

SDS 385: Stat Models for Big Data

Lecture 9: KD trees

Purnamrita Sarkar
Department of Statistics and Data Science
The University of Texas at Austin
<https://psarkar.github.io/teaching>

Background

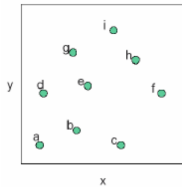
- Has a long history—invented in 1970 by Jon Bentley
- k represents the number of dimensions
- Idea is to partition the data spatially, by using only one dimension at any level.
- While searching, this helps pruning most of the search space.

- Cycle through the dimensions for each level
- Call this cut-dim (cutting dimension)
- Node in tree contains $P = (x, y)$
- So, to find a point, only need to compare the cutting dimension.

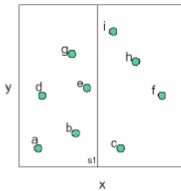
- If there is one point, just form a leaf node
- Otherwise divide the points in half along the cutting axis
 - Find the axis with the widest spread
 - divide in alternative/round robin fashion
- recursively build kdtrees from each half
- Complexity $dn \log n$

Insert

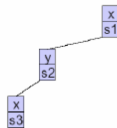
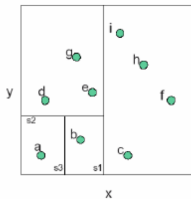
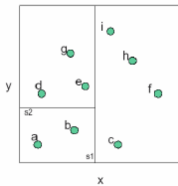
9n 0 h1 12



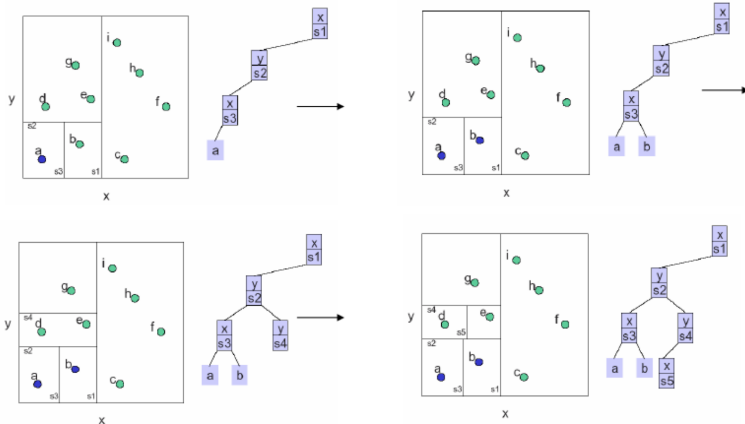
divide perpendicular to the widest spread.



