

Bachelor of Science in Computer Science & Engineering



**A Framework for Identifying Influential People by
Analyzing Social Media Data**

by

Md. Sabbir Al Ahsan

ID: 1504066

Department of Computer Science & Engineering
Chittagong University of Engineering & Technology (CUET)
Chattogram-4349, Bangladesh.

April, 2021

A Framework for Identifying Influential People by Analyzing Social Media Data



Submitted in partial fulfilment of the requirements for
Degree of Bachelor of Science
in Computer Science & Engineering

by

Md. Sabbir Al Ahsan

ID: 1504066

Supervised by

Dr. Mohammad Shamsul Arefin

Professor

Department of Computer Science & Engineering

Chittagong University of Engineering & Technology (CUET)

Chattogram-4349, Bangladesh.

The thesis titled ‘**A Framework for Identifying Influential People by Analyzing Social Media Data**’ submitted by ID: 1504066, Session 2019-2020 has been accepted as satisfactory in fulfilment of the requirement for the degree of Bachelor of Science in Computer Science & Engineering to be awarded by the Chittagong University of Engineering & Technology (CUET).

Board of Examiners

Chairman

Dr. Mohammad Shamsul Arefin

Professor

Department of Computer Science & Engineering

Chittagong University of Engineering & Technology (CUET)

Member (Ex-Officio)

Dr. Asaduzzaman

Professor & Head

Department of Computer Science & Engineering

Chittagong University of Engineering & Technology (CUET)

Member (External)

Dr. Mohammed Moshiul Hoque

Professor

Department of Computer Science & Engineering

Chittagong University of Engineering & Technology (CUET)

Declaration of Originality

This is to certify that I am the sole author of this thesis and that neither any part of this thesis nor the whole of the thesis has been submitted for a degree to any other institution.

I certify that, to the best of my knowledge, my thesis does not infringe upon anyone's copyright nor violate any proprietary rights and that any ideas, techniques, quotations, or any other material from the work of other people included in my thesis, published or otherwise, are fully acknowledged in accordance with the standard referencing practices. I am also aware that if any infringement of anyone's copyright is found, whether intentional or otherwise, I may be subject to legal and disciplinary action determined by Dept. of CSE, CUET.

I hereby assign every rights in the copyright of this thesis work to Dept. of CSE, CUET, who shall be the owner of the copyright of this work and any reproduction or use in any form or by any means whatsoever is prohibited without the consent of Dept. of CSE, CUET.

Signature of the candidate

Date:

Acknowledgements

The success and the outcome of this thesis required a lot of guidance and assistance from many people and I am extremely privileged to have got this through the completion of my thesis. It has been an enriching experience, professionally and personally. All that I have done is only due to such supervision and assistance. Above all, Thanks to almighty Allah for enabling me to complete this thesis successfully. Thereafter, I would like to express my deep gratitude to my honorable thesis supervisor Dr.Mohammad Shamsul Arefin, Professor, Department of Computer Science & Engineering, Chittagong University of Engineering & Technology (CUET) for his guidance, encouragement, and continuous support during my thesis work. I am thankful for his many critical questions, forcing me to see things from different perspectives, and his magnificent support throughout the entire time.

Finally, I want to thank my father and mother for their unconditional love, support, encouragement, and contribution throughout my life and academic career in every aspect over the years.

Abstract

In this paper, we introduce a new framework for identifying the most influential people from social sensor networks. Selecting influential people from social networks is a complicated task as it depends on many metrics like the network of friends, followers, reactions, comments, shares, etc. (e.g., friends-of-a-friend, friends-of-a-friend-of-a-friend). Data on social media are increasing day-by-day at an enormous rate. It is also a challenge to store and process these data. Towards this goal, we use Hadoop to store data and Apache Spark for the fast computation of the data. To select influential people, we apply the mechanisms of skyline query and top- k query. To the best of our knowledge, this is the first work to apply the Apache Spark framework to identify influential people on social sensor network, such as online social media. Our proposed mechanism can find influential people very quickly and efficiently on the data pattern of Facebook.

Keywords— social sensor network; influential person identification; Hadoop; Apache Spark; skyline query; top- k query

Table of Contents

Acknowledgements	iii
Abstract	iv
List of Figures	vii
List of Tables	viii
1 Introduction	1
1.1 Introduction	1
1.2 Framework Overview	1
1.3 Difficulties	2
1.4 Applications	2
1.5 Motivation	3
1.5.1 Problem Statement	3
1.6 Contribution of the thesis	3
1.7 Thesis Organization	4
1.8 Conclusion	4
2 Literature Review	5
2.1 Introduction	5
2.2 Related Literature Review	5
2.2.1 Work Related to Social Networks	5
2.2.2 Work Related to Apache Spark	6
2.2.3 Work Related to Skyline Query	6
2.2.4 Work Related to the Top- k Query	7
2.3 Conclusion	7
3 Methodology	8
3.1 Introduction	8
3.2 Overview of Framework	9
3.2.1 Start HDFS and Spark	10
3.2.2 Accessing HDFS using Python	10
3.2.3 Submitting Job to Spark	11
3.3 Detailed Explanation	11

3.3.1	Dataset Description	11
3.3.2	Performance Metrics	14
3.3.2.1	Friend Score (F_s)	14
3.3.2.2	Follower Score (Fl_s)	15
3.3.2.3	Reaction Score (R_s)	15
3.3.2.4	Comment Score (C_s)	15
3.3.2.5	Group Score (G_s)	15
3.3.2.6	Share Score (S_s)	16
3.3.2.7	Pages Liked Score (P_s)	16
3.3.2.8	Check-in Score (Ch_s)	16
3.3.2.9	Mention Score (M_s)	17
3.3.2.10	Event Score (E_s)	17
3.3.3	Table Construction	17
3.3.4	Influential People Selection	20
3.3.4.1	Skyline Query	20
3.3.4.2	Top- k Query	21
3.4	Conclusion	23
4	Results and Discussions	24
4.1	Introduction	24
4.2	Impact Analysis	24
4.2.1	Social and Environmental Impact	24
4.2.2	Ethical Impact	24
4.3	Evaluation of Performance	25
4.3.1	Experiment Setup	25
4.3.2	Performance on Score Calculation	25
4.3.3	Performance on Skyline Queries	26
4.3.4	Effect of New Metrics in Skyline Computation	27
4.3.5	Performance on Top- k Queries	27
4.4	Conclusion	28
5	Conclusion	29
5.1	Conclusion	29
5.2	Future Work	30

List of Figures

3.1	System architecture for finding influential people	9
3.2	Spark Architecture	10
3.3	Skyline computation.	22
4.1	Performance on score calculation.	26
4.2	Performance on skyline queries.	26
4.3	Effect of new metrics in skyline computation.	27
4.4	Performance on top- k queries.	28

List of Tables

3.1	Facebook data.	12
3.2	Group influence.	13
3.3	Page influence.	13
3.4	Check-in influence.	13
3.5	Event influence.	13
3.6	Symbols and their meanings.	14
3.7	Friends data set in key-value storage (kvs) format.	18
3.8	Groups dataset in kvs format.	18
3.9	Score calculation for each metric.	19
3.10	Candidate domain values.	21
3.11	Top- K query result.	23

Chapter 1

Introduction

1.1 Introduction

Currently, the importance of social media is beyond question. People can connect anywhere in the world at an instant due to the blessings of social media on the Internet [1]. Social media influence is an important term as influencers have a great impact on social media platforms. Every moment, an enormous amount of data is generated on social media like Facebook, Twitter, Instagram, etc. Statistics show that around four petabytes of data are generated per day on Facebook alone, and Apache Spark along with Hadoop can play a big role in processing this huge amount of data [2, 3].

