

AN IMPROVED SELF ORGANIZING MAP USING JACCARD NEW MEASURE  
FOR TEXTUAL BUGS DATA CLUSTERING

ATTIKA AHMED

A thesis submitted in  
fulfillment of the requirement for the award of the  
Degree of Master of Information Technology



Faculty of Computer Science and Information Technology  
Universiti Tun Hussein Onn Malaysia

JANUARY, 2018

This research work I dedicated to my late mother, who passed away on April 2016  
and to my brother Syed Nabeel Ahmed.



PTTA UTHM  
PERPUSTAKAAN TUNKU TUN AMINAH

## ACKNOWLEDGEMENTS

In the name of Allah, Most Gracious, Most Merciful. All Praise is to Allah, the Cherisher and Sustainer of the worlds. Oh Allah, send Grace and Honor on Prophet Muhammad (S.A.W) and on the family and true followers of Prophet Muhammad (PBUH).

I would like to express my appreciation and gratitude to my down to earth supervisor Associate Professor Dr. Rozaida Binti Ghazali for her motivation and support to pursue my Masters. I would also like to thank her for the providing me the grant for publishing my paper in IEEE. I am indeed very grateful. I would also like to thank my family for their support and motivation through the difficult times in during my Masters research work journey. And to the faculty members especially FSKTM postgraduate student's society, I say thank you. Thank you for the good company and friendship you have provided me during my Masters journey.

Finally, I would like thank to my very good friends Abdulkadir Hassan Disina and Dr. Zaharuddin Pindar for their unconditional support specially for helping me in writing research papers and to University Tun Hussein Onn Malaysia for providing the supportive environment to conduct my Masters research.



## ABSTRACT

In software projects there is a data repository which contains the bug reports. These bugs are required to carefully analyze to resolve the problem. Handling these bugs humanly is extremely time consuming process, and it can result the delaying in addressing some important bugs resolutions. To overcome this problem researchers have been introduced many techniques. One of the techniques is the bug clustering. For the purpose of clustering, a variety of clustering algorithms available. One of the commonly used algorithm for bug clustering is K-means, which is considered a simplest unsupervised learning algorithm for clustering, yet it tends to produce smaller number of cluster. Considering the unsupervised learning algorithms, Self-Organizing Map (SOM) considers the equally compatible algorithm for clustering, as both algorithms are closely related but different in way they were used in data mining. This research attempts a comparative analysis of both the clustering algorithms and for attaining the results, a series of experiment has been conducted using Mozilla bugs data set. To address the data sparseness issue, the experiment has been performed on textual bugs' data by using two different distance measure which are Euclidean distance and Jaccard New Measure. The research results suggested that SOM has a limitation of poor performance on sparse data set. Thus, the research introduced the improved SOM algorithm by using a Jaccard NM (SOM-JNM). The SOM-JNM produced significantly better results therefore; it can be consider a challenging approach to address the sparse data problems.



## ABSTRAK

Dalam projek perisian terdapat satu simpanan data yang mengandungi laporan pepijat. Pepijat ini diperlukan untuk menganalisis dengan teliti bagi menyelesaikan masalah ini. Mengawal pepijat ini secara semulajadi memakan masa yang sangat lama, dan ia boleh menyebabkan kelewatan dalam menyelesaikan beberapa isu penting pepijat. Untuk menyelesaikan masalah ini, penyelidik telah diperkenalkan dengan pelbagai teknik. Salah satu dari teknik ini adalah pengelasan pepijat. Untuk tujuan pengelasan, pelbagai algoritma tersedia. Salah satu algoritma yang biasa digunakan untuk pengelasan pepijat adalah K-means, yang mana ia dianggap sebagai algoritma pembelajaran tanpa pengawasan yang paling mudah untuk teknik pengelasan, namun ia cenderung untuk menghasilkan bilangan kelas yang lebih sedikit. Memandangkan algoritma ini tidak diawasi, Peta Penyusun Sendiri (SOM) boleh dianggap algoritma yang sesuai untuk pengelasan, kerana kedua-dua algoritma ini berkait rapat tetapi berbeza dengan cara mereka digunakan dalam perlombongan data. Kajian ini cuba untuk membandingkan kedua-dua kelas algoritma dan untuk mencapai keputusan yang lebih tepat, satu siri eksperimen telah dijalankan menggunakan set data pepijat Mozilla untuk menangani isu kejarangan data. Eksperimen telah dilakukan pada data pepijat berbentuk teks dengan menggunakan dua jarak yang berbeza iaitu jarak Euclidean dan Jaccard New Measure. Hasil kajian mencadangkan bahawa SOM mempunyai batasan prestasi yang lemah pada set data jarang. Oleh itu, kajian memperkenalkan algoritma SOM yang lebih baik dengan menggunakan Jaccard NM (SOM-JNM). SOM-JNM telah menghasilkan keputusan yang jauh lebih baik. Oleh itu, SOM-JNM boleh dipertimbangkan sebagai pendekatan yang baik untuk menangani data jarang.

