering the spatial and thematic properties of these communities (Bakillah et al., 2015).

### 3.3 Similarity Parameters

This section outlines the basis/features/parameters used in detecting communities. It is divided into two subsections. The first subsection deals solely with sentiment analysis. The second subsection deals with other network and node parameters not related to sentiment analysis.

#### 3.3.1 Sentiment Analysis

Zhang et al. (2012) provided a formula to determine similarity in terms of text.

$$sim_{text}(i, j) = \frac{1}{\sqrt{D_{js}(i, j)}} \quad (3.10)$$

$D_{js}$  is the Jensen-Shannon Divergence between the two users topic probability distribution given by

$$D_{js}(i, j) = \frac{D_{kl}(UT_i || M) + D_{kl}(UT_j || M)}{2} \quad (3.11)$$

where  $M = \frac{UT_i + UT_j}{2}$  and  $D_{kl}(P || Q) = \sum_{i \in topics} P(i) \log \frac{P(i)}{Q(i)}$  and  $UT_i$  is the probability distribution of user  $i$  for all topics.  $UT_i[t]$  is the probability distribution for user  $i$  on topic  $t$  (Zhang et al., 2012).

Bakillah et al. (2015) enhanced the FGM algorithm with a similarity measure. A threshold  $T$  for text similarity is used to determine whether two communities are similar enough to increase the priority of merging them. 0.2-0.3 was used as the value of  $T$ . The cosine similarity measure is used to compute similarity between the communities' set of terms:

$$Cosinesimilarity = \frac{A \cdot B}{||A|| ||B||} = \frac{\sum_{k=1}^l (A_k \times B_k)}{\sqrt{\sum_{i=1}^n (A_i)^2 \times \sum_{j=1}^m (B_j)^2}} \quad (3.12)$$

$A$  and  $B$  represent the frequency of a term in the set of terms associated with the first and second community, respectively. The similarity value ranges from 1, meaning dissimilarity, to 1, meaning exact similarity.

This measure can be used as a means for getting the similarity between different communities set of words when merging similar communities will be relevant to the proposed research (Bakillah et al., 2015).

A practical learning method that can be used for sentiment analysis is called the naive Bayes classifier. The naive Bayes classifier can be used for predicting the classification of a new instance based on classification done beforehand on a set of training data. Before the prediction is to be done, a set of training data will be the basis for the classification. Each instance in the set is described by a conjunction of attribute values and any target value from a finite set  $V$  can be taken by the instance. The attribute values are assumed to be conditionally independent. The goal is to find the most probable target value,  $v_{NB}$ , given the attribute values  $(a_1, a_2, \dots, a_n)$  in order to classify the new instance.

The naive Bayes classifier can be given as:

$$v_{NB} = \arg \max_{v_j \in V} P(v_j) \prod_i P(a_i | v_j) \quad (3.13)$$

where  $v_{NB}$  is the max target value given by the classifier,  $v_j$  is an element of the set of target values in  $V$ ,  $P(v_j)$  is the probability of  $v_j$ , the conjunction is the product of attributes given that  $v_j$  occurred or is the actual target value (Mitchell, 1997). For sentiment analysis, the different sentiments, like positive or negative, or subjective or objective can be used as the set of target values. The words from the training data for the sentiment analysis will be the attribute values for the target value it can take (Deitrick & Hu, 2013).

### 3.3.2 Other Parameters

URL similarity is given by the same formula as text similarity in section 3.3.1, only using links instead of topics (Zhang et al., 2012).

Hashtag similarity is given by

$$sim_{hashtag}(i, j) = \sum_{k=1}^n \left(1 - \left| \frac{N_{ik}}{|H_i|} - \frac{N_{jk}}{|H_j|} \right| \right) \frac{N_{ik} + N_{jk}}{|H_i| + |H_k|} \quad (3.14)$$

where  $N_{ik}$  is the number of times user  $v_i$  used the hashtag  $k$  while  $H_i$  is the total hashtags used by  $v_i$  (Zhang et al., 2012).

