

**Privacy Preserving Data Mining Using Random Rotation Based Data
Perturbation Technique**Asha Katharotiya¹ Pankaj Kawadkar²¹Department of Computer science & Engineering²H.O.D. computer science dept., PIES, Bhopal,
^{1,2}Patel Institute Of Engineering & Science, Bhopal, India

Abstract---To preserve privacy of data, privacy preserving data mining is the study of valid mining patterns and models which mask private information. There are many privacy preserving data mining techniques which have been studied. One crucial concept about existing data mining privacy preserving techniques are suitable and designed for static databases and not suitable for data streams. Recently, data streams are introduced as a new type of data which differ from traditional static data. Various features of data streams are: with time, data distribution changes constantly; data is having time preferences; amount of data is extensive; flow of data with fast speed; requirement of immediate response. When has been modified, it would be necessary to rescan whole database, so it leads to more computation time and inability to respond the user fastly. Further, it is observed that accuracy of data is decreases when transformation is carried out on data. So, there has been need to develop the system which preserve privacy along with accuracy. So privacy preserving on data stream mining is very crucial issue.

Keywords: Privacy; Data Streams; K-means clustering

I. INTRODUCTION

In recent years, data mining is shown as a powerful data analysis tool has made remarkable contributions in many areas and has the wide applications viewpoint. With the development of database technology and network technology, a large number of useful data, which contains so much individual privacy information, has been amassed in various fields, such as patient's condition information, preferences to customer, personal background information, account information etc. Once the information leaked, it will be unsafe to individual. If they give the actual data directly to the prospectors, it will predictably produce private information disclosure. As the field of data mining technology extending, privacy disclosure problem becomes worse, causing the attention of phases of industry and social. So, how to do data mining under the circumstances of privacy preserving has become a hot spot in data mining, so privacy preserving data mining (PPDM) is introduced. Securing against unauthorized accesses has been a long-term objective of the database security, and the government research statistical agencies and research community. Solutions to such a difficulty require combining several techniques and mechanisms. In a situation where data have different sensitivity stages, this data may be classified at different levels, and it has made available only to those subjects with an appropriate consent.

Technique of Clustering is a well-known problematic situation in statistics and engineering, namely, how to arrange a set of measurements into a number of clusters. Clustering is an significant area of application for a variety of areas including data mining, vector quantization and statistical data analysis. The problem has been framed in various ways in the, pattern recognition, machine learning, optimization and statistics literature. The fundamental clustering problem is that of assembling together data items that are related to each other. Given a set of data items, clustering algorithms can group analogous items together. Clustering has many applications, such as analysis of customer behavior, targeted marketing, forensics, and bioinformatics.

By mining sensitive characteristics from the original database, Reconstruction based approaches generate privacy aware database. These approaches have been generated smaller side effects in database than heuristic approach. Reconstruction based methods perturb the original data to achieve privacy preserving. The perturbed data would meet the two circumstances. First, an attacker cannot determine the real original data from the issuance of the alteration data. Secondly, the altered data is still preserving some statistical properties of the original data, namely some of the information derived from the partial data are equivalent to data acquired from the original information. Perturbing the data for preventing privacy is very fertile technique used by many researchers. It is also capable to reconstruct the distributions at an cumulative level in order to perform the mining.

There are three types of data perturbation approaches: Rotation Perturbation, Projection Perturbation and Geometric Data Perturbation.

II. ROTATION PERTURBATION

In order to preserve data privacy, a rotation perturbation method is used where the data matrix is multiplied by a random rotation matrix before publishing. There is one advantage of this approach is that it can preserve the geometric properties of the data matrix, so few categories of classifiers which are based on the geometric properties of the data can be achieved comparable accuracy on the altered data as that on the original data. Suppose original dataset has column and N records represented as $X_{d \times N}$, the rotation perturbation of the dataset X will be defined as $G(X) = RX^{[1]}$, here, random rotation orthonormal matrix is $R_{d \times d}$. A key feature of rotation transformation is preserving the Euclidean distance, geometric shape in a multi-dimensional space and inner product.

PCA (Principal component analysis) is a technique that is used to decompose the multidimensional data into lower dimensions. PCA assumes that all the inconsistency in a process should be used in the analysis so it becomes difficult to distinguish the crucial variable from the less vital. Subsequently, Principle component analysis

replaces the original variables of a dataset with the small number of uncorrelated variables called the principal component [3].

