

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*

Abstract

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

Keywords: intrinsic dimension, complexity of searching, metric spaces

1. Introduction

12 Similarity search in metric spaces has received much attention due to its
13 applications in many fields, ranging from multimedia information retrieval to
14 machine learning, classification, and searching the Web. While a wealth of
15 practical algorithms exist to handle this problem, it has been often noted that
16 some datasets are intrinsically harder to search than others, no matter which
17 search algorithms are used. An intuitive concept of “curse of dimensionality”
18 has been coined to denote this intrinsic difficulty, but a clear method to
19 measure it, and thus to predict the performance of similarity searching in a
20 space, has been elusive.
21

[☆]Partially funded by basal funds FB0001, Conicyt, Chile and Fondecyt grant 1131044.
Preprint submitted to Information Systems, May 23, 2016
Chile.

*Corresponding author

Email addresses: gnavarro@dcc.uchile.cl (Gonzalo Navarro),
raparede@utalca.cl (Rodrigo Paredes), nreyes@unsl.edu.ar (Nora Reyes),
cjbustos@unsl.edu.ar (Cristian Bustos)

22 The similarity between a set of objects \mathbb{U} is modeled using a *distance*
23 *function* (or *metric*) $d : \mathbb{U} \times \mathbb{U} \mapsto \mathbb{R}^+ \cup \{0\}$ that satisfies the properties of
24 triangle inequality, strict positivity, reflexivity, and symmetry. In this case,
25 the pair (\mathbb{U}, d) is called a *metric space* [1, 2, 3, 4].

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

38 Similarity queries are usually of two types. For a given database $S \subseteq \mathbb{U}$
39 with size $|S| = n$, $q \in \mathbb{U}$ and $r \in \mathbb{R}^+$, the *range query* $(q, r)_d$ returns all the ob-
40 jects of S at distance at most r from q , formally $(q, r)_d = \{x \in S, d(x, q) \leq r\}$;
41 whereas the *nearest neighbor query* $kNN_d(q)$ retrieves the k elements of S that
42 are closest to q , that is, $kNN_d(q)$ is a set such that for all $x \in kNN_d(q)$ and
43 $y \in S \setminus kNN_d(q)$, $d(q, x) \leq d(q, y)$, and $|kNN_d(q)| = k$.

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

54 In vector spaces with uniformly distributed data, the *curse of dimension-*
55 *ality* describes the well-known exponential increase of the cost of all existing
56 search algorithms as the dimension grows. Non-uniformly distributed vector
57 spaces may be easier to search than uniform ones, despite having the same ex-
58 plicit dimensionality. The phenomenon also extends to general metric spaces
59 despite their absence of coordinates: some spaces are intrinsically harder to

60 search than others. This has led to the concept of *intrinsic dimensionality*
61 (*IDim*) of a metric space, as a measure of the difficulty of searching it. A
62 reliable measure of IDim has been elusive, despite the existence of several
63 formulae.

64 Computing the IDim of a metric space is useful, for example, to determine
65 whether it is amenable to indexing at all. If the IDim is too high, then we
66 must just resort to brute-force solutions or to approximate search algorithms
67 (which do not guarantee to find the exact answers). Even when exact index-
68 ing is possible, the IDim helps decide which kind of index to use and how to
69 tune it. For example, in lower dimension spaces, a pivot-based method works
70 fine using a small set of pivots; whereas in higher dimensions we need to use
71 a large set of pivots [1], which also implies a large amount of memory for the
72 index. Alternatively, if we do not have enough extra memory for the index,
73 we can switch to the *List of Clusters* [11], which has reasonable performance
74 in high dimension spending little space in the index.

75 In this work we aim to empirically study the fitness of various IDim
76 measures to predict the search difficulty of metric space searching. Some
77 measures were specifically developed for metric spaces, based on different
78 theories related to the distribution of distances between objects. Others
79 were originally designed for vector spaces and have then been adapted to
80 general metric spaces. We chose various synthetic and real-life metric spaces
81 and four indexing methods that are representatives of the major families of
82 indices: two based on pivots and two based on compact partitions. Our
83 comparison between real and estimated search difficulty yields that *Distance*
84 *Exponent* [12, 13] and *Correlation* [14] are currently the best predictors in
85 practice, however all the estimators behave relatively well.

86 The rest of this paper is organized as follows. In Section 2, we review
87 some relevant issues of IDim estimators for vector spaces. Next, in Section 3,
88 we survey four methods for estimating IDim in vector spaces and show how
89 to adapt them to the metric case. We also include three new IDim estimators
90 for general metric spaces. The experimental evaluation for the seven methods
91 is presented in Section 4. We finally draw our conclusions and future work
92 directions in Section 5. An early version of this work appeared in [15].

93 **2. Intrinsic Dimension Estimators for Vector Spaces**

94 There are several interesting applications where the data are represented
95 as D -dimensional vectors in \mathbb{R}^D . For instance, in pattern recognition appli-

96 cations, objects are usually represented as vectors [16]. Therefore, data are
97 embedded in \mathbb{R}^D , even though this does not imply that its *intrinsic* dimension
98 is D .

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

105 Even in vector spaces, there are many reasons to estimate the IDim of a
106 dataset. Using more dimensions (more coordinates in the vectors) than nec-
107 essary can bring several problems. For example, the space to store the data
108 may be an issue. A dataset $\mathbb{X} \subseteq \mathbb{R}^D$ with $|\mathbb{X}| = n$ requires to store $n \times D$
109 real coordinates. Instead, if we know that the IDim of \mathbb{X} is $M \leq D$, we
110 can map the points to \mathbb{R}^M and just store $n \times M$ real coordinates. The CPU
111 cost to compute a distance is also reduced. This can in addition help iden-
112 tify the important dimensions in the original data. Also, as the amount of
113 available information increases, compressing the data storage becomes even
114 more important. Secondly, as the asymptotic complexity of the algorithms is
115 monotonically increasing with respect to the dataset dimensionality, a dimen-
116 sionality reduction (to the actual dataset IDim) can produce an important
117 CPU time reduction. For instance, in the case of data classification or pat-
118 tern recognition, producing reliable classifiers is difficult when the dataset
119 dimensionality is high (*curse of dimensionality* [19]); and according to the
120 theoretical approximation of statistical learning [20], the classifier general-
121 ization capability depends on the IDim of the space.

122 There are two approximations to estimate the IDim of a vector space [16,
123 17], namely, *local* and *global* methods. The local ones make the estimation
124 by using the information contained in sample neighborhoods, avoiding the
125 data projection over spaces of lower dimensionality. The global ones deploy
126 the dataset over an M -dimensional space using all the dataset information.
127 Unlike the local methods that only use the information contained in the
128 neighborhood of each data sample, global methods use whole information of
129 the dataset.

130 In this work we focus on global IDim estimators. That is, we consider
131 all the dataset information to estimate the IDim as accurately as possible.
132 Global methods can be split into three families: projection techniques, mul-
133 tidimensional scaling methods, and fractal based methods. The last two are

134 more suitable to extend to metric spaces, so we have selected and adapted
135 some representatives of these groups.

136 **3. Intrinsic Dimension Estimators for Metric Spaces**

137 In general metric spaces, since the curse of dimensionality severely affects
138 the performance of the search process, knowing the IDim can help choose a
139 metric index appropriate to the space dimension and also give some insight on
140 the specific index tuning. For instance, in low IDim spaces, where searching
141 is easier, pivot based indices usually perform better, even when using a small
142 set of pivots. However, they can fail in high IDim spaces, or hard spaces,
143 as a large set of pivot is needed to preserve the performance at the cost of
144 an excessive amount of space for the index. Alternatively, if there is little
145 amount of extra memory, we can use the *List of Clusters* (LC) [11], which is
146 a very appropriate, RAM economical index.

147 Hence, a proper estimation of the operating dataset IDim is very impor-
148 tant, as it helps improve the time and memory costs of the selected solution.

149 There are few IDim estimators that apply directly in general metric
150 spaces. The IDim estimators that are proper to metric spaces can only
151 consider the dataset objects and their distances between each other.

152 In this section we analyze various methods to estimate the IDim of vector
153 spaces and others to general metric spaces. We discuss how to adapt the
154 former to the case of general metric spaces. Note that, since multidimensional
155 spaces are a particular case of metric spaces, our estimators can also be
156 applied to obtain the IDim of D -dimensional vector spaces.

