

# **Big Data Behaviour Modelling and Visual Analytics**

**Jinson Zhang**

A thesis submitted for the degree of  
Doctor of Philosophy in Computing Sciences

**Supervisor by: A/Professor Mao Lin Huang**

May 2017  
Faculty of Engineering & Information Technology  
University of Technology Sydney  
Australia

## Certificate of Authorship/Originality

---

Data:           **October 2016**  
Author:       **Jinson Zhang**  
Title:         **Big Data Behaviour Modelling and Visual Analytics**  
Degree:       **Doctor of Philosophy in Computer Sciences**

---

I certify that the work in this thesis has not previously been submitted for a degree nor has it been submitted as 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 STUDENT

---

## Acknowledgements

Firstly, I would like to thank my supervisor A/Professor Mao Lin Huang for giving me the opportunity to do my PhD within the Faculty of Engineering & Information Technology at the University of Technology Sydney, and for his enormous encouragement, advice, and academic guidance during my part-time research period. Dr. Mao Lin have given me not only academic advice and freedom in my research area, but had also provided me with teaching opportunities over the last 15 years across a broad range of different subjects, including in Internet Programming and Data Visualization.

Secondly, I would like to thank our Visual Analytics Team, consisting of Dr. Jie Liang, Dr. Jie Hua, Dr. Liang Fu Lu, and Dr. Wen Bo Wang, for their knowledge sharing and information exchange. Their focus and patience have encouraged me to finish my work, and I have greatly increased my skills and visualization techniques during our fortnightly team meetings. Their humor and trust have benefited me not only through academic partnerships, but also through friendships outside of research.

Thirdly, I would like to thank the Library for their support during my research period. As a full-time Library staff and part-time researcher, I have been given the ability to work flexibly in order in allowing me to attend conferences and other research/teaching activities. Without their support, I would not have been able to finish my research. The Library has also provided me with enormous online materials, and has helped me with referencing works.

Furthermore, I would like to thank the School of Software, Faculty of Engineering & Information Technology, and the Graduate Research School for their help in providing the Faculty Travel Fund and Vice-Chancellor's Conference Fund to attend national and

international conferences. This has been enormously beneficial for part-time PhD researchers such as myself.

Finally, I would like to take this opportunity to thank my wife Lisa and son Edmund for their enormous help and support. Having started my PhD as part-time researcher, their understanding and belief in my capability has provided enormous encouragement during my research period. They have provided continuous support and help in making this PhD thesis possible, and I am thankful for having my family behind me every step of the way.

# Table Contents

<b>Certificate of Authorship/Originality.....</b>	<b>ii</b>
<b>Acknowledgements.....</b>	<b>iii</b>
<b>Table Contents.....</b>	<b>v</b>
<b>List of Figures.....</b>	<b>viii</b>
<b>List of Tables .....</b>	<b>x</b>
<b>List of Equations .....</b>	<b>xi</b>
<b>Abstract.....</b>	<b>1</b>
<b>Author’s Seventeen published papers for PhD .....</b>	<b>3</b>
<b>Chapter 1: Introduction .....</b>	<b>6</b>
1.1.    Big Data Characteristics.....	6
1.1.1.    Volume.....	7
1.1.2.    Variety.....	8
1.1.3.    Velocity.....	8
1.1.4.    Value and Veracity.....	9
1.2.    Big Data Visual Analytics.....	10
1.2.1.    Dataset Visualization .....	11
1.2.2.    Data-Type Visualization .....	11
1.2.3.    Particular Topic Visualization .....	12
1.3.    Big Data Behaviours .....	12
1.4.    Visualization Techniques.....	14
1.4.1.    Visualization Techniques for Low-Dimensional Data.....	14
1.4.1.1.    Bubble Charts/Scatter Plots .....	14
1.4.1.2.    Word Clouds .....	16
1.4.1.3.    Heat Maps .....	17
1.4.1.4.    Geography Maps.....	19
1.4.1.5.    Data Histograms.....	20
1.4.2.    Visualization Techniques for High-Dimensional Data .....	22
1.4.2.1.    Tree Maps .....	22
1.4.2.2.    Data Clustering .....	24

