



# **Analysis of MST based clustering algorithm with different threshold values**

## **Dissertation**

### **Submitted To:**

Central Department of Computer Science & Information Technology

Tribhuvan University

Kirtipur, Kathmandu

Nepal

In partial Fulfillment of the requirements for the Degree of Master of Science in  
Computer Science & Information Technology

Submitted by:

**Lalit Pant**

February, 2016

Supervisor

**Prof. Dr. Subarna Shakya**

Co-supervisor

**Mr. Arjun Singh Saud**



## **Tribhuvan University**

### **Institute of Science and Technology**

#### **Central Department of Computer Science and Information Technology**

### **Student's Declaration**

I hereby declare that I am the only author of this work and that no sources other than the listed here have been used in this work.

.....  
**Lalit Pant**



**Tribhuvan University**

**Institute of Science and Technology**

**Central Department of Computer Science and Information Technology**

**Supervisor's Recommendation**

I hereby recommend that the dissertation prepared under my supervision by **Mr. Lalit Pant** entitled “**Analysis of MST based clustering algorithm with different threshold values**” be accepted as in fulfilling partial requirement for the completion of Masters Degree of Science in Computer Science & Information Technology.

-----  
**Prof. Dr. Subarna Shakya**

Department of Electronics & Computer Engineering,  
Institute of Engineering,  
Pulchowk, Nepal



**Tribhuvan University**

**Institute of Science and Technology**

**Central Department of Computer Science and Information Technology**

### **LETTER OF APPROVAL**

We certify that we have read this dissertation work and in our opinion it is appreciable for the scope and quality as a dissertation in the partial fulfillment of the requirements of Masters Degree of Science in Computer Science & Information Technology.

#### **Evaluation Committee**

---

**Asst. Prof. Nawaraj Paudel**  
**Head of Department**  
Central Department of Computer Science  
& Information Technology  
Tribhuvan University  
Kirtipur

---

**Prof. Dr. Subarna Shakya**  
**(Supervisor)**  
Department of Electronics & Computer  
Engineering, Institute of Engineering,  
Pulchowk, Nepal

---

**(External Examiner)**

---

**(Internal Examiner)**

## ACKNOWLEDGEMENT

I would never have been able to finish my dissertation without the guidance, support and encouragement of numerous people including my supervisor, co-supervisor, my friends, colleagues and support from my family. At the end of my thesis I would like to thank all those people who made this thesis possible and an unforgettable experience for me.

First, I would like to express my gratitude to my supervisor **Professor Dr. Subarna Shakya**, Institute of Engineering, Pulchowk Campus for his continuous support without which the thesis wouldn't have been possible to complete. His suggestions, guidance, thorough knowledge and expertise helped me immensely in understanding and developing this thesis. Also, the credit of the success of this dissertation work goes to my co-supervisor **Mr. Arjun Singh Saud**. I appreciate for his supervision and guidance.

Most importantly I would like to thank to respected Head of Department of Central Department of Computer Science and Information Technology, **Asst. Prof. Nawaraj Paudel** for his kind support, help and constructive suggestions. I am very much grateful and thankful to all the respected teachers Professor Dr. Shashidhar Ram Joshi, Mr. Dhiraj Kedar Pandey, Mr. Sarbin Sayami, Mrs. Lalita Sthapit, Mr. Bikash Balami and Mr. Jagdish Bhatta for providing me such a broad knowledge and inspirations. I am so much thankful to Mr. Deep Sharma for his continuous support throughout the thesis work and also like to thank my dear friends Ishwari, Rajendra, Dinesh, Chakra, and Bhim, for their cooperation.

Special thanks to my family for their endless motivation, constant mental support and love which have been influential in whatever I have achieved so far. All my class fellows are worthy of my gratefulness for their direct or indirect support in completion of my dissertation.

I have done my best to complete this research work. Suggestions from the readers are always welcomed, which will improve this work.

# ABSTRACT

Clustering analysis has been an emerging research issue in data mining due to its variety of applications. Many algorithms are proposed so far, however each algorithm has been its own merits and demerits and cannot work for real situation. The MST based clustering algorithms have been widely used due to their ability to detect cluster with irregular boundaries. In this dissertation the clustering algorithm is inspired by MST.

