

Faculty of Engineering and Information Technology  
University of Technology Sydney

Data Mining of Classification  
for  
Sybil User Detection

Thesis submitted in partial fulfilment of  
the requirements for the degree of  
**Master of Analytics by Research**

By

Anand Arun Chinchore

June 2016

*Dedicated*

*To my late father....*

*Mr Arun Yashwant Chinchore*



## CERTIFICATE OF AUTHORSHIP / ORIGINALITY

I certify that the work in this thesis has not previously been submitted for a degree nor has it been submitted as a part of requirements for a degree except as fully acknowledged within the text.

I also certify that the thesis has been written by me. Any help that I have received in my research work and the preparation of the thesis itself has been acknowledged. In addition, I certify that all information sources and literature used are indicated in the thesis.

Signature of Candidate

-----



## Acknowledgements

I would like to express my special appreciation and thanks to my supervisor Dr Guandong Xu. You have shown tremendous patience with me as one of my mentors. I would like to thank you for encouraging my research and for allowing me to grow as a researcher in analytics. Your advice on both my research and my career have been priceless.

I would also like to extend my appreciation to my co-supervisor Dr Frank Jiang for providing me with continuous support throughout my Master of Analytics by research and study. Without your professional guidance, persistent help, patience, motivation and immense knowledge, this thesis would not have been possible. Words are not enough.

I thank the management team members of the Advanced Analytics Institute, Mr Colin Wise and Mr Austin Andrade, for their encouragement and guidance towards the completion of my research journey. All of you have been there to support me, when required, over the last two and half years.

I would also like to thank the many other faculties, staff, lecturers, professors and doctors of the University of Technology Sydney, who guided and helped me along the path of my research within UTS.

I thank my workplace manager, and my colleague, Mrs Roula Christodoulides who helped and supported me to balance my time in work and study.

A special thanks to my mother Smt Vasudha Arun Chinchore and my family. Words cannot express how grateful I am to you for all of the sacrifices that you've made on my behalf. Your prayers for me have sustained me thus far. I would also like to thank all of my friends who have supported me in writing and given me incentive to strive towards my goal.

Anand Arun Chinchore

June 2016 @ UTS



# Table of contents

<b>Certificate .....</b>	<b>4</b>
<b>Acknowledgements.....</b>	<b>6</b>
<b>List of figures .....</b>	<b>12</b>
<b>List of tables .....</b>	<b>14</b>
<b>List of publications .....</b>	<b>16</b>
<b>Abstract.....</b>	<b>18</b>
<b>Chapter 1 Introduction</b>	
1.1 Background .....	22
1.1.1 Mobile social networking .....	23
1.1.2 Social media and data mining .....	24
1.1.3 The threats in social networking sites .....	24
1.2 History of the term sybil .....	25
1.2.1 Sybil users in computer terminology .....	25
1.2.2 Honest users and sybil users .....	25
1.2.3 Sybil user communities and group forming .....	26
1.2.4 Sybil user, sybil behaviours and sybil attacks .....	26
1.3 Classification and regression for data mining.....	27
1.3.1 C4.5 algorithm.....	30
1.3.2 Entropy.....	30
1.3.3 Information gain .....	31
1.3.4 Random forest .....	32
1.4 Research aims .....	32
1.5 Research contributions .....	34
1.6 Thesis structure.....	35



## **Chapter 2 Literature review and foundation**

2.1 Directed social networks .....	38
2.2 Defending against sybil attacks in mobile social network .....	39
2.3 Sybil attack in MSN .....	40
2.4 Sybil defender .....	42
2.5 Sybil limit .....	45
2.6 Sybil guard.....	47
2.7 Decentralized defense, sybil attacks, friends and foes and social ties .....	48
2.8 Random-walks: the evolution of sybil defense protocols .....	51
2.9 Decision-making processes .....	53
2.10 Social graph mining .....	54
2.11 Security measures .....	56
2.12 Sybil attacks and defenses in the Internet of Things (IoT).....	56
2.13 Exploiting mobile social behaviours for sybil detection .....	59
2.14 Other researches view points .....	62
2.15 Summary.....	63

## **Chapter 3 Graph-Based behaviour analysis using connectivity between Nodes**

3.1 Background .....	66
3.1.1 Scenarios .....	68
3.1.2 The research process of identification.....	69
3.1.3 Graph theory.....	70
3.2 Research methodology .....	70
3.2.1 Objectives .....	70
3.2.2 Processing methods.....	71
3.3 Conclusion.....	78

