



## Exhaustive k-NN Algorithm

exhaustive k-NN algo: given query pt.  $q$ ,

- 1) Scan through all  $n$  sample pts and compute squared distances to  $q$
- 2) Maintain a max-heap w/  $k$ -shortest distances so far
  - remove top pt. and insert new pt. when we discover a sample pt. closer to  $q$  than top of heap

\* heap not necessary but more efficient for large  $k$  values

- $O(n)$  training time for algorithm
  - all computation done during query
- $O(nd + n \log k)$  query time
  - $O(nd + k \log n \log k)$  expected query time for random pt. order
    - ↳ slightly better runtime  $\rightarrow$  only  $O(k \log n)$  expected heap operations occur
      - Not recommended in practice Since you'd lose more from cache misses than you'll gain from fewer heap operations

(Can an algorithm preprocess training pts. to get sublinear query time?)

- Voronoi diagrams for 2-5 dimensions
- k-d trees better for  $d \rightarrow 30$
- exhaustive k-nn best for  $d > 30$ 
  - PCA or random projection also work
  - locality-sensitive hashing works but not widely used

Cont. below

## Voronoi Diagrams

Let  $X$  be a pt. set

Then the Voronoi cell of  $\vec{w} \in X$  is  $\text{Var}(\vec{w}) = \left\{ \vec{p} \in \mathbb{R}^d \mid \forall \vec{v} \in X, \|\vec{p} - \vec{w}\|_2 \leq \|\vec{p} - \vec{v}\|_2 \right\}$

always convex polyhedron or polytope

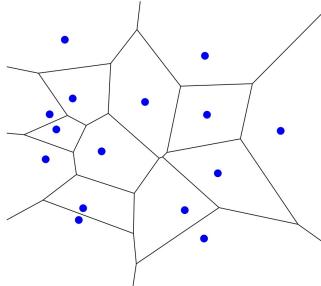


Figure 25.1: Example Voronoi diagram for 15 points

- Size of Voronoi diagram (# Vertices) is  $O(n^{[d/2]})$
- $d = \text{dimension of pts.}$
- tight Upper bound when  $d$  is a small constant
- as  $d$  grows, tightest asymptotic bound is smaller than  $O(n^d)$ , but Complexity still grows exponentially w/  $d$

We can use Voronoi diagrams for querying:

- given query pt.  $\vec{q} \in \mathbb{R}^d$ ,
  - find  $\vec{v} \in X$  for which  $\vec{q} \in \text{Var}(\vec{v})$
  - how data structure needed for this to be efficient
- In 2D, takes  $O(n \log n)$  time to get the Voronoi diagram & trapezoidal map for pt. location
  - query then takes  $O(\log n)$  after map computed
  - much better than linear query.
- In  $d$  dimensions, use a **binary space partition tree** (BSP tree) for pt. location
  - reasonably fast in 3-5 dimensions

\*

Voronoi diagram only supports 1-nearest neighbor queries

- for  $k$ -nearest neighbors, use order- $k$  Voronoi diagram
  - cell for each possible  $k$ -nearest neighbors of  $\vec{q}$
  - **Nobody uses THESE**
    - $\hookrightarrow O(k^2 n)$  for order- $k$  diagram in 2D, worse in higher dimensions
    - $\hookrightarrow$  no software available to compute this
- V-diagrams...
  - exist for other distance metrics
  - good for 1-NN in 2-3D or 4-5D but k-d trees much simpler and faster for  $d > 6$

## k-d trees: decision trees for nearest neighbor search

Differences when compared to decision trees:

