

**LEBANESE AMERICAN UNIVERSITY**

An Optimized Influencer Rating Model using a Joint Event and  
Theme based Approach

By

Ali Srour

A Thesis

Submitted in partial fulfillment of the requirements  
for the degree of Master of Science in Computer Science

School of Arts and Sciences

December 2020

© 2020

Ali Srour

All Rights Reserved

## THESIS APPROVAL FORM

Student Name: Ali Srour I.D. #: 201500067

Thesis Title: An Optimized Influencer Rating Model Using a Joint Event and Theme Based Approach

Program: Masters of Science in Computer Science

Department: Computer Science and Mathematics

School: Arts and Sciences

The undersigned certify that they have examined the final electronic copy of this thesis and approved it in Partial Fulfillment of the requirements for the degree of:

Masters of Science in the major of Computer Science

Thesis Advisor's Name: Azzam Mourad

Signature:  Date: 28 / 12 / 2020  
Day Month Year

Committee Member's Name: Ramzi Haraty

Signature:  Date: 28 / 12 / 2020  
Day Month Year

Committee Member's Name: Abdul-Nasser Kassar

Signature:  Date: 28 / 12 / 2020  
Day Month Year

## THESIS COPYRIGHT RELEASE FORM

### LEBANESE AMERICAN UNIVERSITY NON-EXCLUSIVE DISTRIBUTION LICENSE

By signing and submitting this license, you (the author(s) or copyright owner) grants the Lebanese American University (LAU) the non-exclusive right to reproduce, translate (as defined below), and/or distribute your submission (including the abstract) worldwide in print and electronic formats and in any medium, including but not limited to audio or video. You agree that LAU may, without changing the content, translate the submission to any medium or format for the purpose of preservation. You also agree that LAU may keep more than one copy of this submission for purposes of security, backup and preservation. You represent that the submission is your original work, and that you have the right to grant the rights contained in this license. You also represent that your submission does not, to the best of your knowledge, infringe upon anyone's copyright. If the submission contains material for which you do not hold copyright, you represent that you have obtained the unrestricted permission of the copyright owner to grant LAU the rights required by this license, and that such third-party owned material is clearly identified and acknowledged within the text or content of the submission. IF THE SUBMISSION IS BASED UPON WORK THAT HAS BEEN SPONSORED OR SUPPORTED BY AN AGENCY OR ORGANIZATION OTHER THAN LAU, YOU REPRESENT THAT YOU HAVE FULFILLED ANY RIGHT OF REVIEW OR OTHER OBLIGATIONS REQUIRED BY SUCH CONTRACT OR AGREEMENT. LAU will clearly identify your name(s) as the author(s) or owner(s) of the submission, and will not make any alteration, other than as allowed by this license, to your submission.

Name: Ali Srour

Signature: 

Date: 28 / 12 / 2020

Day

Month

Year

## PLAGIARISM POLICY COMPLIANCE STATEMENT

I certify that:

1. I have read and understood LAU's Plagiarism Policy.
2. I understand that failure to comply with this Policy can lead to academic and disciplinary actions against me.
3. This work is substantially my own, and to the extent that any part of this work is not my own I have indicated that by acknowledging its sources.

Name: Ali Srour

Signature: 

Date: 28 / 12 / 2020  
Day Month Year

## **ACKNOWLEDGMENT**

First of all, I would like to express my sincere gratitude to my thesis supervisor, Dr.Azzam Mourad for his great support and advice, and jury members Dr.Ramzi Harati and Dr.Abdul-Nasser Kassar for their followup and guidance. I would also like to thank my wife and my little son for their inspirational support and love! Not to forget to thank my family, my colleagues, my friends, and anyone who believes in my goals.

# An Optimized Influencer Rating Model using a Joint Event and Theme based Approach

ALI SROUR

## ABSTRACT

The continuous development of social media platforms and exponential level of users' engagement are playing a key role in turning social media platforms into a data source that is indispensable for understanding the behavior of people. Yet this comes along with all types of challenges that needs efficient solutions and big data analysis techniques to allow capturing the different dimensions evolving over social media. Moreover, the importance of detecting influencers over social media platforms has become one of the most challenging research topics given that influencers are at the core of decision-making strategies and leading events' directions on Social Media. Generally, determining influencers can increase revenue and utility; however, measuring the influence of users is essential for determining influencers. Influencers should be credible, reliable, trustworthy, knowledgeable in the domain being discussed and have a high impact that derives people's opinions and lead them towards the proper decisions. However, influencers might play different roles when speaking about misinformation and conspiracy during sensitive and trending event. While different techniques were developed to select influencers over social networks, identifying influencers remains an evolving topic due to the dynamic nature of the social media users and use in addition to their extremely increasing growth which attracts high research interests across multiple disciplines including data science, psychology, sociology, and different human sciences all of which have been studying the topic from multiple angles. In this thesis, we further aim at identifying influencers through proposing influence rates calculation mechanism to find real and highly influential users at a certain event

over Twitter using a mixed theme and event base approach with an emphasis of maximizing accuracy calculation by integrating historical influence rates in addition to content and profiles reputation. We further apply our approach on a global pandemic, the novel Coronavirus, and then we provide results and performance analysis. Finally, we conclude our work by summarizing our major findings and discussing future work.

### **Keywords**

Theme-Event Approach, Social Network Analysis, Influence Rating, User Credibility, User Impact, User Reputation, Natural Language Processing, Big Data, Data Science, Cloud Computing, Social Media, COVID-19, Infodemic.



# TABLE OF CONTENTS

<b>1. Introduction</b>	<b>1</b>
1.1. Motivation and Problem Statement . . . . .	4
1.2. Methodology and Contributions . . . . .	4
1.3. Thesis Organization . . . . .	6
<b>2. Background and Literature Review</b>	<b>7</b>
2.1. Background . . . . .	7
2.1.1. Twitter API . . . . .	7
2.1.2. Python Scripting . . . . .	8
2.1.3. Google Cloud Platform . . . . .	8
2.1.4. Data Analysis Techniques . . . . .	10
2.1.5. Calculation Metrics . . . . .	13
2.1.6. User Social Metrics . . . . .	14
2.2. Literature Review . . . . .	16
2.2.1. Spam and Misleading Posts Detection . . . . .	16
2.2.2. User-Centered and Content-Based Reputation and Credibility Analysis . . . . .	18
2.2.3. Influence Ranking in Social Networks . . . . .	18
2.2.4. Impact of Social Networks . . . . .	20
<b>3. A Proposed Influencer Rating Model over Twitter Social Media Platform</b>	<b>22</b>
3.1. Introduction . . . . .	22
3.2. Approach Overview . . . . .	26

3.3. Analysis Roadmap . . . . .	28
3.4. Joint Theme and Event Model for Maximizing the Influence Ratio Prediction . . . . .	30
3.4.1. Accuracy Maximization . . . . .	31
3.4.2. Reputation Maximization . . . . .	36
3.5. Conclusion . . . . .	38
<b>4. Implementation and Experimental Results</b>	<b>39</b>
4.1. Implementation . . . . .	39
4.1.1. Topic Selection and Methodology Overview . . . . .	39
4.1.2. Data Processing . . . . .	40
4.2. Results and Analysis . . . . .	42
4.2.1. Event Based Analysis . . . . .	43
4.2.2. Theme Based Analysis . . . . .	44
4.2.3. Joint Event and Theme Based Analysis . . . . .	45
4.2.4. Reputation Based Analysis . . . . .	46
4.3. Discussion . . . . .	48
<b>5. Conclusion and Future Work</b>	<b>51</b>
<b>References</b>	<b>52</b>

## LIST OF FIGURES

<b>Figure</b>	<b>Page</b>
1. Example of a Social Network [1] . . . . .	11
2. Example of an Ontology [2] . . . . .	12
3. A flowchart that shows the steps of our proposed influence rating and identification methodology . . . . .	27
4. Topic Selection, Data Collection and Content Labeling Methodol- ogy Flowchart . . . . .	28
5. Venn diagram of the multiple measures that the UIR consists of: User Meta, User Presence, and Tweets Meta . . . . .	32
6. Event Influencers using INR Proposed Calculations . . . . .	43
7. Top 1000 Influencers . . . . .	43
8. Selecting Influencers using the General Approach . . . . .	44
9. Top 1000 General Influencers . . . . .	45
10. List of Influencers with Joint Approach . . . . .	46
11. Credibility Analysis and Aggregations for Selected Influencers of Joint Approach . . . . .	46
12. Comparison between list of selected influencers with influencer gen- eral reputation and Influencer Reputation Calculation . . . . .	47
13. List of Influencers with Reputation Maximization . . . . .	47
14. Multiple bar chart showing Influencers with Reputation Maximiza- tion . . . . .	48

15. Credibility Analysis and Aggregations for List of Influencers with Maximized Reputation . . . . .	48
--	----

## LIST OF ABBREVIATIONS

ABBREVIATION	MEANING
ACC	Accuracy
API	Application Programming Interface
FFr	Followers to Followees Ratio
GCP	Google Cloud Platform
GIR	General Influence Rate
INGIMPr	Influencer General Impact Rate
INGREPr	Influencer General Reputation Rate
INGtCRDr	Influencer General Tweets Credibility Rate
INIMPr	Influencer Impact Rate
INpCRDr	Influencer Profile Credibility Rate
INR	Influence Rate
INREPr	Influencer Reputation Rate
IRR	Irrelevance
NLP	Natural Language Processing
NLU	Natural Language Understanding
REP	Reputation
REPR	Reputation Rate
Tcr	Tweets Count Rate
UER	User Engagement Rate
UMR	User Meta Rate
UPER	User Presence Engagment Rate
UTMR	User Tweets Meta Rate

# Chapter One

## Introduction

Social media platforms have been one of the most prominent ways for connecting people to each other, offering a medium over which information pertaining to individual experiences and stories can be shared with the public. Further, social media platforms are able to enhance the quality of lives of people as they offer easy accessibility to information and provide social support. Users with different backgrounds and ages can discuss trending topics, thoughts and initiatives and hence, social media platforms are no longer only a hub for social interaction but also a hub for sharing health information, market news and other type of information. Further, social media platforms are distinguished with various features and components which make them tailored for different purposes. For example, Facebook is used for non-formal and social communication while Twitter is typically used for business and professional goals. This makes research firms and organizations focus on developing different strategies for every platform to to understand how users interact over it and to assess the reliability of information being disseminated. Further, generated data by users over social media platforms evolve in realtime being noisy, unstructured, and dynamic. Having such characteristics impose serious challenges in front of data analysis techniques requiring solutions with high accuracy and adaptability. Developed algorithms should be able to remove redundancy and noise out of data being collected to have more accurate results and should be customized to fit various problems.

Moreover, news creation and consumption has been changing since the advent of social media. An estimated 2.95 billion people in 2019 used social media worldwide. Most platforms are used to transmit relevant news, guidelines and precautions to people. However, according to WHO, in the context of epidemics, uncontrolled conspiracy theories and propaganda are spread faster than the pandemic events themselves, creating an infodemic and thus causing psychological panic, misleading medical advises, and economic disruption. Accordingly, determining trusted sources is very essential to determine valid information. While medical and research centers can provide major insights about the flow of a certain infectious disease, the opinions and discussions of users and particularly influencers in pandemics are very essential to observe as it helps in determining how users are being sociologically and psychologically affected during the occurrence of this event. Infectious diseases impose a serious health threat and affect societies at a global scale. However, capturing this public unfiltered opinion in the ocean of data available on social media presents a key challenge through which fine-grained filtering and proper extraction should be made to detect influencers. For this reason, in this thesis, we also present a case study that assess the impact of credibility of users on Twitter during a major infectious disease outbreaks COVID-19.

Further, influencers over social platforms are considered an important group that can drive users' opinion and convince them with certain topics. Predicting election candidates, finding most influencers in a certain crisis, and predicting event's user influence ratios are very critical topics discussed nowadays. Identifying and predicting influencers in social networks have many applications including viral marketing [3], searching [4], and expert recommendation [5]. The ability of influencers to spread information has been well studied on Twitter where influence was used for different goals such as human mobility [6], rumor spreading [7] and epidemiology [8] among others. Twitter has gained high popularity, and recently, it has witnessed a spike in its number of users. Starting with eight million users

in 2009, this number has increased to reach more than 300 million in 2019 [9]. Twitter is characterized with its short messages limit and thus people tend to focus on the reflecting their ideas as concisely as possible. This lead research organizations and digital agencies, with the current technological advancements, to study Twitter data due to the various insights that can be derived including social and psychological behavior of people. The authors in [10] define three main types of users on Twitter: information seeker, friends and information source. While information seekers tend to search for knowledge and are typically less active on Twitter, an information source has continuous activity and usually has a huge number of followers. Business professionals and organizations use twitter to market their ideas and products and increase their profits through attracting the highest possible number of people interested. Being a free marketing service, all organizations and agencies aim at devising innovative strategies to attract user's attention to their content and products. The problem of identification of influencers and ranking them on Twitter is still under research and development. The quantitative assessment of influencers is considered a key challenge whereby different solutions and approaches are being proposed. For instance, predicting influencers in micro-blogging e-crimes requires measurements that accurately capture this influence and does not tolerate errors. The definition of influence varies from one context to another; for instance, the influence of a user can be defined as the user's capacity to influence others with his/her opinion, and it also can be defined as the ability of a user to spread a certain message to others that he/she typically interacts with them [11]. While detecting influencers in a social network can be approached at a global or a general level, in this thesis, we focus on providing a solution that discovers influencers in particular events with the objective of jointly maximizing accuracy and minimizing irrelevance. Our suggested model combines multiple influence ratio calculations to achieve this objective. We define an event as an incident that occurs in a certain time frame and specified location that lead users to tweet and initiate various hashtags and



trends while the event is ongoing. We further propose a second step to verify our main approach through which we aim at maximizing reputation of our selected list of influencers.

## **1.1 Motivation and Problem Statement**

The recent technological advancements (e.g cloud computing, machine learning, natural language processing etc. ) are extending the capabilities of social media platforms to support massive number of users while offering smart services and accordingly, turning social media platforms to an attractive source of data with variety of topics generated by users with different profiles and ages. Analyzing such data can present a key enabler for driving essential information about users interests and thus generally provisioning goals that achieve higher utility. However, this comes with all types of challenges as deriving significant information within a certain context necessitates careful examination for data collection, processing and developing solid methodology that critically assess the data available and addresses the given problem with the maximum possible accuracy. For this reason, after finding that the accuracy of identifying the real event influencers will help in re-orienting and managing the spread of a possible misleading content and misinformation different countermeasures by analyzing the impact of social networks in global event like COVID-19, we extend on this study to provide an influence based approach for influencers identification and rating over Twitter events.

