

学校编码：10384
学号：23220131153333

分类号 _____ 密级 _____
UDC _____

厦 门 大 学

硕 士 学 位 论 文

局部敏感哈希改进算法研究

Research on Improved Locality Sensitive Hashing Algorithm

岑伟

指导教师姓名：缪克华 副教授

专业名称：模式识别与智能系统

论文提交日期：2016 年 月

论文答辩时间：2016 年 月

学位授予日期：2016 年 月

答辩委员会主席：

评 阅 人：

2016 年 月

厦门大学博硕士学位论文摘要库

厦门大学学位论文原创性声明

本人呈交的学位论文是本人在导师指导下，独立完成的研究成果。本人在论文写作中参考其他个人或集体已经发表的研究成果，均在文中以适当方式明确标明，并符合法律规范和《厦门大学研究生学术活动规范（试行）》。

另外，该学位论文为（）课题（组）的研究成果，获得（）课题（组）经费或实验室的资助，在（）实验室完成。（请在以上括号内填写课题或课题组负责人或实验室名称，未有此项声明内容的，可以不作特别声明。）

声明人（签名）：

年 月 日

厦门大学学位论文著作权使用声明

本人同意厦门大学根据《中华人民共和国学位条例暂行实施办法》等规定保留和使用此学位论文，并向主管部门或其指定机构送交学位论文（包括纸质版和电子版），允许学位论文进入厦门大学图书馆及其数据库被查阅、借阅。本人同意厦门大学将学位论文加入全国博士、硕士学位论文共建单位数据库进行检索，将学位论文的标题和摘要汇编出版，采用影印、缩印或者其它方式合理复制学位论文。

本学位论文属于：

1. 经厦门大学保密委员会审查核定的保密学位论文，
于 年 月 日解密，解密后适用上述授权。
2. 不保密，适用上述授权。

（请在以上相应括号内打“√”或填上相应内容。保密学位论文应是已经厦门大学保密委员会审定过的学位论文，未经厦门大学保密委员会审定的学位论文均为公开学位论文。此声明栏不填写的，默认为公开学位论文，均适用上述授权。）

声明人（签名）：

年 月 日

摘要

“大数据”时代给数据检索带来了新的挑战，相似性检索显得尤为重要。局部敏感哈希算法是相似性检索中最流行的一种，该算法是建立在哈希的基础上的一种近似最近邻算法，它能将检索时间复杂度缩减到线性。与其它基于 Tree 的数据结构相比，局部敏感哈希算法能较好的处理数据在高维空间中的检索问题。注意到高维数据利用该算法检索得到候选集后需要进行相似度计算，而这一部分所耗时间的占整个数据检索时间的比重非常大，所以该算法在处理大规模数据时的性能仍然需要提高。

针对局部敏感哈希算法，考虑到数据维数和数据量的急剧增加而检索时间还不能满足需求情况下，提出了在不同场合下两种改进的局部敏感哈希算法并应用于图像，本文在局部敏感哈希算法的基础上做了以下工作：

1. 对利用局部敏感哈希算法查找前 k 个最相似 (Top- k) 的问题，提出了一种基于次数排序的局部敏感哈希算法 (ST-LSH)。该算法将唯一化前的索引号按出现次数进行排序来输出，避免了相似度计算花费大量时间的问题。实验结果显示，与改进前的算法相比，在保证识别率基本不变的情况下，改进后的算法大大减少了数据检索时间。

2. 在利用局部敏感哈希算法进行数据类别的查找过程中，为了弥补容易受噪声点的影响，提出了一种基于 k 近邻分类 (k -nearest neighbor) 局部敏感哈希算法 (KNN-LSH)。KNN-LSH 充分利用了唯一化后的候选集的信息，利用候选集唯一化后的索引号所属的类别进行 k 近邻分类。即统计当前候选集所属类别的信息，将查询对象归为所属类别最多的那一类，从而避免了大量的相似度计算。实验结果证明，与改进前的算法相比，KNN-LSH 对分类问题的识别率要好于传统方法，具有稳定性，同时也缩短了数据检索时间。