Following similarity is given by

$$sim_{follow}(i, j) = \frac{c_{friend}}{\sqrt{|Friend_i| |Friend_j|}} + \frac{c_{follower}}{\sqrt{|Follower_i| |Follower_j|}} \quad (3.15)$$

$|Friend_i|$  is the total number of users  $v_i$  follows.  $|Follower_i|$  is the total number of users that follow  $v_i$ .  $c_{friend}$  represents the two users common friends.  $c_{follower}$  represents the two users common followers (Zhang et al., 2012).

Retweeting similarity is given by

$$sim_{retweet}(i, j) = \frac{c_{retweet}}{\sqrt{|R_i| |R_j|}} + \frac{n_{ij} + n_{ji}}{|R_i| |R_j|} \quad (3.16)$$

$R_i$  is the number of users whom  $v_i$  retweet.  $c_{retweet}$  is the number of users both  $v_i$  and  $v_j$  retweet.  $n_{ij}$  is the number of times  $v_i$  retweeted  $v_j$  and  $n_{ji}$  is the inverse case (Zhang et al., 2012).

### 3.4 Community Evaluation Metrics

The average number of mutual following links per user per community (FPUPC) can be used to evaluate communities (Zhang et al., 2012).

A measure called modularity is a method used for measuring the strength of communities through comparing the strength of connections within a community to the strength of random connections between vertices. (Newman & Girvan, 2004) This makes use of the formula indicated in section 3.2.

Modularities approaching 1 indicate strong community structures; these values usually range from 0.3 to 0.7. However, a modularity of 0 would indicate communities that are only as good as randomly produced ones (Newman & Girvan, 2004).

## Chapter 4

# The System Model, Algorithm, and Design

### 4.1 System Overview

The system is a web application that generates communities from data retrieved from Twitter. On startup, the user will select an algorithm. Once an algorithm is selected, the user then selects a similarity parameter that is compatible with the chosen algorithm. Finally, the user chooses to generate communities from the data. The system displays the graphical representation of the communities as well as the community evaluation metrics that can be applied to the communities based from the chosen algorithm and parameter. The system only has one type of user, which can select an algorithm, a parameter, and generate and view communities.

### 4.2 System Objectives

The system must be able to generate communities using a specific algorithm and similarity parameter. The specific objectives are as follows:

- To collect data from Twitter
- To clean data collected from the social networks
- To detect communities using the supported algorithms and parameters

- To display a visualization of the communities
- To display the evaluation metrics for the communities

## 4.3 System Scope and Limitations

Data collection is necessary in order to detect the communities from the social network. Data collected from the social networks will only include the user posts, the follow network (or friend network if from Facebook), hashtag usage, and mentions.

Cleaning the data is necessary in order to reduce noise in the data as well as to ensure the veracity of the data. This would improve the detected communities.

The only algorithms to be considered in the system are divisive and agglomerative hierarchical clustering and fast greedy optimization of modularity (*FGM*). For the first two algorithms, the only parameters to be considered are the positive/negative valence of posts, the subjective/objective valence of posts, the cosine similarity of frequency of topics mentioned, the follow network, hashtag usage, and mentions.

The visualization is necessary in order to provide a more intuitive representation of the communities. The visualization is limited to a graphical visualization to be displayed via HTML5.

Displaying the evaluation metrics is necessary to show the comparative effectiveness of different algorithm and parameter combinations. The evaluation metrics to be considered are the modularity of the communities and the average mutual following links per user per community (*FPUPC*).

## 4.4 Architectural Design

The system will follow the MVC architecture. The view module will use HTML5, CSS3, and Javascript. The controller and model will be implemented as code.

The model module will include a data collection class, representative models for the data collected, and a computational model. The UML class diagram for the model module is found in figure 4.1.

