

ONLINE MASTERS IN DATA SCIENCE

DSC 255 - MACHINE LEARNING FUNDAMENTALS

SPEEDING UP NEAREST NEIGHBOR METHODS

SANJOY DASGUPTA, PROFESSOR

UC San Diego

COMPUTER SCIENCE & ENGINEERING
HALICIOĞLU DATA SCIENCE INSTITUTE

Naive search takes time $O(n)$ for training set of size n : slow!

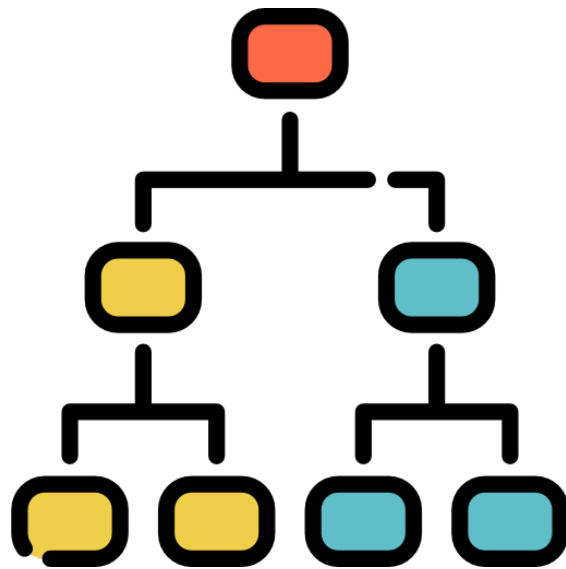


Naive search takes time $O(n)$ for training set of size n : slow!

There are data structures for speeding up nearest neighbor search, like:

1. Locality sensitive hashing
2. Ball trees
3. K -d trees

These are part of standard Python libraries for NN, and help a lot.

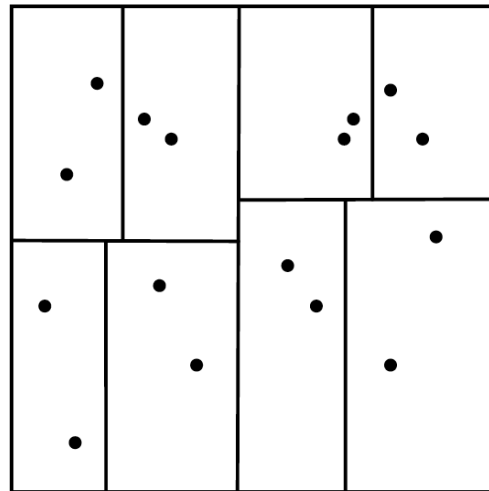


Example: k -d trees for NN search

hierarchical, rectilinear spatial partition.

For data set $\mathcal{S} \subset \mathbb{R}^d$:

1. Pick a coordinate $1 \leq i \leq d$.
2. Compute $\mathbf{v} = \text{median}(\{\mathbf{x}_i : \mathbf{x} \in \mathcal{S}\})$.
3. Split \mathcal{S} into two halves:
 - $\mathcal{S}_L = \{\mathbf{x} \in \mathcal{S} : x_i < v\}$
 - $\mathcal{S}_R = \{\mathbf{x} \in \mathcal{S} : x_i \geq v\}$
4. Recurse on $\mathcal{S}_L, \mathcal{S}_R$

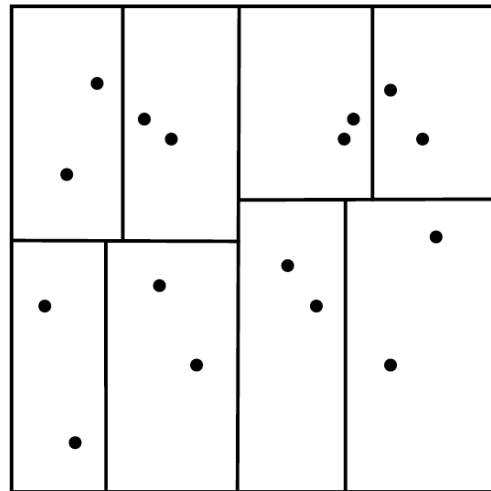


Example: k -d trees for NN search

hierarchical, rectilinear spatial partition.

For data set $\mathcal{S} \subset \mathbb{R}^d$:

1. Pick a coordinate $1 \leq i \leq d$.
2. Compute $\mathbf{v} = \text{median}(\{\mathbf{x}_i : \mathbf{x} \in \mathcal{S}\})$.
3. Split \mathcal{S} into two halves:
 - $\mathcal{S}_L = \{\mathbf{x} \in \mathcal{S} : x_i < v\}$
 - $\mathcal{S}_R = \{\mathbf{x} \in \mathcal{S} : x_i \geq v\}$
4. Recurse on $\mathcal{S}_L, \mathcal{S}_R$



Two types of search, given a query $\mathbf{q} \in \mathbb{R}^d$:

1. **Defeatist search:** Route \mathbf{q} to a leaf cell and return the NN in that cell.
This might not be the true NN.
2. **Comprehensive search:** Grow the search region to other cells that cannot be ruled out using the triangle inequality.