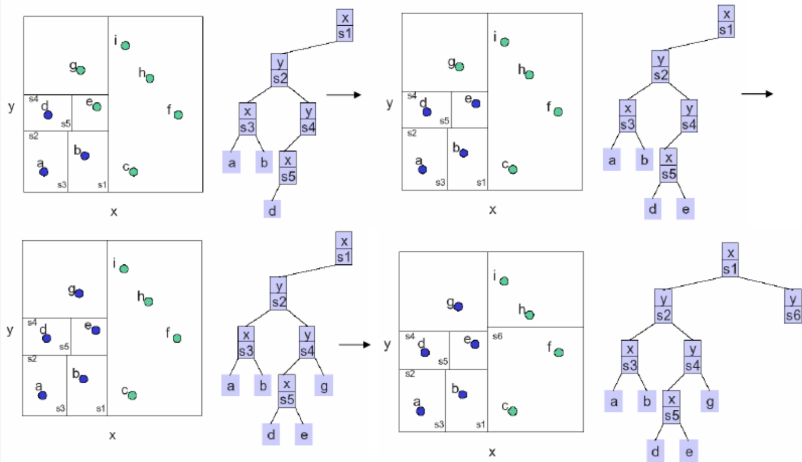
x
s1



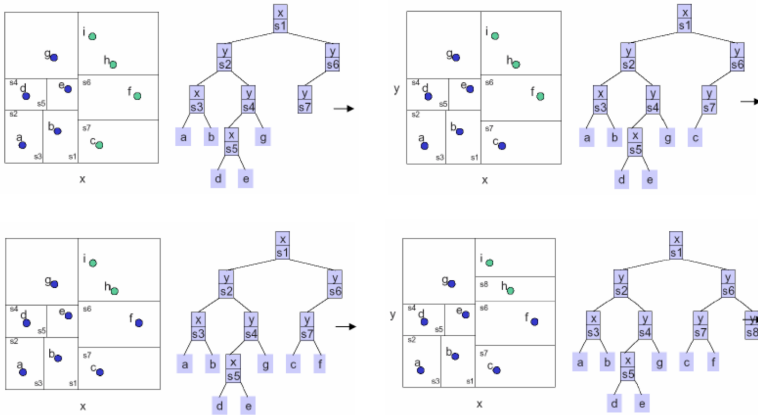
Insert



Insert



Insert



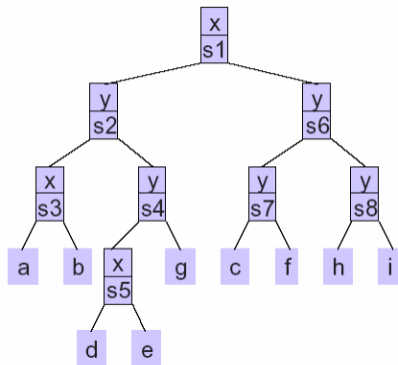
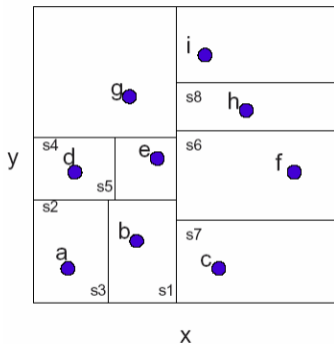
Find point with the smallest element in dimension a

- If cutdim at current node equals a ,
 - the min cannot be in the right subtree
 - recurse on the left subtree

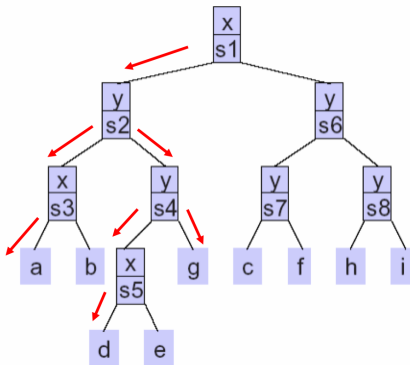
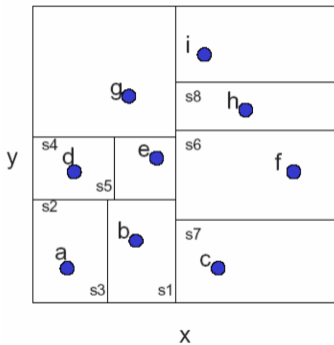
Base case: if there are no left children, stop and return current point.

- Otherwise
 - the min could be in either
 - recurse on both left and right subtrees

Find point with the smallest element in dimension y



Find point with the smallest element in dimension y



Nearest neighbor queries

- Given point Q , find the closest point R
- Have to be careful, because its possible that two points are far away in the tree but close in the Eucidean space.
- For each node store a bounding box

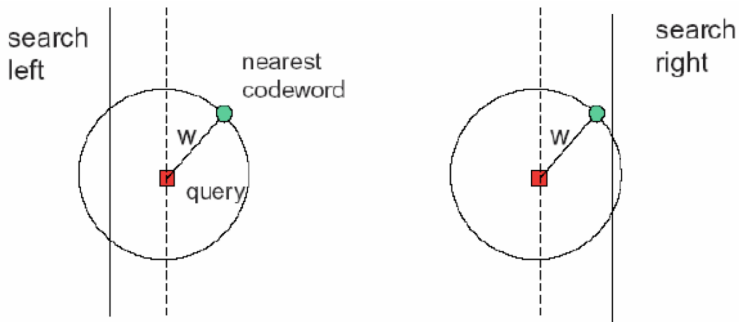
Nearest neighbor queries

- Given point Q , find the closest point R
- Have to be careful, because its possible that two points are far away in the tree but close in the Eucidean space.
- For each node store a bounding box
- Remember the closest point to Q seen so far (call this R')

Nearest neighbor queries

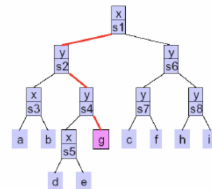
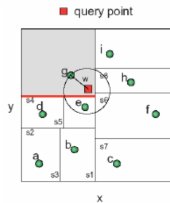
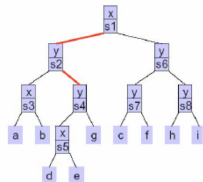
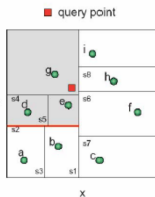
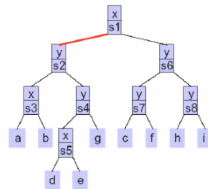
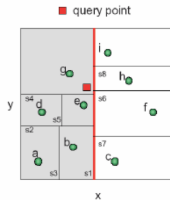
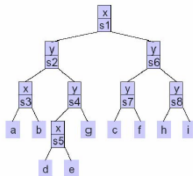
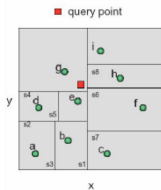
- Given point Q , find the closest point R
- Have to be careful, because its possible that two points are far away in the tree but close in the Eucidean space.
- For each node store a bounding box
- Remember the closest point to Q seen so far (call this R')
- Prune subtrees where bounding boxes cannot contain R'

