

An analysis of Projection Based Multiplicative Data Perturbation for K-means Clustering

Bhupendra Kumar Pandya, Umesh Kumar Singh, Keerti Dixit

*Institute of Computer Science
Vikram University, Ujjain*

Abstract: Random Projections is a very simple yet powerful technique for dimensionality reduction. In this method the data is projected on to a random subspace, which preserves the approximate Euclidean distances between all pairs of points after the projection. It can be proved that the inner product and Euclidean distance are preserved in the new data in the expectation. And many important Data Mining algorithm (e.g., K-means Clustering, KNN Classification etc.) can be efficiently applied to the transformed data and produce expected result. In this research paper we analysis Projection Based Multiplicative data perturbation for k-means Clustering as a tool for privacy-preserving data mining.

Keywords:- Random Projection, K-means Clustering.

1. INTRODUCTION :

Random projection method is very simple and computationally efficient techniques to reduce dimensionality for learning from high dimensional data. This approach is fundamentally based on the Johnson-Lindenstrauss lemma [1], which notes that any set of m points in n -dimensional Euclidean space can be embedded into an $O(\ln m/e^2)$ dimensional space such that the pair wise distance of any two points is maintained with a high probability. Therefore, by projecting the data onto a lower dimensional random space, we can dramatically change its original form while preserving much of its distance-related characteristics. This research paper presents extensive theoretical analysis and experimental results on the accuracy and privacy of the random projection-based data perturbation technique.

1.1 Definition and Fundamental Properties

Random projection refers to the technique of projecting a set of data points from a high dimensional space to a randomly chosen lower dimensional space. Mathematically, let $X \in \mathbb{R}^{n \times m}$ be m data points in n -dimensional space. The random projection method multiplies X by a random matrix $R \in \mathbb{R}^{k \times n}$, reducing the n dimensions down to just k . It is well known that random projection preserves pairwise distances in the expectation. This technique has been successfully applied to a number of applications, for example, VLSI layout [2], nearest-neighbor search [3, 4], image and text clustering [5], distributed decision tree construction [6], motifs in bio-sequences [7] discovery, high-dimensional Gaussian mixture models learning [8], half spaces and intersections of half spaces learning [9].

The following are the steps to reduce the dimensionality of the data by random projections: Suppose that we have a data set $X = \{x_1, \dots, x_n\}$ where each data point is a p

dimensional vector such that $x_i \in \mathbb{R}^p$ and we need to reduce the data to a q dimensional space such that $1 \leq q < p$.

- 1) Arrange the data into a $p \times n$ matrix where p is the dimensionality of the data and n is the number of data points.
- 2) Generate a $q \times p$ random projection matrix R^* using the MATLAB `randn(q, p)` function.
- 3) Multiply the random projection matrix with the original data in order to project the data down into a random projection space.

$$X_{q \times n}^* = R_{q \times p}^* * X_{p \times n}$$

Thus we can see that transforming the data to a random projection space is a simple matrix multiplication with the guarantees of distance preservation.

2. K-MEANS CLUSTERING BY RANDOM PROJECTION

k-means clustering is a data mining/machine learning algorithm used to cluster observations into groups of related observations without any prior knowledge of those relationships. The k-means algorithm is one of the simplest clustering techniques and it is commonly used in medical imaging, biometrics and related fields.

2.1 The k-means Algorithm

The k-means algorithm is an evolutionary algorithm that gains its name from its method of operation. The algorithm clusters observations into k groups, where k is provided as an input parameter. It then assigns each observation to clusters based upon the observation's proximity to the mean of the cluster. The cluster's mean is then recomputed and the process begins again. Here's how the algorithm works:

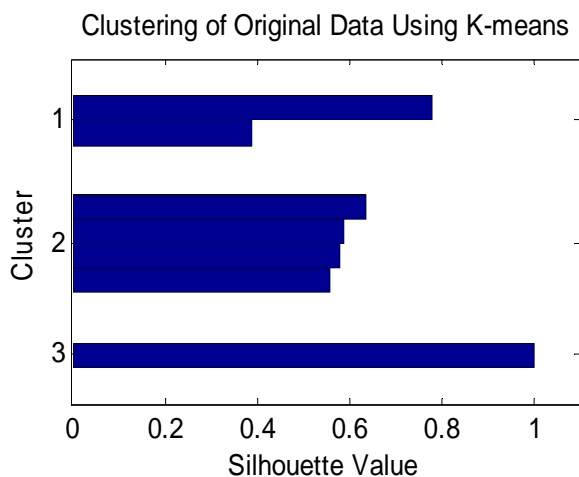
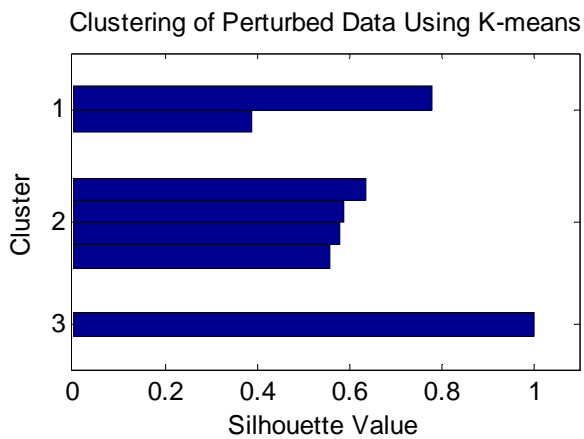
1. The algorithm arbitrarily selects k points as the initial cluster centers ("means").
2. Each point in the dataset is assigned to the closed cluster, based upon the Euclidean distance between each point and each cluster center.
3. Each cluster center is recomputed as the average of the points in that cluster.
4. Steps 2 and 3 repeat until the clusters converge. Convergence may be defined differently depending upon the implementation, but it normally means that either no observations change clusters when steps 2 and 3 are repeated or that the changes do not make a material difference in the definition of the clusters.