1.2 Framework Overview

Considering these facts, we use Hadoop to store the data, as well as Apache Spark to process the data. The Apache Hadoop [4] software library is a framework that allows for the distributed processing of large data sets across clusters of computers using simple programming models. It was created by Doug Cutting in 2005 when he was working for Yahoo at the time for the Nutch search engine project. Hadoop has two major components named the HDFS (Hadoop Distributed File System) [5] and the MapReduce [6] framework.

Apache Spark was originally developed at UC Berkeley in 2009 in order to process a large amount of data (e.g., IoT streaming data or social media streaming data processing from multiple sources [7, 8]), and Apache Spark can serve as a unified analytics engine through machine learning [9] and deal with big data [10]. It can run ten times faster than on disk and a hundred times faster in memory than

Hadoop. Spark runs everywhere, can access diverse data sources, and has very powerful caching. Most importantly, it supports multiple platforms such as R, Scala, Python, and SQL shells [11]. In this paper, we use the Python language in our project to interact with Spark.

1.3 Difficulties

The major difficulties for operating Hadoop and Apache Spark given below:

1. Setting up Infrastructure

Apache Spark and Hadoop both works on parallel processing. To run program fluently we need a group of high configuration computers. So it is a big challenge to set up the infrastructure.

2. Managing Hadoop

Hadoop is hard to manage as well as it often requires extensive internal resources to maintain which is also difficult.

3. Scalability Challenge

As data is growing rapidly nowadays, the data storage and analytics demands are also going to be increased. In a consequence, more server would need which is expensive as well as time consuming.

1.4 Applications

Social Media influencers are keeping a great impact nowadays. They play a huge role in social both social and economical sectors. Social media influencers, ranging from well-known YouTubers to Twitter users with a wide following, are a perfect way to increase brand awareness and engagement, boost sales, and create leads for more potential clients. It's also a great way to use collaborations to promote your brand. So, People will get benefited a lot if they can efficiently find influential people from the huge database of social media.

1.5 Motivation

Social media data are growing exponentially currently. Social influencers have a great impact on these platforms. It is a difficult task to find influential people as this depends on many metrics. If we want to analyze data from this huge source of social media data, we will be faced with mainly two problems:

- the large data set storage problem and
- the large data set processing problem.

Therefore, we use Hadoop and the Apache Spark framework to solve these issues.

1.5.1 Problem Statement

In this paper, we consider the problem of selecting influential people from social networks. We consider the data pattern of Facebook as it is the most popular social media platform currently. We consider ten different metrics to select influential people. The metrics are friends, followers, reactions, comments, groups, shares, pages liked, check-ins, mentions, and events attended.

To select influential people, we apply skyline queries and top- k queries. In a dimensional data set, a point p dominates another point q only if it is better than or equal to q in all dimensions and better than q in at least one. The points that are not dominated by any other point are called the skyline. Top- k queries produce results that are ordered by some computed score. These queries retrieve the k best results according to the user's satisfaction.

1.6 Contribution of the thesis

The major contributions of this article are listed as follows:

- We develop a mechanism to identify influential people.
- We use an efficient method to store and process large data.
- We effectively apply the skyline query and top- k query to handle Facebook data.

1.7 Thesis Organization

The rest of the paper is organized as follows.

- In Chapter 2 a brief discussion on the previous work related to "selecting influential person on social media" is provided.
- In Chapter 3 we provide our methodology, describe our data set in detail, introduce the performance metrics, and then, show the selection procedure to find influential people.
- In Chapter 4 the experimental results and performance evaluations are presented.
- Finally, in Chapter 5 we conclude the paper and outline future research directions.

1.8 Conclusion

In this chapter, the overview of our work has been discussed. The difficulties, applications, and motivation of our work also have been discussed throughout this chapter. We also mentioned the contributions of our work. This chapter gives a summary of our work.

Chapter 2

Literature Review

2.1 Introduction

This section divides the state-of-the-art literature into four sections. We have reviewed the work related to social media in Section 2.2.1 and the work related to Apache Spark in Section 2.2.2. To retrieve the selected people, skyline queries and top- k queries are used. We have reviewed both of these queries in Sections 2.2.3 and 2.2.4.

2.2 Related Literature Review

2.2.1 Work Related to Social Networks

Asif Zaman et al. [12] addressed finding key people on social media using the MapReduce framework. They considered five metrics and applied the skyline query for selecting the key people. However, the MapReduce framework is not very time efficient, and the performance degrades when the data set is large. J. Qiu et al. [13] addressed social influence prediction with deep learning. They proposed a framework that takes a user's local network as the input to a graph neural network for learning his/her latent social representation. They showed that their proposed model significantly outperforms traditional feature engineering-based approaches. Cao et al. [14] addressed the jury selection problem for decision-making tasks on micro-blog services. To avoid the exponentially increasing calculation of the jury error rate, they introduced two efficient algorithms for choosing jury members. By evaluating a Twitter graph, they predicted the error rate of each user. DeMartini et al. [15] suggested a systematic model for search entities. They introduced a vector space model for ranking entities with the application

to find experts. T. Lappas et al. [16] addressed the problem of forming a team of skilled individuals to perform a given task. They mainly focused to minimizing the communication cost among team members.

2.2.2 Work Related to Apache Spark

M. Zaharia et al. [17] first introduced the architecture of Apache Spark with RDDs (resilient distributed datasets), parallel computing, etc. They implemented RDDs in Spark and evaluated a variety of user applications. Satis Gopalani et al. [18] gave a comparison of Apache Spark and the MapReduce framework. They analyzed the performance using the K-means clustering algorithm. X. Meng et al. [19] introduced MLib, which is a machine learning library for Apache Spark. Anand Gupta et al. [20] introduced a framework for big data analysis using Apache Spark and deep learning. For this analysis, they used a technique called cascade learning. L. R. Nair et al. [21] analyzed Twitter data using the Apache Spark framework and classified various job categories among millions of streaming tweets by analyzing the real-time data.

2.2.3 Work Related to Skyline Query

The skyline query problems are related to the maximum vector problems [22, 23]. Borzsonyi et al. [24] first addressed the skyline operator. They introduced three algorithms: Block Nested Loops (BNL), Divide and Conquer (D&C), and B-tree-based schemes. For BNL, every item is contrasted with each other item in the database, and the item is accounted for if no other item dominates. Chomicki et al. [25] addressed a variant of BNL by introducing an algorithm named Sort-Filter-Skyline (SFS) as an improved version of BNL. This algorithm requires a pre-sorting data set. Tan et al. [26] proposed “Bitmap” and “Index”, two progressive processing algorithms. In the Bitmap algorithm, every dimensional value of a point is represented by a few bits. It can be determined whether a given data point is in the skyline by applying a bit-wise AND operation on these vectors. In the Index algorithm, the entire data set is organized into some lists, and points are assigned to a list only if the value in the dimension is the

best among all the dimensions. Chan et al. [27] addressed the k -dominant skyline in high-dimensional space. They proposed efficient algorithms to solve the k -dominant skyline problem. Balke et al. [28] addressed the problem of skyline queries in web information systems. They vertically partitioned the data set and retrieved the skyline in a distributed environment.