157 *3.1. Fractal Based Methods*

158 Unlike other families, fractal based methods can estimate non-integer
159 IDim values. The most popular techniques of this family are *Box Counting*
160 [21], which is a simplified version of the *Haussdorff dimension* [22, 23], and
161 *Correlation* [14]. These techniques have been successfully used to estimate
162 the dimensionality of the underlying dynamic systems that generate time
163 series [24].

164 The dimension estimation by Box Counting D_B of a set $\Omega \subseteq \mathbb{R}^D$ is defined
165 as follows: if $v(r)$ is the number of boxes of size r needed to cover Ω , then

$$D_B = \lim_{r \rightarrow 0} \frac{\ln(v(r))}{\ln\left(\frac{1}{r}\right)}. \quad (1)$$

166 In this method, the boxes are multidimensional regions of side r on each
 167 dimension (that is, they are hypercubes of side r). Regrettably, even though
 168 efficient algorithms have been proposed, the Box Counting dimension can be
 169 computed only for low dimension datasets, because its algorithmic complexity
 170 grows exponentially with the dimension.

171 Estimating the dimension by Correlation is an alternative to Box Counting.
 172 It is defined as follows. Let $\Omega = \{x_1, x_2, \dots, x_n\} \subset \mathbb{R}^D$ and the correla-
 173 tion integral

$$C_m(r) = \lim_{n \rightarrow \infty} \frac{2}{n(n-1)} \sum_{1 \leq i < j \leq n} I(\|x_j - x_i\| < r), \quad (2)$$

174 where $I(\cdot)$ is the indicator function. Intuitively, $C_m(r)$ is the fraction of
 175 object pairs whose distance is lower than r . So, the dimension estimation by
 176 Correlation D_C is

$$D_C = \lim_{r \rightarrow 0} \frac{\ln(C_m(r))}{\ln r}. \quad (3)$$

177 3.1.1. Correlation

178 The most popular method to estimate the dimension by Correlation is the
 179 log-log plot. It consists in plotting $\ln(C_m(r))$ versus $\ln(r)$. The dimension by
 180 Correlation is the slope of the linear section of the curve.

181 To illustrate the process of estimating IDim using the Correlation estima-
 182 tor, in Fig. 1 we show an example of its computation on the real world metric
 183 space Histograms (this dataset is described in Section 4.2). The line plotted
 184 by circles corresponds to the curve $\ln(C_m(r))$ versus $\ln(r)$, obtained from the
 185 experimental data. We estimate the IDim of this dataset by computing the
 186 slope of the linear section at the beginning of the plot (drawn with a line).
 187 Note that this is the section of the curve that shows the usual exponential
 188 growth of the fraction $C_m(r)$ with respect to the intrinsic dimensionality of
 189 the space. At the end of the linear section of the plot, $C_m(r)$ almost reaches
 190 its maximum value, so the growth beyond this linear section is very mild.
 191 Hence, we need to neglect this section of the curve, otherwise we can un-
 192 derestimate the space intrinsic dimensionality. To compute the slope we use
 193 linear regression with least squares over the first linear section of the curve.
 194 Note that this procedure allows us to estimate IDim with a dataset of modest
 195 size, because we are only focused on the section of the curve that does reveal
 196 exponential growth.

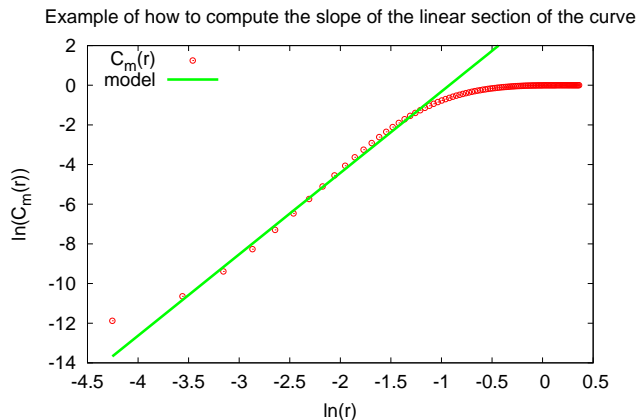


Figure 1: IDim estimation with Correlation by computing the slope of the linear section of the curve.

197 *3.1.2. Ball Counting*

198 Analogously with Correlation, to estimate the dimension by Box Counting, we can compute the slope of the linear section of the curve $\ln(v(r))$ versus $\ln(1/r)$. However, in the general case of metric spaces, we do not have coordinates. Thus, to adapt the Box Counting method, we consider *balls* of radius r , that is, the set of objects within a distance r from a reference object o . We randomly pick the reference objects from the dataset, and count the number $B(r)$ of balls of radius r needed to cover the dataset. To do so, we use the *List of Clusters (LC)* index [11], whose code is available from SISAP [25], with the variant of fixed radius and centers chosen at random. Then, $B(r)$ is just the length of the LC.

208 To estimate the dimension by Box Counting, which in this case is Ball Counting, we replace $\ln(v(r))$ by $\ln(B(r))$, plot $\ln(B(r))$ versus $\ln(\frac{1}{r})$ in log-log and obtain the IDim as the slope of the linear section of the curve by using linear regression with least squares over the experimental data $(\ln(B(r)), \ln(\frac{1}{r}))$, using a procedure analogous to Correlation.

213 *3.2. Distance Exponent*

214 Traina et al. [12, 13] discuss the problem of the selectivity estimation for range queries in real-world metric spaces, including spatial or multidimensional datasets as special cases. It plays an important role when analyzing real metric spaces. Their main finding is that several datasets follow the so-called *Power Law*. They call *Distance Exponent* the exponent of the power

219 law, and show how to use it to derive formulae for estimating the selectivity
 220 of range queries. For instance, the number of objects relevant to the query,
 221 the number of I/Os to answer the query when the data is stored on disk, the
 222 amount of time needed to answer the query, and so on.

223 To find a formula that estimates the number of neighbors of objects within
 224 a distance r in a dataset of n -objects, they introduce the following notions:
 225 (i) the *Distance Plot* of a metric set is the number of object pairs at distance
 226 at most r versus the distance r , and both axes are drawn in logarithmic scale;
 227 and (ii) the *Distance Exponent* is the slope of the line that better fits the
 228 distance plot in case it is linear for a range of scales. Using these two notions,
 229 they define the *Distance Law*.

230 **Definition 1. (*Distance Law*)** Given a dataset of n objects from a metric
 231 space with distance function $d(x, y)$, the average number of distances lower
 232 than a radius r follows a power law; that is, the average number of neighbors
 233 $\overline{nb}(r)$ within a distance r is proportional to $r^{\mathcal{D}}$. Formally,

$$n \cdot \Phi(r) = \overline{nb}(r) \propto r^{\mathcal{D}}, \quad (4)$$

234 where n is the number of objects in the dataset and $\Phi(r)$ is the accumulated
 235 distribution function of the probability of a pair of objects to be within a
 236 distance r .

237 If a dataset has a metric to evaluate the distance between every object
 238 pair, then this plot can always be drawn. They show that the distance plot
 239 has an almost linear behavior for many databases, both real and synthetic.
 240 Building the distance plot requires $O(n^2)$ distance computations. To reduce
 241 this cost, $\overline{nb}(r)$ is estimated using an index [12], in particular the *M-tree* [5].
 242 That is, a way to estimate the distance exponent \mathcal{D} of a dataset stored in a
 243 metric index is by means of the very same index.

244 Since in this work we are only interested in comparing the different IDim
 245 measures, indexing the space is not necessary and we compute $\overline{nb}(r)$ directly,
 246 considering a reference object chosen at random from the dataset. We only
 247 determine the number of elements at distance r from that object. The result
 248 is averaged over various choices for the object.

249 3.3. Fastmap

250 This method arises from the proposal [26] of a fast algorithm to map
 251 objects of any metric space onto points of a k -dimensional space (k being

252 defined by the user), so that the dissimilarities are preserved. Its goal is to
253 speed up searches in traditional or multimedia databases.

254 To do so, the objects are mapped onto the k -dimensional space using
255 k feature extraction functions, provided by domain experts [26]. The main
256 issue is how to define such feature extraction functions. For example, in the
257 metric case of strings with the edit distance [27], it is not clear which features
258 can be considered.

259 For a domain expert, it is generally easier to provide a distance function
260 to compare objects than to provide feature extraction functions. *Fastmap*
261 [28] is a generalization of the original method [26], where the objects are
262 mapped using only a distance function.

263 Fastmap finds, given a dataset of n objects from a metric space (\mathbb{U}, d) , n
264 image points in a k -dimensional target space, such that the distances between
265 the objects in the original space are preserved as much as possible in the
266 target space.