## TABLE OF CONTENTS

<b>DECLARATION</b>	<b>ii</b>
<b>DEDICATION</b>	<b>iii</b>
<b>ACKNOWLEDGEMENT</b>	<b>iv</b>
<b>ABSTRACT</b>	<b>v</b>
<b>ABSTRAK</b>	<b>vi</b>
<b>TABLE OF CONTENTS</b>	<b>vii</b>
<b>LIST OF TABLES</b>	<b>xi</b>
<b>LIST OF FIGURES</b>	<b>xii</b>
<b>LIST OF APPENDICES</b>	<b>xiv</b>
<b>LIST OF PUBLICATIONS</b>	<b>xv</b>
<b>CHAPTER 1 INTRODUCTION</b>	<b>1</b>
1.1 Research Background	1
1.2 Problem Statements	2
1.3 Aim and Objectives	4
1.4 Scope of Research	4
1.5 Significance of Research	5
1.6 Thesis Organization	5
<b>CHAPTER 2 LITERATURE REVIEW</b>	<b>6</b>
2.1 Chapter Introduction	6
2.2 Clustering	7
2.3 Unsupervised Linear and Non-linear Clustering Algorithm	8



2.4	Bug Clustering	9
2.5	Bug Data Clustering Techniques	10
2.6	Neural Network (NN)	13
2.7	Self-Organizing Map (SOM)	15
	2.7.1 SOM Algorithm	18
	2.7.2 SOM Architecture	19
	2.7.3 Best Matching Unit	20
	2.7.4 Applications of SOM	22
2.8	Similarity Measure Methods	24
	2.8.1 Euclidean Distance	24
	2.8.2 Pearson Coefficient	26
	2.8.3 Cosine Similarity	27
	2.8.4 Jaccard Coefficient	29
	2.8.5 Jaccard NM	31
2.9	Advantages of Jaccard NM	32
2.10	Chapter Summary	32
<b>CHAPTER 3 RESEARCH METHODOLOGY</b>		<b>34</b>
3.1	Chapter Introduction	34
3.2	Research Framework	35
3.3	Data Collection and Pre-processing (Phase 1)	37
	3.3.1 DATA Repository	37
	3.3.2 Bugs Data Collection	37
	3.3.3 Data Pre-processing	38
	3.3.4 Metrics Developed	41
3.4	Experimental Setup (Phase 2)	41
	3.4.1 Parameters Settings	42
	3.4.2 Training and Testing Algorithms	44
3.5	Evaluation Comparison (Phase 3)	47
	3.5.1 Internal Evaluation	48
	3.5.2 External Evaluation	50
3.6	Chapter Summary	53
<b>CHAPTER 4 SIMULATION RESULTS AND ANALYSIS</b>		<b>54</b>
4.1	Chapter Introduction	54



4.2	Number of Clusters Determination	55
4.3	Internal Evaluation Results	59
4.3.1	Analyzing K-means	59
4.3.2	Analyzing SOM	61
4.3.3	Analyzing SOM-JNM	62
4.4	External Evaluation Results	63
4.4.1	Analyzing K-means	64
4.4.2	Analyzing SOM	65
4.4.3	Analyzing SOM-JNM	66
4.5	Comparison Based on Internal and External Evaluation	67
4.5.1	Comparison of Internal Evaluation Results	68
4.5.2	Comparison of External Evaluation Results	70
4.6	Parameter Effect on Network Model	72
4.5.1	The Effect of Learning Rate	73
4.5.2	The Effect of Euclidean Distance and Jaccard NM	73
4.7	Evaluation Process	74
4.8	Analysis on Accuracy and Purity	74
4.8.1	Accuracy Comparison	75
4.8.2	Purity Comparison	77
4.9	Chapter Summary	81
<b>CHAPTER 5 CONCLUSIONS AND FUTURE WORKS</b>		<b>82</b>
5.1	Chapter Introduction	82
5.2	Novelty and Research Contribution	83
5.2.1	An Improved SOM Model by Adopting a New Distance Measure Jaccard New Measure (SOM-JNM)	83
5.2.2	To Simulating SOM-JNM on Bug Data Clustering	84
5.2.3	The Evaluation of SOM-JNM Performance Compared to the Benchmark Models	85
5.3	Recommendations for Future Works	85





