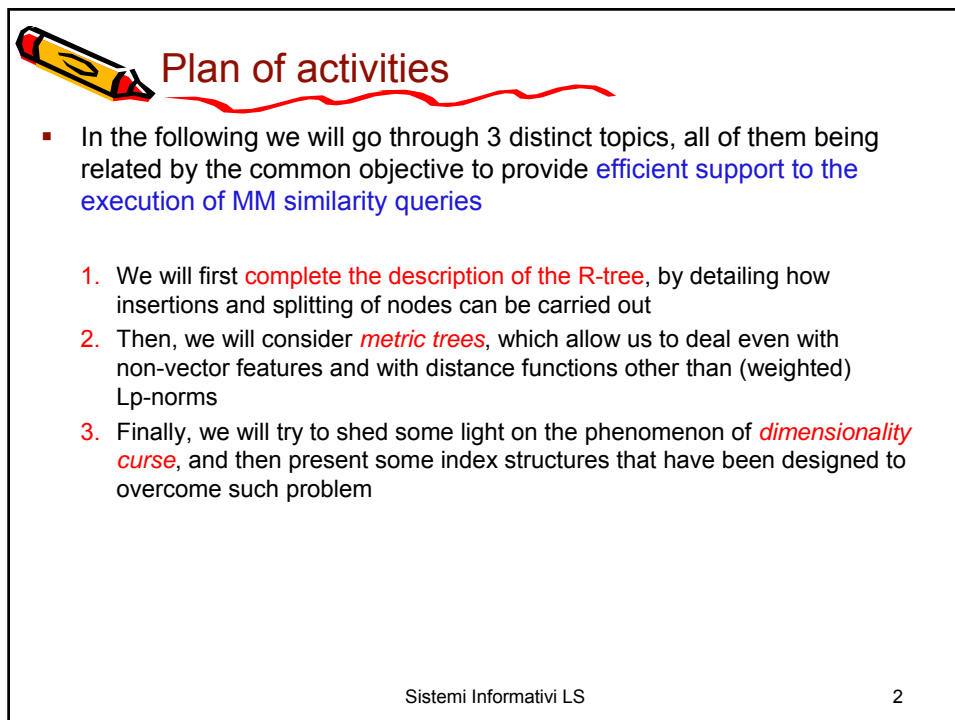


Indexing Multimedia Databases

Prof. Paolo Ciaccia
<http://www-db.deis.unibo.it/courses/SI-LS/>
09_IndexingMMDBs.pdf

Sistemi Informativi LS



Plan of activities

- In the following we will go through 3 distinct topics, all of them being related by the common objective to provide **efficient support to the execution of MM similarity queries**
 1. We will first **complete the description of the R-tree**, by detailing how insertions and splitting of nodes can be carried out
 2. Then, we will consider **metric trees**, which allow us to deal even with non-vector features and with distance functions other than (weighted) Lp-norms
 3. Finally, we will try to shed some light on the phenomenon of **dimensionality curse**, and then present some index structures that have been designed to overcome such problem

Sistemi Informativi LS

2



Back to the R-tree (Guttman, 1984)

- Remind what said some weeks ago:

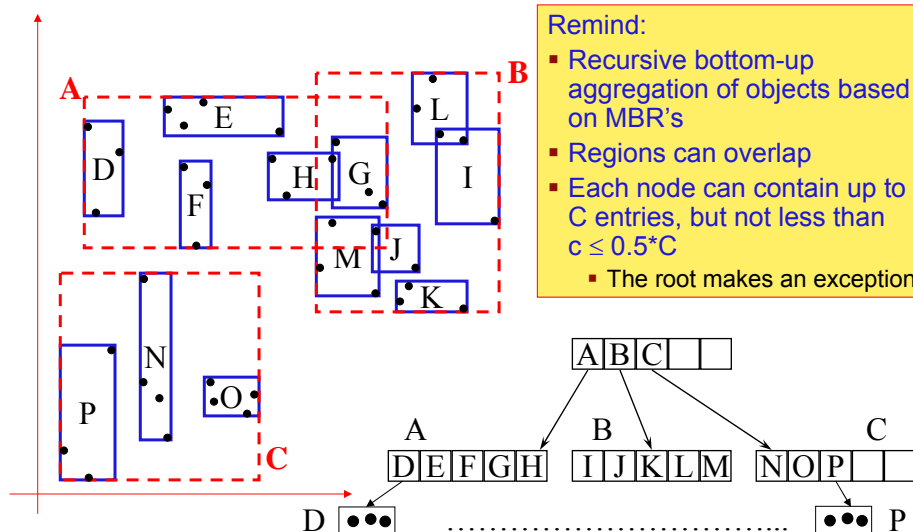


Be sure to understand what the index looks like and how it is used to answer queries; for the moment don't be concerned on how an R-tree with a given structure can be built!

