



## Sensational News And News Maker Using Data Mining And Opinion Extraction

---

Aswathy Dhanapalan

EasyChair preprints are intended for rapid dissemination of research results and are integrated with the rest of EasyChair.

July 1, 2021

# Sensational News And News Maker Using Data Mining And Opinion Extraction

Aswathy Dhanapalan

*School of SCIS*

*UOH*

Hyderabad, India

20MCMT21

aswathydhanapalan2017@gmail.com

**Abstract**—The current outbreak and widespread of microblogging and social networking websites are reshaping many aspects of the modern day social interaction. Social Media today has exceeded the limits of entertainment or simple social interaction contexts. Social media now is better described as a living organism, that has a structure and a soul. The Social Media now reacts the pulse of the people; it reacts to their emotions, and interacts with their opinions. Analyzing and monitoring the content of Social Media can bring about some valuable insights, of which the conventional media means weren't able to convey. Social Networking is the medium which allows million users to share their opinion, ideas, views and huge amount of data is presented for internet users and a lot of data generated too. This project proposes the paradigm to extract the sentiment from famous micro blogging service, called twitter. Twitter is one of the most popular portals, where people post their opinions, views for everything. In this project, data mining techniques are used to automatically classify the sentiments of Tweets taken from Twitter dataset. This project tackles the concept of sensational news and newsmaker detection within Twitter. The project identifies the dimensions of similarity and interaction between news and tweets that provides the tools to calculate sensational news maker and news

**Index Terms**—Data Mining, Sentimental Analysis, Microblogging

## I. INTRODUCTION

This project gives more beneficial solution by providing in depth detailed information of data. In this context this implementation serves processed information of tweets accessed from Twitter Server. This project proposes a new system that delivers large database of Social Networking Site (SNS) called Twitter. Many Third party application are building based on SNS like Twitter, they need to have processed data from their operational purpose. The main stream of the applications is visualization applications. The purpose of this document is to present a detailed description of data mining and opinion extraction used in finding the sensational news and news-maker for a particular period along with the sentiments of people towards that particular news. It explains the purpose and feature of the system, the interfaces of the system what the system will do, the constraints under which it must operate on. This is mainly intended for uses in news medias.

The purpose of this project is to develop an automated way of gauging the sensational news for a particular time period,

news icon for that particular time and the mood of people on social media towards that particular news, specifically Twitter. The software should extract the data from a large database and determine the sentiment of the Twitter community with respect to a given topic. The main scope of data mining are:

- Automated detection of news maker.
- Detection of controversial or tweets with maximum retweets and shares.
- Detecting the statements and policies taken by these newsmakers that has become the topic of discussion.
- Orientation of the people towards these statements.

## II. SYSTEM ANALYSIS

### A. Existing System

The existing system used only structured data analysis. It used phrases and grammar information only to a limited extent. The user has to search a lot in order to find the relevant information from various on-line sites. Analysis were done on small amount of data in these techniques. It does not work in real time environment since it uses hadoop platform for analysis. The existing methods of opinion mining are:

- Product review mining.
- Sentiment analysis of news articles.
- Tracking sentiments toward topics over time.
- Orientation of the people towards these statements.

### B. Proposed System

Proposed system finds the sensational news and news-maker for a particular period along with the sentiments of people towards that particular news. It explains the purpose and feature of the system, the interfaces of the system what the system will do, the constraints under which it must operate on. This is mainly intended for uses in news medias.

## III. SYSTEM DESIGN

Design of the system includes mainly two steps:

- System design
- Detailed design

In System design a structural framework for the entire system is created. It is done in such a way that related part

come under particular groups. Thus after the system design, a network of different groups is obtained. It is the high-level strategy for solving the problem and building a solution. It includes the decision about the organization of the system into subsystems, the allocation of subsystems to hardware and software components, and major conceptual and policy decisions that form the framework for the detailed design.

In detailed design, each group is studied in detail and the internal operations are decided. Based on this, the data structures and the programming language to be used are decided. Apart from detailed design, the system design can be grouped into physical design and structural design. The physical design maps out the details of the physical system and plans the system implementation and specifies the hardware and software requirements.

