# JIMMA UNIVERSITY

# JIMMA INSTITUTE OF TECHNOLOGY

# FACULTY OF COMPUTING AND INFORMATION

# DEPARTMENT OF INFORMATION TECHNOLOGY

## SUMMARIZING SENTIMENTFROM SOCIAL MEDIA DISCOURSES GIVEN IN SPORT DOMAIN USING MACHINE LEARNING

### BY:

### REHIMA NEGASH

**Advisor:**   **Dr.TekluUrgessa**

**Co-Advisor: Mr.HailuBashada**

**Jimma University, Ethiopia**

**September, 2022**

# JIMMA UNIVERSITY

# JIMMA INSTITUTE OF TECHNOLOGY

# FACULTY OF COMPUTING AND INFORMATION

# DEPARTMENT OF INFORMATION TECHNOLOGY


## SUMMERIZING  SENTIMENT FROM SOCIAL MEDIA DISCOURSES GIVEN IN SPORT DOMAIN USING MACHINE LEARNING

### By:


### REHIMA NEGASH


**A Thesis submitted to the college of graduate studies of Jimma University in partial fulfillment of the Requirements for degree of Master of Science in Information Technology**


<div align="right">

**September, 2022**

</div>

**JIMMA UNIVERSITY**

**JIMMA INSTITUTE OF TECHNOLOGY**

**FACULTY OF COMPUTING AND INFORMATION**

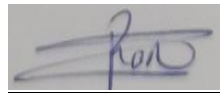**DEPARTMENT OF INFORMATION TECHNOLOGY**

A Thesis Submitted to the School of Graduate Studies of Jimma

University in Partial Fulfillment of the Requirements for the Degree of Master of Science in

Information technology

## Approval sheet

Thesis Title:**Summarizing Sentiment From Social Media Discourses Given In SportDomain**

**Using Machine Learning**          **Submitted by:**

Rehima Negash                                                                    September 23, 2022
Student name                          Signature                              Date
**Approved by:**

### Name and Signature of advisors

Dr. Tekilu Urgesa (PhD)
Principal Advisor                                                             September 23, 2022
                                                                                            Date

Mr. Hilu Beshada (MSc)                                          September 23, 2022
   Co-advisor                            Signature                         Date

### Name and Signature of the Examining Board

_____        _____              _____

Chairman, Examining Board           Signature                            Date

Dr. Kula Kekeba (PhD)                                                     September 23, 2022
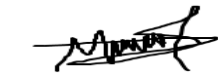
External Examiner                       Signature                            Date

Mr. Mizanu Zelalem(MSc)                                            September 26, 2022

 Internal Examiner                        Signature                            Date

September, 2022

Jimma, Ethiopia

# DEDICATION

I hereby would like to dedicate my thesis to my Mam and Husband. They have provided their unconditional love, support, knowledge, and encouragement throughout this journey. Allah blessed their life in the best of His ways

# ACKNOWLEDGEMENTS

# LIST OF ACRONYMS

AvgMsgOp          Average message opinion

EFF               Ethiopian Football federation

ENFT              Ethiopian National Football team

GUI               Graphical User Interface

NB                Naïve Bayes

NLP               Natural Language Processing

OOG               Opinion Oriented Graph

POS               Part of Speech

SVM               Support Vector Machine

**LIST OF TABLES**

# Table of Contents

**LIST OF FIGURES**

# ABSTRACT

*The sentiment discourse (i.e. responses and replay) on the social media such as, Facebook, Twitter, YouTube, Forums, etc., forms the outbreaks of ample of opinionated thread chains.To drive a complete message, from social Network discourses every single opinionated text under opinion thread has to be seen interdependently. But, the methods of straight forward sentiment mining or computational linguistics being applied, cannot notice the interdependency between two nodes in the absence of opinion oriented graph. In This study, a graph-based opinion summarizing model, whose vertices contain message objects or topic under discussion and its reply nodes that are labeled with opinion polarity is anticipated. The major contribution of this study is the use of back-trace enabled rule based applied on opinion-oriented graph. The total Data set used to undertake this study ware collected from social Network site Facebook, from sport domain and annotated experts. The proposed model extracts the summary of opinions polarity from the corpus of opinion-oriented graph. Hence, it is possible to achieve enhanced decisions by summarized sentiments polarity derived from this graph. Experiments are conducted and have confirmed that the proposed model provided an encouraging result, The result from the model show that the entropy $H(v_x)$ of thread of discussions are between 0 and 0.5 this indicates that the opinion posted were the same type and positive. This reveals that the feedbacks given on sport event are positive feedbacks that encourage National team. However, the graph-based opinion mining model by itself does not automatically identify the orientation of a text. For this, we put forward the automatic sentiment annotation for better performance and the use of separate model for local language for enhanced decisions on domain of the study.*

**Keywords:***Sentiment summary*, ***back-tracing***, ***discourse, opinion oriented-graph; social network; thread analysis***

<div align="center">**CHAPTER ONE**</div>

<div align="center">**INTRODUCTION**</div>

## 1.  Background of the study

With the emerging of the Internet as an important source of information, users are able to express their opinions, emotions, feelings  and their experience toward various topics on social media regarding to products, social, political, sports, religion, businesses and so on [1]. Opinion mining or sentiment analysis is one of the broad concepts in the area of natural language processing and text mining. It is the process of finding users opinion about particular topic from the text expressed as positive, negative and neutral.

In addition to this, it is also the study of semantic relationship between the sentiment expressions of the online users [2]. Online user opinions are two types: Direct opinion and indirect opinion [3]. Direct- opinion gives positive, negative and neutral opinion about the topic directly, whereas indirect opinion is the one that responds, positive or negative opinion depending on the previous opinion given by someone else. For example, for opinion topic "*what do you feel about the new coach?*" if the reply is "*none sense*", it is a direct opinion, whereas, if the reply is "*it is good job*" and" *it is bad to say this*" is indirect opinion*;* here the second two opinions are based on the first opinion. Mining such opinions in online discussions requires an appropriate representation. The task of Opinion mining can be seen in different level of granularity: sentence level, document level and word or feature level [4]. In sentence level, the polarity of each sentence is calculated to identify objective and subjective sentences. The subjective sentences contain opinion words, which help in determining the sentiment about the entity; after which the polarity classification is done into positive and negative classes. In objective sentence the polarity classification is neutral [4]. The document level sentiment classification is concerned with classifying the document based on the overall opinion expressed by the opinion holder as positive, negative and neutral. The Document level analysis is not desirable in forums, Facebook and blogs because the existing data is short text. Feature-based opinion mining also called aspect-based sentiment analysis focuses on the recognition of all sentiment expressions within a given feature (e.g. a product

review) and the features to which they refers **[5].** In this study the researcher uses sentence, word and phrase level direct and indirect opinions text written in, English languages

Currently, Social networks are becoming a great source for opinion discourse analysis. E.g., Facebook is one of the most popular Social Network in the world with more than two billion active users per month **[4]**. It offers the possibility of collecting wealth of information posts appearing in the form of discussions chain, debates, agreement-disagreement discourse **[5] [6]**. As a result, it pays attention for discourse analysis in the area of opinion mining and social networking. The sentiment discourse (i.e. responses and replay) on the social Network sites such as, Facebook, Twitter, YouTube, Forums, etc., forms the outbreaks of ample of opinionated thread chains.

Facebook is the largest social network in the world with over a billion and a half monthly active users (Smith, 2016), and therefore, in this study, we use it as a resource for our dataset and a representative of social media in general. In this platform we observed the reaction text as 16, 43,10,20,25 posts and 22.733, 1.578, 1.434, 36.210, 3.086 Users' Comments that published online regarding Ethiopian National Football team game event from Week 1 to Week 5 of 2021/22 Premier League season respectively.

This ample of information chains appear on such platforms, causes information overload. Hence, it faces difficulty in identifying relevant information for decision making. To get a complete message from threaded opinion appear on the pages, every single opinionated text under these chains has to be seen interdependently. However, it is problematic to get a complete message from such threaded opinion chains, solely by applying the state-of-the art computational linguistic techniques being utilized under opinion mining. In this Study, an opinion-oriented graph-based summarizing model from an opinionated discourse text of social network site is proposed. The major techniques we propose for this study is, the use of back-trace enabled rule based opinion-oriented graph approach and machine learning techniques as input to identify semantic relationship between the sentiment expressions of the online users.

### 1.2. Statement of the Problem

The structure of opinion Information on Social network sites are mainly a kind of discussion oriented or nested opinion. It typically consists of original postings (parent-node) and a plenty of

additional postings (child nodes) which publicly responded to original posts or responses to responses. Hence, as the discussion chain continuously increases it will create complex nested replies. From this, nested opinion text structure making decisions is become quite complex and difficult for individuals to tell the total tendency of a discussion threads solely from state-of-the artwork computational linguistic techniques being utilized under opinion mining. Moreover, it is also time consuming and boring to analyze which opinion is provided with what replay and what are the corresponding opinion relations. Understanding and back-tracing the meaning of such type of opinion thread profoundly require focused attention and more time. This study narrates the overall sentiment summary of the discourse posted on publicly available message about specific posts on specific events. Recently a plenty of discourse analysis research in the area of social networking and opinion mining had been done and are being undertaken using machine learning and social network techniques. From the task of machine learning, Opinion Mining or Sentiment Analysis is a process of automatic extraction of opinion described in the form of positive, negative or neutral about a text in discussion topic [7].

The opinion mining methods are explored in recent Literature mainly rooted in, natural language processing, computational linguistics, and text mining to determine the sentiment polarity of a text at a different level of granularity namely sentence, document, word or phrase and aspect [8], [9], [10], [11]. A typical approach to these methods is to use frequencies of positive and negative words in order to determine whether a discourse is predominantly positive or negative [12], [13]. Such an approach ignores the hierarchical structure of discourse text, whereas this hierarchical structure of discourse posts carries valuable information that tells the interrelationship of posts. As a result, the individual opinion polarity identified by the methods does not carry out complete topic information, out of discourse text. Yet, using knowledge obtained from this hierarchical structure of discourse text is a relatively unexplored direction of discourse sentiment.

**E.g2.** *Consider discourse reviewed regarding Ethiopian national Football team* taken from *https://www.facebook.com/ethiosports/*

1. *Initial post 1. "Ethiopian national Football team has received new coach".*
1.1.   *Comment 1.1. "I am happy with it, it is timely decision."*
1.1.1. *Reply 1.1.1. "I don't think that any change will come by this."*

3

### 1.1.2. *Reply 1.1.2. "sure"*

### 1.1.2.1. *replay 1.1.2.1 "it is bad to say this; the candidate is man of change"*

In this discourse: the reply on the order lists 1.1.1 And 1.1.2 have shown disagreement to the initial Reply 1.1; However, the opinions given in on Reply 1.1.2.1 shows disagreement with the previous two replies, as a result, this reply shows indirectly agreement (positive) with initial post 1. But if we consider this reply message out of discourse it shows disagreement (negative) message. Here, a sequence of a positive message followed by a negative does not necessarily show a negative message between comments and vice versa.

This problem cannot be handled with the state of the art computational linguistic or text mining techniques, as the method does not consider the relationships between opinions orientations from parent to the child node, rather it determines a polarity of the single opinion post. Thus, it is better to explore a new method which examine and determines the total sentiment relationship in a discourse discussion. The new method requires adequate rules offered by experts and annotated text polarity from machine learning. The focus of this study is more of the Summarization of identified opinion polarity by expert and machine learning as input and develop opinion-Summary model that minimize the complexity of opinion thread in discourse. Basically, the method of opinion-oriented graph based analysis requires proper data structure, storage, and representation technique for back-tracing sentiment dependency. It is easy to determine the overall summary of discourse sentiment from a huge amount of collected text from Facebook reviews, prior to making a decision.

Therefore, this study investigates and aims to develop, graph-based sentiment summarizing model from an opinionated discourse text of social network site to have an advantage of both techniques. Based on the above stated problems, the researcher addresses the followings research Questions. These are

- What are the structure and measure techniques an opinions discourse expressed on Facebook Page?
- How can we determine the polarity of sentiment under discussion chain?
- What is the explicit representation and storage of opinion thread structure in the analysis made in social media?

● How can we retrieve the summary of sentiment polarity from threads of discussion to make significant decision?

## 1.3. Objective of the study

### 1.3.1. General Objective

The general objective of this study is to Summarize Sentiment from Social Media Discourses Given in Sport Domain Using Machine Learning.

### 1.3.2. Specific objectives

The specific objectives of the study show the way or techniques that achieves the general objective. The specific objectives of the study are:

▪ To determine the basic opinion measure techniques in opinion-oriented Graph.
▪ To develop a model that shows the work of proposed system
▪ To build the system that computes the summary of opinions from thread structure in terms of positive, negative and neutral.
▪ To evaluate the performance of the model designed with new data set.

## 1.4. Scope of the study

Opinion mining is a complex and challenging field in Social Medias, which requires the effective analysis of discussion chains in thread structures. The scope of the study is concerned with opinions text written in English and Amharic language posted on event posted about Ethiopian national football Team (ENFT) FACEBOOK Page. We select This Platform mainly for the popularity of platform and the availability of ocean of opinion posts from different users as compared to other social media's like tweeter. As Tweeter comment is limited to only 140 characters, it limits user opinion. In this study, we select sport domain because, sport is one of the hottest topics where peoples can express, share their opinion, regarding to players, team, coach, and seasonal status and overall performance of the team. The study used an opinion-oriented graph-based summarizing model from an opinionated discourse text posted on Public Facebook page. The polarities of words, phrases, and sentence are determined manually by

experts from discourses published online between Facebook users dated between March 27, 2022, and April 3, 2022.

