

UNIVERSITY OF TECHNOLOGY SYDNEY
Faculty of Engineering and Information Technology

**Efficient Query Processing and Analytics on High
Dimensional Data**

by

Wanqi Liu

A THESIS SUBMITTED
IN PARTIAL FULFILLMENT OF THE
REQUIREMENTS FOR THE DEGREE

Doctor of Philosophy

Sydney, Australia

2020

Certificate of Authorship/Originality

I, Wanqi Liu declare that this thesis, is submitted in fulfilment of the requirements for the award of Doctor of Philosophy, in the Faculty of Engineering and Information Technology at the University of Technology Sydney.

This thesis is wholly my own work unless otherwise reference or acknowledged. In addition, I certify that all information sources and literature used are indicated in the thesis.

This document has not been submitted for qualifications at any other academic institution.

This research is supported by the Australian Government Research Training Program.

Candidate signature

Production Note:
Signature removed prior to publication.

Date: 28/03/2020

ABSTRACT

Efficient Query Processing and Analytics on High Dimensional Data

by
Wanqi Liu

As a fundamental problem in query processing, similarity search has been applied in many fields including multimedia, machine learning, database, recommendation systems and so on. Generally, it will be challengeable when it comes to the high-dimensional space due to "the curse of dimensionality". Since it would be too expensive to find exact results, approximate solutions have been studied in many papers. There are various distance metrics to evaluate the similarity between the query object and other points in a dataset. In this thesis, we focus on some well-known distance metrics including Euclidean distance and inner product. Except for Euclidean space, we also study query processing on graphs and propose a novel distance metric on graph. The thesis contains four similarity search problems regarding to different distance metrics, which are Approximate Nearest Neighbour, Approximate Inner Product Search, Approximate Furthest Neighbour and Skyline Nearest Neighbour. Given a query point, the four problems all focus on retrieving a set of "similar" points from the dataset.

Given a set of d -dimensional data points, and a query point q , Approximate Nearest Neighbour Search (ANNS) aims to find the approximate closest object to q in the set. More specifically, in this thesis, we focus on c -ANNS problem, which means given a constant c , the purpose is to find a result whose distance is not larger than c times of the exact smallest distance with a certain possibility. Even though this problem has been researched for a long time, there are still some shortage of current algorithms. We studied the existed works, and proposed a novel I/O efficient algorithm to solve c -approximate nearest neighbour problem in external memory, which can dramatically reduce I/O cost and provide rigorous proof of its correctness.

Maximum Inner Product Search (MIPS) is another valuable problem. It returns an object with maximum inner product value to query point q . There are hundreds of solutions for MIPS but still short of comprehensive evaluation and analysis of these methods' performance. In this thesis, we chose several state-of-the-art algorithms of MIPS using different techniques, and conducted a set of comprehensive experiments to evaluate their performance fairly.

Approximate Furthest Neighbour Search is an opposite problem of Nearest Neigh-

bour Search. It finds the furthest object to query point q in a dataset instead of the closest one. Since most recent works for approximate furthest neighbour search in external memory are only suitable for low-dimensional data, we proposed a new I/O efficient technique to achieve a better performance on I/O cost.

In addition to Euclidian space, similarity search is also a fundamental problem in other spaces like graphs. Considering real-world applications, the multi-layer graph model is extensively studied to reveal the multi-dimensional relations between the graph entities. In this thesis, we formulated a new problem called skyline nearest neighbor search on multi-layer graphs, and proposed a baseline algorithm, and two optimizations instead of naively adopting the traditional skyline procedure as a subroutine. We also investigated the rule to optimize search order in the algorithm so that further improve the algorithmic efficiency.

Acknowledgements

Firstly, I wish to express my deepest gratitude to my supervisor Prof. Ying Zhang who is professional, diligent and patient, for offering the opportunity to study at UTS, and for the continuous support, patient guidance and enthusiastic encouragement he has provided through the three years. Discussions with him always enlighten me to improve my works in this thesis. He also helped me to extend my knowledge and gave me many inspirations. It would never have been possible for me to complete this work without his incredible support and valuable ideas.

Secondly, I wish to express my sincere appreciation to Prof. Wei Wang for his advice, and for his brilliant ideas which inspired me deeply. Additionally, his valuable comments from his abundant research experience helped me on my research work. Acknowledgement also goes to Prof. Xumin Lin and my co-supervisor Dr. Lu Qin for their supporting and guidance.