PCA is a standard tool in modern data analysis task in various fields from neuroscience to computer graphics. Since it is a simple, non-parametric method for extracting relevant information from mystifying datasets. With minimal effort PCA provides a road map for how to reduce a composite dataset to a lower dimension to reveal the sometimes simplified and hidden structures that often underlie it. Principal Component Analysis (PCA) is suitable for transforming the multidimensional data into lower dimensions. It is a standard tool in modern data analysis. PCA assumes that all the inconsistency in a process should be used in the analysis so it becomes challenging task to differentiate the important variable from the less important for that PCA is most appropriate for usual distributions (where linear Principle Component Analysis approach provides the best possible solution). Accordingly, Principle Component Analysis replaces the original variables of a dataset with a lesser number of uncorrelated variables called as the "principal components". If the original dataset of dimension D contains highly associated variables, then there is an operational dimensionality exist as, $d < D$, that explains all of the data. The presence of only a few components of d makes which is easier to label each dimension with an intuitive meaning. Furthermore, it is more effective to operate on fewer variables in subsequent analysis.

III. RESULTS AND DISCUSSION

In order to evaluate the clustering accuracy, Series of trials were performed over different sliding window size (w). Our evaluation approach focused on the inclusive quality of generated clusters after dataset perturbation. Experiment was based on following steps:

1. In MOA framework, Set up each dataset as stream.
2. To evaluate measures and cluster membership matrix, defines sliding window (w) over the data stream.
3. Altered all the occurrences in sliding window by applying proposed data perturbation method to protect the sensitive characteristic value.
4. To find the clusters for our performance evaluation, K-Means clustering algorithm is used. Our selection was influenced by (a) K-Means is one of the best known clustering Algorithm and it is also scalable.
 (b) The Number of cluster found from original and perturbed dataset was taken as a measure of cluster.
 Make Comparison that how closely each cluster in the perturbed dataset matches its corresponding cluster in the Original dataset. By computing the F-measure, we expressed the quality of the generated cluster.

To measure accuracy while protecting sensitive data, experiments were performed. Here we have presents two different results, one is analogous to clustering accuracy in terms of membership matrix which was manually plagiaristic from clustering result and another represent the equivalent graph for $F1_P$ (precision) and $F1_R$ (Recall) measures.

Dataset Name	Total instances	Instances processed	Attributes protected
Account Management	42210	43k	Balance, Age Duration

Table 1.1: Dataset configuration to determine accuracy based on Membership Matrix

To determine the accuracy of our proposed method, Table 1.1 shows datasets configuration. To determine set of 3 and 5 clusters using K-Means clustering algorithm, We configured each dataset.

Table 1.2, 1.3 shows the membership matrix acquired while clustering the perturbed attributes of Account Management dataset respectively. Each Matrix representing 3 and 5 clusters scenario for true dataset and discompose dataset.

True dataset clustering provides information about no. of instances are actual classified in each cluster whereas perturb dataset clustering showing result of accurate assignments after attributes data perturbation and percentage of accuracy achieved.

Table 1.2: Resultant accuracy of 5 Cluster

Dataset	Attributes	No. of Cluster	Stream Data	K-Means
Bank Management	Age	5	2000	85.21%
	Balance			89.39%
	Duration			86.81%
	Age	3000	3000	82.96%
	Balance			83.64%
	Duration			81.02%

Table 1.3: Resultant accuracy of 3 Cluster

Dat as et	Attribute s	No. of Cluster Cluster	Stream Data	K-Means
Bank Management	Age	3	2000	88.13 %
	Balance			92.60 %
	Duration			89.42%
	Age	3000	3000	85.59%
	Balance			91.22%
	Duration			89.64%

--	--	--	--	--

For each modified attribute, Results are presented in terms of graphs. Here each graph comprises the measure we obtained when original data is processed without applying privacy preserving method and K-Means is applied in order to evaluate both cases by keeping number of clusters fix ($K=5$, $K=3$), when data is undergone through our proposed privacy preserving method. In defined sliding window size, Instances are processed. Here we representing the accuracy of our method by calculating the precision of individual cluster. F1_R measure determine the recall of system, which take into account the clustering measure provided with MOA framework. We focused on two important measures F1_R and F1_P. F1_P.F1_P measure determine the precision of system by considering the precision of individual cluster. F1_R measure determine the recall of system, which take into account the recall of each cluster