### 1.5. Significance of the study

In the current competitive domain sport, knowing what other people think is a determinant factor in decision-making. The study has an immense value for the sports federation to make efficient and effective decision. For instance, the sport federation has a keen interest for identifying and analyzing the strength and weakness of the national team from the opinions generated by Ethiopian people on social media. In identifying this strength and weakness, the federation encourages the national team by providing different rewards in order to improve their performance and find out new techniques that improve their problems. For example, the federation can change the coach of the team if the summary of opinions regarding the coach is more of negative. It also focuses on selection of players if many users raise negative opinion regarding selection of players. The coacher wants to know user's opinion concerning the coacher performance, players' performances and all events occur during different match of the game. Not only coachers but also important for players of the team to know their performances from the opinions summary of the communities. In general, the summary of users' opinions is so important for the sport federation to know critical problems that have to give emphasis. The research will be used as stepping stone for others researcher who wants to conduct research on related area of study.

### 1.6. Literature Reviews

In this section we will review Opinion mining literatures from different sources such as published papers, journal articles, conference papers, books, and other Internet sources in detail to have better understanding of the problem and to have detail knowledge on the various techniques of opinion mining. We reviewed various related works from the concept of both machine learning and social network techniques for opinion mining of discourse reviews.

### 1.7. Research Methodology

Methodology is a way to systematically solve the research problem. This research is conducted in order to figure out challenges of implementing opinion mining system from social networks.

In order to achieve the main objectives we will follow step by step procedures discussed in chapter three.

## 1.8.    Thesis organization

This thesis report is organized in to five chapters.  Chapter one deals with introduction and over all back ground information of study concepts has been discussed.  Chapter two introduces review of related literature and overview of opinion mining (sentiment mining), different techniques and approach used in sentiment mining researches both from computational linguistics and social network. Moreover, the general steps in sentiment mining are also discussed in this chapter. Chapter three describes the general architecture or frame work of the proposed model, the developed model, algorithms and its implementation has been presented. Chapter four, states the experimental result and evaluations of the model. Chapter five, states the future works, the contribution of the work, recommendations and conclusions.

## CHAPTER TWO

## LITERATURE REVIEW

In reality, textual information is categorized into two types: facts and opinion. Fact is objective sentence that expresses some realistic information about the world, while a subjective sentence expresses some personal feelings or beliefs [14], Example of objective sentence: "This book costs 10$ onAmazon.com!", whereas an example of subjective sentence "This book is amazing"; here the subjective sentence contains adjective words like amazing.    Sentiment analysis, also known as opinion mining as well, takes the written text and translates it into different contexts, such as positive, negative or neutral to identify the orientation of the text.

## 2.1    Opinions or sentiments and their types

Opinions are a private state that is not open to objective observation or verification.   It is defined as a person's idea and thought towards something and it is an assessment, judgment or evaluation of something [15]. In social network opinion of one individual is posted independently or depends on others, which means it is posted based on the previous situations [16].  In opinion mining task it identified that orientation of opinion by the holder towards an object which can be a set of components or attributes.

"*What other people thinks*" has always been an important factor in decision making process. Sentiment Analysis or opinion mining is the process to automatically determine the sentiments expressed in a piece of plain text of the holder's as positive, negative and neutral regard toward the claim about the topic [17].  In general, the term Sentiment is very broad and it constitutes emotions, opinions, moods, personal experiences, feeling, agreement and disagreement etc.  In this study we mainly focused on two opinion sentences in terms of agreement and disagreement detection.

As Bing Liu [17] stated, opinions are classified into two types: regular opinion and comparative opinion.

### 2.1.1   Regular opinion

Regular opinions often referred to as an opinion stated in many literatures and it has two main sub-types [18].

 Direct opinion: A direct opinion refers to an opinion expressed directly about an entity or an entity aspect, e.g., "The picture quality is great."

Indirect opinion: An indirect opinion is an opinion that is expressed indirectly on an entity or aspect of an entity based on its effects on some other entities.  This sub-type often occurs in the medical domain.  For example, the sentence "After injection of the drug, my joints felt worse" describes an undesirable effect of the drug on "my joints", which indirectly gives a negative opinion or sentiment to the drug. In the case, the entity is the drug and the aspect is the effect on joints.

Much of the current research focuses on direct opinions [19].  They are simpler to handle. Indirect opinions are often harder to deal with.  For example, in the drug domain, one needs to know whether some desirable and undesirable state is before or after using the drug.  For example, the sentence "Since my joints were painful, my doctor put me on this drug" does not express a sentiment or opinion on the drug because "painful joints" (which is negative) happened before using the drug [20].

### 2.1.2 Comparative opinion

A comparative opinion expresses a relation of similarities or differences between two or more entities and/or a preference of the opinion holder based on some shared aspects of the entity [21]. For example, the sentences, "Coke tastes better than Pepsi" and "Coke tastes the best" express two comparative opinions. A comparative opinion is usually expressed using the comparative or superlative form of an adjective or adverb, although not always (e.g., prefer). Comparative opinions also have two types, Explicit and implicit opinions [22].

Explicit opinion: An explicit opinion is subjective statement that gives a regular or comparative opinion, e.g., "Coke tastes great," and "Coke tastes better than Pepsi."

Implicit (or implied) opinion: An implicit opinion is an objective statement that implies a regular or comparative opinion. Such an objective statement usually expresses a desirable or undesirable fact.

e.g.:

*"I bought the mattress a week ago, and a valley has formed," and*
*"The battery life of Nokia phones is longer than Samsung phones."*

Explicit opinions are easier to detect and to classify than implicit opinions. Much of the current research has focused on explicit opinions. Unlike regular opinions, it does not make much sense to perform sentiment classification to a comparative opinion sentence as a whole because such a sentence does not express a direct positive or negative opinion. Instead, it compares multiple entities by ranking the entities based on their shared aspects to give a comparative opinion. That is, it expresses a preference order of the entities using comparison. Since most comparative sentences compare two sets of entities, the analysis of an opinionated comparative sentence means to identify the preferred entity set. However, for application purposes, one may assign positive opinions to the aspects of the entities in the preferred set, and negative opinions to the aspects of the entities in the not preferred set. Note that like regular sentences, it is still meaningful to classify whether a comparative sentence expresses an opinion or not, but little research has been done on such classification. In our work both types of opinions are used.

## 2.2    Over view of opinion mining

As presented in figure 2.1, Opinion mining task involves three main steps, opinion retrieval, opinion classification and opinion Summarization [23],

Opinion Retrieval is the process of collecting review text from review websites.    Different review websites contain reviews for products, movies, hotels and news. Information retrieval techniques such as web crawler can be applied to collect the review text data from many sources and store them in database.  This step involves retrieval of reviews, micro blogs, and comments of user.

The primary step in sentiment analysis is opinion classification of review text. Given review documents, D = $\{d_1...d_i\}$ and a predefined category set C= {positive, negative}, sentiment classification attempts to classify each $d_i$ in D, with a label expressed in C.    The approach involves classifying review text into two forms; namely, positive and negative [24].

Opinion Summarization is another major part in opinion mining process.  Summary of reviews provided should be based on features or subtopics that are mentioned in reviews.    Therefore, feature extraction [25] and opinion summarization are key issues. Many researchers worked on summarization of product reviews [26].   The opinion summarization process mainly involves the following two approaches. *Feature based* summarization involves the finding of frequent terms (features) that are appearing in many reviews.   The summary is presented by selecting sentences that contain particular feature information [21].    Sentences in which feature and opinion words are present are displayed in summary of reviews. *Term Frequency based summarization* involvesa count of term occurrences in a document.   If a term has higher frequency it means that the term is more import for summary presentation. In many product reviews certain product features appear frequently and associated with user opinions about it. In this method sentences are scored by term frequency [27].   The summary is presented by selecting sentences that are relevant and contain highest frequency terms. Figure 2.1 the reviewed Architecture of opinion mining which shows how the input is being classified on various steps to summarize the reviews.

Figure 2.1: Reviewed Architecture of Opinion Mining (Loret, et. al, 2012)

## 2.3    Levels of opinion mining

As Liu, [20] defines the task of opinion mining or analysis has been mainly investigated at three levels of classification namely document-level, sentence level and phrase-level (future based)

### 2.3.1. Document level Opinion mining

Opinions analysis has been done at a document level for movies, book reviews and spam detection, etc., starting from the assumption that each document (or review) focuses on a single object (product, topics) and contains opinion from a single opinion holder [28].  This setting is true for reviews, but does not hold for Facebook or blog posts due to the short length of the posted text and each sentence or phrase may indicates different opinions about the topic.   In document Level opinion mining the whole document is written about only one topic and only by one person at a time.   In this study, it is interested in knowing so many peoples' opinion so it is useless for this study.

### 2.3.2. Sentence level Opinion Mining

In the same to the others Opinion Mining, the polarity of each sentence is calculated. The same document level classification methods can be applied to the sentence level classification problem also but Objective and subjective sentences [29] must be found out. The subjective sentences contain opinion words which help in determining the sentiment about the entity. After which the polarity classification is done into positive and negative classes.

### 2.3.3. Phrase level Opinion Mining

Document-level and sentence-level analyses do not discover what exactly people liked and did not like. Phrase-level opinion mining performs finer-grained analysis and directly looks at the opinion [17].The goal of this level of analysis is to discover sentiments on aspects of items. Phrase-Level Opinion Mining was earlier called Aspect level or feature level (feature-based opinion mining and summarization).

### 2.4     Opinion mining Task

Opinion mining task has three main components [22]. Opinion Holder is the person or organization that expresses the opinion. In case of user-generated content (reviews, blogs, social media updates) opinion holders are generally the authors of the post. On the other hand, Opinion Object **is** the feature about which the opinion holder is expressing his/her opinion i.e. opinion object is the center point of the user post. It can be a news, event, product, movie, location, hotel, sport etc. *Opinion Orientation* indicates whether the opinion provided by the person is positive or negative about an object, for example "This team has excellent players". In this review, Opinion Holder is the user who has written this review. Opinion object here is the team and the opinion word is "excellent" which is positively orientated.

### 2.5     Opinion mining Classification approach

There are many approaches for sentiment classification in opinion text. One of the most widely used methods involves classifying a single word or phrase with sentiment, and then calculating an overall sentiment score for a target document [30]. The most commonly applied techniques for opinion mining classification are described as follows.

### 2.6  Machine learning approach

The machine learning methods treat the sentiment classification problem as a topic-based text classification problem. The machine learning approach uses supervised learning method for classification of review text. It requires a corpus containing a wide number of manually tagged examples.    The first step is to train a classifier using sample of reviews with its class (positive/negative) then the built model of trained classifier is used to predict category of new text reviews. Popular machine learning classifiers for text categorization are Support Vector Machines (SVM) and Naive Bayes (NB).  SVM identify a hyper-plane that separates two classes of data. The chosen hyper-plane creates the largest margin between the two classes to make the points belonging to different classes and also make those points away from the hyper-plane as far as possible [31]. On the other hand a  NB classifier  is  a  probabilistic  classifier  based probability  model  that  incorporate  strong  independence assumptions among the features [32]. According to the work of *Pang et al.* [33] the experimental results produced are 81.0% for (NB), and 82.9% for (SVM). They used two classes: positive  and  negative,  and  worked  using product  reviews  that  are  longer  texts than Tweeter and Facebook .   However  the machine learning approach is not in line with short text collected from Facebook and Tweeter.   In spite of the fact that, people use a single sentence points out judgment about entities without considering sentiment score in the whole review about the domain topics.  In Social Web e.g. blogs, forums, reviews, micro blogs and Facebook post the vocabulary used between users is more out of annotated  corpus. The "distinguishing  words"  are  meaningless  for annotated  corpus (e.g. misspellings, ungrammaticality, shortening of words and repetition of letters and punctuation signs, the use of colloquial expressions, acronyms). Hence with machine learning approach a confident result is not produced for decision making.

### 2.6.1    Natural Language processing approach

Natural Language Processing (NLP), also known as computational linguistics, is a field of computer science that studies interactions of human languages with computers [34]. The main goal of NLP is to enable effective human-machine communication. The most important application of a NLP technique used in text mining is Part-of-Speech tagging (POS). It can be seen as the process of assigning a part of speech or other lexical class marker to each word in a sentence [35]. The parts of speech tags divide the words into different categories based on

different roles they play in one sentence. The traditional English language grammar classifies parts-of-speech in the following categories: verb, noun, adjective, adverb, pronoun, preposition, conjunction and article. The reason why POS tagging is so important to information extraction is the fact that each category plays a specific role within a sentence (e. g. Adjectives are good indicators of subjectivity a sentence (e. g. good, bad, worst), but other parts of speech also contribute to the judgment of subjectivity, e.g. verb ("like", "hope", "hate").

### 2.6.2   Lexicon-based techniques

This technique uses sentiment and subjective lexicon of terms.  The basic idea behind this system is to classify reviews based on how many positive and negative terms are present in the document [36].    This is based on a rule-based classifier, where if there are more positive than negative terms then it is considered to be positive. If there are more negative than positive terms then it is considered to be negative. If there is equal number of positive and negative term then it is neutral. The lexicon-based approach also has its own shortcomings: it is hard to use to find domain- or context-dependent orientations of sentiments. In other words, the sentiment orientations of words identified this way are domain- and context-independent [37]. However, there are many sentiments that have context dependent orientations, e.g., 'quiet' is negative for a speaker phone, but it is positive for a vacuum cleaner. Based on common Sense it is also difficult to determine the sentiment polarity, e.g., "this washer uses a lot of water" is with negative polarity [38].

### 2.7  Opinion lexicon generation

Opinion words are employed in many sentiment classification tasks.  As we explore in this research literature opinion words are also known as polar words, opinion-bearing words and sentiment words. To compile or collect the opinion words list, three main approaches have been investigated: manual approach, dictionary based approach and corpus-based approach [14].

### 2.7.1        Manual approach

The manual approach is labor intensive and time consuming, and is thus not usually used alone but combined with automated approaches as the final check; because automated methods can frequently make mistakes.

### 2.7.2    Dictionary Based Approach

This approach extracts the polarity of each sentence in a document. Afterwards, the sense of the opinion words in the phrase is analyzed in order to classify the sentiment in the text. Generally speaking, the techniques that follow this approach are based on lexicons, and use a dictionary of words mapped to their semantic value which labeled as positive and negative. As in [39] stated that the most common lexicon resource for English language is WordNet (which is a semantic lexicon where words are grouped into sets of synonyms (called synsets)) and SentiWordNet which is an extension of Word Net. This one is a sentiment lexicon that represents an index of sentiment words, and it has the polarity information of the relevant word irrespective of whether it carries a positive sentiment or a negative one.

