



Privacy Preserving Data Mining using Attribute Encryption and Data Perturbation

Meenakshi Vishnoi¹, Seeja K.R²

Department of Computer Science

Jamia Hamdard University
New Delhi, India

¹meenakshi.vishnoi85@gmail.com

²seeja@jamiahamdard.ac.in

ABSTRACT

Data mining is a very active research area that deals with the extraction of knowledge from very large databases. Data mining has made knowledge extraction and decision making easy. The extracted knowledge could reveal the personal information, if the data contains various private and sensitive attributes about an individual. This poses a threat to the personal information as there is a possibility of misusing the information behind the scenes without the knowledge of the individual. So, privacy becomes a great concern for the data owners and the organizations as none of the organizations would like to share their data. To solve this problem Privacy Preserving Data Mining technique have emerged and also solved problems of various domains as it provides the benefit of data mining without compromising the privacy of an individual. This paper proposes a privacy preserving data mining technique the uses randomized perturbation and cryptographic technique. The performance evaluation of the proposed technique shows the same result with the modified data and the original data.

Keywords

Data mining, Privacy Preserving, Data Perturbation, Cryptography.

Council for Innovative Research

Peer Review Research Publishing System

Journal: INTERNATIONAL JOURNAL OF COMPUTERS & TECHNOLOGY

Vol 6, No 3

editor@cirworld.com

www.cirworld.com, member.cirworld.com



1. INTRODUCTION

Data mining tool has made the process of extraction of knowledge and information very easy. Most of the organizations now [1][2] depend on data mining results for providing better services, achieving greater profit, and better decision-making. For these purposes organizations collect huge amount of data [3]. For example, to achieve the big profit and to apply the best business strategies business organizations [4] collect data about the consumers for marketing purposes, to provide the best treatment and medical research medical organizations collect medical records. As data mining includes large database which may consist of some sensitive and private information which can lead to the information loss behind the scenes of the data owner and this is one of the major problem. So, none of the organizations and data owner would like to share the data due to data loss. To solve this problem the concept of Privacy Preserving data mining (PPDM) has emerged as an active area of research which preserves both private and sensitive data. PPDM has provided the benefit of data sharing without the misuse of data and also provides the confidentiality of the data and data mining result. Various techniques such as randomization, perturbation, anonymization, swapping, etc are available for PPDM. If only perturbation is done on the data for privacy purpose, it could lead to the information loss for very critical dataset [5]. Thus; we need to work towards minimizing both privacy loss and information loss. We have proposed here an approach of perturbation technique combined with the cryptography. In perturbation some noise has been added to the original data is modified the data and then the attribute name has been encrypted in such a manner that the result obtained from original data and the modified data remains the same providing the better accuracy.

2. RELATED WORK

Various techniques have been developed over past few years for PPDM. The data transformation based approach modifies sensitive data in such a way that it loses its sensitive meaning. In this process statistical properties of interest can be retained but exact values cannot be determined during the mining process [3]. Various data modification techniques such as noise addition [1] it applies random noise to the data sets without considering different privacy requirements of the different Users. One interesting reconstruction approach is proposed by Kargupta et al. in [3]. By using random matrix properties, Kargupta et al. [8] was able to reconstruct the approximate original data by separating the random noise from it. Data shuffling [7] technique which was used to maintain the confidentiality of the data but retains the analytical value of the confidential data. Data transformation [3] in which data is transformed preserving the privacy of sensitive attribute. Role based access of data [5] where each user has control access to data depending on their role. Data masking [6] technique provides protection of sensitive data and hiding of confidential data by modifying the sensitive data to create life-like false values. In Cryptographic techniques [3] the data is encrypted with encryption methods and a set of protocols are used to allow the data mining operation.

2.1 Randomization perturbation technique

In randomization perturbation approach the privacy of the data can be protected by perturbing [9] sensitive data with randomization algorithms before releasing it to the data miner. The perturbed data version is send to the miner to mine the patterns. The algorithm is chosen in such a manner that sensitive data is modified and it remains no longer a sensitive data and preserving the confidentiality. The original data is distorted through adding the noise component to the data which is obtained through randomization. Here each individual entry is added with the noise component. In a set of data records denoted by $X = \{x_1 \dots x_N\}$. For each record x_i a noise component is added which is drawn from the randomization method. The noise components are $y_1 \dots y_N$. The distorted records are $x_1 + y_1 \dots x_N + y_N$. So, this new record is denoted as $z_1 \dots z_N$. This method is very simple and also provides the confidentiality of data.