## **Chapter 4 A classified model for sybil detection**

4.1 Introduction .....	80
4.1.1 The C4.5 algorithm.....	83
4.2 Decision trees, entropy and gain .....	85
4.2.1 Entropy.....	89

4.2.2 C4.5 processing .....	91
4.2.2.1 The issues.....	92
4.2.2.2 Overcoming the issues.....	92
4.2.2.3 The splitting process .....	93
4.2.3 Binary count of columns .....	93
4.2.3.1 Total count for each node and connection.....	93
4.2.3.2 Time difference count for each node and connection .....	94
4.2.3.3 IndexValOne count for each node and connection .....	94
4.2.3.4 ConnIDZero count for each node and connection.....	94
4.2.4 Entropy analysis and generation .....	95
4.2.4.1 Analysis of user vs TimeDiff .....	95
4.2.4.2 First row mismatch issue in dataset .....	98
4.2.4.3 Last row mismatch issue in dataset.....	99
4.2.4.4 Analysis of TimeDiff vs IndexValOne .....	99
4.2.4.5 Analysis of IndexValOne vs ConnID .....	100
4.2.5 Information gain Analysis and generation.....	101
4.2.6 maxGain analysis.....	102
4.3 Random forests and the Gini index.....	103
4.3.1 Variable importance .....	106
4.3.2 Confusion matrix.....	106
4.3.3 Decision tree pruning.....	107
4.4 Conclusion.....	109
<b>Chapter 5 Conclusions and future work</b>	
5.1 Conclusion.....	111
5.1 Future work.....	113
<b>Appendix - List of symbols.....</b>	<b>131</b>
<b>Bibliography .....</b>	<b>133</b>



# List of figures

1.1 Decision tree map node distribution Infocom06 dataset .....	29
1.2 The work flow in this thesis .....	36
2.1 System model of research as explained in Chang W. et al's 2013 research .....	41
2.2 Generation of a social distrust profile, as explained in Chang W. et al's 2013 research .....	41
2.3 Research system model (Zhang K. et. al., 2015) .....	59
2.4 Observation on contact information and pseudonyms changing between normal users and sybil attackers (Zhang K. et. al., 2015).....	61
2.5 Contact rate distribution charts (Zhang K. et. al., 2015 .....	61
3.1 Count of connections between User 1 and all other users in the Infocom06 dataset. ....	72
3.2 User 2 vs User 1 suspicious identity with high and low connections .....	73
3.3 Bar graph: User 1 vs high short-time connectivity between various connections ....	75
3.4 Highest connecting User 1 vs User 2.....	76
3.5 Boxplot User 1 vs User 2 shows the user gap between one to other.....	76
3.6 Bar Diagram for User 2 vs High short-time connectivity .....	78
4.1 Split of training the dataset .....	82
4.2 Decision tree mapping using C4.5 algorithm .....	84
4.3 Model diagram .....	85
4.4 Generalised decision tree map for Infocom06 .....	87
4.5 Decision tree - node1 generated in WEKA .....	91
4.6 Example of one-attribute entropy calculation .....	95-96
4.7 sample of two or more attribute entropy calculation .....	98
4.8 Error rate across decision trees .....	105
4.9 Decision trees pruning .....	109
4.10 User 1 decision tree .....	122
4.11 User 4 decision tree .....	123

4.12 User 8 decision tree .....	124
4.13 User 12 decision tree .....	125
4.14 User 16 decision tree .....	126
4.15 User 18 decision tree .....	127
4.16 User 19 decision tree .....	128
4.17 User 42 decision tree .....	129
4.18 User 66 decision tree .....	130