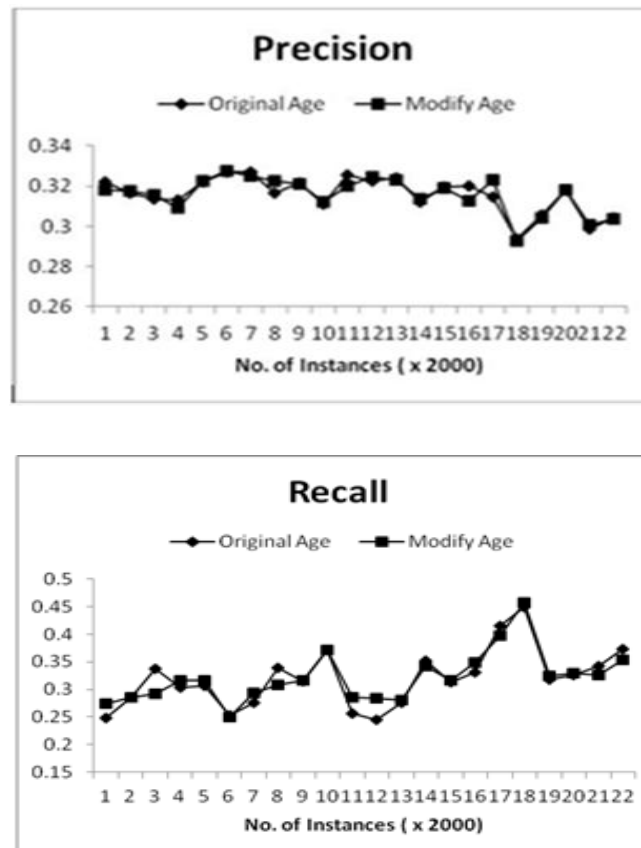
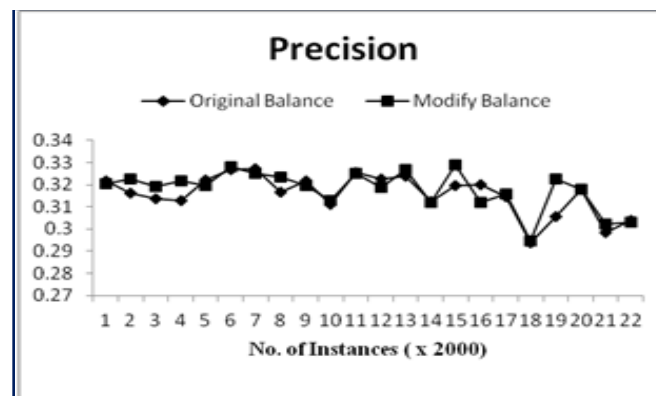