In this dissertation the MST based clustering algorithm has been analyzed using different threshold value on MST and measured by validity index. Given the MST over data set, select or reject the edges of MST in process of forming the clusters, depending on the threshold value. Validity index is the ratio of intra cluster distance and inter cluster distance. Thresholds are taken by mean, standard deviation and mean + standard deviation of MST. These thresholds are evaluated by validity index. Smallest value of validity index is select for best clustering and best threshold value. The algorithm has been tested on the randomly generated data sets and as well as real world data sets.

**Keywords: Clustering Algorithm, MST, Validity Index, Threshold Values**

# CONTENTS

<b>DETAILS</b>	<b>PAGE NO</b>
<b>CHAPTER 1</b>	
<b>INTRODUCTION</b>	
1.1 Introduction.....	1
1.2 Problem Definition.....	2
1.3 Objective.....	2
1.4 Motivation.....	2
1.5 Report Organization.....	2-3
<b>CHAPTER 2</b>	
<b>LITERATURE REVIEW</b>	
2.1 Literature Review.....	4-7
<b>CHAPTER 3</b>	
<b>RESEARCH METHODOLOGY</b>	
3.1 Data Collection.....	8
3.2 Data Analysis .....	8
3.3 Performance Metrics.....	8-9
3.3.1 Validity Index.....	8-9
3.3.2 Threshold Value.....	9

**CHAPTER 4**  
**ALGORITHMS**

4.1 Prim's algorithm.....10-14  
4.2 MST based clustering algorithm.....14-17

**CHAPTER 5**  
**IMPLEMENTATION**

5.1 Tools used.....18  
5.2 Programming language.....18  
    5.2.1 Python IDLE.....18-19  
    5.2.2 Scipy.....19  
    5.2.3 NumPy .....19-20  
    5.2.4 Python list.....20

**CHAPTER 6**  
**RESULT AND ANALYSIS**

6.1 Data set 1:.....21-23  
6.2 Data set 2: .....23-26  
6.3 Data set 3: .....26-35  
6.4 Result.....35



## **CHAPTER 7**

### **CONCLUSION AND FUTURE WORK**

7.1 Conclusion.....	36
7.2 Future Work.....	36
<b>References:</b> .....	37-38
<b>Bibliography:</b> .....	39

## LIST OF FIGURES

FIGURES	PAGE NO
Figure1:- Graph.....	10
Figure1.1:-Constructing MST using prim's algorithm.....	11
Figure1.2:-Constructing MST using prim's algorithm.....	11
Figure1.3:-Constructing MST using prim's algorithm.....	12
Figure1.4:-Constructing MST using prim's algorithm.....	12
Figure1.5:-Constructing MST using prim's algorithm.....	13
Figure1.6:-MST obtained by Prim's algorithm.....	13
Figure2:- Minimum Spanning Tree.....	15
Figure2.1:- Number of cluster obtained by MST based clustering algorithm .....	15
Figure3:-Validity Index VS. Threshold Values for data set1.....	22
Figure3.1:-Validity Index VS. Threshold values for data set2.....	26
Figure3.2:-Validity Index VS. Threshold values for data set3.....	35

## LIST OF TABLES

<b>TABLES</b>	<b>PAGE NO</b>
Table1:- V.I and clusters with different threshold values for data set1 .....	22
Table2:- V.I and clusters with different threshold values for data set2.....	26
Table3:- V.I and clusters with different threshold values for data set3.....	35

## LIST OF ABBREVIATIONS

MST	:	Minimum Spanning Tree
EMST	:	Euclidean Minimum Spanning Tree
SEMST	:	Standard Euclidean Minimum Spanning Tree
HEMST	:	Hierarchical Euclidean MST
MSDR	:	Maximum Standard Deviation Reduction
DHCA	:	Divisive Hierarchical Clustering Algorithm
MDHCA	:	Multi Divisive Hierarchical Clustering Algorithm.
LM	:	Larzlo and Mukherjee
SI	:	Sequential Initialization
VI	:	Validity Index