

COMS 4771 Lecture 3

1. Nearest neighbor classification.

NEAREST NEIGHBOR CLASSIFICATION

THE OPTIMAL CLASSIFIER

Let $(X, Y) \sim P$.

The classifier $f: \mathcal{X} \rightarrow \mathcal{Y}$ with the smallest prediction error

$$\text{err}(f) = \Pr[f(X) \neq Y]$$

is the **Bayes classifier**

$$f^*(x) = \arg \max_{y \in \mathcal{Y}} \Pr[Y = y \mid X = x].$$

THE OPTIMAL CLASSIFIER

Let $(X, Y) \sim P$.

The classifier $f: \mathcal{X} \rightarrow \mathcal{Y}$ with the smallest prediction error

$$\text{err}(f) = \Pr[f(X) \neq Y]$$

is the **Bayes classifier**

$$f^*(x) = \arg \max_{y \in \mathcal{Y}} \Pr[Y = y \mid X = x].$$

Of course, we don't know P , and hence can't generally get a handle of f^* .

THE OPTIMAL CLASSIFIER

Let $(X, Y) \sim P$.

The classifier $f: \mathcal{X} \rightarrow \mathcal{Y}$ with the smallest prediction error

$$\text{err}(f) = \Pr[f(X) \neq Y]$$

is the **Bayes classifier**

$$f^*(x) = \arg \max_{y \in \mathcal{Y}} \Pr[Y = y \mid X = x].$$

Of course, we don't know P , and hence can't generally get a handle of f^* .

What can we do?

- **Last lecture:** use “generative” models to approximate $\Pr[Y = y \mid X = x]$.

THE OPTIMAL CLASSIFIER

Let $(X, Y) \sim P$.

The classifier $f: \mathcal{X} \rightarrow \mathcal{Y}$ with the smallest prediction error

$$\text{err}(f) = \Pr[f(X) \neq Y]$$

is the **Bayes classifier**

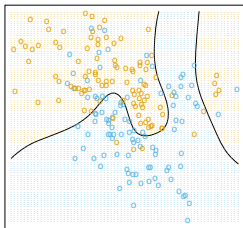
$$f^*(x) = \arg \max_{y \in \mathcal{Y}} \Pr[Y = y \mid X = x].$$

Of course, we don't know P , and hence can't generally get a handle of f^* .

What can we do?

- ▶ **Last lecture:** use “generative” models to approximate $\Pr[Y = y \mid X = x]$.
- ▶ **This lecture:** directly approximate the **decision boundaries** of f^* .

discriminative model



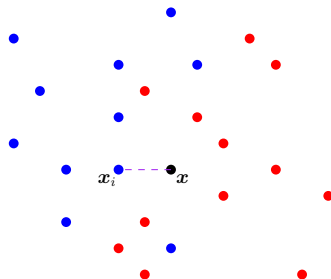
Text

NEAREST NEIGHBOR (NN) CLASSIFIER

Given **training data** $(x_1, y_1), (x_2, y_2), \dots, (x_n, y_n) \in \mathcal{X} \times \mathcal{Y}$, construct $\hat{f}: \mathcal{X} \rightarrow \mathcal{Y}$ as follows:

On input x ,

1. Let x_i be the point among x_1, x_2, \dots, x_n that is closest to x .
2. Return y_i .

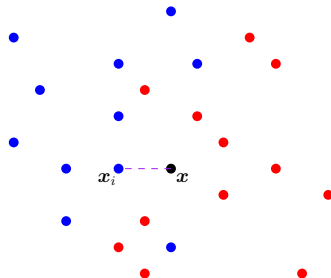


NEAREST NEIGHBOR (NN) CLASSIFIER

Given **training data** $(x_1, y_1), (x_2, y_2), \dots, (x_n, y_n) \in \mathcal{X} \times \mathcal{Y}$, construct $\hat{f}: \mathcal{X} \rightarrow \mathcal{Y}$ as follows:

On input x ,

1. Let x_i be the point among x_1, x_2, \dots, x_n that is closest to x .
2. Return y_i .



Question: how should we measure distance between points in \mathcal{X} ?

DISTANCES

A default choice of distance for data in \mathbb{R}^d :

Euclidean (ℓ_2) distance : $\|\mathbf{u} - \mathbf{v}\|_2 := \sqrt{\sum_{i=1}^d (u_i - v_i)^2}.$

DISTANCES

A default choice of distance for data in \mathbb{R}^d :

Euclidean (ℓ_2) distance : $\|\mathbf{u} - \mathbf{v}\|_2 := \sqrt{\sum_{i=1}^d (u_i - v_i)^2}.$

But there are many other (and often better) options ...

