PREDICTING INFLUENCERS FROM TWEETS DATA

NOREZZATI BINTI MD. NOR HAYATI

UNIVERSITI TEKNIKAL MALAYSIA MELAKA

C Universiti Teknikal Malaysia Melaka

BORANG PENGESAHAN STATUS TESIS

JUDUL: PREDICTING INFLUENCERS FROM TWEETS DATA

SESI PENGAJIAN: <u>2014/2015</u>

Saya NOREZZATI BINTI MD. NOR HAYATI (HURUF BESAR)

mengaku membenarkan tesis (PSM/Sarjana/Doktor Falsafah) ini disimpan di Perpustakaan Fakulti Teknologi Maklumat dan Komunikasi dengan syarat-syarat kegunaan seperti berikut:

- 1. Tesis dan projek adalah hak milik Universiti Teknikal Malaysia Melaka.
- 2. Perpustakaan Fakulti Teknologi Maklumat dan Komunikasi dibenarkan membuat salinan untuk tujuan pengajian sahaja.
- 3. Perpustakaan Fakulti Teknologi Maklumat dan Komunikasi dibenarkan membuat salinan tesis ini sebagai bahan pertukaran antara institusi pengajian tinggi.
- 4. ****** Sila tandakan (/)

	SULIT	(Mengandungi	maklumat	yang
		berdarjah kesela	imatan atau keper	ıtingan
		Malaysia seper	ti yang termak	tub di
		dalam AKTA R	AHSIA RASMI	1972)
	TERHAD	(Mengandungi yang telah organisasi/badan dijalankan)	maklumat TE ditentukan 1 di mana penye	RHAD oleh lidikan
	TIDAK TERHAD			
(TANDATANGAN	PENULIS)	(TANDA	ATANGAN PEN	YELIA)

Alamat tetap:	<u>24 Jalan TU38, Taman Tasik</u>	PM. DR. CHOO YUN HUOY	
<u>Utama, 75450</u>), Ayer Keroh, Melaka.	Nama Penyelia	
Tarikh :		Tarikh :	
CATATAN:	* Tesis dimaksudkan sebagai Laporan Projek Sarjana Muda (PSM)		
	**Jika Tesis ini SULIT atau TER	HAD, sila Lampirkan surat daripada	

C Universiti Teknikal Malaysia Melaka

pihak berkuasa.

PREDICTING INFLUENCERS FROM TWEETS DATA

NOREZZATI BINTI MD. NOR HAYATI

This report is submitted in partial fulfillment of the requirements for the Bachelor of Computer Science (Artificial Intelligence)

FACULTY OF INFORMATION AND COMMUNICATION TECHNOLOGY UNIVERSITI TEKNIKAL MALAYSIA MELAKA 2015

C Universiti Teknikal Malaysia Melaka

DECLARATION

I hereby declare that this project report entitled

PREDICTING INFLUENCERS FROM TWEETS DATA

is written by me and is my on effort and that no part has been plagiarized without citations.

STUDENT	:	Date :
	(NOREZZATI BINTI MD. NOR HAYATI)	
SUPERVISO	R:	Date :
	(PM. DR CHOO YUN HUOY)	

DEDICATION

To my beloved parents ..

ACKNOWLEDGEMENT

First of all, I would like to thank my parents for being very supportive. Their supports give me strength to complete my project successfully.

I also would like to thank to my beloved supervisor, PM. Dr. Choo Yun Huoy. Without her support, I will never get through this project successfully. And I appreciate her time for every consultations and progress.

C Universiti Teknikal Malaysia Melaka

ABSTRACT

Twitter generates such overwhelming information every day and allows developers to get the public stream for analytic purposes. Influencer analysis is thus important to provide insight on market segmentation. However, there are not many researches done on identifying influencer because there exist no objective quantification to determine influencing factor. This project aims to find the influential factors and to cluster into the influencer and non-influencer clusters. The experimental was collected from Twitter on #NepalQuake topic. The attribute proposed for influencer prediction includes retweets, favorites, followers and friends. Simple scoring algorithm and Fuzzy C-Means cluster were use in the project. Data crawling was done through Twitter application in Python. Comparison was made to compare the performance of the simple scoring algorithm versus the Fuzzy C-Means cluster in generating influencer list. Both proposed algorithm is able to list and produce the most influential user. The result of influencer was obtained after the weightage is set to each attribute to represent the importance of attributes. The scoring algorithm will list the influential user with their rank. There is only one user for the influencer cluster. This is because only one user has obvious number of retweets, followers, favorites and friends. In future work, the algorithm can be improved to be commercialized. The measuring on how important the attribute need to be made and also the algorithm will be apply to the unstructured data.