- It's now time to discuss **how an R-tree can be effectively built**
- It has to be considered that many "*R-tree variants*" exist, and it's not our intention to go through their details
- It just suffices to say that one of such variants leads to what is known as the **R*-tree** [BKS+90], which is the commonest version in use
- With respect to the original proposal [Gut84], the R*-tree adds smarter **insertion and split heuristics**, plus a so-called "forced reinsert" technique that we do not consider here



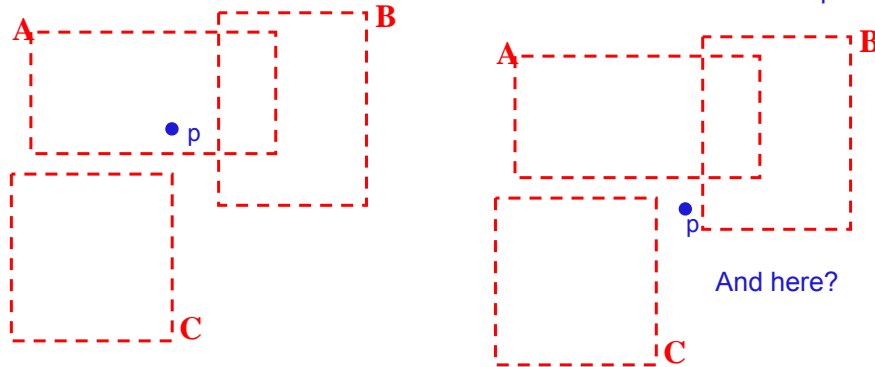
R-tree: how it looks like





R-tree: insertion of a new object

- We start from the root and move down the tree one step at a time, trying to find a “nice place” where to accommodate the new object p
 - For simplicity, we assume that indexed objects are points, similar arguments apply if we index (hyper-)rectangles (MBR's)
- At each step we have a same question to answer: Which child node is the most suitable to accommodate p ?



Sistemi Informativi LS

5



R-tree: the ChooseSubtree method

- The recursive algorithm that descends the tree to insert a new object p , together with its TID, is called **ChooseSubtree**

ChooseSubtree($E_p=(p,TID),ptr(N)$)

1. Read(N);
2. If N is a leaf then: return N // we are done
3. else: { choose among the entries E_c in N
the one, E_{c^*} , for which **Penalty**(E_p, E_{c^*}) is minimum;
4. return **ChooseSubtree**($E_p, E_{c^*}.ptr$) } // recursive call
5. end.

- We invoke the method on the index root
- The specific criterion used to decide “how bad” an entry is, should we choose it to insert p , is encapsulated in the **Penalty** method
 - Variants of the R-tree differ in how they implement **Penalty**
- This insertion algorithm is the one used by most multi-dimensional and metric trees

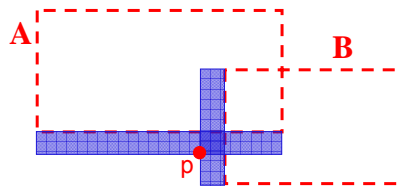
Sistemi Informativi LS

6



R-tree: the Penalty method

- If point p is inside the region of an entry E_c , then the penalty is 0
- Otherwise, Penalty can be computed as the **increment of volume (area) of the MBR**
 - However, if E_c points to a leaf node, then [BKS+90] shows that it's better to consider the **increment of overlap with the other entries**
- Both criteria aim to obtain trees with better performance:
 - **Large area**: increases the number of nodes to be visited by a query
 - **Large overlap**: also degrades performance



B is better than A

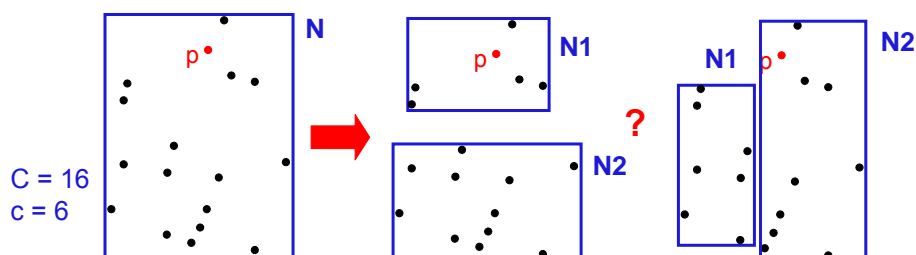
Sistemi Informativi LS

7



R-tree: splitting of a leaf node

- When p has to be inserted into a leaf node that already contains C entries, an overflow occurs, and **N has to be split**
- For leaf nodes whose entries are points the solution aims to split the set of $C+1$ points into 2 subsets, each with at least c and at most C points
- Among the several possibilities, one could consider the choice that leads to have a **minimum overall area**
 - However, this is an **NP-Hard problem**, thus heuristics have to be applied



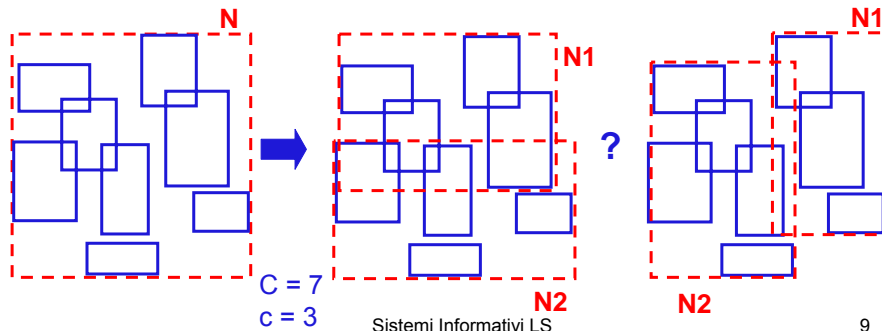
Sistemi Informativi LS

8



R-tree: splitting of a non-leaf node

- As in B+-trees, splits propagate upward and can recursively trigger splits at higher levels of the tree
- The problem to be faced now is how to split a set of $C+1$ (hyper-)rectangles
 - Note that this applies also to leaf nodes if they store MBR's
- The original proposal just aims to minimize the sum of resulting areas
- The R*-tree implements a more sophisticated criterion, which takes into account the areas, overlap, and perimeters of the resulting regions