1.4.2.3.	Parallel Coordinates .....	26
1.4.2.3.1	Visual clustering .....	27
1.4.2.3.2	Axes rotation and reordering.....	28
1.4.2.3.3	Reducing dimensions and feature selection .....	30
1.5.	Research Challenges and Motivations .....	31
1.6.	Author's Contributions in the Thesis .....	34
1.7.	Thesis Organisation .....	35
<b>Chapter 2: Data Behaviour Visual Analytics .....</b>		<b>37</b>
2.1.	Multidimensional Data and Attributes.....	37
2.2.	5Ws Dimension and Behaviours Pattern.....	39
2.3.	5Ws Parallel Coordinates.....	41
2.4.	5Ws Dimension Clustering .....	42
2.5.	5Ws Shrunk Attributes.....	44
2.6.	Noise Attributes .....	45
<b>Chapter 3: Pair-Density Parallel Coordinates .....</b>		<b>48</b>
3.1.	Pair-Density Algorithm.....	49
3.2.	Pattern and Pair-Density .....	51
3.3.	$SD_{(p, pat(x, y, z))}$ via $RD_{(q, pat(x, y, z))}$ for $pat(x, y, z)$ .....	53
3.4.	$SD_{(p, pat(x, y, q))}$ via $PD_{(z, pat(x, y, q))}$ with $pat(x, y, q)$ .....	56
3.5.	$SD_{(p, pat(x, z, q))}$ via $TD_{(y, pat(x, z, q))}$ with $pat(x, z, q)$ .....	59
3.6.	$SD_{(p, pat(y, z, q))}$ via $CD_{(x, pat(y, z, q))}$ with $pat(y, z, q)$ .....	61
3.7.	$CD_{(x, pat(p, y, z))}$ via $RD_{(q, pat(p, y, z))}$ with $pat(p, y, z)$ .....	63
3.8.	$CD_{(x, pat(p, y, q))}$ via $PD_{(z, pat(p, y, q))}$ with $pat(p, y, q)$ .....	65
3.9.	$CD_{(x, pat(p, z, q))}$ via $TD_{(y, pat(p, z, q))}$ with $pat(p, z, q)$ .....	67
3.10.	$TD_{(y, pat(p, x, z))}$ via $RD_{(q, pat(p, x, z))}$ with $pat(p, x, z)$ .....	70
3.11.	$TD_{(y, pat(p, x, q))}$ via $PD_{(z, pat(p, x, q))}$ with $pat(p, x, q)$ .....	72
3.12.	$PD_{(z, pat(p, x, y))}$ via $RD_{(q, pat(p, x, y))}$ with $pat(p, x, y)$ .....	74
<b>Chapter 4: Case Study .....</b>		<b>77</b>
4.1.	Case One: Visual Estimate of US 2008 Flight Delay Patterns .....	77
4.1.1.	5Ws Pattern for US Flight Dataset.....	79
4.1.2.	The Delay Pattern Combined Airline and Flight Distance .....	81
4.1.3.	Clustered Delay Pattern in $SD_{( )}$ via $RD_{( )}$ Parallel Coordinates.....	82
4.1.4.	Airline Flight Pattern between Origin and Destination Airport.....	83

4.2.	Case Two: Visual Analysis of 2009 Spam Email .....	85
4.2.1.	5Ws Pattern for Virus Email .....	87
4.2.2.	Virus Pattern $pat_{(x, y, z)}$ in $SD_{()}$ via $RD_{()}$ Parallel Coordinates .....	88
4.2.3.	Clustered Transferring Pattern $pat_{(p, y, q)}$ .....	91
4.3.	Case Three: Visual Detect DDoS Attacks in ISCX2012 Dataset .....	92
4.3.1.	5Ws Pattern for ISCX2012 Dataset .....	94
4.3.2.	Network Pattern $pat_{(x, y, z)}$ in $SD_{()}$ via $RD_{()}$ Parallel Coordinates .....	95
4.3.3.	Attack Pattern $pat_{(x="attack", y, z)}$ between Attacker and Victim .....	97
4.4.	Reduction of Data Overcrowding .....	99
<b>Chapter 5: Conclusions and Future Work .....</b>		<b>102</b>
5.1.	Conclusions.....	102
5.2.	Future Works .....	105
<b>Bibliography .....</b>		<b>106</b>