2.2.4 Work Related to the Top- k Query

Several works have been done on top- k queries. Marian et al. [29] addressed an algorithm to efficiently evaluate top- k queries over web-accessible databases, and they assumed that only random access was available for web sources. Chang et al. [30] addressed the problem of evaluating ranked top- k queries with expensive predicates. They proposed an algorithm for minimizing probe accesses as much as possible. They assumed an arranged access on one of the traits, while different scores were obtained through testing or executing some user-defined function on the rest of the attributes. Li et al. [31] introduced a framework that gives a methodical and principled structure to help with efficient assessments of top- k queries in relational database systems by expanding the relational algebra. Chaudhuri et al. [32] studied the advantages, as well as the limitations of handling a top- k query by interpreting it in a single range query that conventional relational DBMSs (database management systems) can process efficiently. D. Donjerkovic et al. [33] proposed a probabilistic approach to query optimization.

2.3 Conclusion

In this chapter, a detailed literature review is discussed. We have given a brief description of previous work related to social networks, Apache Spark, skyline query, and top- k query.

Chapter 3

Methodology

3.1 Introduction

The proposed system architecture consists of four modules. They are the storage module, calculation module, selection module, and output module. Figure 3.1 depicts the system architecture of our proposed method. First, we have to store all our data sets in the HDFS (Hadoop Distributed File System), which is the storage module. We have to fetch all our data from here for all the calculations. In the Calculation module, first, we need to write our program for the score calculation. The program was written in Python and PySpark.

We first operate our program individually to generate each metric score in the PySpark shell. After testing all the code, we create an environment for submitting our entire code to Spark. Then, we submit the code to Spark. It will process the data, and the output is generated after processing. In the Selection module, the skyline query and top- k query are used to retrieve the influential people. After retrieving the output data sets, we have to store them again in the HDFS.

In the output module, the results can be easily moved into the HDFS storage. Then, we submit our job to Spark. It will process the data, and the output is generated. In the Selection module, the skyline query and top- k query are used to retrieve the influential people. The output module shows the results retrieved from both queries and stores them in the HDFS.

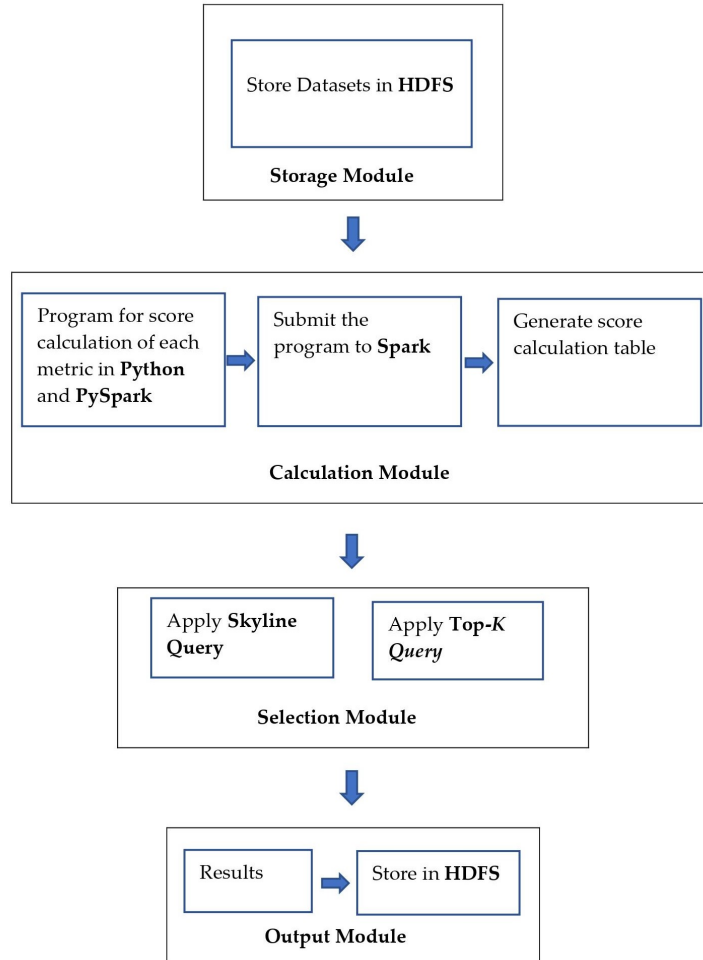


Figure 3.1: System architecture for finding influential people

3.2 Overview of Framework

The Spark follows the master-slave architecture. Its cluster consists of a single master and multiple slaves. The Spark architecture is depicted in fig. 3.2. The overview of how Spark works given below:

- The application is hosted in the driver first.
- Then Spark Context is initiated which is the gateway of Spark.
- Then the Cluster manager allocates resources to Driver program to run the tasks.
- The task is then partitioned and allocated to workers or executors.

3.2.1 Start HDFS and Spark

To store datasets into HDFS and before running our program we have to start HDFS and Apache Spark at first. To do that we have to start Namenode, Datanode, Master, and Worker. We can see if all of our nodes are running by the Jps command.

3.2.2 Accessing HDFS using Python

To access HDFS using python we need to follow some steps. We used the standard method for accessing HDFS. The steps are given below:

- Install HDFS package
- Define hdfscli.cfg in User home directory
- Import Config or Insecure Client in the python program
- Access HDFS

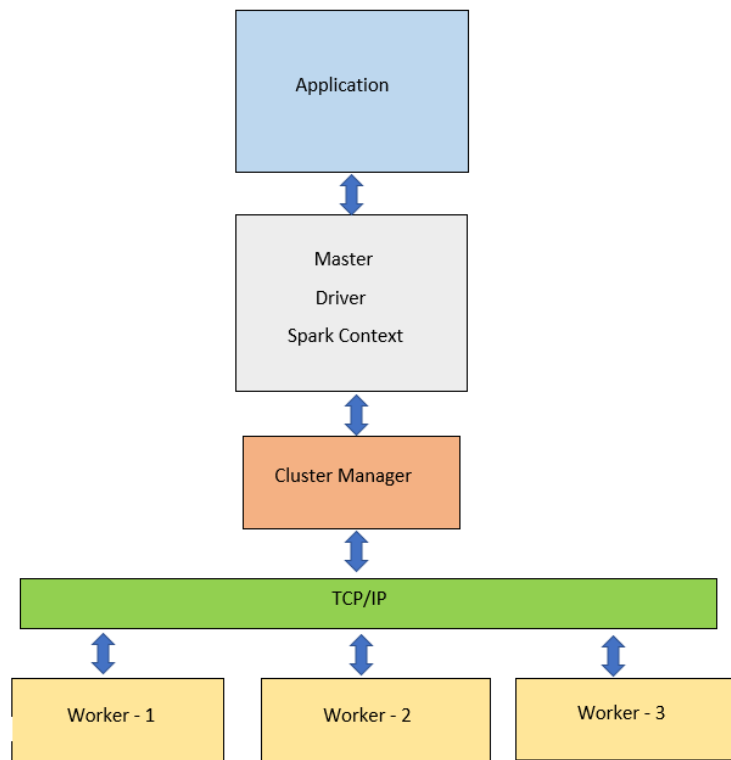


Figure 3.2: Spark Architecture

3.2.3 Submitting Job to Spark

After completing the program, we need to submit the job to the Spark. Then the driver converts the program into tasks. We used Standalone mode for our work.