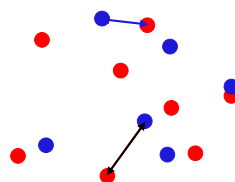


Beyond vector spaces

- It's a matter of fact that vector spaces, equipped with some (weighted) Lp-norm, are not general enough to deal with the whole variety of feature types and distance functions needed in MMDB's

Example:

given 2 sets of points $s1$ and $s2$, their Hausdorff distance is defined as follows:



- 1 \forall (red) point of $s1$ find the closest (blue) point in $s2$
Let $h(s1, s2)$ be the maximum of such distances
- 2 \forall (blue) point in $s2$ find the closest (red) point in $s1$
Let $h(s2, s1)$ be the maximum of such distances
- 3 Let $d_{\text{Haus}}(s1, s2) = \max\{h(s1, s2), h(s2, s1)\}$

Used for matching shapes



Another example: set similarity

- We have **logs of WWW accesses**, where each log entry has format like:

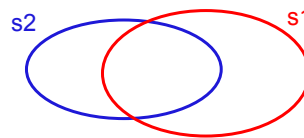
```
www-db.deis.unibo.it pciaccia -
[11/Jan/1999:10:41:37 +0100]
"GET /~mpatella/ HTTP/1.0" 200 1573
```

- Log entries are grouped into **sessions** (= sets of visited pages):

$S = \langle \text{ip_address}, \text{user_id}, [\text{url}_1, \dots, \text{url}_k] \rangle$

and we want to compare "**similar sessions**" (i.e., similar sets), using:

$$d_{\text{setdiff}}(s1, s2) = \frac{|s1 - s2| + |s2 - s1|}{|s1| + |s2|}$$



Another example: edit distance

- A common distance measure for **strings** is the so-called **edit distance**, defined as the **minimum number of characters that have to be inserted, deleted, or substituted** so as to transform a string $s1$ into another string $s2$

$d_{\text{edit}}(\text{'ball'}, \text{'bull'}) = 1$ $d_{\text{edit}}(\text{'balls'}, \text{'bell'}) = 2$ $d_{\text{edit}}(\text{'rather'}, \text{'alter'}) = 3$

- The edit distance is also commonly used in **genomic DB's** to compare **DNA sequences**. Each DNA sequence is a **string over the 4-letters alphabet of bases**:

a: adenine

c: cytosine

g: guanine

t: thymine

$d_{\text{edit}}(\text{'gatctggtgg'}, \text{'agcaaatacag'}) = 7$

g	a	t	c	t	g	g	t	g	-	g
1	=	2	=	3	4	5	=	6	7	=
-	a	g	c	a	a	a	t	c	a	g

The edit distance can be computed using a dynamic programming procedure, similar to the one seen for the DTW



Computing the Edit Distance

- The cost matrix is used to incrementally build the new matrix d_{edit} , whose elements are recursively defined as:

$$d_{edit;i,j} = cost_{i,j} + \min\{d_{edit;i-1,j}, d_{edit;i,j-1}, d_{edit;i-1,j-1}\}$$

r	1	1	1	1	1	0
e	1	1	1	1	0	1
h	1	1	1	1	1	1
t	1	1	1	0	1	1
a	1	0	1	1	1	1
r	1	1	1	1	1	0
	0	1	1	1	1	1
cost		a	l	t	e	r

s2 ↑

→ s1

r	6	5	5	5	4	3
e	5	4	4	4	3	4
h	4	3	3	3	3	4
t	3	2	3	2	3	4
a	2	1	2	3	4	5
r	1	1	2	3	4	4
	0	1	2	3	4	5
dedit		a	l	t	e	r

s2 ↑

→ s1

Sistemi Informativi LS

13



Metric spaces

- A metric space $M = (U, d)$ is a pair, where
 - U is a domain ("universe") of values, and
 - d is a distance function that, $\forall x, y, z \in U$, satisfies the **metric axioms**:

$$\begin{aligned} d(x, y) &\geq 0, d(x, y) = 0 \Leftrightarrow x = y && \text{(positivity)} \\ d(x, y) &= d(y, x) && \text{(symmetry)} \\ d(x, y) &\leq d(x, z) + d(z, y) && \text{(triangle inequality)} \end{aligned}$$

- All the distance functions seen in the previous examples are metrics, and so are the (weighted) Lp-norms
- The only distance we have seen so far that does not fit the metric framework is the **DTW**

Metric indexes only use the metric axioms to organize objects, and exploit the triangle inequality to prune the search space

Sistemi Informativi LS

14



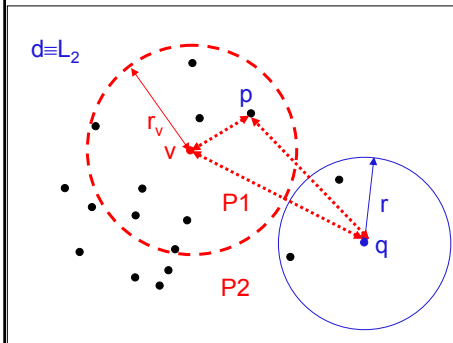
Principles of metric indexing (1)

- Given a “metric dataset” $P \subseteq U$, one of the two following principles can be applied to partition it into two subsets