ABSTRAK

Twitter menghasilkan maklumat yang banyak setiap hari dan membolehkan pemaju untuk mendapatkan maklumat awam untuk tujuan analisis. Walau bagaimanapun, tidak banyak kajian yang dilakukan mengenai pengenalpastian pengaruh dalam Twitter kerana tidak wujud kuantifikasi bagi menentukan faktor pengaruh. Projek ini bertujuan untuk mencari faktor-faktor pengaruh dan untuk menentukan pengaruh tersebut kepada kumpulan yang berpengaruh dan bukan pengaruh. Kajian ini telah dikumpulkan daripada Twitter menerusi topik #NepalQuake. Antara perkara yang dicadangkan untuk meramal pengaruh adalah 'retweet', 'favorite', 'followers', dan 'friends'. Algoritma pemarkahan dan penggunaan Fuzzy C-Means untuk kelompok adalah teknik yang digunakan dalam kajian ini. Python digunakan untuk mendapatkan data melalui Twitter. Perbandingan telah dibuat untuk membandingkan prestasi algoritma yang dicadangkan berbanding teknik Fuzzy C-Means dalam menghasilkan senarai pengaruh. Kedua-dua algoritma yang dicadangkan dapat menyenaraikan pengguna yang paling berpengaruh. Hasil pengaruh telah diperolehi selepas pemberat yang ditetapkan kepada setiap atribut untuk mewakili kepentingan setiap sifat-sifat. Algoritma pemarkahan dapat menyenaraikan pengguna yang berpengaruh dengan kedudukan mereka. Hanya ada satu pengguna bagi kelompok yang berpengaruh. Ini kerana hanya seorang pengguna mempunyai bilangan yang ketara. Bagi kerja-kerja masa hadapan, algoritma boleh diperbaiki untuk dikomersialkan. Algoritma yang dicadangkan juga perlu dikaji supaya dapat bekerjasama denga data yang tidak berstruktur.

TABLE OF CONTENTS

CHAPTER	SUBJECT	PAGE
	DECLARATION	i
	DEDICATION	ii
	ACKNOWLEDGEMENTS	iii
	ABSTRACT	iv
	ABSTRAK	V
	TABLE OF CONTENTS	vi
	LIST OF TABLES	ix
	LIST OF FIGURES	Х
	LIST OF ABBREVIATIONS	xii
	LIST OF APPENDICES	xiii
CHAPTER I	INTRODUCTION	
	1.1. Project Background	1
	1.2. Problem Statement	3
	1.3. Objectives	3
	1.4. Scope	4
	1.5. Project Significance	4
	1.6. Expected Output	5
	1.7. Summary	5
CHAPTER II	LITERATURE REVIEW	
	2.1. Introduction	7
	2.2. Social Networks	8
	2.2.1. Twitter	8
	2.2.2. Facebook	10

C Universiti Teknikal Malaysia Melaka

	2.2.3. Instagram	11
	2.3. Influencer Factor Analysis	11
	2.4. Computational Technique	13
	2.5. Available Social Network Analysis Tools	14
	2.6. Summary	21
CHAPTER III	PROJECT METHODOLOGY AND	
	TECHNIQUE	
	3.1. Project Methodology	23
	3.1.1. Requirement Analysis	24
	3.1.2. Data Crawling	26
	3.1.3. Extracting Useful Features from Online	27
	Twitter Data	
	3.1.4. Influencer Modelling	27
	3.1.5. Result Analysis	28
	3.2. Project Requirement	29
	3.2.1. Software Requirement	29
	3.3. Project Schedule and Milestone	30
	3.4. Summary	32
CHAPTER IV	INFLUENCER MODELLING	
	4.1. Introduction	33
	4.2. Data Crawling Using Python 2.7.9	34
	4.2.1. Tweepy	34
	4.2.2. CSV	35
	4.3. Influencer Modelling Using Proposed	35
	Algorithm	
	4.4. Fuzzy C-Means Clustering Technique	43
	4.5. Summary	47
CHAPTER V	RESULTS AND DISCUSSION	
	5.1. Introduction	48

vii

viii

5.2. Comparison Analysis	49
5.3. Result Analysis	49
5.4. Discussion	52
5.5. Summary	54

CHAPTER VI CONCLUSION

6.1. Introduction	55
6.2. Weaknesses and Strength	55
6.3. Contribution	56
6.4. Future Work	56

REFERENCES 58

APPENDICES

APPENDIX A	62
APPENDIX B	63
APPENDIX C	64
APPENDIX D	65
APPENDIX E	66
APPENDIX F	67

LIST OF TABLES

TA	BLE
----	-----

TITLE

PAGE

2.1	Comparison between Twitter and Facebook	10
2.2	Comparison of several existing tools for	20
	influencer	
3.1	List of software requirement	29
3.2	Project milestone	30
4.1	Comparison between Python 2.7.9 and Python	34
	3.4.3	
4.2	Entity and its representation	37
4.3	Centers from two clusters	46
5.1	Comparison of proposed method to existing	49
	tool	

