

## Scaling up k-NN search using KD-trees

### Complexity of brute force search

Given a query point, scan through each point

- $O(N)$  distance computations per 1-NN query!
- $O(N \log k)$  per k-NN query!

What if  $N$  is huge??  
(and many queries)

### KD-tree representation

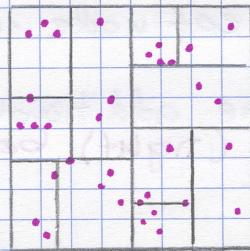
#### KD-trees

Structured organization of documents

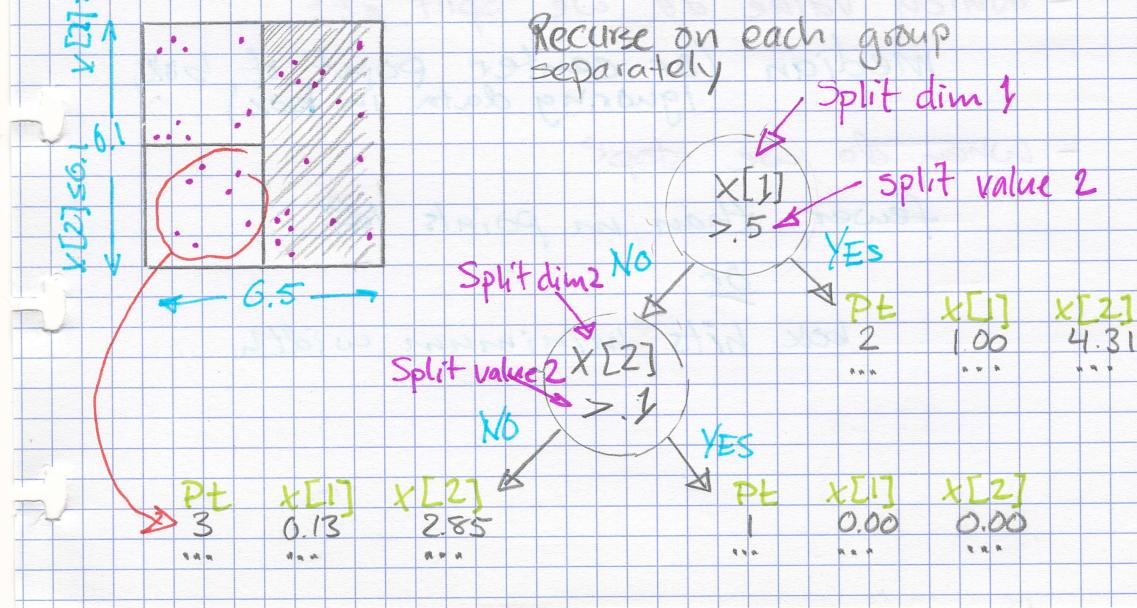
- Recursively partitions points into axis aligned boxes

Enables more efficient pruning of search space

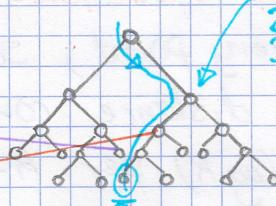
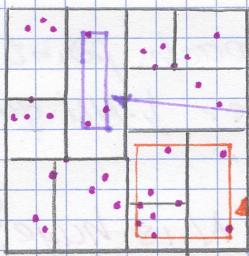
Works "well" in "low-medium" dimensions  
- we'll get back to this -



### KD-tree construction



## KD tree construction



Store:

1. split dim
2. Split value
3. bounding box that is as small as possible while containing points

points here satisfy all conditions to down the tree to this leaf

Continue splitting points at each set

- Creates a binary tree structure

Each leaf node contains a list of points

Keep one additional piece of info at each node:

- #3 - The (tight) bounds of points at or below node

## KD-tree construction choices

Use heuristics to make splitting decisions:

- Which dimension do we split along? widest (or alternate)

- Which value do we split at?

median (or center point of box, ignoring data in box)

- When do we stop?

fewer than n points left

or

box hits minimum width

## NN search with KD-trees

Nearest neighbor with KD-trees



query point  
 $x[1] > 0.5$

distance to NN found so far  
distance to NN is smaller than before

Does nearest neighbor have to dive at leaf node containing query point?