267 For evaluating the dissimilarity preservation in the target space, a *stress*
268 function is defined as follows,

$$stress^2 = \frac{\left(\sum_{i,j} (\hat{d}_{ij} - d_{ij})^2\right)}{\left(\sum_{i,j} d_{ij}^2\right)}, \quad (5)$$

269 where d_{ij} is the dissimilarity measure (the distance of the original space)
270 between objects o_i and o_j , and \hat{d}_{ij} is the Euclidean distance between their
271 respective images p_i and p_j . The stress function gives the relative error that
272 the distances in the target space suffer on average after the transformation.
273 Fastmap begins with an estimation that is iteratively improved, until no
274 additional improvement is possible.

275 In the metric case, we can assume that we have the $n \times n$ matrix of
276 distances between all the dataset objects, and Fastmap must find n points in
277 the k -dimensional space whose Euclidean distances are close to the original
278 matrix of $n \times n$ distances. The crux is to assume that objects are points
279 in some m -dimensional space, with unknown m , and to project these points
280 over k mutually orthogonal directions. The challenge is to compute all these
281 projections using only the distance matrix. Fastmap projects the objects
282 over carefully selected lines. It chooses two objects o_a and o_b , and considers
283 the “line” passing through them in the original space. The projections x'_i of
284 the objects over this line are obtained using the *cosine law*:

285 **Theorem 1. (*Cosine Law*)** Any triangle $o_a \hat{o}_i o_b$ satisfies:

$$d(o_b, o_i)^2 = d(o_a, o_i)^2 + d(o_a, o_b)^2 - 2x'_i d(o_a, o_b). \quad (6)$$

286 Eq. 6 can be solved for x'_i to compute the projection of object o_i with the
287 formula

$$x'_i = \frac{d(o_a, o_i)^2 + d(o_a, o_b)^2 - d(o_b, o_i)^2}{2d(o_a, o_b)}. \quad (7)$$

288 Thus, the input of Fastmap is a set S of size n and, in each iteration, it
289 computes the coordinates of all the n objects over the new axis. So, after
290 k iterations, it produces a k -dimensional target space S' where each object
291 $o_i \in S$ is mapped to a k -coordinate vector $p_i = (x'_{i,1}, x'_{i,2}, \dots, x'_{i,k}) \in S'$,
292 where $x'_{i,j}$ is the j -th projection of the image p_i of the object o_i .

293 In our case, we want to estimate the number of projections needed so
294 that the target space reaches a mapping with a small enough *stress*, that
295 is, preserving the distances sufficiently well. Thus, we modify the Fastmap
296 algorithm so that it computes the *stress* of the target space after each new
297 dimension is added. If the difference between the current and the previous
298 *stress* values is significant, we compute another projection (thus increasing
299 the dimensionality of the target space). Otherwise, the current dimension
300 of the target space is reported as the estimation of the IDim of the original
301 metric space.

302 3.4. *Principal Component Analysis*

303 *Principal Component Analysis (PCA)* [29] is a statistical procedure that
304 projects the data onto new axes, called the principal components, where the
305 axes are ordered by maximum to minimum variance. As the first components
306 accumulate most of the variance, the original data can be projected using
307 the first components controlling how much information we want to preserve
308 or lose (and we can use more components if we want to preserve a larger
309 amount of the data information).

310 A common application of PCA is to reduce the dimensionality of a vector
311 dataset by neglecting the components with small variance, as they have min-
312 imum impact in the amount of information that the projected dataset will
313 have. So, we can identify the number of selected components as the IDim of
314 the dataset.

315 The crux of PCA is that it finds a set of basis vectors, where the first
316 component follows the maximum variance direction, the second follows the

Table 1: Verification of the pivot table IDim for synthetic spaces.

	C5	C10	C15	C20	G5	G10	G15	G20	G101
IDim Space	5.08	8.40	12.13	15.80	4.73	8.83	13.46	15.80	0.92
IDim PT	4.70	7.45	11.17	13.92	4.83	8.01	11.98	13.40	0.96

317 next variance direction, and so on. That is, each component accounts for as
 318 much of the variability in the data as possible. The resulting vectors, called
 319 principal components, are an uncorrelated orthogonal basis set, because they
 320 are the eigenvectors of the data covariance matrix.

321 In the metric case, we do not necessary have objects with coordinates.
 322 So, the first step is to represent the object from a given space as a vector. For
 323 this purpose, we simply select a set of random pivots and compute the pivot
 324 table. Thus, each object is represented as a row in the table, a vector, where
 325 its components are the distances to every pivot. After this, we compute the
 326 principal components of the pivot table. In the process, the components are
 327 sorted by their importance, that is, in decreasing proportion of the variance
 328 of the data. So the first components should accumulate most of the variance.

329 We start the computation with an educated guess. In the synthetic metric
 330 spaces, we use a set of pivots two times bigger than the representational
 331 dimension. The underlying intuition is that we should not estimate an IDim
 332 bigger than the representational dimension, but, as we use random pivots,
 333 we grant the pivot set the chance of incorporating the maximum of distance
 334 information between the objects in the dataset. On the other hand, in the
 335 case of real world metric spaces we use a set of 20 pivots, as our experience
 336 says that the real world metric space under study has a much lower IDim.

337 Since we need to fix a threshold, we use the number of components that
 338 accumulates 90% of the variance of the dataset as the IDim estimation of
 339 the space. This threshold gave us the best results in our experiments. To
 340 perform the statistics computation we use R [30].

341 It is important to verify whether the IDim of the given metric space is
 342 preserved after representing the objects as vectors. To do so, we carry out a
 343 preliminary study which consists in calculating the intrinsic dimensionality
 344 computed with the estimator Distance Exponent over the pivot table. The
 345 results of the study are shown in Tables 1 and 2, where IDim Space is the
 346 IDim estimation that Distance Exponent suggests for the respective space
 347 and IDim PT is the estimation for the vectors in the pivot table.

348 As can be seen, the IDim computed with Distance Exponent over the

Table 2: Verification of the pivot table IDim for real world spaces.

	ENG	NASA	DOCS	HIS
IDim Space	4.96	3.70	3.52	4.70
IDim PT	9.29	3.02	0.57	3.54

349 pivot table is similar to the one computed on the original dataset, with
 350 the exception of the space of strings (ENG) and documents (DOCS). These
 351 results validate our application of PCA. The mismatch in the case of strings
 352 can be due to the discrete nature of the space, and in the case of documents,
 353 to the extremely concentrated histogram of distances. For example, several
 354 other estimators yield similar numbers for these two spaces.

355 3.5. *Intrinsic Search Difficulty*

356 Chávez et al. [1] introduced a measure of the intrinsic complexity of
 357 searching in general metric spaces, which is easy to estimate and is inde-
 358 pendent of the search algorithm.

359 Several authors [31, 32, 33] have proposed to use the *distance histogram*
 360 to characterize the hardness of searching in arbitrary metric spaces, yet the
 361 cost was tailored to a specific index. This measure [1] is the first quantitative
 362 definition. It depends only on the histogram and not on any assumption on
 363 the indexing method.

364 The intuition behind this measure is that, in random vectors in D di-
 365 mensions, the histogram has a larger mean μ and a smaller variance σ^2 as
 366 D increases. In fact, it holds $D = c \cdot \mu^2 / \sigma^2$ for some constant c [34]. Thus,
 367 the same formula could be used to estimate a dimension D from the mean
 368 and variance of the histogram of distances in a general metric space. We do
 369 not have the histogram of the whole universe \mathbb{U} , but we can approximate it
 370 using the histogram of the dataset $S \subset \mathbb{U}$, considering the set S as a random
 371 sample of \mathbb{U} .

372 **Definition 2. (*Intrinsic Search Difficulty*)** Let μ be the mean and σ^2
 373 be the variance of the histogram of distances of a metric space. Then, its
 374 *intrinsic search difficulty* is

$$\rho = \frac{\mu^2}{2\sigma^2}. \tag{8}$$

375 An obvious advantage of this measure, which has contributed to its pop-
376 ularity, is that ρ is easy to compute from a reasonable sampling of pairs in
377 S . Other techniques require more complex and expensive computations.

378 Pestov [35] presents a system of three axioms an intrinsic dimension func-
379 tion must satisfy. He proves that the intrinsic dimension measure ρ satisfies
380 a weak version of these axioms. Later [36], he introduces a set of goals an
381 intrinsic dimension function should fulfill, and ρ achieves many of them.