EXAMPLE: OCR WITH NN CLASSIFIER

- ▶ **Handwritten digits data:** grayscale 28×28 images, treated as **vectors in \mathbb{R}^{784}** , with labels indicating the **digit they represent**.

0 1 2 3' 4 5 6 7 8 9

EXAMPLE: OCR WITH NN CLASSIFIER

- ▶ **Handwritten digits data:** grayscale 28×28 images, treated as **vectors** in \mathbb{R}^{784} , with labels indicating the **digit they represent**.

0 1 2 3' 4 5 6 7 8 9

- ▶ Split into training data S (60000 points) and test data T (10000 points).

EXAMPLE: OCR WITH NN CLASSIFIER

- ▶ **Handwritten digits data:** grayscale 28×28 images, treated as **vectors** in \mathbb{R}^{784} , with labels indicating the **digit they represent**.

0 1 2 3' 4 5 6 7 8 9

- ▶ Split into training data S (60000 points) and test data T (10000 points).
- ▶ **Training error:** $\text{err}(\hat{f}, S) = 0$

EXAMPLE: OCR WITH NN CLASSIFIER

- ▶ **Handwritten digits data:** grayscale 28×28 images, treated as **vectors** in \mathbb{R}^{784} , with labels indicating the **digit they represent**.

0 1 2 3 4 5 6 7 8 9

- ▶ Split into training data S (60000 points) and test data T (10000 points).
- ▶ **Training error:** $\text{err}(\hat{f}, S) = 0$
Test error: $\text{err}(\hat{f}, T) = 0.0309$
- ▶ Examples of mistakes (test point in T , nearest neighbor in S):

2 8 3 5 5 4 4 1

EXAMPLE: OCR WITH NN CLASSIFIER

- ▶ **Handwritten digits data:** grayscale 28×28 images, treated as **vectors in \mathbb{R}^{784}** , with labels indicating the **digit they represent**.

0 1 2 3' 4 5 6 7 8 9

- ▶ Split into training data S (60000 points) and test data T (10000 points).
- ▶ **Training error:** $\text{err}(\hat{f}, S) = 0$
Test error: $\text{err}(\hat{f}, T) = 0.0309$

- ▶ Examples of mistakes (test point in T , nearest neighbor in S):

2 8 3 5 5 4 4 1

- ▶ **Observation:** First mistake (correct label is '2') might've been avoided by looking at three nearest neighbors (whose labels are '8', '2', '2') ...

2 8 2 2

test point three nearest neighbors

k -NEAREST NEIGHBORS CLASSIFIER

Given **training data** $(\mathbf{x}_1, y_1), (\mathbf{x}_2, y_2), \dots, (\mathbf{x}_n, y_n) \in \mathcal{X} \times \mathcal{Y}$, construct $\hat{f}_k: \mathcal{X} \rightarrow \mathcal{Y}$ as follows:

On input \mathbf{x} ,

1. Let $\mathbf{x}_{i_1}, \mathbf{x}_{i_2}, \dots, \mathbf{x}_{i_k}$ be the k points among $\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_n$ that are closest to \mathbf{x} .
2. Return the plurality of $y_{i_1}, y_{i_2}, \dots, y_{i_k}$.

(Break ties in both steps arbitrarily.)

k -NEAREST NEIGHBORS CLASSIFIER

Given **training data** $(\mathbf{x}_1, y_1), (\mathbf{x}_2, y_2), \dots, (\mathbf{x}_n, y_n) \in \mathcal{X} \times \mathcal{Y}$, construct $\hat{f}_k: \mathcal{X} \rightarrow \mathcal{Y}$ as follows:

On input \mathbf{x} ,

1. Let $\mathbf{x}_{i_1}, \mathbf{x}_{i_2}, \dots, \mathbf{x}_{i_k}$ be the k points among $\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_n$ that are closest to \mathbf{x} .
2. Return the plurality of $y_{i_1}, y_{i_2}, \dots, y_{i_k}$.

(Break ties in both steps arbitrarily.)