I am also sincerely grateful to Mr. Hanchen Wang as most of the works in this thesis are conducted collaborating with him. I really appreciate for his patient to assist me in proofreading and revising the drafts. He is always able to find out grammar mistakes which I missed. My sincere thanks also goes to Dr. Dong Wen who shows advanced research skills and helps me to refine my work.

I would also like to thank Mr. Kai Wang, Mr. Mingjie Li, Mr. Bohua Yang, Ms. Conggai Li, Mr. Boge Liu, Mr. Wentao Li, Dr. Yang Yang, Dr. Haida Zhang, Dr. Fan Zhang, Dr. Ouyang Dian, Mr. Junda Lu, Mr. Yufei Wang for giving me a pleasant time when working with them.

Last but not least, I would like to express my great appreciation to my parents for raising me with care and love, providing me with unfailing support and continuous encouragement throughout my years of study and through the process of researching and writing this thesis, and to my grandparents for their understanding and love. Thanks to my boyfriend Chenhao Pan for his love, understanding and continuing support to complete the research work.

Wanqi Liu
Sydney, Australia, March 2020.

List of Publications

- **Wanqi Liu**, Hanchen Wang, Ying Zhang, Wei Wang and Lu Qin, I-LSH: I/O efficient c -Approximate Nearest Neighbor Search in High-dimensional Space, published in ICDE 2019
- **Wanqi Liu**, Dong Wen, Hanchen Wang, Fan Zhang and Xubo Wang, Skyline Nearest Neighbor Search on Multi-Layer Graphs, published in ICDE 2019 workshop
- **Wanqi Liu**, Hanchen Wang, Ying Zhang, Luqin and Wenjie Zhang, I/O efficient algorithm for c -approximate furthest neighbor search in high-dimensional space, will appear in DASFAA 2020
- **Wanqi Liu**, Hanchen Wang, Ying Zhang, Wei Wang, Lu Qin and Xuemin Lin, EI-LSH: An Early-Termination Driven I/O Efficient Incremental c -Approximate Nearest Neighbor Search, Received by VLDBJ

Contents

Certificate	ii
Abstract	iii
Acknowledgments	v
List of Publications	vi
List of Figures	xi
Abbreviation	xiii
1 Introduction	1
1.1 Background	2
1.1.1 Approximate Nearest Neighbour Search	2
1.1.2 Approximate Maximum Inner Product Search	3
1.1.3 Approximate Furthest Neighbour Search	4
1.1.4 Skyline Nearest Neighbour Search	4
1.2 Motivations	5
1.2.1 Approximate Nearest Neighbour Search	5
1.2.2 Approximate Maximum Inner Product Search	6
1.2.3 Approximate Furthest Neighbour Search	7
1.2.4 Skyline Nearest Neighbour Search	8
1.3 Contributions	8
1.3.1 Approximate Nearest Neighbour Search	9
1.3.2 Maximum Inner Product Search	10
1.3.3 Approximate Furthest Neighbour Search	10
1.3.4 Skyline Nearest Neighbour Search	11
2 Literature Survey	12
2.1 Approximate Nearest Neighbour Search	12
2.1.1 List-based LSH algorithm	12

2.1.2	Tree-based LSH method	14
2.1.3	Other relevant work	15
2.2	Approximate Maximum Inner Product Search	16
2.3	Approximate Furthest Neighbour Search	17
2.4	Skyline Nearest Neighbour Search	18
3	Approximate Nearest Neighbour Search	20
3.1	Overview	20
3.2	Motivation	20
3.3	Preliminary	21
3.3.1	Problem Definition	23
3.3.2	LSH with bucket partitioning	23
3.4	Our Approach	25
3.4.1	Motivation	26
3.4.2	Incremental LSH with new ET	28
3.4.3	Incremental LSH with traditional ET	32
3.4.4	Extension for c - k -ANN problem	32
3.5	Analysis	34
3.6	Performance Studies	44
3.6.1	Experiment Setup	44
3.6.2	Evaluate Index Size	47
3.6.3	Evaluate Index Building time	48
3.6.4	Evaluate I/O costs	49
3.6.5	Evaluate Accuracy	50
3.6.6	Effect of the approximate ratio c	52
3.6.7	Effect of the value of P_{ET}	54
3.6.8	Large dataset	54
3.6.9	Summary	55
3.7	Conclusion	57
4	Approximate Maximum Inner Product Search	58
4.1	Overview	58

