



PERFORMANCE EVALUATION OF SIMILARITY MEASURES FOR K-MEANS CLUSTERING ALGORITHM

Usman, D.^{1*} and Sani, S.F.²

¹Department of Mathematics and Computer Science, Faculty of Natural and Applied Sciences
Umaru Musa Yar'adua University, Katsina-Nigeria

²Department of Business Administration, Faculty of Social and Management Sciences
Al-Qalam University, Katsina-Nigeria

*Corresponding author: dausman@gmail.com

ABSTRACT

Clustering is a useful technique that organizes a large quantity of unordered datasets into a small number of meaningful and coherent clusters. Every clustering method is based on the index of similarity or dissimilarity between data points. However, the true intrinsic structure of the data could be correctly described by the similarity formula defined and embedded in the clustering criterion function. This paper uses squared Euclidean distance and Manhattan distance to investigate the best method for measuring similarity between data objects in sparse and high-dimensional domain which is fast, capable of providing high quality clustering result and consistent. The performances of these two methods were reported with simulated high dimensional datasets.

Keywords *k-means clustering, similarity measures, squared euclidean distance, manhattan distance.*

INTRODUCTION

Clustering is a process of grouping a set of physical objects into classes of similar objects and is a most interesting concept of data mining in which it is defined as a collection of data objects that are similar to one another. The purpose of Clustering is to catch fundamental structures in a data and classify them into meaningful group. One of the top most popular clustering methods is the K-Means algorithm due to its simplicity, understandability, and scalability. Hartigan and Wang (1979) opined that cluster analysis is one tool that is used in the exploration of data in which the interactions among patterns are assessed by placing them into groups with unique and distinct characteristics. Guojun et al. (2007) defined cluster analysis as a technique for creating groups of objects such that each cluster contains points that are similar and unique.

The objective is targeted at finding the best grouping for which the observations or objects found in within each cluster are the same. More accurately, cluster analysis consists of a series of processes that partition a given data set $X = \{\bar{x}_1, \bar{x}_2, \dots, \bar{x}_n\} \subset \mathcal{R}^D$ into clusters such that the data points in a cluster are more similar to

each other than points in different clusters (Moses et al., 1999). Thus the principal interest in the clustering process is the revelation of sensible groups or patterns, which allow for the discovery of similarities and dissimilarities so that useful conclusions can be reached. Yet, the standard K-Means method suffers a few drawbacks when clusters are of differing sizes, densities and non-globular shape.

Basically, there is an implicit assumption that the true intrinsic structure of data could be correctly described by the similarity formula defined and embedded in the clustering criterion function. Hence, effectiveness of clustering algorithms under this approach depends on the appropriateness of the similarity measure to the data at hand. The work in this paper is motivated by investigations from the above and similar research findings. It appears to us that the nature of similarity measure plays a very important role in the success or failure of a clustering method. Hence, our objective is to check the best method for measuring similarity between data objects in sparse and high-dimensional domain which is fast, capable of providing high quality clustering result and consistent performance.

MATERIALS AND METHODS

Before clustering the objects, a similarity/distance measure must be determined. The measure reflects the degree of closeness or separation of the target objects and should correspond to the characteristics that are believed to distinguish the clusters embedded in the data. In many cases, these characteristics are dependent on the data or the problem context at hand, and there is no measure that is universally best for all kinds of clustering problems. Moreover, choosing an appropriate similarity measure is also crucial for cluster analysis, especially for a particular type of clustering algorithms. For example, the density-based clustering algorithms, such as DBScan rely heavily on the similarity computation.

Therefore, understanding the effectiveness of different measures is of great importance in helping to choose the best one. However, not every distance measure is a metric. Also to qualify as a metric, a measure d must satisfy the following four conditions: Let x and y be any two objects in a set and $d(x, y)$ be the distance between x and y .

- i. The distance between any two points must be nonnegative, that is, $d(x, y) \geq 0$.
- ii. The distance between two objects must be zero if and only if the two objects are identical, that is, $d(x, y) = 0$ if and only if $x = y$.
- iii. Distance must be symmetric, that is, distance from x to y is the same as the distance from y to x , ie. $d(x, y) = d(y, x)$.
- iv. The measure must satisfy the triangle inequality, which is $d(x, z) \leq d(x, y) + d(y, z)$.

Similarity Measures

Similarity measures quantify how “similar” two patterns are. In most cases we have to ensure that all selected features contribute equally to a similarity measure and there are no features that dominate others. Similarity is fundamental to the definition of a cluster; a measure of the similarity between two patterns drawn from the same feature space is essential to most clustering procedures. It is most common to calculate the dissimilarity between two patterns using a distance measure defined on the feature space. Because of the variety of feature types and scales, the distance measures must be chosen carefully.