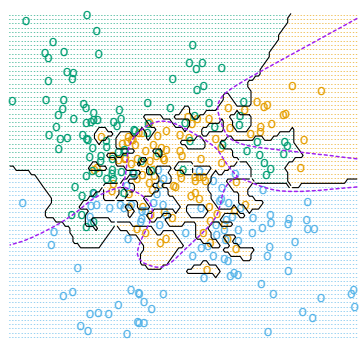
Example: OCR with k -NN classifier

k	1	3	5	7	9
$\text{err}(\hat{f}_k, T)$	0.0309	0.0295	0.0312	0.0306	0.0341

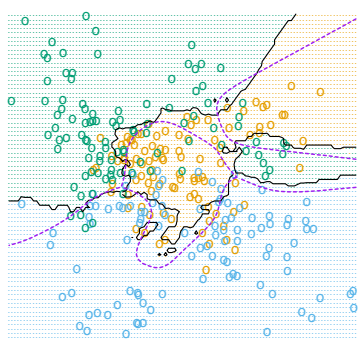
EFFECT OF k

In general:

- ▶ Smaller $k \Rightarrow$ smaller training error. ($k = 1 \Rightarrow$ has zero training error.)
- ▶ Larger $k \Rightarrow$ predictions are more “stable” due to voting.



1-NN



15-NN

Purple dotted lines: Bayes classifier's decision boundaries.

Black solid lines: k -NN's decision boundaries.

CHOOSING k

Question: how do we choose k (say, from some subset of \mathbb{N} like $\{1, 3, 5, 7, 9\}$)?

CHOOSING k

Question: how do we choose k (say, from some subset of \mathbb{N} like $\{1, 3, 5, 7, 9\}$)?

- ▶ **Minimizer of test error:** $\hat{k} := \arg \min_k \text{err}(\hat{f}_k, T)$.

CHOOSING k

Question: how do we choose k (say, from some subset of \mathbb{N} like $\{1, 3, 5, 7, 9\}$)?

- ▶ **Minimizer of test error:** $\hat{k} := \arg \min_k \text{err}(\hat{f}_k, T)$.

Caveat: $\hat{f}_{\hat{k}}$ is no longer independent of T !

CHOOSING k

Question: how do we choose k (say, from some subset of \mathbb{N} like $\{1, 3, 5, 7, 9\}$)?

- ▶ **Minimizer of test error:** $\hat{k} := \arg \min_k \text{err}(\hat{f}_k, T)$.

Caveat: $\hat{f}_{\hat{k}}$ is no longer independent of T !

\Rightarrow Test error is not an unbiased estimate of true error of $\hat{f}_{\hat{k}}$.

CHOOSING k

Question: how do we choose k (say, from some subset of \mathbb{N} like $\{1, 3, 5, 7, 9\}$)?

- ▶ **Minimizer of test error:** $\hat{k} := \arg \min_k \text{err}(\hat{f}_k, T)$.

Caveat: $\hat{f}_{\hat{k}}$ is no longer independent of T !

\implies Test error is not an unbiased estimate of true error of $\hat{f}_{\hat{k}}$.

- ▶ **Better alternatives:** For any set of labeled examples $A \subseteq \mathcal{X} \times \mathcal{Y}$, define $\hat{f}_{(A,k)}$ to be the k -NN classifier that searches for neighbors in A .

1. **Minimizer of hold-out error:** fix $H \subseteq S$,

$$\hat{k} := \arg \min_k \text{err}(\hat{f}_{(S \setminus H, k)}, H).$$

CHOOSING k

Question: how do we choose k (say, from some subset of \mathbb{N} like $\{1, 3, 5, 7, 9\}$)?

- **Minimizer of test error:** $\hat{k} := \arg \min_k \text{err}(\hat{f}_k, T)$.

Caveat: $\hat{f}_{\hat{k}}$ is no longer independent of T !

\implies Test error is not an unbiased estimate of true error of $\hat{f}_{\hat{k}}$.

- **Better alternatives:** For any set of labeled examples $A \subseteq \mathcal{X} \times \mathcal{Y}$, define $\hat{f}_{(A,k)}$ to be the k -NN classifier that searches for neighbors in A .

1. **Minimizer of hold-out error:** fix $H \subseteq S$,

$$\hat{k} := \arg \min_k \text{err}(\hat{f}_{(S \setminus H, k)}, H).$$

2. **Minimizer of leave-one-out cross-validation error:**

$$\hat{k} := \arg \min_k \frac{1}{|S|} \sum_{(x,y) \in S} \text{err}(\hat{f}_{(S \setminus \{(x,y)\}, k)}, \{(x,y)\}).$$

CHOOSING k

Question: how do we choose k (say, from some subset of \mathbb{N} like $\{1, 3, 5, 7, 9\}$)?