Fig 1.1: Accuracy on attribute Age in Bank Management with 5-Cluster



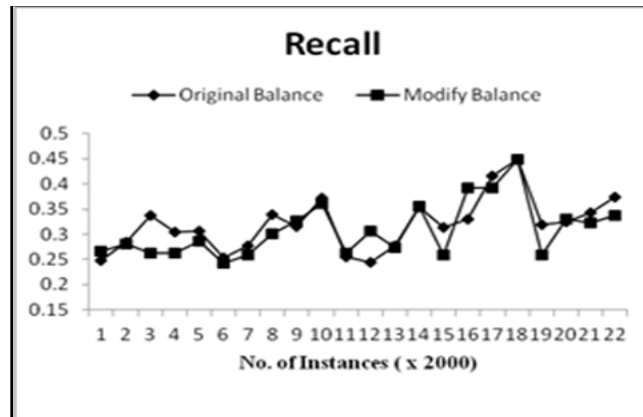


Fig.1.2: Accuracy on attribute Balance in Bank Management with 5-Cluster

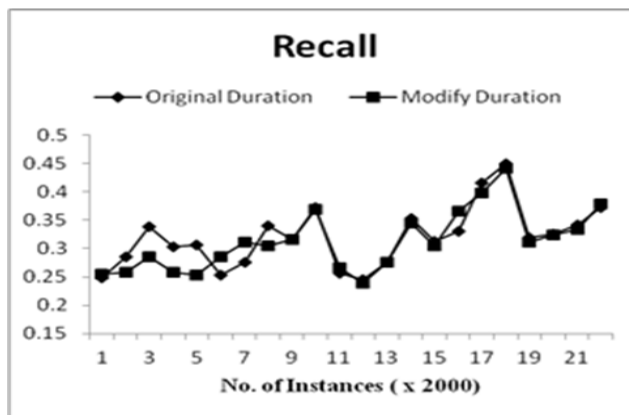
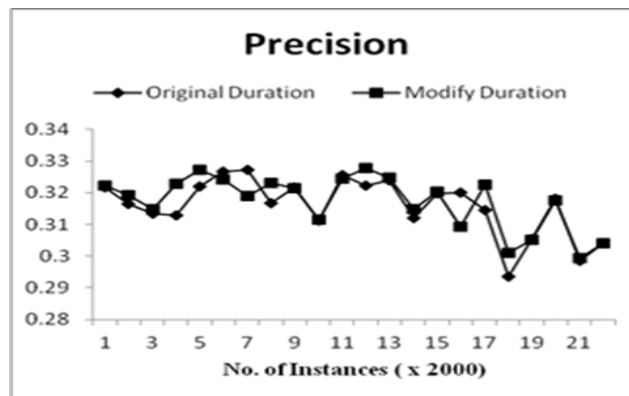


Fig 1.3: Accuracy on attribute Duration in Bank Management with 5-Cluster

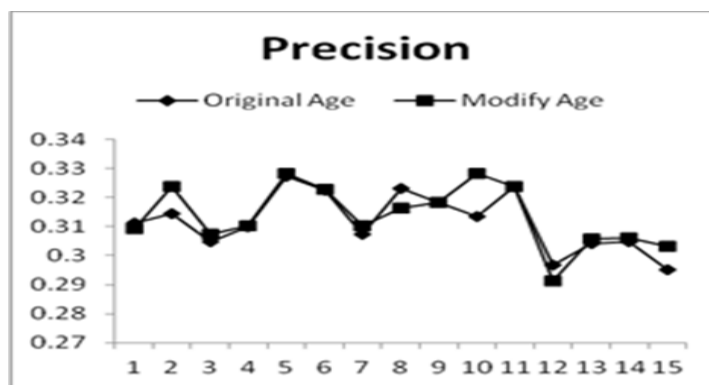
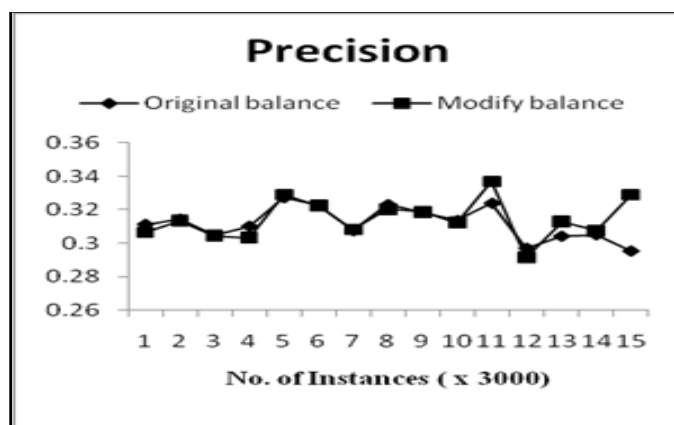