3. EXPERIMENTAL RESULTS

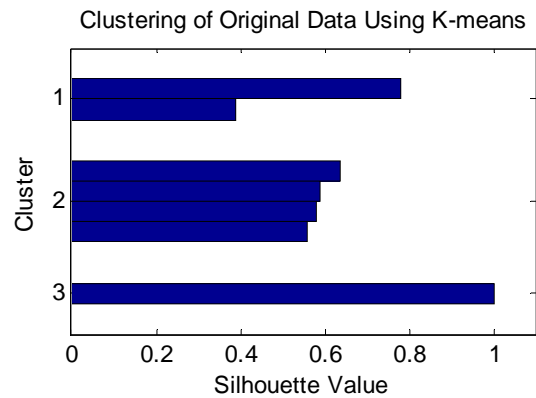
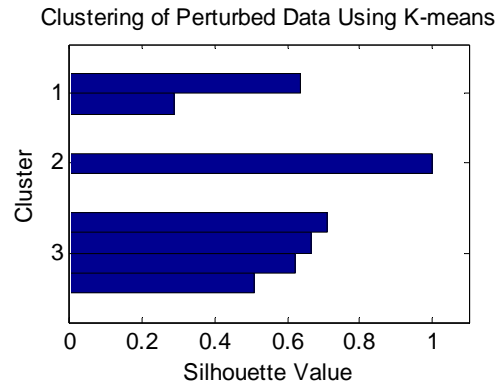
In this study we have Students result database of Vikram University, Ujjain. I randomly selected 7 rows of the data with only 7 attributes (Marks of Foundation, Marks of Mathematics, Marks of Physics, Marks of Computer Science, Marks of Physics Practical, Marks of Computer Science Practical and Marks of Job Oriented Project).

With this data we have generated noise matrices with the help of different different projections and these resultant noise data sets are multiplied with the original data set to form the perturb data sets. We have evaluated Euclidean Distance of original and perturbed data sets with pdist() fuction of Matlab. According to the expectation the Euclidean Distance among the data records are preserved after perturbation. With the Original data we have generated 3 clusters from the kmeans() function of matlab. And similarly we have generated 3 clusters by using the same function with the perturbed data sets. We have used silhouette function for plotting graph of the clustered data generated by the original data and also for plotting graph of the clustered data generated by perturbed data sets.

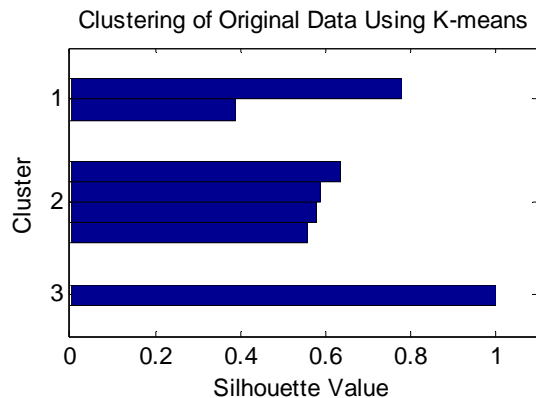
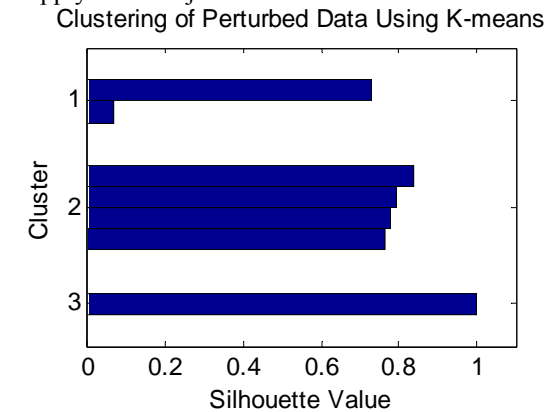
K-means Clustering of Original Data and Perturbed Data after apply 85% Projection -



