

1 A log square average case algorithm to make insertions
2 in fast similarity search

3 Luisa Micó*, Jose Oncina*

4 *Dept. Lenguajes y Sistemas Informáticos*
5 *Universidad de Alicante, E-03071 Alicante, Spain*

6 **Abstract**

7 To speed up similarity based searches many indexing techniques have been
8 proposed in order to address the problem of efficiency. However, most of the
9 proposed techniques do not admit fast insertion of new elements once the index
10 is built. The main effect is that changes in the environment are very costly to
11 be taken into account.

12 In this work, we propose a new technique to allow fast insertions of elements
13 in a family of static tree-based indexes. Unlike other techniques, the resulting
14 index is exactly equal to the index that would be obtained by building it from
15 scratch. Therefore there is no performance degradation in search time.

16 We show that the expected number of distance computations (and the aver-
17 age time complexity) is bounded by a function that grows with $\log^2(n)$ where
18 n is the size of the database.

19 In order to check the correctness of our approach some experiments with
20 artificial and real data are carried out.

21 *Keywords:* similarity search, metric space, dynamic index, insertions

22 **1. Introduction**

23 The similarity search problem can be stated as follows: given a finite data
24 set of objects D , a dissimilarity measure d and a query object q find the set of
25 elements in the data set ($P \subset D$) that is the most similar to the query (minimise
26 a dissimilarity measure). Depending on the amount and type of information

*Corresponding Author

Email addresses: micó@dlsi.ua.es (Luisa Micó), oncina@dlsi.ua.es (Jose Oncina)
Preprint submitted to *Discrete Applied Mathematics* (November 9, 2011)

27 required, several similarity search techniques can be stated: nearest neighbour
 28 search (only the nearest object in the database is retrieved: $p \in P \iff \forall r \in$
 29 $D, d(q, p) \leq d(q, r)$), range search (all the objects in the database nearest to
 30 the query than a value h are retrieved: $p \in P \iff d(q, p) \leq h$), reverse
 31 nearest neighbour search (the elements in the dataset that have the query as
 32 their nearest element: $p \in P \iff \forall r \in D, d(p, q) \leq d(p, r)$), etc.

33 Over the time, these techniques have been applied to databases increasingly
 34 large making their execution times become real bottlenecks.

35 In order to speed up these techniques, fast similarity search methods have
 36 to exploit some property of the search space. Metric space searching techniques
 37 assume that the dissimilarity function ($d(\cdot, \cdot)$) defines a metric over the repre-
 38 sentation space E , that is:

- | | | |
|----|--|---------------------|
| 39 | 1. $\forall x, y \in E, d(x, y) \geq 0$ | non-negativity |
| 40 | 2. $\forall x, y \in E, d(x, y) = d(y, x)$ | symmetry |
| 41 | 3. $\forall x \in E, d(x, x) = 0$ | identity |
| 42 | 4. $\forall x, y, z \in E, d(x, z) \leq d(x, y) + d(y, z)$ | triangle inequality |

43 One of the main characteristics of metric space searching is that no assump-
 44 tion about the structure of the objects (points) is necessary. Some examples
 45 of objects can be: protein sequences (represented by strings) (Lundsteen et al.,
 46 1980), skeleton of images (trees or graphs)(Carrasco and Forcada, 1995)(Es-
 47 colano and Vento, 2007), histograms of images(Cha and Srihari, 2002), etc.

48 At present, many communities have paid great attention to these techniques
 49 because of the need for handling large amounts of data. Then, many metric
 50 space indexes designed to speed up searches have been proposed (some reviews
 51 can be found in (Chávez et al., 2001)(Hjaltason and Samet, 2003)(Zezula et al.,
 52 2006)). These indexes have proved to be very effective in many applications such
 53 as content based image retrieval (Giacinto, 2007), person detection or automatic
 54 image annotation (Torralba et al., 2008), texture synthesis, image colourisation
 55 or super-resolution (Battiato et al., 2007).

56 Unfortunately, most of these indexes are static (Yianilos, 1993)(Brin, 1995)(Micó

57 et al., 1994)(Navarro, 2002). That is, the insertion or deletion of an object re-
58 quires a complete rebuilding of the index. This is very expensive and discourages
59 its use in interactive or on-line training systems.

60 In this work, we propose technique to allow fast insertions. The performance,
61 in search time, of the index does no degrade with the insertions and we show that
62 the expected number of distance computations is bounded by $\log^2(n)$ where n is
63 the size of the database. This result compares very favourably with the number
64 of distance computations needed in a whole rebuild ($n \log(n)$).