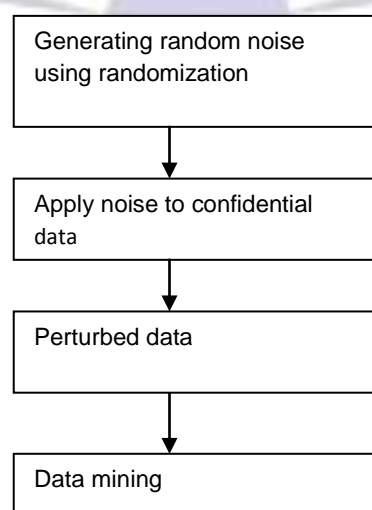




Fig1: Framework for perturbation technique

2.2 Cryptographic technique

Cryptography, the science of communication and computing in the presence of a malicious [3] adversary extends from the traditional tasks of encryption and authentication to protocols for securely distributing computations among a group of mutually distrusting parties. Cryptographic technique is used to encrypt the data to provide the security of data and hiding the individual information. In this paper we have used symmetric algorithm. This algorithm is faster and more secure. A DES algorithm has been used as a symmetric algorithm. We have used this algorithm to encrypt the attribute name of the dataset for preserving the privacy. The key generated is kept with the owner. The data miner has to work on the perturbed data and the encrypted attribute.

3. PROPOSED APPROACH

In this model we have generated a random noise using randomization method. This noise is then added to the original data. The data is modified in such a way that it becomes hard for third party to guess the original data. The field names of the dataset are encrypted and the data miner works on the perturbed and encrypted data.

3.1 ALGORITHM

Step 1

A random noise y_i is generated using randomized method. The original data is distorted using the noise component y_i .

Step 2

This noise component generated in step 1 is added to the original data. Once the noise component y_i is added the data becomes $z_1 = x_1 + y_1$

Step 3

Using cryptographic technique all attributes names is encrypted. Symmetric key encryption technique i.e. DES algorithm has been used to encrypt the attributes name.

Step 4

The key generated during encryption process remains with the sender and is not received by the miner

Step 5

The perturbed and encrypted data is then send to the data miner where the miner applies the mining algorithm. It can be any mining algorithm.

Step 6

The result generated is then send back to the owner in the encrypted format.

Step 7

Data owner after receiving the result decrypts the data using the key which was generated earlier during encryption process and removes the noise.

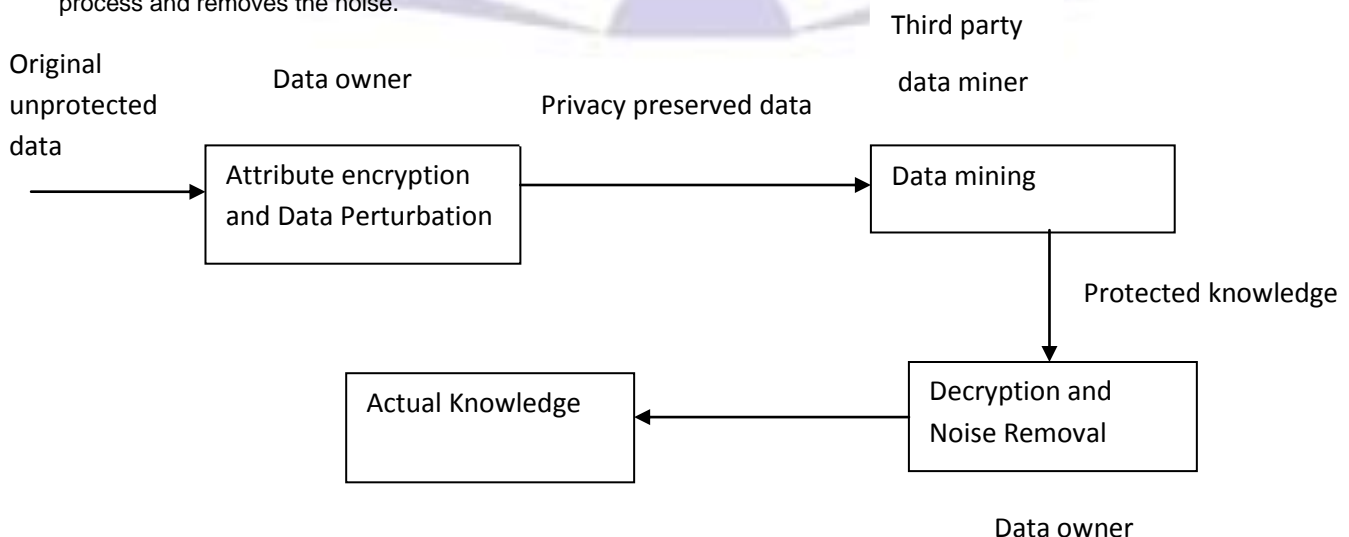




Fig 2: Framework for sharing the data using encryption and perturbation technique