关键字：相似性检索；局部敏感哈希；高维数据；

Abstract

"Big Data" brings new challenges to data retrieval, and similarity search is particularly important. Locality sensitive hashing algorithm is one of the most popular algorithms in similarity search, it is an approximate nearest neighbor algorithm that based on the establishment of hashing, it reduces the complexity of retrieval time to linear. When compared with other data structure algorithms that based on the Tree , locality sensitive hashing can handle data better in high-dimensional space retrieval. Noting that after getting retrieved candidate set for high dimensional data, there will need similarity computation, this part takes a big proportion of time-consuming in the whole retrieval time, so the algorithm performances still needs to improve when dealing with large-scale data.

For locality sensitive hashing algorithm, taking into account the data dimension and a sharp increase in the amount of data retrieval time still can't satisfy the demand, put forward two improved locality sensitive hashing algorithms in different applications and applied them to the image, this paper do the following work:

1. When using locality sensitive hashing algorithm to find the k most similar (Top-k) objects, put forward a locality sensitive hashing algorithm based on the sort of times (ST-LSH). In order to avoid the problem large number of similarity computation, this algorithm outputs data according to the number of occurrences before unique the data. Experimental results show that compared with traditional locality sensitive hashing, improved algorithm ensures the recognition rate nearly unchanged, and greatly reduces data retrieval time.

2. When using locality sensitive hashing algorithm to find data label , in order to compensate its susceptible to noise points, propose locality sensitive hashing algorithm based on k-nearest neighbor classifier (k-nearest neighbor) (KNN-LSH) . KNN-LSH takes full advantage of the information after unique the candidate set, use categories of candidate set after unique technology to k-nearest classification. which means statistic the information of current candidate set belongs to a certain category,

the query object is classified to the largest number of the category , thus avoiding a lot of similarity computation. Experimental results show that when compared with the pre-modified algorithm, KNN-LSH classification performs better recognition rate than the traditional methods, with stability, but also shorten the data retrieval time.

Keywords: similarity search; locality sensitive hashing; high-dimensional data;

厦门大学博硕士论文摘要库

厦门大学博硕士学位论文摘要库

目 录

第一章 绪论	1
1.1 课题背景及研究意义	1
1.2 国内外研究现状	4
1.3 主要研究内容及组织结构	5
第二章 局部敏感哈希算法的基本思想	7
2.1 距离测量度量	7
2.1.1 距离测度的含义	7
2.1.2 欧式距离	7
2.1.3 Jaccard 距离	8
2.1.4 海明距离	8
2.2 最近邻和近似最近邻	8
2.2.1 最近邻	10
2.2.2 近似最近邻	10
2.2.3 K 最近邻	11
2.3 局部敏感哈希算法的简介	12
2.3.1 算法描述	12
2.3.2 实现细节	14
2.3.3 局部敏感哈希函数族	20
2.3.4 局部敏感哈希应用	22
2.4 算法衡量指标	22
2.4.1 查准率和查全率	22
2.4.2 时间效率	23
2.4.3 空间效率	23
2.4.4 高维数据适应性	23
2.5 本章小结	26
第三章 几种改进的局部敏感哈希算法	27

3.1 基于海明距离的局部敏感哈希算法	27
3.1.1 算法简介	27
3.1.2 算法描述	27
3.2 基于 P 稳定分布局部敏感哈希算法	30
3.2.1 算法简介	30
3.2.2 算法描述	31
3.3 基于多探寻局部敏感哈希算法	33
3.3.1 基于熵的局部敏感索引	34
3.3.2 多探寻局部敏感索引	34
3.4 本章小结	36
第四章 基于次数排序的局部敏感哈希算法	39
4.1 改进的依据	39
4.2 改进算法描述	41
4.2.1 哈希表的构建	41
4.2.2 候选集的获取	41
4.2.3 数据处理	43
4.3 实验数据及实验环境	44
4.3.1 实验数据集	44
4.3.2 实验环境	45
4.4 实验结果	45
4.5 本章小结	49
第五章 基于 k 近邻分类的局部敏感哈希算法	51
5.1 改进的依据	51
5.2 改进算法描述	51
5.2.1 候选集的获取	51
5.2.1 数据处理	53
5.3 实验结果	54
5.4 本章小结	56