65 In order to check the correctness of our approach some experiments with
66 artificial data (Euclidean distance in 5, 10 and 15 dimensional spaces) and real
67 data (Euclidean distance in an image database and edit distance in handwritten
68 digits contour strings and English words) have been carried out.

69 Section 2 describes related work and introduces the main ideas in our ap-
70 proach. Section 3 introduces the static index in wich our approach is based, and
71 Section 4 describes our inserting algorithm. Section 5 is devoted to analyse the
72 insertion cost. This analysis is followed by experimental results using artificial
73 and real data in Section 6. Finally Section 7 describes the conclusions drawn
74 from the results and summarises our contribution.

75 **2. The approach**

76 A number of proposals to allow object insertion/deletion operations have
77 been made for metric space indexes (Fu et al., 2000)(Navarro and Reyes, 2008).
78 In some cases dynamic approaches were proposed as a completely new algorithm
79 to allow cheap insertions and deletions such as the M -tree (Ciaccia et al., 1997),
80 and, in other cases, as a modification of previously existing static indexes (Fu
81 et al., 2000)(Navarro and Reyes, 2002)(Procopiuc et al., 2003)

82 Usually, static search methods are faster searching that dynamic indexes and
83 static methods degrades when they are adapted to allow insertions.

84 The main problem when adapting static methods to allow insertions is the
85 need of a reorganization when an insertion is performed. To avoid this overhead,

86 some authors (Navarro and Reyes, 2008) propose the use of buckets in selected
87 places of the index to store the new objects in such a way they can be located
88 easily and does not harm very much the performance of the index. Despite of
89 that, the nearest neighbour search performance is degraded as the size of the
90 buckets increases. To avoid such degradation a rebuilding of the index is forced
91 when the size of a bucket exceeds a threshold. A trade-off between insertion
92 performance and search performance should be established.

93 In our proposal the index obtained after the insertion is the same as the
94 (static) one obtained if a complete rebuild would be made, without adding
95 buckets or any type of additional information to the index. As a consequence,
96 no insertion/search performance ratio should be adjusted and there is no degra-
97 dation of search performance.

98 The idea of the strategy is quite simple: go ahead with the insertion unless a
99 modification in the index is necessary; otherwise, rebuild completely the affected
100 part of the index.

101 Although this strategy can be applied to many indexing techniques, it is
102 specially effective when is applied to Most Distant to the Father (MDF) tree
103 index. This tree based indexing is used in some state of the art searching
104 techniques (Micó et al., 1996)(Gómez-Ballester et al., 2006).

105 The properties that make this structure so effective are:

- 106 1. the structure is based on the use of objects in very low probability regions
- 107 2. the rebuilding of the index section corresponding to one branch of the tree
108 is independent of the other branches.

109 **3. The Most Distant to the Father tree index**

110 One of the most successful methods for reducing the search time (by reducing
111 the average number of distance computations) is the mb-tree (or monotonous
112 bisector tree). This method was originally intended to be used with vector-
113 data and Minkowski metrics although it can be used with arbitrary metrics
114 and then, with complex objects. The mb-tree was proposed by Noltemeier et

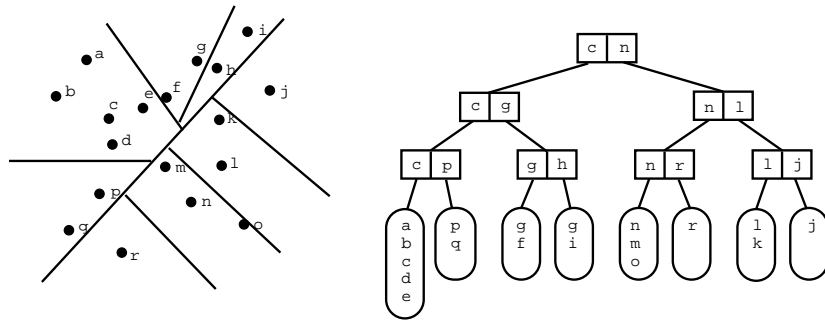


Figure 1: Example of space partitioning produced by a mb-tree in a two-dimensional space (left) and the mb-tree derived from it (right).

115 al. (Noltemeier et al., 1992) to modify the definition of the bisector tree (a
 116 tree that uses generalized hyperplane partitioning augmented by including for
 117 each pivot the maximum distance to an object in its subtree, (Kalantari and
 118 McDonald, 1983)) so that one of the two pivots in each nonleaf node is inherited
 119 from its parent (see Figure 3). This strategy allows to reduce the number of
 120 distance computations during the search (only a new distance, instead of two,
 121 is necessary to compute every time a new level is explored in the tree). But
 122 this is at the cost of a worse partitioning, obtaining a deeper tree. This general
 123 approach allows many different configurations in the selection of pivots.