3.2 Example

Let's take an example of a dataset Weather which is taken from UCI repository. This table contains 5 attribute and 14 records as shown in table 1. The dataset has been modified according to the proposed approach as follows.

Table 1: Original dataset of Weather taken from UCI repository

Outlook	Temperature	Humidity	Windy	Play
Sunny	85	85	False	No
Sunny	80	90	True	No
Overcast	83	86	False	Yes
Rainy	70	96	False	Yes
Rainy	68	80	False	Yes
Rainy	65	70	True	No
Overcast	64	65	True	Yes
Sunny	72	95	False	No
Sunny	69	70	False	Yes
Rainy	75	80	False	Yes
Sunny	75	70	True	Yes
Overcast	72	90	True	Yes
Overcast	81	75	False	Yes
Rainy	71	91	True	No

The above dataset will be modified according to the proposed approach.

- 1) The noise is generated using randomization method and original data is perturbed with the generated noise.
- 2) The attribute names are encrypted using DES.
- 3) The modified dataset is shown in table 2.

Table 2: perturbed and encrypted dataset

Xyfmmm	WQQulo	5b8gaQ	oOqwEv	FoHP/A
Sunny	110	110	False	No
Sunny	105	115	True	No
Overcast	108	111	False	Yes
Rainy	95	191	False	Yes
Rainy	93	105	False	Yes
Rainy	90	105	True	No
Overcast	89	90	True	Yes
Sunny	97	120	False	No
Sunny	94	95	False	Yes



Rainy	100	105	False	Yes
Sunny	100	95	True	Yes
Overcast	97	115	True	Yes
Overcast	106	100	False	Yes
Rainy	96	116	True	No

The above dataset is perturbed and encrypted dataset which is then send to the third party data miner who performs the data mining. The knowledge extracted is then send to the owner who then decrypts and removes noise from it.

It is found that the knowledge extracted from original data as well as Privacy Preserved data are found to be same as shown in table 3.

Table 3: Table displaying the result of both original and perturbed dataset.

Original Data	No. of leaves	Size of tree	Correctly classified instances	Incorrect instances
	5	8	9	5
Privacy Preserved Data	5	8	9	5
Knowledge Extracted from original data Knowledge extracted from Privacy Preserved data are same				

4. RESULTS AND DISCUSSIONS

During this whole process data remains preserved at the owner site. The privacy is achieved because the data remains perturbed and encrypted during the mining process and key of encryption remains with the sender. Here we have taken a diabetes dataset from UCI repository. The dataset consist of 768 instances and 9 attributes. All the data except the class values are perturbed and all of the attribute names are encrypted using DES algorithm which has been implemented in java and data perturbation is done using randomized method. In this approach Weka tool has been used to obtain the result. Two types of result are displayed.

- i) Pre-processing of attributes.
- ii) Decision tree generated from J48 classifier for both perturbed and original dataset.