3.3 Detailed Explanation

3.3.1 Dataset Description

The symbolic representation of Facebook data is shown in Table 3.1. To find influential people, we selected the 10 criteria given in the figure. Here, we show the symbol names for all metrics for simplicity. The User ID column shows the list of all users on Facebook. Here, we describe all the columns for User A from the User ID column. The Friends column demonstrates that User A has friends like U, V, X, etc., and followers like S, D, R, etc. In the Like column, (1, T) means that User T has given a reaction to the first status of User A. In the Comment column, (3, Z, Comment5) means that User Z has given a comment, which is Comment1 to the third status of User A. The Groups column shows that User A is connected to different groups like G1, G3, etc. As all the groups do not have the same influence, we assume that the group that is more active is more important. We use Table 3.2 to determine the activity of the group that we called influencers. In the Share column, (1, E) means User E has shared the first status or post of User A. The Pages liked column shows that User A has liked pages P1, P2, etc. The importance of all pages is not the same. To determine the page influence, we use Table 3.3. The Check-ins column shows that User A has checked-in different places named C1, C2, etc. A place is said to be more popular on Facebook depending on some criteria, which we show in Table 3.4. In the Mentioned by column, (B,6) means that User A was mentioned or was tagged by User B in his/her sixth status. The Events attended column shows that User A has attended various events named E1, E2, etc. However, the importance of all events is not the same. The metrics are shown to determine event influence in Table 3.5.

Table 3.1: Facebook data.

User ID	Friends	Followers	Reactions	Comments	Groups	Shares	Pages Liked	Check-ins	Mentioned by	Events Attended
A	U, V,	S, D,	(1, T),	(3, Z, Comment5),	G1,	(1,E),	P1,	C1,	(B, 6),	E1, E2,
	X, Z,	R, Y,	(3, I),	(6, E, Comment4),	G3,	(8, C),	P2,	C2,	(M, 7),	
	G5,	P3,	C5,	
	
B	C, M,	S, P,	(2, F),	(5, E, Comment9),	G3,	(2, F),	P8,	C3,	(E, 2),	E3, E5,
	L, S,	W, X,	(9, C),	(2, P, Comment3),	G7,	(3, D),	P5,	C5,	(C, 5),	
	G9,	P7,	C9,	
	
C	A, C,	M, J,	(16, I),	(7, A, Comment8),	G12,	(7, S),	P2,	C1,	(A, 12),	E4, E6,
	Y, U,	K, L,	(21, J),	(1, F, Comment2),	G15, G19,	(5, N),	P4,	C9,	(B, 17),	
	I,	P7,	C11,	
	

Table 3.2: Group influence.

Group	Total Members	Members Added the Last 7Days	Total Posts Last 30 Days
G1	1270	E,F,G..	200
G2	2100	D,E,F..	134
G3	1700	M,N,O..	305
...

Table 3.3: Page influence.

Page	Followers	Comments	Reviews
P1	A,B,C...	(5,Q,Comment1)...	(R,4)...
P2	A,C,E...	(2,P,Comment3)...	(S,5)...
P3	L,M,P...	(1,T,Comment9)...	(U,3)...
....

Table 3.4: Check-in influence.

Place	Followers	Comments	Reviews
C1	B,C,D...	(9,S,Comment2)...	(Q,5)...
C2	W,E,R...	(7,E,Comment1)...	(A,3)...
C3	T,Y,U...	(3,Q,Comment5)...	(C,5)...
....

Table 3.5: Event influence.

Event	Going	Interested
E1	180	350
E2	305	751
E3	200	500
...

Table 3.2 consists of the group name, total members, members added the last seven days, and total posts the last 30 days. We consider the more active group to be more important. The group activity is given in Table 3.2. Using these three metrics, we can calculate the group influence on Facebook. Table 3.3 consists of the page name, followers, comments, and reviews. On every page, there is a review section where people can give a review of that page like a four-star review, five-star review, etc. We store this in the Review column. We assume a Facebook page to be influential considering all these metrics in Table 3.3. Here, the Comment column has the same format as described in Table 3.1.

In the Reviews column, (R,4) means that User 4 has given a four-star rating to that particular page. Table 3.4 is for determining the check-in influence. It consists of the name where a user checked-in, the followers, the comment section, and the review section of that place. The structure of Tables 3.3 and 3.4 is the same. Table 3.5 consists of the event, going, and interested column. A popular event always gets a huge response from people. Therefore, we consider how many people are interested in that event, as well as how many people are going to that event to determine the event’s popularity.

3.3.2 Performance Metrics

In this section, we give the basic definitions, as well as the formulas for social network metrics that are used in this work. As we considered a total of 10 metrics to find influential people, we assume that each metric’s data set is in kvs (key-value storage) format. Then, we calculate the metric score for every metric separately in the PySpark [34] shell. Table 3.6 shows the considered metrics with their symbols.

Table 3.6: Symbols and their meanings.

Considered Metrics	Symbols
Friend Score	F_s
Follower Score	Fl_s
Reaction Score	R_s
Comment Score	C_s
Group Score	G_s
Share Score	S_s
Pages Liked Score	P_s
Check-in Score	Ch_s
Mention Score	M_s
Event Score	E_s

3.3.2.1 Friend Score (F_s)

A friend score is the total number of friends a user has on social media [12]. It is defined as:

$$F_s = \text{Total number of friends} \quad (3.1)$$

3.3.2.2 Follower Score (Fl_s)

A follower score is the total number of friends a user has on social media [12]. It is defined as:

$$Fl_s = \text{Total number of followers} \quad (3.2)$$

3.3.2.3 Reaction Score (R_s)

The term reaction refers to the like, sad, happy, angry, etc, reaction that a user has to someone's post. The reaction score is the average number of the "reaction" count that a user has achieved. To determine the reaction score, we define the following rules:

$$R_s = (\text{Total reaction received from all posts} / \text{Total posts}) \quad (3.3)$$

3.3.2.4 Comment Score (C_s)

In every post, many relevant and irrelevant comments are made by users. To determine the comment score of a user, a comment bias function [12] is used to determine the polarity of the comment. The function receives comments and returns a value for each comment. We give +1 for the positive comments, -1 for negative comments, and 0 for neutral comments. For simplicity, we considered all complicated comments as neutral comments. After determining each comment bias, we replace each comment name with its corresponding value. Then, we use the formula:

$$C_s = (\text{Total comment bias} / \text{Total posts}) \quad (3.4)$$

There are mainly three types of sentiments that can be classified: positive, negative, and neutral. There are many features that are often used in this problem [35] like terms and their frequency, parts of speech, opinion words and phrases, rules of opinions, etc. We used opinion words and phrases to detect comment bias. Words that are commonly used to express positive or negative sentiments are called opinion words.

3.3.2.5 Group Score (G_s)

There are many groups in which a user is involved. To obtain a group score, we first calculate the group influence, which is depicted in Table 3.2. We assume a group to be more active, as well as influential by three criteria: total members, members added the last 30 days, and total posts the last 30 days. Group influence

is the summation of all these metrics for each group. It is defined as follows:

$$\text{Group influence} = (\text{Total members} + \text{members added the last 30 days} + \text{total posts the last 30 days}) \quad (3.5)$$

The group score is the summation of all the group influence of a user that belongs to the social network. It is defined as:

$$G_s = \text{Total group influence} \quad (3.6)$$

3.3.2.6 Share Score (S_s)

A Facebook user's post is shared by many people according to the quality of that posts. The share score is the average number of "share" counts of that user. It is defined as:

$$S_s = (\text{Total shares} / \text{Total posts}) \quad (3.7)$$

3.3.2.7 Pages Liked Score (P_s)

To obtain a page score, we first calculate the page influence, which is depicted in Table 3.3. We assume a page to be more influential based on total followers, the comment score, and the review score. Total followers are the total number of followers following a page. To determine the comment score, we use the same method as discussed earlier. To determine the reviews, we use the following formula:

$$\text{Total reviews} = (\text{Total review score} / \text{Total reviewers}) \quad (3.8)$$