124 The MDF tree is a binary indexing structure based on a hyperplane parti-
 125 tioning approach (Micó et al., 1996)(Gómez-Ballester et al., 2006) with similar
 126 properties to the mb-trees. The main difference is related to the selection of the
 127 representatives (pivots) for the next partition (branch of the tree).

128 In the MDF-tree firstly a pivot is randomly selected as the root of the tree
 129 (first level). Secondly, the most distant point from the root is selected as a new
 130 pivot, and the remaining points are distributed according to to the closest pivot.
 131 This procedure is recursively repeated until each leaf node has only one object
 132 (see Figure 3).

133 For each node, the covering radius (the distance from the pivot to the most
 134 distant point in the subspace) is computed and stored in the respective node.

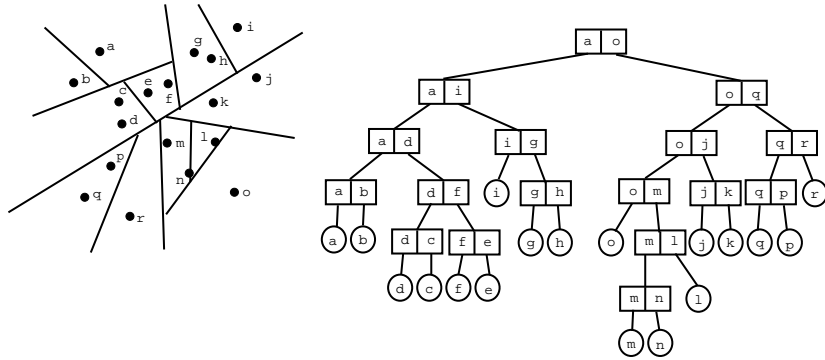


Figure 2: Example of space partitioning produced by a MDF-tree in a two-dimensional space(left) and the MDF-tree derive from it (right)

135 This procedure is described in algorithm 1.

136 The function `build_tree(ℓ, S)` takes as arguments the future representative
 137 of the root node (ℓ) and the set of objects to be included in the tree (exclud-
 138 ing ℓ) and returns the MDF-tree that contains $S \cup \{\ell\}$. The first time that
 139 `build_tree(ℓ, S)` is called, ℓ is selected randomly among the data set. In the
 140 algorithm, M_T is the pivot corresponding to T , r_T is the covering radius, and
 141 T_L (T_R) is the left (right) subtree of T .

142 It is easy to see that the space complexity of the index is $O(n)$, with n being
 143 the number of points, and the time complexity is $O(hn)$ where h is the depth
 144 of the tree.

145 4. Incremental tree

146 In this work we focus on a procedure to obtain, each time an insertion is
 147 performed, the exact index that will be obtained if a complete rebuild of the tree
 148 was made (a preliminary version of this idea, with no theoretical guarantees, was
 149 presented in Micó and Oncina (2009)). Note that in such case the performance
 150 of the search algorithm that uses the MDF index is exactly the same as in
 151 the case when the index is build from scratch. Then no further research or
 152 experiments in search degradation performance is needed.

Algorithm 1: build_tree(ℓ, S)

Data:

$S \cup \{\ell\}$: set of points to include in T ;

ℓ : future left representative of T

create MDF-tree T

if S is empty **then**

$M_T = \ell$

 // New representative of T

$r_T = 0$

else

$r = \operatorname{argmax}_{x \in S} d(\ell, x)$

$r_T = d(\ell, r)$

$S_\ell = \{x \in S \mid d(\ell, x) < d(r, x)\}$

$S_r = \{x \in S \mid d(\ell, x) \geq d(r, x)\} - \{r\}$

$T_L = \text{build_tree}(\ell, S_\ell)$

$T_R = \text{build_tree}(r, S_r)$

end

return T

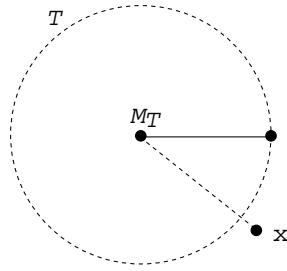


Figure 3: Case when a complete rebuild is needed

153 The main idea consists on comparing the object to be inserted with the pivot
 154 (representatives) on every recursive call to check if it is farther than the farthest
 155 so far. If it is farther (it has to be a pivot in the resulting tree), the affected
 156 part of the tree is completely rebuilt. Otherwise, the insertion is made in the
 157 subtree whose pivot is nearest.