5.4	Concluding Remarks	86
	<b>REFERENCES</b>	<b>87</b>
	<b>APPENDIX</b>	<b>91</b>
	<b>VITA</b>	<b>92</b>



**PTTA UTHM**  
PERPUSTAKAAN TUNKU TUN AMINAH

## LIST OF TABLES

2.1	List of unsupervised linear and non-linear clustering algorithms.	9
2.2	The summarization of some researchers clustering technique	12
2.3	Contingency Table	31
3.1	Internal and external evaluation measures	47
3.2	The results interpretation for Silhouette average width	49
4.1	K-means internal validity results	59
4.2	SOM validity index	61
4.3	SOM-JNM validity index	62
4.4	K-means external validity results	64
4.5	SOM external validity results	65
4.6	SOM-JNM external validity results	66



## LIST OF FIGURES

2.1(a)	Before clustering, the set of round shape balls	7
2.1(b)	After clustering the empty, shaded and fully colored balls grouped	8
2.2	An example of simple neural network	14
2.3	Illustration of neural network learning rules	15
2.4	SOM grid representation	16
2.5	Structure of SOM	20
2.6	An example of distance of one dimension SOM	21
2.7	An example of distance of two dimension SOM	2Error!
	<b>Bookmark not defined.</b>	
2.8	Euclidean distance representation	25
2.9	The best fit line using positive and negative slope	26
2.10	Cosine similarity Equation representation	28
2.11	Vector representation of similarity using cosine similarity	28
2.12	Jaccard similarity with two data sets	29
2.13	Jaccard coefficient equation elaboration	30
3.1(a)	Proposed Methodology	35
3.1(b)	Proposed Methodology frame work	36
3.2	Canopy Enthought parsing illustration	39
3.3	The obtained data from tokenizing process	40
3.4	An illustration of parse tree	40
3.5	Frame work for preprocessing of data set	41
3.6	Proposed SOM-JNM algorithm using Jaccard NM approach	46
4.1	K-means clusters determination based on Silhouette method	59
4.2	SOM cluster determination based on Silhouette method	57
4.3	SOM JNM clusters determination based on Silhouette method	58
4.4	Internal validity indexes for K-means	60



PTTA  
PERPUSTAKAAN TUN AMINAH

- 4.5 Internal validity indexes for SOM61
- 4.6 Internal validity indexes for SOM-JNM63
- 4.7 External validity results representation for K-means64
- 4.8 External validity results representation for SOM65
- 4.9 External validity representation for SOM-JNM67
- 4.10 Silhouette results comparison for all three network models68
- 4.11 Homogeneity and Separation comparison for all the selected algorithms69
- 4.12 External validity measures comparison for all three algorithms 71
- 4.13 The comparison for accuracy of all the three models75
- 4.14 FM index comparisons for accuracy of all the three models76
- 4.15 Accuracy comparison of Adjusted Rand index for all algorithms77
- 4.16 Homogeneity and Separation percentage comparison in SOM-JNM 78
- 4.17 Homogeneity and Separation comparison of SOM79
- 4.18 Homogeneity and Separation comparison for K-means 80



**LIST OF APPENDIX**

<b>APPENDIX</b>	<b>TITLE</b>	<b>PAGE</b>
A	MATLAB coding for implimentation SOM-JNM	2



**PTTA UTHM**  
PERPUSTAKAAN TUNKU TUN AMINAH

**LIST OF PUBLICATIONS**

- i. Ahmed, A., & Ghazali, R. (2016, October). An improved self-organizing map for bugs data clustering. In *Automatic Control and Intelligent Systems (I2CACIS), IEEE International Conference on* (pp. 135-140). IEEE.



PTTA UTHM  
PERPUSTAKAAN TUNKU TUN AMINAH



# PTTA UTHM

PERPUSTAKAAN TUNKU TUN AMINAH



# PTTA UTHM

PERPUSTAKAAN TUNKU TUN AMINAH



## CHAPTER 1

### INTRODUCTION

#### 1.1 Research Background

The word 'bug' actually is short for Bugbear. Sometimes found as Bugaboo. Its meaning is much closer to 'Gremlin', "where the people who worked on engineering prototypes often grew to suspect that the problems were due to malicious spooks" (Guo *et al.*, 2011). In 1947 at the computational laboratory of Harvard faculty a moth trapped in a relay which caused an error in the Mark II computer, this bug carefully removed and taped to the log book which still exists (Guo *et al.*, 2011). Since then the term "bugs" established for software error.