# List of Figures

Figure 1.1	3Vs model of Big Data.....	7
Figure 1.2	5Vs model of Big Data.....	9
Figure 1.3	Behaviour analysis by string matching (Alam et al 2016).....	13
Figure 1.4	Hans Rosling's bubble chart .....	15
Figure 1.5	Example of word cloud.....	17
Figure 1.6	Heat map of Bank World's activity .....	18
Figure 1.7	Uncertainty ribbon in geography map.....	20
Figure 1.8	HOG histogram bin with different GPR downtime .....	21
Figure 1.9	Large file system with the different shapes in treemaps .....	23
Figure 1.10	Matrix structure and transformation used in treemaps.....	24
Figure 1.11	Group with high attribute value in red rectangle .....	25
Figure 1.12	Blog network clustering by ADraw .....	26
Figure 1.13	Curve edges to reduce the visual clutter .....	28
Figure 1.14	Axes re-ordering to reduce the visual clutter .....	29
Figure 1.15	DNA Microarray data with Scattering Points in Parallel Coordinates.....	31
Figure 2.1	Big Data 5Ws pattern.....	40
Figure 2.2	Example of 5Ws parallel coordinates .....	41
Figure 2.3	Example of Big Data in 5Ws pattern crossing multiple datasets .....	42
Figure 2.4	Tree structure of 5Ws pattern for Big Data.....	43
Figure 2.5	Example of clustered 5Ws parallel coordinates .....	44
Figure 2.6	Example of SA in 5Ws parallel coordinates .....	45
Figure 2.7	Example of noise data in 5Ws parallel coordinates .....	46
Figure 3.1	Example of relationship for $SD_{(p, pat(x, y, z))}$ via $RD_{(q, pat(x, y, z))}$ .....	53
Figure 3.2	Example of 5Ws parallel coordinates with $SD_{(p, pat(x, y, z))}$ via $RD_{(q, pat(x, y, z))}$ .....	55
Figure 3.3	Example of 5Ws parallel coordinates with $SD_{(p, pat(x, y, q))}$ via $PD_{(z, pat(x, y, q))}$ .....	57
Figure 3.4	Example of 5Ws parallel coordinates with $SD_{(p, pat(x, z, q))}$ via $TD_{(y, pat(x, z, q))}$ .....	60
Figure 3.5	Example of 5Ws parallel coordinates with $SD_{(p, pat(y, z, q))}$ via $CD_{(x, pat(y, z, q))}$ .....	62
Figure 3.6	Example of 5Ws parallel coordinates with $CD_{(x, pat(p, y, z))}$ via $RD_{(q, pat(p, y, z))}$ .....	64
Figure 3.7	Example of 5Ws parallel coordinates with $CD_{(x, pat(p, y, q))}$ via $PD_{(z, pat(p, y, q))}$ .....	66
Figure 3.8	Example of 5Ws parallel coordinates with $CD_{(x, pat(p, z, q))}$ via $TD_{(y, pat(p, z, q))}$ .....	68
Figure 3.9	Example of 5Ws parallel coordinates with $TD_{(y, pat(p, x, z))}$ via $RD_{(q, pat(p, x, z))}$ .....	71
Figure 3.10	Example of 5Ws parallel coordinates with $TD_{(y, pat(p, x, q))}$ via $PD_{(z, pat(p, x, q))}$ .....	73
Figure 3.11	Example of 5Ws parallel coordinates with $PD_{(z, pat(p, x, y))}$ via $RD_{(q, pat(p, x, y))}$ .....	75
Figure 4.1	US 2008 flight patterns in 5Ws parallel coordinates.....	80
Figure 4.2	Delay pattern $pat_{(x>600, y, z)}$ in $SD_{()}$ via $RD_{()}$ parallel coordinates .....	81
Figure 4.3	Clustered delay pattern $pat_{(x>600, y, z)}$ in $SD_{()}$ via $RD_{()}$ parallel coordinates .....	83
Figure 4.4	Airline flight pattern $pat_{(p, y, q)}$ in $CD_{()}$ via $PD_{()}$ parallel coordinates .....	84
Figure 4.5	Example of an email incident.....	85
Figure 4.6	2009 email virus pattern in 5Ws parallel coordinates .....	88
Figure 4.7	Virus pattern $pat_{(x=virus, y, z)}$ in $SD_{()}$ via $RD_{()}$ parallel coordinates with SA.....	90
Figure 4.8	Virus transferring pattern $pat_{(p, y, q)}$ in $CD_{()}$ via $PD_{()}$ parallel coordinates.....	91