Distances and similarities play an important role in cluster analysis (Jain and Dubes, 1988; Anderberg, 1973). In the literature of data clustering, similarity measures, similarity coefficients, dissimilarity measures, or distances are used to describe quantitatively the similarity

or dissimilarity of two data points or two clusters.

In general, distance and similarity are reciprocal concepts. Often, similarity measures and similarity coefficients are used to describe quantitatively how similar two data points are or how similar two clusters are: the greater the similarity coefficient, the more similar are the two data points. Dissimilarity measure and distance are the other way around: the greater the dissimilarity measure or distance, the more dissimilar are the two data points or the two clusters.

Every clustering algorithm is based on the index of similarity or dissimilarity between data points Jain and Dubes (1988). If there is no measure of similarity or dissimilarity between pairs of data points, then no meaningful cluster analysis is possible. A distance metric is a real-valued function d , such that for any points x, y and z :

$$d(x, y) \geq 0, \text{ and } d(x, y) = 0 \text{ if and only if } x = y$$

$$d(x, y) = d(y, x) \tag{2.1}$$

$$d(x, z) \leq d(x, y) + d(y, z) \tag{2.2}$$

First property, positive definiteness, assures that distance is always a nonnegative quantity, so the only way distance can be zero is for the points to be the same. The second property indicates the symmetry nature of distance. The third property is the triangle inequality, according to which introducing a third point can never shorten the distance between two points (Larose, 2005). There are several measures of distance which satisfy the metric properties, some of which are:

Euclidean Distance

The Euclidean distance is the most common distance metric used in low dimensional data sets. It is also known as L_2 norm. The Euclidean distance is the usual manner in which distance is measured in real world. In this sense, Manhattan distance tends to be more robust to noisy data.

$$d_{\text{euclidean}}(X, Y) = \sqrt{\sum_i (x_i - y_i)^2} \tag{2.4}$$

where X and Y are m -dimensional vectors and denoted by $X = (x_1, x_2, x_3, \dots, x_m)$ and $Y = (y_1, y_2, y_3, \dots, y_m)$ represent the m attribute values of two records (Larose, 2005). While Euclidean metric is useful in low dimensions, it doesn't work well in high dimensions. The drawback of Euclidean distance is that it ignores the similarity between attributes. Each attribute is treated as totally different from all of the attributes Ertoz et al. (2003).

Manhattan Distance

This metric is also known as L_1 norm or the rectilinear distance. This is also a common distance metric and gets its name from the rectangular grid patterns of streets in midtown Manhattan. Hence, another name for the distance metric is also city block distance. It is defined as the sum of distances travelled along each axis.

The Manhattan distance looks at the absolute differences between the coordinates. In some situations, this metric is more preferable to Euclidean distance, because the distance along each axis is not squared so a large difference in one dimension will not dominate the total distance Berry and Linoff (1997).

$$d_{\text{manhattan}}(X, Y) = \sum_i^m |x_i - y_i| \tag{2.5}$$

Experimental Results and Discussion

It is very difficult to conduct a systematic study comparing the impact of similarity metrics on cluster quality, because objectively evaluating cluster quality is difficult in itself. In practice, manually assigned category labels are usually used as baseline criteria for evaluating clusters. As a result, the clusters, which are generated in an unsupervised way, are compared to the pre-defined category structure, which is normally created by human experts. This kind of evaluation assumes that the objective of clustering is to replicate human thinking, so a clustering solution is good if the clusters are consistent with the manually created categories. However, in practice datasets often come without any manually created categories and this is the exact point where clustering can help. Therefore, measures like cluster coherence in terms of the within-cluster distances and the

well-separateness between clusters in terms of between-cluster distances were used for evaluation in this paper.

The two metric distance functions discussed in section 2 are analysed and compared. The K-Mean clustering algorithm was implemented using each of the metric distance functions: Squared Euclidian and Manhattan distance measures. The cluster formations, error sum of squares, as the smaller the error sum of squares the better cluster formation and the running time required for the two approaches were used to measure the clustering quality among the two approaches. A simulation experiment is conducted with the pairs $(p, n) = (20, 500), (50, 500)$, where p refers to the number of variables and n is the sample size. The data was generated from multivariate normal distribution $N_p(0, I_p)$ with covariance matrix $\Sigma = b\Gamma$, $b > 0$, and Γ is a symmetric matrix of size $(p \times p)$ with all diagonal elements equal 1 and all off diagonal elements equal ρ where $\rho = 0$ and $b = 1.2$ as in Mason *et al.* (2009). The $\rho = 0$ values is a representative of no correlation. For $b = 1.2$ the covariance matrix for the $\rho = 0$ value is:

$$\Sigma = \begin{bmatrix} 1.2 & 0 & \dots & 0 \\ 0 & 1.2 & \dots & 0 \\ \vdots & \vdots & \ddots & \vdots \\ 0 & 0 & \dots & 1.2 \end{bmatrix};$$

$\rho = 0$

In order to make the advantage of the two approaches very clear, show its separation and compactness the paper consider three and five centroids. The running time required by each experiment and their error sum of squares for the two approaches are presented in Table 1. The cluster formations are also shown in Figure 1 to 8 respectively.