- choose splitting feature w/ largest width, i.e. feature  $i$  in  $\max_{i,j,k} (x_{ji} - x_{ki})$
  - nearest neighbor search doesn't care about entropy
  - Want sphere around query pt. to not intersect many boxes of decision tree instead  
→ helps if boxes nearly cubical rather than long & thin
  - cheap alternative: rotate through features:
    - Split on 1st feature at depth 1
    - Split 2nd at depth 2, etc
    - faster by  $O(d)$
  - Splitting Value is median pt. for feature  $i$ , or  $\frac{x_{ji} + x_{ki}}{2}$ 
    - median guarantees  $\lceil \log_2 n \rceil$  tree depth &  $O(nd\log n)$  tree building time
    - $O(n\log n)$  time if rotating through features
  - alternative to median: Split at box's center
    - improves aspect ratio of boxes
    - Might make tree unbalanced
  - or, median at odd depth & box center at even depth
    - $O(\log n)$  depth
  - each internal node of k-d tree stores a sample pt. that lies in node's box
    - this pt. is usually the splitting pt.
    - Some implementations have pts. only at leaves
      - preferred w/ internal pts. → when searching through tree, we often stop searching sooner.
- Also: given query  $\vec{q}$ , find pt.  $\vec{v}$  s.t.  $\|\vec{q} - \vec{v}\| \leq (1 + \epsilon) \|\vec{q} - \vec{s}\|$ ,  $\vec{s}$  is closest sample pt.
- $\epsilon = 0 \Rightarrow$  exact nearest neighbors
- $\epsilon > 0 \Rightarrow$  approx. nearest neighbors

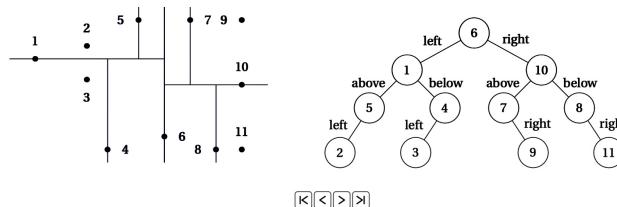


Figure 25.2: Example of a k-d tree. To step through the animation, you must use a PDF reader that supports animations, eg. Adobe Acrobat or Okular.

query algo. maintains 2 things:

- 1) nearest neighbor found so far (or k nearest)
    - successive improvements decrease w/ distance
  - 2) maintains binary min-heap of unexplored subtrees keyed by dist. from  $\vec{q}$ 
    - successive pops increase min. distance in heap
- each subtree represents an axis-aligned box
- query searches for boxes closest to  $\vec{q}$  first
    - avoid searching unnecessary boxes
  - measure distance from  $\vec{q}$  to a box
    - distance is key for subtree in heap
- Search stops when distance from  $\vec{q}$  to kth nearest neighbor found so far is less than distance from  $\vec{q}$  to nearest unexplored box  
( $\times$  times  $1 + \epsilon$  if approx. NN)

## 1-NN query algo:

```
1 Q = heap containing root node with key zero
2 r = infinity
3 while Q is not empty and  $(1 + \epsilon) * \min(Q) < r$ 
4   // pop min and update nearest neighbor
5   B = Q.removeMin()
6   w = sample point at B
7   r = min(r, dist(q, w))  Which  $P \in [1, \infty]$ 
8
9   // add child boxes to heap
10  B1, B2 = child boxes of B
11  if  $(1 + \epsilon) * \text{dist}(q, B1) < r$  then
12    Q.insert(B1, dist(q, B1)) // key for B1 is dist(q, B1)
13  if  $(1 + \epsilon) * \text{dist}(q, B2) < r$  then
14    Q.insert(B2, dist(q, B2)) // key for B2 is dist(q, B2)
```

- If k-NN, replace r w/ max-heap that holds k nearest neighbors  
(like exhaustive search algo.)

## Summary:

- Start w/ heap initialized w/ root of k-d tree
- Pop from min heap & compare distance from  $\vec{q}$  to p. in current node
  - update nearest neighbor if necessary
- look at shortest distance from  $\vec{q}$  to children boxes of current node
  - if these distances to boxes are less than current nearest neighbor, then node in box might be new nearest neighbor → add child to heap
- repeat until min. element in heap is too far away ( $\vec{q} \rightarrow$  nearest neighbor  $\vec{q} \rightarrow$  <sup>closest</sup> unexplored box)

Why  $\epsilon$ -approx.  $NN^2$ ?

→ don't need to visit  
each box