158 This may seem a quite expensive strategy, but as the pivots are very unusual
 159 objects (the farthest of its sibling pivot), and the sizes of the subtrees decrease
 160 very quickly, big reconstructions of the tree seldom happens. The high cost of
 161 big rebuilds is compensated by its low probability.

162 Let T be the MDF tree built using a database D . Let x be the new object
 163 to be inserted in the index, and let T' the MDF tree built using the data set
 164 $D \cup \{x\}$. The algorithm detects and rebuilds the subtree of T that is different
 165 from T' .

166 Let us denote by M_T the representative of the root node of a subtree T of
 167 the MDF tree, let r_T be its covering radius, and let T_L (T_R) be the left (right)
 168 MDF subtree of T .

169 We have several cases:

170 **C 1.** If $d(M_T, x) > r_T$, T' differs from T in the root node because the object x
 171 is selected in T' as the representative of the right node. Then the whole
 172 tree T is rebuilt in order to include x (see fig. 3).

173 **C 2.** Otherwise, the roots of the trees T and T' are identical. Then we have
 174 two cases (see fig. 4) :

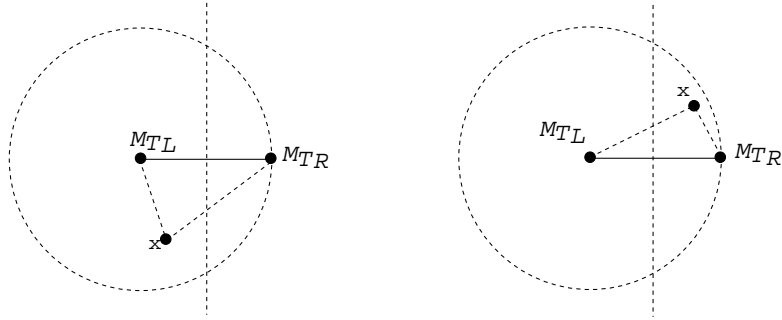


Figure 4: Inserting in the left (right) subtree

175 **C 2.1.** if $d(M_{T_L}, x) < d(M_{T_R}, x)$ the object x should be inserted in the
 176 left tree T_L and then the tree T'_R is identical to T_R .

177 **C 2.2.** Conversely if $d(M_{T_L}, x) \geq d(M_{T_R}, x)$ the object should be inserted
 178 in T_R and the tree T'_L is identical to T_L .

179 Algorithm 2 shows the insertion procedure.

180 5. Average time complexity

181 An MDF tree is generally unbalanced and, in the worst case, it can be fully
 182 degenerated.

183 We introduce a parameter α to measure the inbalance of a tree. Let $\alpha \in$
 184 $[0.5, 1.0[$ be defined such that for all node T in a MDF tree, where T_1 and T_2
 185 are its two children and where $|T_1| \leq |T_2|$

$$|T_1| \leq \alpha |T|$$

$$|T_2| \geq (1 - \alpha) |T|$$

186 An upper bound to the depth (h) of a α -unbalanced tree can be easily computed
 187 taking into account that, in the worst case, the size of the bigger child decreases
 188 at least a factor α in each level until we arrive to a leaf (size 1).

Algorithm 2: insert_tree(T, x)

Data:

T : MDF-tree

x : object to be inserted in T

if $d(M_T, x) > r_T$ **then**

| $T = \text{build_tree}(M_T, \{s | s \in T\} \cup \{x\} - \{M_T\})$

else if T_L is empty **then**

| $T = \text{build_tree}(M_T, \{x\})$

else

| $d_\ell = d(M_{T_L}, x)$ // this distance has been previously computed

| $d_r = d(M_{T_R}, x)$

| **if** $d_\ell < d_r$ **then**

| | insert_tree(T_L, x)

| **else**

| | insert_tree(T_R, x)

| **end**

end

189 That is, $\alpha^h n = 1$, and then $h = -\frac{\log(n)}{\log(\alpha)}$. For example, if the tree is balanced
 190 ($\alpha = 0.5$) $h = \log_2(n)$.

191 Then, the number of distance computations required to build a α -unbalanced
 192 MDF tree of size n is upper bounded by $nh = n(-\frac{\log(n)}{\log(\alpha)})$.

193 In the following, we are going to obtain an upper bound of this function.
 194 Let we denote by $E(n)$ the expected number of distance computations when
 195 inserting an object in a α -unbalanced MDF tree of size n .