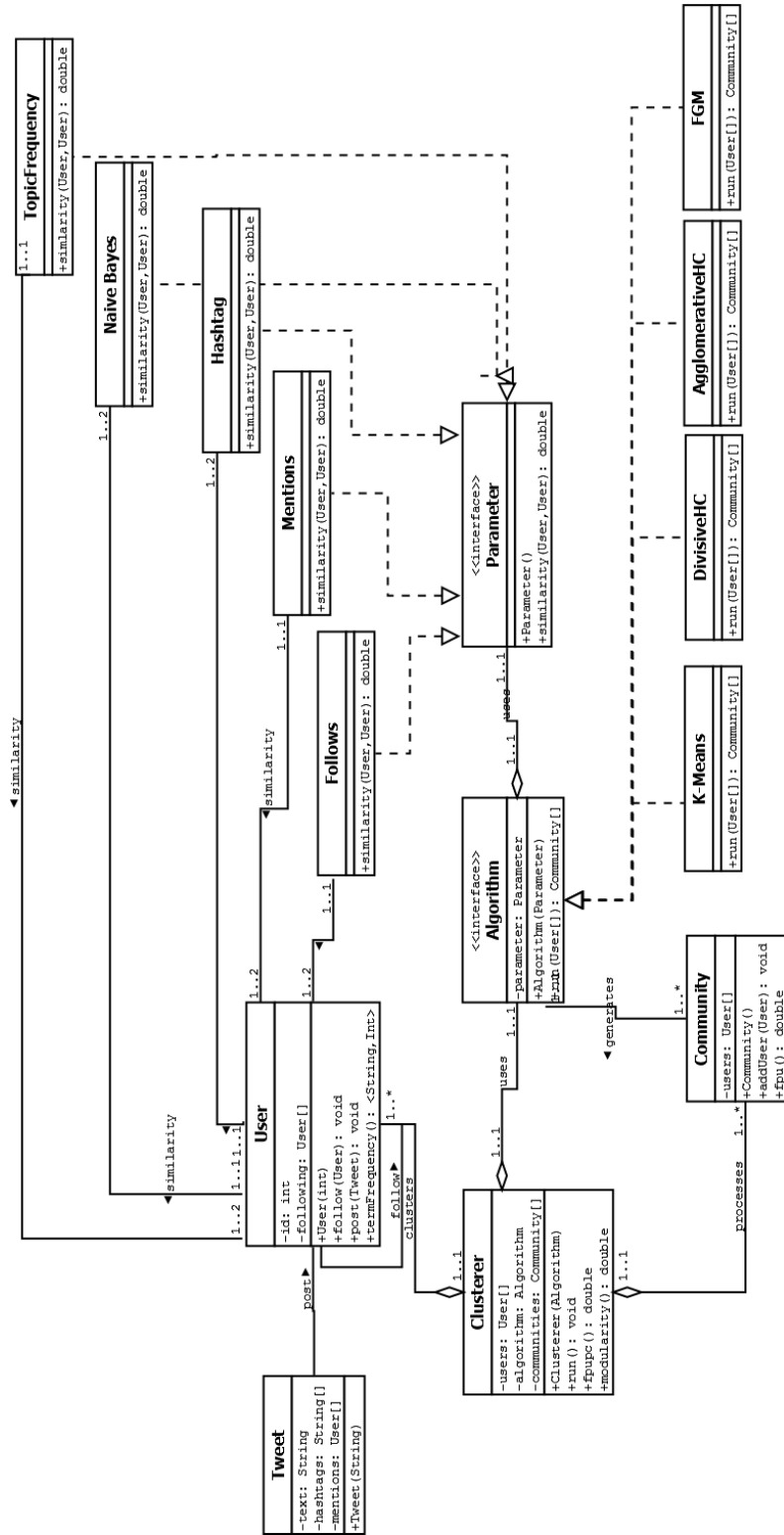


Figure 4.1: UML Diagram of the Model Component of the Proposed System

Data collection is done via the tweepy API. The data collection methodology is as follows - three public twitter accounts with a significant amount of followers will be selected. For each account, the first 1000 most recent retrievable tweets will be obtained. Next, a breadth-first search (BFS) will be done for that account's followers, and their 500 most recent tweets will be retrieved as well.