## **1.2 Methodology and Contributions**

Our suggested approach entails selection of influencers and assessing their level of influence based on two main contexts: theme and event. We aim at capturing all event influencers and at the same time all theme influencers (for similar time, location, and sample size criteria) and then correlating and analyzing his-

torical influence and other attributes for our selected event influencers to obtain maximum accuracy. The first step in our approach aims at calculating influencers ratings and then selecting those influencers in the context of an event while jointly maximizing accuracy and minimizing irrelevance by joining similar calculations and findings about the selected theme that the event belongs to. In the second step, we aim at selecting, out of the resulting selected list of influencers in the first step, the set of influencers while maximizing their reputations by analysing their content and profiles credibility and impact. Thus we measure the influence of each user from multiple angles and using a combination of influence ratios to achieve the aforementioned objectives.

We summarize our main contributions as follows:

- Using twenty two different data points and 12 calculated ratios to calculate the Event's Influence Rate and calculate users' historical influence rates to maximize the influence rates calculations of a certain event on Twitter.
- Using Lexicon-based approach to measure both content and profile credibility of a selected list of influencers based on their profiles and their tweeting activities.
- Optimizing the accuracy of influence rating through a Joint Theme and Event based Model and then disregard and eliminate Event based irrelevant users using content analysis. Further, we combine popularity and reputation based influence ratio through a Joint Theme and Event base Model
- Elaborating Both computing and non-computing findings, implications, social networks management strategies and research directions supported with thorough literature review for a field to become of great importance in the near future

## 1.3 Thesis Organization

This thesis consists of 5 chapters. In chapter 2, we introduce background information important for understanding the rest of our work and we illustrate the different proposed approaches and a wide literature review about both impact of social networks and event influence rating on Twitter. In chapter 3, we discuss our proposed approach and its corresponding components in addition to illustrating the mathematical models and the optimization problems. Chapter 4 explains the detailed implementation steps and analysis results, in addition to discussing the experimental results and proving the approach accuracy. And finally, in chapter 5 we conclude the whole work and we explain some of the possible future enhancement and directions. Portions of this thesis have been published in [12].

# Chapter Two

## Background and Literature Review

### 2.1 Background

#### 2.1.1 Twitter API

Twitter provides two publicly available types of Application Programming Interfaces (API) for developers that allows access to collect data from the social network. In this section, we describe each available type as follows:

##### **API REST**

API REST allows developers to use simple HTTP primitives (e.g. GET, POST) to read and write Twitter data [13]. For example, we can using GET issue a query and obtain the response to this query over Twitter in JSON format. However, API REST limits the number of times each request done within a certain range of time using working windows and tokens. Thus a developer can initiate requests according to the amount of tokens allowed for this type of request and wait until the second working window starts to start again with the same amount of tokens limit to initiate requests. Hence, representing a full social network can take a long time and this presents an important limitation of the API REST.

## Streaming API

The Streaming API offers developers accessibility to the overall flow of tweets and events in real time (e.g. deletions, replies) [14]. The Streaming API allows developers to establish a connection that supports continuous data streaming along with multiple filters that customize collections to the developer's preferences. However, since 2012, new updates on API allowed restricted access to data, offering other amount of data for a certain profit.

### 2.1.2 Python Scripting

Python scripting involves developing a collection of commands to be executed. Unlike usual written programs, this script is typically contains a set of functions and modules which is run using command line or from within a Python interactive shell to perform a specific task [15].

### 2.1.3 Google Cloud Platform

Google Cloud Platform (GCP) <sup>1</sup> is a cloud solution developed by google that offers different products for developers to build various types of programs including websites, distributed applications and others [16]. GCP provides these services on the same internal infrastructure used for hosting Google main end-user solutions such as Google Search. Being reliable, different companies has been hosting their development solutions on GCP including Spotify, Airbus and The New York Times among others.

GCP offers a set of physical assets (e.g., processors, hard disks) and virtual resources (e.g., virtual machines, kubernetes instances) which are hosted in Google's data centers in multiple regions and zones. Further, Google offers through GCP more than 90 products categorized into nine groups:

- Compute: which offers multiple technologies for execution and computation

---

<sup>1</sup><https://cloud.google.com/>

intensive tasks.

- Storage and Databases: which provides storage resources with customized sizes ranging from data for small businesses to large business
- Networking: which supports connecting virtual instances to each other over cloud using different technologies including Virtual Private Cloud (VPC) network and cloud balancing tools.
- Big Data: which offers data storage, processing, and analytics on a more scalable, flexible, cost-effective which are essential when data volumes grows exponentially.
- Cloud AI: which aims at providing a code-based data science development environment that allows data scientists work the full project from start to end over cloud.
- Management Tools: which focuses on simplifying the work over the cloud through allowing developers to specify all needed resources in a declarative format supporting mobile applications that enables close monitoring of work and cost incurred by developed work.
- Identity and Security: which offers a platform for securely managing identity, access, application and endpoints over cloud.
- IoT: which provides a set of utilities that enable connecting, processing, storing and analyzing data in the edge and cloud offering an integrated software stack with machine learning capabilities and scalable cloud services.
- API Platform: which offers a wide range of APIs pertaining to storage, machine-learning API and other tools that can be added to every project when needed.

## 2.1.4 Data Analysis Techniques

In this section, we provide key concepts about some important data analysis techniques and methodologies that are used through our work:

### Content Analysis

Content Analysis is a methodology that aims at compressing themes, paragraphs, and semantics into fewer content categories based on some rules and identify and eventually identify findings. Scientists analyze the meanings of these words using frequencies and relationships to derive better inferences about them. Content analysis can be used for multiple goals: 1) identify intentions and trends of a certain group of users 2) determine emotional and psychological state of users 3) analyze patterns in the communicated content and 4) understand the responses of users towards certain topics.

One of the main approaches to apply content analysis is to analyze most frequent words and keyword frequencies and derive the proper conclusions based on this type of information [17]. Generally, accurate approaches exclude stop words, word variations and other types of words that increases the margin of error and confuse the proposed models. Further, Lexicons with the same meaning and synonyms terms are mapped to the same semantic concept in order to depict the actual content meaning with higher accuracy.

### Social Network Analysis

Social Network Analysis (SNA) is defined as applying different network theories to study social relationships over social networks. Users represent nodes in the network and they are connected with edges that resemble relationships (e.g friends). Different SNA approaches are being used to study the interactions between nodes and examining the relationships between these entities. SNA has been used for multiple research areas including: influence propagation, expert finding, link prediction, community or group detection, recommender systems,

predicting trust and distrust between users, behavior analysis and opinion mining [18]. Figure 1 depicts an example on how social network analysis maps and measures relationships flows between people, groups, or organizations. This can allow better understanding of communities and identify the major characteristics for each group. Different metrics are used to assess these relationships among people including: degree centrality, betweenness centrality, closeness centrality, network centralization, network reach and network integration.



Figure 1: Example of a Social Network [1]

## Ontology

Ontologies are defined as a set of relationships and descriptions associated with a certain field. An ontology depicts the taxonomy related to a theme or topic. Building ontologies require extracting the proper terminology related to a certain domain, depicting concepts and learning the non-taxonomic relations and rules. An ontology defines the knowledge about a certain domain by presenting a set of vocabulary and terms that are usually shared in the context of this domain. An ontology is built of two main components i) a set of classes that defines ideas about the selected domain and ii) a set of relationships between these classes.



There are four key components that comprise an ontology: classes, properties, restrictions, and individuals [19].

- A class depicts a concept in a certain domain. For instance, both Pizza and Toppings are example on classes in the figure below 2. As the figure shows, each class can have subclasses that usually belong to this class but have additional information and attributes. This relationship is represented using a directed edge between class and sub-class nodes. For instance, MushroomPizza is a sub-class that belongs to class Pizza.
- The property defines the relationship between two given classes [2]. For example, Pizza and Topping classes are related through hasTopping which illustrates the type of relationships between them.
- Since a class may have some constraints that are needed, a restriction specify the set of rules imposed by a certain class. The example in fig. 2 shows restrictions imposed by Mashroom that should be satisfied to be done: 1) it should hasTopping[some] Mushroom and 2) hasTopping[some] Mozzarella.
- An individual is an instance of a given class. Usually, individuals belonging to a certain class would have the same properties and restrictions.

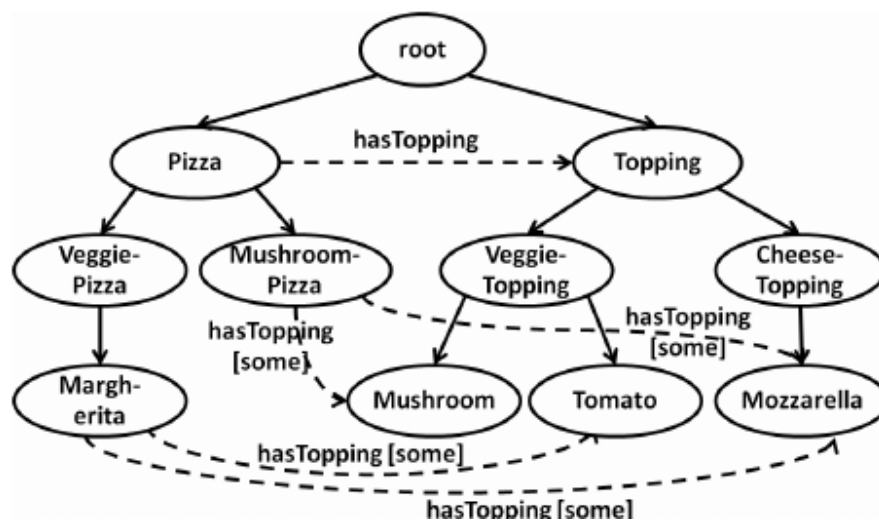


Figure 2: Example of an Ontology [2]

## **Sentiment Analysis**

Sentiment Analysis or opinion mining is a branch of Natural Language Processing (NLP) which focuses on analyzing user's text to predict its sentiment. Sentiment analysis over social media is being a highly demanded task. For example, companies are interested in determining people's reactions over social media regarding their products and accordingly, they tend to customize their products based on users' sentiment.

## **Natural Language Processing (NLP)**

Natural Language Processing (NLP), or computational linguistics, aims at using computational models to help computers understand human languages. Key applications that uses NLP solutions to solve problems includes translation of text from one language to another, automatic replies on asked questions, and spelling correction and grammar checking. NLP combines machine learning and linguistics to solve such problems.

NLP is data-driven and thus the data chosen can greatly challenge the effectiveness of NLP methodologies. Further, building the datasets incorporates defining the problem to be solved and the appropriate measurements to be taken. Another important challenge that NLP faces is developing models that can understand words in variable situations as each word, despite having a unique meaning, differ in its meaning from one context to another making it ambiguous to understand at different levels. [20]

### **2.1.5 Calculation Metrics**

In this section, we will define the term "Accuracy" which is the main metrics of our work.

## **Accuracy**

The term accuracy in social media influence rating is related to the result of a mathematical or arithmetic calculation for users influence rates. We can describe the approach that is used to measure and calculate users influence rates as accurate or non-accurate, and we can also give an accuracy threshold or percentage based on the result's correctness. The lower the result's error margin, the higher the accuracy of the used approach. Moreover, to maximize the accuracy of finding influencers in a social media event, that means we need to find the actual event influencers based on a different approach and mathematical model that can reduce the gap between the user's calculated influence rates and their actual influence. Finally, influence rating accuracy can be clearly defined as the level of match between a user's given influence rate and his actual influence value.

### **2.1.6 User Social Metrics**

A user is considered a key driver in social media platforms and hence, studying user's behavior entails a set of metrics that define this behavior. In this section, we describe important social related metrics to users as follows:

#### **Reputation**

Reputation can be defined as an attitude constructed on two main components: an emotional component and a rational component [21]. Reputation is used for different purposes on Twitter, such as political activities, human mobility and epidemiology, among others. Reputation management has been helping in understanding reputation on social media as they capture feedback of users analyzing multiple metrics [22].

#### **Influence**

Influence can be generally defined as the effect induced by a certain person on the ideas, thoughts or behavior of other people [23]. In the context of social media

platforms, various definitions for influence has been defined and measured in research [24]. In this thesis, influence as the effect induced by a user on other users that cause the propagation of ideas, thoughts and behavior over societies [25]. Katz et al. [26] explain that influencers are able to produce, using word-of-mouth, a chain-reaction of influence resulting in high reach.

## **Credibility**

Credibility can be defined as trustworthiness and inherent persuasiveness. Utilizing credible sources and information over social media is essential for deriving accurate conclusions. However, evaluating credibility is a challenging problem sometimes users are not well known and thus no guarantees or responsibilities about the content. Twitter provides the username as the only information about this source of information and thus this profile may be a fake profile generating false information. Other accounts are verified accounts and refer to legitimate source that is authoring the account's tweets <sup>2</sup>. However, this is only a small group of users whose accounts are verified. Hence, measuring the credibility of a certain social media user is essential to depict the credibility of the given piece of information. Different levels of credibility are defined and various research interests are being explored to measure credibility for every level:

- **Post Credibility:** which defines whether a certain post represent relevant and accurate information about a specific topic [27].
- **User Credibility:** which quantifies the reliability of a certain user and is typically associated with a certain score [28].
- **Topic Credibility:** which corresponds to the acceptance of a certain topic or event [29].

---

<sup>2</sup><https://twitter.com/verified>

## 2.2 Literature Review

In this section, we provide a literature review in relation to influence maximization and highlight the different approaches proposed. In addition to a literature review about the different techniques used to measure the impact of social network in public event on Twitter.

### 2.2.1 Spam and Misleading Posts Detection

To detect spam users in social networks, it is required to first analyze message content [30–35]. However, many approaches enhance their calculations by further analysing user profiles and connections [36]. In [37], the authors detected spam users using a semi-supervised method by investigating tweets’ contents and classifying maliciousness levels. While, The authors of [38] used URLs and domain names specified in the user profiles in order to classify and measure their potential maliciousness. Similarly, Lee et al. [39] proposed a real-time approach to detect redundant URLs. In addition, Guille et al. [40] in his approach, took the advantage of URLs being reused to spot any malicious intent. Amleshwaram et al. [41] distinguished spam users using a supervised model by measuring hyperlinks and other URLs features. Moreover, authors of [42] proposed a classification approach to classify users between promoters, spammer and legitimated from their published content (especially videos). Using supervised machine learning on top of manually selected and studied users, some authors were able to classify spam and malicious users. SHEN et al. [43] in his approach, tried to classify users based on their tweets content by analysing their behaviors. However, for a large datasets and thus potential large error margins, it is important to say that manual work is a must when it comes to feeding all the mentioned supervised and semi-supervised learning algorithms.