During the software development or maintenance process the bugs introduce due to many reasons such as misunderstanding of software architectural design, poor strategy of design, implementation, error handling methodology missing in coding, continuous changes in code after development, etc. In software projects, there is a database which use as bug repository since it is used to collect and manage the bugs reports from users and developers as well as with any other technical support or team member. These bugs are required to analysis very carefully that whether the bugs are valid or invalid (occurred due to operating system, but not by product), duplicate or unique (reported earlier or not), important or unimportant (are they really bugs or something else) and who will resolve it. This process is called bug triaging and the person who performs this job is called triage. The trigger manage the bugs repository in a way that only real bugs should contain and so that the important bugs can be addressed quickly. Every day large number of bugs reports, according to Jadhav *et al.*,(2015) almost 20~30 reports received in software system on an average per day (Jadhav *et al.*, 2015). However, in larger projects such as Mozilla (Mozilla, 2010) more than 300 bugs reports on an average per day (Jadhav *et al.*,2015). Handling

these bugs humanly is extremely time consuming process and it can result the delaying in addressing some important bugs resolutions.

To resolve this problem, bugs reports contain a field where the reporter asked to assign bug priority. The reporter needs to report based that how important it is, but sometimes the field left empty or the reporter may not assign the correct priority. As its possible that the reporter opinion about the bug priority is different with the triager. Since the reporter can be a naïve user of the software and the level of software knowledge is insufficient while the triager has more information of the software as a whole. To assign the correct priority to the bug reports, traiger analyzed the content of report using his own knowledge as well which he gained during triagging and than send it to the developer where the developer assigned the priority to the bug depends on its seviourity.

As the bugs testing is one of the important aspects of software maintenance and testing, research has been growing significantly, but still bugs are one of the major concerns for developers. This study focuses on providing the alternate to the developers in order to minimize the bugs' problem. This can be significantly useful as it provides the better accuracy and purity to coup with the duplicate bugs' issues and weak bugs' prioritization problem.

## 1.2 Problem Statement

Software bugs have created much concern among the researchers due to catastrophic consequences which led to many life claiming incidents and financial losses. Significantly the research has been grown in software bugs resolving issue but still bugs are a major concern for both researchers and developers. Considering the conventional approach to overcome the bugs' problems, one of the strategies follows is called reassignment. Developers often used the bug tracking system in order to find the expertise in search of queries, which is the best to resolve the issue in the particular software or sub module so that the task can be reassign to him/her. It is simply because lack of time to investigate deeply the bug and genuine attempt to find a person with better expertise. It happens mostly because of five reasons; finding the root cause, determining ownership, poor bug report quality, hard to determine proper fix, and workload balancing. As the number of reassignment increase, the time

required to fix the bug increase too. On contrary to popular belief, reassignments aren't necessarily 'bad', since it does take a few reassignments to find the true cause of a bug and who to properly fix it. On the other hand, if there were few reassignments but the optimal bug fixer was not identified, then that could lead to a low quality or faulty fixings (Guo *et al.*, 2011). These reassignments can be happens not only for the search of expertise in that particular software or sub module, but also because of disagreements between developer teams.

As the reassignment strategy if not failed, but challenge due to time consuming issue, researchers stress for an alternative approach. Over the last decades there has been some encouraging research for addressing this issue, but still bug resolving issue a big concern for the software industry.

One of the non-conventional approaches to resolve the bugs' issues used by researchers are Neural Networks (NN), often refers as Artificial Neural Networks (ANN). NN is the mathematical alternative approach inspired by the biological nervous system. The NN has adaptability and can handle complex system same like human brain such as images, pattern, speech recognition, etc. NN is the non liner statistical data modeling tool (Mishra, 2015) used by researchers for bugs' data clustering. The advantage of NN is that it is adaptive and does not need any understanding of input data. They are self-adaptive techniques. The drawback is that NNs are not easy to represent and fewer statistical techniques can be applied (Zeng & Rine, 2004).

One of the statistical approaches is K-means. Russ *et al.*,(2009) chose the bug cluster to experiment on the fact that they describe the same defect "We regard a bug and its duplicates as a cluster" (Russet *al.*, 2009). This is one of the issues among the bug clustering problems. To avoid the bug duplication researchers emphasize on bug prioritization. Bug prioritization strategy is very helpful to address the bugs' issues. Many clustering algorithms including K-means, used for bug prioritization but the problem arise when the algorithms space complexity increases due to data sparseness (NRC, 2012).