No.	1:preg	2:plas	3:pres	4:skin	5:triu	6:mass	7:pedi	8:age	9:class
1	6.0	148.0	72.0	35.0	0.0	33.6	0.627	50.0	tested_positive
2	1.0	85.0	66.0	29.0	0.0	36.6	0.701	31.0	tested_negative
3	8.0	183.0	64.0	0.0	0.0	23.5	0.672	32.0	tested_positive
4	1.0	89.0	66.0	23.0	0.0	28.1	0.167	21.0	tested_negative
5	0.0	137.0	40.0	35.0	168.0	43.1	2.288	33.0	tested_positive
6	5.0	136.0	74.0	0.0	0.0	25.6	0.201	30.0	tested_negative
7	3.0	78.0	50.0	32.0	88.0	31.0	0.248	26.0	tested_positive
8	10.0	115.0	0.0	0.0	0.0	0.0	0.134	29.0	tested_negative
9	2.0	197.0	70.0	45.0	340.0	30.5	0.158	53.0	tested_positive
10	8.0	125.0	96.0	0.0	0.0	0.0	0.232	54.0	tested_positive
11	4.0	130.0	92.0	0.0	0.0	27.8	0.191	30.0	tested_negative
12	10.0	168.0	74.0	0.0	0.0	38.3	0.337	34.0	tested_positive
13	10.0	139.0	80.0	0.0	0.0	27.1	1.442	57.0	tested_negative
14	1.0	189.0	62.0	23.0	846.0	30.1	0.398	59.0	tested_positive
15	5.0	166.0	72.0	39.0	175.0	25.8	0.587	51.0	tested_positive
16	7.0	100.0	0.0	0.0	0.0	30.0	0.484	32.0	tested_positive
17	0.0	118.0	84.0	47.0	230.0	45.8	0.551	31.0	tested_positive
18	7.0	107.0	74.0	0.0	0.0	29.0	0.254	31.0	tested_positive
19	1.0	103.0	30.0	36.0	83.0	43.5	0.183	33.0	tested_negative
20	1.0	115.0	70.0	30.0	96.0	34.6	0.529	32.0	tested_positive
21	3.0	126.0	88.0	41.0	235.0	39.3	0.704	27.0	tested_negative
22	8.0	99.0	84.0	0.0	0.0	25.4	0.388	50.0	tested_negative
23	7.0	186.0	90.0	0.0	0.0	29.8	0.451	41.0	tested_positive
24	9.0	119.0	80.0	35.0	0.0	29.0	0.283	29.0	tested_positive
25	11.0	143.0	94.0	33.0	148.0	36.6	0.254	51.0	tested_positive
26	10.0	125.0	70.0	26.0	115.0	31.1	0.205	41.0	tested_positive
27	7.0	147.0	76.0	0.0	0.0	29.4	0.257	43.0	tested_positive
28	1.0	97.0	66.0	15.0	140.0	23.2	0.487	22.0	tested_negative
29	13.0	145.0	82.0	39.0	110.0	22.2	0.245	57.0	tested_negative
30	5.0	117.0	92.0	0.0	0.0	24.1	0.337	30.0	tested_negative
31	5.0	109.0	75.0	26.0	0.0	36.0	0.546	60.0	tested_negative
32	3.0	158.0	76.0	36.0	245.0	31.6	0.851	28.0	tested_positive
33	3.0	88.0	58.0	11.0	54.0	24.8	0.267	22.0	tested_negative
34	6.0	92.0	92.0	0.0	0.0	29.9	0.188	28.0	tested_negative
35	10.0	122.0	78.0	31.0	0.0	27.8	0.512	45.0	tested_negative
36	4.0	103.0	62.0	12.0	102.0	24.0	0.966	33.0	tested_negative
37	11.0	138.0	76.0	0.0	0.0	33.2	0.42	35.0	tested_negative
38	9.0	102.0	76.0	37.0	0.0	32.9	0.665	40.0	tested_positive
39	2.0	90.0	68.0	42.0	0.0	38.2	0.502	27.0	tested_positive
40	4.0	111.0	72.0	47.0	207.0	27.1	1.38	56.0	tested_positive
41	3.0	180.0	64.0	23.0	76.0	24.0	0.272	26.0	tested_negative

