

Optimization of Parallel Coordinates for Visual Analytics

By

Liang Fu Lu

Supervisor: Mao Lin Huang

*A Thesis submitted in Fulfillment for the Degree of Doctor of
Philosophy*

in

**Faculty of Engineering and IT
University of Technology, Sydney Australia**

March 2016

CERTIFICATE OF AUTHORSHIP/ORIGINALITY

UNIVERSITY OF TECHNOLOGY SYDNEY

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

How time flies. It has been nearly five years since for my research study I first took my feet on the grand and attractive university-University of Technology, Sydney, where I had been dreaming for years to go. Fortunately, I acquired a variety of knowledge in the short but precious semesters like a tiny boat loading to full capacity of great achievements, expectations and dreams was sailing backward in the sun shine through the ordeal from the sea storm. At the tranquil and solemn campus, everywhere is surrounded by the strong and profound academic atmosphere, resounded with vigorous study aspirations and left the track of our familiar figure shadows and good memories, which arouses a strong sense of gratitude in my heart.

Firstly, I would like to gratefully acknowledge Mao Lin Huang, an A/Prof in the field of information visualization, who has who has earned a great reputation at home and abroad for his numerous contributions to scientific researches and practical achievements, for the enthusiastic supervision in the process of the research. This thesis grew out of a series of dialogues with him. Each step of completing the thesis including selection, proposal, research and writing is dissolved in Prof. Huang's hard work and guidance. What was deeply impressed me is not the toil and sweat he scattered in numerous readings and modifications, but the wise academic vision, noble life faith and rigorous scholarship attitude he had, which will benefit me a lot in my rest whole life. He brought me closer to the reality I had initially perceived, and eventually enabled me to grasp its rich complexity.

Furthermore, I owe sincere thankfulness to the members of Visualization Team and fellow researchers and the professors of iNext Research centre. Their scholarly attainments and instructions in humor, serious and patient way gave me great help and

inspirations, and that besides gaining knowledge also benefitted me many like friendship, concern, trust and so on.

This research also benefited tremendously from many researchers and staffs in the University of Technology, Sydney. In addition, thank you to the participants in the usability study for the cooperation and valuable feedbacks which are taken as a key factor in providing a necessary guarantee on the datum I collected in accuracy and precision and the evidences in objectivity and scientificity for the thesis.

Last but not least, I owe my deepest gratitude to Huan Xu, my wife, who bravely took the burden of supporting the family alone without any complaints and totally poured her love onto my daughter and me and my whole family. All of her toil and sweat, tolerance and understanding, selfless support and devotion set up the strong will and belief for me in my study. I really appreciated sincerely my extended family for their continuous encouragement and support to make this PhD thesis possible. Thanks my daughter Wenxuan Lu for constantly delivering me angel smiles which accompanied me all the way. Her sweet smiles once encouraged me stand up again after fall a million times and then bravely and firmly faced the difficulties I had encountered. You are all my love most worthy of spending my lifetime loving.

Contents

Contents	v
Figure List	vii
Table List	viii
Equation List	ix
Algorithm List	x
Abstract	xi
Chapter 1. INTRODUCTION	1
1.1 INFORMATION VISUALIZATION	1
1.2 HIGH DIMENSIONAL DATA VISUALIZATION	5
1.3 PARALLEL COORDINATE PLOTS	9
1.4 RESEARCH CHALLENGES.....	13
1.5 RESEARCH OBJECTIVES.....	16
1.6 CONTRIBUTIONS.....	19
1.7 THESIS ORGANIZATION	23
Chapter2. VERTICES OPTIMIZATION IN PARALLEL COORDINATE PLOTS	27
2.1 CLUTTER DESCRIPTION IN PCP.....	27
2.2 NEW ALGORITHM FOR CLUTTER REDUCTION	30
2.3 CASE STUDIES	35
2.4 SUMMARY	45
Chapter 3. NEW AXES RE-ORDERING METHOD IN PARALLEL COORDINATE PLOTS	47
3.1 SIMLARIITY MEASURE AND DIMENSION RE-ORDERING METHODS	47
3.1.1 SIMLARIITY MEASURE.....	48
3.1.2 DIMENSION RE-ORDERING METHODS.....	49
3.2 NEW APPROACH FOR DIMENSION RE-ORDERING	53
3.2.1 Linear/Nonlinear Correlation.....	53
3.2.2 Similarity-based Reordering.....	56
3.3 CASE STUDIES.....	58