Adding all the measures, we obtain the page influence. It is defined as:

$$\text{Page influence} = (\text{Total followers} + \text{comment score} + \text{review score}) \quad (3.9)$$

The pages liked score is the summation of all users that have liked the pages. It is defined as:

$$P_s = \text{Total page influence} \quad (3.10)$$

3.3.2.8 Check-in Score (Ch_s)

A user checks into different places on Facebook. To determine the importance of these places, we consider the follower score, comment score, and review score of that page and call it check-in influence. We calculate the follower score, comment

score, and review score by the same method used previously. It is defined as:

$$\text{Check-in influence} = (\text{follower score} + \text{comment score} + \text{review score}) \quad (3.11)$$

The check-in score is the summation of all the check-in influence of a user that belongs to the social network. It is defined as:

$$Ch_s = \text{Total check-in influence} \quad (3.12)$$

3.3.2.9 Mention Score (M_s)

A Facebook user is mentioned by several users in their posts, photos, etc. The mention score is the number of times in total a user ID is get mentioned or tagged by another user. It is defined as:

$$M_s = \text{Total get mentioned} \quad (3.13)$$

3.3.2.10 Event Score (E_s)

A person attends several events, and to determine event importance, we assume two criteria: how many people are going and how many people are interested in that event. It is defined as:

$$\text{Event influence} = (\text{Total going} + \text{Total interested}) \quad (3.14)$$

The event score is the summation of the event influence of a user of the social network. It is defined as:

$$E_s = \text{Total event influence} \quad (3.15)$$

3.3.3 Table Construction

As we mentioned before, all our data set is in kvs (key-value storage) format. We first run all our data sets individually according to the formula given in Section ?? to obtain the individual score of each data set, like the friend score, follower score, reaction score, comment score, group score, share score, pages liked score, check-in score, mention score, and event score. However, for every data set that has a comment column, we use the CommentBias() [12] function to determine the

polarity and replace the value with the corresponding comment. Then, we follow our algorithms for further processing.

We show the friend data set and group data set in kvs format in Tables 3.7 and 3.8. The rest of the data sets also follows the same format. To avoid redundancy, we show here only two data sets.

By using Algorithm 1, the friend score, follower score, reaction score, comment score, share score, and mention score are calculated. However, we shown the friend score calculation here. Other metric scores are calculated by following the same procedure described in Algorithm 1. By using Algorithm 2, the group score, pages liked score, check-in score, and event score are calculated. We shown only the group score calculation here. Other metric scores are calculated by following the same procedure described in Algorithm 2.

The constructed table is shown in Table 3.9. We use this table for further calculations to find the influential people.

Table 3.7: Friends data set in key-value storage (kvs) format.

User ID	Friends
A	U,V,X,Z,...
B	C,M,L,S,...
C	A,C,Y,U,...
D	M,K,S,O,...
....

Table 3.8: Groups dataset in kvs format.

User ID	Groups
A	G1,G3,G5,...
B	G3,G7,G9,...
C	G12,G15,G19,...
D	G2,G4,G6,...
....

Algorithm 1 Friend score calculation.

```
1: procedure FRIEND SCORE
2:   initialize Spark context
3:   load data set into a data frame  $df$ 
4:   createOrReplaceTempView of  $df$ 
5:   perform SQL operation for score calculation
6:   store result in  $df\_friends$ 
7: end procedure
```

Algorithm 2 Group score calculation.

```
1: procedure GROUP SCORE
2:   initialize Spark context
3:   load group data set into a data frame  $df\_1$ 
4:   load group influence data set into data frame  $df\_2$ 
5:   createOrReplaceTempView of  $df\_2$ 
6:   perform SQL operation for score calculation, and store in  $df\_3$ 
7:   inner join  $df\_1$  and  $df\_3$ 
8:   store result in  $df\_groups$ 
9: end procedure
```

Table 3.9: Score calculation for each metric.

User ID	F_s	Fl_s	R_s	C_s	G_s	P_s	S_s	Ch_s	M_s	E_s
A	1002	200	103	200	302	409	70	54	19	7
B	905	125	95	120	504	206	90	65	12	9
C	507	50	53	75	210	115	105	43	20	5
D	1506	500	150	240	270	101	320	29	45	11
.....

3.3.4 Influential People Selection

Selecting influential people from the huge database of social media is also a difficult task. There are many efficient methods for selecting influential people. We have used the skyline query and top- k query for our desired people selection. The most significant disadvantage of using skyline query is that it can yield either too few result or too large result collection. When the result set is too small, the user will not benefit from the computation because the resulting data set will already be filled or not be willing to serve. When the result set is too big, it may make it difficult for the result seeker to make a decision.

3.3.4.1 Skyline Query

We first perform descending sort on the user ID according to the friend score, follower score, reaction score, comment score, group score, share score, pages liked score, check-in score, mention score, and event score. Then, we join all the tables. Then, we select the candidate list from among these sorted user IDs. As we have to select influential people based on ten criteria, an influential person ten opportunities to be selected as the best person. Therefore, for candidate selection, we programmed a Python function that maintains a counter for each user and appends a list, and if a user is retrieved ten times, then the calculation stops. Thus, we will find our candidate list. The total procedure of candidate selection is given in [12].

Here, the list returns S,A,B,Y,E,N,X,C,Q as the candidate list. Each candidate list contains separate domain values, which are given in Table 3.10. We have to perform a dominance test to retrieve the influential people. The dominance test is performed using a simple comparison algorithm. After the dominance test, we get our desired influential people. The procedure of the skyline query is shown graphically in Figure 3.3.

Table 3.10: Candidate domain values.

User ID	F_s	Fl_s	R_s	C_s	G_s	P_s	S_s	Ch_s	M_s	E_s
S	701	305	80	107	192	446	92	60	28	11
A	1002	200	103	200	302	409	70	54	19	14
B	905	125	95	120	504	206	90	65	12	9
Y	1501	550	152	243	267	103	229	49	55	15
E	1056	210	109	254	304	550	96	65	67	12
N	654	120	19	63	113	90	65	38	4	1
X	700	99	40	45	101	112	70	56	8	2
C	507	50	53	75	210	115	105	43	20	5
Q	877	376	40	83	120	154	68	44	29	6

3.3.4.2 Top- k Query

There are many methods to run top- k queries efficiently. The top- k query aims to retrieve only the k best results from a large data set. We use the naive method as Spark uses the distributed computation technique, and the naive method gives efficient results in the Spark framework. We can use SQL queries very easily in PySpark.

The steps of the naive approach of top- k query are given below.

1. Compute the total score by adding all the attribute values (given in Table 3.9) using the SQL query in PySpark.
2. Sort all user IDs based on their scores in descending order.
3. Return the first k highest scored user IDs.

We set the limit to 10 users, so we got the top 10 influential people. The output of the top- k query is shown in Table 3.11.