382 As the measure ρ has been designed for general metric spaces, we use it as
383 is. We consider the dataset S and we compute all the distances $d(x, y), \forall x, y \in$
384 S . Then we compute the average $\mu = \frac{1}{n^2} \sum_{x, y \in S} d(x, y)$ and the variance

385 $\sigma^2 = \frac{1}{n^2} \sum_{x, y \in S} (d(x, y) - \mu)^2$. Finally, we obtain the value of $\rho = \frac{\mu^2}{2\sigma^2}$ and

386 report it as the IDim of the metric space.

387 3.6. Sparse Spatial Selection Method

388 This method is based on a very simple pivot selection strategy [37], called
389 *Sparse Spatial Selection (SSS)*. This strategy has the advantage that it is not
390 required to know the number of pivots in advance.

391 Initially, the set of pivots contains only the first object of the collection.
392 Then, for each element $x \in S$, x is chosen as a new pivot if its distance to
393 every pivot in the current set of pivots is equal or greater than αd^+ , being
394 α a constant parameter (for indexing purposes this constant α takes values
395 around 0.5) and d^+ the maximum distance between two objects in the space.
396 That is, an object in the collection becomes a new pivot if it is located at
397 more than a fraction of the maximum distance d^+ with respect to all the
398 current pivots. For example, if $\alpha = 0.5$ an object is chosen if it is located
399 further than a half of the maximum distance d^+ from the already selected
400 pivots.

401 Therefore, the selected pivots will not be too close to each other. Forcing
402 the distance between two pivots to be greater or equal than αd^+ , it is ensured
403 that they are well distributed in the whole space. It is important to take into
404 account that these pivots are not very far away from each other, neither very
405 far from the rest of the objects in the collection (i.e., they are not necessarily
406 *outliers*), but they are well distributed and cover the whole space.

407 An distinguishing feature of SSS is that an element x is compared against
408 the pivots already selected and it becomes a new pivot if needed. In this way
409 the number of pivots does not depend on the collection size but on its intrinsic

410 dimensionality. This number of pivots is very similar to the optimum number
411 for other strategies.

412 We note that surrounding a pivot, there is a ball of radius αd^+ . Also,
413 the method produces pivots as long as there is some part of the space that is
414 not covered by any previous pivot. This resembles the fractal methods. So,
415 we can use a similar technique to estimate IDim. Let $P(\alpha)$ be the number
416 of pivots produced by SSS for a given value of α . So, we plot $\ln P(\alpha)$ versus
417 $\ln \frac{1}{\alpha}$ and obtain the slope of the linear section of the curve by using linear
418 regression with least squares over the experimental data $(\ln(P(\alpha)), \ln(\frac{1}{\alpha}))$.

419 4. Experimental Evaluation

420 We evaluate experimentally the seven IDim estimators described on general
421 metric spaces. We consider two kinds of metric spaces, depending on
422 the data source:

423 **Synthetic:** these are spaces generated artificially so that they present some
424 interesting characteristic to be evaluated. For instance, uniformly dis-
425 tributed vectors in \mathbb{R}^D with known dimension.

426 **Real world:** these are metric spaces obtained from real-world applications.
427 For instance, a feature vector space of images obtained from a NASA
428 image set.

429 4.1. Synthetic Metric Spaces

430 These are vector spaces with Euclidean distance. They are treated as
431 metric spaces, as we do not consider the coordinate information. A first set
432 is formed by vectors with uniform distribution, so that the representational
433 dimension matches the IDim. Here, we can test the estimators in a case
434 where the IDim is known. A second set is formed by vectors with Gaussian
435 distribution, so that the representational dimension is greater than the IDim
436 (the more clustered is the space, the lower is the IDim). The distance is also
437 Euclidean. Here, we aim to check whether the estimators give lower values
438 as the IDim decreases.

439 4.1.1. Uniformly Distributed Vectors with Euclidean Distance

440 We generate four datasets of 100,000 uniformly distributed vectors in the
441 unitary cube $[0, 1)^D$, with $D = 5, 10, 15,$ and 20 . The spaces are called C5,
442 C10, C15, and C20, respectively.

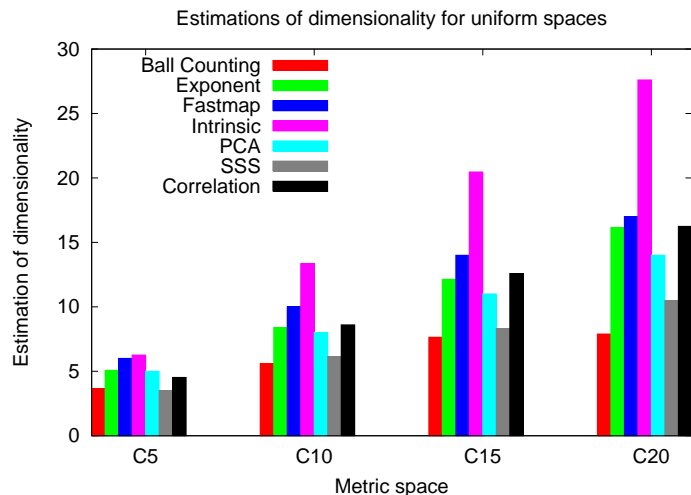


Figure 2: Comparison of dimensionality estimations for uniform spaces.

443 Fig. 2 depicts the estimations obtained with the seven IDim estimators,
 444 namely, Ball Counting, Distance Exponent, Fastmap, Intrinsic Search Diffi-
 445 culty, PCA, SSS, and Correlation, for these four metric spaces. As it can be
 446 seen, Ball counting becomes insensitive to the correct dimension in C20. The
 447 other six methods increase proportionally with D , but with different slopes.
 448 Fastmap is the one with the best fit, matching D almost perfectly, closely
 449 followed by Correlation and Distance Exponent. Intrinsic Search Difficulty
 450 based estimator shows a consistent factor multiplying D . Both PCA and
 451 SSS fit well in C5, but as IDim grows, the fit loses precision.

452 *Search degradation as IDim grows.* To verify that the dataset IDim is re-
 453 sponsible of the search degradation, we pick C5 and extend its vectors with
 454 zeroes to produce spaces with 10, 15, and 20 representational dimensions,
 455 and study the search performance over it.

456 We perform 10 executions of the algorithms, building the index with
 457 90% of the database elements and reserving the remaining 10% (chosen at
 458 random) for the queries. So, the query objects do not belong to the index.
 459 We average the results over the 10 executions. In each execution, the objects
 460 in the metric space are permuted at random. Therefore, each of the 10 indices
 461 uses a different dataset S , and the query objects are also different.

462 We use two pivot indices and two compact partition indices. For the pivot
 463 index family, we use the generic pivot algorithm and the Vantage Point Tree

464 index [38, 39].

465 In the case of the generic pivot algorithm, we choose at random a set of
466 pivots $\mathcal{P} = \{P_1, P_2, \dots, P_k\} \subset S$ of size $|\mathcal{P}| = k$. We store the kn distances
467 between pivots and objects, and use them to filter out candidates using the
468 triangle inequality. For each space, we experimentally determine the number
469 of pivots that obtains the best search performance. Thus, the results shown
470 for each case correspond to the best possible ones for this method, in the
471 corresponding metric space.

472 On the other hand, the Vantage Point tree (VPT) is a tree recursively
473 built by taking an arbitrary element p as the root. The distances from the
474 root to every object in the database are computed $\{d(p, u), u \in \mathbb{U}\}$. Let M
475 be the median of those distances. All objects such that $d(p, u) \leq M$ are
476 assigned to the left node and the rest to the right node. Then, we recurse
477 until the number of elements is smaller than a certain bucket size. To solve a
478 query in the VPT the query ball is tested to see if there could be candidates
479 in the left and right nodes. It is possible to enter both subtrees.

480 For the case of compact partition based algorithms, we consider the LC
481 [11], which is one of the best indexes for medium and high dimensions, and the
482 Spatial Approximation Tree [40]. We use the LC variant that has a maximum
483 size for each cluster. For each metric space considered, we experimentally
484 determine the cluster size whose performance is the best, and this is the result
485 shown in the plots.

486 Finally, the Spatial Approximation Tree (SAT) is a data structure aim-
487 ing at approaching the query spatially by starting at the root and getting
488 iteratively closer to the query by navigating the tree. The SAT is built as
489 follows. An element a is selected as the root, and is connected to a set of
490 *neighbors* $N(a)$, defined as a subset of elements x in the dataset such that x
491 is closer to a than to any other element in $N(a)$. The other elements (not in
492 $N(a) \cup \{a\}$) are assigned to their closest element in $N(a)$. Each element in
493 $N(a)$ is recursively the root of a new subtree containing the elements assigned
494 to it.