Fig 3: original dataset of diabetes taken from UCI repository

No.	1:p2d1	2:162dH	3:161eH	4:161eH	5:162dH	6:162dH	7:161eH	8:162dH	9:161eH
1	16.0	184.0	82.0	46.0	0.0	43.6	10.627	60.0	tested_positive
2	11.0	95.0	76.0	29.0	0.0	36.6	10.701	41.0	tested_negative
3	18.0	193.0	74.0	16.0	0.0	33.3	10.672	42.0	tested_positive
4	11.0	89.0	76.0	33.0	304.0	38.1	10.167	31.0	tested_negative
5	10.0	147.0	90.0	46.0	176.0	53.1	12.288	43.0	tested_positive
6	15.0	126.0	84.0	10.0	0.0	35.6	10.201	40.0	tested_negative
7	15.0	88.0	60.0	42.0	96.0	41.0	10.248	36.0	tested_positive
8	20.0	123.0	80.0	18.0	0.0	45.3	10.134	36.0	tested_negative
9	12.0	207.0	80.0	55.0	853.0	48.5	10.195	63.0	tested_positive
10	16.0	159.0	96.0	18.0	0.0	33.0	10.232	64.0	tested_positive
11	14.0	129.0	102.0	18.0	0.0	47.6	10.191	40.0	tested_negative
12	26.0	138.0	84.0	18.0	0.0	48.0	10.527	44.0	tested_positive
13	20.0	148.0	90.0	10.0	0.0	37.1	11.442	67.0	tested_negative
14	11.0	184.0	70.0	33.0	856.0	48.1	10.388	60.0	tested_positive
15	15.0	176.0	82.0	29.0	283.0	35.8	10.587	61.0	tested_positive
16	17.0	130.0	25.0	18.0	0.0	40.0	10.484	45.0	tested_positive
17	10.0	128.0	94.0	57.0	240.0	60.8	10.551	41.0	tested_positive
18	17.0	117.0	84.0	18.0	0.0	39.6	10.254	41.0	tested_positive
19	11.0	113.0	40.0	48.0	93.0	63.2	10.183	43.0	tested_negative
20	11.0	125.0	80.0	40.0	306.0	44.6	10.529	42.0	tested_positive
21	13.0	136.0	96.0	51.0	245.0	48.3	10.704	37.0	tested_negative
22	18.0	109.0	94.0	10.0	0.0	45.4	10.388	60.0	tested_positive
23	17.0	206.0	100.0	10.0	0.0	48.9	10.451	51.0	tested_positive
24	18.0	128.0	90.0	48.0	0.0	39.0	10.283	29.0	tested_positive
25	21.0	153.0	104.0	40.0	156.0	46.6	10.254	61.0	tested_positive
26	20.0	135.0	80.0	36.0	125.0	41.1	10.205	51.0	tested_positive
27	17.0	157.0	86.0	18.0	0.0	48.4	10.257	53.0	tested_positive
28	11.0	187.0	76.0	25.0	150.0	33.2	10.487	32.0	tested_negative
29	23.0	155.0	92.0	29.0	200.0	32.2	10.245	67.0	tested_negative
30	15.0	137.0	102.0	10.0	0.0	44.1	10.337	46.0	tested_negative
31	15.0	119.0	85.0	38.0	20.0	46.0	10.546	70.0	tested_positive
32	13.0	168.0	86.0	46.0	255.0	41.6	10.811	38.0	tested_positive
33	13.0	98.0	68.0	21.0	04.0	34.8	10.287	32.0	tested_negative
34	14.0	103.0	102.0	10.0	0.0	29.9	10.180	38.0	tested_negative
35	20.0	152.0	86.0	41.0	30.0	37.6	10.512	55.0	tested_negative
36	14.0	113.0	70.0	40.0	202.0	34.0	10.966	43.0	tested_negative
37	21.0	148.0	86.0	10.0	0.0	43.2	10.42	45.0	tested_negative
38	18.0	112.0	86.0	47.0	0.0	43.9	10.665	56.0	tested_positive
39	12.0	180.0	76.0	52.0	0.0	48.2	10.503	37.0	tested_positive
40	14.0	121.0	82.0	57.0	217.0	47.1	11.39	66.0	tested_positive
41	13.0	180.0	74.0	25.0	80.0	44.0	10.272	26.0	tested_negative

Fig 4: Perturbed and encrypted dataset

Below are the statistics which shows all the information of the values stored in each attribute of the dataset. The statistics differ according to the data type of the attributes. If the attribute is nominal, the list consists of each possible value for the attribute along with the number of instances that have that value. If the attribute is numeric, the list gives four statistics

describing the distribution of values in the data—the minimum, maximum, mean and standard deviation. The color coded statistics is displayed only if the attribute is nominal. There are two class values 'tested-positive' and 'tested-negative'. The red color is for 'tested-negative' class and blue is 'tested-positive' class. According to the range of values for each attribute the class values differ. For example for attribute 'preg' when the value is 0 or 1 the class value for tested_positive and tested_negative is 246 .