In [44], authors were able to select influencers based on some extracted features from their activity for a period of six months including but not limited to: active days and times and actual retweet delay in addition to other tweets

and user profile meta data like followers and followees. The mentioned authors applied four different algorithms (Bagging, Random Forest, Extra Trees, and Voting) in order to analyze and classify tweets and users. While, authors of [45] propose a method to predict election results using Twitter through extracting sentimental information and content sentiment analysis to predict the popularity of candidates. The authors then propose another approach in [46] to predict election results through calculating their popularity using content sentiment analysis based on tweets content terms weighting approach. In [47], the authors aim at finding tweets popularity based on time warping and sequence clustering algorithm using retweets and replies. Tweets are divided into time-based groups and then the popularity of each tweet post is calculated using sequential clustering. The centroids of each cluster are depicted using the Barycenter for every cluster. These centroids are utilized to generate popularity profile templates and then those templates are used in assessing new tweets. The authors in [48] explore the most important features in tweets that relates to financial popularity and develop a prediction model was created using binary logistic regression. They calculate popularity ratio for each tweet based on a descriptive analysis list for all tweets like (number of retweeted tweets over all tweets, number of liked tweets over all tweets, and so on) and calculate mean values for popular and non-popular tweets. In [49], the authors propose a sentiment analysis based approach to calculate the reputation of tweets and users using collected Saudi tweets in Twitter. The authors propose an Arabic Dataset Saudi dialect lexicon to detect the polarity of tweets and then use beta probability density functions to derive reputation scores. The authors in [50] propose a ranking mechanism based on social user communications and activity.

## 2.2.2 User-Centered and Content-Based Reputation and Credibility Analysis

In addition to finding spam and malicious users, many authors work on sorting and ranking users based on their influence ratios by using different techniques and approaches based on two major categories: focusing on tweets content in order to calculate and allocate reputation level using ML techniques [51], [52], [53], while the other set of approaches focused on the network analysis of each users' graph network to assess their level of influence [54], [55], [52]. But at the same time, some authors provided a mix between the two mentioned categories or sets of approach targeting higher accuracy. Following are some highlights from the mentioned approaches. Jain et al. [56] used a kind of graph theory algorithm in order to score users based on their centrality in order to classify and identify universal leaders opinions in later stages. But in terms of social distancin, authors of [57] provided an important analysis on how this can affect trust and trustworthiness among users. Mohammadinejad et al. [58] in his article presented a framework derives the users' influence levels from their analysed personalities. The authors of. [59] identified user behaviors based on their frequent activities, and then scoring their credibility. However, Wang et al. [60] correlated tweets credibility to user profiles. In addition to that, Tsikerdekis et al. [61] highlighted more important user habits or actions like creating multiple accounts. Also, Ahmad et al. [62] article presented a survey on different approaches used for detecting social networks rumors. Buzz et al. [63] in his article, used sentiment analysis for Arabic language tweets in order to classify content and users based on their sentiment score. While, Alrubaian et al. [64] deduced user credibility in order to measure the level of fake and malicious news spreading.

## 2.2.3 Influence Ranking in Social Networks

Influence detection in events and public conversations has become one of the most challenging research direction in the field of Social Network Analysis especially in

Twitter event. But finding influencers and ranking users based on their influence ratio can be done using different methods and techniques. Authors in [65], [66], [67–70] found that user meta data like number of followers, number of tweets, and number of followees in addition to tweets meta data like retweet and favorite counts are enough to find users influence ratios. On the other hand, authors in [71] rank user influence based on their actual relationships. While [72] links users influence with their social activity during the selected event. The authors of [73] calculated a potential social networking ratio (SNP) for each user in among a number of most popular Twitter accounts in Austria using their accounts meta data like followers and followees counts. While Bakshy, Eytan et. al [74] used diffusion trees techniques in order to calculate the influence ratios of users whom their tweets has URLs by calculating the reach of those URLs in other platforms. In [75], M. Anjaria and R. M. R Guddeti used Incremental Learning algorithms with NLTK sentiment analysis to predict the presidential elections in the US. While other authors such as in [76] classified and categorized users into four influence levels categories based on their interactions and activities. In [77], Y. Mei, Y. Zhong and J. Yand in order to calculate the popularity for each user, they used eight data points in addition to NewFollowers and NewMentions features to find the top hundred users in Australia. Riquelmea et al. [51] targeted the propagation ratios of users' tweets using two different linear centrality approaches. Similarly, Li et al. [78] proposed an eigenvector centrality based approach to measure users' influence rate as well. Lahuerta-Otero et al. [79] applied behavioral analysis techniques among special kind of twitter users in which those techniques can increase user's influence ratios. However, Sharma et.al [80] proposed a novel approach that combines users tweets and their trend scores in order to elect potential influence ratios. Moreover, Huynh et al [81] focused on analysing tweets tags and their correlation with the speed of their propagation.



## 2.2.4 Impact of Social Networks

The authors in [82] highlight that existing semantic textual similarity (STS) datasets and models fail to provide good performance results to specific domains such as COVID-19. They introduce CORD19STS dataset which includes 13,710 annotated sentence pairs collected from COVID-19 open research dataset challenge and then generate one million sentence pairs out of them. Further, the authors use a finetuned BERT-like model, Sen-SCI-CORD19-BERT, to create a balanced dataset with a total of 32K sentence pairs. The authors in [83] propose COVID-Twitter-BERT (CT-BERT), a transformer-based model, pretrained on a large corpus of Twitter messages on the topic of COVID-19. The proposed model is compared with other transformer-based models showing better performance and can be used for COVID related natural language processing tasks, including classification and question-answering. The authors in [84] study the diffusion of COVID-19 related information using massive data analysis on Twitter, Instagram, YouTube, Reddit and Gab. They analyze engagement and interest of users in COVID-19 topic and then differentially assess the progress of the discourse on a global scale for each platform and their users. The authors trace different volumes of misinformation in each platform and then provide platform-based results about rumors' amplification. In [85], the author study 43.3M English tweets about COVID-19 and conclude that humans were mostly focusing on public health concerns while bots were used to promote political conspiracies in the United States. The authors in [86] estimate low credibility information on Twitter during the outbreak, and the role of bots in spreading this type of information based on content analysis. Similar to [85], the authors conclude the role of social bots posting and amplifying low credibility information while the majority of generated data is by humans. The authors in [87] aim to explore evolving discourse about COVID-19 pandemic using Natural Language Processing, Text Mining, and Network Analysis. The authors identify most common responses to the pandemic and the evolution of information and misinformation

over Twitter and then analyze how each COVID-19 topic was changing with time. A dataset of tweets from all over the world is collected in multiple languages at the time when the total number of COVID-19 cases reported was below 600 worldwide. The authors in [88] curate a set of 1000 day-scale time series of 1-grams across 24 languages on Twitter and derive some conclusions including a comparison to numbers of confirmed deaths due to COVID-19 over time. The authors then host a time series on Github with daily updates with the aim of help researchers in their analysis and investigations. The authors in [89] propose a methodology that aims at understanding conversations about COVID-19 through examining discussed themes related to COVID-19 taking place on Twitter and relating each theme to other information regardless of their quality through shared URL links. The authors depict a spatio-temporal relationship between the flow of information and new revealed cases of COVID-19. In [90], the authors propose Weighted Correlated Influence (WCI) approach which combines the relative impact of timeline-based and trend-specific features of social media users. The proposed approach merges both the profile activity and underlying network topology to calculate the influence score for each user. The authors identify different trends related to COVID-19 over Twitter to generate their results. In [91], the authors identify information influencers during the COVID-19 pandemic and present analysis over a dataset of collected Arabic tweets during the COVID-19 pandemic. The study uses the network topology to identify influencers during the month of March, 2020 and then implement both HITS and PageRank algorithms to analyze and compare the ranking among users. The results show that both HITS and PageRank algorithms have 40% similar influencers. The authors in [92] aim to identify the main topics posted by Twitter users related to the COVID-19 pandemic. The authors identify 12 topics related to COVID-19 composed of four main themes: the origin of the virus; the multiple sources of the virus; the impact of the virus on people, countries, and the economy; and the possible ways to mitigate the risk of being infected.

# Chapter Three

## A Proposed Influencer Rating Model over Twitter Social Media Platform

### 3.1 Introduction

Social media platforms are playing a major role in spreading information, thoughts, and discussions among users. Being supported by various devices with different computing and storage resources, people are able to communicate to each other at anytime and place making it an active and real-time hub for different topics and events. Social media platforms have also attracted users with different backgrounds to talk in any scope of interest; for example, average citizens started to discuss different ideas and thoughts about political decisions [93]. This interconnectedness between users has triggered dynamic changes in societies as people form and share opinions over social media [94]. The age of ideas however varies from one context to another; some ideas appear and then die quickly while others amplify and remain existing. However, understanding the reason behind the adoption of certain ideas necessitates determining the key aspects that affect the behavior of users, one which is influencers in those social networks. Social media influencers have a key impact on the interaction of users and directing their

attention towards some topics. In fact, organizations have been leveraging this group when preparing for any public relations plan [95]. People may change their original opinions according to the opinion of their influencer due to his/her credibility [96]. Influencers may promote a certain brand for an organization's in order to maximize its popularity and revenue, post case studies related to health and medicine or discuss sociological and economical information. However, sometimes, influencers may lead to a great loss as they may pass false information and fake news which consequently decrease their influence. Hence, it is very important to select influencers that may achieve the highest levels of impact and effectively shape communities. However, identifying the most relevant influencers over social media platforms is coupled with various challenges. First, detecting influencers in a social network necessitates a scalable and big data solution that captures the increasing volume of social network interactions generated by the massive amount of users. Second, the variety of themes and events discussed over social media requires data mining techniques that captures the influencers evolving within a certain event. Third, the volatile relations and ratings for each user is also a key challenges as influencers may change with time and thus solutions should consider the dynamic nature of user profiles. Finally, the veracity and reliability of data posted by users is of key importance to consider when detecting influencers. Thus, selecting influencers entails various metrics including but not limited to users tweets, followers, meta data and profile among others. The term influencer has indeed encompassed different meaning based on its definitions. For instance, an influencer is considered as a user with a powerful network, extreme activity and significant impact [97]. Further, other work consider influencers as users with large number of followers while others define them as users with extreme prolific blogs [98]. Hence, there is no agreement on the definition of an influencer. Defining influence remains an emerging topic through which different influence measurements are offered and contributed. Now to realize the importance of influencers on social media with massive number of users, selecting

influencers should be limited to a certain context or event in order to achieve the needed goal for the correct groups of audience. However, a critical aspect of this selection is considering the different underlying metrics that maximizes the relevance of selected influencers. While most influencer rating algorithms and approaches are built on top of well-known metrics such as users metadata (i.e. Number of Followers, Number of Tweets, ..) to calculate their influence ratios, adopting these metrics may fail when selecting influencers in events, as a user can be the most followed one in the country but is not involved in this particular event and hence, considering such metrics that depicts this types of influencers in the overall influence calculation decreases the accuracy and hinders capturing the actual influencers in this event. Further, measuring influence depending only on the number of followers may be misleading as followers can be associated with spammers or fake accounts. While some assumptions may claim that an influential user is the mostly-connected person, different work showed that this claim can not be generalized [99]. For example, the authors in [100] showed that people usually gain influence on certain topics such as sports and politics but not on a global scale. Hence, a thorough and quantitative analysis that considers the network and user specific metrics are needed to detect the most influential people during a certain event. This result leads us to explore more on topical authorities. Our proposed approach considers two different activity dimensions for users and their tweets, in other words, covering tweets and their meta data in the first hand, and users with their networks and other meta data in the other hand. In our work, influence maximization based on reputation is also introduced to increase the accuracy and verify the joint theme and event based models.

With more than 300 million active users per month, Twitter has become an attractive platform for organizations and individuals with strong political, social and economic interest and seek enhancing their popularity and reputation. For this reason, we propose a joint theme and event model with the objective of maximizing the influencers ratio and obtaining the most accurate influencers in

a certain event on Twitter. Our novel methodology is based on multiple factors encapsulated within different influence ratio calculations that aim at defining the influence score for every user. This strategies applied allow us to select top  $k$  most influential users involved in a certain event. The influence ratio calculations are modeled on a joint event and theme basis to select most influential users involved in a certain event. Our main contributions underlay in capturing the productivity and activity of users' under time and place constraints which define the event being under our study to find top influencer. Our proposed model uses 12 calculated ratios to derive the user-event influence rate and and then maps it with users' historical influence rate with the objective of optimizing the accuracy of influence rating. We use various user specific data points and network features and the we apply a joint theme and event based model to derive the corresponding influencers in specified event. The outcomes of this approach are verified and validated using a second layered reputation based approach. First, by calculating potential influence for all Twitter users in the selected event from different stand points (i.e. event based and theme based), we seek to empirically list the set of influencers available jointly at theme level and event levels. The user influence combines different influence scores that assess engagement of user in this event and community through weighted distribution that captures the importance of each ratio. This ratios can be easily calculated and exploits user and tweet levels to find out influence. Further, a case study is conducted where we apply our approach on a sample of users and their corresponding tweets during COVID-19 pandemic and we then provide performance analysis to users.

We list our main contributions in the following:

- The proposed approach integrates user data, tweets, user and tweets meta-data, and network features in order to identify event based influencers enhancing the performance in comparison with other approaches.
- We propose a joint theme and event based approach for influencers with the objective of maximizing accuracy and minimizing irrelevance. We follow

this step with a second step that targets maximizing reputation of the selected list of influencers by measuring both popularity and credibility of both user profiles and tweets within and before the event.

- We provide a case study related to COVID-19 using the proposed approach of ranking.
- Our proposed approach can be further extended to include new features due to its low-level complexity. It also can be applied on any event with its corresponding chosen theme.

The remainder of the work in this chapter is organized as follows. A full approach diagram and flowchart is shown in section 3.2. We draw an analysis roadmap in the section 3.3. In section 3.4, we define the main problem and then discuss the proposed methodology. After that, we write a conclusion for what we have done in this chapter in section 3.5

## 3.2 Approach Overview

After selecting a specific event on Twitter, our approach can be used to find top influencers for the mentioned event. By mixing influence calculations about both the selected event with its corresponding location and timeframe and the topic or the theme that this exact event derived from (Super Category) with the same location and timeframe, our approach can run simultaneously between both levels and in parallel to maximize the accuracy of finding the real and the most influential users for the selected event. In addition, adding a layer of reputation maximization is a major enhancement and contribution that can improve the final rating results using text analysis technique for both credibility and impact measurements and findings. Figure 3 shows the whole approach in one flowchart describing all implementation and analysis steps and headlines. Figure 4 shows a flowchart of the topic selection and data collection mechanism. More details are illustrated in section 3.3.