Generally, this method works as follows: A small set of sentiment words (seeds) with known positive or negative orientations is first collected manually, which is very easy. The algorithm then increases this set by searching in the Word Net or another online dictionary for their synonyms and antonyms. The newly found words are added to the seed list. The next iteration begins. The iterative process ends when no more new words can be found. The dictionary based approach and opinion words collected have some shortcomings. The approach is unable to find opinion words with domain specific orientations, which is quiet common. For example, for a speakerphone, if it is "**quiet**", it is usually negative. However, for a car, if it is "quiet", it is positive.

### 2.7.3   Corpus based approach

Popular corpus-driven method is to determine the emotional affinity of words which is meant to learn their probabilistic affective scores from large corpus. The method in the Corpus based approach relies on syntactic or co-occurrence patterns and also a seed list of opinion words to find other opinion words in a large corpus. The technique starts with a list of seed opinion adjectives words, and uses a set of linguistic constraints of conventions on connectives to identify additional adjective opinion words and their orientations. One of the constraints is about conjunction (AND), which says that conjoined adjectives usually have the same orientation. For example, in the sentence "this car is beautiful and spacious," if "beautiful" is known to be positive, it can be inferred that "spacious" is also positive. This is so

because people usually express the same opinion on both sides of a conjunction. The following sentence is rather unnatural, "**this car is beautiful and difficult to drive", if it is changed to "this car is beautiful but difficult to drive",** it becomes acceptable.  Rules or constraints can also be designed for other connectives; OR, BUT, EITHER-OR, and NEITHER-NOR. This idea is called sentiment consistency [14].   The corpus based approach has also its own drawback, it requires a large corpus to get a large set of opinion words but it is a best approach for domain or context dependent meanings.

### 2.7.4   Ontology based approach

Ontology defines the common words and concepts (the meaning) used to describe and represent an area of knowledge. This definition has two parts: describing and representing an area of knowledge, defining the common words and concepts of the description [40]. Ontology appears specially promising for sentiment mining.  The use of ontology has the potential to refine and improve the process of sentiment mining by identifying specific properties of a domain as well as relationships between different concepts from that domain[41]. Ontology itself is an explicitly defined reference model of application domain with the purpose of improving information consistency and knowledge sharing. It describes the semantics of a domain both in human-understandable and computer processable way. In general, opinion mining is quite context sensitive, and at a coarser granularity, quite domain dependent. As a result a fine grain approach for opinion mining is needed [42]

### 2.8  The Applications of Opinion Mining

Information concerning people's opinions can be a very important component for more accurate decision making in a number of domains. Companies, for instance, have a keen interest in finding out what are their customers' opinions on a new product launched on a marketing campaign [43]. As well consumers on the other hand would benefit from accessing other people's opinions and reviews on a given product they are intending to purchase, as recommendations from other users tend to play a part on influencing purchasing decisions (For example, items that receive a lot of negative feedback will not be recommended.) [44]. Knowledge of other people's opinions is also important in the political realm, where for instance, one could find out the sentiment towards a new piece of legislation, or an individual such as a

politician or protester [45]. In the financial industry, sentiment information present on financial news has been studied to assess its impact on the performance of securities [46].

## 2.9 Reviews of Related Literature

We reviewed various related works from the concept of machine learning and social network techniques for opinion mining from discourse posts. For the task of machine learning approach, Opinion Mining or Sentiment Analysis from discourse text is defined as a process of extraction of sentiment described in the form of positive, negative or neutral about a particular topic or problem **[7].** As stated in **[11], [12], [13]** Sentiment analysis techniques can be roughly divided into the lexicon-based methods and machine-learning methods **[14]**. Lexicon-based methods rely on a sentiment lexicon, a collection of known and pre-compiled sentiment terms.

Machine learning approaches make use of syntactic or linguistic features [15], [16]to find out sentiment statics of opinions sentence. In these, we believe that determining the meanings of discourses opinions are challenging for the method as the sentiments or opinion meanings are determined by back-tracing of each opinion post. When we examine opinion message in the discourse shown in eg2,in section 1.2 some of the sentiments are ambiguous from the view of computational linguistic compared to the actual meaning throughout discourses. Because it does not mean that all positive opinions are given to imply only for a positive message and a negatives opinion are given to imply only for a negative message. In the works [17], [18], [19], [20], [21] determine the meanings of opinion in discourse posts using connectives, i.e. cue words and phrases [22], [23] discussed probabilistic models for identifying elementary discourse units at the clausal level and generating trees at the sentence level, using lexical and syntactic information from a discourse-annotated corpus. However, most of these discourse-based works narrow their scope to detect the sentiment polarity of a single unit. However, discourse analysis work concerned with the actual meaning of a message at the whole discourse remains unexplored dimension.

The intention of this work is to tide and compute more than one discourses unit (sentences, posts) relationship in order to determine the final sentiment meaning used for the decision.

### 2.9.1   Structure of social network

In this portion we aimed to investigate current social network techniques for opinion analysis, structure of opinion posting in social network and graph based representation of opinion are reviewed from related study.

According to [47] online forums, discussions are typically structured as **threads**, which is tree shaped structures in which multiple posts can share the same parent or ancestor post. In fact, it is often the case that a single post may bring out many comments, which can either respond to the initial poster or to one of the comments on the replay.   This comment of opinion threads has two types of sentences: the *antecedent sentence* and *reaction sentence.   The antecedent sentence is* the root sentence or the post which makes many reactions text. Reaction sentence is a reply (claim) made in an ancestor post.   Each opinion sentence can be shown of either agreement (positive) or disagreement (negative).   In [48] two sentences are in agreement mode when they show positive opinion between *antecedent and reaction* and shows disagreement when they shows negative opinion on the same argument.

For example, the message "*I do not agree with the  new coaching style chosen for the national team* " that receives as a reply the message «*Yeah, the players are not familiar with this style*» is a sequence of two negative opinions that shows agreement between them. On the contrary, the message "*It is bad to say this*" as a reply to the message "*it is foolish game*" does not point out agreement even though both messages express a negative opinion and they are connected. Similarly, a sequence of a positive message followed by a negative one does not necessarily show disagreement between the messages. Thus it is better to explore critical techniques that handle such kinds of opinion thread. In our case we reviewed related articles that contribute to the basic idea from current social network techniques for opinion analysis as follow.

### 2.9.2   . Current Social Network techniques

The Social Network Analysis deals with the analysis of the relationships that exists between entities in a social network.   For instance, in a social network of people, the analysis can include who is friend with whom, who can influence which group of people, who can have access to the information that goes through this network, etc.   Being [17] analyzes newsgroups by applying Social Network techniques and they interpret online communities by assigning roles to the members of the groups.   This is done by observing how people relate to each other in a graph-

based model of post-reply relations. They notice that short discussion threads point out question-answer exchanges and longer threads indicate proper discussions.

In [18], authors represent a newsgroup as a user-based graph and they base their analysis on the "reply-to" links between the users. They do not consider the content of each text because they claim that the statistical methods do not work for small messages where users use similar vocabulary. Our approach is more similar with this two papers, the little difference is we used data structure concept in which unstructured data is first structured and labeled in tree structure for analysis purpose, hence, we didn't consider the opinion text or string instead opinion target and their relation-ships.

Our approach also differs in that, we didn't use lexicon resource or dictionaries to find grammatical meaning of the sentences instead manually annotated corpus that opinions are retrieved to be summarized from tree structures.

In [47] sentiment classification based on the opinion frame annotation is performed. The classification algorithm used was collective classification which performs classification on a graph. The nodes are sentences (or other expressions) that need to be classified, and the links are relations. In the discourse context, they are sentiments related discourse relations. These relations can be used to generate a set of relational features for learning.

Each node itself also generates a set of local features. The relational features allow the classification of one node to affect the classification of other nodes in the collective classification scheme.

Zhou et al [49] used the discourse information within a single compound sentence to perform sentiment classification of the sentence. For example, the sentence "**Although Fujimori was criticized by the international community, he was loved by the domestic population because people hated the corrupted ruling class**" is a positive sentence although it has more negative opinion words.

Other paper is the work of *Abbott et al* [50]that explore several forms of agreement and disagreement and ask the annotator to take into account the context of the phrases by providing the entire document. In other related work, *Bender et al* [51]annotate Wikipedia discussion forums for positive and negative alignment moves which express agreement and disagreement respectively between the source and target. Their annotation includes praise, doubt, and

sarcasm in addition to explicit agreement and disagreement. They did not have an annotation tool, but simply they manually annotate the documents directly.

*Helander et al* [52], that analyzes the Innovation Jam 2006 among IBM employees and external contributors. The representation of the discussion is seen from the point of view of posts rather than users. The difference from our work is that our objective is not to find out the degree of innovation of a discussion but to identify opinions. Moreover, in our case, the participants of the opinion Discussion come from different backgrounds and they have different concepts and beliefs. Also, while in the IBM Innovation Jam the users are known since they are specific IBM employees, in our work users remains anonymous. The anonymity allows people to express more honestly how they feel about a certain issue.

In *Maurel et al.* [53], forums in the domain of tourism have been analyzed and they have extracted information regarding user sentiments and tourist destinations. They apply syntactic and semantic processing techniques and they adapt the grammar rules or the opinion words they try to identify according to the domain. They do not, though, represent the discussion as a graph.

As in [16], links in social network are implicit opinion in the postings which is the form of "in response to" tags they retained from the raw data only, those postings that contained both the author and the person whom the author was responding to. Each such posting yields a link. They were used 65% to 80% of the three datasets from the archives of the Usenet postings and they consider the remaining posting as either not opinion responses, they experiment with two commonly used text classification methods: Support Vector Machines and Naive Bayes classifiers. Both SVM and Naive Bayes were unable to distinguish the two classes, for any of the datasets. The reason they investigate for the low accuracy is that in a newsgroup discussion the vocabulary used by two sides is quite similar. Meaningful words are contained equally frequently in both the positive and negative classes, and the "distinguishing words" are meaningless. Finally, they conclude that extracting useful information from a newsgroup, forum discussion and similar contents using conventional text mining techniques has been hard because the vocabulary used in the two sides of an issue is generally identical and because individual postings tend to be less.

Lastly, local study has been conducted for sentiment classification of Amharic [8], in this the lexica of Amharic sentiment terms are developed to identify and assign initial polarity value to the sentiment terms detected. They use the NLP techniques; normalization and tokenization to

detect movie reviews and Newspapers. However, these techniques analysis text written in straight Amharic text without considering its relations with other sentiment written in different sentences

In another hand Most of existing Social Network opinion mining work deals with the analysis of the relationships between entities in a social network like who is friends with whom, who are experts and who post reaction text, what is the central node, etc. In Fisher et al. [24] analyze newsgroups by applying Social Network techniques and they interpret online communities by assigning roles to the members of the groups. This is done by observing how people relate to each other in a graph-based model of post-reply relations.

In [25] the authors represent a newsgroup as a user-based graph and they base their analysis on the "reply-to" links between the users. The most related study with our work is the work of Stavrianou et al [26]. They propose a framework for discussion analysis by combining Social Network and Opinion Mining techniques and they study the structure of an online debate and analyze the user reactions, preferences and opinions on a certain subject, by combining user-based graph and opinion-based graph. The main objective of the author was to enhance the user-based graph with additional opinion information.

The authors also proposed a measure to analyze discussion thread mainly considering relationship resides between entities rather than opinion chain. They did not consider the relationship between message nodes resides in the threads. Moreover, they also did not consider the intra-relationship present between thread chains which are important to express the aggregate opinion polarity. Our work is the extension of the work of author's [26] by providing more improved opinion thread measure techniques and a well-defined set of rules that can enhance a thread discussion analysis. In general, the following aspects were considered in this study:

- We incorporate more detail data structure concepts for the data storage and representation of opinion polarities which is important in the opinion-oriented graph model.
- An attempt was made to develop a graph-based opinion summary model from opinionated discourse posts.
- We proposed a measure and rules that consider the additional variables, the relationship between message contents (opinion polarity presented between each node).

- We considered also the intra-relationship resides between thread chains that enable us to overlook the total summary of threaded opinion polarity

## 2.9. Summary

This chapter reviewed different research attempts to solve the problem of opinion mining from computational linguistic, machine learning and social network Techniques. In social network; posts are appears in terms of agreement and disagreement, discussion thread, link analysis and discourse information. We observed that extracting useful information from Facebook, forum discussion and similar contents using straight forward text mining techniques is not significant to decision making. The reviews show that *opinion based graph* representations are a good technique to obtain better result for assertive decision making in the proposed social networks domain of study.

# CHAPTER THREE

# RESEARCH METHODOLOGY

### 3. Introduction

In this section we will introduce the way to systematically solve our research problem and the way to achieve the main objectives of the study. For this we will follow step by step procedures enable as to achieve the result. First we start from identifying the data source for the study. The data set is secondary data; it is opinionated discourse text that will be collect from the national football team Facebook page. The data-sets used for the experiments will be crawling threads of discussions, comments and plenty of replies that are, posted on specific event. We will use uses the Facebook Graph API to crawl public posts and its' associated public comments and replies. However, before opinion summarization by the use of graph theory (tree structure), for crawled labels of opinion text, a very basic preprocessing phase has been applied to the corpus before linguistic annotation. The ambiguous text expressed in idioms, slang, Misspellings, Laughter, neutral opinions has to be ignored In order to make the annotation simple for linguistic men. In order to represent these thread chains, data structure tree node is used. An opinion-oriented graph-based summarizing model will be follow. The Graph API allows us to navigate through the graph of the social network, which is organized into tree nodes. The crawled datasets are labels of text available in XML file format. The labeled XML file is converted into data structure nodes to indicate parent-child relationships. We will use python programming language as Implementation tools. The developed opinion oriented model will be evaluated with complimentary information provided by computational linguistic or text mining to that of graph-based summary.