The `UserCrawler.py` file has a function `get_all_tweets(username)`, which is responsible for retrieving the most recent 500 tweets of a particular user, as specified by the `username` parameter. Initially, three accounts (seeds) will be chosen. Each seed needs to have at least 200 followers. These three seeds are then pushed into a queue. An iterative process then begins, which is described as follows:

At the start of each iteration, a user is popped from the head of the queue. This user's most recent 500 tweets are collected, and the list of users, representing the account's followers, is obtained. For each user in that list, the algorithm checks if it has seen that particular user previously. If not, the user is added to the end of the queue.

This iterative process continues until either a depth of 3 in the BFS tree is reached, or a count of 500,000 tweets is reached.

In parallel with the above algorithm, a separate csv file is also generated. This csv file will contain information about the users' following network. Each row of the csv file will begin with the user id of the particular user, followed by the user id's of all the accounts that user is following.

To anonymize the data, all json elements that may point back to the user are removed. Specifically, these elements are "name", "description", "profile\_image\_url", "profile\_image\_url\_https", "profile\_background\_image\_url\_https", "profile\_banner\_url", and "profile\_background\_image\_url". Additionally, each user id is replaced with an automatically incrementing integer id (starting at 1). This is because the user id can still be mapped to the actual user profile.

For the data representation, this comprises a `User` class and a `Tweet` class.

The computational model will have a central class `Clusterer` which takes a list of users as an attribute, as well as an `Algorithm` object. The algorithm follows the Strategy design pattern, which abstracts the algorithms in concrete implementations of the `Algorithm` interface. The `Algorithm` Interface contains an instance of the `Parameter` interface which also follows the Strategy design pattern. In this case, retrieving the similarities be-

tween two User objects is abstracted, to be implemented by the concrete realizations of the Parameter interface. By setting the algorithm and parameter, the Clusterer objects runs the algorithm and produces a list of Community objects, which have two derived attributes *modularity* and *FPUPC*. The Clusterer caches these Communities and can also return the aggregate *modularity* and *FPUPC* of the set of detected communities.

The controller simply interfaces all the model submodules. It checks on startup if the data collection module is done collecting and cleaning data. It passes the data collected, represented as User and Post classes, to the Clusterer class, after setting an algorithm and parameter. This produces the Community objects that the controller then passes as part of the HTTP response.

The view module will be implemented with HTML5, CSS3, and Javascript.

## 4.5 System Functions

This section outlines the different functions of the system.

### 4.5.1 Select Algorithm

The user selects an algorithm in order to set a required parameter for generating communities. Figure 4.2 shows a sample screen for this user story.

Select an Algorithm ▼	Select a Parameter ▼	Generate Communities
K-Means Divisive Hierarchical Clustering Agglomerative Hierarchical Clustering Fast Greedy Optimization of Modularity		

Figure 4.2: Sample Screen of User Story 4.5.1

Pre-condition: The system has data from Twitter. The main menu is the current screen.



Scenario:

1. The user chooses to select an algorithm.
2. The system displays all the algorithms (K-Means, Divisive Hierarchical Clustering, Agglomerative Hierarchical Clustering, Fast Greedy Optimization of Modularity) supported by the system.
3. The user selects one of the algorithms.
4. The user confirms their choice.
5. The system redirects to the main menu and displays the selected algorithm as part of the system status.

Post-condition: The system has an algorithm to use for generating communities. The “Select Parameter” option is now enabled in the main menu. The allowed parameters have been enabled in the Select Parameter screen.

Acceptance Criteria:

- Test if the algorithm selected is valid.
- Test if only one algorithm is selected.
- Test if the algorithm is displayed as part of the system status.
- Test if the system redirects to the main menu.
- Test if the “Select Parameter” option is now enabled.