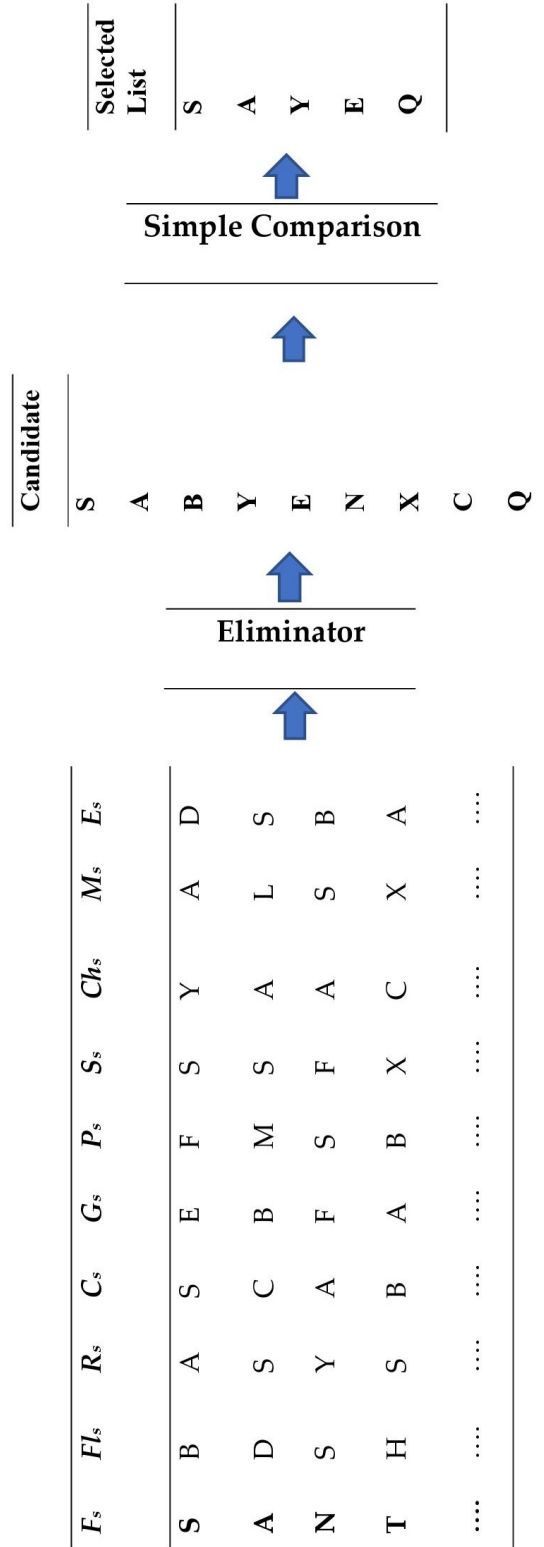


Figure 3.3: Skyline computation.

Table 3.11: Top- K query result.

User ID	Influential Score
S	9006
G	8987
J	8970
H	8901
L	8890
T	8889
K	8850
M	8820
Z	8809
X	8805

3.4 Conclusion

In this chapter, we have discussed the overall methodology of our proposed framework. We have discussed briefly the working procedure of Hadoop as well as Spark. We have discussed all the performance metrics to select the prominent people from social media. We have also depicted the dataset pattern in this section. Selecting influential people from the huge source of the database is also a difficult task. We have used skyline query and top- k query which is also briefly described in this section.

Chapter 4

Results and Discussions

4.1 Introduction

In this chapter, we will evaluate the results and performance of our proposed framework.

4.2 Impact Analysis

The impact analysis has divided into two parts and they are briefly discussed in section 4.2.1 and section 4.2.2.

4.2.1 Social and Environmental Impact

Influencers on social media now keeping a great impact in both social and economic spheres. They are a great way to boost brand awareness, increase sales, and create leads for more potential clients. It's also an excellent way to promote the brand through collaborations. People would greatly benefit if they could quickly recognize influential individuals from the massive social media database.

4.2.2 Ethical Impact

There are some ethical impacts as one can select a social media influencer either for good deeds or bad deeds.

4.3 Evaluation of Performance

This section provides the experimental settings, as well as evaluates the performance of our methodology.

4.3.1 Experiment Setup

We set up hadoop cluster with given specifications:

1. Number of commodity PCs: 4
2. Memory of each PC: 8
3. Operating System: CentOS 7
4. Hadoop version: 2.7.1
5. Python version: 3.5.2
6. Java version: JDK 7u79
7. Apache Spark version: 2.0.0
8. The replication factor of the Hadoop configuration: 2

4.3.2 Performance on Score Calculation

The score calculation table is shown in Table 3.6. We compared this result with the Hadoop MapReduce framework. There was a little time variation after submitting the job every time. Therefore, we repeated our job five times and took the average result. For automatic repetition, we used a crontab [36] schedule. The result is shown in Figure 4.1. We see that there is a constant time delay with respect to the total data volume from 50K to 500K. However, if the data volume becomes larger, there would be a slight increase in the processing time.

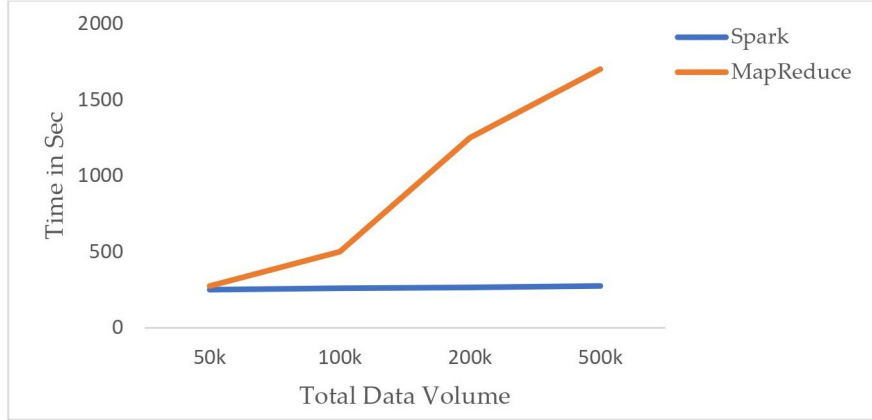


Figure 4.1: Performance on score calculation.

4.3.3 Performance on Skyline Queries

We used the naive approach of the skyline query. Generally, the naive approach of the skyline query takes too much time. However, in the Spark framework, we did not face such a problem. We also compared our work with the MapReduce framework. However, the time difference was huge. This is shown in Figure 4.2. Here, we also see that the time delay remains constant for the total data volume from 50K to 500K. The processing time would also vary here when the data volume increases.

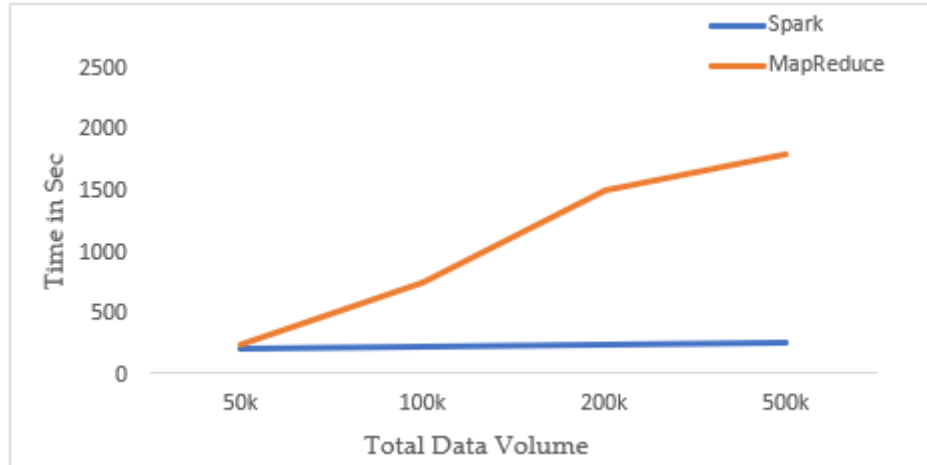


Figure 4.2: Performance on skyline queries.