- **Minimizer of test error:** $\hat{k} := \arg \min_k \text{err}(\hat{f}_k, T)$.

Caveat: $\hat{f}_{\hat{k}}$ is no longer independent of T !

\implies Test error is not an unbiased estimate of true error of $\hat{f}_{\hat{k}}$.

- **Better alternatives:** For any set of labeled examples $A \subseteq \mathcal{X} \times \mathcal{Y}$, define $\hat{f}_{(A,k)}$ to be the k -NN classifier that searches for neighbors in A .

1. **Minimizer of hold-out error:** fix $H \subseteq S$,

$$\hat{k} := \arg \min_k \text{err}(\hat{f}_{(S \setminus H, k)}, H).$$

2. **Minimizer of leave-one-out cross-validation error:**

$$\hat{k} := \arg \min_k \frac{1}{|S|} \sum_{(x,y) \in S} \text{err}(\hat{f}_{(S \setminus \{(x,y)\}, k)}, \{(x,y)\}).$$

\implies Test error is an unbiased estimate of true error of $\hat{f}_{\hat{k}}$.

CHOOSING k

Question: how do we choose k (say, from some subset of \mathbb{N} like $\{1, 3, 5, 7, 9\}$)?

- **Minimizer of test error:** $\hat{k} := \arg \min_k \text{err}(\hat{f}_k, T)$.

Caveat: $\hat{f}_{\hat{k}}$ is no longer independent of T !

\implies Test error is not an unbiased estimate of true error of $\hat{f}_{\hat{k}}$.

- **Better alternatives:** For any set of labeled examples $A \subseteq \mathcal{X} \times \mathcal{Y}$, define $\hat{f}_{(A,k)}$ to be the k -NN classifier that searches for neighbors in A .

1. **Minimizer of hold-out error:** fix $H \subseteq S$,

the difference is at
k's value!!!

$$\hat{k} := \arg \min_k \text{err}(\hat{f}_{(S \setminus H, k)}, H).$$

hold one
point out for
test

2. **Minimizer of leave-one-out cross-validation error:**

$$\hat{k} := \arg \min_k \frac{1}{|S|} \sum_{(x,y) \in S} \text{err}(\hat{f}_{(S \setminus \{(x,y)\}, k)}, \{(x,y)\}).$$

\implies Test error is an unbiased estimate of true error of $\hat{f}_{\hat{k}}$.

More on this later in the course.

CONSISTENCY OF k -NN

Say a learning algorithm is **consistent** if

$$\lim_{n \rightarrow \infty} \mathbb{E} \left[\text{error of learned classifier with training sample size } n \right] = \text{err}(f^*).$$

CONSISTENCY OF k -NN

Say a learning algorithm is **consistent** if

the best classifier!

$$\lim_{n \rightarrow \infty} \mathbb{E} \left[\text{error of learned classifier with training sample size } n \right] = \text{err}(f^*).$$

k -NN is consistent provided that $k := k_n$ is chosen as an increasing but sublinear function of n :

$$\lim_{n \rightarrow \infty} k_n = \infty, \quad \lim_{n \rightarrow \infty} \frac{k_n}{n} = 0$$

(some other mild conditions might also have to hold).

as the size of training sample grows up. the error of the learned classifier getting closer to $\text{err}(f)$.

CONSISTENCY OF k -NN

Say a learning algorithm is **consistent** if

$$\lim_{n \rightarrow \infty} \mathbb{E} \left[\text{error of learned classifier with training sample size } n \right] = \text{err}(f^*).$$

k -NN is consistent provided that $k := k_n$ is chosen as an increasing but sublinear function of n :

$$\lim_{n \rightarrow \infty} k_n = \infty, \quad \lim_{n \rightarrow \infty} \frac{k_n}{n} = 0$$

(some other mild conditions might also have to hold).

1-NN is not consistent unless $\text{err}(f^*) = 0$, although

$$\lim_{n \rightarrow \infty} \mathbb{E} \left[\text{err}(\hat{f}_1) \right] \leq 2 \text{err}(f^*) \cdot \left(1 - \frac{K}{2(K-1)} \text{err}(f^*) \right).$$

NEAREST NEIGHBOR SEARCH

- ▶ Naïve implementation of NN classifiers uses n distance computations to compute $\hat{f}_k(\mathbf{x})$ for any test point $\mathbf{x} \in \mathcal{X}$.

