

Nearest Neighbor Search in High Dimensional Space

by

Mingjie Li

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

Doctor of Philosophy

Centre for Artificial Intelligence

Faculty of Engineering and Information Technology

University of Technology Sydney

December, 2020

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 part of requirements for a degree at any other academic institution 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.

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

Production Note:
Signature: Signature removed
prior to publication.

Date: 10/12/2020

ACKNOWLEDGEMENTS

Foremost, I would like to express my sincere gratitude to my supervisor Prof. Ying Zhang for his continuous support and guidance throughout my PhD career. Ying is professional, patient and kind. He introduced me to the research area of nearest neighbor search and further provided constant motivation which kept me going. His valuable ideas and suggestions always guided me and broadened my knowledge and horizon in the related areas. I am very thankful to him for his patience and confidence during my PhD study. Even experiencing failures and difficulties in research and life, his constant encouragement and support always kept me optimistic and positive. Additionally, Ying is a good mentor and friend for me. He gave me many useful advices on planing career development.

I would like to thank Prof. Wei Wang for his constructive ideas and invaluable suggestions on the works in this thesis. The insightful discussions with Prof. Wang always gave me many inspirations. In addition, I would like to thank Prof. Ivor W. Tsang and A/Prof. Lu Qin for the discussions and suggestions on the related topic.

I would like to thank Prof. Xuemin Lin and A/Prof. Wenjie Zhang for supporting the works in this thesis, as most of the works were conducted

in collaboration with them. I thank Prof. Lin for offering a rigorous while interesting research environment.

I am thankful to the faculties and staffs at the school of computer science at University of Technology Sydney. They were very helpful and supportive throughout my PhD study. It was indeed a pleasure to be a part of such an exciting community.

I would like to thank Dr. Xin Cao, Dr. Lijun Chang, Dr. Xiaoyang Wang, Dr. Shiyu Yang, Dr. Zengfeng Huang, and Dr. Yixiang Fang for sharing the ideas and experiences. Thanks to the members from database groups at UNSW and UTS, including Dr. Long Yuan, Dr. Longbin Lai, Dr. Xiang Wang, Dr. Xing Feng, Dr. Xubo Wang, Dr. Haida Zhang, Dr. Yang Yang, Mr. Xuefeng Chen, Dr. You Peng, Dr. Boge Liu, Ms. Xiaoshuang Chen, Mr. Yuren Mao, Dr. Fan Zhang, Dr. Dong Wen, Dr. Dian Ouyang, Ms. Wanqi Liu, Mr. Hanchen Wang, and Mr. Yuanhang Yu, for creating a dynamic and vibrant atmosphere in the labs and in life. I would also like to thank Mr. Daniel Ouyang, Mr. Peng Zhang, and Mr. Xunxiang Yao for their selfless help and care during my PhD life.

Last but not least, I would like to thank my father Mr. Zheng Li, my mother Mrs. Fengying Zhu, my brother Mr. Mingchang Li and my sister Mrs. Jenly Li for their continuous support, encouragement and love during my entire PhD journey and in my life. I am greatly indebted to them.

PUBLICATIONS

- **Mingjie Li**, Ying Zhang, Yifang Sun, Wei Wang, Ivor W. Tsang, Xuemin Lin. An Efficient Exact Nearest Neighbor Search by Compounded Embedding. DASFAA 2018. (Chapter 4)
- **Mingjie Li**, Ying Zhang, Yifang Sun, Wei Wang, Ivor W. Tsang, Xuemin Lin. I/O Efficient Approximate Nearest Neighbour Search based on Learned Functions. ICDE 2020. (Chapter 5)
- Wen Li, Ying Zhang, Yifang Sun, Wei Wang, **Mingjie Li***, Wenjie Zhang, Xuemin Lin. Approximate Nearest Neighbor Search on High Dimensional Data - Experiments, Analyses, and Improvement. TKDE 2019. (*Corresponding Author)(Chapter 6)