Fig.1.4: Accuracy on attribute Age in Bank Management with 5-Cluster



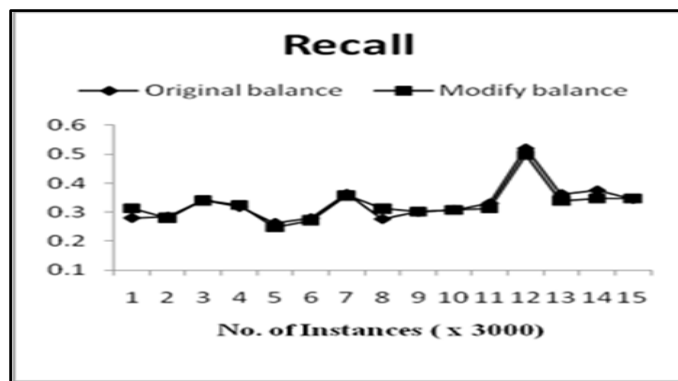


Fig. 1.5: Accuracy on attribute Balance in Bank Management with 5-Cluster

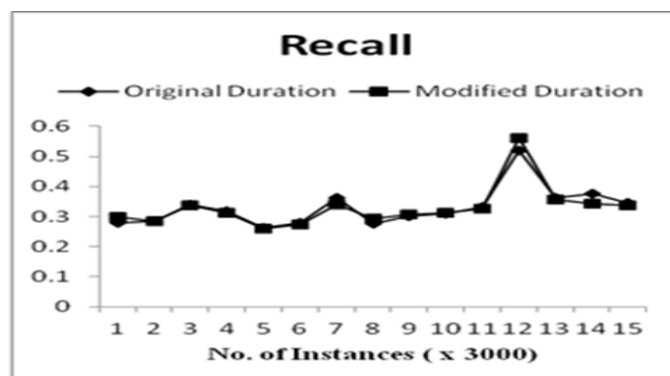
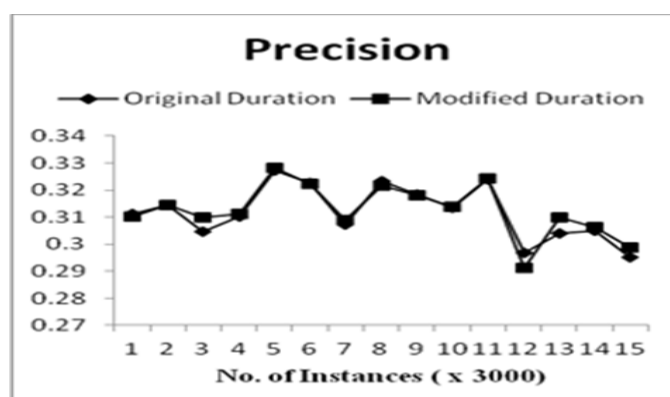


Fig. 1.6: Accuracy on attribute Duration in Bank Management with 5-Cluster

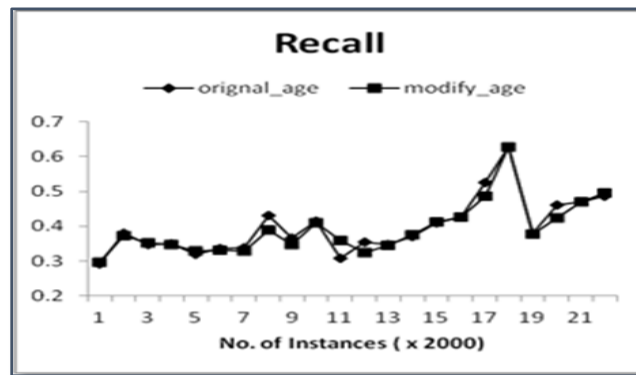
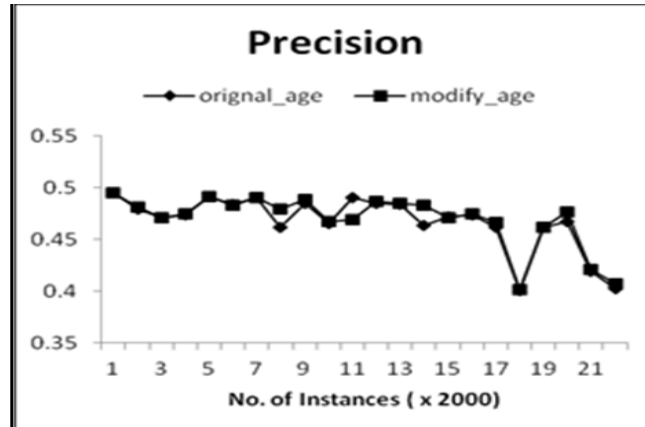
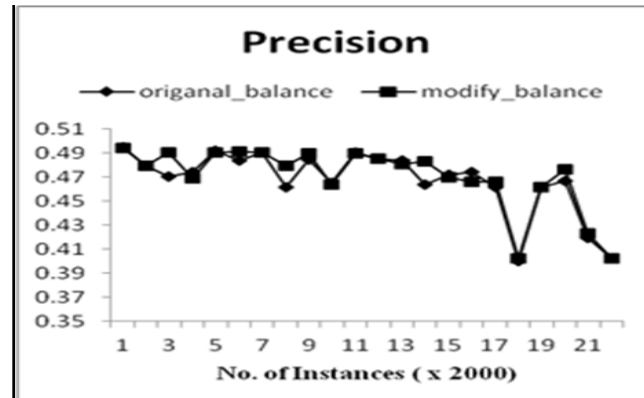