Therefore, this research proposed the Self-Organizing Map using Jaccard New Measure (SOM-JNM) which is a type of NN. SOM is considered good for large data clustering without any external training. Beside the advantage of SOM over other algorithms to cluster the large data set, the main disadvantage of SOM is that it is a slow learner, as some researcher describe it "Time consuming algorithm, this is



because as the number of neurons activation at same time slower the performance of the algorithm. And as the number increases the computation increases which results in increasing computational time (Patole *et al.*, 2010) . The main reason of being time consumption is the limitation of Euclidean distance of working on multi-dimensional or sparse data, high sensitivity to noise and outliers (especially for sparse data),which results the too much time and resources consumption. Due to Euclidean distance limitation of computing the distance on sparse data set, this research suggested to use another distance measure Jaccard New Measure (Jaccard NM). Since the Jaccard NM provides a room to address the data sparseness issue and to compute the distance with null values. It motivates this research to substitute the Euclidean distance with Jaccard NM. By combining the properties of SOM with Jaccard NM, the clustering ability of the algorithm enhanced and overcome the weak performance over the sparse data set.

### 1.3 Aim and Objectives

The aim of this research is to cluster bug data using the proposed SOM-JNM. Therefore in order to achieve the research aim, few objects have been set:

- (i) To propose an improved Self Organizing Map (SOM) by adopting a new distance measure, the Jaccard NM.
- (ii) To simulate (i) on bug data.
- (iii) To evaluate (ii) based on its cluster purity and accuracy, and benchmark it against the ordinary SOM with other clustering technique K-means.

### 1.4 Scope of Research

The research focuses on improvising SOM by adopting a new distance measure Jaccard NM and to implement SOM-JNM for clustering the Bugzilla Mozilla data (Mozilla , 2010). Mozilla data has been used and the preprocessing has been applied to the data in order to bring the data in a form where clustering can be applied. To simulate the proposed SOM-JNM, bug data from Bugzilla Mozilla data has been normalized and used to evaluate the simulation based on the cluster purity and accuracy. The simulated results compared against the ordinary SOM with other

benchmark algorithms K-means. To evaluate the performance, the research focuses on internal and external measures which are Silhouette, Homogeneity, Separation, Adjusted Rand, Jaccard and Fowlkes Mellow respectively.

### **1.5 Significance of Research**

The research provides an improved approach to solve the bug issues facing by the software industry. SOM is considered one of the best approach for categorize and textual data both as it has the quality of self-organizing against the large data and cluster it. Thus, this research can be beneficial for addressing the bugs' issues into large data repository and the bug prioritization for large data repository will become less time consuming and less costly. Moreover, the study provides the steady and reliable approach for addressing the database bug repository issues to deal in a proper manner so that it will be beneficial for developer and user both as well as the technical support team. On the larger scale it has been beneficial for the software industry to decrease the financial cost needed to spend on bugs resolving issues at the same time for the economy in general.

### **1.6 Thesis Organization**

This thesis has been organized in 5 chapters. The first chapter discusses the existing research gaps in using SOM along with aim and objective of this research work. The second chapter is dedicated to the literature review related to the previously done work over SOM and K-means clustering. The research methodology of this research work has been discussed, to improvise the ordinary SOM by adopting a Jaccard NM and introduce SOM-JNM in third chapter. Fourth chapter is about the experimental results of this research work, where all the experiment and its results discussed in detail. The final chapter of this thesis is chapter five, which is the discussion over the research novelty and the objectives achieved in this research work. It also discussed this research work limitations and the future recommendation for the relevant research.



## CHAPTER 2

### LITERATURE REVIEW

#### 2.1 Introduction

Over the last fifteen years the interest in Artificial Intelligence (AI) in the computer science field has been immensely increased and some very amazing applications of bioscience and engineering field have been designed based on AI. AI basically deals with intellectual mechanism, which includes Neural Network (NN), offers a great advantage over conventional modeling, including the ability of large amount of noisy data from dynamic and nonlinear processes where nonlinearities and variables plays a vital role. NN's are at the forefront of computational systems designed to produce, or at least mimic, intelligent behavior. Unlike classical AI systems that are designed to directly emulate rational, logical reasoning, neural networks aim at reproducing the underlying processing mechanisms that give rise to intelligence as an emergent property of complex, adaptive systems (Heger, 2005). The NN has adaptability and can handle complex system same like human brain; such as image, pattern, speech recognition, etc. In AI most tasks that require intelligence, require an ability to induce new knowledge from experiences. Thus, a large area within AI is machine learning. Machine learning involves the study of algorithms that can extract information automatically. Similarly to AI, machine learning is very broad and can include almost everything, as long as there is some inductive component to it. Data mining is an area that has taken much of its inspiration and techniques from machine learning. Data mining is carried out in specific situation, on a particular data set, with a goal in mind. Moreover, data mining procedures could be either unsupervised



(when the answer is unknown-discovery) or supervised (the answer is known-prediction). Since the goal and data set for data mining can be different so the techniques also varies to extract the results. Common data mining techniques are including clustering, classification, regression trees, and NN (Han & Kamber, 2000). Classification is the task where the answer is known and called as supervised learning, while clustering is the task which is unsupervised learning because the answer is unknown. Clustering procedure used different algorithms to perform the clustering task. One of the commonly used clustering algorithms is K-means (Baçã, Lobo, & Painho, 2005). Other than K-means, SOM is another algorithm used for large data clustering. Beside its limitation of slow learning, SOM considered as good algorithm for large data clustering (Patole *et al.*, 2010).