4.3.4 Effect of New Metrics in Skyline Computation

The main problem of skyline computation is that either it may give a very small result set or a very large result set. It will create a great problem for a user if this condition occurs. Growing the size of social media matrices can limit the chance of experiencing such an issue. In [12], they used 5 matrices. As we use 10 metrics, we obtained more precise results, which are shown graphically in Figure 4.3.

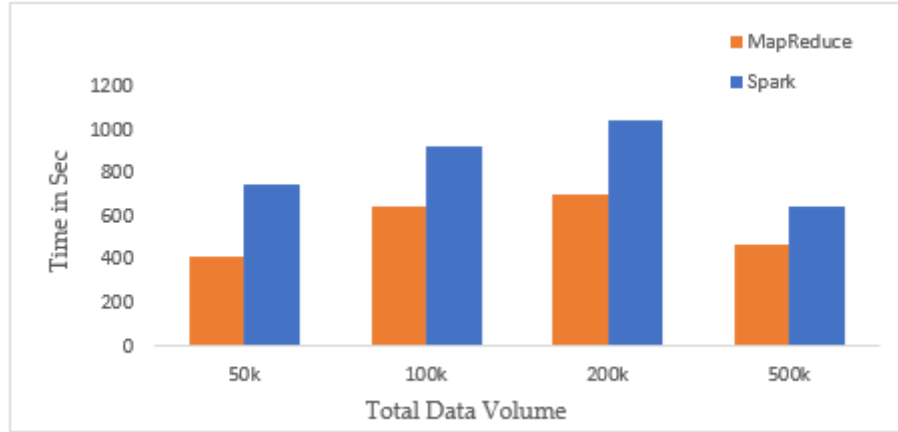


Figure 4.3: Effect of new metrics in skyline computation.

4.3.5 Performance on Top- k Queries

The naive approach of top- k queries generally takes much time. However, within the Spark framework, it gives time-efficient results. We have compared it with the MapReduce framework. We can see that Spark is far better than MapReduce. The time comparison is shown in Figure 4.4.

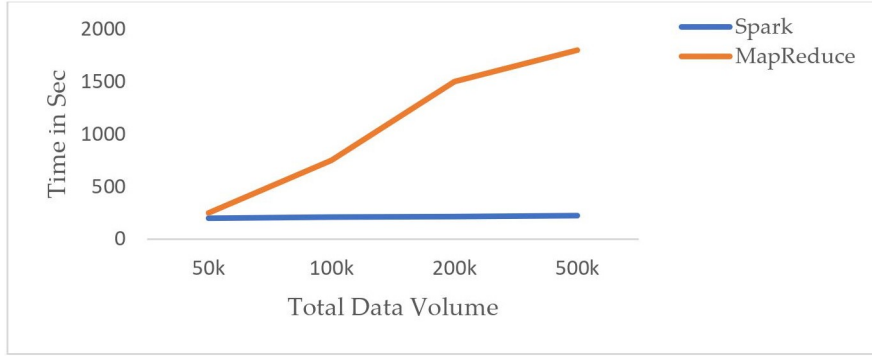


Figure 4.4: Performance on top- k queries.

4.4 Conclusion

In this chapter, the result and analysis of our proposed framework have been shown. We have seen from the graphs that, the Spark framework has outscored the MapReduce framework in all the sectors discussed above.

Chapter 5

Conclusion

5.1 Conclusion

As the nature of social networks is becoming more complicated day by day, social media data are also increasing at an exponential pace. Influencers on social media are currently receiving even more exposure as well. Nowadays, social media influencers have a significant effect. They play a significant role in both the social and economic sectors. Influencers on social media, ranging from well-known YouTubers to Twitter users with a large following, are an excellent way to raise brand awareness, increase revenue, and generate leads for more potential clients. It's also a perfect way to market your brand through partnerships. People would benefit greatly if they can easily identify powerful people from the vast social media database. In this article, as we consider Facebook in our work, we suggest an effective method of identifying prominent individuals from the social network. Our method can obtain data from both skyline queries and top- k queries efficiently. The key issue with skyline computation is that it can produce either a very small or a very large result collection. If this condition happens, it can cause a significant problem for the customer. Increasing the size of social media matrices will reduce the likelihood of such a problem occurring. So we use top- k queries along with skyline queries. We use Hadoop as the storage to manage the big data and Apache Spark to process the data. The algorithms used for calculating the score of different metrics and submitting the Spark job are simple, easy to implement, and very fast in obtaining the result as well. The effectiveness of our calculation is also demonstrated by multiple experiments. In this work, we use synthetic data considering the data patterns of Facebook.

5.2 Future Work

In the future, we plan to use Twitter data and data from other social networking sites to further verify the performance and effectiveness of our proposed system. In addition, this work can be further expanded by adding more parameters to get more precise results.

References

- [1] S. Mahbub, E. Pardede, A. Kayes and W. Rahayu, ‘Controlling astroturfing on the internet: A survey on detection techniques and research challenges,’ *International Journal of Web and Grid Services*, vol. 15, no. 2, pp. 139–158, 2019 (cit. on p. 1).
- [2] L. Lu, H. Dong, C. Yang and L. Wan, ‘A novel mass data processing framework based on hadoop for electrical power monitoring system,’ in *2012 Asia-Pacific Power and Energy Engineering Conference*, IEEE, 2012, pp. 1–4 (cit. on p. 1).
- [3] D. Q. Tu, A. Kayes, W. Rahayu and K. Nguyen, ‘Iot streaming data integration from multiple sources,’ *Computing*, pp. 1–31, 2020 (cit. on p. 1).
- [4] *Apache hadoop*, <https://hadoop.apache.org/> (accessed on 09 September 2020) (cit. on p. 1).
- [5] K. Shvachko, H. Kuang, S. Radia and R. Chansler, ‘The Hadoop distributed file system,’ in *2010 IEEE 26th Symposium on Mass Storage Systems and Technologies, MSST2010*, 2010, ISBN: 9781424471539. DOI: 10.1109/MSST.2010.5496972 (cit. on p. 1).
- [6] J. Dean and S. Ghemawat, ‘MapReduce: Simplified data processing on large clusters,’ in *OSDI 2004 - 6th Symposium on Operating Systems Design and Implementation*, 2004, pp. 137–149. DOI: 10.21276/ijre.2018.5.5.4 (cit. on p. 1).
- [7] D. Q. Tu, A. Kayes, W. Rahayu and K. Nguyen, ‘Isdi: A new window-based framework for integrating iot streaming data from multiple sources,’ in *International Conference on Advanced Information Networking and Applications*, Springer, 2019, pp. 498–511 (cit. on p. 1).
- [8] Q.-T. Doan, A. Kayes, W. Rahayu and K. Nguyen, ‘Integration of iot streaming data with efficient indexing and storage optimization,’ *IEEE Access*, vol. 8, pp. 47 456–47 467, 2020 (cit. on p. 1).
- [9] I. H. Sarker, A. Kayes, S. Badsha, H. Alqahtani, P. Watters and A. Ng, ‘Cybersecurity data science: An overview from machine learning perspective,’ *Journal of Big Data*, vol. 7, no. 1, pp. 1–29, 2020 (cit. on p. 1).