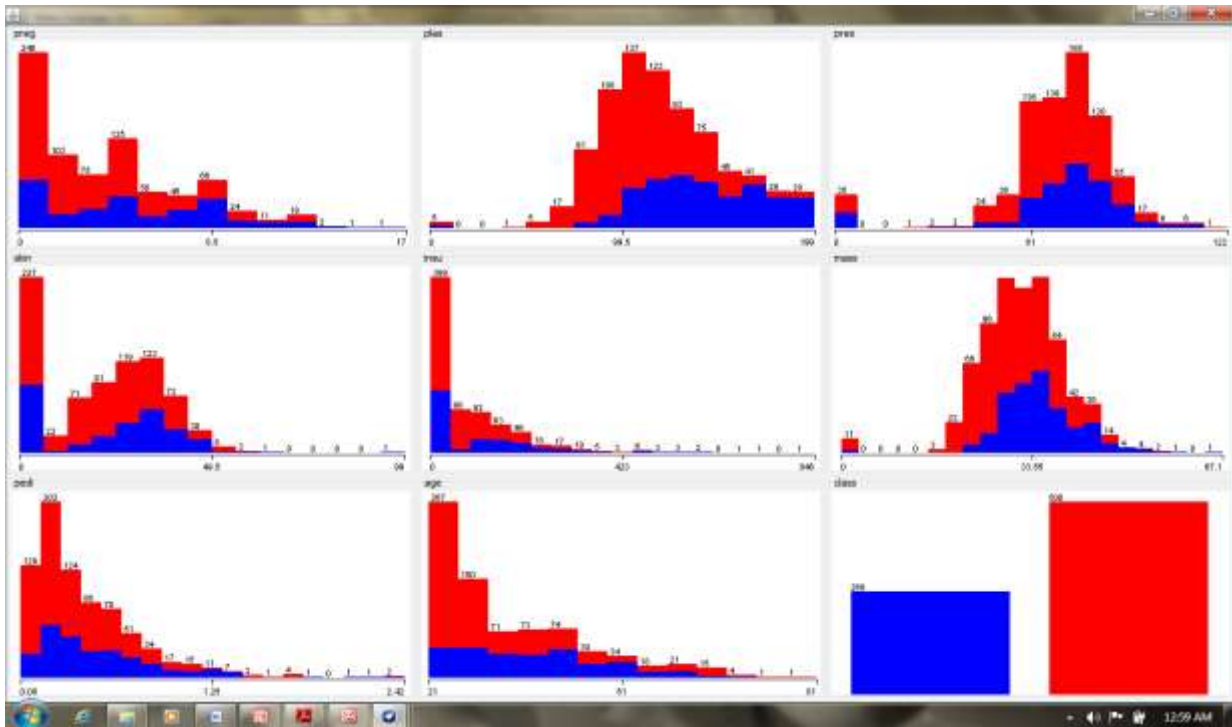


Fig 5: Result of original diabetes dataset

Below are the statistics which shows all the information of the values stored in each attribute of the dataset the values are perturbed and attribute names are encrypted the statistics generated is similar to the statistics of original dataset.