Figure 4.9	ISCX2012 network pattern in 5Ws parallel coordinates.....	95
Figure 4.10	Network pattern $\text{pat}_{(x, y, z)}$ in $\text{SD}_{()}$ via $\text{RD}_{()}$ parallel coordinates.....	96
Figure 4.11	Attack pattern $\text{pat}_{(x=\text{"attack"}, y, z)}$ between attackers and victims .....	97
Figure 4.12	DDoS attack pattern in $\text{SD}_{()}$ via $\text{RD}_{()}$ parallel coordinates .....	98
Figure 4.13	Reduction of data cluttering (a) .....	100
Figure 4.14	Reduction of data cluttering (b) .....	100

## List of Tables

Table 1.1.	Dataset Volume Size .....	7
Table 3.1.	5Ws Pair-Density and Patterns .....	52
Table 4.1.	5Ws classification for US 2008 flight dataset.....	78
Table 4.2.	5Ws dimension for US 2008 flight dataset .....	79
Table 4.3.	5Ws classification for 2009 email dataset.....	86
Table 4.4.	5Ws dimension for 2009 virus email .....	87
Table 4.5.	5Ws dimension for ISCX2012 network dataset.....	93
Table 4.6.	5Ws pattern for ISCX2012 network traffic.....	94
Table 4.7.	Data pattern for three cases.....	99

# List of Equations

Equation 2. 1 .....	37
Equation 2. 2 .....	37
Equation 2. 3 .....	38
Equation 2. 4 .....	39
Equation 2. 5 .....	40
Equation 2. 6 .....	40
Equation 2. 7 .....	43
Equation 2. 8 .....	45
Equation 2. 9 .....	46
Equation 3. 1 .....	49
Equation 3. 2 .....	49
Equation 3. 3 .....	50
Equation 3. 4 .....	51
Equation 3. 5 .....	54
Equation 3. 6 .....	56
Equation 3. 7 .....	59
Equation 3. 8 .....	61
Equation 3. 9 .....	63
Equation 3. 1 0 .....	65
Equation 3. 1 1 .....	67
Equation 3. 1 2 .....	70
Equation 3. 1 3 .....	72
Equation 3. 1 4 .....	74

## Abstract

Big Data is composed of text, images, video, audio, mobile or other forms of data collected from multiple datasets, and is rapidly growing in both size and complexity. This has created a huge volume of multidimensional data within a very short time period. Big Data is therefore too big, too complex and moves too fast for us to analyze using traditional methods. Big Data behaviour is considered as a set of concepts and categories that describes Big Data's acts towards others. The challenges facing Big Data analysis and visualization include: 1) how to classify Big Data across multiple datasets and different forms of data, 2) how to visualize structured and unstructured Big Data behaviour patterns for multidimensional data, 3) how to display Big Data behaviour patterns with very large volumes onto a normal-sized screen, 4) how to visualize Big Data behaviour patterns without the loss of information.

Big Data visualization normally requires optimized solutions through using different visual techniques for integrating display and exploration. To illustrate the huge amount of multidimensional data within a standard-size screen, visualization needs to find an efficient classification method for multiple datasets across any form of data. The current data interactive exploration has normally optimized data for visualization by excluding some pieces of information, resulting in missing information. Big Data visualization also suffers from visual cluttering and data overcrowding problems, whilst dealing with huge amounts of multidimensional data.