NEAREST NEIGHBOR SEARCH

- ▶ Naïve implementation of NN classifiers uses n distance computations to compute $\hat{f}_k(\mathbf{x})$ for any test point $\mathbf{x} \in \mathcal{X}$.
 - ▶ If using Euclidean distance in \mathbb{R}^d , then each distance computation is $O(d)$ operations.
- $\implies O(dn)$ operations per test point.

NEAREST NEIGHBOR SEARCH

- ▶ Naïve implementation of NN classifiers uses n distance computations to compute $\hat{f}_k(x)$ for any test point $x \in \mathcal{X}$.
 - ▶ If using Euclidean distance in \mathbb{R}^d , then each distance computation is $O(d)$ operations.

$\implies O(dn)$ operations per test point.

- ▶ **Alternatives:**
 1. Settle for an approximate nearest neighbor using *locality sensitive hash functions*.
 2. Store the n training data in a geometric data structure that permits fast NN queries.

LOCALITY SENSITIVE HASH FUNCTIONS

(Informally:) A family \mathcal{H} of hash functions from \mathbb{R}^d to \mathbb{Z} is a **locality-sensitive hash family** if

- For any points $\mathbf{a}, \mathbf{b}, \mathbf{c} \in \mathcal{X}$ with $\|\mathbf{a} - \mathbf{b}\|_2 \ll \|\mathbf{a} - \mathbf{c}\|_2$,

$$|\{h \in \mathcal{H} : h(\mathbf{a}) = h(\mathbf{b})\}| \gg |\{h \in \mathcal{H} : h(\mathbf{a}) = h(\mathbf{c})\}|.$$

LOCALITY SENSITIVE HASH FUNCTIONS

(Informally:) A family \mathcal{H} of hash functions from \mathbb{R}^d to \mathbb{Z} is a **locality-sensitive hash family** if

- For any points $\mathbf{a}, \mathbf{b}, \mathbf{c} \in \mathcal{X}$ with $\|\mathbf{a} - \mathbf{b}\|_2 \ll \|\mathbf{a} - \mathbf{c}\|_2$,

$$|\{h \in \mathcal{H} : h(\mathbf{a}) = h(\mathbf{b})\}| \gg |\{h \in \mathcal{H} : h(\mathbf{a}) = h(\mathbf{c})\}|.$$

It turns out there are such hash families!

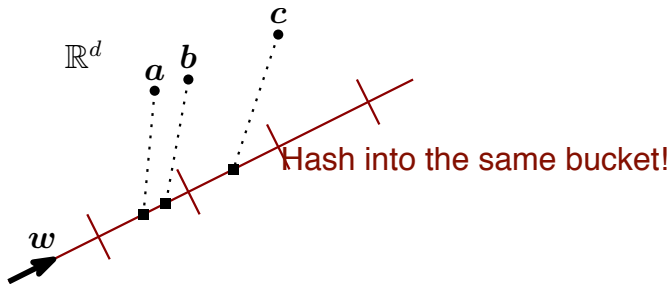
A LOCALITY SENSITIVE HASH FAMILY

A LSH family based on projections to one-dimensional subspaces

For $w \in \mathbb{R}^d$ with $\|w\|_2 = 1$, $r \in \{2^i : i \in \mathbb{Z}\}$, $s \in [0, r]$:

$$h_{w,r,s}(x) := \left\lfloor \frac{w^\top x + s}{r} \right\rfloor.$$

- ▶ w determines the one-dimensional subspace,
- ▶ r determines a distance resolution, and
- ▶ s determines a shift of the bucket boundaries.



LOCALITY SENSITIVE HASH FUNCTIONS

(Informally:) A family \mathcal{H} of hash functions from \mathbb{R}^d to \mathbb{Z} is a **locality-sensitive hash family** if

- ▶ For any points $\mathbf{a}, \mathbf{b}, \mathbf{c} \in \mathbb{R}^d$ with $\|\mathbf{a} - \mathbf{b}\|_2 \ll \|\mathbf{a} - \mathbf{c}\|_2$,

$$|\{h \in \mathcal{H} : h(\mathbf{a}) = h(\mathbf{b})\}| \gg |\{h \in \mathcal{H} : h(\mathbf{a}) = h(\mathbf{c})\}|.$$

Procedure:

- ▶ Select a hash function $h \in \mathcal{H}$ at random.
(In practice, some parameters of the hash function, like r and s , may be tuned via hold-out or cross-validation.)
- ▶ Create pointer from buckets $j \in \mathbb{N}$ to points $\mathbf{x} \in S$ such that $h(\mathbf{x}) = j$.
- ▶ Given test point \mathbf{x} , search bucket $h(\mathbf{x})$ for nearest neighbor.
(The bucket will generally contain far fewer than n points.)