K-means Clustering of Original Data and Perturbed Data after apply 70% Projection -

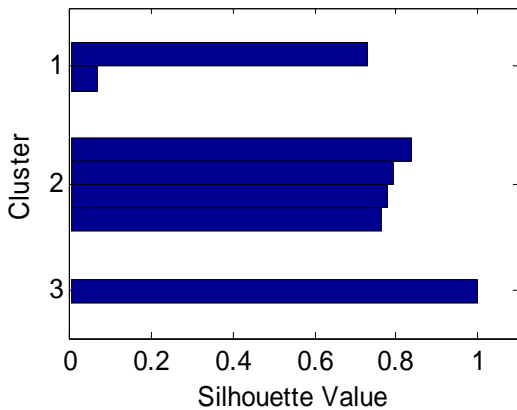


K-means Clustering of Original Data and Perturbed Data after apply 60% Projection -

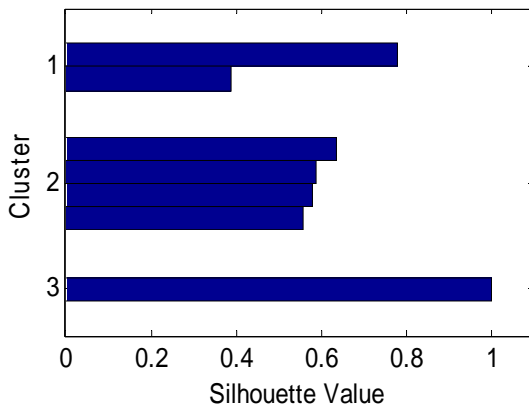


K-means Clustering of Original Data and Perturbed Data after apply 45% Projection –

Clustering of Perturbed Data Using K-means



Clustering of Original Data Using K-means



4. DISCUSSION

It is proved by the experimental result that we get the expected result after applying clustering to the perturbed data as after applying clustering to the original data. Hence we can say that data perturbed by this technique can be used in clustering techniques and we can work with high dimensional data and large datasets. So we can use the perturbed data in various data mining applications like marketing, organization, land use, insurance, city planning etc.

5. CONCLUSION

In this research paper, we have analyzed the effectiveness of Projection based perturbation and we considered the use of this technique as a data perturbation technique for privacy preserving data mining. This technique is quite useful as it allows many interesting data mining algorithms to be applied directly to the perturbed data and produce expected result, e.g., K-means clustering, with little loss of accuracy.

The tremendous popularity of K-means algorithm has brought to life many other extensions and modifications. Euclidean distance is an important factor in k-means clustering. In Distance preserving perturbation technique the Euclidean distance is preserved after perturbation. Hence the data perturbed by this technique can be used in various clustering techniques.

REFERENCE

- [1] H. Kargupta, S. Datta, Q. Wang, and K. Sivakumar, "On the privacy preserving properties of random data perturbation techniques," in Proceedings of the IEEE International Conference on Data Mining, Melbourne, FL, November 2003.
- [2] Z. Huang, W. Du, and B. Chen, "Deriving private information from randomized data," in Proceedings of the 2005 ACM SIGMOD Conference, Baltimore, MD, June 2005, pp. 37–48.
- [3] S. Guo and X. Wu, "On the use of spectral filtering for privacy preserving data mining," in Proceedings of the 21st ACM Symposium on Applied Computing, Dijon, France, April 2006, pp. 622–626.
- [4] N. R. Adam and J. C. Worthmann, "Security-control methods for statistical databases: a comparative study," ACM Computing Surveys (CSUR), vol. 21, no. 4, pp. 515–556, 1989.
- [5] G. T. Duncan and S. Mukherjee, "Optimal disclosure limitation strategy in statistical databases: Dtering tracker attacks through additive noise," Journal of The American Statistical Association, vol. 95, no. 451, pp. 720–729, 2000.
- [6] R. Gopal, R. Garfinkel, and P. Goes, "Confidentiality via camouflage: The cvc approach to disclosure limitation when answering queries to databases," Operations Research, vol. 50, no. 3, pp. 501–516, 2002.
- [7] D. Agrawal and C. C. Aggarwal, "On the design and quantification of privacy preserving data mining algorithms," in Proceedings of the twentieth ACM SIGMOD-SIGACT-SIGART symposium on Principles of Database Systems, Santa Barbara, CA, 2001, pp. 247–255.
- [8] S. Guo, X. Wu, and Y. Li, "On the lower bound of reconstruction error for spectral filtering based privacy preserving data mining," in Proceedings of the 10th European Conference on Principles and Practice of Knowledge Discovery in Databases (PKDD'06), Berlin, Germany, 2006.
- [9] K. Muralidhar, R. Parsa, and R. Sarathy, "A general additive data perturbation method for database security," Management Science, vol. 45, no. 10, pp. 1399–1415, 1999.