$x[2] > v$

Update distance bound when new nearest neighbor is found

Traverse tree looking for nearest neighbor to query point

1. Start by exploring leaf node containing query point
2. Compute distance to each other point at leaf node
3. Backtrack and try other branch at each node + visited

Use distance bound and bounding box of each node to prune parts of tree that cannot include nearest neighbor

## Complexity of NN search with KD-trees

Complexity

For (nearly) balanced, binary trees...

- Construction
  - Size:  $2N - 1$  nodes if 1 data point at each leaf  $\rightarrow O(N)$
  - Depth:  $O(\log N)$
  - Median + Send points left/right:  $O(N)$  at every level of tree
  - Construction time:  $O(N \log N)$
- 1-NN query
  - Traverse down tree to starting point:  $O(\log N)$
  - Maximum backtrack and traverse:  $O(N)$  in worst case
  - Complexity range:  $O(\log N) \rightarrow O(N)$

Under some assumptions on distribution of points, we get  $O(\log N)$  but exponential in d

complexity for N queries

- Ask for nearest neighbor to each doc

N queries

- Brute force 1-NN:

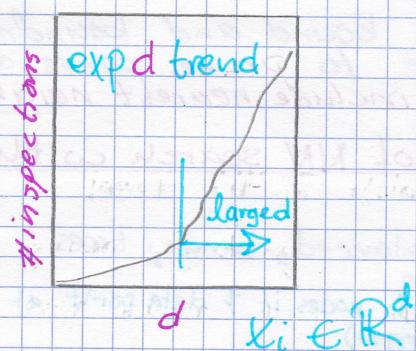
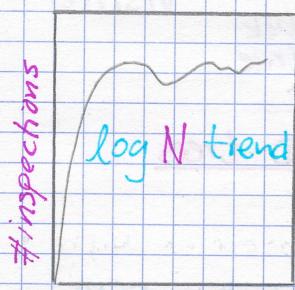
$$O(N^2)$$

- kD-trees:

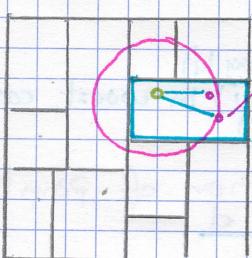
$$O(N \log N) \rightarrow O(N^2)$$

potentially  
very large  
savings for  
large N!

Inspections vs. N and d



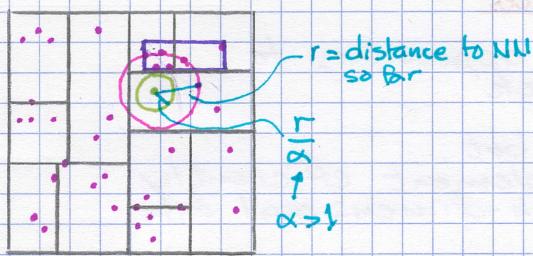
k-NN with KD-trees



distance to 2<sup>nd</sup> nearest neighbor  
in 2-NN example

Exactly the same algorithm, but  
maintain distance to furthest of  
current k nearest neighbors

## Approximate k-NN search using KD-trees



Saves lots of  
search time at  
little cost in  
quality of NN!

Before: Prune when distance to bounding box  $> r$

Now: Prune when distance to bounding box  $> r/\alpha$

Prunes more than allowed, but can **guarantee** that if we return a neighbor at distance  $r$ , then there is no neighbor closer than  $r/\alpha$

Bound loose... In practice, often closer to optimal

### Closing remarks on KD-trees

Tons of variants of kd-trees

- On construction of trees (heuristics for splitting, stopping, representing branches...)
- Other representational data structures for fast NN Search (e.g., ball trees)

### Nearest neighbor search

- Distance metric and data representation crucial to answer returned

For both, high-dim spaces are hard!

- Number of kd-tree search can be exponential in dimension
  - o Rule of thumb  $N \gg 2^d$  Typically useless for large  $d$ .
- Distances sensitive to irrelevant features
  - o Most dimensions are just noise  $\rightarrow$  everything is far away
  - o Need technique to learn which features are important to given task