

# Privacy Preservation in Data Mining Through Noise Addition

**Md Zahidul Islam**

A thesis submitted in fulfillment of the requirements for the degree of  
Doctor of Philosophy



THE UNIVERSITY OF  
**NEWCASTLE**  
AUSTRALIA

School of Electrical Engineering and Computer Science  
University of Newcastle  
Callaghan  
New South Wales 2308  
Australia

November 2007

# Certificate of Originality

*I hereby certify that the work embodied in this thesis is the result of original research and has not been submitted for a higher degree at any other University or Institution.*

(Signed) \_\_\_\_\_

Md Zahidul Islam

# Acknowledgements

I would like to thank my supervisor A/Prof. Ljiljana Brankovic who is something more than just a supervisor to all her students. Whenever I was in a trouble she was there with her genuine suggestions and dependable directions. She introduced this research area to me. If I have learnt anything on how to do research then it is due to her wise supervision. She always led us to be independent researchers having high ethics and moral values.

I would also like to thank my co-supervisor Professor A.S.M. Sajeev for his support, encouragement and wisdom. I am also grateful to Dr Regina Berretta, Dr Michael Hannaford, Dr Alexandre Mendes, Professor Mahbub Hassan, Professor M. F. Rahman, Professor Mirka Miller and Professor Elizabeth Chang for their support and encouragement.

I would give my special thanks to my friends Helen, Mousa, Mouris and Tanya for their enormous moral support throughout my study. My thanks also to all Faculty and Staff members of the School of Electrical Engineering and Computer Science and all postgraduate students of the school during my study for being so kind and friendly to me.

Last but not least, I would like to thank my wife Moonmoon for her patience, care, support, trust, love and encouragement. My special thanks to my children Abdellah, Saifullah and Mahir for their love and support. I would like to thank my parents Harun and Nilu, my father in law, mother in law, sister and brother in law for their encouragement. They have been a very supportive family all the way.

*This thesis is gratefully dedicated to*

**My Family:**

**Moonmoon**, my wife

**Abdellah, Saifullah and Mahir**, my sons

**My Parents, My Sister, My Father in law, My Mother in law**  
and **All Relatives**

*for their patience, their unwavering support and their faith.*

*Say: “If the ocean were ink (wherewith to write out) the words of my Lord. Sooner would the ocean be exhausted than would the words of my Lord, even if we added another ocean like it, for its aid.” (Qur’an 18:109)*

# List of publications arising from this thesis

1. M. Z. Islam, and L. Brankovic, Privacy Preserving Data Mining: A Framework for Noise Addition to all Numerical and Categorical Attributes, In *Data Mining and Knowledge Discovery*. (In Preparation)
2. M. Z. Islam, and L. Brankovic, Privacy Preserving Data Mining: Noise Addition to Categorical Values Using a Novel Clustering Technique, In *IEEE Transactions on Industrial Informatics*, 2007. (Submitted on the 3rd September, 2007)
3. L. Brankovic, M. Z. Islam and H. Giggins, Privacy-Preserving Data Mining, Security, Privacy and Trust in Modern Data Management, Springer, Editors Milan Petkovic and Willem Jonker, ISBN: 978-3-540-69860-9, Chapter 11, 151-166, 2007.
4. M. Z. Islam, and L. Brankovic, DETECTIVE: A Decision Tree Based Categorical Value Clustering and Perturbation Technique in Privacy Preserving Data Mining, In *Proc. of the 3rd International IEEE Conference on Industrial Informatics*, Perth, Australia, (2005).
5. M. Z. Islam, and L. Brankovic, A Framework for Privacy Preserving Classification in Data Mining, In *Proc. of Australian Workshop on Data Mining and Web Intelligence (DMWI2004)*, Dunedin, New Zealand, CRPIT, **32**, J. Hogan, P. Montague, M. Purvis and C. Steketee, Eds., *Australian Computer Science Communications*, (2004) 163-168.
6. M. Z. Islam, P. M. Barnaghi and L. Brankovic, Measuring Data Quality: Predictive Accuracy vs. Similarity of Decision Trees, In *Proc. of the 6th International Conference on Computer & Information Technology (ICCIT 2003)*, Dhaka, Bangladesh, Vol.2, (2003) 457-462.
7. M. Z. Islam, and L. Brankovic, Noise Addition for Protecting Privacy in Data Mining, In *Proc. of the 6th Engineering Mathematics and Applications Conference (EMAC*

2003), Sydney, Australia, (2003) 457-462.