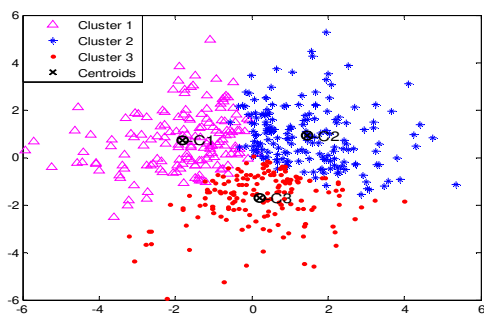


Figure 1: K-Means clustering with SED

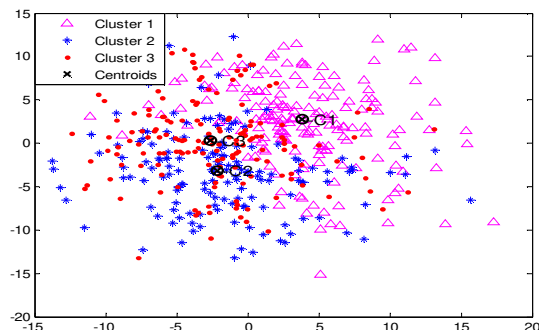


Figure 2: K-Means clustering with MD

Figure 1 and 2 gives the results of the K-Means clustering using Squared Euclidean distance (SED) and Manhattan distance (MD) with simulated dataset containing 500 sample size and 20 variables. Their error sums of squares are 14567.2, 35928.9 and the time taken for execution equal 9.63 and 10.45 respectively.

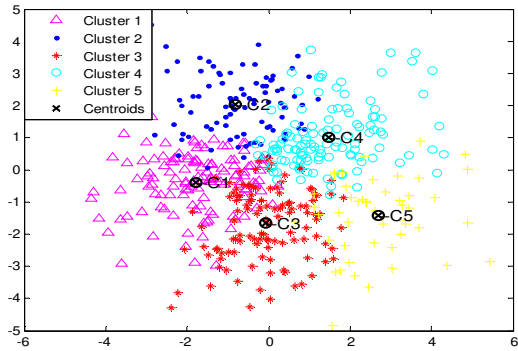


Figure 3: K-Means clustering with SED

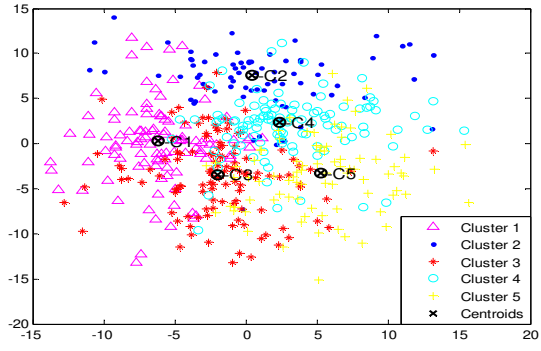


Figure 4: K-Means clustering with MD

Figure 3 and 4 gives the results of the K-Means clustering using Squared Euclidean distance (SED) and Manhattan distance (MD) with simulated dataset containing 500 sample size and 20 variables. Their error sums of squares are 13948.5, 34918.5 and the time taken for execution equal 6.74 and 7.16 respectively.

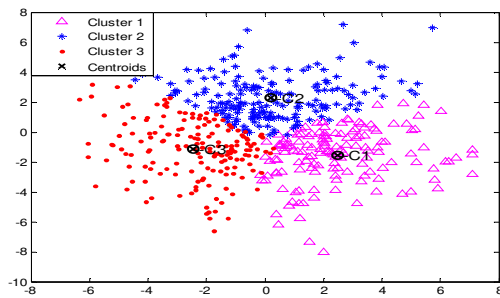


Figure 5: K-Means clustering with SED

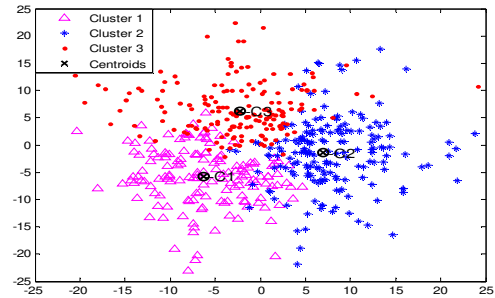


Figure 6: K-Means clustering with MD

Figure 5 and 6 gives the results of the K-Means clustering using Squared Euclidean distance (SED) and Manhattan distance (MD) with simulated dataset containing 500 sample size and 50 variables. Their error sums of squares are 28581.4, 61354.5 and the time taken for execution equal 07.86 and 09.34 respectively.

