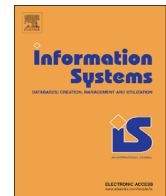




ELSEVIER

Contents lists available at ScienceDirect

Information Systems

journal homepage: www.elsevier.com/locate/infosys

An empirical evaluation of intrinsic dimension estimators [☆]

Gonzalo Navarro ^a, Rodrigo Paredes ^{b,*}, Nora Reyes ^{c,*}, Cristian Bustos ^c

^a Center of Biotechnology and Bioengineering, Department of Computer Science, University of Chile, Chile

^b Departamento de Ciencias de la Computación, Universidad de Talca, Chile

^c Departamento de Informática, Universidad Nacional de San Luis, Argentina

ARTICLE INFO

Article history:

Received 16 December 2015

Accepted 10 June 2016

Keywords:

Intrinsic dimension

Complexity of searching

Metric spaces

ABSTRACT

We study the practical behavior of different algorithms and methods that aim to estimate the intrinsic dimension (IDim) in metric spaces. Some of them were specifically developed to evaluate the complexity of searching in metric spaces, based on different theories related to the distribution of distances between objects on such spaces. Others were originally designed for vector spaces only, and have been extended to general metric spaces. To empirically evaluate the fitness of various IDim estimations with the actual difficulty of searching in metric spaces, we compare two representatives of each of the broadest families of metric indices: those based on pivots and those based on compact partitions. Our conclusions are that the estimators Distance Exponent and Correlation fit best their purpose.

© 2016 Elsevier Ltd. All rights reserved.

1. Introduction

Similarity search in metric spaces has received much attention due to its applications in many fields, ranging from multimedia information retrieval to machine learning, classification, and searching the Web. While a wealth of practical algorithms exist to handle this problem, it has been often noted that some datasets are intrinsically harder to search than others, no matter which search algorithms are used. An intuitive concept of “curse of dimensionality” has been coined to denote this intrinsic difficulty, but a clear method to measure it, and thus to predict the performance of similarity searching in a space, has been elusive.

The similarity between a set of objects \mathfrak{U} is modeled using a *distance function* (or *metric*) $d: \mathfrak{U} \times \mathfrak{U} \rightarrow \mathbb{R}^+ \cup \{0\}$

[☆] Partially funded by basal funds FB0001, Conicyt, Chile and Fondecyt grant 1131044, Chile.

* Corresponding authors.

E-mail addresses: gnavarro@dcc.uchile.cl (G. Navarro), rapared@utalca.cl (R. Paredes), nreyes@unsl.edu.ar (N. Reyes), cjbustos@unsl.edu.ar (C. Bustos).

<http://dx.doi.org/10.1016/j.is.2016.06.004>

0306-4379/© 2016 Elsevier Ltd. All rights reserved.

that satisfies the properties of triangle inequality, strict positivity, reflexivity, and symmetry. In this case, the pair (\mathfrak{U}, d) is called a *metric space* [1–4].

In some applications, the metric spaces are of a particular kind called “vector spaces” of finite *explicit* or *representational* dimension, where the elements consist of D coordinates of real numbers. In this case, we can use some Minkowski metric or any other metric appropriate to the specific case (for instance, the cosine distance) as the dissimilarity measure between two objects. Many works exploit the geometric properties of vector spaces, but they usually cannot be extended to general metric spaces, where the only available information is the distance between objects. Since in most cases the distance is very expensive to compute, the main goal when searching in metric spaces is to reduce the number of distance evaluations. In contrast, vector space operations tend to be cheaper and the primary goal when searching them is to reduce the CPU cost or the number of I/O operations carried out.

Similarity queries are usually of two types. For a given database $S \subseteq \mathfrak{U}$ with size $|S| = n$, $q \in \mathfrak{U}$ and $r \in \mathbb{R}^+$, the *range query* $(q, r)_d$ returns all the objects of S at distance at most r

1 from q , formally $(q, r)_d = \{x \in S, d(x, q) \leq r\}$; whereas the
 2 nearest neighbor query $kNN_d(q)$ retrieves the k elements of
 3 S that are closest to q , that is, $kNN_d(q)$ is a set such that for
 4 all $x \in kNN_d(q)$ and $y \in S \setminus kNN_d(q)$, $d(q, x) \leq d(q, y)$, and
 5 $|kNN_d(q)| = k$.

6 A naïve way to answer similarity queries is to compare
 7 all the database elements with the query q and return
 8 those elements that are close enough to q . This *brute force*
 9 approach is too expensive for real applications. Research
 10 has then focused on ways to reduce the number of distance
 11 computations performed to answer similarity queries. There has been significant progress around the idea of
 12 building an *index*, that is, a data structure that allows
 13 discarding some database elements without explicitly
 14 comparing them to q . Moreover, there are some relatively
 15 recent works [5–10] that try to get jointly the goals of
 16 reducing the number of distance evaluations and the
 17 number of I/O operations performed.

18 In vector spaces with uniformly distributed data, the
 19 *curse of dimensionality* describes the well-known expo-
 20 nential increase of the cost of all existing search algo-
 21 rithms as the dimension grows. Non-uniformly distributed
 22 vector spaces may be easier to search than uniform ones,
 23 despite having the same explicit dimensionality. The
 24 phenomenon also extends to general metric spaces despite
 25 their absence of coordinates: some spaces are intrinsically
 26 harder to search than others. This has led to the concept
 27 of *intrinsic dimensionality* (IDim) of a metric space, as
 28 a measure of the difficulty of searching it. A reliable measure
 29 of IDim has been elusive, despite the existence of several
 30 formulae.