# List of other publications during the candidature

1. M. Alfalayleh, L. Brankovic, H. Giggins, and M. Z. Islam, Towards the Graceful Tree Conjecture: A Survey, In *Proc. of the 15th Australasian Workshop on Combinatorial Algorithms (AWOCA 2004)*, Ballina, Australia, (2004).

# Abstract

Due to advances in information processing technology and storage capacity, nowadays huge amount of data is being collected for various data analyses. Data mining techniques, such as classification, are often applied on these data to extract hidden information. During the whole process of data mining the data get exposed to several parties and such an exposure potentially leads to breaches of individual privacy.

This thesis presents a comprehensive noise addition technique for protecting individual privacy in a data set used for classification, while maintaining the data quality. We add noise to all attributes, both numerical and categorical, and both to class and non-class, in such a way so that the original patterns are preserved in a perturbed data set. Our technique is also capable of incorporating previously proposed noise addition techniques that maintain the statistical parameters of the data set, including correlations among attributes. Thus the perturbed data set may be used not only for classification but also for statistical analysis.

Our proposal has two main advantages. Firstly, as also suggested by our experimental results the perturbed data set maintains the same or very similar patterns as the original data set, as well as the correlations among attributes. While there are some noise addition techniques that maintain the statistical parameters of the data set, to the best of our knowledge this is the first comprehensive technique that preserves the patterns and thus removes the so called Data Mining Bias from the perturbed data set.

Secondly, re-identification of the original records directly depends on the amount of noise added, and in general can be made arbitrarily hard, while still preserving the original patterns in the data set. The only exception to this is the case when an intruder knows enough about the record to learn the confidential class value by applying the classifier. However, this is always possible, even when the original record has not been used in the training data set. In other words, providing that enough noise is added, our technique makes the records from the training set as safe as any other previously unseen records of the same kind.

In addition to the above contribution, this thesis also explores the suitability of prediction accuracy as a sole indicator of data quality, and proposes technique for clustering both categorical values and records containing such values.

# Contents

<b>List of Figures</b>	<b>xi</b>
<b>List of Tables</b>	<b>xiv</b>
<b>1 Introduction</b>	<b>1</b>
<b>2 Data Mining</b>	<b>5</b>
2.1 Introduction to Data Mining . . . . .	5
2.1.1 Definition . . . . .	5
2.1.2 Comparison with Traditional Data Analyses . . . . .	6
2.1.3 Data Mining Steps . . . . .	7
2.2 Data Mining Tasks . . . . .	9
2.3 Applications of Data Mining . . . . .	14
2.3.1 Usefulness in General . . . . .	14
2.3.2 Some Applications of Data Mining Techniques . . . . .	15
2.4 Privacy Issues Related to Data Mining . . . . .	17
2.5 Conclusion . . . . .	20
<b>3 Privacy Preserving Data Mining - A Background Study</b>	<b>21</b>
3.1 Classification Scheme and Evaluation Criteria . . . . .	21
3.2 Data Modification . . . . .	28
3.2.1 Noise Addition in Statistical Database . . . . .	28
3.2.2 Noise Addition in Data Mining . . . . .	32
3.2.3 Data Swapping . . . . .	39
3.2.4 Aggregation . . . . .	40
3.2.5 Suppression . . . . .	40
3.3 Secure Multi-Party Computation . . . . .	41
3.4 Comparative Study . . . . .	46
3.5 Conclusion . . . . .	46
<b>4 Class Attribute Perturbation Technique</b>	<b>48</b>
4.1 The Essence . . . . .	48
4.2 Noise Addition to Class Attribute . . . . .	50
4.3 The Experiment . . . . .	57