My approach includes two parts: Big Data behaviour modelling and Big Data visualization. I have firstly established the 5Ws dimensions for Big Data classification, based on data behaviour ontologies, that can be applied to multiple datasets and to any

form of data. Each data incident contains these 5Ws dimensions, which are posed as a set of concepts and categories that describes Big Data acts for; **When** did the data occur, **Where** did the data come from, **What** did the data contain, **How** was the data transferred, **Why** did the data occur, and **Who** received the data. Secondly, I have introduced Pair-Density algorithms to measure Big Data behaviour patterns, which enables comparison and analysis between any two dimensions of behaviours. Two non-dimensional axes in parallel coordinates have then been created by using Pair-Density to measure and compare visual patterns for Big Data visualization. Finally, Shrunk Attributes has been deployed into Pair-Density parallel coordinates. This not only narrows down Big Data patterns for better understanding, but also dramatically reduces data cluttering and overcrowding in Big Data visualization.

Three different datasets with a combined total of more than 2.5 million data incidents have been implemented for measuring and visualizing different data patterns, including both numerical and non-numerical dimensions. The experimental results have shown that my new approach has significantly improved the accuracy of Big Data visualization, reduced data cluttering by more than 80% without the loss of information. The use of 5Ws dimensions and Pair-Density parallel coordinates therefore has large potential benefits and applications across both the business and research fields.

This thesis contains the research approach and implementation results obtained by the author during his Ph.D period. The majority of methods and results have been published in **Seventeen** research papers in journals and conference proceeding by May 2016.

## Author's Seventeen published papers for PhD

- 1). J. Zhang, ML. Huang, (2016), "2D Approach Measuring Multidimensional Data Pattern in Big Data Visualization", IEEE first International Conference on Big Data Analysis (ICBDA2016), In Proceeding of 2016 IEEE International Conference on Big Data Analysis, IEEE Computer Society, pp. 194-199, March 2016, DOI: 1109/ICBDA.2016.7509823
- 2). LF Lu, ML. Huang, and J. Zhang, (2016), "Two Re-ordering Methods in Parallel Coordinates Plots", Journal of Visual Languages and Computing, Elsevier, Vol. 33, pp. 3-12, April 2016, DOI: 10.1016/j.jvlc.2015.12.001
- 3). J. Zhang and M.L. Huang, (2016), "Data Behaviours Model for Big Data Visual Analytics", International Journal of Big Data Intelligence, InderScience Publishers, Vol. 3, No. 1, pp. 1-17, Dec 2015, DOI: 10.1504/IJBDI.2016.073899
- 4). J. Zhang, M.L. Huang, and Z. Meng, (2015), "Visual Analytics for BigData Variety and Its Behaviours", International Journal of Computer Science and Information Systems, ComSIS Consortium, Vol. 12, No. 4, pp. 1171-1191, Nov 2015, DOI: 10.2298/CSIS141122050Z
- 5). WB. Wang, M.L. Huang, J. Zhang, and W. Lai, (2015), "Detecting Criminal Relationships Through SOM Visual Analytics", 19th International Conference on Information Visualization (IV2015), In Proceeding of 19th International Conference on Information Visualization, IEEE Computer Society, pp. 316-321, July 2015, DOI:10.1109/IV.2015.62
- 6). J. Zhang and M.L. Huang, (2015), "A New Analytics Model for Large Scale Multidimensional Data Visualization", Cloud Computing and Big Data (CloudCom-Asia 2015), Lecture Notes in Computer Science, Vol. 9106, pp. 55-71, June 2015, DOI: 10.1007/978-3-319-28430-9\_5
- 7). J. Zhang, M.L. Huang, W.B. Wang, L.F. Lu, and Z.P. Meng, (2014), "Big Data Density Analytics using Parallel Coordinate Visualization", The 13th