196 Let x be the object to be inserted and assume that all the elements in $D \cup \{x\}$
 197 where extracted i.i.d. from an unknown probability distribution. Following
 198 alg. 2 we have four possibilities:

- 199 1. if $n = 1$ then $E(n) = 1$
- 200 2. if $d(M_T, x) > r_T$ a rebuilt of the subtree is necessary. Its cost is upper
 201 bounded by $n(-\frac{\log(n)}{\log(\alpha)})$
- 202 3. if $d(M_{T_L}, x) < d(M_{T_R}, x)$, x is inserted in the left subtree
- 203 4. if $d(M_{T_L}, x) \geq d(M_{T_R}, x)$, x is inserted in the right subtree

204 Note that since we are assuming that all the points are extracted i.i.d., all
 205 the points have the identical probability of being the new pivot (fulfilling con-
 206 dition 2) and then, the probability of this event is $\frac{1}{n}$. Therefore, the probability
 207 of event 3 or event 4 is $\frac{n-1}{n}$.

208 Moreover, the action taken by the algorithm in such cases is to make an
 209 insertion in one of its children. Since the tree is α -unbalanced the cost of each
 210 of this actions are bounded by worst case: $E(\alpha n)$.

Now, expressing all that in an equation we have the upper bound:

$$E(n) \leq 1 + \frac{1}{n} n \left(-\frac{\log(n)}{\log(\alpha)} \right) + \frac{n-1}{n} E(\alpha n) \quad (1)$$

211 This equation is composed by three terms. The first term takes into account
 212 the distance computation needed to know the distance from the sample to the
 213 pivot. Second an third terms takes into account the possibility of the new
 214 sample being farther (or not) than the present representative in the right node.
 215 Second term is the probability that the new sample is farther than the present

216 representative ($\frac{1}{n}$), multiplied by the cost of rebuilding the subtree. Third
 217 term is the probability that the new sample is not the farther representative,
 218 multiplied by the expected number of distance computations of inserting in the
 219 bigger of the two children (size at most αn).

220 If we unfold equation 1 we have:

$$E(n) \leq h - \frac{1}{\log(\alpha)} (\log(n) + \log(n\alpha) + \log(n\alpha^2) + \dots \textit{h times})$$

221 where we have taken into account that $\frac{n-i}{n} < 1, \forall i < n$.

$$E(n) \leq -\frac{\log(n)}{\log(\alpha)} - \frac{1}{\log(\alpha)} \log \left(\prod_{i=0}^h n\alpha^i \right)$$

222 and using some properties of the log function:

$$E(n) \leq \frac{\log^2(n)}{2\log^2(\alpha)} - \frac{3\log(n)}{2\log(\alpha)} \quad (2)$$

223 This upper bound shows that, in the worst case, the expected number of
 224 distance computations grows with $\log^2(n)$. Very far of the worst case ($n \log(n)$).

225 6. Experimental results

226 The experiments were done using artificial and real data represented as vec-
 227 tors or strings. For artificial data, the datasets were generated using a uniform
 228 distribution in the 5, 10 and 15 dimensional unit hypercube.

229 Three real data databases are used:

230 **NASA:** is a collection of 40 150, 20-dimensional vectors obtained from NASA
 231 video and image archives. The authors of the database (Katayama and
 232 Satoh, 1999) divided each image in four regions, nine color histograms
 233 for each region were computed. The features are a PCA projection to a
 234 20 dimensional space. The Euclidean distance was used as dissimilarity
 235 measure (more details can be found in <http://www.sisap.com>).

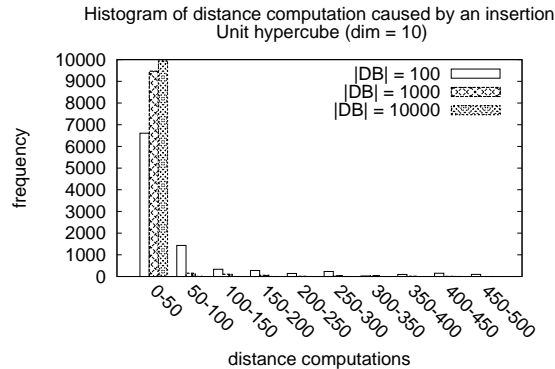


Figure 5: Histogram of the distance computations caused by an insertion for 100, 1000, and 10 000 artificial database sizes in dimension 10. Only the first 10 bins are showed

236 **English:** is an English dictionary of 69 069 words extracted from the dictionary
 237 of the GNU spell checker (<ftp://ftp.gnu.org/gnu/aspell/dict/0index.html>).
 238 In this case the edit distance was used as dissimilarity measure.