LIST OF FIGURES

FI	GI	JR	ES
----	----	----	----

TITLE

PAGE

1.1	Overall approach for the project	2
2.1	Twitter functionality	9
2.2	Concept of Klout algorithm	16
2.3	Output for Keyhole tool	17
2.4	Example of visualization types of Keyhole	18
	tool	
2.5	Output example for Traackr tool	19
3.1	Flow of the methodology phases	24
3.2	Twitter application	25
3.3	Twitter application keys	25
4.1	OAuth authentication is used to request	35
	the Twitter stream	
4.2	Python code for csv instance	35
4.3	Output for raw twitter streams	38
4.4	Python code for one day stream	39
4.5	Python code for one week stream	40
4.6	Python code for one hour stream	41
4.7	Example of dataset for the first run	42
4.8	Python implementation for merging the	43
	files	
4.9	GUI for clustering function in Matlab	44
4.10	Center for both cluster	45
4.11	Sample illustration for the implementation	46
5.1	The only row that return a contra value	50
5.2	Data for user Stephenfry that matches the	51

plotted point

5.3	Data cursor shows point for influencer's	51
	value (Retweets vs Favorites)	
5.4	Data cursor shows point for influencer's	52
	value (Followers vs Friends)	
5.5	Data cursor shows point for influencer's	53
	value (Retweets vs Followers)	

LIST OF ABBREVIATIONS

CSV	-	Comma separated values
XLS	-	Excel file format
API	-	Application programming interface
URL	-	Uniform resource locator
FCM	-	Fuzzy C-Means

xii

LIST OF APPENDICES

APPENDIX

TITLE

PAGE

A	Sample dataset for first run	62
B	Sample dataset for rank collected	63
С	Sample dataset for merged files with ordered	64
	rank	
D	Sample dataset used for FCM clustering	65
	method	
Е	Sample of clustered dataset	66
F	Sample of clustered result	67

CHAPTER I

INTRODUCTION

1.1 Project Background

Social media represent how important information being produced and consumed. User generates content in form of comments, tweets and blog posts establish a connection between producers and the consumers of information. The social media users generate such overwhelming volume of information. This kind of information will be useful to the various marketers around the world. In this project, the information is collected specifically from Twitter data. By searching and predicting the influencers in social media data, it will help us to improve the lifestyle, business and also the authorities.

Influencers are the person or authorities that have a power to be a trendsetter, leader, or someone who sets the rule. These people have a big number of followers, and also can affect others to change what they think. Therefore finding influencers are important nowadays because the market, news media, fashion brand or any company always wanted to know their market and how to approach customers. Especially by searching related topics using a hashtags that Twitter had provided.

Twitter is an online social network also known as a microblogging service that allows user to post 140 characters of short messages. It is called 'tweets'. This microblogging service started in 2006. Twitter users will also allow to read, post, comments, and 'retweet' or also known as share the post they interested in. The information is now expanded in large volume as well as the users. The advantage of using Twitter in this project is that Twitter provides an application programming interface (API) to crawl the public information. Twitter data is accessible for public and useful for data analytic purposes.

In this project, few sizes of data were crawled from Twitter. One topic is selected to make a crawl on the information of users. There are four attributes selected to make a crawl on, which is the number of retweets, number of favorites, number of follower, and number of friends. Details about the functionalities will discuss in the next chapter. The approach of this project is explained as in Figure 1.1.



Figure 1.1 : Overall approach for the project

An algorithm is proposed to find the scoring value of the influencers. In addition to that, the influential percentage is made. The scoring is assigned to every data collected and comparison is made. There are several existing tools that can find and predict influential user, therefore these tools will be set as a benchmark for the scoring value and used in comparison phase. More details about the algorithm and implementation will be explained in the following chapter.

Furthermore, the clustering method will be applied in these dataset collected and will compare the result and analysis topic. In this project, fuzzy c-means clustering method is used as the clustering method. The method should be able to cluster the dataset into two clusters which are the influencer group and noninfluencer group. This report will explain all steps and technique includes preliminary studies.

1.2 Problem Statement

The problem statement for this project is some authorities or company might want to find the most influential user or a group of influencers in certain topics that relate to their research. There are no objective quantifications for measuring the influencers in Twitter. The measurement of influencer is important and need to find the correct measurement in order to get the most logical and accurate influencer. There are many ways to measure the influencer, but it will return different kind of influencers. For example, the followers count only will return the famous user in Twitter. Thus in this project, the weightage is proposed to each factors of influencers.