### 3.1 Research Design

This section describes the design of the proposed opinion mining model for exploring opinion information in detail and presents the general architecture of the study. It is the conceptual framework within which the research is going to be conducted. The research design is more of experimental and explanatory research and used fully secondary data sources.

### 3.2. General system architecture

The general architecture of the proposed model (sentiment or opinion mining model for opinionated text reviews in English from social Network Facebook site as shown in figure 3.1. As it depicted from the figure, the architecture contains different components based on the processes required. These components are: review of opinion subjects from the site, Structuring opinion information, annotation of pinion information, constructing opinions based graph, measuring opinion polarity, and summary of opinion polarity.



**Figure 3.1: Proposed architecture of Graph-based opinion summarizing model**

 (Source: Own)

### 3.2. Review of Data Source

This step extracts opinions Threads that contain a plenty of comments and replies posted by different opinion holders, on the page regarding Ethiopian national football team sport events. The datasets used for the experiments were crawled from national football team Facebook public page. We crawls a total of 49 threads of discussions, containing 1120 comments and replies. In order to make the annotation simple for linguistic men, the irrelevant comments were removed.

The web has suddenly changed the way how people express their views and opinions. They can now post reviews on various websites, participate in discussion on various forums, write a blog describing their experience, update their status on social websites like face book, twitter, blogs etc. This data on review websites, discussion forums, blogs, and social networks can be collectively called as user-generated content. Each of this user-generated content has their unique property.

In this study, we select only comments posted on Facebook. Because the opinion posted on Facebook is a kind of conversation or deluge. In this Deluge any Users who are eligible to take part in this media have tendency to generate many opinions without limitation of space. Peoples are more reflective of idea when they are looking others opinions. In Facebook the user can view others post and make response. The remaining user-generated content like twitter and blog post are less reflective of opinion because more of the opinions of the author are posted regardless of the many previous posted opinions.

Twitter allows only 140 character status updates. And also in blog only restricted authorized users can post blogs. Due to the above reasons we select Facebook as the main source of our data set to conduct this study. When we observe the above social media's there are a lot of posts concerning the problem domain. We observe that different controversy issue, debate as well as replies between users opinion show agreement and disagreement. This creates opinions chain. In such opinion chain the researcher faces difficulty in identifying the overall outcome or final conclusion Even though we can identify individual meaning of opinion sentence.

So that these interesting sources of data motivated us to build a model that enable us to know the overall outcome of the event by designing and measuring opinions information. The dataset collected for the experiment contains mixed opinions written by more of Amharic, English and Afan Oromo. As we know Ethiopia is known by multicultural and multilingual country, As a result all opinion holders can give their opinion in two or more than two languages. The opinions text written in the Amharic language and English is selected since many users post in these two languages. These opinions posted are written in different language and posted on two different pages Ethiopian national football team (walya) and Ethiopian football Federation, the messages they provides are about the same event.

## 3.1 Structuring opinion information

Different Facebook users can posts thousands of messages, Design scheme is required to make sense of all the information exchanged between the Facebook users. The basic organization unit of the Facebook post is the Thread of conversation. A thread is a set of messages which address the question or discussion topic announced in the first message which is called thread head (Zhang, Ackerman, &Adamic, 2007). The thread head is the initial message or issue, followed by the holder of the message. The remaining messages are all follow-ups to the thread heads, i.e. replies or comments from other participants.

Different Facebook users can post thousands of messages; a design scheme is required to make sense of all the information exchanged on the Facebook page. The basic organizational unit of the Facebook post is the Thread of discussion. A thread is a collection of messages which address discussion topic declared in the first message called thread head [15]. In order to represent these thread chains, data structure tree node is used. This method is appropriate to show whether each message is a follow-up to the original post or to one of the replies arises from it. We used the Facebook Graph API to crawl public posts and its' associated public comments and replies. The Graph API allows us to navigate through the graph of the social network, which is organized into tree nodes. The crawled datasets are labels of text available in XML file format. The labeled XML file is converted into data structure nodes to indicate parent-child relationships.

This node is directly mapped to a number list to refer the parent of the replies (child node). The parent-child node has created a forest tree structure, for this we applied the data structure storage technique linked list, to summarize threaded opinion polarity. However, before opinion summarization by the use of graph theory (tree structure), for crawled labels of opinion text, a very basic preprocessing phase has been applied to the corpus before linguistic annotation. The ambiguous text expressed in other language, in idioms, slang, Misspellings, Laughter, neutral opinions are removed and then only the structured text of thread opinion has been given to linguistic expertise to annotate the thread texts as positive, negative and neutral. However, after annotation, we exempt neutral opinions, because, it has no any influence on decision going to be made so that it is not considered in the model as input. But the neutral opinion is seen as an equal number of positive and negative opinions extracted in opinion chain while

summarization. To illustrate this, Table 3.1 represents the structure of the annotated opinion thread.

| Opinions | Order lists (node) | Opinion sentences | Polarity |
|---|---|---|---|
| Thread 1 Post 1 ) | 1 | *Ethiopian National team lost score* | **Neutral** |
| Comment | 1.1 | **I think, it is a coach problem.** | **Negative** |
| Replay | 1.1.1 | **You are right, the coach acts over confidence.** | **Negative** |
| Replay | 1.1.1.1 | *ትክክል ነህ አሠልጣኙ በበፊቱ ብቃቱ ላይ አይ ገኝም::* | **Positive** |
| Replay | 1.1.2. | **I don't think, he knows well what the modern football is.** | **Positive** |
| Comment | 1.2 | **To me the problem is with goal keeper.** | **Negative** |
| Replay | 1.2.1 | *እንዴ ነከ ሀ ነ ግንግ ብ ጠባቂ ው ስህተት አለበት ብየ አሳም ንም::* | **Positive** |
| Replay | 1.2.1.1 | *አ ይደ ለም አን ቴ ጨ ዋታ በ ትክክል አለ የም::* | **Negative** |
| Replay | 1.2.1.1.1 | **በ ር ግ ጥ በ ረ ኛ ው ት ል ቅ ስ ህ ተ ት ሰ ር ቶ ል::** | **Negative** |
| Comment | 1.3 | **The problem is at defense position.** | **Negative** |
| Replay | 1.3.1 | **You are right; the great mistake is there.** | **Positive** |
| Replay | 1.3.1.1 | *የሚ ገ ርም ቡ ዲ ን አለ ኝ!! ግ ን ተ ከ ላ ካ ይ መ ስ መ ር ች ግ ር ይ ተ ያ ል* | **Neutral** |
| Replay | 1.3.2 | **It is the gap that has to be filled.** | **Negative** |
| Replay | 1.3.3 | **Yes, they are not fit for the position.** | **Positive** |
| Replay | 1.3.4 | **For sure we have incredible team if the defiance line is has solved.** | **Positive** |
| Thread 2 | 2 | What do you feel about national team game? | **Neutral** |
| Comment | 2.1 | It is letting me down; they can't even defeat this ordinary team. | **Negative** |
| Comment | 2.2 | The referee clearly changed Ethiopia's game plan, too many wrong decisions observed. | **Negative** |
| Comment | 2.3 | **GO…GO.. ETHIOPIA!!!!! It's our time to shine in the African football!!!** | **Positive** |
| Comment | 2.4 | **I am proud of Ethiopia National team I love you my mama Ethiopia!!** | **Positive** |
| Comment | 2.5 | **Professionals sport team, we support them for ever** | **Positive** |
| Comment | 2.6 | **Oh God! What an amazing emotional moment at the National Anthem!** | **Positive** |

**Table.3.1 Structure of opinion thread in Facebook**

The table above is the layout of the Facebook user's post and exchange opinion information under opinion thread. In this opinion thread each number list represents one opinion sentence.

The numbers list represents opinion sentence or none (initial post, comment, and reply)

E.g.1**. Ethiopian National team lost score** and

**2.What do you feel about national team game?**

Are two different opinion issues or topics, indicate thread head *(root node).* It is the initial message from which the remaining comments are arisen.   The rest of the number lists are indicate towhicheach participant posted or an idea of sentiment flow belongs to. This method is a suitable method to design the idea of opinions for proper data structure sentiment analysis.   In order to design the opinions in thread the researcher consider the following information.

To structure the corpus of opinions thread, the following three-fold hierarchical structure of Facebook discussion were considered.

**Post:** It is the main message objects or issue about which user's posts or it is the topic/issue that requires different comments or feedback from different users to make decisions. It is clear that the suspense content of the post has been often neutral, but this post can inspire the user to create long discourse.

**Comment:** Commentis a reaction text written on the front page of a public Facebook page for the initial message post.

**Users:** A community of people that responds to each other (friendship, interest in the same topic or not). They are the participants of the discussions identified by a user name and they can participate by either writing a new message or replying to an already posted one. However, information of users are not our objective instead their opinions.

**Reply**: The relations between the exchanged messages point out which message replies to what message. It enables us to keep the flow of opinion chain in discussion thread (e g. in the above Table 3.1 the direct replies-to number list, 1→{1.1, 1.2}, 1.1→{1.1.1}, {1.1.2.}, 1.2→{1.2.1}, {1.2.1.1}, {1.2.1.1.1} and etc.

**Opinion:** the sentiment behavior towards a previous post or parent post either in agreement (positive) or disagreement (negative). E.g. the message "I think, it is a coach problem" conveys negative opinion to the root message or opinion thread head. The message "*you are right, the coach acts over confidence*" suggest positive opinion or agreement to this opinion but implies negative opinion for the root message (thread head). Another issue related to this concept such as share,

like, unlike, friendship and tags are points that will be considered in our new dimension of future research work.

### 3.1.1 Opinion-oriented model

In this section we present the opinion-oriented model we propose for the representation of the structure of Facebook discussions. We consider that the participants of the discussions identify themselves by a user name and they can participate by either writing a new message or replying to an already posted one. The relations "reply-to" between the exchanged messages point out which message replies to what and they are considered to be known. The model is based on a graph-based representation. Most graph-based existing works consider users to be the vertices of the graph. In this study, we propose to use polarity of message objects as the vertices. Because in Data structure vertices are represents the node, this node used to store different values of variables.

### 3.1.2 . Opinion oriented graph

In this section we present a framework which achieves a Facebook discussion representations corresponding to the domain of study. The new representation allows us to exploit the structural characteristics of a Facebook discussion and analyze it from a semantics-oriented point of view. Most of graph-based existing works, consider users to be the vertices of the graph [8]. In this study, we suggest using polarity of message objects as the vertices. Because in Data structure, vertices are representing the nodes, these nodes are used to store opinion polarity. We represent this framework as "opinion-oriented graph", whose definition is as follows.

**Definition1.** Opinion-oriented graph (**OOG**) is a graph G = $(Vn, R_n)$ with a set of $Vn$ vertices and a set of $R_n$ reply. Each vertex $v_i$ represents a "message object" or thread head and its weight values are (0). Each reply $R_{nij} = (r_{ni}; r_{nj})$ points out direction from $r_{ni}$ to $r_{nj}$, and it is weighted by a value that represents the opinion polarities expressed in the message object $r_{ni}$ as a reply to what has been posted in the message object $r_{nj}$.

The weight is a function w: R → Z and it takes negative values when the opinion polarity is negative (-1), and a positive (1) value when a positive opinion is expressed. An opinion-oriented

graph (**OOG**) consists of opinion orientation which allows us to define opinion measures that extract useful information from it. In this study, we present two basic components; the discussion threads and the discussion chains. The distinction between a discussion thread and a discussion chain becomes apparent from Figure 3.2.



**Figure3.2. Opinion discussions threads**

In this figure, the first thread consists of 7 discussion chains:

1. $\{\,post[1],\ rep[1.1],\ rep\ [1.1.1]\}$

2. $\{\,post[1],\ rep\ [1.1],\ rep\ [1.1.2],\ rep\ [1.1.2.1]\}$

3. $\{\,post\ [1],\ rep\ [1.2],\ rep\ [1.2.1],\ rep\ [1.2.1.1],\ rep\ [1.2.1.1.1]\}$

4. $\{\,post\ [1],\ rep\ [1.3],\ rep\ [1.3.1],\ rep\ [1.3.1.1],\ rep\ [1.3.1.1.1]\}$

5. $\{\,post\ [1],\ rep\ [1.3],\ rep\ [1.3.2],\ rep\ [1.3.2.1],\ rep\ [1.3.2.2.1]\}$

6. $\{\,post\ [1],\ rep\ [1.3],\ rep\ [1.3.2],\ rep\ [1.3.2..1],\ rep\ [1.3.2.2.2]\}$

7. $\{\,post\ [1],\ rep\ [1.3],\ rep\ [1.3.2..1]\}$

The set of the discussion threads in an opinion-oriented graph G is the union of all the maximal connected components of G. The discussion chains consist of the paths in the graph whose starting node is a root and ending node is a leaf node when we invert the direction of edges. The chains are important in an opinion-oriented graph. The longest discussion chain can point out the

longest exchange of messages in a Facebook post discussions and it can be measured by the maximum number of edges that start from a root node and end up to a leaf node.

### 3.1.3   Opinion Relation measures

After having defined the opinion-oriented model, we now present some opinion relation measures that enable us to determine the relation of opinions in discussion and the opinion polarity of the participants. We have separated the opinion relation measures into three categories according to whether they characterize per node, per discussion chain and per discussion thread or as a whole. In order to define the opinion relation measures we first clearly defined what is an opinion Relation in our model?

The opinion $(V_X, V_y)$ is denoted by the weight of a reply $r_{x,} = (V_X, V_y)$, $r_x \in E$ and it expresses the opinion polarity which is present in node $v_x$ which is a reply to the node $v_{y,}$. It takes values in {-1; 0; 1} if the opinion expressed in the message object $v_{x,}$ is negative, neutral (i.e. if $V_X$ is root node) and positive respectively. Following this representation, opinion Measure techniques, are done as we can see in the next section.