239 **Contour:** is a set of 10 000 8-directions contour strings extracted from the
 240 NIST Special Database 3 (Garris and Wilkinson, 1994). NIST database
 241 contains 128×128 black and white (bilevel) images of handwritten digits
 242 that was collected among Census Bureau employees.

243 First, in order to study the distribution of the number of distance computa-
 244 tions needed to rebuild the index when an object is inserted, 10 000 insertions of
 245 an object over a fixed MDF-tree with sizes 100, 1000, and 10 000 was made. The
 246 number of distance computations were counted and its histogram is depicted in
 247 Figure 5. Only the case for the uniform distribution in a 10 dimensional space
 248 is plotted, the other cases show a similar behaviour.

249 It can be observed that independently of the size of the tree, almost all the
 250 insertions cause very few distance computations (left side). On the other hand,
 251 there are very few insertions that cause a large number of distance computa-
 252 tions (right side). In all the experiments, rarely more than a hundred distance
 253 computations have been made in one insertion.

254 Next experiments are intended to study the behaviour of the algorithm when
255 inserting objects in a database. For that, 10 000 different databases were gen-
256 erated for each size varying from 100 to 10 000 in steps of 100. Each point in
257 the plots is the average number of the distance computations provoked by an
258 insertion in each of the 10 000 MDF indexes build from these databases. The
259 95% percentile was also computed and plotted in the following experiments.
260 That means that the 95% of the cases make a number of distance computations
261 under this curve.

262 Moreover, the theoretical upper bound was plotted to check experimentally
263 its validity. In order to do that, for each MDF tree, the values of the unbalance
264 factor α was computed for each node of the tree. In order to meet the definition,
265 the α for a tree should be the maximum α of its nodes. Instead of that, the
266 95% percentile of the node α 's was computed to avoid pathological high values
267 of α . Note that doing so the predicted values for the upper bound are going to
268 be lower than if the maximum would be computed. The figures also show the
269 value of the α factor for the corresponding experiment. The results are showed
270 in Figure 6 for the artificial data and Figure 7 for the real data experiments.

271 The experimental results fits very well with the theoretical prediction. It can
272 be seen that in all the cases the distance computations caused by an insertion
273 seems to grow very slowly with the database size. Moreover, the 95% percentile
274 decreases as the database size increases. This effect is due to the fact that the
275 cost of the worst case increases much faster that the average case. Then, in
276 order to compensate for the few worst case events, many events have to be very
277 cheap.

278 7. Conclusions

279 In this work we have proposed a simple but efficient algorithm to insert
280 objects in a MDF-tree. This algorithm, unlike others, has the property that the
281 index obtained after the insertion is the same as the one obtained if a complete
282 rebuild would be made. Then, the search efficiency does not degrade when

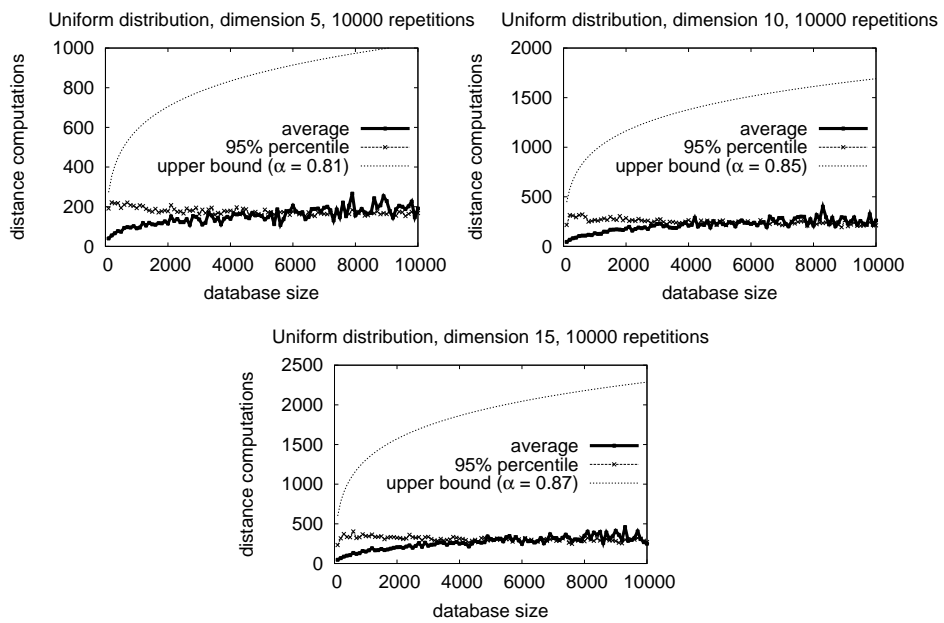


Figure 6: Experimental average number and theoretical upper bound of the distance computations caused by an insertion for increasing size database sets and for dimensions 5, 10 and 15.