31 Computing the IDim of a metric space is useful, for
 32 example, to determine whether it is amenable to indexing
 33 at all. If the IDim is too high, then we must just resort to
 34 brute-force solutions or to approximate search algorithms
 35 (which do not guarantee to find the exact answers). Even
 36 when exact indexing is possible, the IDim helps decide
 37 which kind of index to use and how to tune it. For
 38 example, in lower dimension spaces, a pivot-based
 39 method works fine using a small set of pivots; whereas
 40 in higher dimensions we need to use a large set of pivots
 41 [1], which also implies a large amount of memory for the
 42 index. Alternatively, if we do not have enough extra
 43 memory for the index, we can switch to the *List of Clusters*
 44 [11], which has reasonable performance in high dimension
 45 spending little space in the index.

46 In this work we aim to empirically study the fitness of
 47 various IDim measures to predict the search difficulty of
 48 metric space searching. Some measures were specifically
 49 developed for metric spaces, based on different theories
 50 related to the distribution of distances between objects.
 51 Others were originally designed for vector spaces and have
 52 then been adapted to general metric spaces. We chose
 53 various synthetic and real-life metric spaces and four
 54 indexing methods that are representatives of the major
 55 families of indices: two based on pivots and two based on
 56 compact partitions. Our comparison between real and
 57 estimated search difficulty yields that *Distance Exponent*
 58 [12,13] and *Correlation* [14] are currently the best pre-
 59 dictors in practice, however all the estimators behave
 60 relatively well.

61 The rest of this paper is organized as follows. In [Section](#)
 62 [2](#), we review some relevant issues of IDim estimators for
 63 vector spaces. Next, in [Section 3](#), we survey four methods
 64 for estimating IDim in vector spaces and show how to
 65 adapt them to the metric case. We also include three new
 66 IDim estimators for general metric spaces. The experi-
 67 mental evaluation for the seven methods is presented in
 68 [Section 4](#). We finally draw our conclusions and future work
 69 directions in [Section 5](#). An early version of this work
 70 appeared in [15].

2. Intrinsic dimension estimators for vector spaces

71 There are several interesting applications where the
 72 data are represented as D -dimensional vectors in \mathbb{R}^D . For
 73 instance, in pattern recognition applications, objects are
 74 usually represented as vectors [16]. Therefore, data are
 75 embedded in \mathbb{R}^D , even though this does not imply that its
 76 *intrinsic* dimension is D .

77 There are many definitions of IDim. For instance, the
 78 IDim of a given dataset is the minimum number of free
 79 variables needed to represent the data without loss of
 80 information [17]. In general terms, a dataset $\mathbf{x} \subseteq \mathbb{R}^D$ has
 81 IDim $M \leq D$, if its elements fall completely within an M -
 82 dimensional manifold of \mathbb{R}^D [18]. Another intuitive notion
 83 is the logarithm of the search cost, as in many cases this
 84 cost grows exponentially with the dimension.

85 Even in vector spaces, there are many reasons to esti-
 86 mate the IDim of a dataset. Using more dimensions (more
 87 coordinates in the vectors) than necessary can bring sev-
 88 eral problems. For example, the space to store the data
 89 may be an issue. A dataset $\mathbf{x} \subseteq \mathbb{R}^D$ with $|\mathbf{x}| = n$ requires to
 90 store $n \times D$ real coordinates. Instead, if we know that the
 91 IDim of \mathbf{x} is $M \leq D$, we can map the points to \mathbb{R}^M and just
 92 store $n \times M$ real coordinates. The CPU cost to compute a
 93 distance is also reduced. This can in addition help identify
 94 the important dimensions in the original data. Also, as the
 95 amount of available information increases, compressing the
 96 data storage becomes even more important. Secondly,
 97 as the asymptotic complexity of the algorithms is mono-
 98 tonically increasing with respect to the dataset dimen-
 99 sionality, a dimensionality reduction (to the actual dataset
 100 IDim) can produce an important CPU time reduction. For
 101 instance, in the case of data classification or pattern
 102 recognition, producing reliable classifiers is difficult when
 103 the dataset dimensionality is high (*curse of dimensionality*
 104 [19]); and according to the theoretical approximation of
 105 statistical learning [20], the classifier generalization cap-
 106 ability depends on the IDim of the space.

107 There are two approximations to estimate the IDim of a
 108 vector space [16,17], namely, *local* and *global* methods. The
 109 local ones make the estimation by using the information
 110 contained in sample neighborhoods, avoiding the data
 111 projection over spaces of lower dimensionality. The global
 112 ones deploy the dataset over an M -dimensional space
 113 using all the dataset information. Unlike the local methods
 114 that only use the information contained in the neighbor-
 115 hood of each data sample, global methods use whole
 116 information of the dataset.

Download English Version:

<https://daneshyari.com/en/article/4945155>

Download Persian Version:

<https://daneshyari.com/article/4945155>

[Daneshyari.com](https://daneshyari.com)