To measure the follow of an idea a concept of antecedents were used. This consists of the set of reply nodes towards a vertex $V_X$ and according to the theory of graphs, it is defined as:

Reply $(V_X) = \{V_X, \in V \mid (V_y; V_X) \in R\}$………………………………………….. (1)

#### 3.1.3.1 Opinion Relation measures at vertex (node)

 A Message Object  $v_{x,} \in V$ may be replied to during the discussion through posts. These posts may contain the opinions of the responder expressed by positive or negative opinions. Opinion polarity measure can be determined at; three level i.e. at thread level, at thread chains, and at the node (vertex) level. Opinion measure per vertex (node) is determined by identifying the level of opinion chain. This level is defined from $(L_1, to\ L_{n,})$ which implies the flow of opinion polarity from root to leaf node. The opinion information contains two types of opinions direct and indirect opinions. The opinion found on first level of vertex is a type of direct opinions.

Direct opinion information is the direct reply for the message object. Level one (L1) Opinions are a direct post for root node (thread head).  But the replies posted from the first level of opinion chains to n (leaf node) are not direct. Instead they reinforce or antagonize indirectly about the

previous post. The opinion measure in the first level is straight forward, it could be right for computational linguistic to analysis; However the number of this opinions are few in number, so it is insignificant to undertake language study.

This level opinions measures identify individual opinion polarity; it is objective is to determine only sentiment polarities of the nodes as only positive and negative. The relation is going to be determined at chin level. But we do not have neutral in this level because the opinion sentence we employed is all of subjective and neutral were employed at the opinion chain.

### 3.1.3.2 Opinion Relations measures between chains

The second opinion measure task is done the through sway of opinion chain. The objective of this measure is to determine the sentiment relations throughout discussion chains. This measure is identified starting from level one (L1 − Ln). A discussion chain $G_{C,} = (V_C; R_C), R_C \in E_C$ in the graph G is a path whose starting node is a root and ending node is a leaf when we inverse the direction of the edges. In this we defined the opinion received by a message object $v_{x,}$ from the root node to the leaf node as:

$$Reply\ (V_{X,}r) = r_1 * r_2 \dots r_i,\ where\ r_i \in R_{C,},\ and\ R_C \in E_C \text{'}\dots\dots\dots\dots\dots (2)$$

In general, to determine the final summary of opinion chain computed from these relations, we identified the Rules come after the following figure 3.3

**Figure 3.3 opinion oriented Graph**

**Rule 1**:

If all reply in the opinion chain is positive, the product of reply in a sequence of opinions chain will give positive opinion summary (+). It points out an agreement between the reply's nodes regarding to the initial post or discussion topic.

**Example:** from the above opinion thread 1, the opinion chain { *rep* [1.2.], *rep* [1.2.1], *rep* [1.2.1.1], and *rep* [1.2.1.1.1]} indicates positive reply which means there is no disagreement between the replies posts throughout the chain. In this our objective is to obtain the final result obtained throughout the chain. So the products in this sequence give positive opinion result. So the decision made in this is relevant.

**Rule 2:**

If the number of reply in a sequence of opinion chain is odd and their annotated polarity is negative the product of reply will be negative.

For all Replies, $R_{i=-ve}$ the $\prod_{i=1}^{n} Ri = \{-ve\ for\ n\ is\ odd\ node\ +ve\ for\ n\ is\ even\ node\}$

Example: from the above opinion thread 1, the opinion chain { $rep$ [1.3], $rep$ [1.3.2], and $rep$ [1.3.2.1]} are three replies annotated as negative. But it does not mean that all replies imply negative opinion about the initial post. The $rep$ [1.3] is reply for negative, $rep$[1.3.2] is reply for positive it indirectly supports the initial post, $rep$[1.3.2.1] is reply for negative it indirectly antagonizes the initial post. In general the reply message that indicates the negative idea about the topic is greater than that positive one. The product between these opinions also gives negative result. Therefore, the assumption made by the rule is significant for decision making.

**Rule 3:**

If the first reply in the sequence of opinion chain is negative and the remaining reply in the chain is annotated as positive, the product of all replies throughout the chain will be negative.

**Example:** in the opinion chain { $rep$ [1.3], $rep$ [1.3.1], $rep$ [1.3.1.1], and $rep$ [1.3.1.1.1]}. The first reply $rep$ [1.3] is posted as having negative opinion polarity to the initial post. The remaining replies are annotated as positive which means they are the supporter of the first reply. . As a result they convey indirectly negative opinion to the root message. Thus the product throughout opinion chain in this assumption will give negative result so this assumption is convenient in deciding the final conclusion as negative.


**Rule 4:**

If the number of reply in opinion chain is even and there annotation polarities are negative the product of all reply throughout the chain will be positive'.

**Example:** The opinion chain { $rep$ [1.3], $rep$ [1.3.2], $rep$ [1.3.2.2], $rep$ [1.3.2.2.1]} are four replies annotated as negative. In this even though they annotated as negative the opinion message they convey is different. Here the $rep$ [1.3] and $rep$ [1.3.2.2] show negative idea about the message object but the $rep$ [1.3.2] and [1.3.2.2.1] show positive idea about the message object. In this rule the product of all replies throughout the chain will give positive conclusion but the numbers of replies that convey positive and negative opinion are the same so it is inconvenient to take it as positive conclusion.    Due to this rule and other combinations of opinion chains we improve the opinion relation measure techniques from the product to the average sum of the

product. In this the new opinion relation measure summary made at this rule is taken as neutral reply.

For this we further defined opinion relation measure techniques as the average opinion received by a message object $v_x$, from the root node to the leaf node as follows:

$$avgMsgOp\ (v_x) = \frac{\sum_{i=1}^{n}\ (v_{x,}\ (\pi R_i)}{n} \ \dots\dots\dots\dots\dots\dots\dots\dots\dots\dots\dots\dots\dots\ (3.1)$$

Where $(\pi R_i)$ indicates the Reply of the node $v_{x,}$ and it shows the sum of product of replies across a series of opinion chain from first *reply* to last reply (root node to leaf node), n is the total number of replies travels from root node to leaf node. It is equal to the number of opinion across discussion chain. It depends on the maximum size of opinion travels on one direction in the discussion chain. In this study 100 size of n is taken by a researcher as default to take account the problem may extends to beyond existing size. In fact limiting the size of opinion chain is; limiting user opinion, this can make biased in research design, so that it is better to determine based on the problem going to be addressed.

The average opinion towards a message object is an indication of the polarity of the discussion chain towards the specific post.

$$avgMsgOp\ (vx) = \{-1, if\ avgMsgOp < 0\ \ 0, if\ avgMsgOp = 0 + 1, if\ avgMsgOp > 0\}(3.2)$$

If the average opinion is 0, there is a balance between positive and negative opinions. Similarly if the *AvgMsg Op* is less than zero the Summary of opinion polarity is negative. Finally if the *AvgMsg Op* is greater than zero the Summary of opinion polarity is positive.

For instance, in the opinion-oriented graph of Figure 3.3, we can see that the thread 1, thread 2 and thread 3 have 5, 6 and 6 opinion chain respectively the average message opinion (*AvgMsg Op*) and its polarities are as shown in Table 3.2 below

| Msg.no | Thread no | Chain no | AvgMsg Op | Polarities |
|--------|-----------|----------|-----------|------------|
| 1 | 1 | 1 | 0 | 0 |
| 2 | 1 | 2 | 1 | + |
| 3 | 1 | 3 | -1 | - |
| 4 | 1 | 4 | -0.33 | - |
| 5 | 1 | 5 | 0 | 0 |
| 6 | 2 | 1 | -1 | - |
| 7 | 2 | 2 | 1 | + |
| 8 | 2 | 3 | 0 | 0 |

| 9 | 2 | 4 | -1 | - |
|---|---|---|---|---|
| 10 | 2 | 5 | 0 | 0 |
| 11 | 2 | 6 | -1 | - |
| 12 | 3 | 1 | -1 | - |
| 13 | 3 | 2 | -1 | - |
| 14 | 3 | 3 | -1 | - |
| 15 | 3 | 4 | -1 | - |
| 16 | 3 | 5 | -1 | - |
| 17 | 3 | 6 | -1 | - |

**Table 3.2 sample of $avgMsgOp$ measures applied to the short discussion.**

### 3.1.3.3. Opinion measure at the thread level

The last opinion measure in our model is the measure for the whole chain which involves measure at thread level. This measure is used to identify interesting opinion information at global views. This opinion information is measured by entropy H( $v_x$) and we define the amount of opinion information held by a node $v_x \in$ V (that has been replied to), as:

$$H(v_{x,}) = - \sum_{i=-1,0,1}^{n} (\frac{AvgMsgOp(v_{x,i})}{|replay(v_{x,})|} log \frac{AvgMsgOp(v_{x,i})}{|replay(v_x)|}) \dots \dots \dots \dots . (3.2)$$

The opinion information is an indication of the variety of opinions received by a node. If, for instance, a node has received reply posts that are all of the same opinion orientation, then the entropy will be 0. This information can be interpreted as: there is common opinion regarding the message expressed by the particular node.

**Example:** As indicated in table 3.2 above there are three opinion threads thread1, thread 2 and thread 3,their entropy $H(v_x)$ is depicted in table 3.3.

| Thread No | $AvgMsgOpinion(v_{x,1})$ | $AvgMsgOpinion(v_{x,-1})$ | $AvgMsgOpinion(v_{x,0})$ | $H(v_{x,})$ |
|---|---|---|---|---|
| Thread 1 | 1 | 2 | 2 | 0.4582 |
| Thread 2 | 1 | 3 | 2 | 0.3976 |
| Thread 3 | 0 | 6 | 0 | 0 |

**Table 3.3: sample of entropy measure for thread discussion**

From this table, we can understand that the Thead1 has the highest entropy of all. This is the message that has received replies with the highest variety of opinions polarity. Thread3 has zero entropy which shows lack of opinion variety in the replies it has received. This measure is used to rank opinion information from the most to the least interesting. A message that has received few but varied positive and negative opinions can be more interesting than one that has received plenty of messages that are all positive or negative. In this model, we measure this by the entropy. This measure allows the selection of the right issue to focus on. After the theoretical framework for opinion measures completed we are going to discuss the proposed algorithms that enable us to extract, classify and summarizes its polarity.

### 3.1.4 Proposed approach

In the proposed approach, we develop an algorithm for retrieving opinion thread structure from annotated corpus. Thread structure definitely plays an important role in identifying the parent child relationships between long discussion chains in this algorithm. This thread structure is the process where by a parent message is explicitly linked to one or more responding child messages. The approach consists of two steps. We used data structure node that is represented by a sequence of number list. In this case the first number list is the root node. The possible numbers lists that represent parent (root node) are (1, 2, 3…...n) and (1.1, 2.1, 3.1 etc.)

Each element of the number list (i.e. child node) contains a polarity of opinion between message chains, which measures value as positive (+1), negative (-1) and neutral (0). The second step is to determine parent-child relationships. These relations are obtained by opinion measure techniques in the section 3.3.3.3 which is called $AvgMsgOpinion$ if the measure of $AvgMsgOpinion\ is\ greater\ than\ 0$ the relation is positive and if the $AvgMsgOpinion\ is\ less\ than\ 0$ The relation is negative and if $AvgMsgOpinion\ is\ equal\ to\ 0$ the relation is neutral which means equal number of positive and negative. This base line merely used as threshold to identify relation along opinion chains. The Graph API allows us to navigate through the graph of the social network, which is organized into tree nodes. The crawled datasets are labels of text available in XML file format. The labeled XML file is converted into data structure nodes to indicate parent-child relationships. This node is directly mapped to a number list to refer to the parent of the reply (child node). The parent-child node has created a forest tree

structure, for this we applied the data structure storage technique linked list, to summarize threaded opinion polarity. Then the summarized opinions result is retrieved from corpus saved in XML file format. The following section illustrates the brief description of proposed framework.



**Figure 3.5: Flow chart of the proposed system**

The detail description of the data structures of framework is described as follow.

**Step1**. Create a user defined data type

**Step2**. Create Class **forest_tree** -->defines structure, attributes, and different functions

**Step3**. Class Forest_tree**Structtree_node**--> defines attributes /*such as attribute of sentiments

**[Explanation]** - file; /* which is used to store consecutive number list, sentiment data and its polarity

**Step4**. Class forest_tree**Structtree_nodedata**; /*allows sentiment polarity flow in tree structure

**Step5**. Class multi_tree**Structtree_noderesult**;/ *To determines whether the node is support /*or not to the previous node

**Step6**. Class forest_treeStructtree_node**level**;/*Determines level of the node

**Step7**. Class forest_treeStructtree_node**pnodenum[]**; /* Contains chain of the node or the //flow of the sentiments

**Step8**. Class forest_treeStructtree_node ***node [];/***Contains arrays of pointer variable that is Self referential

**[Explanation]** - file; /* which is used to store sentiment data and polarity is a set of /*consecutive number list.

**Step9**. Class forest_treeStructtree_node**nxtindex**; /*it indicates index of next empty child /*node of the current node

**Step10**. Class forest_treeStructtree_node**avgMsgOpinion**; /*it tells the polarity of each chain

**Step11**. Class forest_tree ***root**;/*Starting node of the thread

**Step12**. Class forest_tree ***curr**;/* current node of the thread

**Step13**. Class forest_tree**pos**; /* total positive polarity chain of the thread

**Step14**. Class forest_tree**neg**; /* total negative polarity chain of the thread

**Step15.** Class forest_tree**neut**; /* total neutral polarity chain of the thread

/* Descriptions of Eight different Function Prototype are used in class forest_tree to solve different functionality

**Step16.isEmpty()**

/*isEmpty function is used to check if the thread is empty

**Step17**. **empty()**

/*empty function is used to reset the thread to empty

**Step18.Insert()**

/* Insert function is used to enter the sentiment value into the tree node

**Step19. call_countLeave()**

/*call_countLeave is used to call countLeave function

**Step20**. **countLeave()**

/*countLeave function is used to count total positive, negative or neutral polarity of the thread chain

**Step21**. **setCurrent();**

/*set current function is used to set the current node by calling findCurrent function /*iteratively and sending the root node.

**Step22**. **findCurrent()**

/*find current function is used to find the current node by accepting root node from /*setCurrentfunctionr.