Ball decomposition: take a point v (“vantage point”), compute the distances of all other points p w.r.t. v , $d(p, v)$, and define

$$P1 = \{p : d(p, v) \leq r_v\} \quad P2 = \{p : d(p, v) > r_v\}$$

If r_v is chosen so that $|P1| \approx |P2| \approx |P|/2$ we obtain a balanced partition



Consider a range query $\{p : d(p, q) \leq r\}$

If $d(q, v) > r_v + r$ we can conclude that no point in $P1$ belongs to the result

Proof:

we show that $d(p, q) > r$ holds $\forall p \in P1$.

$$\begin{aligned} d(p, q) &\geq d(q, v) - d(p, v) && \text{(triangle ineq.)} \\ &> r_v + r - d(p, v) && \text{(by hyp.)} \\ &\geq r_v + r - r_v && \text{(by def. of P1)} \\ &\geq r \end{aligned}$$

Similar arguments can be applied to $P2$

Sistemi Informativi LS

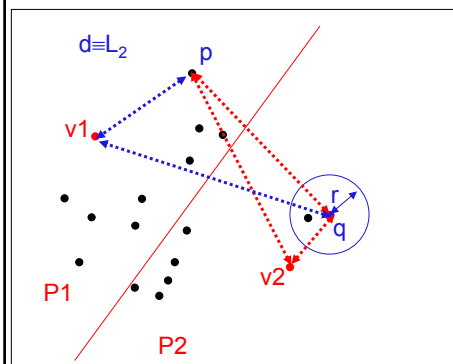
15



Principles of metric indexing (2)

Generalized Hyperplane: take two points $v1$ and $v2$, compute the distances of all other points p w.r.t. $v1$ and $v2$, and define

$$P1 = \{p : d(p, v1) \leq d(p, v2)\} \quad P2 = \{p : d(p, v2) < d(p, v1)\}$$



Consider a range query $\{p : d(p, q) \leq r\}$

If $d(q, v1) - d(q, v2) > 2*r$ we can conclude that no point in $P1$ belongs to the result

Proof:

we show that $d(p, q) > r$ holds $\forall p \in P1$.

$$\begin{aligned} d(q, v1) - d(p, q) &\leq d(p, v1) && \text{(triangle ineq.)} \\ d(p, v1) &\leq d(p, v2) && \text{(def. of P1)} \\ d(p, v2) &\leq d(p, q) + d(q, v2) && \text{(triangle ineq.)} \end{aligned}$$

Then:

$$\begin{aligned} d(q, v1) - d(p, q) &\leq d(p, q) + d(q, v2) \\ d(p, q) &\geq (d(q, v1) - d(q, v2))/2 \\ &> r && \text{(by hyp.)} \end{aligned}$$

Sistemi Informativi LS

16



The M-tree (Ciaccia, Patella & Zezula, 1997)

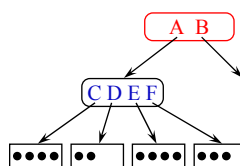
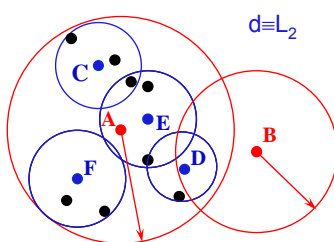
- The M-tree has been **the first dynamic, paged, and balanced metric index**
- Intuitively, it **generalizes “R-tree principles” to arbitrary metric spaces**
 - The M-tree treats the distance function as a “black box”
- Since 1997 [CPZ97] it has been used by several research groups for:
 - Image retrieval, text indexing, shape matching, clustering algorithms (including the WWW log example), fingerprint matching, DNA DB's, etc.
 - [CNB+01] and [HS03] are both excellent surveys on searching in metric spaces
- C++ source code freely available at <http://www-db.deis.unibo.it/Mtree/>



Remind: at a first sight, the M-tree “looks like” an R-tree.
However, remember that **the M-tree only “knows” about distance values**, thus **it ignores coordinate values and does not rely on any “geometric” (coordinate-based) reasoning**

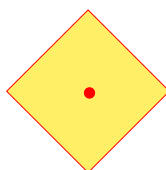


M-tree: how it looks like

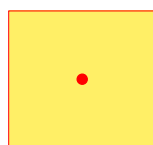


- Recursive bottom-up aggregation of objects based on regions
- Regions can overlap
- Each node can contain up to C entries, but not less than $c \leq 0.5 \cdot C$
 - The root makes an exception

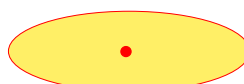
- Depending on the metric, the “shape” of index regions changes



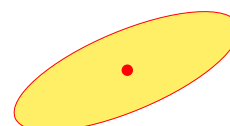
L_1



L_∞



Weighted Euclidean
Sistemi Informativi LS



quadratic distance



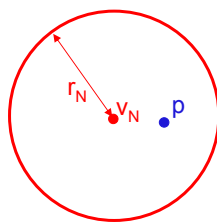
The M-tree regions

- Each node N of the tree has an associated region, $\text{Reg}(N)$, defined as

$$\text{Reg}(N) = \{p: p \in U, d(p, v_N) \leq r_N\}$$

where:

- v_N (the “center”) is also called a *routing object*, and
- r_N is called the *(covering) radius* of the region
- The set of indexed points p that are reachable from node N are guaranteed to have $d(p, v_N) \leq r_N$



- This immediately makes it possible to apply the pruning principle:
If $d(q, v_N) > r_N + r$ then prune node N:



Entries of leaf and internal nodes

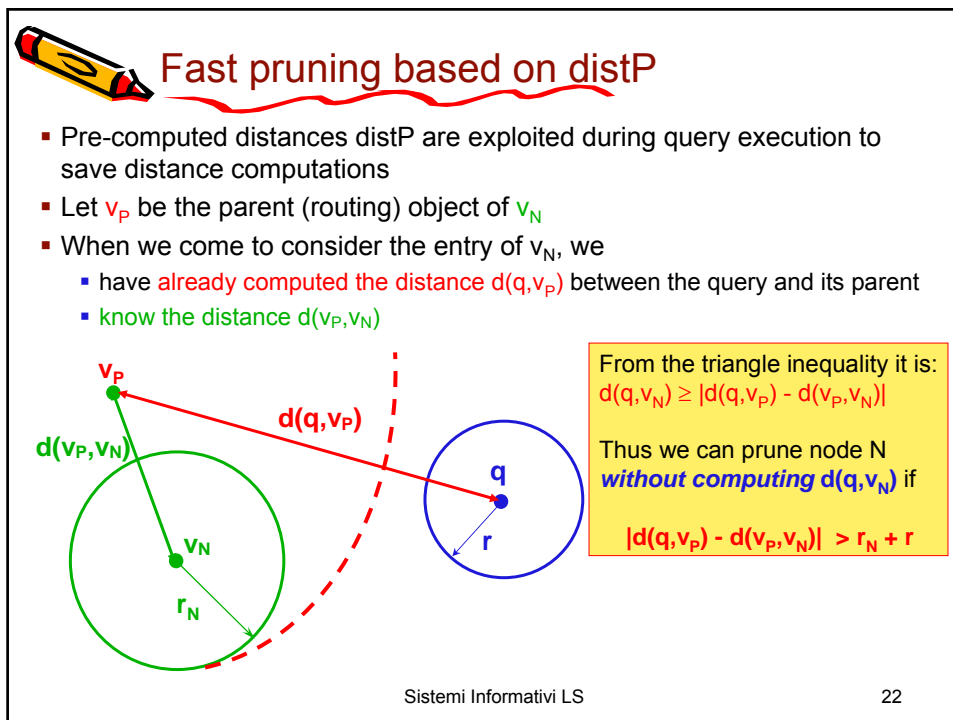
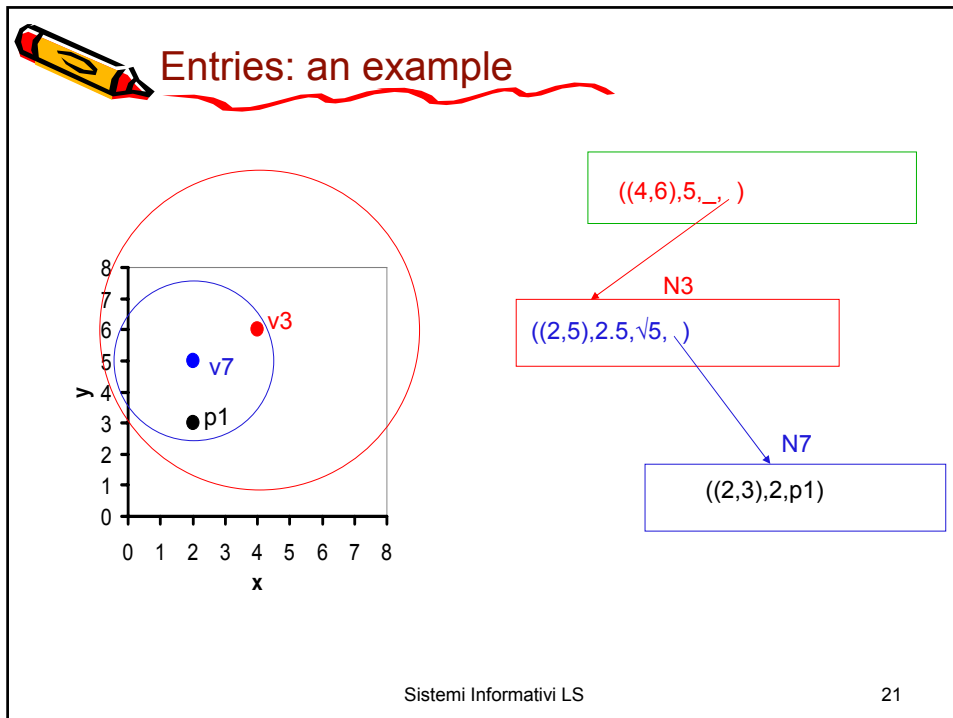
- Each node N stores a variable number of *entries*

Leaf node:

- An entry E has the form $E = (\text{ObjFeatures}, \text{distP}, \text{TID})$, where
 - ObjFeatures* are the feature values of the indexed object
 - distP* is the distance between the object and its parent routing object (i.e., the routing object of node N)

Internal node:

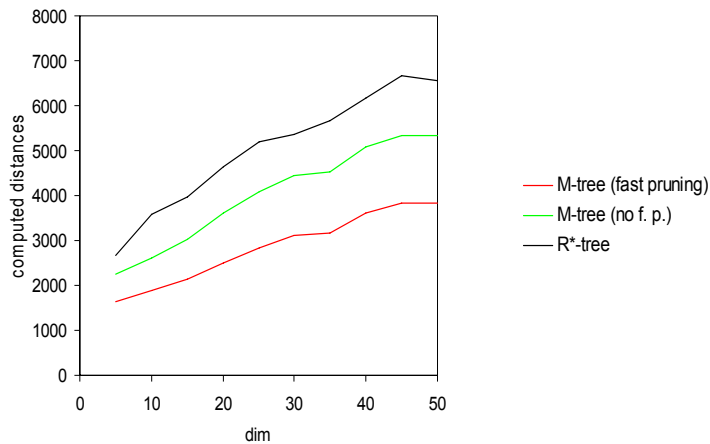
- An entry E has the form $E = (\text{RoutingObjFeatures}, \text{CoveringRadius}, \text{distP}, \text{PID})$, where
 - RoutingObjFeatures* are the feature values of the routing object
 - CoveringRadius* is the radius of the region
 - distP* is the distance between the routing object and its parent routing object (this is undefined for entries in the root node)





Experimental results

- Synthetic datasets (10 Gaussian clusters)
- Up to 40% cost reduction with fast pruning



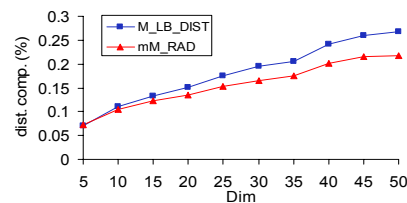
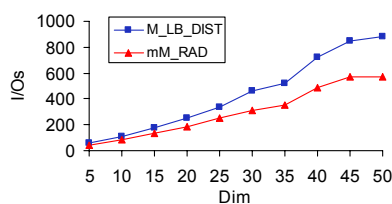
Sistemi Informativi LS

24



Insertion and split (sketch)

- The procedure to insert a new object is based on the ChooseSubtree method
- The Penalty method considers the increase of the covering radius needed to accommodate the new object
 - Remind: no "volume" can be computed!
- For managing a split, there are several alternatives, among which [CPZ97]:
 - mM_RAD minimize the maximum of the two resulting radii
 - M_LB_DIST choose the closest and the farthest object from v_N
- Experiments demonstrate that mM_RAD is the best



Sistemi Informativi LS

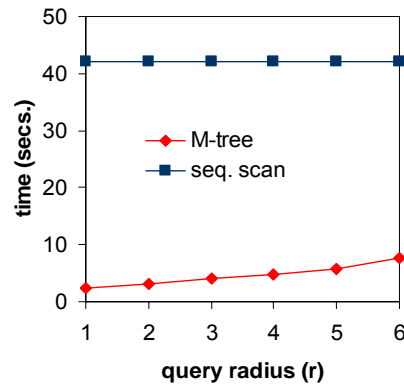
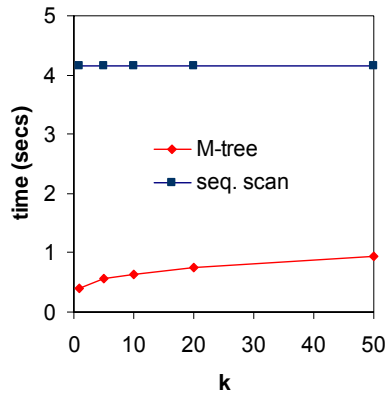
25



Experiments (k-NN and range queries)

- 68,000 color images
- 32-dim (color histograms), L_2
- 161,212 text rows
- Edit distance

The logic of search algorithms is the one already seen for range and k-NN queries with the R-tree



Sistemi Informativi LS

26



High-dimensional spaces (1)

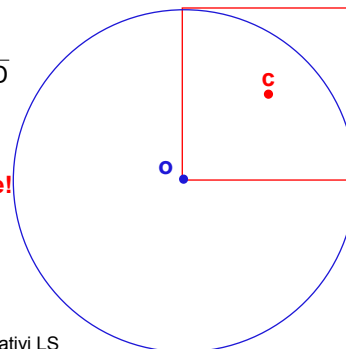
- The geometry of high-dimensional spaces is intriguing, since our common-sense intuitions fail, as the following examples show

1st example: "is the center in the sphere?"

- Consider the unitary hypercube $[0,1]^D$ with center $c = (0.5, \dots, 0.5)$, and the D -dimensional hypersphere S^D centered in the origin $o = (0, \dots, 0)$ and with radius $r = 1$.
- Our intuition, and the figure as well, confirms that for $D=2$ c is inside S^D
- Let's see what happens when D grows:

$$L_2(c, o) = \sqrt{\sum_{i=1, D} 0.5^2} = \sqrt{D \times 0.5^2} = 0.5 \times \sqrt{D}$$

- Thus, when $D > 4$ c is external to the sphere!



Sistemi Informativi LS

27



High-dimensional spaces (2)

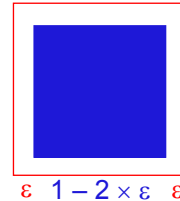
2nd example: "where are the points?"

- Consider again the unitary hypercube $[0,1]^D$
- Now, take a hypercube B of side $1 - 2 \times \varepsilon$ and center $c = (0.5, \dots, 0.5)$
- The volume of B grows like

$$\text{Vol}(B) = (1 - 2 \times \varepsilon)^D$$

- As the table shows, even for (very) small ε values, $\text{Vol}(B)$ sharply reduces

$\varepsilon \setminus D$	2	50	100	500	1000
0.1	0.64	1.43E-05	2.04E-10	3.51E-49	1.23E-97
0.05	0.81	0.01	2.66E-05	1.32E-23	1.75E-46
0.01	0.96	0.36	0.13	4.10E-05	1.68E-09



- If we have N points uniformly distributed over $[0,1]^D$, then only a fraction equal to $\text{Vol}(B)$ will be contained, on the average, in B
- Thus, **all points are close to the surface of $[0,1]^D$!**

Sistemi Informativi LS

28

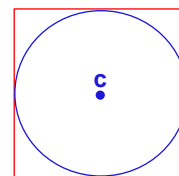


High-dimensional spaces (3)

3rd example: "How big a sphere is?"