495 In Fig. 3, we show the cost of range queries retrieving 0.01%, 0.1% and
496 1% of the vector dataset per query, using the generic pivot index, the LC,
497 the VPT, and the SAT, see Figs. 3(a), (b), (c), and (d), respectively. These
498 results are compared with the ones for searching C10, C15, and C20. The four
499 plots show that the search effort performed by the four tested indices remain
500 almost unaltered when working on the four spaces of IDim 5 (independently
501 of the representational dimension of the space), while the curves for C10,

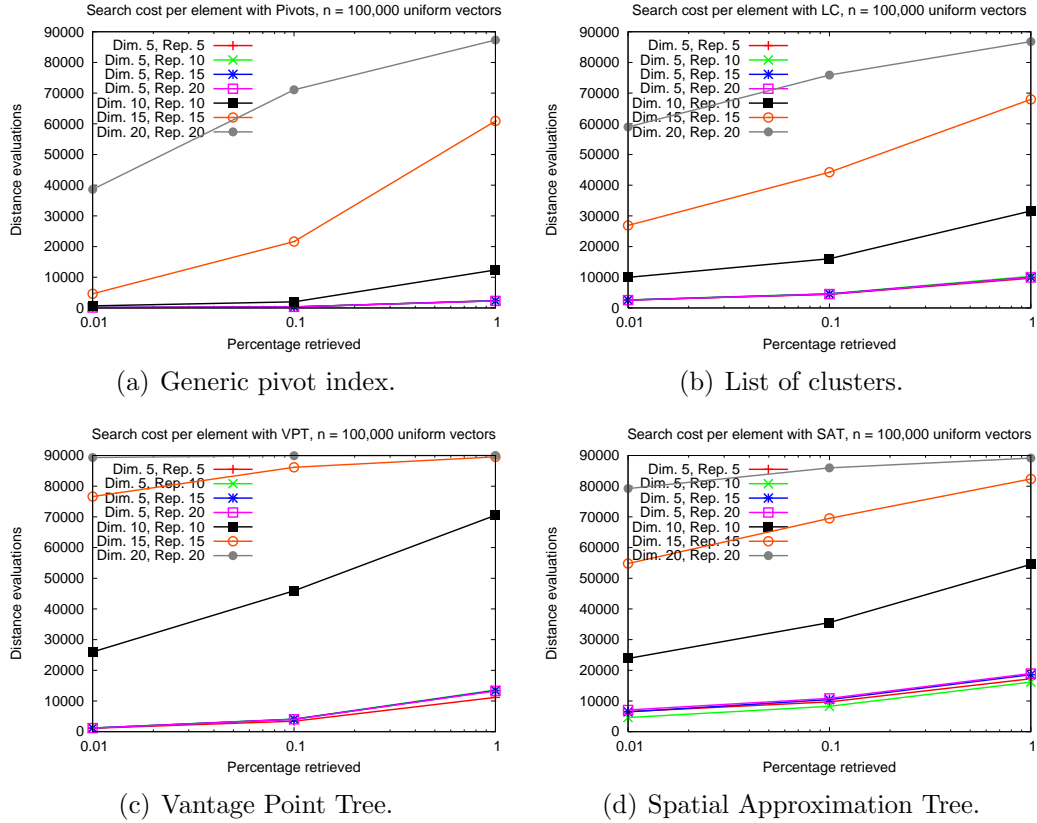


Figure 3: The search effort does not vary when the IDim of the space does not change.

502 C15, and C20 show the usual degradation.

503 4.1.2. Gaussian Distributed Vectors with Euclidean Distance

504 We generate 100,000 vectors in \mathbb{R}^D , where each coordinate has mean
 505 $\mu = 1$ and variance $\sigma^2 = 0.1$, for $D = 5, 10, 15,$ and 20 . In these spaces,
 506 there are no, a priori, clusters of elements. These spaces are called G5, G10,
 507 G15, and G20. Note that the object are not confined in the unitary cube.

508 We also generate 100,000 vectors in \mathbb{R}^{101} with 200 clusters. In this space,
 509 the first 100 coordinates of each vector follow a $\mathcal{N}(\mu = 1, \sigma^2 = 0.1)$ distri-
 510 bution. The 101-th coordinate stores the cluster identifier. So, the cluster
 511 centers are essentially uniformly distributed in the last coordinate. Geomet-
 512 rically speaking, one can imagine that this space is a sequence of 200 crisp
 513 clusters in a line immersed in \mathbb{R}^{101} . This space is called G101.

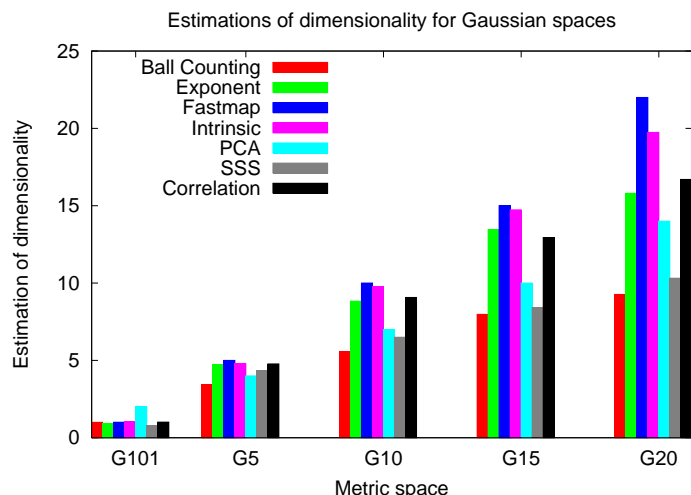


Figure 4: Comparison of dimensionality estimations for Gaussian spaces.

514 Fig. 4 shows the estimations obtained with the seven IDim estimators
 515 for these metric spaces. As can be seen, all the methods give increasing
 516 IDim values as the representational dimension grows (G5 through G20) as
 517 expected, but with different behaviors. Ball Counting and Distance Exponent
 518 become less sensitive to high dimensionality, from G15 to G20 they show a
 519 small increment. Intrinsic, SSS and Correlation give IDim values that grow
 520 steadily through G5 through G20, and we note that the Intrinsic Search
 521 Difficulty gives markedly lower values than in the uniform case. Fastmap and
 522 PCA show a large increment from G15 to G20, and we note that the IDim
 523 estimated by Fastmap in G20 is higher than in C20, which is unexpected.

524 We note that the lower the IDim, the more similar the dimensionality
 525 estimation. In fact, all the measures estimate the IDim of G5 around 4
 526 to 5 and the IDim of G101 around 1. This last result is very interesting.
 527 We prepare the dataset G101 with the purpose of having a space with high
 528 representational dimension but with very low intrinsic dimension and all the
 529 estimators detect this fact.

530 4.2. Real Metric Spaces

531 We pick four spaces from the Metric Library [25] ¹ in order to estimate
532 their IDims with the seven IDim estimators. The selected spaces are varied:

533 **Dictionary:** this is a dictionary of 69,069 English words. In this space, we
534 use a discrete function, the *Edit Distance* or *Levenshtein Distance* [27].

535 **NASA:** this is a set of 40,700 images from NASA, represented as feature
536 vectors of 20 real coordinates per vector, under the Euclidean distance.
537 They were generated from images downloaded from the NASA photo
538 and video archive site, used in contests conducted by the *Center for*
539 *Discrete Mathematics and Theoretical Computer Science* (DIMACS) ².
540 To obtain images from the videos, cuts are detected based on the tran-
541 sition of the color histogram and then representative images are ex-
542 tracted when the changes in the histogram reach a given threshold.
543 Later, the images are split into four sub-regions, i.e., upper-left, upper-
544 right, lower-left and lower-right, and histograms of the subregions are
545 calculated in order to take account of the composition of the image.
546 The four histograms are concatenated to compose a 36-dimensional
547 feature vector. Finally, using principal component analysis the feature
548 vectors are reduced to 20-dimensional vectors ³.

549 **Histograms:** this is a dataset of 112,682 histograms of medical images, each
550 one composed by 8-D color histograms of 112 real components ⁴. As
551 any quadratic form function can be used as the distance in this case, we
552 also have chosen the Euclidean distance, as is the simplest alternative.