4.4	Conclusion . . . . .	63
<b>5</b>	<b>Non-class Numerical Attributes Perturbation Technique</b>	<b>71</b>
5.1	Introduction . . . . .	71
5.2	The <i>Leaf Innocent Attribute Perturbation Technique</i> . . . . .	74
5.3	The <i>Leaf Influential Attribute Perturbation Technique</i> . . . . .	76
5.4	The <i>Random Noise Addition Technique</i> . . . . .	77
5.5	Conclusion . . . . .	77
<b>6</b>	<b>Non-class Categorical Attributes Perturbation Technique</b>	<b>79</b>
6.1	Introduction . . . . .	79
6.2	Background . . . . .	80
6.3	An Overview of Existing Categorical Attribute Clustering Techniques . . . . .	82
6.4	<i>DETECTIVE</i> : A Novel Categorical Values Clustering Technique . . . . .	91
6.4.1	The Preliminaries . . . . .	91
6.4.2	<i>DETECTIVE</i> . . . . .	91
6.4.3	The Essence . . . . .	92
6.4.4	Illustration . . . . .	94
6.4.5	The Similarities . . . . .	95
6.4.6	The Difference . . . . .	96
6.4.7	<i>EX-DETECTIVE</i> . . . . .	97
6.5	CAPT: Categorical Attributes Perturbation Technique . . . . .	101
6.6	Experimental Results . . . . .	102
6.6.1	Experiments on <i>DETECTIVE</i> . . . . .	104
6.6.2	Experiments on <i>CAPT</i> . . . . .	106
6.7	Properties of Synthetic Data Sets . . . . .	111
6.7.1	Properties of Credit Risk (CR) Data Set . . . . .	111
6.7.2	Properties of Customer Status (CS) Data Set . . . . .	112
6.8	Conclusion . . . . .	114
<b>7</b>	<b>The Framework and Experimental Results</b>	<b>116</b>
7.1	The Framework . . . . .	116
7.2	The Extended Framework . . . . .	117
7.3	Experiments . . . . .	119
7.3.1	Experiments on the Adult Data Set . . . . .	122
7.3.2	Experiments on Wisconsin Breast Cancer Data Set . . . . .	137
7.4	Conclusion . . . . .	148
<b>8</b>	<b>Measuring of Disclosure Risk</b>	<b>155</b>
8.1	<b>Measuring Disclosure Risk</b> . . . . .	156
8.2	Conclusion . . . . .	164

<b>9</b>	<b>Data Quality</b>	<b>165</b>
9.1	Motivation . . . . .	165
9.2	Our Work . . . . .	167
9.3	Experimental Results . . . . .	169
9.4	Conclusion . . . . .	173
<b>10</b>	<b>Conclusion</b>	<b>175</b>
	<b>Bibliography</b>	<b>179</b>

# List of Figures

2.1	An example of a decision tree. Squares represent internal nodes, the unshaded circle represents homogeneous leaf where all records have the same class value and shaded circles represent heterogeneous leaves. . . . .	11
2.2	Main Clustering Methods. . . . .	13
3.1	Classification of Data Sets Based on Distribution. . . . .	22
3.2	A Classification of Privacy Preserving Techniques. . . . .	25
4.1	An example of a decision tree classifier. . . . .	51
4.2	The decision tree obtained from 300 records of the original <i>BHP</i> data set. .	59
4.3	The decision tree obtained from the 1st of the five <i>BHP</i> data sets that have been perturbed by the <i>RPT</i> . . . . .	60
4.4	The decision tree obtained from the 2nd of the five <i>BHP</i> data sets that have been perturbed by the <i>RPT</i> . . . . .	61
4.5	The decision tree obtained from the 3rd of the five <i>BHP</i> data sets that have been perturbed by the <i>RPT</i> . . . . .	62
4.6	The decision tree obtained from the 4th of the five <i>BHP</i> data sets that have been perturbed by the <i>RPT</i> . . . . .	64
4.7	The decision tree obtained from the 5th of the five <i>BHP</i> data sets that have been perturbed by the <i>RPT</i> . . . . .	65
4.8	The decision tree obtained from one of the ten <i>BHP</i> data sets that have been perturbed by the <i>PPT</i> . . . . .	66
4.9	The decision tree obtained from another <i>BHP</i> data set that has been perturbed by the <i>PPT</i> . . . . .	67
4.10	The decision tree obtained from a 3rd <i>BHP</i> data set that has been perturbed by the <i>PPT</i> . . . . .	68
4.11	The decision tree obtained from one of the ten <i>BHP</i> data sets that have been perturbed by the <i>ALPT</i> . . . . .	68
4.12	The decision tree obtained from another <i>BHP</i> data set that has been perturbed by the <i>ALPT</i> . . . . .	69
4.13	The decision tree obtained from a 3rd <i>BHP</i> data set that has been perturbed by the <i>ALPT</i> . . . . .	70