1.3 Objectives

The purpose of this research is to create a model that can predict influencer from a big dataset. The objectives of this research are as follow:

a) To identify significant influential factor

The dataset need to be crawled and. From the information collected, the influential factors need to be identified and justify from the previous studies.

b) To propose a scoring algorithm to quantify influential value of Twitterers

- The algorithm should be able to return the value of influencer's score. This algorithm need to add the weightage in order to be fair to all attributes and able to sort the influencer and non-influencer
- c) To propose a clustering method to cluster influencers
 - Define one clustering method and cluster the influencer and non-influencers.

1.4 Scope

This project focuses on hashtags topics that exist in Twitter application. Twitter allow public stream extraction for one week duration, thus in this project the chosen topic is '#NepalQuake'. The factor of influencer is selected based on the preliminary studies to do the data crawling and also the proposed technique will add the weightage that is total of 10 is divided to the influential factors. Furthermore, the data will be cluster to two clusters.

1.5 **Project Significance**

The project significance can be divided into two sub topic which is academic significant and industrial significant.

Academically, the project proposed an algorithm to be used for influencer searching and in future for predicting and decision making purposes.

Industrial significance can be explained for future business decision making. By having this type of model or application, the prediction can be made in various types of data. Otherwise, in future, everybody can predict whether a user is an influencer by only reading a single tweet. This model will help the marketer around the world to improve their business. The company will find their target by finding this influencer and also to check their own market influence.

The authorities will be using this model for future decision making especially if the data is collaborated with the unstructured data. Also the model is useful when the data is getting bigger in future. Furthermore, the model can be useful in terms of disaster issues. Such as nature disaster, this model could give alert to certain rescue to the place needed, or the media company can cover as many news at the same time.

1.6 Expected Output

In the end of this research, the suggested algorithm could be able to describe and determine the influencer from huge set of data collection. This means that the proposed model will capable to work with big size of data or simply called the big data.

On the other hand, the proposed algorithm will display both scoring values and predict the accuracy of the model created. All of the crawled data will be saved into the spreadsheet for monitoring purposes. In addition to that, the clustering method will be able to display a group of influencers and non-influencers. At the end of both techniques, the expected results are that it will produce the same username or same information.

1.7 Summary

As the conclusion to this chapter, the proposed technique will be developed in order to make a prediction. Prediction model will ease the human power to make calculation, observation and intuitions. Therefore this model is useful and after it is specialized in predicting the influencer, it will specifically be more useful to certain authorities to improve their business. The report is organized in the following order, literature review in chapter 2, project methodology and techniques in chapter 3, influencer modelling in chapter 4, results and analysis in chapter 5 and conclusion in chapter 6.

CHAPTER II

LITERATURE REVIEW

2.1 Introduction

This chapter will explain the literature review with respect to this project. All the findings and collection of information regarding the topics which include the twitter information will explain in this chapter. Literature review is the phase of finding the previous research and related technical papers that is related to the project. It defines the process of searching related project, collect the information, analyzing the information to related scope, and determine the conclusion. Furthermore, to evaluate other people studies to find the related works and methodology.

In this chapter, the research was made on several related topics. Earlier studies were made on how to measure the influencer or specifically stated as the scoring measurement. The next issue is focused on how to crawl the data from Twitter. More studies were made on the existing tools and comparison of these tools functionalities.

2.2 Social Networks

Nowadays, there are too many social networks that allows users to describe their opinions and at the same time share with their friends around the world. The content that all users are sharing is called a social media. Social media is part of the image, blogs, songs, video, or text that user upload in their social networks, (Social Media vs Social Networking, 2013). The example of social networks in this era is Facebook, Twitter, LinkedIn, Instagram, Flickr, Tumblr and many more. The users of these social networks are from students, teenagers, adults, news media, authorities and also celebrities around the world. In this topic, the analysis of the social network is made.

2.2.1 Twitter

This microblogging service is using the alias "@username" to mention other user to join the conversation and each user can post 140 characters per time (Boyd, Golder, & Lotan, 2010). On one day, user have the limitations of tweeting 1000 post of messages or also known as messages. On the other hand, twitter also provide a famous hashtags function "#topic". This hashtags or trending section will list down the most mentioned topic at time. Every user who uses this hashtag will be assumed to have a homophily relationship or topical similar relationship, (Weng, Lim, & Jiang, 2010). Twitter also having the friends and followers function which is user can decide whether they want to be followed or having friends, and otherwise follow other. The other user who follows the owner of the profile is called followers. While if other user being followed by the profile owner is called following or friends. Each user can follow 2000 users. Another unique function of Twitter is the retweet function. Retweet is also known as sharing. Because of the character limit for posting tweets, Twitter provide retweet function for sharing or embed the post to own post. Favorite is also another function in Twitter, but this function almost being ignored by people. The functionality is same as the well-known 'like' function on Facebook. But differ to Facebook, by favorited a post, the post will be saved in user's profile.

8