

# 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

- 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.

- Say you have some algorithm to decide which dimension to split on and which value to split on (call this cut-val).
- Call this cut-dim (cutting dimension)
- Node in tree is described by (cut-dim, cut-val)
- 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





y  $\begin{array}{c|c} & i & & \\ \hline & & & \\ & & & \\ \hline & & & \\ & & & \\ \end{array}$ 





5



6







- 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 x



# Find point with the smallest element in dimension x



- $\bullet\,$  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

- $\bullet\,$  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')

- $\bullet\,$  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.





















y







#### 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





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

# Timing vs dimensions





Number of inspections graphed against tree dimension. In these experiments the points had an underlying distribution with the same dimensionality as the tree. • Algorithm 1: find nearest neighbor with above approach, report label

- Algorithm 1: find nearest neighbor with above approach, report label
  - Backtracking still takes time
- Algorithm 2: do approximate nearest neighbor, dont backtrack, report label
  - Still need to search brute force in the leaf node

- Algorithm 1: find nearest neighbor with above approach, report label
  - Backtracking still takes time
- Algorithm 2: do approximate nearest neighbor, dont backtrack, report label
  - Still need to search brute force in the leaf node
- Algorithm 3: no brute force search. Can you tell me what to do?

### **Curse of dimensionality**

- What happens when you have high dimensional data?
- How about when data lies in a lower dimensional manifold?



Figure 1: A curled plane: the swiss roll.

- How to fix this?
  - Projections?

• Pick direction of most variability

- Pick direction of most variability
- PCA?

- Pick direction of most variability
- PCA?
- Fast algorithm:
  - Pick point at random, call it X
  - Pick point far away from X, call it Y
  - Pick point far away from Y, call it Z
  - Y Z gives you a good proxy for the direction of most variability

• Pick direction of most variability

- Pick direction of most variability
- Project all points on that direction
- Split by median
- In each split
  - Compute center
  - Compute distance of center to furthest point
  - Continue recursively

# Ball trees



- Traverse tree in depth first order
- Check if distance of pointx to current node B is smaller than distance to current nearest neighbor (call this  $n_x$ )
  - If no, move on.
  - If yes, and *B* is a leaf node, then find nearest neighbor in *B*, and update the nearest neighbor *n*<sub>x</sub>
  - If yes, and *B* is an internal node, recursively search both children of node *B*. Search child whose center is closest, first.





(b)

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