Thus, this research is focused the bugs data clustering which is usually a large data set and the poor clustering results would have negative effect when the bugs prioritization applied to that data set. This chapter also discussed the related works that in line with the problem under study, the bug data clustering.

## 2.2 Clustering

Clustering is the division of data into groups of similar objects. In clustering, some details are marginalized in exchange for data simplification. Clustering can be regarded as data modeling technique that offers concise summaries of the data. Clustering is therefore related to many disciplines and plays an essential role in a broad range of applications. The set of different balls has been used to elaborate the clustering process. The data set before clustering as shown in Figure 2.1(a), while Figure 2.1(b) is the visualization of the after clustering where, three types of balls has been clustered.

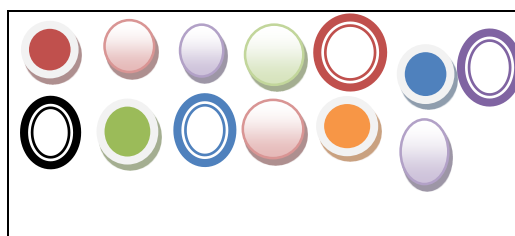


Figure 2.1(a): Before clustering, the set of round shape balls.

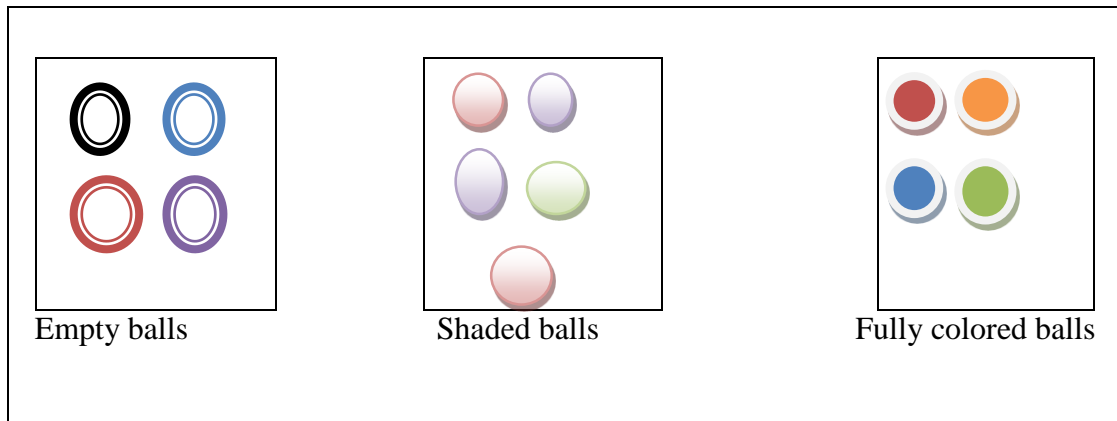


Figure 2.1(b): After clustering the empty, shaded and fully colored balls grouped.

Clustering algorithms are required to classify some conditions in order to be considered efficient algorithms. These conditions are as follows:

- 1) Data must be scalable otherwise it may get the wrong result.
- 2) Clustering algorithm must be able to deal with different types of attributes.
- 3) Clustering algorithm must be able to find clustered data with the arbitrary shape.
- 4) Clustering algorithm must be insensitive to noise and outliers.
- 5) Clustering algorithm must be able to deal with data set of high dimensionality.

Clustering algorithms broadly classified into two categories: unsupervised linear clustering algorithms and unsupervised non-linear clustering algorithms.

### 2.3 Unsupervised Linear and Non-linear Clustering Algorithm

In clustering, unsupervised linear algorithms are usually based on sequential search where the data is linearly separable. For instance, the points belonging to cluster 1 can be separated from the points belonging to cluster 2 by a hyper plane. While in nonlinear algorithms, due to increased data dimension the computation required between two data points. Some of the examples of these algorithms are listed in Table 2.1.