Fig. 1.7: Accuracy on attribute Age in Bank Management with 3-Cluster



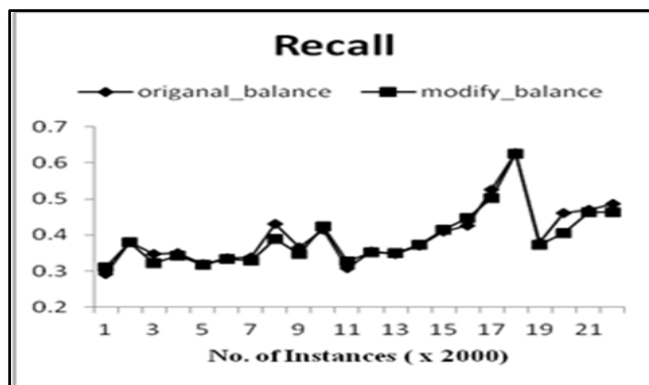


Fig. 1.8: Accuracy on attribute Balance Bank Management with 3-Cluster

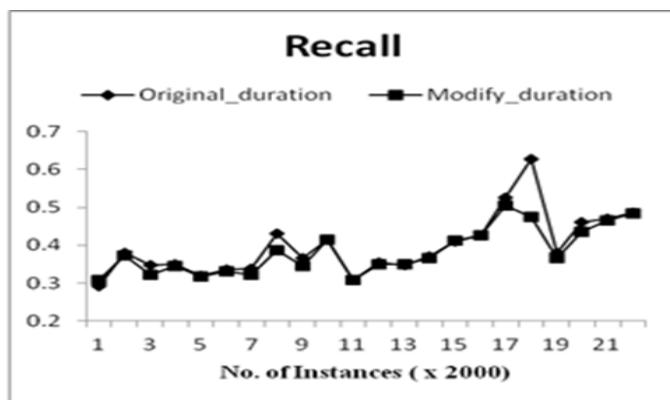
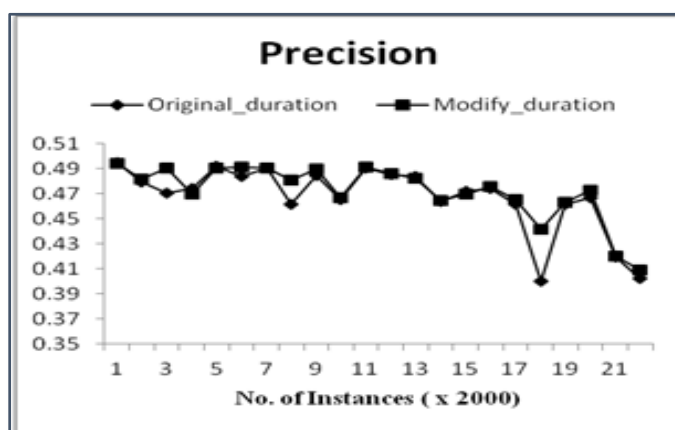


Fig. 1.9: Accuracy on attribute Duration in Bank Management with 3-Cluster

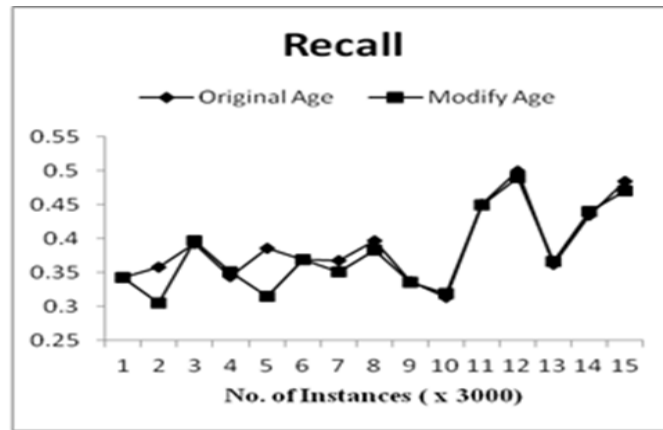
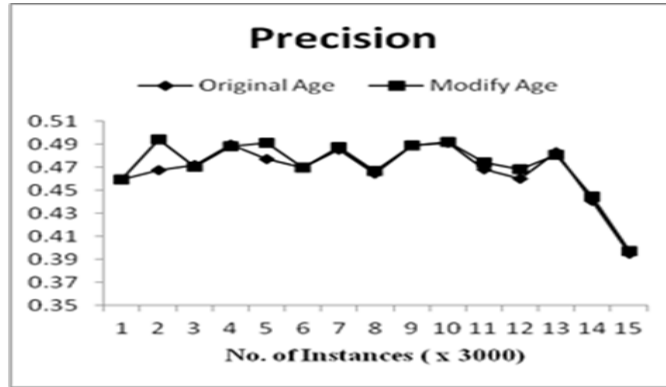
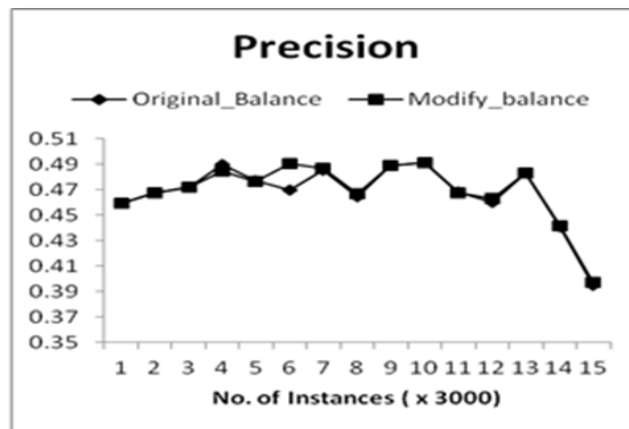