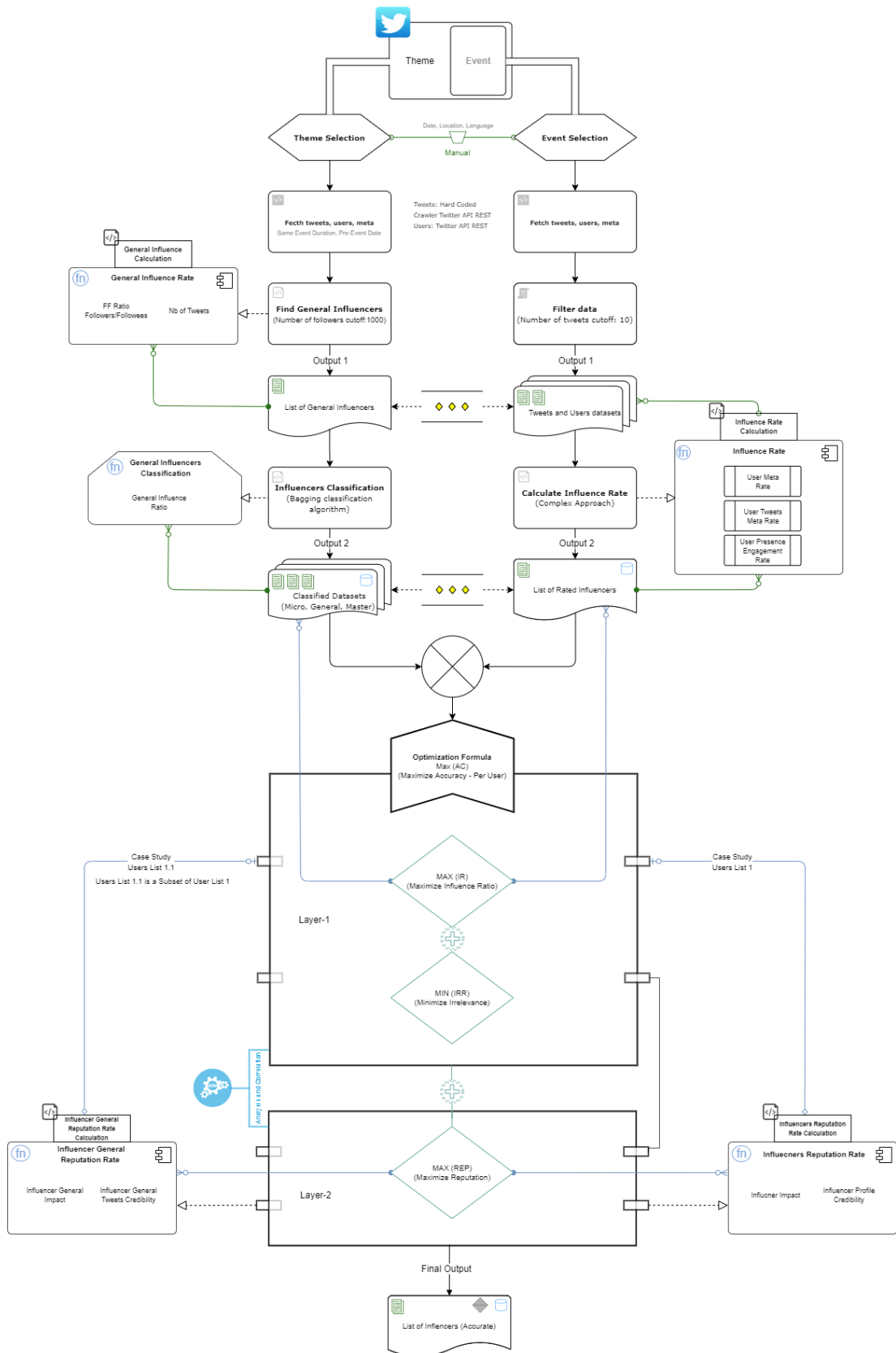


Figure 3: A flowchart that shows the steps of our proposed influence rating and identification methodology



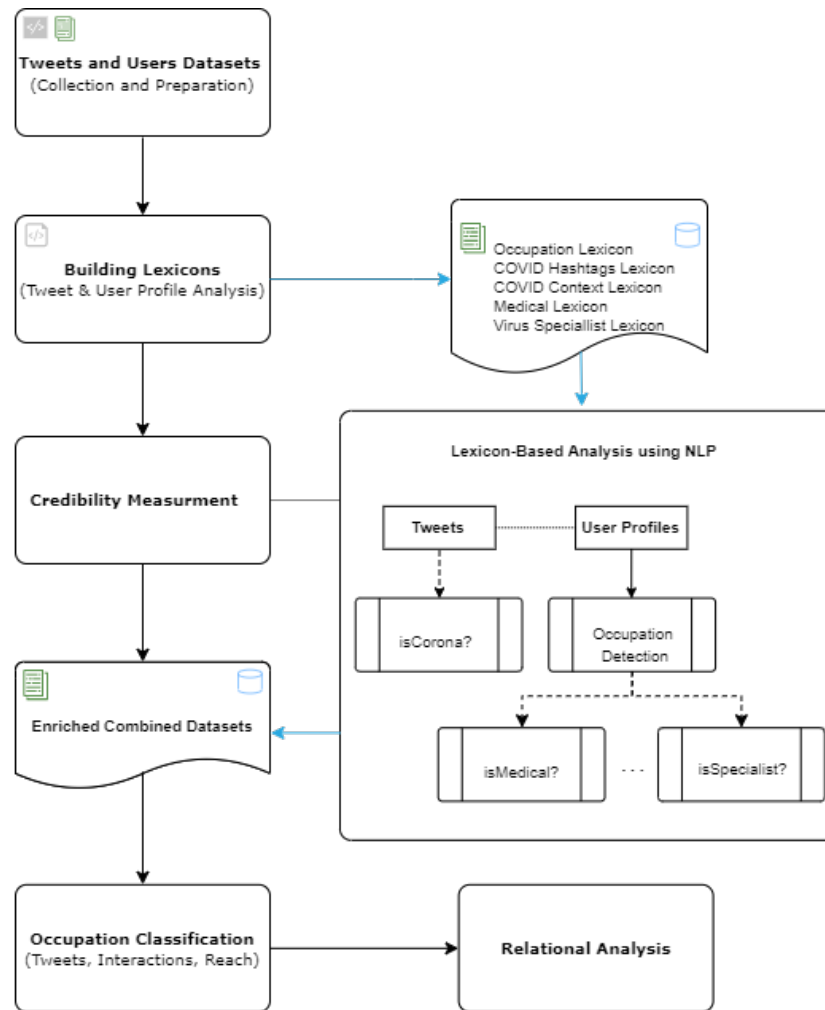


Figure 4: Topic Selection, Data Collection and Content Labeling Methodology Flowchart

### 3.3 Analysis Roadmap

In this section, we will state the analysis roadmap as a list of all approach analysis steps shown in Figure 3 that their results will be illustrated in details in Chapter 4 (Given that the datasets are already collected, cleaned and prepared for analysis for both event and theme):

1. *Engagement Calculation*: Calculate the Engagement Rate of all users of the selected event (the sample that we have previously selected)
2. *Event Influencers Selection*: Select the top 1000 event influencers based on the calculated Engagement rates
3. *Theme Influencers Selection*: Sort and select the top 1000 influencers from

among the Theme users dataset based on their number of followers

4. *Event Influence Rate Calculation:* Calculate and sort Event Influence Rate for the selected 1000 event users
5. *Theme General Influence Rate Calculation:* Calculate and sort Theme General Influence Rate for the selected 1000 theme users
6. *Event/Theme Influencers Classification:* Classify both event and theme influencers lists into Master, General, and Micro influence categories
7. *Event Influence Rating Maximization:* Maximize Influence Rate calculation by combining both Event and Theme results (finding those event influencers who have general influence)
8. *Event Influence Irrelevance Minimization:* Minimize irrelevance by re-arranging users in both event and theme influencers lists who might appear as spam, malicious, or blacklisted ones
9. *Event/Theme Activity Calculation:* Calculate and differentiate activity ratios for both event and theme influencers
10. *Event Influencers Credibility Aggregation:* Find credibility aggregations for event influencers based on their profiles, tweets, and type
11. *Joint Influencers Selection for Reputation:* Select top 100 influencers from the joint result of theme and event influencers/users lists
12. *Influencers Impact and Profile Credibility Measurement:* Calculate the impact and profile credibility ratios for the selected 100 influencers based on their profile biographies
13. *Influencers General Impact and Content Credibility Measurement:* Calculate the general impact and content credibility ratios for the selected 100 influencers based on their historical activity

14. *Influencers Reputation Calculation:* Calculate the influencers reputation rate and the general reputation rate for the selected 100 influencers
15. *Influencers Reputation Maximization:* Combine the previously calculated rates to maximize the reputation rate for the selected 100 influencers
16. *Influencers Sorting:* Sort and display the result of the maximized reputation rates
17. *Influencers Exported List:* The result shows the top 100 influencers in the selected event using the joint approach model and methodology

### 3.4 Joint Theme and Event Model for Maximizing the Influence Ratio Prediction

Given a set of users  $U$  with corresponding tweets  $T$ , followers  $F$ , and followees  $F_e$ , select the influential list of users  $L_u \in U$  that maximizes the accuracy  $ACC$  and reputation  $REP$  of each influencer  $u \in U$ .

In this work, our objective is to maximize: 1) the accuracy  $ACC$  achieved for each influencer 2) the reputation  $REP$  obtained for each influencer based on joint theme and event based models:

$$\text{maximize}(ACC) + (REP) \tag{3.1}$$

Our proposed methodology is depicted in Fig.3. We consider a sample of users  $U$  engaging in a certain event and we are interested in identifying influencers evolving in this event. We assume that the events relate to a certain country and hence, we constrain the sample of users to one country. Our aim is to rank this sample of users to get the most influential users among them. To achieve this goal, we use a joint theme and event based approaches for ranking the given sample of users. While the event based approach aims at rating the influence of

each user, the theme based approach gets general list of influencers and apply classification algorithm to get general influencers. We account for irrelevance in terms of spam accounts and we try to eliminate irrelevance as much as possible. We further optimize the list of influencers, crossed between both approaches, to get the over list of influencers  $L_e^t$ . We then optimize the reputation of those influencers through selecting top influencers in  $L_e^t$  and applying a reputation based approach as a second step to validate and enhance the output of list of influencers. Our problem is composed of two main steps 1 and 2 through which we maximize the accuracy  $ACC$  for each influencer and then optimize according to the reputation  $REP$  objective. We use both event and theme model to depict the list of influencers. To achieve this objective, we jointly study a theme and event based selection models with different influence rate calculations. Obviously, the theme and event should relate to each other as this is critically important in order to achieve the aforementioned objective. We should note that the same sample of users should be used for this joint approach. We further discuss the aforementioned objective in the following:

### 3.4.1 Accuracy Maximization

We aim at maximizing the accuracy for each selected influencer in a certain event. We define our objective of accuracy maximization in terms of influence rate  $INR$  and irrelevance  $IRR$  where we aim at maximizing the influence rate while minimizing the irrelevance as follows:

$$max_{IRC}(ACC) = max(INR) + MIN(IRR) \quad (3.2)$$

To achieve the aforementioned objective, we follow a joint approach composed of two main components: 1) Theme based Model 2) Event based Model. While theme-based selection strategy focuses on obtaining influencers with general and historical perspective in this theme, the event-based selection strategy aims at ob-

taining influencers particularly related to this special event using multiple metrics. We further describe the steps of each model in what follows:

1) *Theme General Influence Rate Calculations*: After fetching the sample of users  $U$  to using Twitter API, we aim at identifying general influencers with at least 1000 followers within the given theme. Our general influence rate  $GIR$  model entails the historical influence achieved for each influencer as follows:

$$GIR = w_1(FFr) + w_2(Tcr) \quad \forall w \in [0, 1] \quad \sum_{i=1}^2 w_i = 1 \quad (3.3)$$

where  $FFr$  corresponds to the followers followee rate and  $Tcr$  corresponds to the Tweet count rate. Due to the generality of this approach, multiple influential users with different backgrounds and roles are obtained within the given theme. As a result, we apply bagging algorithm in order to classify this list of influencers into different groups (Micro, General, Master). Our general influence classification criteria is as follows:

$$C(GIR) = \begin{cases} C_{mic} = 26 < GIR_{u=n}^m < 50 \\ C_{Gen} = 51 < GIR_{u=n}^m < 75, & \text{where } 1 \leq n \leq m \\ C_{Mas} = 76 < GIR_{u=n}^m < 100 \end{cases}$$

The list of influencers  $L_t$  is then obtained for this theme.

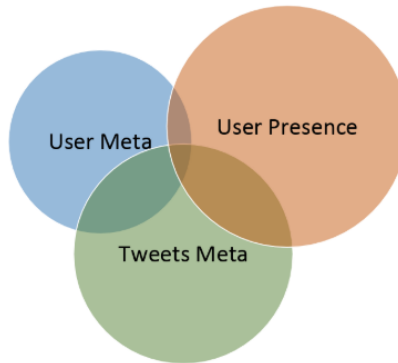


Figure 5: Venn diagram of the multiple measures that the UIR consists of: User Meta, User Presence, and Tweets Meta

2) *Event Influence Rate Approach*: The goal of event-based model is to select

the list of influencers  $L_e$  with the highest influence in a certain event. Fig. 5 depicts the different components of this approach and shows the general importance of each component. However, the weights may differ from one event to another. User-centric measures are of key importance as they essentially help in defining user’s engagement. We define two types of user-centric measures: User Meta Rate ( $UMR$ ) which is expressed in terms of time, number of followers and followees and User Presence Engagement Rate ( $UPER$ ) which captures the engagement of user in this particular event in terms of tweeting frequency and duration. Hence, for the same sample of users  $u$ , we aim at detecting the set of users  $u$  that are mostly engaging in this particular event who may not have been influencers at a previous stage before this event. We fetch all necessary data corresponding to tweets, users and meta and then filter this data according to a certain tweet cutoff.

We calculate  $UMR$ ,  $UPER$ , and  $UTMR$  as follows:

- i) User Meta Rate (UMR): While analyzing the content of the tweets play an important role in understanding the influence of users, it also very important to capture other related metrics pertaining to the user him self to better assess his/her influence. Our User Meta Rate  $UMR$  aims at capturing the time, followers and following rate in addition to the user listed rate to calculate this overall measurement as follows:

$$UMR = \alpha U_{SY} + \beta U_{Fr} + \gamma U_F + \kappa U_L \quad (3.4)$$

where  $\alpha$ ,  $\beta$ ,  $\gamma$  and  $\kappa$  correspond to the weights given to each term respectively.  $U_{SY}$  corresponds to the User Since Years Rate,  $U_{Fr}$  corresponds to Users Followers Rate,  $U_F$  corresponds to the User Following Rate and  $U_L$  corresponds to the User Listed Rate. We calculate each of the  $U_{SY}$ ,  $U_{Fr}$  and  $U_F$  as follows:

- User Since Year Rate ( $U_{SY}$ ): which depicts how long this user has been

using this account. A user with a long period of time may have a higher level of engagement or a larger number of connections than a user with a new registered account, and thus this user can be more influential than new registered user. To calculate the User Since Year Rate  $U_{SY}$  of a certain user, we consider the number of years starting the registration of this user over the total number of years starting Twitter platform as follows:

$$U_{SY} = \frac{\text{Number of Years Since User Registration}}{\text{Number of Years Since 2006}} * 100 \quad (3.5)$$

- User Followers Rate ( $U_{Fr}$ ): Generally, the number of followers can be considered as a key indicator for measuring the influence of a certain user. Since we aim at minimizing irrelevance, which corresponds to eliminating spam and fake accounts as much as possible, we use User Followers Rate  $U_{Fr}$  in our influence calculations whereby we consider the number followers for this particular user out of the total number of followers in this particular event:

$$U_{Fr} = \frac{\text{User Number of Followers}}{\text{Total Number of Followers for All Users}} * 100 \quad (3.6)$$

- User Following Rate  $U_F$ : We calculate the User Following Rate  $U_F$  in terms of the number of followings and the total number of followings as follows:

$$U_F = \frac{\text{User Number of Followings}}{\text{Total Number of Followings for All Users}} * 100 \quad (3.7)$$

The list of Rate Influencers  $L_e$  is then derived based on this approach. To obtain the overall list  $L_u$ , the two lists  $L_e$  and  $L_t$  are crossed to derived the list of influencers.