3.3.1	Cars dataset	59
3.3.2	Liver disorders dataset	61
3.4	SUMMARY	65
Chapter 4. USING ARCED AXES IN PARALLEL		
COORDINATE GEOMETRY.....		67
4.1	ANALYSIS OF PCP	67
4.2	OVERVIEW OF APPROACHES ON PCP	69
4.3	ARC-BASED PARALLEL COORDINATES GEOMETRY	74
4.3.1	Optimizing Length of Arced Axis.....	74
4.3.2	Arc-Coordinate Geometry	76
4.3.3	Contribution-Based Layout.....	80
4.4	CASE STUDIES	82
4.4.1	Random and Car Datasets	82
4.4.2	Case Study in Network Security Domain	84
4.5	SUMMARY	87
Chapter 5. CONCLUSION AND FUTURE WORK		89
PUBLICATION LIST		93
APPENDIX		95
REFERENCES		96

Figure List

FIGURE 1. SCIENTIFIC VISUALIZATION FOR STANFORD BUNNY(WIJK 2002).	2
FIGURE 2. PARALLEL COORDINATE PLOTS FOR CAR DATASET.	2
FIGURE 3. CLASSIFICATION OF INFORMATION VISUALIZATION.(KEIM 2002).....	4
FIGURE 4. SPACE-SCALE DIAGRAMS ILLUSTRATING TWO EXPERIMENTAL COMPARISONS OF ZOOMING INTERFACES WITH DISPLAY SIZE VARIED BETWEEN D AND 2D. (JAKOBSEN AND HORNBAEK 2013).....	5
FIGURE 5. CHERNOFF FACE VISUALIZATION ON HIGH DIMENSIONAL DATA(KABULOV AND TASHPULATOVA 2010).....	6
FIGURE 6. SCATTERPLOT MATRIX FOR A 7-DIMENSIONAL CAR DATASET.(ÉLMQVIST, DRAGICEVIC ET AL. 2008)	7
FIGURE 7. 3D SCATTERPLOT MATRIX SHOWING THE 8D “OLIVE OIL” DATA SET. (SANFTMANN AND WEISKOPF 2012).....	8
FIGURE 8. FIVE CLUSTERS IN 2D VISUALIZATION OF 100,000 ARTIFICIALLY GENERATED DATA ITEMS(KEIM AND KRIGEL 1994).9	
FIGURE 9. EXAMPLES OF VISUAL CLUTTER IN PARALLEL COORDINATES VISUALIZATION. SEE REGIONS BOUNDED BY TWO ELLIPSES THAT CONTAIN A LARGE NUMBER OF EDGE CROSSINGS.	29
FIGURE 10. THE EXAMPLE DISPLAY OF UNCERTAIN VALUES VISUALIZED IN PARALLEL COORDINATES. THE LINES IN RED AND GREEN BEHAVE REALISTIC AND UNCERTAIN DATA RESPECTIVELY. THE DUMMY VERTICES ARE SHOWN BY GREEN CIRCLE. .	32
FIGURE 11: CASE 1 - RANDOM DATA IN PARALLEL COORDINATES: (A) THE INITIALIZATION OF INCOMPLETE DATA ITEMS; (B) VISUALIZATION OF SUBOPTIMUM POSITIONS OF UNCERTAIN VALUES; (C) THE OPTIMAL POSITIONS OF VERTICES.	37
FIGURE 12. CASE 2 - AN INCOMPLETE DATASET AMEX A VISUALIZED IN PARALLEL COORDINATE VISUALIZATION: (A) THE INITIAL DRAWING OF THE GIVEN DATA WITH TEN UNCERTAIN VALUES; (B) THE NEW DRAWING OF THE SAME GIVEN DATA AFTER THE IMPLEMENTATION OF OUR OPTIMIZATION METHOD. THE DATA SOURCE IS AVAILABLE AT: HTTP://DAVIS.WPI.EDU/XMDV/DATASETS/AMEXA.HTML	38
FIGURE 13: CASE 3 - FORBES 94, A DATASET WITH 5 VARIABLES VISUALIZED IN PARALLEL COORDINATE VISUALIZATION: (A) ORIGINAL PLOT; (B) AFTER CLUTTER REDUCTION. DATA FROM HTTP://WWW-STAT.WHARTON.UPENN.EDU/WATERMAN/FSW/DATASETS/TXT/FORBES94.TXT	44
FIGURE 14. CARS DATASET VISUALIZATION IN PARALLEL COORDINATES.	63
FIGURE 15. AXES REORDERING VISUALIZATION OF LIVER DISORDERS DATASET.	64
FIGURE 16. DUALITY PROPERTY BETWEEN POINTS AND LINES IN CARTESIAN AND PARALLEL COORDINATE PLOTS(WEGMAN 1990).	71
FIGURE 17. THE RATIONALE OF ARC COORDINATES PLANE.	76
FIGURE 18. RANDOM DATA REPRESENTED IN TWO DIFFERENT COORDINATES SYSTEMS	83
FIGURE 19. CAR DATASET VISUALIZED IN PCP AND ACP RESPECTIVELY	86
FIGURE 20. DETECTING DDoS ATTACKS USING ACP: RED AND GREEN LINES DESCRIBE THE SMURF AND NEPTUNE ATTACKS RESPECTIVELY	86