5.1	The decision tree obtained from 349 records of the original (unperturbed) <i>WBC</i> data set. . . . .	73
6.1	The basic concept of similarity, of two values belonging to a categorical attribute, in <i>CACTUS</i> . . . . .	83
6.2	An illustration of the correlation analysis by <i>CORE</i> . . . . .	84
6.3	An example showing a limitation of <i>CORE</i> . . . . .	86
6.4	Representation of a data set as a hyper-graph. . . . .	90
6.5	A section of the decision tree built on the <i>CR</i> data set. The tree considers attribute <i>City</i> as class attribute. . . . .	95
6.6	Basic steps of <i>EX-DETECTIVE</i> - for clustering records based on attributes A and B. . . . .	98
6.7	Clustering records based on the attributes A, B and C. . . . .	99
6.8	Clustering records of a data set having numerical attribute/s along with categorical attribute/s. . . . .	100
6.9	Details of a decision tree built from the unperturbed <i>CS</i> data set. The tree considers attribute <i>Car Make</i> as class attribute. This tree is used for clustering values of the attribute <i>Car Make</i> . . . . .	105
6.10	A decision tree $T_o(status)$ , built on the original <i>CS</i> data set. The tree considers the natural class attribute <i>Status</i> as class attribute. . . . .	107
6.11	A decision tree $T_p(status)$ , built on a perturbed <i>CS</i> data set. The tree considers the attribute <i>Status</i> as class attribute. . . . .	108
6.12	A decision tree built on a total perturbed <i>CS</i> data set. The tree considers attribute <i>Car Make</i> as class attribute. . . . .	109
6.13	A decision tree built on another total perturbed <i>CS</i> data set. The tree considers attribute <i>Car Make</i> as class attribute. . . . .	109
7.1	The decision tree $DT_{training}$ obtained from 25,600 records of the training <i>Adult</i> data set. . . . .	124
7.2	The decision tree obtained from a data set perturbed by the <i>Framework</i> . . .	125
7.3	The decision tree obtained from a data set perturbed by the <i>Framework</i> . . .	125
7.4	The decision tree obtained from a data set perturbed by the Random Framework. . . . .	128
7.5	The decision tree obtained from a data set perturbed by the Random Framework. . . . .	129
7.6	The decision tree obtained from a data set perturbed by the <i>Extended Framework</i> . . . . .	130
7.7	The decision tree obtained from a data set perturbed by the <i>Extended Framework</i> . . . . .	131
7.8	The decision tree obtained from a data set perturbed by Random Extended Framework. . . . .	135
7.9	The decision tree obtained from a data set perturbed by Random Extended Framework. . . . .	136
7.10	The decision tree obtained from the training <i>WBC</i> data set. . . . .	139
7.11	The decision tree obtained from a <i>WBC</i> data set perturbed by the <i>Framework</i> . . .	140

7.12	The decision tree obtained from a <i>WBC</i> data set perturbed by the <i>Framework</i> .	141
7.13	The decision tree obtained from a data set perturbed by Random Technique.	142
7.14	The decision tree obtained from another data set perturbed by Random Technique. . . . .	143
7.15	The decision tree obtained from a data set perturbed by the <i>Extended Framework</i> . . . . .	144
7.16	The decision tree obtained from a data set perturbed by the <i>Extended Framework</i> . . . . .	145
7.17	An example of how the information gain of an attribute can increase due to a noise addition. . . . .	151
8.1	The probability distribution of a perturbed record originating from the target record $x$ . . . . .	160
8.2	Entropies of a perturbed data set calculated for each original record. . . . .	161
9.1	A decision tree obtained from the training <i>BHP</i> data set having 300 records. Squares represent internal nodes, unshaded circle represents homogeneous leaf and shaded circles represent heterogeneous leaves. . . . .	168

# List of Tables

2.1	Harris Poll Survey: Privacy Consciousness of Adults [98] . . . . .	19
3.1	Privacy Preserving Techniques - Comparative Study . . . . .	47
6.1	Time taken by the whole program in seconds on the mushroom data set. . .	88
6.2	Time taken by the whole program in seconds on the landslide data set. . .	89
6.3	The cluster produced by <i>CACTUS</i> from the <i>CS</i> data set. . . . .	105
6.4	Similarities of perturbed trees with corresponding original trees. . . . .	110
6.5	Similarities of perturbed trees with corresponding original trees - using J48.	111
7.1	Prediction Accuracy of the Classifiers Obtained from the Unperturbed and Various Perturbed Adult Data Sets. . . . .	152
7.2	Prediction Accuracy of the Classifiers Obtained from the Unperturbed and Various Perturbed WBC Data Sets. . . . .	153
7.3	Prediction Accuracy of the Classifiers Obtained from WBC Data Sets Perturbed by GADP technique only. . . . .	154
8.1	A Compare of Entropies for the Cases Where the Intruder Has Access to the Original and the Perturbed Data Set. . . . .	163
9.1	Experimental Results of Decision Tree Classifiers on the BHP Data Set. . .	170
9.2	Experimental Results of Decision Tree Classifiers on the WBC Data Set. . .	171
9.3	Experimental Results of Neural Network Classifier on BHP Data Set. "Diff." in Col. G Means That There are 5 Different Values in All 5 Experiments and Hence There is No Single Mode Value. . . . .	172
9.4	Experimental Results of Neural Network Classifier on WBC Data Set. . . .	172