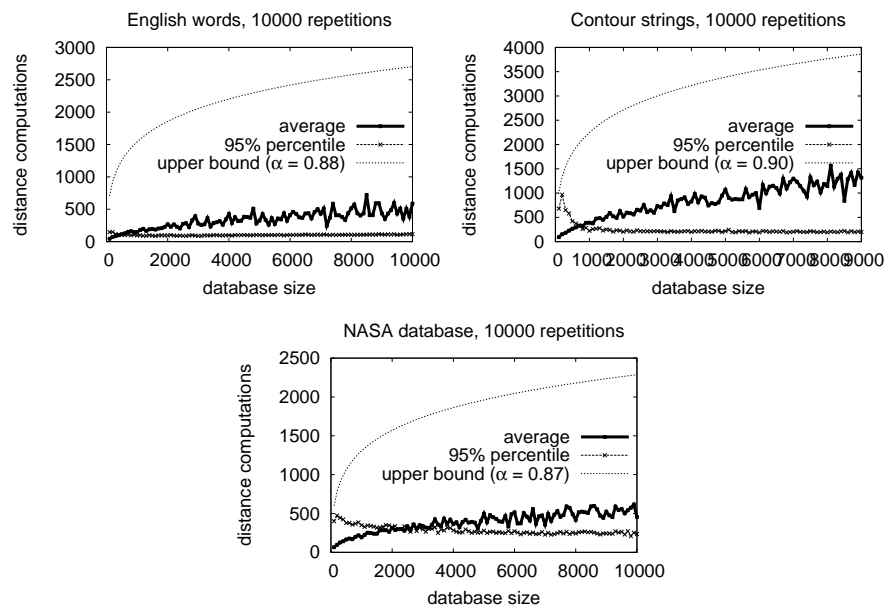


Figure 7: Experimental average number and theoretical upper bound of the distance computations caused by an insertion for increasing size database sets and for English, NASA and Contour databases.

283 insertions are incrementally done.

284 We have shown that the average number of distance computations (and the
285 average complexity) is bounded by a function that grows with the square of the
286 logarithm of the size of tree ($\log^2(n)$). This is a big improvement if compared
287 with the "naïve" approach that grows with $n \log(n)$.

288 Moreover, we have tested this upper bound with several artificial and real
289 data experiments. These experiments, as well as confirming the theoretical
290 results, also shows that the 95% percentile decreases when the database size
291 increases. That means that when the database size increases, pathological in-
292 sertions, which provokes wide reconstructions of the tree, becomes very uncom-
293 mon.

294 **Acknowledgements**

295 This work has been supported in part by grants TIN2009-14205-C04-01 from
296 the Spanish CICYT (Ministerio de Ciencia e Innovación), the IST Programme
297 of the European Community, under the Pascal Network of Excellence, IST-2002-
298 506778, and the program CONSOLIDER INGENIO 2010 (CSD2007-00018).

299 Battiato, S., Blasi, G. D., Reforgiato, D. Advanced indexing schema for imaging
300 applications: three case studies. *Image Processing, IET* 1 (3), 249–268.

301 Brin, S., 1995. Near neighbor search in large metric spaces. In: *Proceedings of*
302 *the 21st International Conference on Very Large Data Bases*. pp. 574–584.

303 Carrasco, R. C., Forcada, M. L., 1995. A note on the nagendraprasad-wang-
304 gupta thinning algorithm. *Pattern Recognition Letters* 16, 539–541.

305 Cha, S.-H., Srihari, S. N., 2002. On measuring the distance between histograms.
306 *Pattern Recognition* 35, 1355–1370.

307 Chávez, E., Navarro, G., Baeza-Yates, R., Marroquin. Searching in metric
308 spaces. *ACM Computing Surveys* 33 (3), 273–321.