553 **Documents:** this space has 1,265 documents, represented as vectors accord-
554 ing to the vector model of documents used in the Information Retrieval
555 field. To compare documents we use the *cosine distance*. Each vector
556 has a coordinate for each vocabulary term in the collection, and docu-
557 ments can be seen as points in a vector space of high representational

¹Available at <http://www.sisap.org/library/dbs/> .

²Available at <http://www.dimacs.rutgers.edu/Challenges/Sixth/software.html> .

³More details on this dataset can be obtained from
<http://www.dimacs.rutgers.edu/Challenges/Sixth/participants.html#KS> .

⁴Available at <http://www.dbs.informatik.uni-muenchen.de/~seidl/DATA/histo112.112682.gz> .

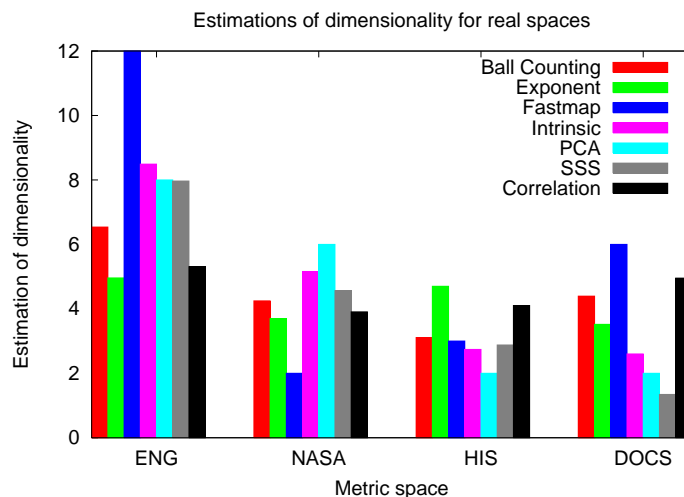


Figure 5: Comparison of dimensionality estimations for real metric spaces.

558 dimension. The documents are files obtained from the TREC-3 collec-
 559 tion ⁵.

560 We start the experimental evaluation in real metric spaces by estimating
 561 their IDims. These results are shown in Fig. 5. As it can be seen, all the
 562 methods shown coincide in that the English dictionary apparently has the
 563 highest IDim. Also, in these real world metric we can detect two groups of
 564 estimators that report similar values for IDims. The first is composed by
 565 Intrinsic Search Difficulty, PCA and SSS, and the second by Ball Counting,
 566 Distance Exponent and Correlation. On the other hand, Fastmap shows an
 567 erratic behavior.

568 To measure the intrinsic hardness of the searching, we consider the same
 569 four indices as before, using range queries:

570 **Dictionary:** As the metric is discrete, we use radii 1, 2, 3, and 4, retrieving
 571 on average about 0.003%, 0.037%, 0.326%, and 1.757% of the database.

572 **NASA:** In this continuous metric we use radii 0.012, 0.285, and 0.53, re-
 573 trievaling on average approximately 0.01%, 0.1%, and 1% of the dataset.

⁵Available at <http://trec.nist.gov>.

574 **Histograms:** This metric is also continuous. To retrieve on average ap-
575 proximately 0.01%, 0.1%, and 1% of the dataset, we use query radii
576 0.051768, 0.082514, and 0.131163.

577 **Documents:** The distance is also continuous. We use query radii 0.14,
578 0.15, and 0.195, which retrieve on average 0.01%, 0.1%, and 1% of the
579 database.

580 Fig. 6 shows the correlation between the search cost with the Pivot index
581 (on the left) and the List of Clusters (on the right), and the estimation
582 reported for each considered IDim estimator, namely, Ball counting, Distance
583 Exponent, Fastmap, Intrinsic Search Difficulty, PCA, SSS, and Correlation.

584 Fig. 7 illustrates the correlation between the search cost with the Vantage
585 Point Tree (on the left) and the Spatial Approximation Tree (on the right)
586 with respect to the seven estimators.

587 We plot the ratio between the logarithm of the search cost, measured
588 with distance computations, and the estimations of the IDim. This mea-
589 sures how close is the logarithm of the actual search costs to the predicted
590 IDim: if the search cost is consistently $s = c^d$, where d is the predicted IDim
591 and c is a constant, then the plots should always be close to $\log c$. Thus
592 the best methods are those that give roughly the same value regardless of
593 the index used. In Table 3, we show the mean and standard deviation of
594 $\log(\text{Search Difficulty})/\text{IDim}$ obtained for each estimator. To compute the ta-
595 ble, we consider the four real world metric spaces, the three radii, and the
596 four indices considered for each estimator.

597 As on the synthetic spaces, Distance Exponent and Correlation turn out
598 to be the best predictors for the four tested indices, as the ratio between
599 $\log(\text{Search Difficulty})$ and the IDim estimation remains similar in all the real
600 world spaces tested considering the different query radii and indices. This
601 can be corroborated in Table 3: Distance Exponent and Correlation have
602 the lowest standard deviations, just 0.135 and 0.200, respectively. As it can
603 be observed, the other methods are not stable enough with respect to search
604 costs. In fact, their standard deviations range from 0.289 for Ball Counting to
605 0.605 for PCA. This is because in some cases they underestimate the search
606 difficulty, and in the others they make an overestimation.

607 The search cost, however, depends both on the database size and on
608 the output size (which in turn depends on the query radius). Therefore,
609 dividing the cost of the search by these measures may give a more stable

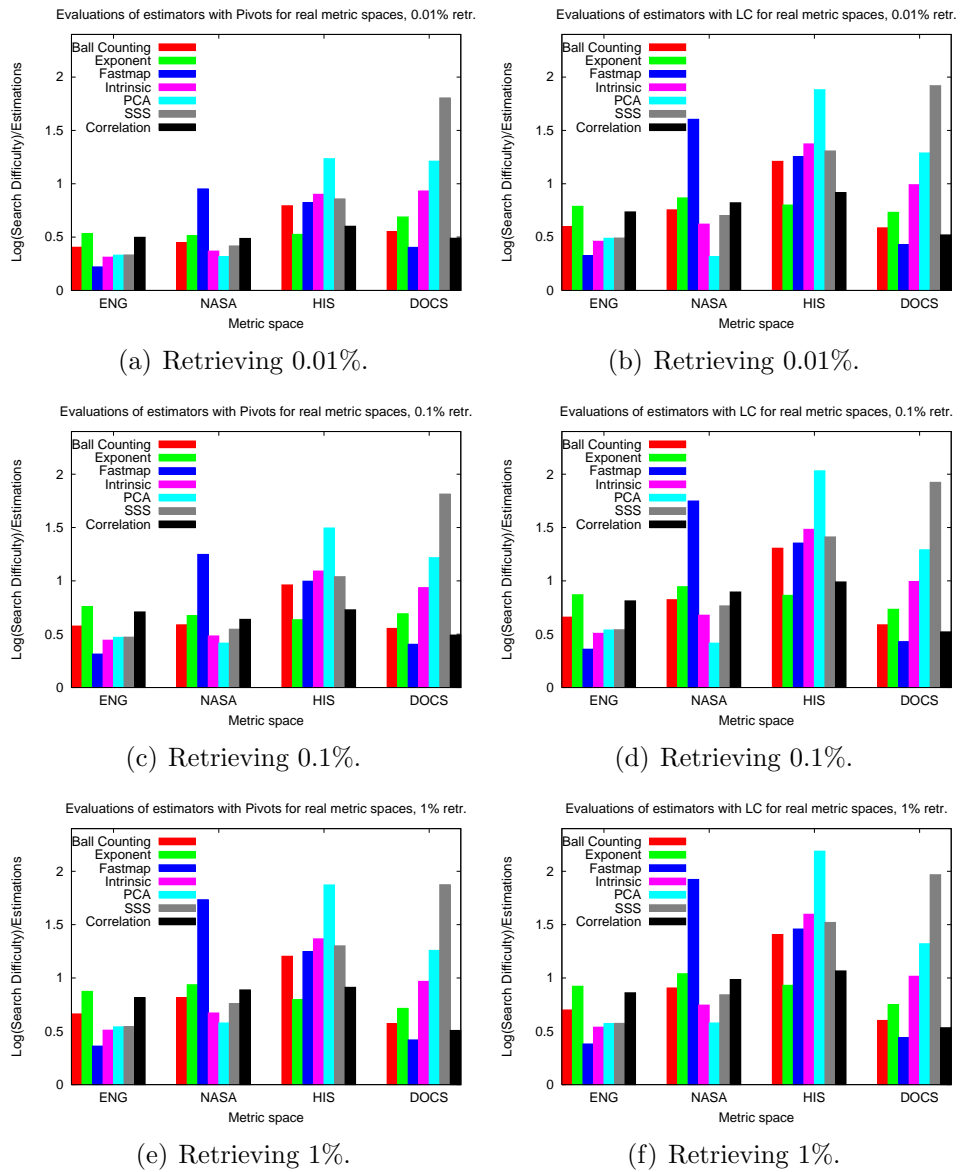


Figure 6: Comparison of Ball Counting, Distance Exponent, Fastmap, Intrinsic Search Difficulty, PCA, SSS, and Correlation IDim estimators, for real metric spaces. On the left, using Pivots. On the right, using List of Clusters.