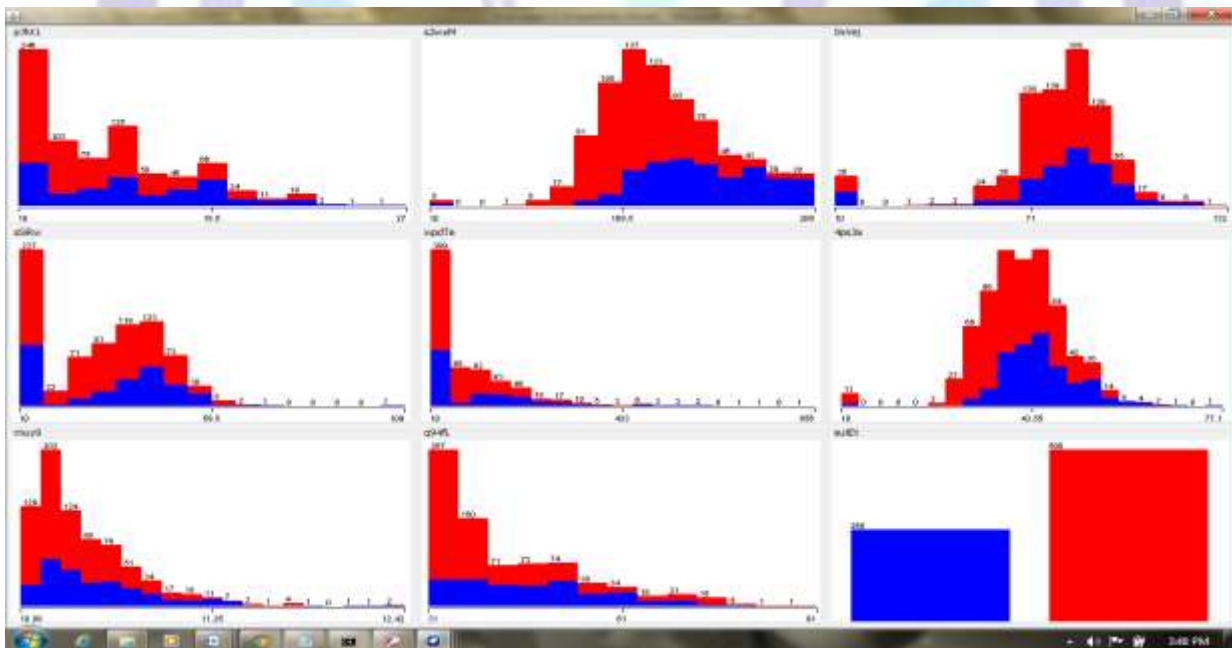


Fig 6: Result of perturb data

Below the table shows the result generated when actual and perturbed data are classified using J48 classifier.

Table 3: Table displaying the result of both original and perturbed dataset.

Actual Data	No. of leaves	Size of tree	Correctly classified instances	Incorrect instances
	20	39	567	201
Perturbed and Encrypted Data	20	39	567	201

Below is the tree generated after classifying the dataset in J48 classifier using Weka tool. The tree generated is a pruned tree .In J48 classifier. Here leaf nodes indicate which class an instance will be assigned to and should that node be reached. We can also see that the numbers are written in brackets after the leaf nodes indicate the number of instances assigned to that node, followed by how many of those instances that are incorrectly classified as a result.

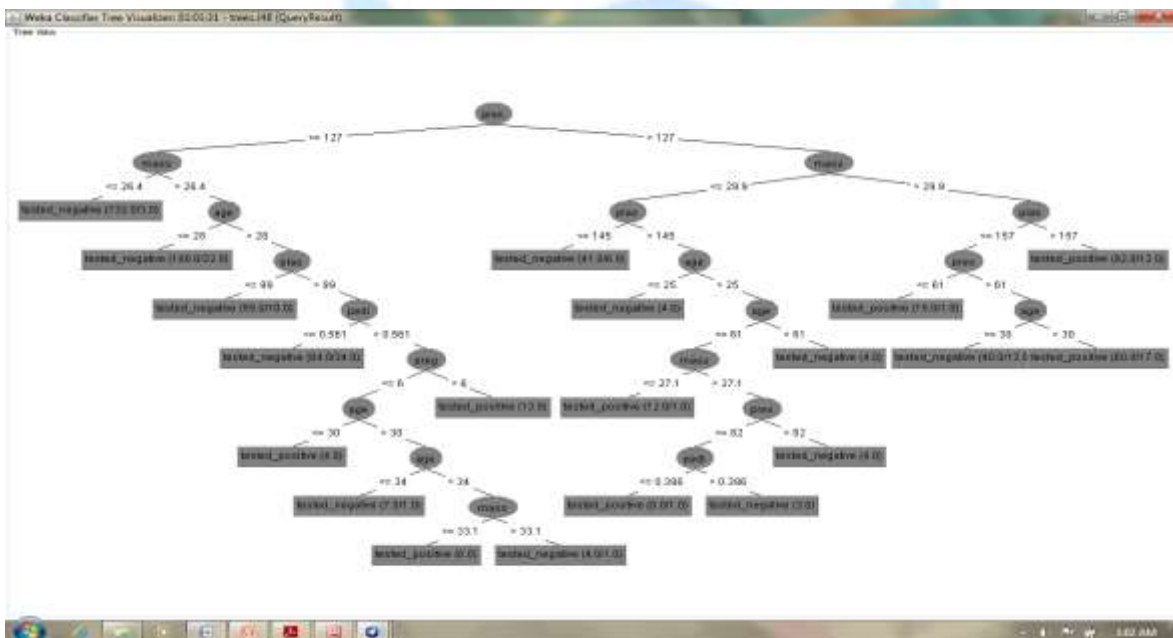


Fig 7: Decision tree generated from original data

Below is the tree generated after classifying the perturbed dataset in J48 classifier using the Weka tool. Here the nodes are the encrypted attribute name. And the tree generated is same as that of the tree generated using actual dataset.