**Step23**. **display()**

/* display function is used to display total sentiments for each thread

/*Two other Function Prototype are used outside of class forest_tree to solve different functionality

**Step24**. **isNumberList()**

/*isNumberList function is used to check a string whether it is in number list format or not

**Step25. subString()**

/* substring function is used to cut a string from Number List to get the current Parent node

**Step26**. **Main function ( )**

/* Main function under main body of the program and performs all Logics

The remaining descriptions and sample source code are listed out in appendix I

In this study using the above approach we implement graph-based opinion summary for analysing opinion given on sport domain in the next chapter

# CHAPTER FOUR

# IMPLIMENTATION AND EXPERIMENTS

This chapter presents the implementation and summary of experiment for graph-based opinion summary system for annotated text from corpus. The chapter includes the following tasks. Review of opinion information, design of opinion information, manual annotation of the whole data set, develop data structure and algorithms used for the summary of opinion information, finally evaluation of the result has been done.

In data collection (review of opinion) phase the identification of data source, that is used as input for the study, consideration and selection of languages, that opinion text were written by is some task. In designing opinion information task we identified the nature of the opinion in social network and prepared opinion thread based on sentiment flow between two opinion sentences. After that the next task is annotating the opinion information. In annotation phase the task done is assignment of sentiment polarity for opinion listed in each tread of discussions. Then assign tree structure based storage location. In opinion oriented graph representation the data structure and tree concept which focus on storage of annotated opinion for construction of tree structures are the critical concept we employed here. The other phase is developing data structure and algorithms. This task is to implement an algorism that extracting opinion polarities from threaded opinion graph and return the summary of opinion polarities. Finally the evaluation of the result has been also presented.

## 4.1. Data set preparation

We created our own data sets from two Facebook page, namely*Ethiopian National Football team (walya)* and *Ethiopian Football Federation* Facebook page. We first crawls 49 threads of discussions, containing **1235** reactions (comments and replies), comments are the first reaction text given to initial post from these **706** posts were written in English, and **529** written in Amharic script and **115** were ambiguous text expressed in idioms, slang, Misspellings, Laughter and Amharic message written in Latin script which is not clear. The total data set used for experiment is 1120 reaction texts

## 4.2. Data cleaning

It is known that the text written on Facebook, blogs, forums, YouTube and other social network comments; contains different structures of writing style. The main obstacles we identified are (e.g. Misspellings, ungrammaticality, shortening of words and/or repetition of letters and punctuation signs, the use of informal expressions, "urban" acronyms "sense of laughter").

**Informal expressions:** e.g. Guys.....got ur idea....thx. (i. e. Guys, I have got your ideas. Thanks!!) .

**Acronym***: e.g. (*OMG, i. e. oh my God)***,** many acronyms exchanged between the participants of the discussion chains are more of subjective or not common for other readers of the posts.

**Laughter**: e.g. *hehehehe, kkkkkkkk,hahaha!!!hihihi,* indicates Laughter. It is a variant indication of laughter that is used by Participants of discussion to express positive laughter or negative laugher.

The users are not always articulating their laugh for positive thing but also they articulate their laugh for negative thing. Thus it makes ambiguity in determining the polarity of such laughter. For this we considered only comments that are well structured and convey a clear message. The rest of user's opinions which are not structured well and not clear had been removed. In general these challenges and others are needs the work of text processing.

## 4.3. Data set

The total data set used for experiments were selected from 15 discussion threads that containing 1120 reaction text written in both language. To make annotation clear for linguistics men's, we purposively selected the most popular and controversial threads of discussions from the above data set. From these 698 of them are labeled as positive (+), and the remaining 422 opinions are labeled as negative, by linguistic men.

## 4.4. Annotation by expert

After the above preprocessing steps had been done the overall corpus of opinion thread chain is annotated as positive polarity (+) and negative polarity (**-**) sign and the polarity of the thread head (root) is annotated as (**0**) which indicates the opinion thread head. It is simply point of discussion and does not express any sentiment polarities. However, the polarity of neutral is not mentioned here; because neutral is taken as equal number of positive and negative in opinion chains. The polarity of the opinion sentence is labeled and annotated by professional staff of

language and literature from Jimma University as illustrated **inTable.1.** Five respondents were selected purposively by researcher for annotation purpose.

### 4.5. The Proposed Algorithms

Given annotated opinion thread, the proposed graph-based opinion mining mode operates in three steps. First, it reads the polarity of annotated opinion thread from the text file.

Then create a tree structure that contains parent-child relationships. Then Insert the polarity of a text into a created tree structure for the opinion thread. Next opinions measure is applied to the automatically created opinion-oriented graph. This is done in two ways. First is done at discussion chain and the second is the measure at the whole thread. Finally, all the polarity of annotated opinion thread is summarized into predefined categories: positive (+), negative (-) or neutral. The following algorithm is the high-level view of algorithms which describe how tree structures are created, how opinion polarities are stored and how sentiment polarity values are summarized into its pre-defined class.

**Algorithm 4.1: Back Tracing Algorithm (Source: own)**

1. *start*
2. **for** every annotated opinion discourse Thread **DT**
3. **for** every opinion polarity OP of discourse thread **Op**
4. **Read** Its OP of **DT**
5. **If** the OP of a text is found in **DT**
   5.1. **create** parent node **pn**
   5.2. **If** $p_n$ has a child node **Cn**
   5.3. **Create** a child node **Cn**
   5.4. **Read** its OP then
   5.4.1. **Insert** OP to **Cn**
   5.5. **Repeat** from step 4.3 to create new **Cn**
   5.6. **if** new **Cn** is created
   5.6.1. **then** Insert OP to new **Cn**
   5.7. **If** the **Cn**is leaf node
   5.7.1. Computes the AvgMsgOP

**5.7.2. If** AvgMsgOP is > 0

**5.7.2.1.** Assign OP to class (positive)

**5.7.2.2. If**AvgMsgOP is <0

**5.7.3.** Assign OP to class(negative)

**5.7.4. If**AvgMsgOP is=0

**5.7.4.1.** Assign AvgMsgOP to class (neutral)

**5.7.5. else** repeat from repeat step 4.3

**5.7.6. End**

The above algorithm implemented to automatically retrieve the final summary of opinion polarity from sample of data set collected from the pages about topic concerning Ethiopian national football Team.

| *Thread no* | Comments and Replies. | Opinion Chains | AvgMsgOP | | |
| --- | --- | --- | --- | --- | --- |
| | | | Aggregate (+) | Aggregate (-) | Neutral (0) |
| 1 | 24 | 8 | 6 | 2 | 0 |
| 2 | 27 | 5 | 4 | 1 | 0 |
| 3 | 96 | 28 | 18 | 8 | 2 |
| 4 | 87 | 24 | 18 | 5 | 1 |
| 5 | 30 | 9 | 2 | 7 | 0 |
| 6 | 39 | 12 | 8 | 3 | 1 |
| 7 | 66 | 17 | 10 | 2 | 7 |
| 8 | 105 | 35 | 31 | 4 | 0 |
| 9 | 48 | 14 | 7 | 7 | 0 |
| 10 | 24 | 6 | 6 | 0 | 0 |
| 11 | 72 | 20 | 20 | 0 | 0 |
| 12 | 6 | 2 | 1 | 1 | 0 |
| 13 | 6 | 1 | 0 | 1 | 0 |
| 14 | 12 | 4 | 3 | 1 | 0 |
| 15 | 198 | 36 | 19 | 12 | 6 |
| **Totals** | **1120** | **224** | **152** | **56** | **16** |

**Table4.1.statics of sample data set and its summary**

The annotated result shows that 224 posts are summarized as positive, 56 as negative and, 16as neutral, where 69% is "positive" 24% is "negative" and 7% is neutral. Here the neutral is computed from the occurrence of equal number of positive and negative opinions in discussion chains.

Applying our model to bigger discussions with hundreds of messages is interesting. The table above is the real discussion of the data set collected from Facebook sites stated before. From this table, we have identified that, the number of discussion threads, Reply-posts and chains that appear in the discussion. In this discussion thread we also observe many replay-posts that can make many chains or reaction with varied opinion polarities to discussion thread. Here as the number of chains increases for the particular thread; determining the most interesting opinion information in opinion-based graph is become quite complex.

So to simplify this complexity some global analysis and computation is important. The table below depicts the real statics of sample data set collected from topic concerning Ethiopian National Football Team and Ethiopian Football Federation. As indicated in table 4.1 the Entropy $H(v_x)$ is computed to know the most vital messages of the discussion at global level or for the entire discussion.

| Thread Post | Chains | (+) | (-) | Entropy $H(v_x)$ |
|---|---|---|---|---|
| Thread Post[1] | 8 | 6 | 2 | 0.244 |
| Thread Post [2] | 5 | 4 | 1 | 0.217 |
| Thread Post [3] | 28 | 18 | 8 | 0.278 |
| Thread Post [4] | 24 | 18 | 5 | 0.235 |
| Thread Post [5] | 9 | 2 | 7 | 0.230 |
| Thread Post [6] | 12 | 8 | 3 | 0.267 |
| Thread Post [7] | 17 | 10 | 2 | 0.255 |
| Thread Post [8] | 35 | 31 | 4 | 0.181 |
| Thread Post [9] | 14 | 7 | 7 | 0.296 |
| Thread Post [10] | 6 | 6 | 0 | 0 |
| Thread Post [11] | 20 | 20 | 0 | 0 |
| Thread Post [12] | 2 | 1 | 1 | 0.301 |
| Thread Post [13] | 1 | 0 | 1 | 0 |

| | | | | |
|---|---|---|---|---|
| *Thread Post* [14] | 4 | 3 | 1 | 0.244 |
| *Thread Post* [15] | 36 | 19 | 12 | 0.306 |

Table 4.2: The Entropy of *Thread Post*

### 4.6. Discussion of experimental results

From the table the summary result shows that 1120 comments 152 were summarized as positive, 56 as negative and, 16 as neutral, where 68% were "positive" 24% were "negative" and 7% were neutral. Here the neutral is computed from the occurrence of an equal number of positive and negative opinions. Here we can observe that about 20% of opinions are summarized between the three classes, it minimize the complexity of opinion chains to a minimum manageable level. The model summarizes more as opinion chain increased or become more complex. We also noticed that the *Thread Post* 15 is the most imperative message post; having **36** replies of which **19** were positive, **15** were negative. We notice that the *Thread Post*[15] has the highest entropy of all; indeed, this is the message post that has received replies with the highest variety of opinions. We also notice that the average opinion of all messages is positive which indicates the general tendency of discussion. In contrast to this we realize that the *Thread Post*[8] is the least imperative one having 35 replies of which 31 were positive, 4 were negative. We notice that the *Thread Post* [8]has the least entropy of all. Here this is the message that has received a plenty of positive replies that shows the general tendency of the discussion as positive.

In general, we conclude that the measure of entropies in the above result resides between 0 and 1. Hence, 0 Entropy indicates all of the opinion posts are from the same polarity either positive or negative. This shows that the opinion of the discussion thread is not debating issue. However, the maximum entropy 1 indicates that the balance number of opinion polarities, this measure indicates that the issue under discussion is a hot topic or debating point.

It is very important for the decision maker to give more emphasis on the issue. In this analysis, we considered only two polarities of opinions, positive and negative for calculating entropy measure because the neutral opinions are not important for this measure as it is already used for the equal number of positive and negative polarity in the model.

In general, the average entropy of the above message objects is below 0.5, it indicates that the opinions are more of the same type, which is more of a positive opinion. As we observe in the

above table, the entropy measures of

the opinion posted for the sport domain indicates users are posted more of positive opinion about the national team.

In general the entropy measure is used for Mining relevant opinion information from discussion thread and transforming it from a complex to a simpler one by extracting only the threads that has plenty of comments and replies from long discussions. Finally, the concept in the above algorithm implements and summarizes the threaded opinion reviews as depicted in the following figure.

The figure 4.1 bellow shows the Snapshot of the summary of opinion polarity in each thread of discussion.

**Figure 4.1: the Snapshot of the summary of opinion polarity for ach discussion thread.**

Here as depicted from the figure above the model could generate the summary of opinion polarity for each opinion thread.

In fact in the absence of such model determining the final tendency (summary) of even a single thread is not easy. Because the controversies, debates, agreement and disagreement seined between users opinion are forming a long discussion chains. Hence, the summarized opinion information can minimize the difficulty in identifying the most imperative information.

As a result concluding the final inclination of opinion polarity throughout the swing (chain) is difficult and confusing for decision makers in the absence of such model. So the model displays the summary of opinions polarity from each opinion thread. The summarized opinion information by the method minimizes difficulty in identifying the most imperative information. To this end the result provided from the model can support decision maker to identify quickly, the main urgent topic of discussion that have to give emphasis from the bulky of information post on the page about the given event.

### 4.7. Experiments

### 4.8. Evaluation method

Precision and recall, are the well-known Performance evaluation parameters of information retrieval (IR), we were used for our model validation. Precision measures the exactness of a classifier. Precision is the ratio of the number of opinion classified correctly to the total number of opinions classified in a given category. A high precision means less false positive, while a lower precision means more false positives.

$$P = \frac{TP}{TP+Fp} \ldots\ldots\ldots\ldots\ldots\ldots\ldots\ldots\ldots\ldots\ldots\ldots\ldots\ldots\ldots\ldots 4.1$$

Where, $TP$ denotes the number of opinions reviewed which are classified correctly and $Fp$ denotes the number of opinions reviews which are classified incorrectly.

Recall measures the completeness or sensitivity of a classifier. It is the ratio of $TP$ and the whole opinion reviews belonging to the category. A high recall means less false negative, while lower recall means more false negatives.

$$R = \frac{TP}{TP+FN} \ldots\ldots\ldots\ldots\ldots\ldots\ldots\ldots\ldots\ldots\ldots\ldots\ldots\ldots\ldots\ldots\ldots 4.2$$

Where, $FN$ denotes the number of opinion reviewed which are missed by the classifier, i.e. neither classified correctly or incorrectly (unclassified category).

There is trade-off between precision and recall. Greater precision decreases recall and greater recall leads to decreased precision.