4.2	Background	58
4.2.1	Problem definition	58
4.2.2	Algorithm scope	59
4.2.3	Categories	60
4.3	Approximate MIPS algorithms with theoretical guarantee	60
4.3.1	Locality-Sensitive Hashing	61
4.3.2	From LSH to MIPS	62
4.3.3	H2ALSH	63
4.3.4	L2-ALSH	64
4.3.5	Sign-ALSH	65
4.3.6	Norm-range LSH	65
4.4	Approximate MIPS algorithms without theoretical guarantee	66
4.4.1	Graph-based algorithms	66
4.4.2	Tree-based algorithms	69
4.5	Experiments	70
4.5.1	Experiment settings	70
4.5.2	Evaluation Metrics	71
4.5.3	Parameter settings.	72
4.5.4	Comparison within each category	73
4.5.5	Second round comparing	79
4.5.6	Conclusion	81
4.6	Summary	84
5	Approximate Furthest Neighbour Search	86
5.1	Overview	86
5.2	Motivation	86
5.3	Preliminary	87
5.3.1	Problem Definition	87
5.3.2	LSH family for furthest neighbor	88
5.3.3	$(c, 1, p_1, p_2)$ -sensitive Reverse LSH	89
5.4	Approach	90

5.4.1	Motivation	90
5.4.2	Approach	91
5.4.3	c - k -AFN	94
5.5	Analysis	94
5.6	Experiments	97
5.6.1	Experiment Setup	97
5.6.2	Index Size	98
5.6.3	I/O costs	99
5.6.4	Running time	100
5.6.5	I/O and ratio	101
5.6.6	Summary	102
5.7	Conclusion	103
6	Skyline Nearest Neighbour Search on multi-layer graph	104
6.1	Overview	104
6.2	Preliminaries	104
6.2.1	Problem definition	104
6.3	Approach	107
6.4	Optimizations	108
6.4.1	An Early-Stop Approach	108
6.4.2	Layer chosen strategy	111
6.5	Performance Studies	114
6.5.1	Experiment setup	114
6.5.2	Running time	115
6.5.3	Effectiveness of the optimizations	115
6.6	Conclusion	116
7	Conclusion	119
7.1	Contributions	119
7.2	Future work	120
	Bibliography	122

List of Figures

1.1	Multi-layer graph example	9
3.1	Motivation for EI-LSH	21
3.2	Example: When ET could fail	26
3.3	Example for incremental search	28
3.4	Example of $f_h(x)$	35
3.5	The sample space \mathcal{S}	37
3.6	Motivation of early termination	39
3.7	Success possibility δ	42
3.8	I/O costs varying k	51
3.9	I/O costs vs. ratio	53
3.10	I/O costs vs. c	53
3.11	I/O costs vs. P_{ET}	54
3.12	I/O cost on difference dataset size ($k = 50$)	56
4.1	recall varying k	76
4.2	running time varying k	77
4.3	running time vs. recall	78
4.4	running time varying k	81
4.5	recall varying k	82
4.6	recall vs. time	83
4.7	recall vs. time	84
4.8	recall vs. speedup	85
4.9	Recall $\geq 95\%$	85
5.1	Example for separation	89

5.2	Motivation for RI-LSH	92
5.3	I/O costs varying k	100
5.4	Running time varying k	101
5.5	I/O vs. ratio	102
6.1	Skyline nearest neighbors of v_0 on graph 1.1	107
6.2	Procedure of optimization	110
6.3	a graph suitable for SNNS-ET2	112
6.4	running time	115
6.5	number of visited vertices	116

Abbreviation

NN: Nearest Neighbour

ANN: Approximate Nearest Neighbour

c -ANN: c -Approximate Nearest Neighbour

MIPS: Maximum Inner Product Search

FN: Furthest Neighbour

AFN: Approximate Furthest Neighbour

c -AFN: c -Approximate Furthest Neighbour

ET: Early-Termination

NT: Normal-Termination

SNNS: Skyline Nearest Neighbour Search