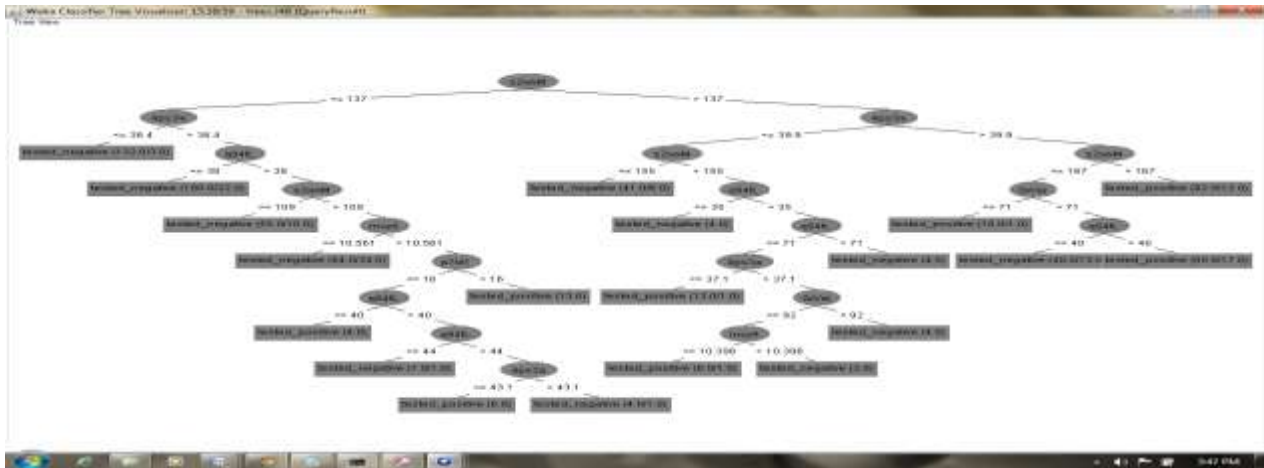


Fig 8: Decision tree of perturbed and encrypted dataset

5. CONCLUSION

Privacy preserving data mining is an emerging subfield of data mining. With the help of Privacy Preserving Data Mining sensitive data like medical data, banking data etc. can be protected without revealing the actual data to the third party data miner. In this paper, a privacy preserved data mining technique using data perturbation and encryption is proposed. The proposed approach is evaluated by using the weather dataset available in UCI repository. It is found that the knowledge extracted from original data as well as the protected data are same. Therefore the proposed approach can be used for mining sensitive data by protecting the individual privacy.

REFERENCES

- [1] R.Agarwal and R.Srikant, "Privacy preserving data mining", In Proceedings of the 19th ACM SIGMOD conference on Management of Data ,Dallas,Texas,USA, May2000.
- [2] J. Canny, "Collaborative filtering with privacy". In IEEE Symposium on security and privacy , pages 45-57 Oakland, May 2002.
- [3] P.Kamakshi , Dr.A.Vinaya Babu." Preserving Privacy and Sharing the Data in Distributed Environment using Cryptographic Technique on Perturbed data", In Journal of Computing , volume 2, April 2010.
- [4] The privacy preserving properties of random data perturbation techniques", In Proceedings of the 3rd IEEE International Conference on Data Mining, pages 99–106, Melbourne, Florida, November 19-22, 2003.
- [5] *Lalanthika Vasudevan , S.E. Deepa Sukanya, N. Aarthi**," Privacy Preserving Data Mining Using Cryptographic Role Based Access Control Approach", in the Proceedings of the International Multi Conference of Engineers and Computer Scientists 2008 Vol I.
- [6] An Oracle white paper, July 2010
- [7] Muralidhar K. and Sarathy R. " Data Shuffling- a new masking approach for numerical data." Management Science, forthcoming, 2006.
- [8] H. Kargupta, S. Datta, Q. Wang, K. Sivakumar, On the privacy preserving properties of random data perturbation techniques, in:
ICDM, IEEE Computer Society, 2003, pp. 99–106.
- [9] R.Agrawal, A.Evfimievski, R.Srikant, "Information sharing across private databases", In Proc.of ACM SIGMOD, 2003.
- [10] www.cs.waikato.ac.nz/ml/weka/
- [11] Mohammad Ali Kadampur, Somayajulu D.V.L.N.," A Noise Addition Scheme in Decision Tree for Privacy Preserving Data Mining", JOURNAL OF COMPUTING, VOLUME 2, ISSUE 1, JANUARY 2010.
- [12] Krishna Priya .J Geetha Mary. A," Perturbation of String Values", In IJCSIT Vol. 2 (3) , 2011