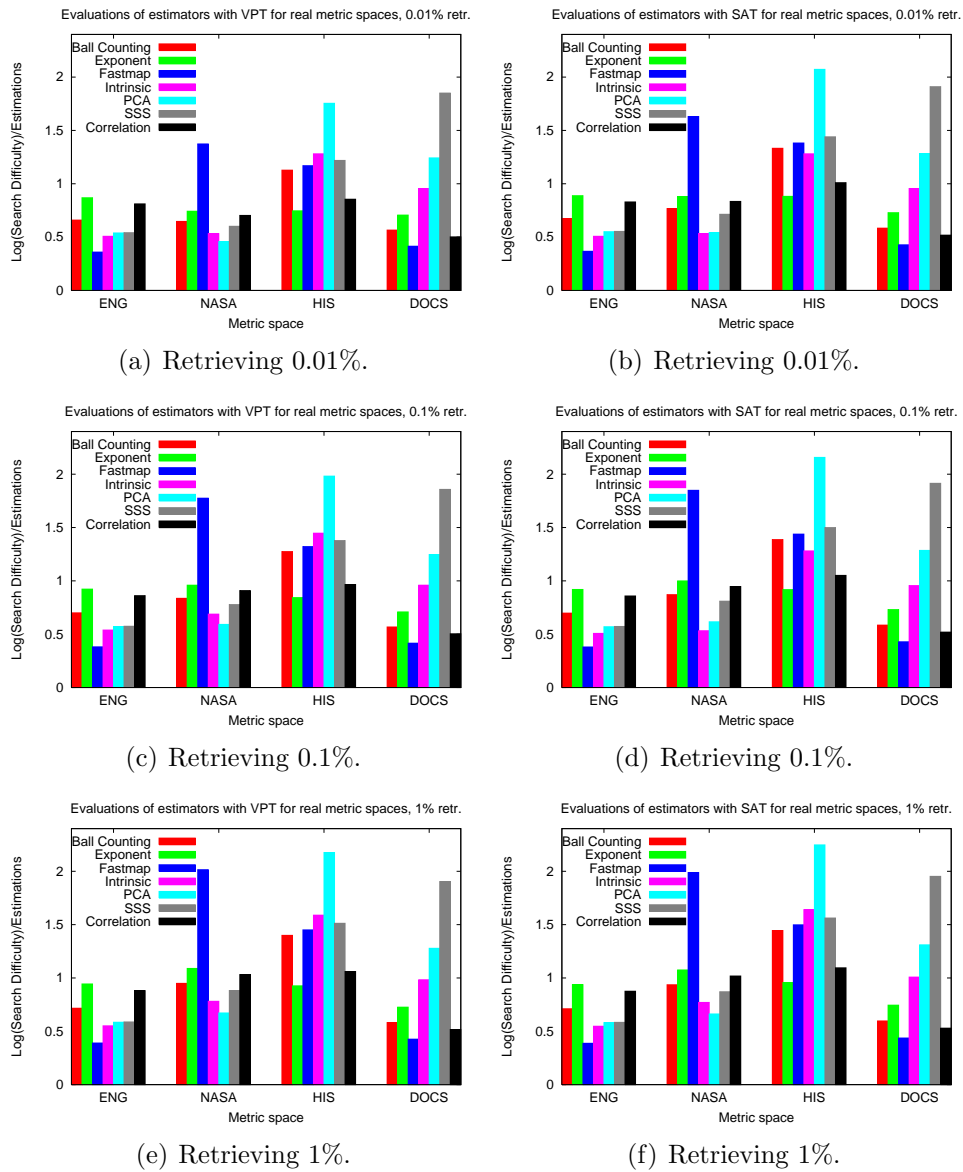


Figure 7: Comparison of Ball Counting, Distance Exponent, Fastmap, Intrinsic Search Difficulty, PCA, SSS, and Correlation IDim estimators, for real metric spaces. On the left, using Vantage Point Tree. On the right, using Spatial Approximation Tree.

Table 3: Comparison of $\log(\text{Search Difficulty})/\text{Estimation}$ for the seven IDim estimators.

IDim Estimator	Mean	Standard Deviation
Ball Counting	0.810	0.289
Distance Exponent	0.822	0.135
Fastmap	0.924	0.592
Intrinsic Search Difficulty	0.861	0.364
PCA	1.068	0.605
SSS	1.121	0.556
Correlation	0.773	0.200

Table 4: Comparison of $\log(\text{Search Difficulty}/\text{DBS})/\text{Estimation}$ for the seven IDim estimators.

IDim Estimator	Mean	Standard Deviation
Ball Counting	0.218	0.175
Distance Exponent	0.212	0.149
Fastmap	0.287	0.302
Intrinsic Search Difficulty	0.255	0.196
PCA	0.279	0.247
SSS	0.283	0.186
Correlation	0.204	0.156

610 measure of search difficulty. We repeat this evaluation considering the re-
611 lation between the fraction of the database visited when solving a query
612 and the estimation of IDim (this is, $-\log(\text{Search Difficulty}/\text{DBS})/\text{IDim}$,
613 where DBS stands for database size) and, on the other hand, the num-
614 ber of distance evaluation per each object in the query answer set (this is,
615 $\log(\text{Search Difficulty}/\text{QOS})/\text{IDim}$, where QOS stands for query output size).
616 We summarize these results in Tables 4 and 5. The reduced standard devi-
617 ations confirm that the estimations are indeed more stable, especially when
618 dividing by the database size. Still, we obtain the same conclusions: The two
619 most stable intrinsic dimensionality estimators are *Distance Exponent* and
620 *Correlation*.

Table 5: Comparison of $\log(\text{Search Difficulty}/\text{QOS})/\text{Estimation}$ for the seven IDim estimators.

IDim Estimator	Mean	Standard Deviation
Ball Counting	0.427	0.173
Distance Exponent	0.448	0.164
Fastmap	0.398	0.281
Intrinsic Search Difficulty	0.491	0.280
PCA	0.599	0.406
SSS	0.681	0.553
Correlation	0.407	0.137

621 5. Conclusions

622 The *Intrinsic Dimension* (IDim) of metric spaces measures their search
623 difficulty, independently of the type of index used. Computing the IDim is
624 useful to determine whether a metric space can be indexed at all (or we must
625 resort to sequential scanning or approximate methods), which kind of index
626 would perform better, and what search performance to expect.

627 We have analyzed seven IDim estimators in a practical perspective. Some
628 were defined for D -dimensional coordinate spaces, and we have adapted them
629 to the more general metric spaces. We compared their predictions with the
630 actual search cost using various synthetic and real-life metric spaces, so as
631 to verify which are better at predicting the search difficulty.

632 Although our results are preliminary, they suggest that all the methods
633 considered obtain appropriate estimations over synthetic metric spaces, be-
634 cause their values grow as the dimension increases. However, if we compare
635 the estimations with the real search costs, the Distance Exponent [12, 13]
636 and Correlation [14] turn out to be more stable. This is corroborated in
637 Table 3, that shows the mean and standard deviation of the ratio between
638 $\log(\text{Search Difficulty})$ and the IDim estimation for the seven estimators. The
639 standard deviations computed for Distance Exponent and Correlation are
640 just 0.135 and 0.200, respectively, and are the two lowest ones. On the other
641 hand, the other estimators, namely, Ball Counting (our adaptation of Box
642 Counting [21]), Fastmap [28], the simple measure proposed by Chávez et
643 al. [1], PCA, and SSS [37] sometimes obtain values less than the logarithm of
644 search costs and other times greater than them. These conclusions are also

645 supported by Tables 4 and 5.

646 As future work, we plan to analyze other estimators. For instance, we
647 can study the concentration dimension [35].

