

Metric Spaces

- Elements:
 - \mathbf{U} : universe of objects.
 - \mathbf{d} : distance function.

- Definition:

Metric Space (\mathbf{U}, \mathbf{d}) .

- Data set $\mathbf{X} \subseteq \mathbf{U}$, $|\mathbf{X}| = N$.
 - Objective: given $\mathbf{q} \in \mathbf{U}$, retrieval all similar objects objects to \mathbf{q} in \mathbf{X} .

Properties

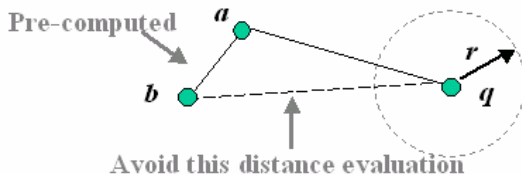
$d: U \times U \rightarrow \mathbf{R}^+$ distance function:

- $\forall x, y \in U, d(a, b) \geq 0$ positiveness
- $\forall x \in U, d(a, a) = 0$ strictly positiveness
- $\forall x, y \in U, d(a, b) = d(b, a)$ symmetry
- $\forall x, y, z \in U, d(a, b) \leq d(a, c) + d(c, b)$

Triangle inequality

Triangle inequality

- Reduce the number of distance evaluations performed during a range search operation:



$$\text{abs}(d(q,a) - d(b,a)) > r$$

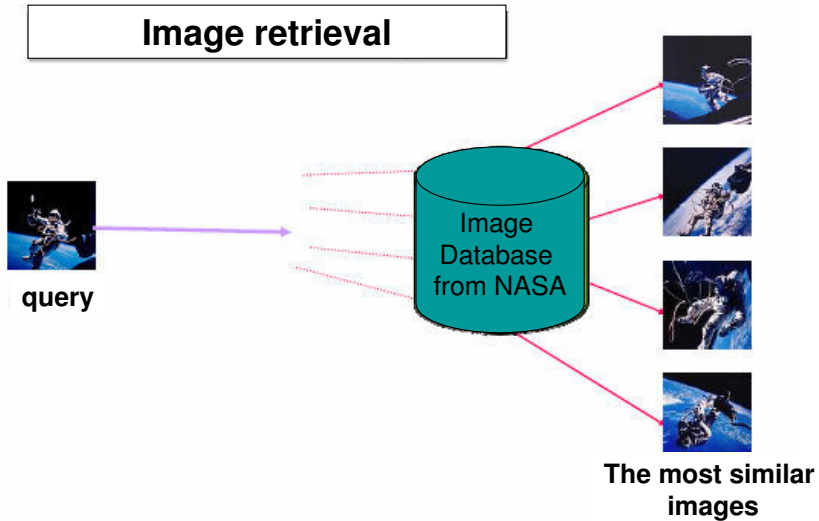
HYPOTHESIS: DISTANCE EVALUATION ARE EXPENSIVE TO COMPUTE
OBJECTIVE: REDUCE THE NUMBER OF DISTANCE EVALUATIONS

Applications

- Information retrieval
- Recognition of images
- Recognition of faces, voice or fingerprint
- Genetic databases
- Medical databases
- etc.



Example

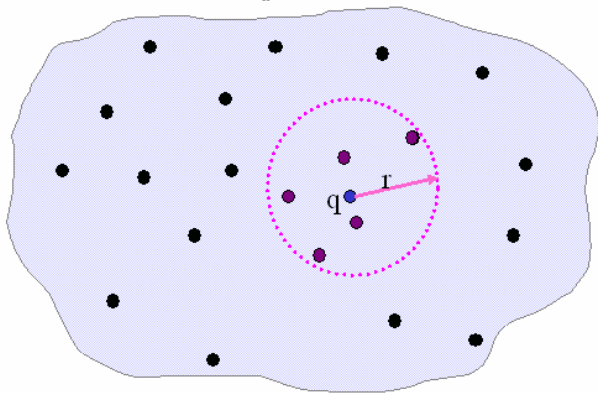


Similarity search

- There are basically two kinds of searches in a Metric Space:
 - Range Search (q, r) : retrieval all objects in X at distance at most r from the query q .
 - k -NN(q) : retrieval the k most similar objects to a query q .