- User Listed Rate  $U_L$ :

$$U_L = \frac{\text{Number of Listed Times}}{\text{Total Number of Listed Times of all Users}} * 100 \quad (3.8)$$

- ii) User Presence Engagement Rate (UPER): Knowing that influencers are usually active but not all active users are influencers, the engagement can be considered as a very important factor that affects the total influence. We develop this metric in order to quantify the productivity of a certain user. Posting a large number of tweets is one of the key indicators about the intense level of engagement of the user in this event. We consider User Presence Engagement Rate  $UPER$  as follows:

$$UPER = \frac{\text{User Number of Tweets} * UER}{\text{Total Number of Tweets for All Users}} \quad (3.9)$$

where  $UER$  corresponds to User Engagement Rate as follows:

$$UER = \frac{\text{User Duration of Tweeting}}{\text{Total Duration of Tweeting for All Users}} * 100 \quad (3.10)$$

Where we calculate User Duration of Tweeting as follows:

$$\text{User Duration of Tweeting} = \text{User Last Tweet Date} - \text{User First Tweet Date} \quad (3.11)$$

And we calculate the total duration of tweeting as follows:

$$\text{Total Duration of Tweeting} = \text{Last Tweet Date} - \text{First Tweet Date} \quad (3.12)$$

- iii) User Tweet Meta Rate (UTMR): Quantifying the importance of tweet generated by a certain user is very critical in determining the influence of this user and hence, we further define User Tweet Meta Rate  $UTMR$  in terms of



retweets and favorites rate achieved by a certain initiated tweet as follows:

$$\text{User Tweet Meta Rate} = 60\%(\text{User Retweets Rate}) + 40\%(\text{User Favorites Rate}) \quad (3.13)$$

Where we calculate the User Retweets Rate and User Favorites Rate respectively as follows:

$$\text{User Retweets Rate} = \frac{\text{User Number of Retweets}}{\text{Total Number of Retweets for All Users}} * 100 \quad (3.14)$$

$$\text{User Favorites Rate} = \frac{\text{User Number of Favorites}}{\text{Total Number of Favorites for All Users}} * 100 \quad (3.15)$$

To calculate the influence rate, we combine  $UMR$ ,  $UPER$  and  $UTMR$  as follows:

$$INR = w_1(UMR) + w_2(UPER) + w_3(UTMR) \quad \forall w \in [0, 1] \quad \sum_{i=1}^3 w_i = 1 \quad (3.16)$$

where  $w_1$ ,  $w_2$  and  $w_3$  corresponds to the weights of each of  $UMR$ ,  $UPER$  and  $UTMR$  respectively. Crossing each of the event-based model and theme-based model leads to a list of influencers  $L_e^t$ .

### 3.4.2 Reputation Maximization

We further consider sampling out of selected influencers  $L_e^t$  to apply a second step with the objective of maximizing the of selected influencers. We select a list of influencers 1.1 (please refer to figure 3) out of the user list  $L_e^t$ . Our objective is represented as follows:

$$REPR = w_1(INGREPr) + w_2(INREPr) \quad \forall w \in [0, 1] \quad \sum_{i=1}^2 w_i = 1 \quad (3.17)$$

where  $INGREPr$  corresponds to the Influencer General Reputation Rate and  $INREPr$  corresponds to the Influencer Reputation Rate. The Influencer General Reputation Rate  $INGREPr$  aims at capturing the reputation rate of the selected

influencer before the occurrence of the event and thus gets a historical reputation about this user. We calculate the *INGREPr* in terms of Influencer General Impact *INGIMPr* and Influencer General Tweets Credibility *INGtCRDr* based on user theme tweets as follows:

$$INGREPr = w_1(INGIMPr) + w_2(INGtCRDr) \quad \forall w \in [0, 1] \quad \sum_{i=1}^2 w_i = 1 \quad (3.18)$$

We calculate *INGIMPr* and *INGtCRDr* as follows:

$$INGIMPr = \frac{T_t * Fr + R_t + F_t}{\text{Total Impact of All users}} * 100 \quad (3.19)$$

where  $T_t$  corresponds to the User Number of tweets per theme,  $Fr$  corresponds to the Number of Followers,  $R_t$  corresponds to User Number of Retweets per theme and  $F_t$  corresponds to User Number of Favorites per theme.

$$INGtCRDr = \begin{cases} 1 * \left( \frac{\text{Number of user's tweets per theme}}{\text{Total Number of Tweets}} \right) * 100 & \text{if credible} \\ 0 & \text{if not credible} \end{cases} \quad (3.20)$$

We also aim at assessing the reputation of this selected influencer in the context of this event and hence, we further calculate the Influencer Reputation Rate *INREPr* as follows:

$$INREPr = w_1(INIMPr) + w_2(INpCRDr) \quad \forall w \in [0, 1] \quad \sum_{i=1}^2 w_i = 1 \quad (3.21)$$

where *INIMPr* corresponds to the Influencer Impact and *INpCRDr* corresponds to the Influencer Profile Credibility.

where *INIMPr* and *INpCRDr* are defined as follows:

$$INIMPr = \frac{T_e * Fr_e + R_e + F_e}{\text{Total Impact of All users}} * 100 \quad (3.22)$$

where  $T_e$  corresponds to the User Number of tweets per event,  $Fr$  corresponds to

the Number of Followers,  $R_e$  corresponds to User Number of Retweets per event and  $F_e$  corresponds to User Number of Favorites per event.

$$INpCRDr = \begin{cases} 1 * \left( \frac{\text{Number of user's tweets per event}}{\text{Total Number of Tweets}} \right) * 100 & \text{if credible} \\ 0 & \text{if not credible} \end{cases} \quad (3.23)$$

### 3.5 Conclusion

In this chapter, we discussed our proposed influencer rating joint approach, displayed the full approach flowchart and diagram, defined our implementation and analysis roadmap for the next chapter, and illustrated all the mathematical models and the optimization calculations and functions that can show the complex formation of our implemented scripts while processing and analysing datasets.

# Chapter Four

## Implementation and Experimental Results

### 4.1 Implementation

#### 4.1.1 Topic Selection and Methodology Overview

The adopted selection methodology, illustrated in Figure 4 in Chapter 3, starts by selecting the event topic and choosing its relevant hashtags and keywords that were used to build the basic twitter search query in order to get the targeted results. After that, the system collects one million unique tweets with all their relevant unique user profiles. After fetching tweets and user profiles, additional aggregations and labels were added to each tweet and/or profile which indicate a few additional attributes and classifications.

- One million public tweets were collected using a hard-coded python script that uses Tweepy [101]. All tweets contains at least one of the following keywords (terms not hashtags) "corona", "covid", or "sarscov2" which are the most used keywords during the COVID-19 global event. After fetching all tweets, all unique users are listed to be collected using Twitter REST API V1 [102] access tokens in order to be analysed and classified in later stages.

- In order to classify tweets whether they are corona related or non corona related tweets (context based), an ontology/lexicon of common related keywords/terms is produced and used within the NLP entity extractor modules in order to perform the mentioned classifications. Below are the different lexicons used during our study to label tweets and profiles based on the NLP results: Corona Top Used Hashtags Lexicon, Corona Social Media Context Lexicon for Tweets, Occupation Lexicon for Grouping Users Based on their Biographic Information, Medical Occupation Lexicon for Users and Virus Specialty Occupation Lexicon for Users.
- After building all the mentioned lexicons, each single tweet was labeled as "isCorona" or "isNotCorona" tweet, and each user profile was labeled as "isMedical" or "isSpecialty" based on the users claimed biography information. In addition to that, a specific occupation field was added to each user based on their claimed biography details as well.
- An enriched dataset was next designed to help aggregating and creating more relationships between all the data points and the available classifications. The obtained dataset is structured as the following: TweetID, TweetHashtagCount, TweetFavorates, TweetRetweets, TweetMentions, TweetTotalInteractions (favorites plus retweets, TweetTotalReach (based on the number of followers for the tweet's user), UserID, UserClaimedLocation, UserOccupation (extracted from the user's profile), isCoronaTweet, isMedicalProfileUsers, and isSpecialtyProfileUser.

### 4.1.2 Data Processing

We present more details about our proposed framework in Figure 3 (Chapter 3) in the following:

1. *Event Selection:* In this step, we select the event parameters such as the event, location, language, duration. We select "COVID-19" as our event

with global location and restrict the language to English Language. Our data was collected from Jan 25 2020 to March 20 2020 with a size equal 1 million unique tweet and 288.4 thousands unique user. (More details about event topic selection in Figure 3 in Chapter 3)

2. *Event Data Collection:* This is done through querying all tweets containing at least one of the following terms: “covid, corona, sarscov” over Twitter public blog.
3. *Event Data Cleaning and Filtering:* We use GCP DataFlow and BigQuery to manage and execute sequence of steps related to cleaning and filtering the tweets.
4. *Event Data Analysis:* We utilize GCP DataFlow, BigQuery Analysis, and custom python scripts to analyze the tweets.
5. *Complex Approach Calculation:* We apply the previously mentioned approach in section 4.3 in terms of calculating Influencer Rating, Classifying Influencers, and Influencers Credibility and Impact Measurements.
6. *Event Exported Datasets:* We then obtain, based on the aforementioned calculations, a dataset of 1 million unique tweets, dataset of top 1000 Influencers (ordered by number of event-tweets), and a sorted list of 1000 Event Influencers (based on the result of the INR).
7. *Theme Selection:* In this step, we select the theme parameters such as the event, location, language, duration. We select “Medical, Virus” as our event with global location and restrict the language to English Language. Our data was collected from Oct 01 2020 to Dec 05 2020 with a size equal 850 thousands unique tweets and 401.7 thousands unique user. (More details about theme topic selection in Figure 3 in Chapter 3)
8. *Theme Data Collection:* We get all tweets containing at least one of the following keywords: “medic, virus, vaccine” over Twitter Public Blog.

9. *Theme Data Cleaning and Filtering:* We use GCP DataFlow and BigQuery to manage and execute sequence of steps related to cleaning and filtering the tweets.
10. *Theme Data Analysis:* We utilize GCP DataFlow, BigQuery Analysis, and custom python scripts to analyze the tweets.
11. *Theme-based Calculations and Findings:* We further, based on the mentioned approach in section 4.3, calculate General Influence Rates, Classify General Influencers, and calculate activity ratios for all theme influencers.
12. *Theme Exported Datasets:* We then obtain, based on the aforementioned calculations, a dataset of 1 million unique tweets, dataset of top 1000 General Influencers (ordered by number of followers) and a sorted list of 1000 General Influencers (based on the result of the GIR).
13. *Joint Theme-Event Exported Datasets:* While calculating influence rates and maximizing reputation for the selected influencers, the following joint theme-event datasets are obtained, a dataset of all event influencers with maximized influence ratios after merging theme-based general influence rates, a dataset of the selected 100 influencers with maximized reputation ratios, and finally, a sorted dataset of 100 influencers (based on the final Reputation Rate REPR) and this is the final exported result that contains the top 100 influencers for the selected event

## 4.2 Results and Analysis

In this section, we provide different performance analysis that highlights numerical evaluation. We choose COVID-19 as a main event. We should note that the steps of data processing are presented in the previous section (Section 4.1.2). We show results of the performance of our solution in each approach alone and then we provide the overall findings and results

## 4.2.1 Event Based Analysis

We first extract around 1 million tweets for a total number of users equal to 288,439, a total number of followers equal to 21,393,493 and a total number of interactions equal to 36,964,431,381. The figure 6 depicts influencers in the context of COVID event. As the figure shows, top influencers include *EclipseMist*, *evankirstel*, and *my\_amigouk* with INR equal to 41.92, 37.21, and 35.27 respectively.

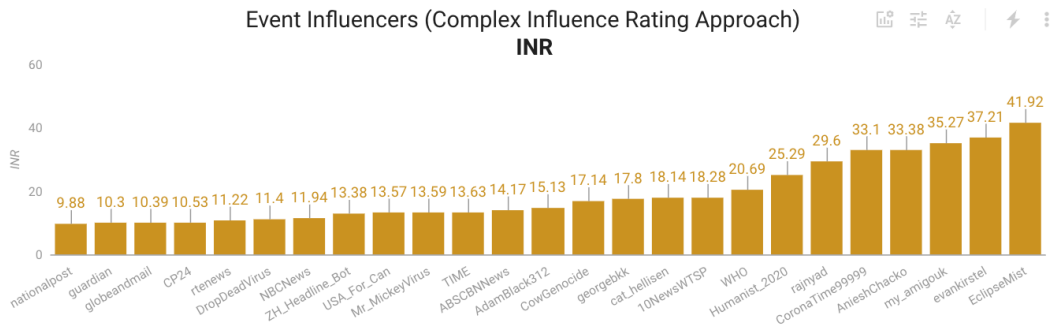


Figure 6: Event Influencers using INR Proposed Calculations

We further in figure 7 show top 1000 influencers with their corresponding UMR, UPER, and UTMR ratios (refer to section 4.3 for more details). The figure illustrates the non-linear and dense relationship among those three components.

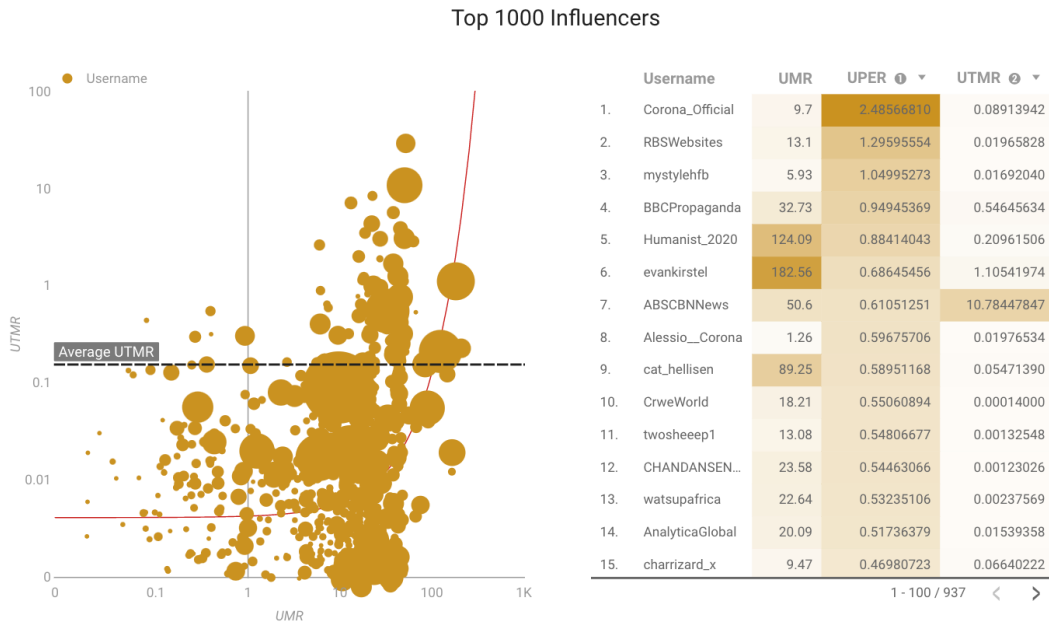


Figure 7: Top 1000 Influencers



As the figure 7 shows, the distribution of the UMR, UPER and UTMR vary from one user to another, recording different username for the three different maximums. For example, UPER is maximum for user numbered 1, the UMR is maximum for user numbered 6 and the UTMR is maximum for user number 7. This due to the fact that each user has different tweeting behavior and characteristics.