## 4.5.2 Select Parameter

The user selects a similarity parameter in order to set a required parameter for generating communities. Figure 4.3 shows a sample screen for this user story.

Pre-condition: The system has finished collecting data from the data collection module. User story 4.5.1 has been finished.

Scenario:

1. The user chooses to select a parameter.

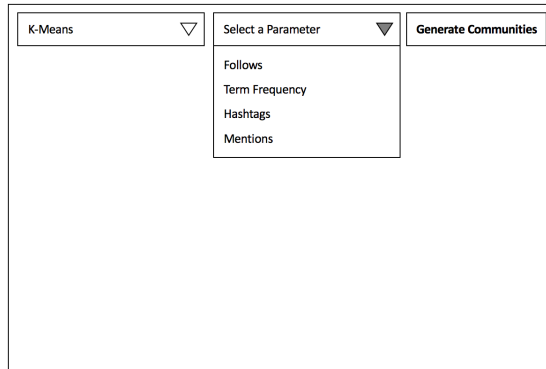


Figure 4.3: Sample Screen of User Story 4.5.2

2. The system displays all parameters compatible with the selected algorithm (Follows, Term Frequency, Hashtags, Mentions).
3. The user selects one parameter to use.
4. The user confirms their choice.
5. The system redirects to the main menu with the “Generate Communities” option enabled.

Post-condition: The system has a similarity parameter to use for generating communities. The “Generate Communities” option has been enabled in the main menu.

Acceptance Criteria:

- Test if the user has already selected an algorithm.
- Test if only the parameters compatible with the selected algorithm are displayed in the “Select Parameter” screen.
- Test if the parameter selected is a valid choice.
- Test if only one parameter is selected.
- Test if the system redirects to the main menu.
- Test if the “Generate Communities” option is enabled.

### 4.5.3 Generate Communities

The user asks the system to generate communities to see the graphical representation of the communities and to see the evaluation metrics of the generated communities for the selected algorithm-parameter combination. Figure 4.4 shows a sample screen for this user story.

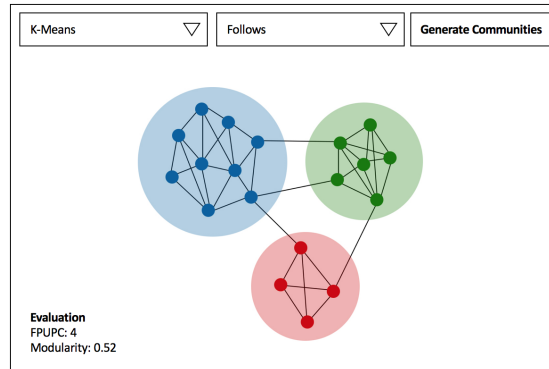


Figure 4.4: Sample Screen of User Story 4.5.3

Pre-condition: The system has finished collecting data from the data collection module. User stories 4.5.1 and 4.5.2 have been finished.

Scenario:

1. The user chooses to generate communities.
2. The system computes for the communities using the selected algorithm and parameter.
3. The system displays the graphical representation of the communities as well as the values of the evaluation metrics.

Post-condition: The system is now displaying a graphical representation of the communities and the value of the evaluation metrics.

Acceptance Criteria:

- Test if the user has already selected an algorithm and parameter.
- Test if the detected communities are appropriate for the chosen algorithm and parameter.
- Test if the evaluation metrics' values are correct.

## 4.6 Physical Environment and Resources

The software was developed using Python as a language and Django as a framework for the web application. Python 3.5 is required to run the code. The web application must be run by a server that supports Django 1.10.1.

# Appendix A

## Research Ethics Documents

This section contains the research ethics documents related to this research proposal.