TABLE OF CONTENT

CERTIFICATE OF AUTHORSHIP/ORGINALITY	iii
ACKNOWLEDGEMENTS	iv
PUBLICATIONS	vi
TABLE OF CONTENT	vii
LIST OF FIGURES	x
LIST OF TABLES	xii
ABSTRACT	xiii
Chapter 1 Introduction	1
1.1 Exact Nearest Neighbor Search	2
1.2 Approximate Nearest Neighbor Search by Learning to hash	5
1.3 Approximate Nearest Neighbor Search: An Experimental Study	8
Chapter 2 Literature Review	12
2.1 Exact Nearest Neighbor Search	12
2.2 Approximate Nearest Neighbor Search	14
2.2.1 Hashing-based Methods	15
2.2.2 Partition-based Methods	20
2.2.3 Graph-based Methods	21
Chapter 3 Problem Statement	23
3.1 Problem Definition	23
3.2 Notations	24
Chapter 4 Exact Nearest Neighbor Search by Compounded Em- bedding	26
4.1 Overview	26

TABLE OF CONTENT

4.2	Embedding and Distance Lower Bound	26
4.2.1	Motivation	27
4.2.2	Embedding Method	28
4.2.3	Correctness of Distance Lower Bound	29
4.2.4	Optimization of Distance Lower Bound	31
4.2.5	Using PCA Technique	33
4.3	Efficient Exact NNS Algorithm	34
4.3.1	Motivation	34
4.3.2	Exact NNS Algorithm	36
4.3.3	Performance Analysis	38
4.3.4	Discussion	39
4.4	Experiments	40
4.4.1	Experimental Settings	40
4.4.2	Performance Evaluation	42
4.5	Conclusion	47
Chapter 5 Approximate Nearest Neighbor Search By Learned Functions		48
5.1	Overview	48
5.2	Our ANNS Framework	49
5.2.1	Our ANNS Solution	49
5.2.2	Performance Analysis	52
5.3	Learning to Index by Linear Hashing	53
5.3.1	Linear Model and Its Objective Function	53
5.3.2	Relaxation and Optimization	55
5.4	Learning to Index by Neural Network	60
5.4.1	DNN Architecture	60
5.4.2	Objective Function	61
5.5	Discussion	62
5.6	Experiments	63
5.6.1	Experimental Settings	63
5.6.2	Parameter Tuning	68
5.6.3	Performance Comparison	70
5.6.4	Summary	74
5.7	Conclusion	76
Chapter 6 Approximate Nearest Neighbor Search: An Experimental Evaluation		77
6.1	Overview	77
6.2	Evaluation Scope	78
6.3	The State-of-the-art ANNS Algorithms	79

6.3.1	LSH-based methods	79
6.3.2	L2H-based methods	80
6.3.3	Partition-based Algorithms	82
6.3.4	Graph-based Algorithms	83
6.4	Diversified Proximity Graph	84
6.4.1	Motivation	85
6.4.2	Diversified Proximity Graph	86
6.5	Experiments	87
6.5.1	Experimental Settings	87
6.5.2	Evaluation Measures	89
6.5.3	Comparison with Each Category	90
6.5.4	Second Round Evaluation	93
6.5.5	Summary	99
6.6	Further Analyses	102
6.6.1	Space Partition-based Approach	102
6.6.2	Graph-based Approach	104
6.7	Conclusion	105
	Chapter 7 Epilogue	107
	REFERENCES	109