LOCALITY SENSITIVE HASH FUNCTIONS

(Informally:) A family \mathcal{H} of hash functions from \mathbb{R}^d to \mathbb{Z} is a **locality-sensitive hash family** if

- ▶ For any points $\mathbf{a}, \mathbf{b}, \mathbf{c} \in \mathbb{R}^d$ with $\|\mathbf{a} - \mathbf{b}\|_2 \ll \|\mathbf{a} - \mathbf{c}\|_2$,

$$|\{h \in \mathcal{H} : h(\mathbf{a}) = h(\mathbf{b})\}| \gg |\{h \in \mathcal{H} : h(\mathbf{a}) = h(\mathbf{c})\}|.$$

Procedure:

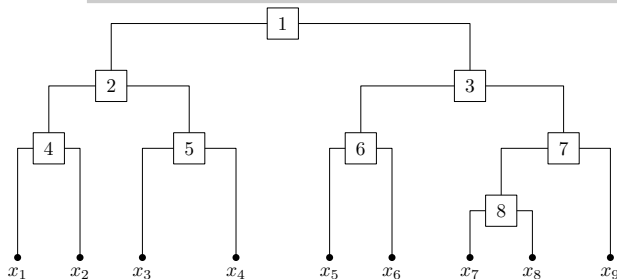
- ▶ Select a hash function $h \in \mathcal{H}$ at random.
(In practice, some parameters of the hash function, like r and s , may be tuned via hold-out or cross-validation.)
- ▶ Create pointer from buckets $j \in \mathbb{N}$ to points $\mathbf{x} \in S$ such that $h(\mathbf{x}) = j$.
- ▶ Given test point \mathbf{x} , search bucket $h(\mathbf{x})$ for nearest neighbor.
(The bucket will generally contain far fewer than n points.)

Do this with several hash functions from \mathcal{H} to boost the chances that you find close neighbors.

TREE STRUCTURES FOR ONE-DIMENSIONAL DATA

A data structure for fast NN search in \mathbb{R}^1

Sort training data so that $x_1 \leq x_2 \leq \dots \leq x_n$, then construct binary tree:

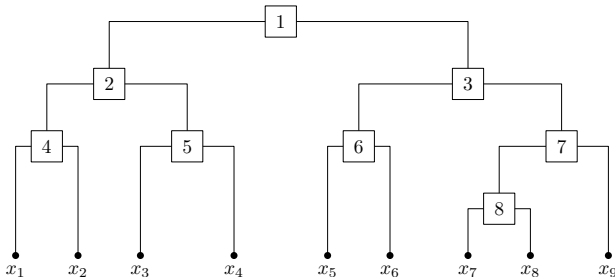


TREE STRUCTURES FOR ONE-DIMENSIONAL DATA

Text

A data structure for fast NN search in \mathbb{R}^1

Sort training data so that $x_1 \leq x_2 \leq \dots \leq x_n$, then construct binary tree:

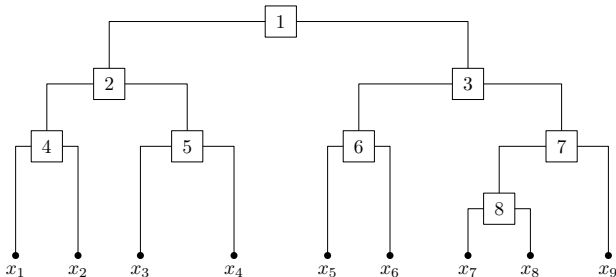


With each tree node, remember **midpoint** between rightmost point in left child, and leftmost point in right child.

TREE STRUCTURES FOR ONE-DIMENSIONAL DATA

A data structure for fast NN search in \mathbb{R}^1

Sort training data so that $x_1 \leq x_2 \leq \dots \leq x_n$, then construct binary tree:

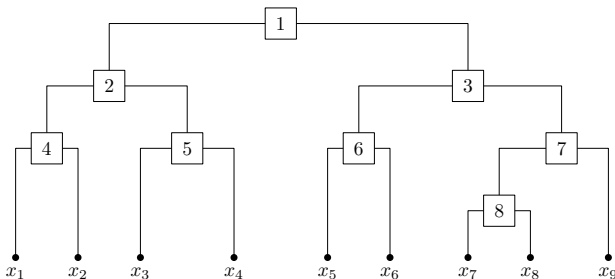


With each tree node, remember **midpoint** between rightmost point in left child, and leftmost point in right child. **This permits very efficient NN search.**

TREE STRUCTURES FOR ONE-DIMENSIONAL DATA

A data structure for fast NN search in \mathbb{R}^1

Sort training data so that $x_1 \leq x_2 \leq \dots \leq x_n$, then construct binary tree:



With each tree node, remember **midpoint** between rightmost point in left child, and leftmost point in right child. **This permits very efficient NN search.**

If tree is (approximately) balanced, then $O(\log(n))$ time to find NN!