## List of Tables

3.1 User 1 vs short time count .....	75
3.2 Attribute ID vs Quantity (highest connection).....	75
3.3 Analyse - Descriptive Statistics – Frequencies (High Frequencies repetition and engagements of User1 with other nodes) .....	77
3.4 Analyse - Descriptive Statistics- Frequencies (High Frequencies repetition and engagements of beginning user of column 2User2 on by their position in between nodes changes dramatically) .....	77
4.1 Gini index calculation.....	106
4.2 Confusion matrix.....	106
4.3 Confusion matrix statistics.....	106
4.4 Confusion matrix and statistics for response variable.....	107
4.5 Total count – frequency count for each pair. ....	115
4.6 TimeDiff count – frequency count for each pair .....	115
4.7 IndexOneVal count – frequency count for each pair .....	116
4.8 ConnIDZero count – frequency count for each pair .....	116
4.9 Entropy calculations – user - TimeDiff datasets for each pair .....	117
4.10 Frequency mismatch – with total count for each pair .....	117
4.11 Frequency mismatch – total count with TimeDiff for each pair .....	118
4.12 TimeDiff vs IndexValOne entropy for each pair.....	118
4.13 IndexValOne vs. ConnIDZero entropy for each pair .....	119
4.14 IndexValOne vs ConnIDZero few entropy for User 55.....	119
4.15 Information gain calculation .....	120
4.16 maxGain calculation .....	120
4.17 Gini index calculation .....	121



## List of publications

### Papers published

Chinchore, A, Jiang, F, Xu, G 2015, 'Intelligent Sybil attack detection on abnormal connectivity behaviour in mobile social networks', in the proceedings of the Springer International Publishing 10th International Conference, KMO (2015), Maribor, Slovenia, 24-28 August, vol. 224, pp. 602-617.





## Abstract

Data analytics and Big Data application research, along with new structures in complex data, are revealing the secrets of their own complexity and patterns with valuable and critical, but challenging, issues through newly designed tools, techniques and models in data science technology. A common example concerns the interconnectivity of social network users on mobiles, involving content and information sharing through mobile social networks.

There have been a large number of studies on mobile networks. Many focus on a variety of secured applications that attempt to exploit social connections, impersonate users or attack social groups. Such applications are often created with the intention of collecting confidential information, laundering money, blackmail or to perform other criminal activities.

Existing methods for identifying such activity, such as distributed systems, social graph-based sybil detection, behaviour classification, and local ranking systems that estimate the trust level between users, rely on the dependencies between random nodes of connection on mobile social networks. These models aim to detect suspicious connections and have the advantage of learning the relationships between nodes and data. However, their detection patterns tend to impose the behavioural patterns typically associated with community-based and external networks.

In data mining, the graph-based and classification models used for pattern collection can accurately predict patterns in data in targeted categories. Decision trees, commonly used for classification, are trees in which each branch represents a choice between a number of alternatives, and each leaf represents a classification, or decision. For example, a decision tree may help an institution decide whether a node in a dataset is suspicious, or considered to be sybil, if a

decision tree can be induced from a set of data about its instances and the - classifications of those instances. It could also provide the flexibility to demonstrate data distribution. Thus, researchers have tried to combine different techniques and methods into network-based models to detect various patterns generated by sybil nodes within a network.

The purpose of this thesis is to abridge existing classification and regression techniques to identify sybil nodes, and the correlation of those nodes with time, to address these research limitations.

Classification and regression techniques predict behaviour based on continuous or categorical responses. For example if the predicted response is continuous, then it is called a regression tree. If the response is categorical, it is called a classification tree. At each node of the tree, the value of one the connected input nodes is checked and a binary answer – yes or no – determines whether one continues to the left or right sub-branch. When a leaf is reached, a prediction follows from a series of entropy calculations and graphing techniques.

This thesis introduces a novel classification model for sybil detection in mobile social behaviour that identifies dependencies using connection duration and other attributes. Roger Quinlan's C4.5 algorithm, its resulting decision tree and a random forest simplify the step-by-step identification process, while maintaining its merits. Partial correlations between nodes are simplified using Rattle programming, and the dataset is divided into majority nodes to assist processing.

This research also includes a behavioural survey of the nodes and an extended analysis using a classification system for sybil detection, with a particular focus on sybil attacks in mobile social network environments. Each sybil node is tracked and identified based on the frequency and duration of its connections with other nodes.

An outline of how the classified model identifies behaviour is also included, along with an explanation of the flow of the decision tree and the C4.5 algorithm process, which press-gangs identified sybil nodes based on the results of entropy calculations and information gain. The calculated entropy for each node connection across the all datasets informs the information gain. The

maxGain calculations for individual node bring the final stage of draw decision tree and helped to predict the sybil nodes, compare and justify the sybil attackers .

These processes and new models applied to sybil detection provide insight into the behaviour of connections, through deep analytics and entropy gain. The evidence gleaned from this research brings significant knowledge to data analytics and data science in the identification of threats on mobile social networks.