- Consider the unitary hypercube $[0,1]^D$ and the D -dimensional hypersphere S^D centered in $c = (0.5, \dots, 0.5)$ and with radius $r = 0.5$
- The volume of S^D can be computed as (D even):
$$\text{Vol}(S^D) = \frac{\pi^{D/2} \times 0.5^D}{(D/2)!}$$
- The following table (from [WSB98]) shows, for various values of D and assuming that points are uniformly distributed over $[0,1]^D$:
 - The volume of S^D , $\text{Vol}(S^D)$
 - The number of points N needed to have, on the average, at least 1 point in S^D (this is just $1 / \text{Vol}(S^D)$)

D	$\text{Vol}(S^D)$	N
2	0.785	1.27
4	0.308	3.24
10	0.002	401.50
20	2.46E-08	40631627
40	3.28E-21	3.05E+20
100	1.87E-70	5.35E+69



- Thus, **the number of points should grow exponentially to have at least 1 point in S^D !**

Sistemi Informativi LS

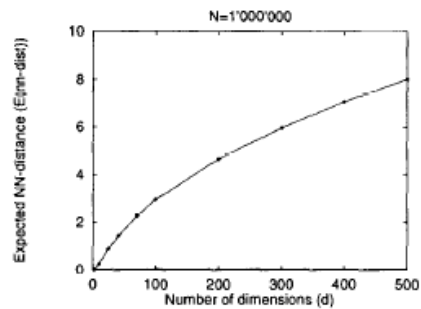
29



High-dimensional spaces (4)

4th example: "How far is the nearest neighbor?"

- Continuing with the previous example, we can compute the expected (Euclidean) distance of the nearest neighbor of the center $c = (0.5, \dots, 0.5)$ of S^D
- The following graph (from [WSB98]) shows how the NN distance grows with D when $N = 10^6$



- Thus, **the closest point is far away!**

Sistemi Informativi LS

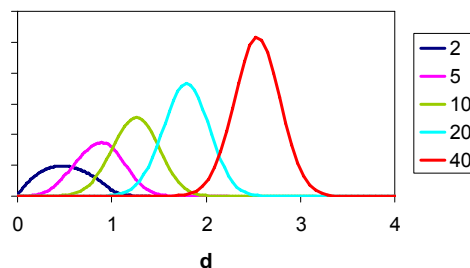
30



High-dimensional spaces (5)

5th example: "How far are the other points?"

- We now plot the **distance distribution of the dataset**, for various values of D
- The distance distribution shows, for a given value of d , which is the percentage of points whose distance is d



- It can be observed that when D grows, **the variance of distances decreases**
- Thus, **in high-dimensional spaces all the points tend to have the same distance from the query!**

Sistemi Informativi LS

31



Basic facts about high-dim. spaces (1)

- The analysis in [WSB98] demonstrates that, **no matter how smart you are in designing a new index structure, there always exists a value of D such that the index performance will deteriorate, and sequential scan will become the best alternative!**
- However, **the analysis applies to uniformly distributed datasets and Euclidean distance...**
- If your data are not uniformly distributed (as it always happens!), then the authors argue that their analysis still applies, provided one considers the **"intrinsic dimensionality"** of the dataset
- The concept of "intrinsic dimensionality" is not precisely definable, intuitively it is the **"true dimensionality"** of our data
 - E.g.: **a line has intrinsic dimensionality 1, regardless of D**
- Some attempts to characterize the intrinsic dimensionality of a dataset have been based on the concept of **fractals** (e.g., see [FK94])

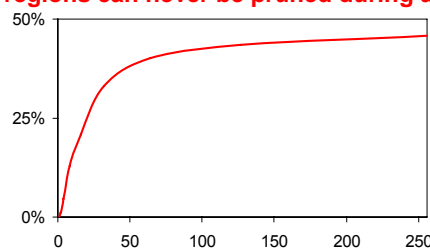
Sistemi Informativi LS

32



Basic facts about high-dim. spaces (2)

- From a more pragmatical point of view, experimental results obtained with both spatial and metric indexes confirm that high-dimensional datasets are often a nightmare!
- This is the so-called **"dimensionality curse"**!
- For the structures we have seen (R-tree and M-tree), what is observed is an **incredible amount of overlap between the regions of index nodes**
 - The graph shows the **percentage of M-tree regions that enclose a query point q**, i.e., those regions for which $d_{\min}(q, \text{Reg}(N)) = 0$
 - Thus, **all such regions can never be pruned during a k-NN search!**



Sistemi Informativi LS

33



Note: Partitioning without overlap

- If we partition the $[0,1]^D$ space into **non-overlapping regions**, similar problems arise
- For instance, consider a uniform distribution of points, and assume we **split a dimension in the mid-point 0.5** (thus, each time we double the number of regions). **We can split at most $D' = \lceil \log_2 N \rceil$ dimensions**
- Consider the region: **$\text{Reg} = [0,0.5] \times \dots \times [0,0.5] \times [0,1] \times \dots \times [0,1]$** whose farthest point is **$q = (1, \dots, 1)$**
- The Euclidean distance of q from Reg is:

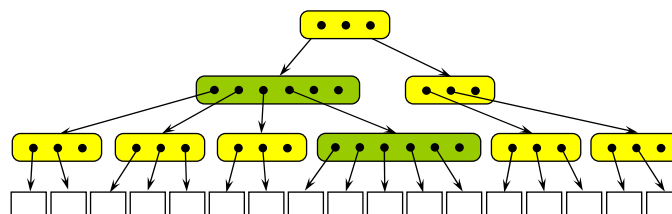
$$L_2(\text{Reg}, q) = \sqrt{\sum_{i=1, D'} (1-0.5)^2} = \sqrt{D' \times 0.5^2} = 0.5 \times \sqrt{D'} = 0.5 \times \sqrt{\lceil \log_2 N \rceil}$$

- With $N = 10^6$ we have $D'=20$ and $L_2(\text{Reg}, q)=2.236$
- Since this is independent of D , whereas the expected NN distance grows with D , **for values of D large enough ($D \geq 80$) Reg will be accessed, and this holds for any other region!**



The X-tree [BKK96]: basic idea

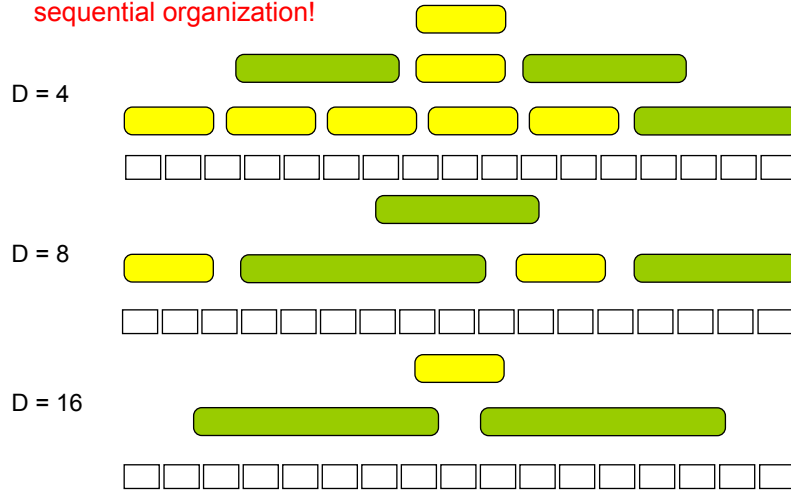
- The **X-tree** is an **evolution of the R-tree**, aiming to deal with the “overlap problem”
- When a node has to be split, if an overlap-free split is possible then it is performed as usual, otherwise a new, larger, **super-node**, is allocated
 - Thus, now we have nodes of variable size
- The price to be paid is that searching within a super-node is more costly than searching within nodes





The X-tree: what happens when D grows

- Although the X-tree performs better than the R-tree for medium values of D, when the dimensionality increases the index degenerates to a sequential organization!

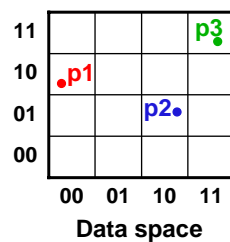


36



The VA-file (Weber, Schek & Blott, 1998)

- The basic idea of the VA-file [WSB98] is to speed-up the sequential scan by exploiting a "Vector Approximation"
- Each dimension of the data space is partitioned into 2^{b_i} intervals using b_i bits
 - E.g.: the 1st coordinate uses 2 bits, which leads to the intervals 00, 01, 10, and 11
- Thus, each coordinate of a point (vector) requires now b_i bits instead of 32
- The VA-file stores, for each point of the dataset, its approximation, which is a vector of $\sum_{i=1,D} b_i$ bits



p1	0.1	0.6	Feature values
p2	0.7	0.4	
p3	0.9	0.3	

p1	00	10	VA-file
p2	10	01	
p3	11	11	

Sistemi Informativi LS

37

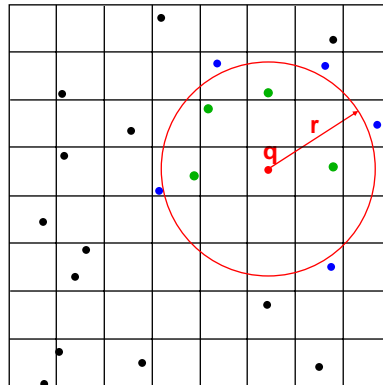


The VA-file: query processing

- Query processing with the VA-file is based on a **filter & refine approach**
- For simplicity, consider a range query

Filter: the VA file is accessed and only the points in the regions that intersect the query region are kept

Refine: the feature vectors are retrieved and an exact check is made



actual results
false drops
excluded points

Sistemi Informativi LS

38



Conclusions (?)

- The issue of efficiently indexing complex datasets is far from having been solved
- Starting from the end of 90's, many solutions have been proposed, and new ideas have emerged
- Unfortunately, **the absence of a well-defined and accepted benchmark makes it almost impossible to compare all such solutions**
- The basic lesson to be learned is that, **no matter how a structure has been cleverly designed, ultimately it has to be contrasted with the sequential scan!**
- Thus, be skeptical if someone claims to have designed an index showing "superior performance" w.r.t. the others: **always look if sequential scan has been taken as a competitor!**

Sistemi Informativi LS

39