TREE STRUCTURES FOR MULTI-DIMENSIONAL DATA

A data structure for fast NN search in \mathbb{R}^d , $d > 1$

Many options, but a popular one is the **K-D tree**.

TREE STRUCTURES FOR MULTI-DIMENSIONAL DATA

A data structure for fast NN search in \mathbb{R}^d , $d > 1$

Many options, but a popular one is the **K-D tree**.

Construction procedure

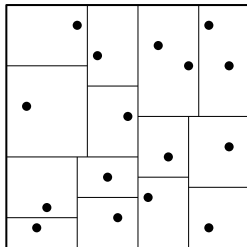
Given points $S \subset \mathbb{R}^d$:

1. Pick a coordinate $j \in \{1, 2, \dots, d\}$.
2. Let m be the median of $\{x_j : \mathbf{x} \in S\}$.
3. Split points into halves:

$$L := \{\mathbf{x} \in S : x_j < m\},$$

$$R := \{\mathbf{x} \in S : x_j \geq m\}.$$

4. Recurse on L and R .



Text

TREE STRUCTURES FOR MULTI-DIMENSIONAL DATA

A data structure for fast NN search in \mathbb{R}^d , $d > 1$

Many options, but a popular one is the **K-D tree**.

Construction procedure

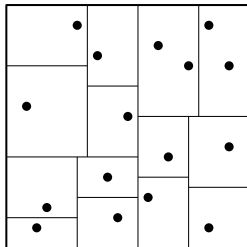
Given points $S \subset \mathbb{R}^d$:

1. Pick a coordinate $j \in \{1, 2, \dots, d\}$.
2. Let m be the median of $\{x_j : \mathbf{x} \in S\}$.
3. Split points into halves:

$$L := \{\mathbf{x} \in S : x_j < m\},$$

$$R := \{\mathbf{x} \in S : x_j \geq m\}.$$

4. Recurse on L and R .



Easy to lookup points in S (in $O(\log(n))$ time), but how about NN search?

TREE STRUCTURES FOR MULTI-DIMENSIONAL DATA

A data structure for fast NN search in \mathbb{R}^d , $d > 1$

Many options, but a popular one is the **K-D tree**.

Construction procedure

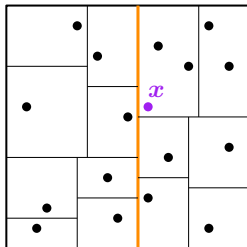
Given points $S \subset \mathbb{R}^d$:

1. Pick a coordinate $j \in \{1, 2, \dots, d\}$.
2. Let m be the median of $\{x_j : \mathbf{x} \in S\}$.
3. Split points into halves:

$$L := \{\mathbf{x} \in S : x_j < m\},$$

$$R := \{\mathbf{x} \in S : x_j \geq m\}.$$

4. Recurse on L and R .



Easy to lookup points in S (in $O(\log(n))$ time), but how about NN search?

Same $O(\log(n))$ -time routing of a test point $\mathbf{x} \in \mathbb{R}^d$ is **overly optimistic**:
might not yield the NN!

SEARCHING GENERAL TREE STRUCTURES

Generic NN search procedure for binary space partition trees

Given a test point \mathbf{x} and a tree node v (initially $v = \text{root}$):

1. Pick most optimistic child L , recursively find NN of \mathbf{x} in L (call it \mathbf{x}_L).
2. Let R be the other child. If

$$\|\mathbf{x} - \mathbf{x}_L\|_2 < \min_{\mathbf{x}' \in R} \|\mathbf{x} - \mathbf{x}'\|_2 \quad (\star)$$

then return \mathbf{x}_L .

3. Otherwise recursively find NN of \mathbf{x} in R (call it \mathbf{x}_R);
return the closer of \mathbf{x}_L and \mathbf{x}_R .