## REFERENCES

- Baçon, F., Lobo, V., & Painho, M. (2005). Self-organizing Maps as Substitutes for K-Means. *Springer* , 476 – 483.
- Bakos, Y. J. (2010). *Mining Similarity Using Euclidean Distance, Pearson Correlation, and Filtering*. Retrieved dec 7, 2015, from humanoriented.com: [http://mines.humanoriented.com/classes/2010/fall/csci568/portfolio\\_exports/mvoget/similarity/similarity.html](http://mines.humanoriented.com/classes/2010/fall/csci568/portfolio_exports/mvoget/similarity/similarity.html)
- Begum, M., & Akthar, M. N. (2013). KSOMKM: An Efficient Approach for High. *International Journal of Electrical Energy*.
- Berglund, E., & Sitte, J. (2013). The Parameter-Less Self-Organizing Map. *IEEE* , 305-316.
- Bloehdorn, S., & Blohm, S. (2006). A self organizing map for relation extraction from wikipedia using structured data representations. Paper presented at the Proc. Int. Workshop on Intelligent Information Access (IIIA-2006).
- Bnanankhah, A., & Nejadkoorki, F. (2012). Artificial Neural Networks: A Non-Linear Tool. *International Conference on Applied Life Sciences (ICALS2012)* (pp. 81-85). Turkey: INTECH.
- Boroš, M. (2011, may 24). *Cluster Analysis*. Retrieved june 16, 2016, from let.rug.nl: <http://www.let.rug.nl/nerbonne/teach/rema-stats-meth-seminar/presentations/Boros-Clustering-2011>
- Cavazos, T. (2000). Using Self-Organizing Maps to Investigate Extreme Climate Events: An Application to Wintertime Precipitation in the Balkans. *AMS* .
- CBS, Department of Systems Biology. (2015). Retrieved from [http://www.cbs.dtu.dk/courses/27618.chemo/27618\\_lecture\\_clustering\\_part2.pdf](http://www.cbs.dtu.dk/courses/27618.chemo/27618_lecture_clustering_part2.pdf)
- Christian, S. P. (2013). *Machine Learning :: Cosine Similarity for Vector Space Models*. Retrieved dec 7, 2015, from pyevolve:



PTTA UTHM  
PERPUSTAKAAN TUNJUKKAN AMANAH

<http://blog.christianperone.com/2013/09/machine-learning-cosine-similarity-for-vector-space-models-part-iii/>

Edhlund, B., & McDougall, A. (2015). *NVivo 10 Essentials*. griffith: Amazon.com.

Fry , Z. P., & Weimer, W. (2013). Clustering Static Analysis Defect Reports to Reduce Maintenance Costs. *IEEE* , 282-291.

Ghaemmaghani, F., & Sarhadi, R. M. (2013). SOMSN: An Effective Self Organizing Map for Clustering. *International Journal of Computer Applications*, (0975 – 8887).

Ghahrai, A. (2008, November 9). *Defect Clustering in Software Testing*. Retrieved may 7, 2015, from testingexcellence.com: <http://www.testingexcellence.com/defect-clustering-in-software-testing/>

Guo, P. J., Zimmermann, T., Nagappan, N., & Murphy, B. (2011). Not my bug! and other reasons for software bug report reassignments. Paper presented at the Proceedings of the ACM 2011 conference on Computer supported cooperative work.

Han, J., & Kamber, M. (2000). *Data Mining: Concepts and Techniques*. Morgan Kaufmann Publishers.

Heger, D. A. (2005). An Introduction to Artificial Neural Networks(ANN)-Methods, Abstraction, and Usage. Retrieved from <http://www.dhtusa.com/media/NeuralNetworkIntro.pdf>

Herbert, J., & Yao, J. T. (2007). GTSOM: GAME THEORETIC SELFORGANIZING. *Springer* , 199-224.

Hynar, M., Burda, M., & Sarmanova, J. (2005). Unsupervised clustering with growing self-organizing neural network--a comparison with non-neural approach. Paper presented at the DATESO.

Ishii, M., Shimodate, T., Kageyama, Y., Takahashi, T., & Nishida, M. (2012). Quantification of Emotions for Facial Expression: Generation of Emotional Feature Space Using Self-Mapping. In *Applications of Self-Organizing Maps* (p. 298). InTech, Chapters.

Jadhav, A., Jadhav, K., Bhalerao, A., & Kharade, A. (2015). A Survey on Software Data Reduction Techniques. *IJCSIT*, 4611-4612.

Kanwal, J., & Maqbool, O. (2012). Bug Prioritization to Facilitate Bug Report Triage. *JOURNAL OF COMPUTER SCIENCE AND TECHNOLOGY* .