LIST OF FIGURES

1.1	Illustration of our idea of index and query processing. Grey points are embedding values of data points, orange points are embedding values of the query point (i.e., q_1^*, \dots, q_m^*).	8
4.1	Motivation of Distance Lower Bound	28
4.2	Illustration of Our Embedding Method	29
4.3	Motivation of Exact NNS using Embedded Space	34
4.4	Search time with respect to m and n	42
4.5	Comparison of search time on all datasets	42
4.6	Pruning performance (lower the better)	44
4.7	Comparison with respect to k	45
4.8	Pre-processing time	46
4.9	Speedup with respect to recall	46
5.1	The Sigmoid Function	56
5.2	The Architecture of Our Non-Linear Hash Learning	60
5.3	The impact of parameters of OPFA on Deep	68
5.4	The impact of parameters of NeOPFA on Deep	69
5.5	I/O Cost with respect to k on all datasets	71
5.6	Ratio with respect to k on all datasets	72
5.7	Recall with respect to k	73
5.8	Search Time with respect to k	73
5.9	Pre-processing Time on All Datasets	75
6.1	Motivation of Diversified Proximity Graph	85
6.2	Speedup vs Recall for LSH-based and L2H-based Methods	91
6.3	Speedup vs Recall for Partition-based and Graph-based Methods	92
6.4	Speedup with Recall of 0.8	94
6.5	Recall with Speedup of 50	94
6.6	Speedup vs Recall on Different Datasets	95
6.7	Recall vs Percentage of Data Points Accessed	96
6.8	Accuracy vs Recall	97
6.9	The Ratio of Index Size and Data Size (%)	98

6.10	Index Construction Time (seconds)	98
6.11	Index Memory Cost (MB)	98
6.12	Precision vs Recall	99
6.13	F1 score vs Recall	99
6.14	Analyses of Space Partitioning-based Methods	103
6.15	minHops Distributions of KGraph and DPG	104

LIST OF TABLES

2.1	Overview of the Indexing and Searching of Existing Exact NNS Algorithms	14
2.2	Overview of the Indexing and Searching of Representative ANNS Algorithms	16
3.1	Summary of Notations	24
4.1	Dataset Summary	41
5.1	Statistics of Datasets	65
5.2	Parameter Settings of OPFA	66
5.3	Index Sizes of All Algorithms (in Megabytes)	74
6.1	Dataset Summary	89
6.2	mAP for each algorithm	100
6.3	Ranking of the Algorithms Under Different Criteria	101

ABSTRACT

Nearest neighbor search (NNS) in high dimensional space is a fundamental and essential operation in applications from many domains, such as machine learning, databases, multimedia and computer vision, to name a few. In this thesis, we investigate both exact and approximate NNS in high dimensional space.

For the exact NNS, we propose an efficient technique which can have a significant speedup over the state-of-the-art exact solutions. Specifically, we first propose a novel compounded embedding technique, by which we achieve a tight distance lower bound for Euclidean distance. Then each point in a high dimensional space can be embedded into a low dimensional space such that the distance between two embedded points lower bounds their distance in the original space. Following the *filter-and-verify* paradigm, we develop an efficient exact NNS algorithm by pruning a large number of candidates using the new lower bounding technique. Comprehensive experiments on many real-world data demonstrate the effectiveness and efficiency of our new algorithm.

In terms of the approximate NNS, we propose an external memory-based approximate NNS algorithm by learning to hash. Specifically, we introduce a novel data-sensitive indexing and query processing framework for approximate NNS with an emphasis on optimizing the I/O efficiency, especially, the sequential

I/Os. The proposed index consists of several lists of point IDs, ordered by values that are obtained by learned hashing functions on each corresponding data point. The functions are learned from the data and approximately preserve the order in the high-dimensional space. We consider two instantiations of the functions (linear and non-linear), both learned from the data with novel objective functions. Comprehensive experiments on six large scale high dimensional datasets show that our proposed methods with learned index structure perform much better than the state-of-the-art external memory-based approximate NNS methods in terms of I/O efficiency and search accuracy.

Although lots of approximate NNS algorithms have been continuously proposed in the literature each year, there is no comprehensive evaluation and analysis of their performance. Therefore, we conduct a comprehensive and systematic experimental evaluation for the state-of-the-art approximate methods. Our study (1) is cross-disciplinary (i.e., including 19 algorithms in different domains, and from practitioners) and (2) has evaluated a diverse range of settings, including 20 datasets, several evaluation metrics, and different query workloads. The experimental results are carefully reported and analyzed to understand the performance results. Furthermore, we propose a new method that achieves both high query efficiency and high recall empirically on majority of the datasets under a wide range of settings.