### 4.2.2 Theme Based Analysis

We extract around 1 million tweets with for a total number of users equal to 391,795, a total number of followers equal to 16,973,449,114 and a total number of interactions equal to 6,950,323. The figure 6 depicts influencers in the context of COVID event. As the figure shows, top influencers include *UberFacts*, *VoceNaoSabiaQ*, and *ANI* with INR equal to 335.6, 90.39 , and 83.35 respectively.

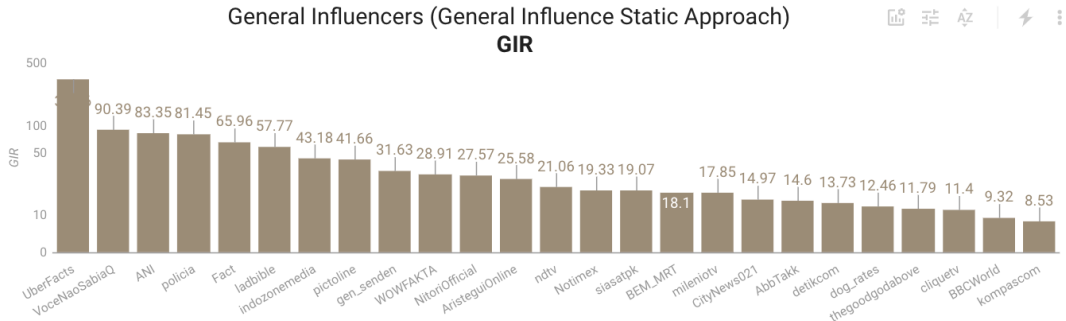


Figure 8: Selecting Influencers using the General Approach

We further in figure 9 shows top 1000 influencers with their corresponding Followers, Tcr and FFr. A non-linear relationship is also depicted when using this approach but with less density.

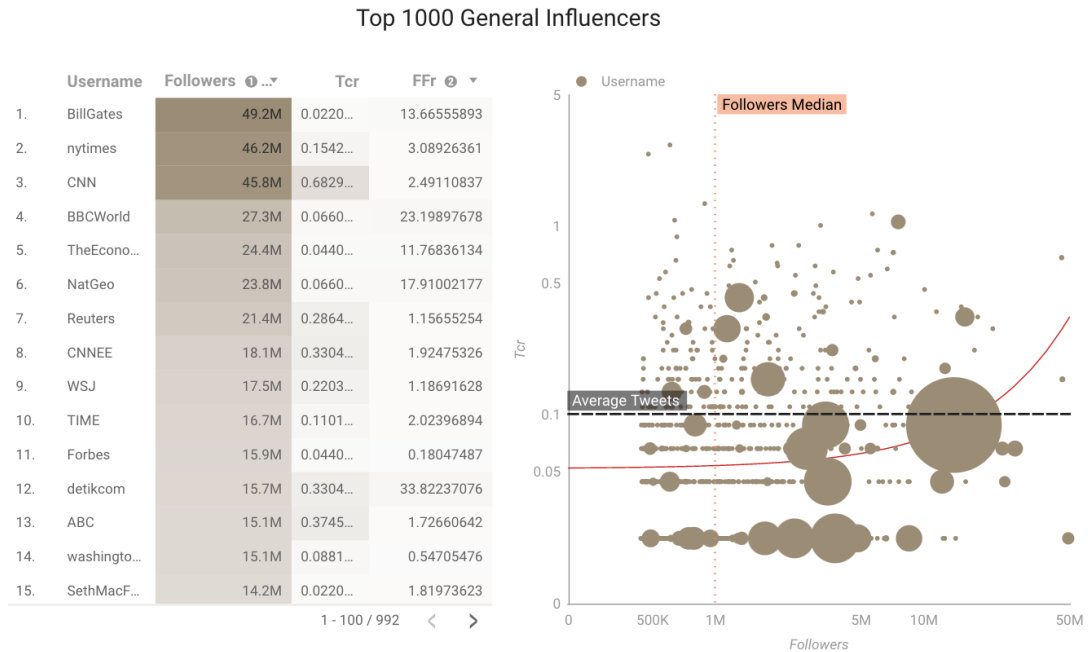


Figure 9: Top 1000 General Influencers

### 4.2.3 Joint Event and Theme Based Analysis

The resulting list of influencers using the joint theme and event based approach is depicted in figure 10. As can be depicted, the GIR alone fails to capture influencers in the context of COVID while the INR achieves a higher level of accuracy. The figure shows further analysis on the scale of tweets achieved by influencers captured by each of the event and theme based approaches individually. Selected influencers in the theme based approach reach an average of 4.54 tweet per influencer (0.45%) which is less than that of the event based approach with 143.27 tweets per influencer (13.56%) and this signifies the impact of influencers selected when tweeting in terms of the total number of users.

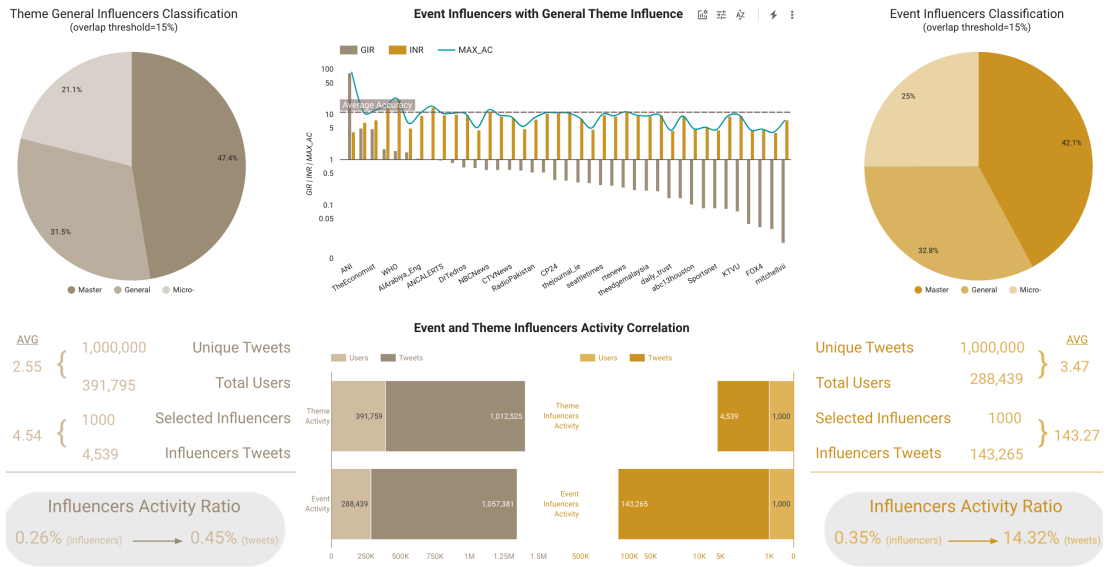


Figure 10: List of Influencers with Joint Approach

Further, the credibility of those top influencers is measured in terms of content credibility, profiles, virus clinical profile and company in figure 11. The percentage of content credibility is low in comparison with the total content released about COVID. This shows that most influencers did not have credible content. Similarly, medical profiles record a low percentage in the set of influencers, this can be due to the fact that most influencers pertain to the group of journalists and news agencies which release information.

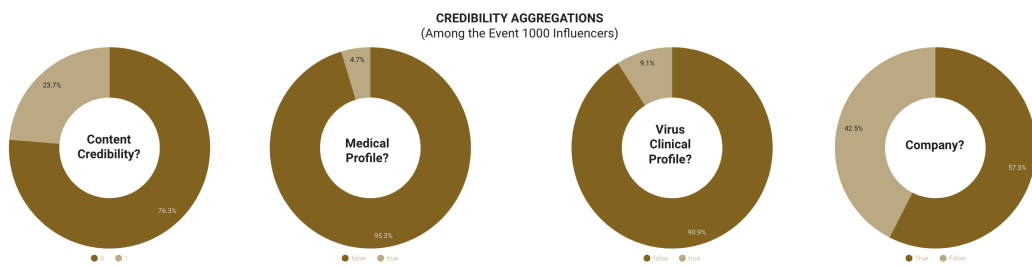


Figure 11: Credibility Analysis and Aggregations for Selected Influencers of Joint Approach

#### 4.2.4 Reputation Based Analysis

In this subsection, we compute the reputation for Top 100 influencers. Figure 12 shows the distribution of reputation for both Influencer General Reputation

## Calculation (INGEPr) and Influencer Reputation Calculation (INREPr).

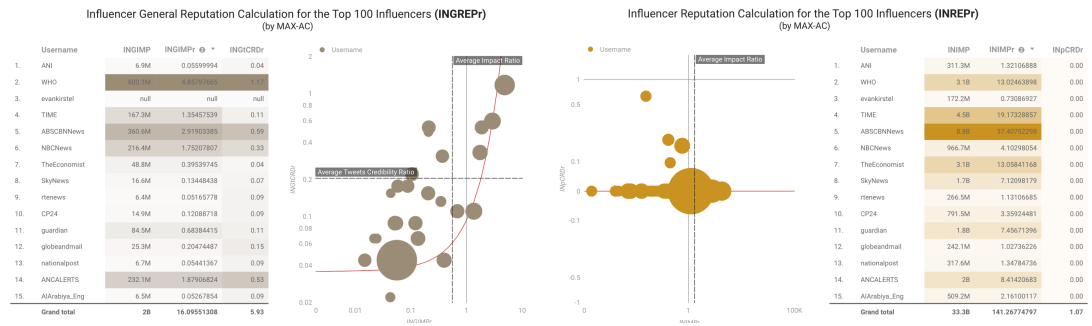


Figure 12: Comparison between list of selected influencers with influencer general reputation and Influencer Reputation Calculation

Figures 13 and 14 further lists main influencers while highlighting major components of the reputation based approach. The additional components added in reputation.

For example, WHO, ANI, and ABS CBL news

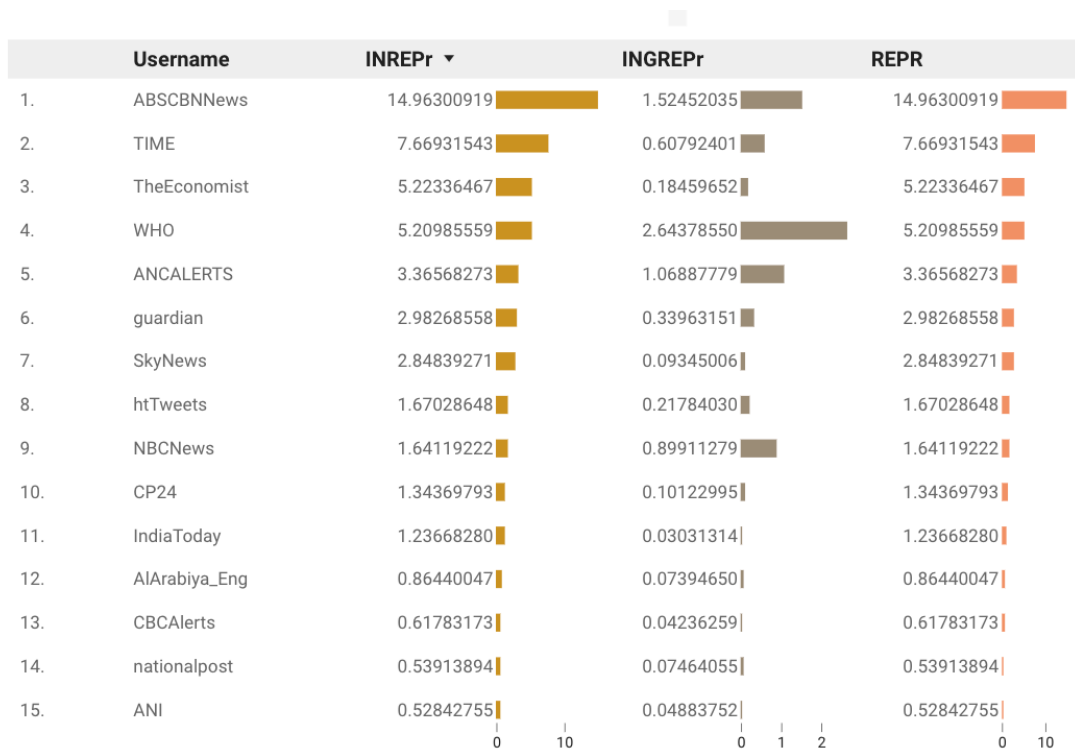


Figure 13: List of Influencers with Reputation Maximization

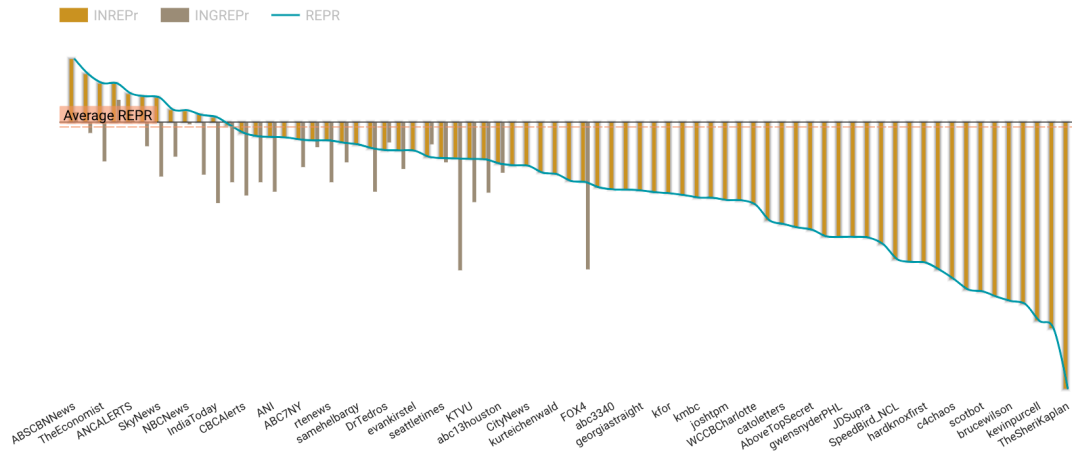


Figure 14: Multiple bar chart showing Influencers with Reputation Maximization

The figure 15 below compares the achieved credibility of top influencers between the event and theme approaches in terms of their provided content, medical profile, virus clinical profile and company. As the figure depicts, the credibility of event-based influencers is much higher than the credibility of theme-based influencers at the four dimensions.