# Appendix B

## Resource Persons

**Ms. Charibeth Cheng**  
Adviser  
College of Computer Studies  
De La Salle University-Manila  
`chari.cheng@delasalle.ph`

# References

- Amor, B., Vuik, S., Callahan, R., Darzi, A., Yaliraki, S. N., & Barahona, M. (2015). Community detection and role identification in directed networks: understanding the twitter network of the care.data debate.
- Bakillah, M., Li, R.-Y., & Liang, S. H. L. (2015, February). Geo-located community detection in twitter with enhanced fast-greedy optimization of modularity: The case study of typhoon haiyan. *Int. J. Geogr. Inf. Sci.*, 29(2), 258–279. doi: 10.1080/13658816.2014.964247
- Bryden, J., Funk, S., & Jansen, V. A. (2013). Word usage mirrors community structure in the online social network twitter. *EPJ Data Science*, 2(1), 1–9. doi: 10.1140/epjds15
- Cao, N., Lu, L., Lin, Y.-R., Wang, F., & Wen, Z. (2015). Socialhelix: visual analysis of sentiment divergence in social media. *Journal of Visualization*, 18(2), 221–235. doi: 10.1007/s12650-014-0246-x
- Clauset, A., Newman, M. E. J., & Moore, C. (2004, Dec). Finding community structure in very large networks. *Phys. Rev. E*, 70, 066111. doi: 10.1103/PhysRevE.70.066111
- Darmon, D., Omodei, E., & Garland, J. (2015, 08). Followers are not enough: A multifaceted approach to community detection in online social networks. *PLoS ONE*, 10(8), 1-20. doi: 10.1371/journal.pone.0134860
- Deitrick, W., & Hu, W. (2013). Mutually enhancing community detection and sentiment analysis on twitter networks. *Journal of Data Analysis and Information Processing*(1), 19-29.
- Java, A., Song, X., Finin, T., & Tseng, B. (2007, August). Why We Twitter: Understanding Microblogging Usage and Communities. In *Proceedings of the joint 9th webkdd and 1st sna-kdd workshop 2007* (p. 56-65). Springer.
- Lam, A. (2016). Improving twitter community detection through contextual sentiment analysis of tweets.
- Lancichinetti, A., Radicchi, F., Ramasco, J. J., & Fortunato, S. (2011). Finding statistically significant communities in networks. *PLoS ONE*, 6(4), 1 - 18.
- Lim, K., & Datta, A. (2012). Following the follower: Detecting communities with common interests on twitter. In *Proceedings of the 23rd ACM conference*

- on *hypertext and social media (ht12)* (Vol. 1, pp. 317–318). Association for Computing Machinery. doi: 10.1145/2309996.2310052
- Lim, K. H., & Datta, A. (2012). Finding twitter communities with common interests using following links of celebrities. In *Proceedings of the 23rd ACM conference on hypertext and social media (ht12)*.
- Mitchell, T. M. (1997). *Machine learning*. New York, NY, USA: McGraw-Hill, Inc.
- Newman, M. E. J., & Girvan, M. (2004). Finding and evaluating community structure in networks. *Physical Review E Phys. Rev. E*, 69(2). doi: 10.1103/physreve.69.026113
- Papadopoulos, S., Kompatsiaris, Y., Vakali, A., & Spyridonos, P. (2012). Community detection in social media. *Data Mining and Knowledge Discovery*, 24(3), 515–554. doi: 10.1007/s10618-011-0224-z
- Tang, L., & Liu, H. (2010). *Community detection and mining in social media*. Morgan & Claypool. doi: 10.2200/S00298ED1V01Y201009DMK003
- West, R., Paskov, H. S., Leskovec, J., & Potts, C. (2014). Exploiting social network structure for person-to-person sentiment analysis.
- Xie, J. (2012). *Agent-based dynamics models for opinion spreading and community detection in large-scale social networks* (Unpublished doctoral dissertation). Troy, NY, USA. (AAI3533361)
- Zhang, Y., Wu, Y., & Yang, Q. (2012). Community discovery in twitter based on user interests. *Journal of Computational Information Systems*, 8(3), 991–1000.