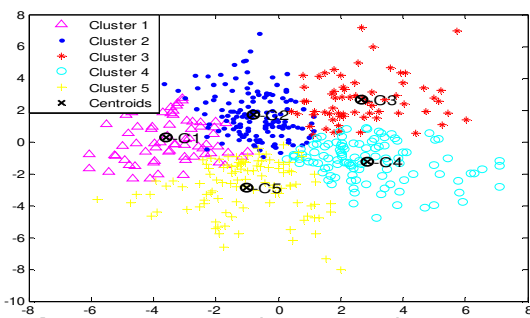


Figure 7: K-Means clustering with SED

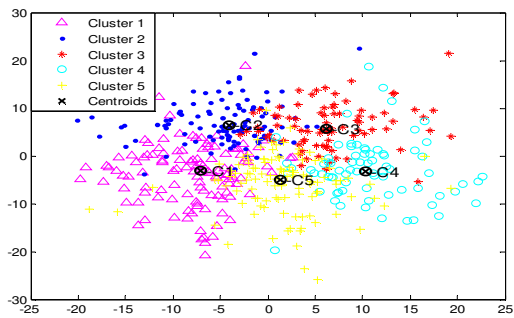


Figure 8: K-Means clustering with MD

Figure 7 and 8 gives the results of the K-Means clustering using Squared Euclidean distance (SED) and Manhattan distance (MD) with simulated dataset containing 500 sample size and 50 variables. Their error sums of squares are 27380.2, 60351.4 and the time taken for execution equal 08.11 and 09.89 respectively.

Table 1: Error Sum of Squares and Time Taken

Method	Error Sum of Squares (20, 500)	Time Taken (20, 500)	Error Sum of Squares (50, 500)	Time Taken (50, 500)
SED 3 Centers	14567.2	9.63	28581.4	07.86
MD 3 Centers	35928.9	10.45	61354.5	09.34
SED 5 Centers	13948.5	6.74	27380.2	08.11
MD 5 Centers	34918.5	7.16	60351.4	09.89

CONCLUSION

A distance measuring function is used to measure the similarity among objects, in such a way that more similar objects have lower dissimilarity value. Several distance measures can be employed for clustering tasks. Each measure has its own merit and demerits. The selection of different measures is a problem dependent. Hence, choosing an appropriate distance measure for K-Mean clustering algorithm can greatly reduce the burden of the

algorithm. The experimental results implies that K-Means method performs very well with Squared Euclidian distance providing better error sum of squares and reduced time taken for the execution as shown in Table 1. However, it was also observed that the clusters are well separated and compact as revealed in Figure 1, Figure 3, Figure 5 and Figure 7. This agrees with the findings of Berry and Linoff (1997) that says, compactness and separation are used to measure the significance of clustering results.

REFERENCES

- Anderberg, M. (1973). Cluster analysis for applications. New York: Academic Press.
- S. Salleh, S. Olariu and B. Sanugi. Single-row transformation of complete graphs. *Journal of Supercomputing*, 31, 265-279, 2005.
- Berry, M. J. A. and Linoff, G. S. (1997). *Data Mining Techniques for Marketing, Sales and Customer Support*. John Wiley & Sons, Inc., New York,
- Ertöz, L., Steinbach, M. and Kumar, W. (2003). Finding Clusters of Different Sizes, Shapes, and Densities in Noisy, High Dimensional Data, *Proceedings of the Third SIAM International Conference on Data Mining*, Volume 3, 2003, San Francisco.
- Guojun, G., Chaoqun, M. and Jianhong, W. (2007). *Data Clustering Theory, Algorithms and Applications*. American Statistical Association and The Society for Industrial and Applied Mathematics.
- Hartigan, J. and Wang, M. (1979). A K-means clustering algorithm. *Appl. Stat.*, 28:100-108.
- Jain, A. and Dubes, R. (1988). *Algorithms for Clustering Data*. Englewood Cliffs, NJ: Prentice-Hall.
- K-median Problem. *Proceedings of the 31st Annual ACM Symposium on Theory of Computing*.
- Larose, D. I. (2005). *Discovering Knowledge in Data: An Introduction to Data Mining*, New Jersey: John Wiley and Sons.
- Mason, R. L., Chaou, Y. M. and Young, J. C. (2009). Monitoring Variation in a Multivariate Process when the Dimension is Large Relative to the Sample Size, *Communication in Statistics. Theory and Methods*, 36:(6), 939-951.
- Moses, C., Guha, S., Tardos, E. and David, B. S. (1999). A Constant-Factor Approximation Algorithm for the K-median Problem. *Proceedings of the 31st Annual ACM Symposium on Theory of Computing*.