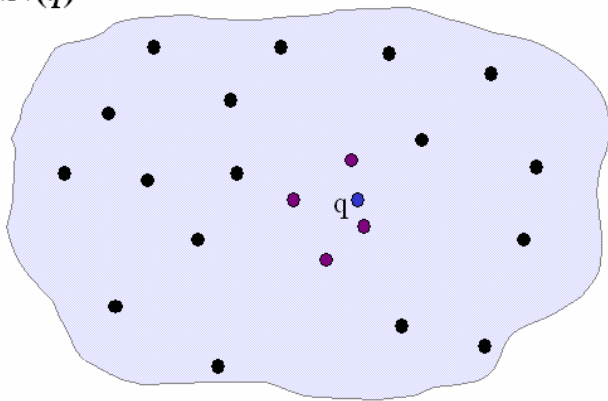
Range Search

Range query $(q, r)_d$

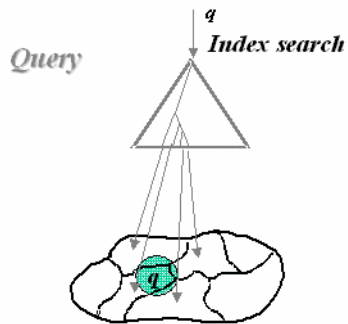
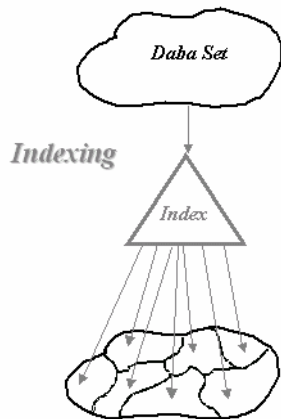


K -NN search

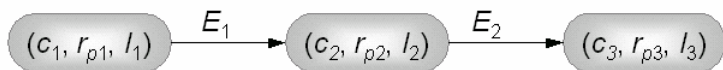
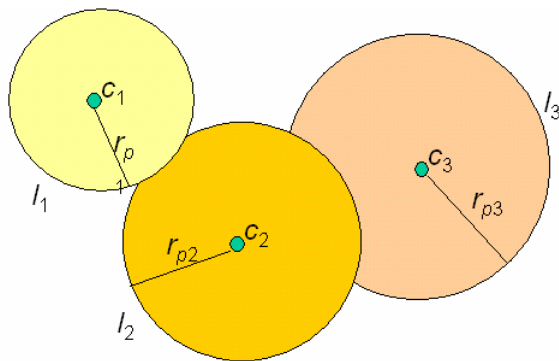
$4-NN(q)$



Indexing



List of Clusters (LC)



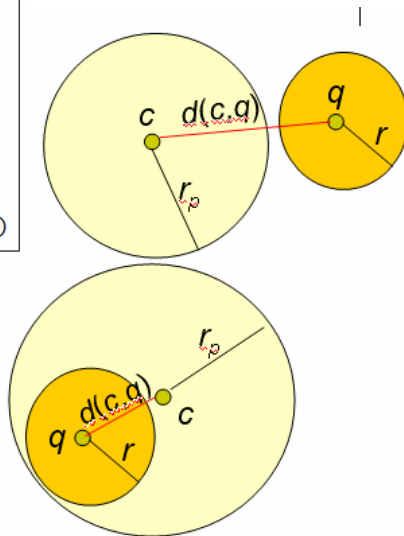
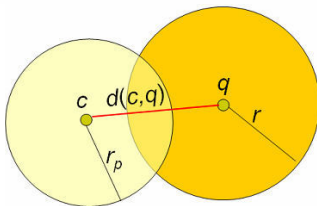
Build(U)

1. If $U = \emptyset$ Then Return an empty list
2. Select a center $c \in U$
3. $I \leftarrow kNN(c)$ in $U - \{c\}$
4. Let be $r_c = \max_{x \in I} d(x, c)$
5. $E \leftarrow U - I$
6. Return $(c, r_c, I) : \text{Build}(E)$