PTTA UTHM  
PERPUSTAKAAN TUNKU TUN AMINAH

- Kaur, M., & Garg, S. K. (2014). Survey on Clustering Techniques in Data Mining for Software Engineering. *International Journal of Advanced and Innovative Research*, 3, 238-243.
- Kohonen, T. (2012). *Self-organization and associative memory* (Vol. 8): Springer.
- Kriesel, D. (2007). *A Brief Introduction to Neural Networks*. Retrieved August, 15, 2011. Bonn, Germany.
- Liu, B., Xia, Y., & Yu, P. S. (2005). Clustering Via Decision Tree Construction Foundations and advances in data mining (pp. 97-124): Springer.
- Lobo, V. J. (2009). Application of Self-Organizing Maps. In *Information Fusion and Geographic Information Systems* (pp. 19-36). Springer.
- McCulloch, W. S., & Pitts, W. (1990). A logical calculus of the ideas immanent in nervous activity. *Bulletin of mathematical biology*, 52(1-2), 99-115.
- Mehta, P. (2016). *Karl Pearson's Formula for Finding the Degree of Correlation*. Retrieved 2016, from economicsdiscussion.net: <http://www.economicsdiscussion.net/correlation/karl-pearseons-formula-for-finding-the-degree-of-correlation/2579>
- Mozilla . (2010). Retrieved from Mozilla: Mozilla. <http://www.mozilla.org>, 2010.
- Naseem, R., Maqbool, O., & Muhammad, S. (2010). An improved similarity measure for binary features in software clustering. Paper presented at the Computational Intelligence, Modelling and Simulation (CIMSIM), 2010 Second International Conference on.
- Naseem, R., Maqbool, O., & Muhammad, S. (2011). Improved similarity measures for software clustering. Paper presented at the Software Maintenance and Reengineering (CSMR), 2011 15th European Conference on
- NRC, N. R. (2012). Model Validation and Prediction. In *Assessing the Reliability of Complex Models: Mathematical and Statistical Foundations of Verification, Validation, and Uncertainty Quantification* (pp. 52-83). The National Academic Press.
- Patel, K., Sawant, P., Tajane, M., & Shankarmani, D. R. (2015). BUG TRACKING AND PREDICTION. *International Journal of Innovative and Emerging Research in Engineering* , 174-179.
- Patole, M. V., Pachghare, M. V., & Kulkarni, D. P. (2010). Self Organizing Maps to Build Intrusion Detection System. *International Journal of Computer Applications* .



PTTA UTHM  
 PERPUSTAKAAN TEKNIKAL AMINAH

- Prajapat, R. (2015, February 24). *DnI Institute*. Retrieved May 14, 2015, from <http://dni-institute.in>: <http://dni-institute.in/blogs/neural-network-tutorial/>
- Rosenberg, A., & Hirschberg, J. (2007). V-Measure: A Conditional Entropy-Based External Cluster Evaluation Measure. *EMNLP-CoNLL Vol. 7*, (pp. 410-420). paruge.
- Rousseeuw, P. J. (1987). Silhouettes: a graphical aid to the interpretation and validation of cluster analysis. *Elsevier* , 53-65.
- Rus, & Shiva. (2009). Clustering of Defect Reports Using Graph Partitioning Algorithms. 442-445.
- Russo, T., Scardi, M., & Cataudella, S. (2014). Applications of Self-Organizing Maps for Ecomorphological Investigations through Early Ontogeny of Fish. *PLOS* .
- Saimadhu. (2015, april 4). *Five most popular similarity measures implementation in python*. Retrieved march 2016, from dataaspirant: <http://dataaspirant.com/2015/04/11/five-most-popular-similarity-measures-implementation-in-python/>
- Sathya, R., & Abraham, A. (2013). Comparison of Supervised a. (*IJARAI*) *International Journal of Advanced Research in Artificial Intelligence*, , 34-38.
- Schmidhuber, J. (2015). Deep learning in neural networks: An overview. *elsvier* , 85-117.
- Stefnoic, P., & Kurasova, O. (2011). Influence of learning rate and nebghiouing functions on Self Organizing Map. In *Advances in Self Organizing Map* (pp. 141-150). Springer Berlin Heidelberg.
- Weber, M., Teeling, H., Huang, S., Waldmann, J., Kassabgy, M., M Fuchs, B., et al. (2011). Practical application of self-organizing maps to interrelate biodiversity and functional data in NGS-based metagenomics. *The ISEM Journal* , 918-928.
- Xu, R., & Wunsch, D. (2005). Survey of clustering algorithms. *Neural Networks, IEEE Transactions on*, 16(3), 645-678.
- Zeng, H., & Rine, D. (2004). Estimation of software defects fix effort using neural networks. Paper presented at the Computer Software and Applications Conference, 2004. COMPSAC 2004. Proceedings of the 28th Annual International.