第六章 总结与展望	57
6.1 论文总结	57
6.2 研究和展望	58
参考文献	59
攻读硕士期间研究成果	63
致 谢.....	64

厦门大学博硕士论文摘要库

厦门大学博硕士学位论文摘要库

Contents

Chapter 1 Introduction.....	1
1.1 Background and Research Meaning	1
1.2 Research Status of Locality Sensitive Hashing.....	4
1.3 Cotents of The Paper	5
Chapter 2 The Basic Idea of Locality Sensitive Hashing	7
2.1 Distance Measurement	7
2.1.1 Introduction of Distance Measurement.....	7
2.1.2 Euclidean Distance.....	7
2.1.3 Jaccard Distance.....	8
2.1.4 Hamming Distance.....	8
2.2 Nearest Neighbor and Approximate Nearest Neighbor.....	8
2.2.1 Nearest Neighbor.....	10
2.2.2 Approximate Nearest Neighbor	10
2.2.3 K-Nearest Neighbor	11
2.3 Introduction to Locality Sensitive Hashing	12
2.3.1 Algorithm Description	12
2.3.2 Implementaion Details	14
2.3.3 Hash Family	20
2.3.4 Applications	22
2.4 Standard of Algorithm.....	22
2.4.1 Accurate	22
2.4.2 Time Efficient	23
2.4.3 Space Efficient.....	23
2.4.4 High Dimensional	23
2.5 Summary.....	26
Chapter 3 Serveral Improvements for Locality Sensitive Hashing	27

3.1 Locality Sensitive Hashing Based on Hamming Distance	27
3.1.1 Algorithm Introduction	27
3.1.2 Algorithm Description	27
3.2 Locality Sensitive Hashing Based on p-Stable Distributions	30
3.2.1 Algorithm Introduction	30
3.2.2 Algorithm Description	31
3.3 Locality Sensitive Hashing Based on Multi-Probe.....	33
3.3.1 Entropy-Based Locality Sensitive Hashing Indexing	34
3.3.2 Multi-Probe Locality Sensitive Hashing Indexing	34
3.4 Summary.....	36
Chapter 4 Locality Sensitive Hashing Based on Frequency	39
4.1 The Background of Improved Algorithm	39
4.2 Description of Improved Algorithm	41
4.2.1 Construction of Hash Table.....	41
4.2.2 Obtain of Candidate sets	41
4.2.3 Data Processing.....	43
4.3 Experimental Data Set and Enviroment.....	44
4.3.1 Experimental Data Set	44
4.3.2 Experimental Environment	45
4.4 Experiment Results.....	46
4.5 Summary.....	46
Chapter 5 Locality Sensitive Hashing Based on K-Nearest Neighbor Classification	51
5.1 The Background of Improved Algorithm	51
5.2 Description of Improved Algorithm	51
5.2.1 Obtain of Candidate sets	51
5.2.1 Data Processing.....	53
5.3 Experiment Results.....	54

5.4 Summary.....	56
Chapter 6 Conclusions and Prospect	57
6.1 Conclusion of The Paper	57
6.2 Prospect of Study.....	57
References	59
Publications.....	63
Acknowledgements	65

厦门大学博硕士学位论文摘要

厦门大学博硕士学位论文摘要库

Degree papers are in the “[Xiamen University Electronic Theses and Dissertations Database](#)”.

Fulltexts are available in the following ways:

1. If your library is a CALIS member libraries, please log on <http://etd.calis.edu.cn/> and submit requests online, or consult the interlibrary loan department in your library.
2. For users of non-CALIS member libraries, please mail to etd@xmu.edu.cn for delivery details.