Fig. 1.10: Accuracy on attribute Age in Bank Management with 3-Cluster



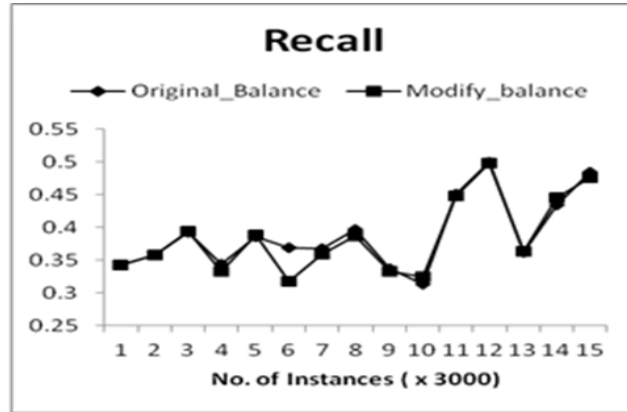


Fig. 1.11: Accuracy on attribute Balance in Bank Management with 3-Cluster

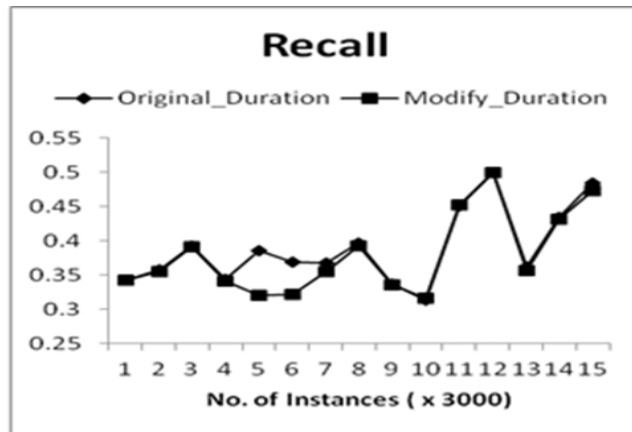
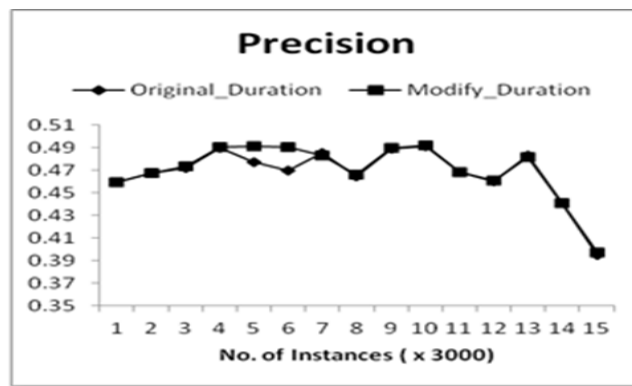


Fig 1.12: Accuracy on attribute Balance in BankManagement with 3-Cluster

IV.CONCLUSION

While presenting on a publicly accessible place like internet, the proposed method can be used to hide sensitive information. The proposed privacy preserving prototype has been successfully implemented in java under Windows 7 operating system and evaluated using Massive Online Analysis (MOA). The arrived results were more substantial and promising.. Additionally, the proposed model can be used to multi party cooperative clustering development. Some of the results of earlier works have been shown, accuracy sometimes suffers as a result of security. However in the proposed method, the accuracy has been conserved and in some cases, the accuracy was almost equal to that of original data set.