The F-measure is the harmonic mean of P and R and takes account of both the measures. As a result, F-measure is defined as follows:

$$F = \frac{2PR}{P+R} \ldots\ldots\ldots\ldots\ldots\ldots\ldots\ldots\ldots\ldots\ldots\ldots\ldots\ldots\ldots\ldots\ldots\ldots\ldots\ldots\ldots\ldots\ldots\ldots\ldots\ldots..4.3$$

The experiment is conducted by comparing the direct meaning of opinions conveys with the message it conveys in discussion chains.

### 4.9. Performance evaluation

All **1120** opinions are compared with the expected meaning of the sentiment mining model for conducting the experiment in the absence of opinions chain. The results measured by, precision, recall and F-measure is presented in table 4.3 as follows.

**Table 4.3 summary of experimental result**

| Reviews | Class | Precision | Recall | F-measure |
|---|---|---|---|---|
| In discussion chain | Positive | 0.789 | 0.873 | 0.828 |
| | Negative | 0.88 | 0.673 | 0.762 |

In this section, we evaluate the model. Typically, the evaluation is done linguistically, by comparing expert's annotated results in the presence of discussion threads.

The following Table 4.4 shows that the performance evaluation is done on 15 opinions thread provided by five annotators using Cohen's Kappa: to determine the agreement disagreement level between annotators.

**Table 4.4: annotation of opinion polarity by linguistic men**

| Reviews | Annotation | Sentiment class | Numbers |
|---|---|---|---|
| In discussion chain | Annotator 1 | Positive | 640 |
| | | Negative | 480 |
| | Annotator 2 | Positive | 692 |
| | | Negative | 428 |

| | Annotator 3 | Positive | 660 |
|---|---|---|---|
| | | Negative | 460 |
| | Annotator 4 | Positive | 460 |
| | | Negative | 456 |
| | Annotator 5 | Positive | 720 |
| | | Negative | 480 |
| | **Cohen's kappa** | Positive | 0.05 |
| | | Negative | 0.058 |

The average annotation result done by all annotators from the total discussion is 60%, 40% positive and negative respectively. The Cohen's kappa result (i.e. 0.05 positive, 0.058 negative) shows from the total dataset, 95% agreement in positive and 94.2% in negative between the annotators.

The evaluation is done by showing the advantages and the corresponding information that can be extracted from an opinion-based graph as compared to that of expected result from the straight forward opinion mining by natural language processing techniques.

Typically, the comparison is done between the expert annotated result categorized by the algorithm in the presence of opinion chains and that of the direct meaning of the opinions obtained (in the absence of opinions chain).

These results were also tested by using machine learning developed online available tool that has best performance. For this, we used the known performance evaluation matrixes, precision, recall, and F-measure as depicted table 4.5.

| Reviews | Confusion matrix | | | performance evaluation | | |
|---|---|---|---|---|---|---|
| Using machine learning **(SVM)** | Class | Positive | negative | Precision | recall | F-measure |
| | Positive | 613 | 103 | 0.84 | 0.85 | 0.84 |
| | Negative | 116 | 288 | 0.73 | 0.71 | 0.71 |

**Table 4.5: performance evaluation result Using SVM**

### 4.10. Experimental Result

The experiment shows that 74% accuracy when we compared these results with our model which works on linguistically annotated data the overall accuracy of the system has shown smaller

amount of result. The model shows that about 67.5% of opinions are annotated as positive and 33.5% is annotated as negative. Whereas experts annotate 60% of it as positive and 40% as negative.

The experimental results show that there is a thread-off among positive and negative polarities compared to experts annotation result. This is due to some of the opinions given to the initial posts are indirect or depend on previous one, the order in which these opinions text comes has an impact on the message it tends to expresses. Consider the following discussion thread.

"የኢትዮጵያእግርኳስፌዴሬሽንየዲሲፕሊንጥፋትፈፀሙዋልባላቸውሁለትተጫዋቾችላይየገነዘብቅጣትአስተላለፈ::"
**1.1** It's not good decision
**1.1.1** you are right
**1.1.2** it's not English premier league men"
**1.1.3** I do not agree with you, They are so crap
**1.2** I think it is better to suspend both players from the national team
**1.2.1** ሳተወቅትናገራሌ
**1.2.2** No የለነሱኢንጎዳሌን
**1.2.2.1** I really impressed by those people who gave their comment wanting the two players back in to the national team

The reply on 1.1 "you are right" is conveys positive opinion when we consider alone in the absence of opinion chain, However in the above dialog or opinion chain it mean that it indirectly implies disagreement with decision taken by EFF (Ethiopian Football federation). In contrast the reply on 1.1.3 '' I do not agree with all of you. They are so crap'' is conveys negative opinion in the absence of opinion chain, However in the presence of above dialog or opinions chains it mean that it indirectly implies agreement with decision taken by EFF.

Due to such opinion text exist in positive but express negative and vice versa the total accuracy of the model is reduced.

From these, we observe that the meaning of opinion and its' relationships in discussion threads are determined through *back tracing of the previous post*. This cannot be resolved from the use of solely straightforward text mining or Computational linguistics in the absence of the opinion oriented graph.

The other challenge that affects the annotation result in both experiments is the use of an informal expression like proverb, pragmatics, idiomatic, slang, semantic, Misspellings, Laughter etc. as annotation result may vary from expert to expert in such ambiguous text.

## 4.11. Sway analysis of discussion chain.

The second analysis method we employed is Sway analysis. Sway analysis is the analysis of discussion chain. This approach consists of three things: For each discussion, we identified the data structure and storage for opinions polarities. Then we find out the type of opinions that appear in the messages of the discussion, and finally we retrieved automatically the summarized result from tree structure. The focus of this evaluation is on the presentation of opinion-based model which facilitates the discussion analysis from information overloading and not on the ways to identify opinion data as subjective or objective in text processing.

In Table 4.1 we give some information of the discussions we analyzed such as the number of reply posts exchanged between nodes. The column «Opinion chains» shows the number of discussion chains which are characterized by a sequence of node labeled with a positive or negative opinion and aggregated positive or negative result.

The experiments with real discussions allowed us to observe the characteristics of online discussions and confirm the importance of opinion-based graphs since they capture information that cannot be provided by opinion retrieval algorithm in machine learning. The machine learning techniques do not handle opinions from the user point of view.

E.g.: for the Amharic proverb"*ያገኑንስርዶበገሩበሬ!*"is the reply for message object, "***Would you believe that, if we outsource the coach? He/she will bring expected result?***" It is clear that the reply show disagreement for the message object. The reply message is a kind of proverb. This is challenging in NLP to determine its meaning. The graph-based opinion mining generates better result for decision maker since the model summarizes the labeled opinions by human experts.

In more detail, from Table 4.1, we see that in the majority of the discussions, the number of Opinion chain is less than the number of reply messages. This illustrate that many messages were raised from disagreement in the topic than an agreement topics. This points out the importance of our proposed model in overwhelming information overload and allows a discussion analyst to concentrate only on the parts of the discussion that contain vital opinions without waste time in analyzing the whole discussion. The graph based opinion enable us to see at a quick look how the opinion flows inside the discussion, how the positive messages alternate with the negative ones. This useful information is not provided by one directional analysis of text mining.

In our experiments we identified the «opinion chains» which are characterized by a series of replies that holds sentiments. Having represented the debate from the point of view of message

objects instead independent analysis allows us to identify quicker interesting opinion discussion chains in decision making. A message that has received few but varied positive and negative opinions can be more interesting than one that has received plenty of messages that are all positive or negative. In our model, we measure this by the entropy. This measure allows the selection of the imperative issues to be given attention.

Furthermore the opinion-based graph extraction allows us to observe the sequences of opinions. A sequence of two positive nodes in a discussion may show agreement between the users. The prerequisite, though, for assuming agreement is that the messages express a positive opinion on the same argument.

For example, the message «*I do not agree with the newly selected coach for the national team*» that receives as a reply the message «Yeah, they are not going with our style» is a sequence of two negative opinions that shows agreement between them.

On the contrary, the message «*It is bad to say this*» as a reply to the message «*This is a boring game*» does not point out agreement even though both messages express a negative opinion and they are connected.

In general, the opinion-based graph differs from NLP or text mining techniques in different ways. Firstly, opinion-based graph tells the relationship between two or more than two sentences whereas NLP or text mining techniques tells only the polarity of a single sentence or the whole review without considering the sentiment relation in discussion thread. It only calculate the whole sentiment score in a given reviews to classify the reviews as positive or negative. Secondly, unlike the NLP techniques, graph-based opinion summarization tells which opinion threads are more important for the purpose of decision making. Finally, the method of computational linguistic is more effective for the direct opinion information than indirect opinion information. However, the opinions of the domain are more of indirect and informal message like proverb. Since the polarity of opinion information is determined by the experts, the graph-based opinion analysis gives a convincing result. The result from the model show that the entropy $H(v_x)$ of thread of discussions is between 0 and 0.5 this indicates that the opinion posted were the same type and positive. This reveals that the feedbacks given on sport event are positive feedbacks that encourage National team.

# CHAPTER FIVE

## CONCLUSION AND RECOMMENDATION

### 5.1. CONCLUSION

The web has dramatically changed the way that people express their views and opinions. They can now post reviews of the event at social Network sites and express their view on almost everything on Facebook, forums and blogs etc. This online explosive posts is considerable source of information for many practical applications, Like sport federation. It could find out the opinions of people and sport-analyst to know the strength and weakness of the national team if they want to make effective decision. However, it is inconvenient for a human reader to make relevant decision. As a result, automated opinion measure and extraction systems are needed.

This research work has strained to go through the techniques of graph-based opinion mining method that summarizes Facebook discussion from the corpus of annotated opinion thread. The technique permits us to determine the sentiment orientation or relation from large discussion chain. The model involves both direct and indirect meaning of opinion sentence wants to convey in the presence of dialog. The proposed opinion measures and extraction algorithm offer a sentiment-oriented analysis and generate aggregate summary result of the discussion.

 The evaluation of the model shows that graph-based opinion summarization provides information that cannot provided by straight forward text mining, as such, it is useful and it has a lot to offer to the discussion analysis, even if the data annotation done linguistically by expert. The evaluation of the model shows that graph-based opinion summarizing model works hundred percent for any annotated opinion text done either by text mining or linguistic men. However, improving the annotation result prior to graph-based summary need further study due to the complexity of social network text. We put forward on the use of machine learning techniques to enhance the efficiency of the developed model.

### 5.2. Contributions of the study
Some of the main contributions of this research work are given below.
- Provides a summary of opinion information from thread structure in discussion.
- Define and present measures that can give important information regarding the opinion flow.

- Developed a model that extracts opinion summary by holding sentiment relation between nodes.
- Overcomes information over load (difficulty in identifying relevant information) while reading long discussion thread.
- The model allows us to identify quicker interesting opinion discussion chains in decision making.

5.3. Major challenges of the study

One of the major challenges confronting us in this study is the nature of data sources. It is known that posts on Face books are written in more than one language. As a result the opinion posts are requires more than one model to classify a post and need more than one language expertise to annotate the whole dataset. Developing separate model might be easy for language study. However, opinions are written in different language they have relationships in discussion chain, though developing a separate model for individual language makes opinion information meaningless or none sense because the given opinion is only meaning full in the presence of discussion chain or dialogues. So that developing a multi-lingual model from the study of computational linguistic or text mining is so important for this problem. Nevertheless, developing a multilingual model is not an easy task. It requires a more researcher time and need further investigation of the languages. In another

One of the major challenges confronting us in this study is the nature of data sources. It is known that posts on Face books are written in more than one language. As a result the opinion posts are requires more than one model to classify a post and need more than one language expertise to annotate the whole dataset. Developing separate model might be easy for language study. However, opinions are written in different language they have relationships in discussion chain, though developing a separate model for individual language makes opinion information meaningless or none sense because the given opinion is only meaning full in the presence of discussion chain or dialogues. So that developing a multi-lingual model from the study of computational linguistic or text mining is so important for this problem. Nevertheless, developing a multilingual model is not an easy task. It requires a more researcher time and need further investigation of the languages. In other hand, detecting the relationships between two opinion sentences is a challenging field from computational linguistics. So the combinations of

computational linguistic and graph-based opinion mining system can produce a convincing result. The other challenge is the data quality in social network. As stated in section 4.5. The text written on Facebook, blogs, forums, reviews, micro blogs, social network comments; includes different structures of writing style. As a result we hope these challenges and many others are needs the combination of text mining and graph theory to achieve a better result.

## 5.4. Recommendation

Even though this study attempts to develop opinions-oriented graph for stated social media, developing a full-fledged, fully functional and a more efficient analysis is still required. We believe that future precise opinion-oriented graph analysis is helpful. More measures need to be defined and more large-scale experiments are needed for the formal validation of the model. One future direction is to combine the user-based and the graph-based opinion mining analysis model for an improved decisions. In this the user-based graphs is helpful to extract the user's behavior in the discussion domain. Furthermore, an interesting future issue is to add the time dimension in our model. This will permit monitoring how opinion changes over time. In this way, we could observe whether people become more satisfied with certain changes about the given topic, or even whether people are finally convinced after a long discussion. The information extracted by the graph-based opinion model can be used in many ways. We have experimented by using it in order to rank Facebook posts from the most to the least interesting. This is a combination of many criteria such as how many chains a message causes, whether it receives chain that contain similar opinions, whether these opinions have the same strength or not. Initial results are hopeful but more extensive experiments are needed. Another strategy that can be considered in the future is Automatic sentiment detection and annotations for opinion oriented graph model. It can be achieved by developing multilingual opinion mining with combination of opinion oriented graph; this is also our future direction to carry out further research in order to find out the model that facilitates this identification.

**References**

[1]. B. Pang and L. Lee, "Opinion Mining and Sentiment Analysis," Foundations and Trends in Information Retrieva, 2008.

[2]. Po-Wei Liang and Bi-Ru Dai. Opinion Mining on Social Media Data, DOI 10.1109/MDM. 2013.73, 978-0-7695-4973-6/13 © 2013 IEEE.

[3]. Dar and A. Jain, "opinion mining through on to tree," vol. 4, 2014.