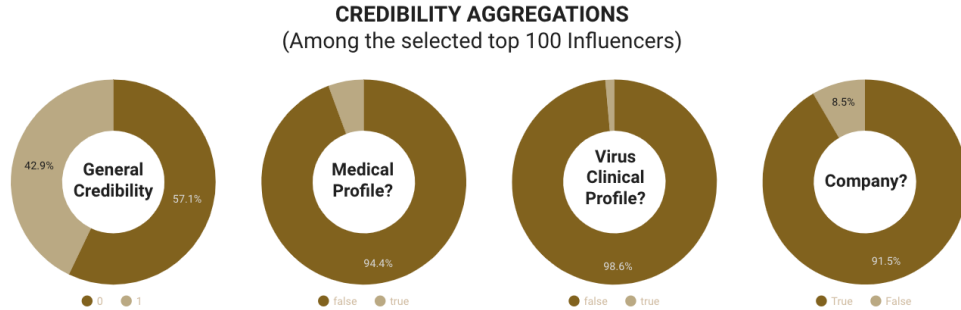


Figure 15: Credibility Analysis and Aggregations for List of Influencers with Maximized Reputation

### 4.3 Discussion

In this section, we discuss the different results showed in section 4.2 while highlighting important aggregations and relationships. Starting from the very beginning of our approach, the selection, selecting hot or trending events in social media is basically related to those people or agents who are called users who tweet and publish content using similar/common terminologies or hashtags and

thus create a topic traffic and interest, later called "Trend". Such event might vary between different types of social media platform "Twitter", "Facebook", "Instagram", "YouTube", or other platforms based on the content format, targeted audience, and other factors. Elections, Crisis, or Protests are some of many different types of events that might happen every day on social media. From the mentioned details, finding event real influencers can help in understanding different user behaviors and thus recommending different strategies and further research enhancements and directions. In our approach, influencers are not just those who have large numbers of followers (even if this have a high impact ratio in general) but that does not mean that they have to be the influencers and/or the leaders of a certain event.

After selecting and sorting out 1000 influencers from the selected event based on their engagement (number of tweets within the selected event time-frame), some of the names that showed at the beginning of the list continued with us and also showed at the final exported list. While at the same time, most of the selected 1000 theme influencers based on their followers count (number of followers) that were at the very top of the sorted list, were disappeared in the next maximization levels, which proves that user number of followers is not necessarily a main factor in the whole calculation. However, when calculating event influence rates, we take into account user meta information and we also add it our calculation but with different weights. Presence Engagement is the most important factor in this part of calculation, which means that the engagement level of the selected user within the event compared to all engagements by all selected users, and this might vary a lot between users for a single event. In addition, Tweets Meta Rate which is the level of interactions with a certain tweet is also reflects the impact of the mentioned tweet, and thus the possible reach that it might create between users and their networks, and that also is a primary factor within the mentioned calculation. One more highlight in this part, is that we can clearly see in the GIR exported results and charts in figure 9 the linear relationship between

followers count (the sorting element) and the FF ratio for example, while there is no relation at all with the tweeting count rate  $T_{cr}$ . At the same time, figure 6 shows the non-relational connection between the three calculated rates: UMR, UPER, and UTMR which is organically true. As a conclusion, each one of the mentioned factors has its own weight that drive the user to the top or the bottom of the final sorted list of influencers.

Adding historical influence rate to the selected event influencers can maximize the accuracy of the assumption, in which, in our selected case study, some of the selected event influencers have historical general influence in similar events or at least in the same field. That way, we can support our finding and maximize the accuracy of the calculation (see figure 10).

At this stage, calculated event influencers with theme general influence can still have high rates than other, while this might be the only accuracy threshold. To eliminate the mentioned threshold, the reputation calculation was necessary and important in which it combines both content and profile credibility in addition to content and profile impact/popularity analysis within the event timeframe and the theme timeframe (which is before the event get started). This calculation and findings proves that highly influential theme and event users may vary upon the credibility and impact findings, and that is exactly the truth. For instance, a journalist who has a huge network of followers and a high level of engagement during the event, and whom his tweets get a large number of interactions, this user is clearly considered as an event influencer, but at the same time, this specific user might be giving advice about COVID-19 during the event for the first time while his profile does not match the content he provides and the event topic which means that his profile is not credible enough to talk about the field like "Medical" or "Healthcare" profiled users. So those tweets which supposed to be tweeted by credible profiled users will have higher content credibility ratios and thus higher profile credibility rates for their users. In this case, the final list of influencers can change based on the reputation rate sorting (see figure 14).

# Chapter Five

## Conclusion and Future Work

In this thesis, we studied the problem of influence maximization for users over Twitter social media platform while maximizing the accuracy and minimizing irrelevance through a joint model combining a theme and event based models and then we, using the list of influencers obtained, applied a maximization for reputation achieved for this set of obtained influencers. We formulated the problem in terms of different influence calculations metrics to capture the various metrics that depicts the behavior of social media users and can be affecting the influence of user in a certain event.

Further, this work can be leveraged in any event to obtain top influencers and depict the credibility at multiple dimensions. Multiple questions can be answered including but not limited to the following:

- If a user is has a very high engagement rate during the event, can this specific user has lower reputation than other users?
- Can an influencer be considered credible in a certain event even if he did not tweet before within the theme of the event?
- How much credible are those influencers in an event in terms of their content? profile? and historical activities?
- How we can defer between general influencers and event influencers?



- Can an influencer who has a very large network and reach level considered as non-influencer even if he has a good participation in the event?
- Can we consider an influencer who has a very small network of followers as an event influencer?
- Can a user who has a very high participation or engagement level during the event to be consider as non-influencer?
- How can we find hidden influencers in a similar event?
- How does time affect on being an influencer?
- How can this combined approach achieve a higher level of accuracy?

Moreover, we discuss further future work foreseen based on this proposed methodology. An interesting extension for this work can be integrating intelligent solutions that capture sentiment meaning in influence calculation. Using Natural Language Processing, the influence of tweets posted by users can be also utilized to depict the impact of each tweet that meets a certain predefined key point indicators, and thus, we can enhance the accuracy of our influence based approach as negative and positive tweets about users play a key role in affecting the total influence of users. Another extension for this work could be to study the influence of each engaged group alone such as organizations and people of different professions and quantify the influence of each group on certain audiences. Moreover, an additional extension for this work is to deeply identify and find the credible voices that can be very helpful for governments as alias to defeat fake news and other misinformation styles that happens during similar critical event. In addition, adding graph network analysis into this approach by applying interconnected relationship influence ratios my be a great enhancement in the future. Further, we can consider more platforms to be studies in the context of a certain event to detect influencers and achieve a higher level of accuracy.

# Bibliography

- [1] “Strengthen grow your network with social network analysis,” 2015.
- [2] M. Kim, J. Cobb, M. Harrold, T. Kurc, A. Orso, J. Saltz, A. Post, S. Malhotra, and S. Navathe, “Efficient regression testing of ontology-driven systems,” *Proceedings of the 2012 International Symposium on Software Testing and Analysis*, pp. 320–330, 07 2012.
- [3] D. Kempe, J. Kleinberg, and E. Tardos, “Maximizing the spread of influence through a social network,” in *Proceedings of the Ninth ACM SIGKDD International Conference on Knowledge Discovery and Data Mining*, KDD ’03, (New York, NY, USA), p. 137–146, Association for Computing Machinery, 2003.
- [4] L. Adamic and E. Adar, “How to search a social network,” *Social Networks*, vol. 27, no. 3, pp. 187 – 203, 2005.
- [5] X. Song, B. L. Tseng, C.-Y. Lin, and M.-T. Sun, “Personalized recommendation driven by information flow,” in *Proceedings of the 29th Annual International ACM SIGIR Conference on Research and Development in Information Retrieval*, SIGIR ’06, (New York, NY, USA), p. 509–516, Association for Computing Machinery, 2006.
- [6] B.-H. J. R. A. González-Bailón, Sandra and M. Yamir, “The dynamics of protest recruitment through an online network,” p. 197, *Scientific Reports*, 2011.

- [7] J. Borge-Holthoefer and Y. Moreno, “Absence of influential spreaders in rumor dynamics,” *Phys. Rev. E*, vol. 85, p. 026116, Feb 2012.
- [8] G. Stilo, P. Velardi, A. E. Tozzi, and F. Gesualdo, “Predicting flu epidemics using twitter and historical data,” in *Brain Informatics and Health* (D. Ślzak, A.-H. Tan, J. F. Peters, and L. Schwabe, eds.), (Cham), pp. 164–177, Springer International Publishing, 2014.
- [9] J. Clement, “Number of monthly active twitter users worldwide from 1st quarter 2010 to 1st quarter 2019,” 2019. Available at <https://www.statista.com/statistics/282087/number-of-monthly-active-twitter-users/>.
- [10] A. Java, X. Song, T. Finin, and B. Tseng, “Why we twitter: Understanding microblogging usage and communities,” *of the 9th WebKDD and 1st SNA*, vol. 43, pp. 56–65, 01 2007.
- [11] V. Nebot, F. Rangel, R. Berlanga, and P. Rosso, “Identifying and classifying influencers in twitter only with textual information,” in *Natural Language Processing and Information Systems* (M. Silberstein, F. Atigui, E. Kornyshova, E. Métais, and F. Meziane, eds.), (Cham), pp. 28–39, Springer International Publishing, 2018.
- [12] A. Mourad, A. Srour, H. Harmanani, C. Jenainatyi, and M. Arafeh, “Critical impact of social networks infodemic on defeating coronavirus covid-19 pandemic: Twitter-based study and research directions,” *arXiv preprint arXiv:2005.08820*, 2020.
- [13] “What is a rest api?.” Available at [https://idratherbewriting.com/learnapidoc/docapis\\_what\\_is\\_a\\_rest\\_api.html](https://idratherbewriting.com/learnapidoc/docapis_what_is_a_rest_api.html).
- [14] “Twitter api tutorial.” Available at <http://socialmedia-class.org/twittertutorial.html>.

- [15] P. scripts and modules, “Python scripts and modules,” Univesity of Wash-  
ington. Available at [https://faculty.washington.edu/rjl/classes/  
am583s2013/notes/python\\_scripts\\_modules.html](https://faculty.washington.edu/rjl/classes/am583s2013/notes/python_scripts_modules.html).
- [16] Google, “Why google cloud platform,” 2014.
- [17] S. Stemler, “An overview of content analysis. - practical assessment, re-  
search evaluation,” 2016.
- [18] G. Nandi, “A survey on using data mining techniques for online social net-  
work analysis,” *International Journal of Computer Science Issues*, vol. 10,  
p. 162, 11 2013.
- [19] N. F. Noy and D. L. McGuinness., “Ontology development 101: A guide  
to creating your first ontology,” *Stanford Knowledge Systems Laboratory  
Technical Report KSL-01-05 and Stanford Medical Informatics Technical  
Report SMI-2001-0880*, 03 2001.
- [20] K. Chowdhary, “Natural language processing,” in *Fundamentals of Artifi-  
cial Intelligence*, pp. 603–649, Springer, 2020.
- [21] C. Fombrun, N. Gardberg, and J. Sever, “The reputation quotientism: A  
multi-stakeholder measure of corporate reputation,” *Journal of Brand Man-  
agement*, vol. 7, 07 2013.
- [22] M. El Marrakchi, H. Bensaid, and M. Bellafkih, “Scoring reputation in on-  
line social networks,” in *2015 10th International Conference on Intelligent  
Systems: Theories and Applications (SITA)*, pp. 1–6, 2015.
- [23] Y. Mei, W. Zhao, and J. Yang, “Influence maximization on twitter: A  
mechanism for effective marketing campaign,” pp. 1–6, 05 2017.
- [24] M. Cha, H. Haddadi, F. Benevenuto, and K. P. Gummadi, “Measuring user  
influence in twitter: The million follower fallacy,” in *ICWSM*, 2010.

- [25] A. Goyal, F. Bonchi, and L. V. S. Lakshmanan, “A data-based approach to social influence maximization,” *Proc. VLDB Endow.*, vol. 5, p. 73–84, Sept. 2011.
- [26] E. Katz and P. Lazarsfeld, *Personal Influence, the Part Played by People in the Flow of Mass Communications*. A Report of the bureau of applied social research Columbia university, Free Press, 1966.
- [27] C. Castillo, M. Mendoza, and B. Poblete, “Information credibility on twitter,” WWW ’11, (New York, NY, USA), p. 675–684, Association for Computing Machinery, 2011.
- [28] M.-A. Abbasi and H. Liu, “Measuring user credibility in social media,” in *Social Computing, Behavioral-Cultural Modeling and Prediction* (A. M. Greenberg, W. G. Kennedy, and N. D. Bos, eds.), (Berlin, Heidelberg), pp. 441–448, Springer Berlin Heidelberg, 2013.
- [29] S. Kwon, M. Cha, K. Jung, W. Chen, and Y. Wang, “Prominent features of rumor propagation in online social media,” in *2013 IEEE 13th International Conference on Data Mining*, pp. 1103–1108, 2013.
- [30] N. L. R. M. Merchant., “Social media firms will use more ai to combat coronavirus misinformation, even if it makes more mistakes.,” 2020. Available at <https://thenextweb.com/neural/2020/03/17/social-media-firms-will-use-more-ai-to-combat-coronavirus-misinformation->
- [31] M. Shirakawa, K. Nakayama, T. Hara, and S. Nishio, “Wikipedia-based semantic similarity measurements for noisy short texts using extended naive bayes,” *IEEE Transactions on Emerging Topics in Computing*, vol. 3, pp. 205–219, 2015.
- [32] C. Sumner, A. Byers, R. Boochever, and G. J. Park, “Predicting dark triad personality traits from twitter usage and a linguistic analysis of tweets,”

*2012 11th International Conference on Machine Learning and Applications*, vol. 2, pp. 386–393, 2012.