Searching on the LC

Search(L, q, r)

1. If L is empty Then Return
2. Let $L = (c, r_c, I) : E$
3. Compute the distance $d(c, q)$
4. If $d(c, q) \leq r$ Add c to the set of results
5. If $d(c, q) \leq r_c + r$ Then Search I exhaustively
6. If $d(c, q) > r_c - r$ Then Search(E, q, r)



Sparse Spatial Selection - SSS

- Maximum distance between two objects in the database

$$M = \max \{d(x, y) / x, y \in X\}$$

Pivots Selection Stage

PIVOTS $\leftarrow \{x_1\}$

for all $x_i \in U$ **do**

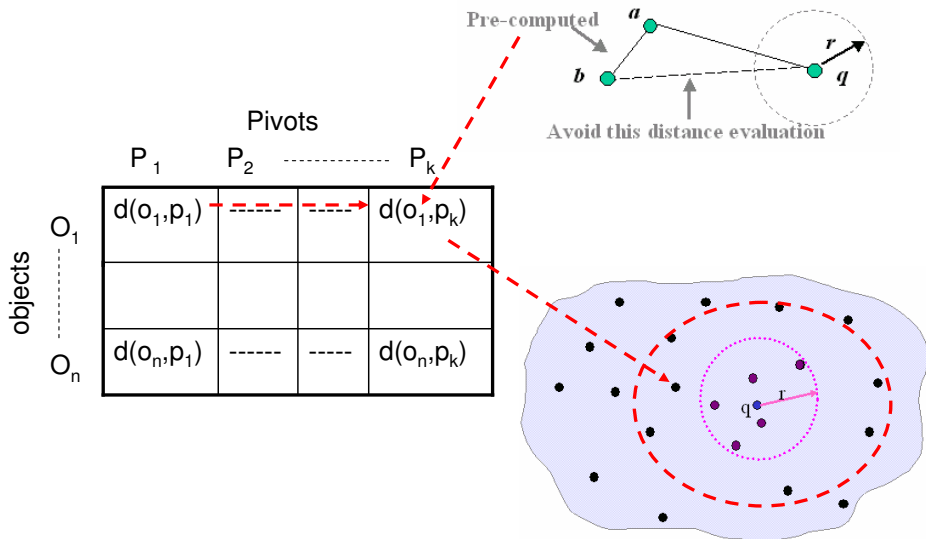
if $\forall p \in \text{PIVOTS}, d(x_i, p) \geq 0.5 * M$ **then**

 PIVOTS $\leftarrow \text{PIVOTS} \cup \{x_i\}$

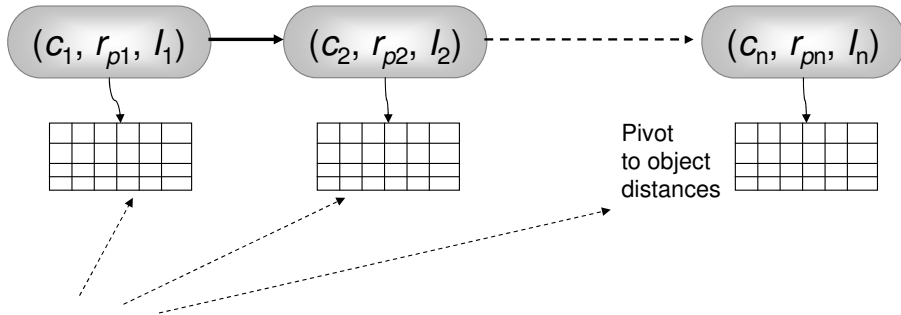
Searching on the SSS

		Pivots			
		P_1	P_2	-----	P_k
objects	O_1	$d(o_1, p_1)$	-----	-----	$d(o_1, p_k)$
	O_n	$d(o_n, p_1)$	-----	-----	$d(o_n, p_k)$

Searching on the SSS



LC-SSS Combination



SSS Pivots

The same SSS pivots in each table (bucket).

**Also the same ordering of pivots in the table columns,
the first two are the most distant ones, and so on.**