- [10] *Apache sparkTM - what is spark*, [Online]. Available: <https://databricks.com/spark/about> (accessed on 15 July 2020) (cit. on p. 1).
- [11] *Apache sparkTM - unified analytics engine for big data*, [Online]. Available: <https://spark.apache.org/> (accessed on: 7 July 2020) (cit. on p. 2).
- [12] A. Zaman, Md. Anisuzzaman Siddique, Annisa and Y. Morimoto, ‘Finding Key Persons on Social Media by Using MapReduce Skyline,’ *International Journal of Networking and Computing*, vol. 7, no. 1, pp. 86–104, 2017, ISSN: 2185-2839. DOI: 10.15803/ijnc.7.1_86 (cit. on pp. 5, 14, 15, 17, 20, 27).
- [13] J. Qiu, J. Tang, H. Ma, Y. Dong, K. Wang and J. Tang, ‘DeepInf: Social influence prediction with deep learning,’ in *Proceedings of the ACM SIGKDD International Conference on Knowledge Discovery and Data Mining*, 2018, pp. 2110–2119, ISBN: 9781450355520. DOI: 10.1145/3219819.3220077. arXiv: 1807.05560 (cit. on p. 5).
- [14] C. C. Cao, J. Shej, Y. Tong and L. Chen, ‘Whom to ask? Jury selection for decision making tasks on micro-blog services,’ *Proceedings of the VLDB Endowment*, vol. 5, no. 11, pp. 1495–1506, 2012, ISSN: 21508097. DOI: 10.14778/2350229.2350264. arXiv: 1208.0273 (cit. on p. 5).
- [15] G. Demartini, J. Gaugaz and W. Nejdl, ‘A vector space model for ranking entities and its application to expert search,’ in *Lecture Notes in Computer Science (including subseries Lecture Notes in Artificial Intelligence and Lecture Notes in Bioinformatics)*, vol. 5478 LNCS, 2009, pp. 189–201, ISBN: 3642009573. DOI: 10.1007/978-3-642-00958-7_19 (cit. on p. 5).
- [16] T. Lappas, K. Liu and E. Terzi, ‘Finding a team of experts in social networks,’ in *Proceedings of the ACM SIGKDD International Conference on Knowledge Discovery and Data Mining*, 2009, pp. 467–475, ISBN: 9781605584959. DOI: 10.1145/1557019.1557074 (cit. on p. 6).
- [17] M. Zaharia, M. Chowdhury, M. J. Franklin, S. Shenker and I. Stoica, ‘Spark: Cluster computing with working sets,’ in *2nd USENIX Workshop on Hot Topics in Cloud Computing, HotCloud 2010*, 2010 (cit. on p. 6).
- [18] S. Gopalani and R. Arora, ‘Comparing Apache Spark and Map Reduce with Performance Analysis using K-Means,’ *International Journal of Computer Applications*, vol. 113, no. 1, pp. 8–11, 2015. DOI: 10.5120/19788-0531 (cit. on p. 6).

- [19] X. Meng, J. Bradley, B. Yavuz, E. Sparks, S. Venkataraman, D. Liu, J. Freeman, D. B. Tsai, M. Amde, S. Owen, D. Xin, R. Xin, M. J. Franklin, R. Zadeh, M. Zaharia and A. Talwalkar, ‘MLlib: Machine learning in Apache Spark,’ *Journal of Machine Learning Research*, vol. 17, 2016, ISSN: 15337928 (cit. on p. 6).
- [20] A. Gupta, H. K. Thakur, R. Shrivastava, P. Kumar and S. Nag, ‘A Big Data Analysis Framework Using Apache Spark and Deep Learning,’ *IEEE International Conference on Data Mining Workshops, ICDMW*, vol. 2017-November, no. 1, pp. 9–16, 2017, ISSN: 23759259. DOI: 10.1109/ICDMW.2017.9. arXiv: 1711.09279 (cit. on p. 6).
- [21] L. R. Nair and S. D. Shetty, ‘Streaming twitter data analysis using spark for effective job search,’ *Journal of Theoretical and Applied Information Technology*, vol. 80, no. 2, pp. 349–353, 2015, ISSN: 18173195 (cit. on p. 6).
- [22] J. L. Bentley, H. T. Kung, M. Schkolnick and C. D. Thompson, ‘On the Average Number of Maxima in a Set of Vectors and Applications,’ *Journal of the ACM (JACM)*, vol. 25, no. 4, pp. 536–543, 1978, ISSN: 1557735X. DOI: 10.1145/322092.322095 (cit. on p. 6).
- [23] H. T. Kung, F. Luccio and F. P. Preparata, ‘On Finding the Maxima of a Set of Vectors,’ *Journal of the ACM (JACM)*, vol. 22, no. 4, pp. 469–476, 1975, ISSN: 1557735X. DOI: 10.1145/321906.321910 (cit. on p. 6).
- [24] S. Borzsonyil and K. Stocker, ‘SkylineOperator,’ *Ieee*, pp. 421–430, 2001 (cit. on p. 6).
- [25] K. Street, N. York and O. M. J. Canada, ‘Jan Chomicki Jarek Gryz Dongming Liang October 2002 Department of Computer Science,’ no. October, 2002 (cit. on p. 6).
- [26] K. L. Tan, P. K. Eng and B. C. Ooi, ‘Efficient progressive skyline computation,’ *VLDB 2001 - Proceedings of 27th International Conference on Very Large Data Bases*, pp. 301–310, 2001 (cit. on p. 6).
- [27] C. Y. Chan, H. V. Jagadish, K. L. Tan, A. K. Tung and Z. Zhang, ‘Finding k-dominant skylines in high dimensional space,’ *Proceedings of the ACM SIGMOD International Conference on Management of Data*, pp. 503–514, 2006, ISSN: 07308078. DOI: 10.1145/1142473.1142530 (cit. on p. 7).
- [28] W. T. Balke, U. Güntzer and J. X. Zheng, ‘Efficient distributed skylining for web information systems,’ *Lecture Notes in Computer Science (including*

subseries Lecture Notes in Artificial Intelligence and Lecture Notes in Bioinformatics), vol. 2992, pp. 256–273, 2004, ISSN: 16113349. DOI: 10.1007/978-3-540-24741-8_16 (cit. on p. 7).

- [29] A. Marian, N. Bruno and L. Gravano, *Evaluating top-k queries over web-accessible databases*, 2004. DOI: 10.1145/1005566.1005569 (cit. on p. 7).
- [30] K. C. C. Chang and S. W. Hwang, ‘Minimal probing: Supporting expensive predicates for top-K queries,’ *Proceedings of the ACM SIGMOD International Conference on Management of Data*, pp. 346–357, 2002, ISSN: 07308078 (cit. on p. 7).
- [31] C. Li, K. C. C. Chang, I. F. Ilyas and S. Song, ‘RankSQL: Query algebra and optimization for relational top-k queries,’ *Proceedings of the ACM SIGMOD International Conference on Management of Data*, pp. 131–142, 2005, ISSN: 07308078 (cit. on p. 7).
- [32] S. Chaudhuri and L. Gravano, ‘Evaluating Top-k Selection Queries,’ *Vldb*, vol. 1, pp. 397–410, 1999 (cit. on p. 7).
- [33] D. Donjerkovic and R. Ramakrishnan, ‘Probabilistic Optimization of Top N Queries,’ *International Conference on Very Large Databases (VLDB)*, vol. 1, pp. 411–422, 1999 (cit. on p. 7).
- [34] *Pyspark*, [Online]. Available: <https://spark.apache.org/docs/2.1.0/api/python/pyspark.html> (accessed on 18 July 2020) (cit. on p. 14).
- [35] B. Pang and L. Lee, ‘Opinion mining and sentiment analysis,’ in *Foundations and Trends in Information Retrieval*, vol. 2, 2008, pp. 1–135, ISBN: 9789380544199. DOI: 10.1561/15000000011 (cit. on p. 15).
- [36] *Cron job*, [Online]. Available: <https://en.wikipedia.org/wiki/Cron> (accessed on 1 October 2020) (cit. on p. 25).