- International Symposium on Pervasive System, Algorithm, and Network (I-SPAN2014). In Proceeding of IEEE 17th International Conference on Computational Science and Engineering (CSE), IEEE Computer Society, pp. 1115-1120, Dec 2014, DOI:10.1109/CSE.2014.219
- 8). W.B. Wang, M.L. Huang, L.F. Lu, and J. Zhang, (2014), “Improving Performance of Forensics Investigation with Parallel Coordinates Visual Analytics”, The 8th International Conference on Frontier of Computer Science and Technology (FCST2014), In Proceeding of IEEE 17th International Conference on Computational Science and Engineering (CSE), IEEE Computer Society, pp.1838-1843, Dec 2014, DOI:10.1109/CSE.2014.337
  - 9). J. Zhang, M.L. Huang, and Z. Meng, (2014), “BigData Visualization: Parallel Coordinates using Density Approach”, The 2nd International Conference on System and Informatics (ICSAI2014), In Proceeding of IEEE 2nd International Conference on System and Informatics, IEEE Computer Society, pp. 1056-1063, Nov 2014, DOI:10.1109/ICSAI.2014.7009441
  - 10). J. Zhang and M.L. Huang, (2014), “Density approach: a new model for BigData analysis and visualization”, Concurrency and Computation: Practice and Experience, Vol. 28, Issue. 3, pp 661-673, First publish online July 2014, DOI: 10.1002/cpe.3337
  - 11). J. Zhang and M.L. Huang, (2013), “5Ws Model for BigData Analysis and Visualization”, The 2nd International Conference on Big Data Science and Engineering (BDSE2013), In Proceeding of IEEE 16th International Conference on Computational Science and Engineering (CSE), IEEE Computer Society, pp. 1021-1028, Dec 2013, DOI:10.1109/CSE.2013.149
  - 12). J. Zhang and M.L. Huang, (2013), “Detecting Flood Attack through New Density-Pattern Based Approach”, The 15th International Conference on High Performance Computing and Communication (HPCC2013), In Proceeding of 15th IEEE International Conference on High Performance Computing and Communication, IEEE Computer Society, pp. 246-253, Nov 2013, DOI: 10.1109/HPCC.and.EUC.2013.44

- 13). M.L. Huang and J. Zhang, (2013), “Visual Analysis and Detection of Network Flood Attack through Two-Layer Density Approach”, The 3rd International Conference on Computer Science and Network Technology (ICCSNT2013), In Proceeding of 3rd IEEE International Conference on Computer Science and Network Technology, IEEE Computer Society, pp. 625-629, Oct 2013, DOI: 10.1109/ICCSNT.2013.6967191
  
- 14). J. Zhang and M.L. Huang, (2013), “Visual Analytics Model for Intrusion Detection in Flood Attack”, The 12th International Conference on Trust Security and Privacy in Computing and Communications (TrustCom2013), In Proceeding of 12th IEEE International Conference on Trust, Security and Privacy in Computing and Communications, IEEE Computer Society, pp. 277-284, July 2013, DOI:10.1109/TrustCom.2013.38
  
- 15). J. Zhang, M.L. Huang, and D.B. Hoang, (2013), “Visual Analytics for Intrusion Detection in Spam Emails”, International Journal of Grid and Utility Computing, InderScience Enterprises, Vol. 4, No. 2/3, pp. 178-186, Jan 2013, DOI:10.1504/IJGUC.2013.056254
  
- 16). J. Zhang, M.L. Huang, and D.B. Hoang, (2011), “Detecting DDoS Attack in Spam Emails using Density-Weight Model”, The 2nd International Conference on Information Theory and Information Security, In Proceeding of 2nd IEEE International Conference on Information theory and Information Security, IEEE Press, Vol. II, pp. 344-352, Nov 2011, <https://opus.lib.uts.edu.au/research/handle/10453/29567>
  
- 17). M.L. Huang, J. Zhang, Q. Nguyen, and J. Wang, (2011), “Visual Clustering of Spam Emails for DDoS Analysis”, The 15th International Conference on Information Visualization (IV2011), In Proceeding of 15th IEEE international Conference on Information Visualization, IEEE Computer Society, pp. 65-72, July 2011, DOI:10.1109/IV.2011.41