Nearest neighbor queries

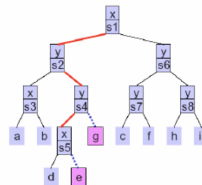
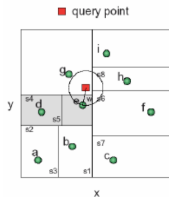
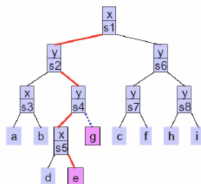
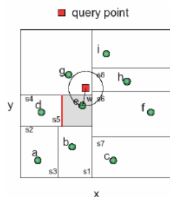
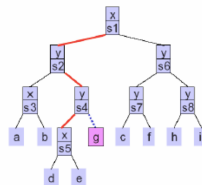
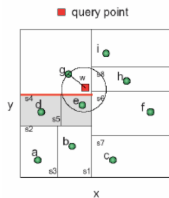
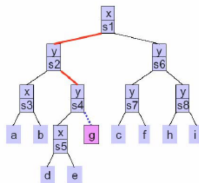
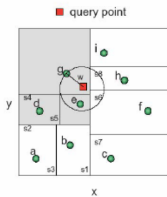


- If circle overlaps with left subtree, search left subtree
- If circle overlaps with right subtree search right subtree
- Has been shown to work in about $O(\log n)$ time.

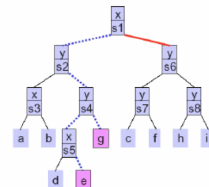
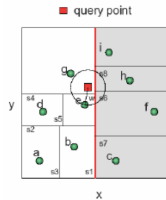
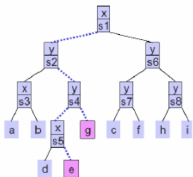
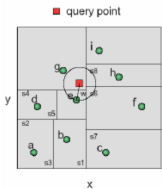
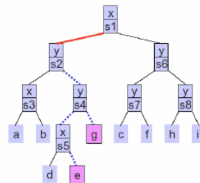
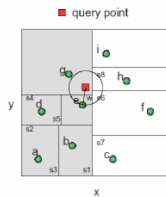
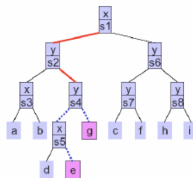
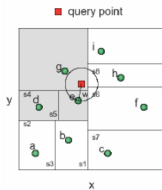
NN search



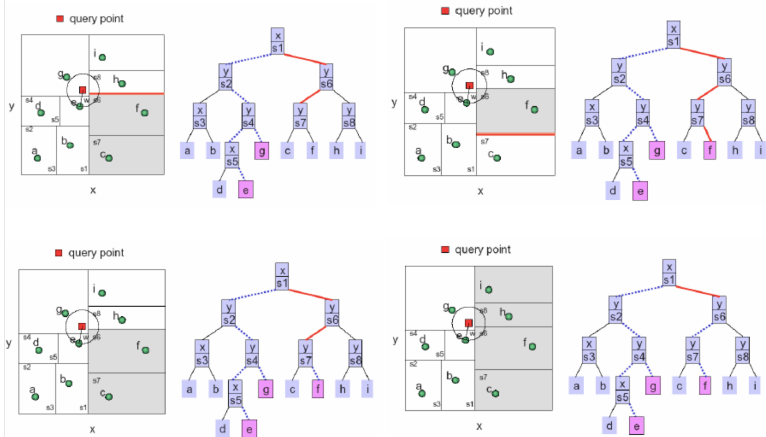
NN search



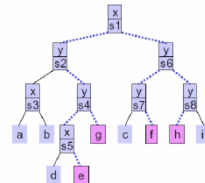
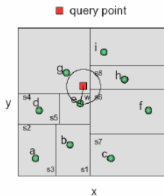
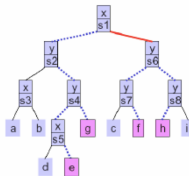
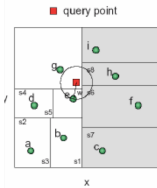
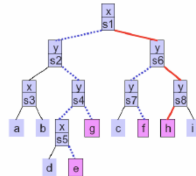
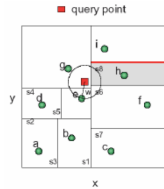
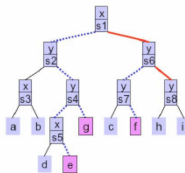
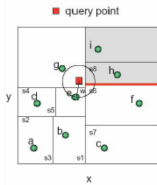
NN search