- 309 Ciaccia, P., Patella, M., Zezula, P., 1997. M-tree: An efficient access method for
310 similarity search in metric spaces. In: Proceedings of the 23rd International
311 Conference on VLDB. Morgan Kaufmann Publishers, Athens, Greece, pp.
312 426–435.
- 313 Escolano, F., Vento, M. (Eds.), 2007. Graph-Based Representations in Pattern
314 Recognition, 6th IAPR-TC-15 International Workshop, GBRPR 2007, Ali-
315 cante, Spain, June 11-13, 2007, Proceedings. Vol. 4538 of Lecture Notes in
316 Computer Science. Springer.
- 317 Fu, A. W., Chan, P. M., Cheung, Y., Moon, Y. S., 2000. Dynamic vp-tree in-
318 dexing for n-nearest neighbor search given pair-wise distances. VLDB Journal
319 9, 154–173.
- 320 Garris, M., Wilkinson, R., 1994. NIST special database 3: Handwritten seg-
321 mented characters.
- 322 Giacinto, G., 2007. A nearest-neighbor approach to relevance feedback in content
323 based image retrieval. In: CIVR '07: Proceedings of the 6th ACM interna-
324 tional conference on Image and video retrieval. ACM, New York, NY, USA,
325 pp. 456–463.
- 326 Gómez-Ballester, E., Micó, L., Oncina. Some approaches to improve tree-based
327 nearest neighbour search algorithms. Pattern Recognition 39 (2), 171–179.
- 328 Hjaltason, G., Samet, H., 2003. Index-driven similarity search in metric spaces.
329 ACM Trans. Database Syst. 28 (4), 517–580.
- 330 Kalantari, I., McDonald, G., 1983. A data structure and an algorithm for the
331 nearest point problem. IEEE Trans. Software Eng. 9 (5), 631–634.
- 332 Lundsteen, C., Philip, J., Granum, E., 1980. Quantitative analysis of 6985 dig-
333 itized trypsin g-banded human metaphase chromosomes. Clin Genet 18 (5),
334 355–70.

- 335 Micó, L., Oncina, J., Carrasco, R., 1996. A fast branch and bound nearest
336 neighbor classifier in metric spaces. *Pattern Recognition Letters* 17, 731–73.
- 337 Micó, L., Oncina, J., Vidal, E., 1994. A new version of the nearest-neighbour ap-
338 proximating and eliminating search algorithm (aesa) with linear preprocessing
339 time and memory requirements. *Pattern Recognition Letters* 15, 9–17.
- 340 Micó, L., Oncina, J., 2009. Experimental analysis of insertion costs in a naïve
341 dynamic mdf-tree. *Lecture Notes in Computer Science* 5524, 402–408.
- 342 Navarro, G., 2002. Searching in metric spaces by spatial approximation. *VLDB*
343 *Journal* 11 (1), 28–46.
- 344 Navarro, G., Reyes, N., 2002. Fully dynamic spatial approximation trees. In:
345 *Proceedings of the 9th International Symposium on String Processing and*
346 *Information Retrieval (SPIRE 2002)*, LNCS 2476. Springer, pp. 254–270.
- 347 Navarro, G., Reyes, N., 2008. Dynamic spatial approximation trees. *J. Exp.*
348 *Algorithmics* 12, 1–68.
- 349 Noltemeier, H., Verbarq, K., Zirkelbach, C., 1992. Monotonous bisector^{*} trees
350 – a tool for efficient partitioning of complex scenes of geometric objects. In:
351 *Data Structures and Efficient Algorithms*. pp. 186–203.
- 352 Procopiuc, O., Agarwal, P. K., Arge, L., Vittery, J. S., 2003. Bkd-tree: A
353 dynamic scalable kd-tree. In: *In Proc. International Symposium on Spatial*
354 *and Temporal Databases*. pp. 46–65.
- 355 Katayama, N., Satoh, S., 1999. Analysis and Improvement of the SR-tree: an
356 Index Structure for Nearest Neighbor Searches. 6th DIMACS Implementation
357 Challenge Workshop: Near Neighbor Searches, Baltimore (Maryland) (USA)
358 (1999)
- 359 Torralba, A., Fergus, R., Freeman, W., 2008. 80 million tiny images: A large
360 data set for nonparametric object and scene recognition. *IEEE Transactions*
361 *on Pattern Analysis and Machine Intelligence* 30 (11), 1958–1970.

362 Yianilos, P., 1993. Data structures and algorithms for nearest neighbor search
363 in general metric spaces. In: Proceedings of the ACM-SIAM Symposium on
364 Discrete Algorithms. pp. 311–321.

365 Zezula, P., Amato, G., Dohnal, V., Batko, M., 2006. Similarity Search: The
366 Metric Space Approach. Springer.