[4]. M. Nidhi, "Classification of Opinion Mining Techniques," International Journal of Computer Applications, vol. 56, no. 13, pp. 1-6, 2012.

[5]. B. Liu, "Sentiment Analysis," in 5th Text Analytics Summit, Boston, 2009.

[6]. M. Sigrid, C. Paolo and D. Luca, A Hybrid Method for Sentiment Analysis, France: Statistical Analysis Software (SAS) press Grenoble, 2008.

[7] M. K. JayashriKhairnar, "Machine Learning Algorithms for Opinion Mining and Sentiment Classification," International Journal f Scientific and Research Publications, vol. 3, no.6, 2013.

[8] B. Liu, "Sentiment Analysis and Opinion Mining. Synthesis Lectures on Human Language Technologies," in Morgan & Claypool,2017.

[9]. B. Pang, L. Lee and S. Vaithyanathan, "Thumbs up? Sentiment classification using machine learning techniques", in In Proceeding of the conference on empirical methods in natural language, pp. 79-86, Philadelphia, US, 2018.

[10]. SwapnaSomasundaran, JanyceWiebe, and Josef Ruppenhofer "Discourse level opinion interpretation", In Proceedings of COLING, 2008.

[11]. M. Taboada, J. Brooke, M. Tofiloski, K. Voll, and M. Stede," Lexicon-based methods for sentiment analysis.Computational

Linguistics", Vol. 1, pp.1–41,2010.

[12]. Popescu and O. Etzioni. "Extracting product features and opinions from reviews", In Proc. of the Conference on Empirical Methods for Natural Language Processing (EMNLP), pp. 339–346, Vancouver, Canada, 2014.

[13]. Boy and M-F. Moens. "A machine learning approach to sentiment analysis in multilingual web texts. Information Retrieval", Vol. 12, no. 5, pp.526–558, 2009.

[14]. Pak and P. Paroubek. "Twitter Based System:Using Twitter for Disambiguating Sentiment Ambiguous Adjectives" In roceedings of the 5th International Workshop on Semantic Evaluation, Association for Computational Linguistics, pp.436-439, 2010b.

[15]. Alec Go, RichaBhayani, and Lei Huang, "Twitter sentiment classification using distant supervision", CS224N Project Report, Stanford,2009.

[16]. Carley, K.M., Diesner, J.: "AutoMap: Software for Network Text Analysis", CASOS (Center for Computational Analysis of Social and Organizational Systems), ISRI, CMU,2015)

[17]. Van Atteveldt, W.H.: "Semantic Network Analysis. Techniques for Extracting, Representing, and Querying Media Content." VrijeUniversiteit, Amsterdam, 2018)

[18]. Pitler, E., Raghupathy, M., Mehta, H., Nenkova, A., Lee, A., and A. Joshi. "Easily identifiable discourse relations." In Proceedings of COLING, pages 87-90, 2008

[19]. Somasundaran, Swapna. "Discourse-level relations for Opinion Analysis." PhD Thesis, University of Pittsburgh. 2010

[20]. Prasad, R., Dinesh, N., Lee, A., Miltsakaki, E., Robaldo, L., Joshi, A. K., and Webber, B. L. "The Penn discourse treebank 2.0." In LREC. European Language Resources Association.,2008

Marcu, Daniel. "The Theory and Practice of Discourse Parsing and Summarization", MIT Press, Cambridge, MA. 2000.

[21]. Asher, Nicholas andBenamara, Farah and Mathieu, Yvette Yannick. "Distilling opinion in

[22] discourse" : A preliminary study. In Proceedings of Computational Linguistics (CoLing). 2008.

[23]. Fisher, D., M. Smith, and H.Welser, "You are who you talk to Detecting roles in Usenet newsgroups", Proceedings of the 39th Annual HICSS. IEEE Computer Society, 2015.

[24]. Jindal and L. Nitin and Bing, "Identifying comparative sentences in text documents: in Proceedings of ACM SIGIR Conf," in on Research and Development in Information Retrieval, SIGIR, 2006,2006a.

[25]. B. Liu, "Sentiment Analysis and Subjectivity," in Handbook of Natural Language Processing, Second Edition ed., 2010.

[26]. X. Feiyu and C. Xiwen, "Opinion Mining," Saarbrucken,Germany, 2007.
A. Rakesh, R. Sridhar, S. Ramakrishnan and X. Yirong, "Mining Newsgroups Using Networks Arising From Social Behavio," IBM Almaden Research Center, no. CA 95120.

[27]. L. Bing, "Sentiment Analysis and Opinion Mining," in Synthesis Lectures on Human Language Technologies, 2012.

[28]. Jindal and L. Nitin and Bing, "Identifying comparative sentences in text documents: in Proceedings of ACM SIGIR Conf," in on Research and Development in Information.

[29]. M. Nidhi, "Classification of Opinion Mining Techniques," International Journal of Computer Applications, vol. 56, no. 13, pp. 1-6, 2012.

[30]. B. Liu, "Sentiment Analysis and Opinion Mining. Synthesis Lectures on Human Language Technologies," in Morgan & Claypool, 2012.

[31]. A. Jacob, R. Sara and M. Kathy, "Annotating Agreement and Disagreement in Threaded Discussion," Columbia University, 2006.

[33]. B. Alexandra, "method and Resources for sentiment Analysis in Multilingual Documents of Different text types," pp. 12-18, 2011.

[34]. G. Angulakshm and D. Manickachozia, "An Analysis opinion miningtechniques and tools," international Journal of Advanced in computer and cominication engineering, vol. 3, no. 7, 2014.

[35]. Ganapathibhotla, Murthy and L. Bing, "Mining opinions in comparative sentences in Proceedings of International Conference," in Computational Linguistics, COLING, 2008.

[36]. L. Chien-Liang, W. Hsaio, C.-H. Lee, G.-C. Lu and J. Emery, "Movie Rating and Review Summarization in Mobile Environment," 2012.

[37]. L. Chien-Liang, W. Hsaio, C.-H. Lee, G.-C. Lu and J. Emery, "Movie Rating and Review Summarization in Mobile Environment," 2016.

[38]. T. Mikalai and P. Themis, " Survey on mining subjective data on the web: Data Mining Knowledge Discovery," Springer, 2012.

[39]. Lloret, A. Balahur, J. M. Gómez, A. Montoyo and M. Palomar, "Towards a unified framework for opinion retrieval; mining and summarization," Journal of Intelligent Information Systems , p. 11 747, 2012.

[40]. m. Hu and B. Liu, "Mining opinion features in customer reviews," Proceedings of AAAI, p. 755–760, 2004.

[41]. N. Mishra, "Classification of Opinion Mining Techniques," International Journal of Computer Applications, vol. 56, no. 13, pp. 1-6, 2018.

[42]. O. Alexander, "Sentiment Mining for Natural Language Documents," Australian National University, 2009.

[43]. T. Joachims, "Text categorization with support vector machines: Learning with many relevant features," in In Proceedings of the 10 European Conference on Machine Learning Springer, Berlin, 1998.

[44].

[45]. Pang, L. Lee and S. Vaithyanathan, "Thumbs up? Sentiment classification using machine learning techniques," in In Proceeding of conference on empirical methods in natural language, Philadelphia,US, 2002..

J. Felipe, P. Almeida and Mattosinho, "Mining Product Opinions and Reviews," 02 09 2013. [Online].

[46]. Jurafsky.D and Martin.H, "Speech and Language Processing: An Introduction to Natural Language Processing,," in Computational Linguistics and Speech Recognition, Prentice Hal, 2019.

[47]. A. Kennedy and D. Inkpen, "Sentiment Classification of Movie and Product Reviews Using Contextual Valence Shifters:Intelligence, Computational Intelligence," 2006.

[48]. L. Bing, "Sentiment Analysis and Opinion Mining," Synthesis Lectures on Human LanguageTechnologies, 2012.

[49]. A. Samaneh and Moghaddam, "Aspect-based opinion mining in online reviews," simonfraser university, 2017.

[50]. B. Esuli and F. Sebastiani, "Determining term subjectivity and term orientation for opinion mining," in In Proceedings EACL-06 the 11rd Conference of the European Chapter of the Association for Computational Linguistics, 2006.

[51]. M. Michae, L. Daconta, J. Obrst, Kevin and T. Smith, "The Semantic Web:," in A Guide to the Future of XML, Web Services, and Knowledge Management, 2003.

[52]. C. Xiwen and X. Feiyu, "Fine-grained Opinion Topic and Polarity Identification," in European Language Resources Association (ELRA), 2006.

[53]. L. Zhao and L. Chunping, "Ontology Based Opinion Mining for Movie reviews .," Springer-Verlag Berlin Heidelberg, 2009.

[54]. J. Tatemura, "Virtual reviewers for collaborative exploration of movie reviews," in

[55]. Proceedings of Intelligent User Interfaces (IUI), p. 272–275, 2000.

L. Terveen, W. Hill, B. Amento, D. McDonald and J. Creter, "PHOAKS: A system for sharing recommendations," Communications of the Association for Computing Machinery (CACM), vol. 40, no. 40, p. 59–62, 1997.

[56]. C. Cardie, C. Farina, T. Bruce and E. Wagne, "Using natural language processing to improve eRulemaking," in Proceedings of Digital Government Research (dg.o), 2006.

[57]. A. Devitt and K. Ahmad, "Sentiment Polarity Identification in Financial News: A Cohesion Based Approach," Proceedings of the 45th Annual Meeting of the Association of Computational Linguistics, p. 984–991, 2007.

[58]. Somasundaran, Swapna, N. Galileo, G. Lise and J. Wiebe, "Opinion graphs for polarity and discourse classification," in Proceedings of the 2009 Workshop on Graph-based Methods for Natural Language Processing, 2018.

[59]. A. Stavrianou, J. Velcin and J.-H. Chauchat, "A combination of opinion mining and social network techniques for discussion analysis," ERIC Laboratoire, Université Lumière Lyon 2, 2012.

[60]. Zhou, Lanjun, L. Binyang, G. Wei, W. Zhongyu and W. KamFai, "Unsupervised discovery of discourse relations for eliminating intrasentence polarity ambiguities," In Proceedings of the Conference on Empirical Methods in Natural Language Processing, 2009.

[61]. M. A. Rob Abbott, F. Jean, E. R. B. Tree and K. Joseph, " " Recognizing disagreement in informal political argument," In Proceedings of the Workshop on Language in Social Media, 2011.

[62]. M. Emily, Bender, T. Jonathan, M. O. Morgan, Z. Mark, B. Hutchinson, A. Marin, B. Zhang and a. M. Ostendorf, "Annotating social acts: authority claims and alignment moves in wikipedia talk pages," In Proceedings of the Workshop on Languages in Social Media, 2011.

| [63]. | M. Helander, R. Lawrence and Y. Liu, "Looking for great ideas: Analyzing the innovation jam. In KDD'07: Proceedings of the ACM SIGKDD International," in Conference on Knowledge Discovery and Data Mining, 2007. |
| [64]. | . S. Maurel, P. Curtoni and L. Dini, "'analyse des sentiments dans les forums," in In Atelier Fouille des Donnéesd''Opinions (FODOP 08), 2008. |
| [65]. | J. Zhang, M. Ackerman and L. Adamic, "Expertise networks in online communities: Structure and algorithms," in In Proc. of the 16th International Conference on World Wide Web, 2007. |

Appendix I

**/\*Description of Basic Variables and others steps**

Step27. b=variable of multi_tree

Step28. tmp=variable of integer that indicates polarity (i.e.( −) or (+))

Step29. c= variable that holds array of word file

Step30. fin = handler of input file

Step31. Open file with handler fin using build in open function

Step32. While file is not end

Step33. Set c to current input word from the file

Step34. IF c is plus OR minus OR (startth is true and c is 0)

Step35. Set nextnumlist true

Step36. IF startth is true and c is 0

Step37. Set startth to false and tmp to 1

Step38.END IF

Step39. IF c is plus

Step40. Set tmp to 1

Step41. ELSE

Step42. Set tmp to -1

Step43. END IF

Step44. Send tmp to insert function

Step45 .END IF

Step46. ELSE IF isNumberList returns true and nextnumlist is true

Step47. Set nextnumlist to 0

Step48. Send nodenum and c to strcpy

Step49. Send chknodenum and c to strcpy

Step50. Send chknodenum to substring

Step51. IF length of chknodenum is greater than 0

Step52. IF chknodenum is different from current node

Step53. Find current

Step54. End If

Step55. Else

Step56. Set startth to 1

Step57. IF first is true //skip displaying before the completion of processing threads

Step58. Set first to 0

Step59 Else

Step60. Set order to nodenum less than 1

//Step61. Count total aggregate result of thread

Step62. END IF

Step63. END IF

Step64. End IF

Step65. END WHILE

//Algorithm to measure and extract sentiment

Step66. Isempty()

/*operation for Isempty function

Step67. Match root with NULL

Step68. Return result

Step69. empty()

/* operation for empty function

Step70. Set root to NULL

Step71. Set pos to 0

Step72. Set neg to 0

Step73. Set nut to 0

Step74. Insert()

/* operation to insert function

Step75. Accepts polarity

Step76. Create new pointer object of tree_node with name t

Step77. Set polarity to data of t(t->data)

Step78. Set nodenum to pnodenum of t(t->pnodenum)

Step79. IF root is empty

Step80. Set t to root

Step81. ELSE

Step82. Set t to next empty node

Step83. END IF

Step84. call_countLeave()

/* operation of call_countLeave

Step85. Call countLeave

Step86. countLeave()

/* operation of countLeave

Step87. Accepts pointer of tree_node

Step88. IF tree_node is NULL

Step89. return 0;

Step90. ELSE

Step91. IF first child node of tree_node is NULL

Step92. IF aveMsgOpinion of tree_node is greater than zero

Step93. pos++

Step94. ELSE IF aveMsgOpinion of tree_node is less than zero

Step95. neg++

Step96. ELSE

Step97. nut ++

Step98. END IF

Step99. ELSE

Step100. WHILE tree_node has a child

Step101. Repeat step 80 by sending child of tree_node

Step102. END WHILE

Step103. END IF

Step104. END IF

**Step105 /* end algorithm**