NN search



NN search



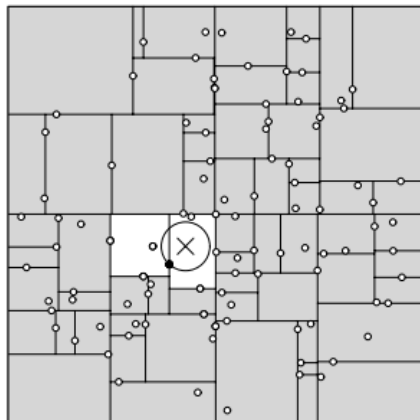


Figure 6.5

Generally during a nearest neighbour search only a few leaf nodes need to be inspected.

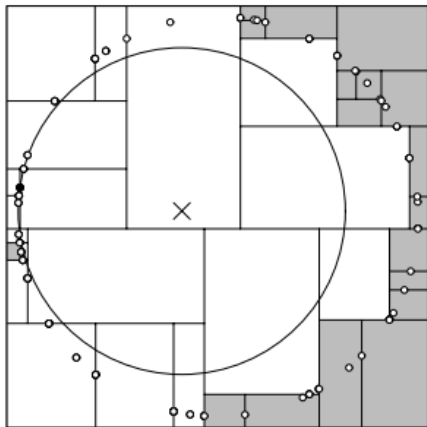


Figure 6.6

A bad distribution which forces almost all nodes to be inspected.

Timing vs tree size

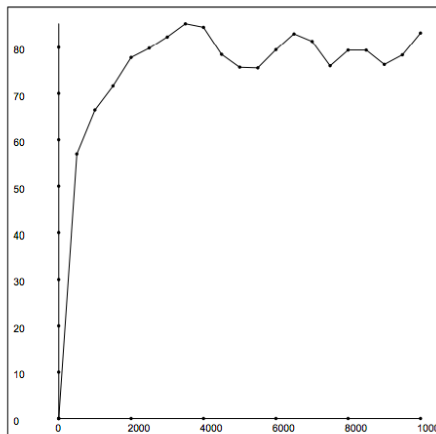


Figure 6.8

Number of inspections
against kd-tree size for an
eight-dimensional tree with
an eight-dimensional un-
derlying distribution.

Timing vs dimensions

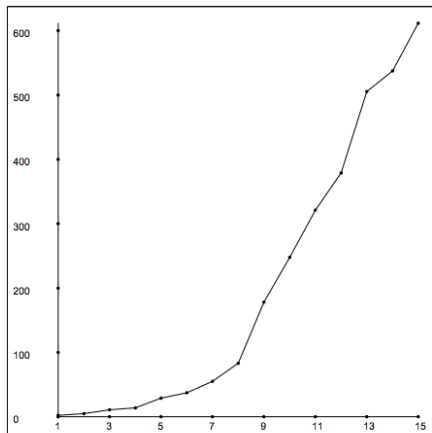
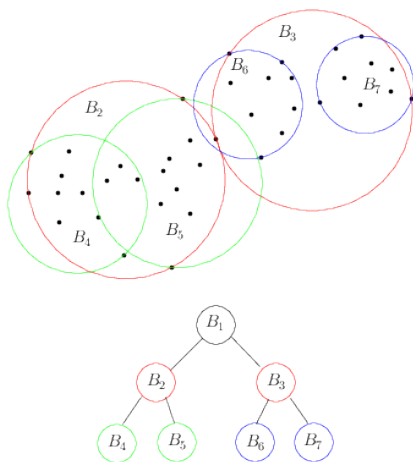


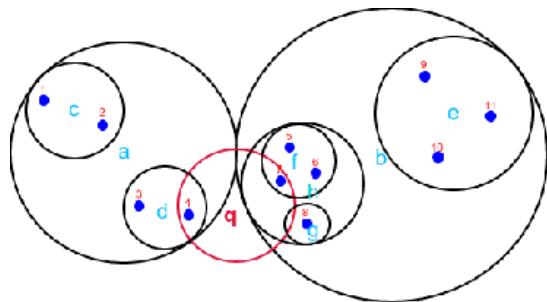
Figure 6.9

Number of inspections graphed against tree dimension. In these experiments the points had an underlying distribution with the same dimensionality as the tree.

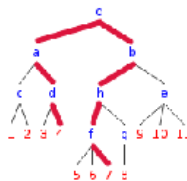
Ball trees



Ball tree search



(a)



(b)

Acknowledgment

- The kdtrees animations were borrowed from
 - Thinh Nguyen's slides
 - Carl Kingsford's slides
- Andrew moore's tutorial