REFERENCES

- [1] Majid Bashir Malik And M. Asger Ghazi And Rashid Ali ;“Privacy Preserving Data Mining Techniques: Current Scenario And Future Prospects”; Third International Conference On Computer And Communication Technology; 978-0-7695-4872-2/12 \$26.00 © 2012 Ieee
- [2] Hitesh Chhinkaniwala And Dr. Sanjay Garg “Privacy Preserving Data Mining Techniques: Challenges & Issues” In Proceedings Of International Conference On Computer Science & Information Technology, Cslt – 2011,P.609
- [3] Chirag N. Modi,UdaiPratapRao And DhirenR.Patel “Maintaining Privacy And Data Quality In Privacy Preserving Association Rule Mining”, In 2010 Second International Conference On Computing, Communication And Networking Technologies
- [4] W.T. Chembian¹, Dr. J.Janet, “A Survey On Privacy Preserving Data Mining Approaches And Techniques”,In Proceedings Of The Int. Conf. On Information Science And Applications Icisa 2010,6 February 2010, Chennai, India
- [5] Xiaolin Zhang And Hongjing Bi; “Research On Privacy Preserving Classification Data Mining Based On Random Perturbation”; International Conference On Information, Networking And Automation (Icina); 978-1-4244-8106-4/\$26.00 © 2010 Ieee
- [6] Ching-Ming Chao, Po-Zung Chen And Chu-Hao Sun ;“Privacy-Preserving Classification Of Data Streams”; Tamkang Journal Of Science And Engineering; Vol. 12, No. 3, Pp. 321_330 (2009)
- [7] M. Naga Lakshmi, K Sandhya Rani;” Privacy
Issn: 2319-8753 Vol. 2, Issue 9, September 2013
- [8] Albert Bifet, Geoff Holmes, Bernhard Pfahringer, Philipp Kranen, Hardy Kremer, Timm Jansen, And Thomas Seidl
“Moa: Massive Online Analysis, A Framework For Stream Classification And Clustering.”
- [9] Vassilios S. Verykios, Elisa Bertino, Igor NaiFovino,LoredanaParasilitiProvenza, YucelSaygin, YannisTheodoridis
“State-Of-The-Art In Privacy Preserving Data Mining”
- [10] Haisheng Li;” Study Of Privacy Preserving Data Mining” ;Third International Symposium On Intelligent Information Technology And Security Informatics; 978-0-7695-4020-7/10 \$26.00 © 2010 Ieee
- [11] Jian Wang ,YongchengLuo ,Yan Zhao Jiajin Le;” A Survey On Privacy Preserving Data Mining”; 2009 First International Workshop On Database Technology And Applications; 978-0-7695-3604-0/09 \$25.00 © 2009 Ieee
- [12] MohammadrezaKeyvanpour, SomayyehSeifiMoradi;” Classification And Evaluation The Privacy Preserving Data Mining Techniques By Using A Data Modification–Based Framework”; International Journal On Computer Science And Engineering (Ijcse); Issn : 0975-3397 Vol. 3 No. 2 Feb 2011
- [13] Keke Chen Ling Liu;” A Random Rotation Perturbation Approach To Privacy Preserving Data Classification ”
- [14] S. Kasthuri, T. Meyyappan;” Detection Of Sensitive Items In Market Basket Database Using Association Rule Mining For Privacy Preserving”; Proceedings Of The 2013 International Conference On Prim, February 21-22; 978-1-4673-5845-3/13/\$31.00©2013 Ieee
- [15] Nikunj H. Domadiya;” Hiding Sensitive Association Rules To Maintain Privacy And Data Quality In Database”; 978-1-4673-4529-3/12/\$31.00_C 2012 Ieee
RahenaAkhter, RownakJahanChowdhury, Keita Emura, Tamzida Islam, Mohammad ShahriarRahman, NusratRubaiyat;” Privacy-Preserving Two-Party K- Means Clustering In Malicious Model”; 2013 Ieee 37th Annual Computer Software And Applications Conference Workshops; 978-0-7695-4987-3/13 \$26.00© 2013 Ieee
- [17] Jaideep Vaidya , BasitShafiq ;” A Random Decision Tree Framework For Privacy-Preserving Data Mining”; 1545-5971/13/\$31.00 © 2013 Ieee
2013