648 **Acknowledgements**

649 We gratefully acknowledge the anonymous referees who helped to improve
650 the presentation.

651 **References**

- 652 [1] E. Chávez, G. Navarro, R. Baeza-Yates, J. Marroquín, Searching in
653 metric spaces, *ACM Computing Surveys* 33 (3) (2001) 273–321.
- 654 [2] G. Hjaltason, H. Samet, Index-driven similarity search in metric spaces,
655 *ACM Trans. on Database Systems* 28 (4) (2003) 517–580.
- 656 [3] P. Zezula, G. Amato, V. Dohnal, M. Batko, *Similarity Search: The Met-*
657 *ric Space Approach*, Vol. 32 of *Advances in Database Systems*, Springer,
658 2006.
- 659 [4] H. Samet, *Foundations of Multidimensional and Metric Data Structures*
660 *(The Morgan Kaufmann Series in Computer Graphics and Geometric*
661 *Modeling)*, Morgan Kaufmann Publishers Inc., San Francisco, CA, USA,
662 2005.
- 663 [5] P. Ciaccia, M. Patella, P. Zezula, M-tree: an efficient access method for
664 similarity search in metric spaces, in: *Proc. 23rd Conf. on Very Large*
665 *Databases (VLDB)*, 1997, pp. 426–435.
- 666 [6] V. Dohnal, C. Gennaro, P. Savino, P. Zezula, D-index: Distance search-
667 ing index for metric data sets, *Multimedia Tools and Applications* 21 (1)
668 (2003) 9–33.
- 669 [7] T. Skopal, J. Pokorný, V. Snásel, PM-tree: Pivoting metric tree for sim-
670 ilarity search in multimedia databases, in: *ADBIS (Local Proceedings)*,
671 2004.

- 672 [8] G. Navarro, N. Reyes, Dynamic spatial approximation trees for massive
673 data, in: T. Skopal, P. Zezula (Eds.), Proc. 2nd Intl. Workshop on
674 Similarity Search and Applications (SISAP), IEEE CS Press, 2009, pp.
675 81–88.
- 676 [9] G. Navarro, R. Uribe, Fully dynamic metric access methods based on
677 hyperplane partitioning, *Information Systems* 36 (4) (2011) 734–747.
- 678 [10] G. Navarro, N. Reyes, Dynamic list of clusters in secondary memory,
679 in: Proc. 7th Intl. Workshop on Similarity Search and Applications
680 (SISAP), LNCS 8821, 2014, pp. 94–105.
- 681 [11] E. Chávez, G. Navarro, A compact space decomposition for effective
682 metric indexing, *Pattern Recognition Letters* 26 (9) (2005) 1363–1376.
- 683 [12] C. Traina Jr., A. J. M. Traina, C. Faloutsos, Distance exponent: a new
684 concept for selectivity estimation in metric trees, Research Paper 99-
685 110, School of Computer Science, Carnegie Mellon University (03/1999
686 1999).
- 687 [13] C. Traina Jr., A. J. M. Traina, C. Faloutsos, Distance exponent: A
688 new concept for selectivity estimation in metric trees., in: Proc. 16th
689 Intl. Conf. on Data Engineering (ICDE), 2000, p. 195.
- 690 [14] F. Camastra, A. Vinciarelli, Estimating the intrinsic dimension of data
691 with a fractal-based method., *IEEE TPAMI* 24 (10) (2002) 1404–1407.
- 692 [15] C. Bustos, G. Navarro, N. Reyes, R. Paredes, An empirical evaluation
693 of intrinsic dimension estimators, in: Proc. 8th Intl. Conf. on Similarity
694 Search and Applications (SISAP), LNCS 9371, Springer, 2015, pp. 125–
695 137.
- 696 [16] A. K. Jain, R. C. Dubes, Algorithms for clustering data, Prentice-Hall,
697 Inc., Upper Saddle River, NJ, USA, 1988.
- 698 [17] F. Camastra, Data dimensionality estimation methods: a survey, *Pat-
699 tern Recognition* 36 (12) (2003) 2945–2954.
- 700 [18] K. Fukunaga, Introduction to statistical pattern recognition (2nd ed.),
701 Academic Press Professional, Inc., San Diego, CA, USA, 1990.

- 702 [19] R. Bellman, Adaptive control processes: a guided tour, Princeton Uni-
703 versity Press Princeton, N.J, 1961.
- 704 [20] V. N. Vapnik, The nature of statistical learning theory, Springer-Verlag
705 New York, Inc., New York, NY, USA, 1995.
- 706 [21] B. Mandelbrot, Fractals: Form, Chance and Dimension, W. H. Freeman,
707 San Francisco, 1977.
- 708 [22] J. P. Eckmann, D. Ruelle, Ergodic theory of chaos and strange attrac-
709 tors, Rev. Mod. Phys. 57 (1985) 617.
- 710 [23] E. Ott, Chaos in dynamical systems, Cambridge University Press, Cam-
711 bridge, New York, 1993.
- 712 [24] D. Kaplan, L. Glass, Understanding Nonlinear Dynamics, Springer-
713 Verlag, New York, 1995.
- 714 [25] K. Figueroa, G. Navarro, E. Chávez, Metric spaces library, available at
715 http://www.sisap.org/Metric_Space_Library.html (2007).
- 716 [26] H. V. Jagadish, A retrieval technique for similar shapes, in: SIGMOD
717 Conference, ACM Press, 1991, pp. 208–217.
- 718 [27] V. I. Levenshtein, Binary codes capable of correcting deletions, inser-
719 tions, and reversals, Soviet Physics Doklady 10 (8) (1966) 707–710.
- 720 [28] C. Faloutsos, K.-I. Lin, Fastmap: A fast algorithm for indexing, data-
721 mining and visualization of traditional and multimedia datasets, in:
722 Proc. 1995 ACM SIGMOD Intl. Conf. on Management of Data, ACM
723 Press, 1995, pp. 163–174.
- 724 [29] I. T. Jolliffe, Principal Component Analysis, 2nd Edition, Springer Series
725 in Statistics, Springer, 2002.
- 726 [30] R Core Team, R: A Language and Environment for Statistical Comput-
727 ing, R Foundation for Statistical Computing, Vienna, Austria (2013).
- 728 [31] S. Brin, Near neighbor search in large metric spaces, in: Proc. 21st
729 Conf. on Very Large Databases (VLDB’95), 1995, pp. 574–584.

- 730 [32] E. Chávez, J. Marroquín, Proximity queries in metric spaces, in: Proc.
731 4th South American Workshop on String Processing (WSP'97), Carleton
732 University Press, 1997, pp. 21–36.
- 733 [33] P. Ciaccia, M. Patella, P. Zezula, A cost model for similarity queries
734 in metric spaces., in: Proc. 17th ACM SIGACT-SIGMOD-SIGART
735 Symp. on Principles of Database Systems (PODS), 1998, pp. 59–68.
- 736 [34] P. Yianilos, Excluded middle vantage point forests for nearest neighbor
737 search, Tech. rep., NEC Research Institute, Baltimore, MD, in *6th DI-*
738 *MACS Implementation Challenge: Near Neighbor Searches Workshop,*
739 *ALLENEX'99* (1998).
- 740 [35] V. Pestov, Intrinsic dimension of a dataset: what properties does one
741 expect?, in: 2007 Intl. Joint Conf. on Neural Networks (IJCNN), 2007,
742 pp. 2959–2964. doi:10.1109/IJCNN.2007.4371431.
- 743 [36] V. Pestov, An axiomatic approach to intrinsic dimension of a dataset,
744 Neural Networks 21 (2–3) (2008) 204–213, advances in Neural Net-
745 works Research: 2007 Intl. Joint Conf. on Neural Networks (IJCNN).
746 doi:http://dx.doi.org/10.1016/j.neunet.2007.12.030.
- 747 [37] N. R. Brisaboa, A. Fariña, O. Pedreira, N. Reyes, Similarity search using
748 sparse pivots for efficient multimedia information retrieval, in: 8th IEEE
749 Intl. Symp. on Multimedia (ISM), IEEE CS, 2006, pp. 881–888.
- 750 [38] P. Yianilos, Data structures and algorithms for nearest neighbor search
751 in general metric spaces, in: Proc. 4th ACM-SIAM Symposium on Dis-
752 crete Algorithms (SODA'93), SIAM Press, 1993, pp. 311–321.
- 753 [39] T. Chiueh, Content-based image indexing, in: Proc. 20th Conf. on Very
754 Large Databases (VLDB'94), 1994, pp. 582–593.
- 755 [40] G. Navarro, Searching in metric spaces by spatial approximation, The
756 Very Large Databases Journal (VLDBJ) 11 (1) (2002) 28–46.