SEARCHING GENERAL TREE STRUCTURES

Generic NN search procedure for binary space partition trees

Given a test point \mathbf{x} and a tree node v (initially $v = \text{root}$):

1. Pick most optimistic child L , recursively find NN of \mathbf{x} in L (call it \mathbf{x}_L).
2. Let R be the other child. If

$$\|\mathbf{x} - \mathbf{x}_L\|_2 < \min_{\mathbf{x}' \in R} \|\mathbf{x} - \mathbf{x}'\|_2 \quad (\star)$$

then return \mathbf{x}_L .

3. Otherwise recursively find NN of \mathbf{x} in R (call it \mathbf{x}_R);
return the closer of \mathbf{x}_L and \mathbf{x}_R .

Note: can't always guarantee $O(\log(n))$ search time due to Step 3.

SEARCHING GENERAL TREE STRUCTURES

Generic NN search procedure for binary space partition trees

Given a test point x and a tree node v (initially $v = \text{root}$):

1. Pick most optimistic child L , recursively find NN of x in L (call it x_L).
2. Let R be the other child. If

$$\|x - x_L\|_2 < \min_{x' \in R} \|x - x'\|_2 \quad (\star)$$

then return x_L .

3. Otherwise recursively find NN of x in R (call it x_R);
return the closer of x_L and x_R .

Note: can't always guarantee $O(\log(n))$ search time due to Step 3.

Question: How do you check if (\star) is true?

- **Note:** it's correct though computationally wasteful to declare "false" in Step 2 even if (\star) turns out to be true.

USING GEOMETRIC PROPERTIES

For K-D trees:

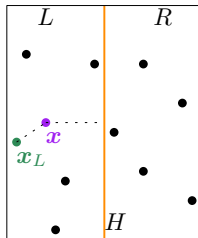
L and R are separated by a hyperplane $H = \{z \in \mathbb{R}^d : z_j = m\}$.

USING GEOMETRIC PROPERTIES

For K-D trees:

L and R are separated by a hyperplane $H = \{z \in \mathbb{R}^d : z_j = m\}$.

Suppose test point x is in L , and the NN of x in L is x_L .



USING GEOMETRIC PROPERTIES

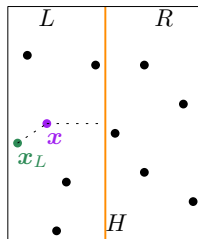
For K-D trees:

L and R are separated by a hyperplane $H = \{z \in \mathbb{R}^d : z_j = m\}$.

Suppose test point x is in L , and the NN of x in L is x_L .

By geometry,

$$\begin{aligned} \min_{x' \in R} \|x - x'\|_2 &\geq \text{distance from } x \text{ to } H \\ &= |x_j - m|. \end{aligned}$$



USING GEOMETRIC PROPERTIES

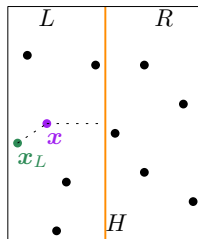
For K-D trees:

L and R are separated by a hyperplane $H = \{z \in \mathbb{R}^d : z_j = m\}$.

Suppose test point x is in L , and the NN of x in L is x_L .

By geometry,

$$\begin{aligned} \min_{x' \in R} \|x - x'\|_2 &\geq \text{distance from } x \text{ to } H \\ &= |x_j - m|. \end{aligned}$$



A valid check: if $\|x - x_L\|_2 < |x_j - m|$, then

$$\|x - x_L\|_2 < \min_{x' \in R} \|x - x'\|_2.$$

once there is
point at left part

In this case, we can skip searching R and immediately return x_L .

EFFICIENT NN SEARCH?

For certain kinds of binary space partition trees (similar to K-D trees), enough pruning will happen so NN search typically completes in $O(2^d \log(n))$ time.

- ▶ Very fast in low dimensions.
- ▶ But can be slow in high dimensions.

EFFICIENT NN SEARCH?

For certain kinds of binary space partition trees (similar to K-D trees), enough pruning will happen so NN search typically completes in $O(2^d \log(n))$ time.

- ▶ Very fast in low dimensions.
- ▶ But can be slow in high dimensions.

But NN search is only means to an end—ultimate goal is good classification.
K-D tree construction doesn't even look at the labels!

EFFICIENT NN SEARCH?

For certain kinds of binary space partition trees (similar to K-D trees), enough pruning will happen so NN search typically completes in $O(2^d \log(n))$ time.

- ▶ Very fast in low dimensions.
- ▶ But can be slow in high dimensions.

But NN search is only means to an end—ultimate goal is good classification.

K-D tree construction doesn't even look at the labels!

Question: Can we use trees to directly build good classifiers? (Next lecture.)