- [33] C. Chen, Y. Wang, J. Zhang, Y. Xiang, W. Zhou, and G. Min, “Statistical features-based real-time detection of drifted twitter spam,” *IEEE Transactions on Information Forensics and Security*, vol. 12, no. 4, pp. 914–925, 2017.
- [34] P. Meel, H. Agrawal, M. Agrawal, and A. Goyal, “Analysing tweets for text and image features to detect fake news using ensemble learning,” in *International Conference on Intelligent Computing and Smart Communication 2019* (G. Singh Tomar, N. S. Chaudhari, J. L. V. Barbosa, and M. K. Aghwariya, eds.), (Singapore), pp. 479–488, Springer Singapore, 2020.
- [35] B. Halawi, A. Mourad, H. Otrok, and E. Damiani, “Few are as good as many: An ontology-based tweet spam detection approach,” *IEEE Access*, vol. PP, pp. 1–1, 10 2018.
- [36] W. Gad and S. Rady, “Email filtering based on supervised learning and mutual information feature selection,” in *2015 Tenth International Conference on Computer Engineering Systems (ICCES)*, pp. 147–152, 2015.
- [37] S. Sedhai and A. Sun, “Semi-supervised spam detection in twitter stream,” *IEEE Transactions on Computational Social Systems*, vol. 5, no. 1, pp. 169–175, 2018.
- [38] B. A. A. Alghamdi, J. Watson, and Y. Xu, “Toward detecting malicious links in online social networks through user behavior,” in *Proceedings of the 2016 IEEE/WIC/ACM International Conference on Web Intelligence Workshops (WIW)* (S. Ferilli, M. Reformat, and S. Belkasim, eds.), pp. 5–8, United States of America: IEEE, 2016.

- [39] S. Lee and J. Kim, “Warningbird: A near real-time detection system for suspicious urls in twitter stream,” *IEEE Transactions on Dependable and Secure Computing*, vol. 10, no. 3, pp. 183–195, 2013.
- [40] A. Guille and C. Favre, “Mention-anomaly-based event detection and tracking in twitter,” *2014 IEEE/ACM International Conference on Advances in Social Networks Analysis and Mining (ASONAM 2014)*, pp. 375–382, 2014.
- [41] A. A. Amleshwaram, N. Reddy, S. Yadav, G. Gu, and C. Yang, “Cats: Characterizing automation of twitter spammers,” in *2013 Fifth International Conference on Communication Systems and Networks (COMSNETS)*, pp. 1–10, 2013.
- [42] F. Benevenuto, T. Rodrigues, J. Almeida, M. Gonçalves, and V. Almeida, “Detecting spammers and content promoters in online video social networks,” in *Proceedings of the 28th IEEE International Conference on Computer Communications Workshops, INFOCOM’09*, p. 337–338, IEEE Press, 2009.
- [43] H. Shen and X. Liu, “Detecting spammers on twitter based on content and social interaction,” in *2015 International Conference on Network and Information Systems for Computers (ICNISC)*, (Los Alamitos, CA, USA), pp. 413–417, IEEE Computer Society, jan 2015.
- [44] V. Visansirikul and S. Kitisin, “Identifying influencers with ensemble classification approach on twitter,” in *2018 22nd International Computer Science and Engineering Conference (ICSEC)*, pp. 1–4, 2018.
- [45] L. Wang and J. Q. Gan, “Prediction of the 2017 french election based on twitter data analysis,” in *2017 9th Computer Science and Electronic Engineering (CEECE)*, pp. 89–93, 2017.

- [46] L. Wang and J. Q. Gan, “Prediction of the 2017 french election based on twitter data analysis using term weighting,” in *2018 10th Computer Science and Electronic Engineering (CEECE)*, pp. 231–235, 2018.
- [47] R. Sermsai and S. Laohakiat, “Analysis and prediction of temporal twitter popularity using dynamic time warping,” in *2019 16th International Joint Conference on Computer Science and Software Engineering (IJCSSSE)*, pp. 176–180, 2019.
- [48] G. Amoudi, “Popularity prediction in twitter during financial events,” in *2018 21st Saudi Computer Society National Computer Conference (NCC)*, pp. 1–6, 2018.
- [49] H. A. Al-Hussaini and H. Al-Dossari, “A lexicon-based approach to build service provider reputation from arabic tweets in twitter,” *International Journal of Advanced Computer Science and Applications*, vol. 8, 2017.
- [50] N. Abirami and K. Manohari, “Measuring user reputation on twitter using page rank algorithm,” 2018.
- [51] F. Riquelme, P. González-Cantergiani, X. Molinero, and M. Serna, “Centrality measure in social networks based on linear threshold model,” *Knowledge-Based Systems*, vol. 140, pp. 92–102, 01 2018.
- [52] A. Gün and P. Karagöz, “A hybrid approach for credibility detection in twitter,” in *Hybrid Artificial Intelligence Systems* (M. Polycarpou, A. C. P. L. F. de Carvalho, J.-S. Pan, M. Woźniak, H. Quintian, and E. Corchado, eds.), (Cham), pp. 515–526, Springer International Publishing, 2014.
- [53] H. Kwak, C. Lee, H. Park, and S. Moon, “What is twitter, a social network or a news media?,” vol. 19, 01 2010.
- [54] Y. Riyanto and J. Yeo, “Directed trust and trustworthiness in a social network: An experimental investigation,” *Journal of Economic Behavior Organization*, vol. 151, no. C, pp. 234–253, 2018.



- [55] X. Li, Y. Liu, Y. Jiang, and X. Liu, “Identifying social influence in complex networks: A novel conductance eigenvector centrality model,” *Neurocomputing*, vol. 210, 06 2016.
- [56] L. Jain, R. Katarya, and S. Sachdeva, “Opinion leader detection using whale optimization algorithm in online social network,” *Expert Syst. Appl.*, vol. 142, 2020.
- [57] Y. E. Riyanto and Y. X. Jonathan, “Directed trust and trustworthiness in a social network: An experimental investigation,” *Journal of Economic Behavior and Organization*, vol. 151, pp. 234–253, 2018.
- [58] M. Brede, “How does active participation affect consensus: Adaptive network model of opinion dynamics and influence maximizing rewiring,” *Complexity*, vol. 2019, 06 2019.
- [59] S. Zhang, Y. Cai, and H. Xia, “A privacy-preserving interactive messaging scheme based on users credibility over online social networks,” *2017 IEEE/CIC International Conference on Communications in China (ICCC)*, pp. 1–6, 2017.
- [60] J. Cetkovic, S. Lakić, M. Lazarevska, M. Žarković, S. Vujošević, J. Cvijović, and M. Gogić, “Assessment of the real estate market value in the european market by artificial neural networks application,” *Complexity*, vol. 2018, pp. 1–10, 01 2018.
- [61] M. Tsikerdekis and S. Zeadally, “Multiple account identity deception detection in social media using nonverbal behavior,” *IEEE Transactions on Information Forensics and Security*, vol. 9, no. 8, pp. 1311–1321, 2014.
- [62] F. Ahmad and S. Rizvi, *Identification of Credibility Content Measures for Twitter and Sina-Weibo Social Networks*, pp. 372–384. 09 2019.
- [63] S. A. Curiskis, B. Drake, T. R. Osborn, and P. J. Kennedy, “An evaluation of document clustering and topic modelling in two online social networks:

- Twitter and reddit,” *Information Processing Management*, vol. 57, no. 2, p. 102034, 2020.
- [64] M. Alrubaian, M. Al-Qurishi, M. M. Hassan, and A. Alamri, “A credibility analysis system for assessing information on twitter,” *IEEE Transactions on Dependable and Secure Computing*, vol. 15, no. 4, pp. 661–674, 2018.
- [65] S. A. Ríos, F. Aguilera, J. D. Nuñez-Gonzalez, and M. Graña, “Semantically enhanced network analysis for influencer identification in online social networks,” *Neurocomputing*, vol. 326-327, pp. 71–81, 2019.
- [66] Y. Liu and J. Cao, “Irank: A novel algorithm for identifying influencers in micro-blog social networks,” *2019 International Conference on Data Mining Workshops (ICDMW)*, pp. 735–740, 2019.
- [67] H. Kwak, C. Lee, H. Park, and S. Moon, “What is Twitter, a social network or a news media?,” in *WWW ’10: Proceedings of the 19th international conference on World wide web*, (New York, NY, USA), pp. 591–600, ACM, 2010.
- [68] M. Cha, H. Haddadi, F. Benevenuto, and K. P. Gummadi, “Measuring user influence in twitter: The million follower fallacy,” in *ICWSM*, 2010.
- [69] M. Luiten, W. A. Kusters, and F. W. Takes, “Topical influence on twitter : A feature construction approach,” 2012.
- [70] E. P. J. J. WENG, Jianshu; LIM and Q. HE, “Twitterrank: Finding topic-sensitive influential twitterers.,” *Proceedings of the Third ACM International Conference on Web Search Data Mining*, 2010.
- [71] B. A. Huberman, D. M. Romero, and F. Wu, “Social networks that matter: Twitter under the microscope,” *CoRR*, vol. abs/0812.1045, 2008.
- [72] R. Cappelletti and N. Sastry, “Iarank: Ranking users on twitter in near real-time, based on their information amplification potential,” in *International*

*Conference on Social Informatics (SocialInformatics)*, (Los Alamitos, CA, USA), pp. 70–77, IEEE Computer Society, dec 2012.

- [73] I. Anger and C. Kittl, “Measuring influence on twitter,” in *Proceedings of the 11th International Conference on Knowledge Management and Knowledge Technologies, i-KNOW '11*, (New York, NY, USA), Association for Computing Machinery, 2011.
- [74] E. Bakshy, J. M. Hofman, W. A. Mason, and D. J. Watts, “Everyone’s an influencer: Quantifying influence on twitter,” in *Proceedings of the Fourth ACM International Conference on Web Search and Data Mining, WSDM '11*, (New York, NY, USA), p. 65–74, Association for Computing Machinery, 2011.
- [75] M. Anjaria and R. R. Guddeti, “Influence factor based opinion mining of twitter data using supervised learning,” pp. 1–8, 01 2014.
- [76] C. Schenk and D. Sicker, “Finding event-specific influencers in dynamic social networks,” pp. 501–504, 10 2011.
- [77] Y. Mei, Y. Zhong, and J. Yang, “Finding and analyzing principal features for measuring user influence on twitter,” in *2015 IEEE First International Conference on Big Data Computing Service and Applications*, pp. 478–486, 2015.
- [78] X. Li, Y. Liu, Y. Jiang, and X. Liu, “Identifying social influence in complex networks: A novel conductance eigenvector centrality model,” *Neurocomputing*, vol. 210, pp. 141 – 154, 2016. SI:Behavior Analysis In SN.
- [79] E. Lahuerta-Otero and R. Cordero-Gutiérrez, “Looking for the perfect tweet. the use of data mining techniques to find influencers on twitter,” *Computers in Human Behavior*, vol. 64, pp. 575 – 583, 2016.

- [80] P. Sharma, A. Agarwal, and N. Sardana, “Extraction of influencers across twitter using credibility and trend analysis,” *2018 Eleventh International Conference on Contemporary Computing (IC3)*, pp. 1–3, 2018.
- [81] T. Huynh, I. Zelinka, X. H. Pham, and H. D. Nguyen, “Some measures to detect the influencer on social network based on information propagation,” in *Proceedings of the 9th International Conference on Web Intelligence, Mining and Semantics, WIMS2019*, (New York, NY, USA), Association for Computing Machinery, 2019.
- [82] X. Guo, H. Mirzaalian, E. Sabir, A. Jaiswal, and W. Abd-Almageed, “Cord19sts: Covid-19 semantic textual similarity dataset,” 2020.
- [83] M. Müller, M. Salathé, and P. E. Kummervold, “Covid-twitter-bert: A natural language processing model to analyse covid-19 content on twitter,” 2020.
- [84] M. Cinelli, W. Quattrocioni, A. Galeazzi, C. M. Valensise, E. Brugnoni, A. L. Schmidt, P. Zola, F. Zollo, and A. Scala, “The covid-19 social media infodemic,” 2020.
- [85] E. Ferrara, “What types of covid-19 conspiracies are populated by twitter bots?,” *First Monday*, May 2020.
- [86] K.-C. Yang, C. Torres-Lugo, and F. Menczer, “Prevalence of low-credibility information on twitter during the covid-19 outbreak,” 2020.
- [87] C. E. Lopez, M. Vasu, and C. Gallemore, “Understanding the perception of covid-19 policies by mining a multilanguage twitter dataset,” 2020.
- [88] T. Alshaabi, J. R. Minot, M. V. Arnold, J. L. Adams, D. R. Dewhurst, A. J. Reagan, R. Muhamad, C. M. Danforth, and P. S. Dodds, “How the world’s collective attention is being paid to a pandemic: Covid-19 related 1-gram time series for 24 languages on twitter,” 2020.

- [89] L. Singh, S. Bansal, L. Bode, C. Budak, G. Chi, K. Kawintiranon, C. Padden, R. Vanarsdall, E. Vraga, and Y. Wang, “A first look at covid-19 information and misinformation sharing on twitter,” 2020.
- [90] S. Jain and A. Sinha, “Identification of influential users on twitter: A novel weighted correlated influence measure for covid-19,” *Chaos, Solitons Fractals*, vol. 139, p. 110037, 2020.
- [91] S. Alqurashi, A. Alashaikh, and E. Alanazi, “Identifying information superspreaders of covid-19 from arabic tweets,” 07 2020.
- [92] A. Abd-alrazaq, D. Alhuwail, M. Househ, M. Hamdi, and Z. Shah, “Top concerns of tweeters during the covid-19 pandemic: Infoveillance study,” *Journal of Medical Internet Research*, vol. 22, 03 2020.
- [93] W. H. Dutton, “The fifth estate emerging through the network of networks,” *Prometheus*, vol. 27, no. 1, pp. 1–15, 2009.
- [94] B. Brereton, N. Hurley, and Dkit, *Review of Prell C. (2012) Social Network Analysis: history, theory and methodology Los Angeles, London, New Delhi, Singapore, Washington DC, Sage Publications Ltd.* 01 2014.
- [95] K. Freberg, K. Graham, K. Mcgaughey, and L. Freberg, “Who are the social media influencers? a study of public perceptions of personality,” *Fuel and Energy Abstracts*, vol. 37, pp. 90–92, 03 2011.
- [96] P. Dahlgren, “The internet, public spheres, and political communication: Dispersion and deliberation,” *Political Communication*, vol. 22, pp. 147–162, 04 2005.
- [97] Z. Li, “Psychological empowerment on social media: Who are the empowered users?,” *Public Relations Review*, vol. 42, no. 1, pp. 49–59, 2016.
- [98] M. D. Veirman, V. Cauberghe, and L. Hudders, “Marketing through instagram influencers: the impact of number of followers and product divergence

on brand attitude,” *International Journal of Advertising*, vol. 36, no. 5, pp. 798–828, 2017.

- [99] M. Kitsak, L. Gallos, S. Havlin, F. Liljeros, L. Muchnik, H. Stanley, and H. Makse, “Identification of influential spreaders in complex networks,” *Nature Physics*, vol. 6, 01 2010.
- [100] Z. Z. Alp and G. Öğüdücü, “Influential user detection on twitter: Analyzing effect of focus rate,” in *2016 IEEE/ACM International Conference on Advances in Social Networks Analysis and Mining (ASONAM)*, pp. 1321–1328, 2016.
- [101] tweepy, “tweepy api reference,” 2020. Available at <http://docs.tweepy.org/en/latest/api.html>.
- [102] twitter, “tweepy api reference.” Available at <https://developer.twitter.com/en/docs/api-reference-index>.