Structured design is an attempt to minimize the complexity and make a problem manageable by subdividing into smaller segments, which is called modularization or decomposition. In this way structuring minimizes intuitive reasoning and promotes maintainable provable of systems. The structured design partitions a program into small, independent modules. They are arranged in a hierarchy that approximates a model of the business and is organized in a top-down manner

Logical design proceeds in a top-down manner. General features, such as reports and inputs are identified first. Then each is studied individually and in more detail. Hence the structured design is an attempt to minimize the complexity and make a problem.

#### IV. SYSTEM WORKING

Based upon the functions, the project had been divided to 6 modules which are mentioned as follows :

- Tweet Retrieval
- Text Preprocessing
- Named Entity Recognition
- Text Summarization
- Sentiment Calculation
- Data Visualization

##### A. *Tweet Retrieval*

In this first part, we will see different options to collect data from Twitter. In order to have access to Twitter data programmatically, we need to create an app that interacts with the Twitter API. Tweets are retrieved from twitter by using this API. Twitter provides REST APIs you can use to interact with their service. There is also a bunch of Python-based clients out there that we can use without re-inventing the wheel. In particular, Tweepy in one of the most interesting and straightforward to use.

##### B. *Text Preprocessing*

I start our analysis by breaking the text down into words. Tokenisation is one of the most basic, yet most important, steps in text analysis. The purpose of tokenisation is to split a stream of text into smaller units called tokens, usually words or phrases. While this is a well understood problem with several

out-of-the-box solutions from popular libraries, Twitter data pose some challenges because of the nature of the language. The tokenisation is based on regular expressions. To improve the richness of your pre-processing pipeline, we can improve the regular expressions, or even employ more sophisticated techniques like Named Entity Recognition.

##### C. *Named Entity Recognition*

Named Entity Recognition (NER) is a subproblem of information extraction and involves processing structured and unstructured documents and identifying expressions that refer to peoples, places, organizations and companies. NER is a fundamental task and it is the core of natural language processing (NLP) system. NER involves two tasks, which is firstly the identification of proper names in text, and secondly the classification of these names into a set of predefined categories of interest, such as person names, organizations (companies, government organisations, committees, etc), locations (cities, countries, rivers, etc), date and time expressions.

For humans, NER is intuitively simple, because many named entities are proper names and most of them have initial capital letters and can easily be recognized by that way, but for machine, it is so hard. One might think the named entities can be classified easily using dictionaries, because most of named entities are proper nouns, but this is a wrong opinion. As time passes, new proper nouns are created continuously.

Therefore, it is impossible to add all those proper nouns to a dictionary. Even though named entities are registered in the dictionary, it is not easy to decide their senses. Most problems in NER are that they have semantic ambiguity, on the other hand, a proper noun has Different senses according to the context.

##### D. *Text Summarization*

Summarization is a process of generating summaries by a computer program. Summarization process involves interpretation, transformation and generation. It is very difficult for human beings to manually summarize large documents of text. There is an abundance of text material available on the internet. However, usually the Internet provides more information than is needed. Therefore, a twofold problem is encountered: searching for relevant documents through an overwhelming number of documents available, and absorbing a large quantity of relevant information. The goal of automatic text summarization is condensing the source text into a shorter version preserving its information content and overall meaning.

##### E. *Sentiment Calculation*

Sentiment Calculation is a Natural Language Processing and Information Extraction task that aims to obtain writers feelings expressed in positive or negative comments, questions and requests, by analyzing a large numbers of tweets. Generally speaking, sentiment calculation aims to determine the attitude of a speaker or a writer with respect to some topic or the overall tonality of a document. In recent years, the exponential



inspiration to improve my work. I thank all the people for their help directly and indirectly to complete my project.

Lastly, I would like to thank each and every person who directly or indirectly helped me in the completion of the project especially my Parents and Peers who supported me throughout my project.

#### REFERENCES

- [1] <http://programminghistorian.org/lessons/output-data-as-html-file>  
Accessed on 25/4/2021
- [2] [2] <https://marcobonzanini.com/2015/03/02/mining-twitter-data-with-python-part-1/> Accessed on 04/3/2021
- [3] [3] <http://hadoopindepth.blogspot.in/> Accessed on 23/4/2021
- [4] <http://glowingpython.blogspot.in/2014/09/text-summarization-with-nltk.html> Accessed on 24/05/2021