Table List

TABLE 1 DETAILS OF DATA SETS USED IN THE THESIS.	25
TABLE 2 CLUTTER REDUCTION USING OPTIMAL ORDERING ALGORITHM FOR ALL CASES.	41
TABLE 3 CLUTTER REDUCTION IN DIMENSION DECOMPOSITION FOR CASE 3 FORBES94.....	42
TABLE 4 THE COMPARISON OF THE SIMILARITY VALUES USING PCC AND NCC TO CARS DATASET.	65

Equation List

EQ. 1	33
EQ. 2	34
EQ. 3	54
EQ. 4	55
EQ. 5	55
EQ. 6	56
EQ. 7	58
EQ. 8	58
EQ. 9	74
EQ. 10	75
EQ. 11	75
EQ. 12	76
EQ. 13	79
EQ. 14	79
EQ. 15	79
EQ. 16	82

Algorithm List

ALGORITHM 1. DETERMINATION OF POSITIONS OF INCOMPLETE DATA	31
ALGORITHM 2. SIMILARITY-BASED REORDERING ALGORITHM	57

Abstract

The visualization and interaction of multidimensional data always requires optimized solutions for integrating the display, exploration and analytical reasoning of data into a kind of visual pipeline for human-centered data analysis and interpretation. However, parallel coordinate plot, as one of the most popular multidimensional data visualization techniques, suffers from a visual clutter problem. Although this problem has been addressed in many related studies, computational cost and information loss still hamper the application of these techniques, which leads to large high dimensional data sets. Therefore, the main goal of this thesis is to optimize the visual representation of parallel coordinates based on their geometrical properties.

At the first stage, we set out to find optimization methods for permuting data values displayed in parallel coordinate plot to reduce the visual clutter. We divide the dataset into two classifications according to the values and the geometric theory of the parallel coordinate plot: numerical data and non-numerical data, and missing data may exist between them occasionally. We apply Sugiyama's layered directed graph drawing algorithm into parallel coordinate plot to minimize the number of edge crossing among polygonal lines. The methods are proved to be valuable as it can optimize the order of missing or non-numerical value to tackle clutter reduction.

In addition, it is true that optimizing the order is a NP-complete problem, though changing the order of the axis is a straightforward way to address the visual clutter problem. Therefore, we try to propose in the research a new axes re-ordering method in parallel coordinate plot: a similarity-based method, which is based on the combination

of Nonlinear Correlation Coefficient (NCC) and Singular Value Decomposition (SVD) algorithms. By using this approach, the first remarkable axis can be selected based on mathematical theory and all axes can be re-ordered in line with the degree of similarities among them. We also propose a measurement of contribution rate of each dimension to reveal the property hidden in the dataset.

In the third stage, we put forward a new projection method which is able to visualize more data items in the same display space than the existing parallel coordinate methods. Moreover, it is demonstrated clearly in the research that the new method enjoys some elegant duality properties with parallel coordinate plot and Cartesian orthogonal coordinate representation. Meanwhile, the mean crossing angles and the amount of edge crossing between the neighboring axes are utilized in this research to demonstrate the rationale and effectiveness of our approaches.