

10-301/601: Introduction to Machine Learning

Lecture 4 – Overfitting & KNNs

Hoda Heidari, Henry Chai & Matt Gormley

1/29/24

Front Matter

- Announcements:
 - HW2 released 1/24, due 2/5 at 11:59 PM

Recall: Decision Tree – Pseudocode

```
def train( $\mathcal{D}_{train}$ ):  
    store root = tree_recurse( $\mathcal{D}_{train}$ )  
def tree_recurse( $\mathcal{D}'$ ):  
    q = new node()  
    base case – if (SOME CONDITION):  
    recursion – else:  
        find best attribute to split on,  $x_d$   
        q.split =  $x_d$   
        for  $v$  in  $V(x_d)$ , all possible values of  $x_d$ :  
             $\mathcal{D}_v = \{(x^{(n)}, y^{(n)}) \in \mathcal{D} \mid x_d^{(n)} = v\}$   
            q.children( $v$ ) = tree_recurse( $\mathcal{D}_v$ )  
    return q
```

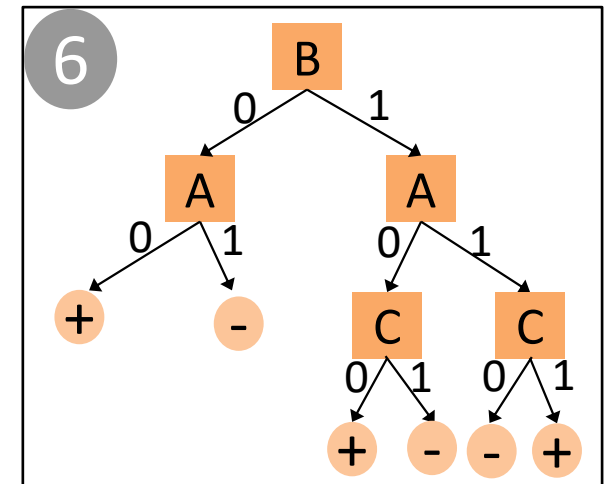
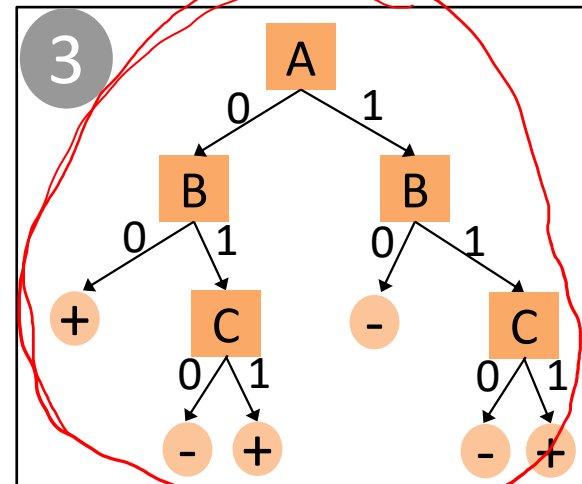
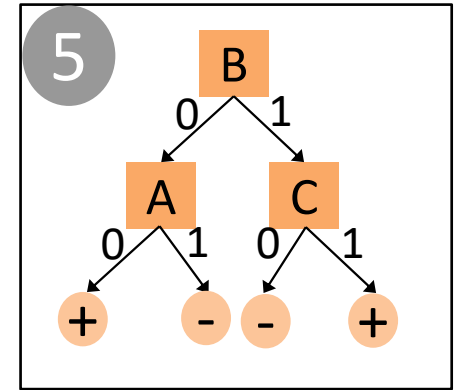
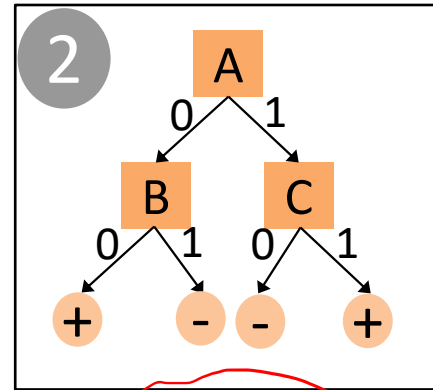
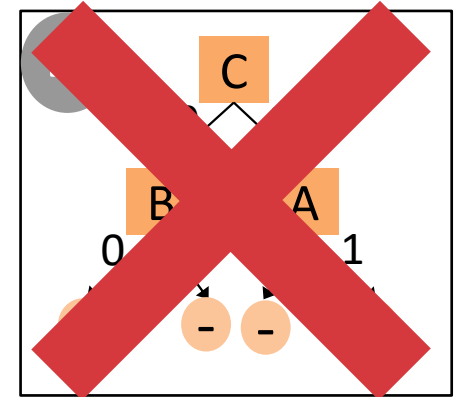
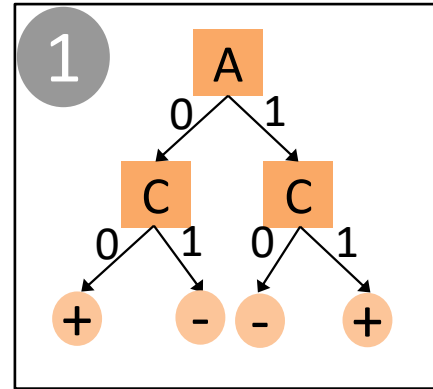
Recall: Decision Tree – Pseudocode

```
def train( $\mathcal{D}_{train}$ ):  
    store root = tree_recurse( $\mathcal{D}_{train}$ )  
def tree_recurse( $\mathcal{D}'$ ):  
    q = new node()  
    base case – if ( $\mathcal{D}'$  is empty OR  
        all labels in  $\mathcal{D}'$  are the same OR  
        all features in  $\mathcal{D}'$  are identical OR  
        some other stopping criterion):  
        q.label = majority_vote( $\mathcal{D}'$ )  
  
    recursion – else:  
        return q
```

Poll Question 1:
Which decision tree would you learn if you used training error rate as the splitting criterion? Break ties alphabetically.

Open the poll, either by clicking the [Poll] link on the schedule page of our course website or going to <http://poll.mlcourse.org>

<i>A</i>	<i>B</i>	<i>C</i>	<i>y</i>
0	0	0	+
0	0	1	+
0	1	0	-
0	1	1	+
1	0	0	-
1	0	1	-
1	1	0	-
1	1	1	+

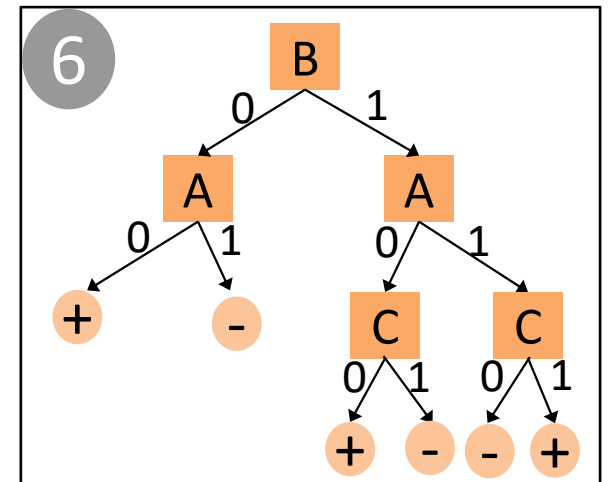
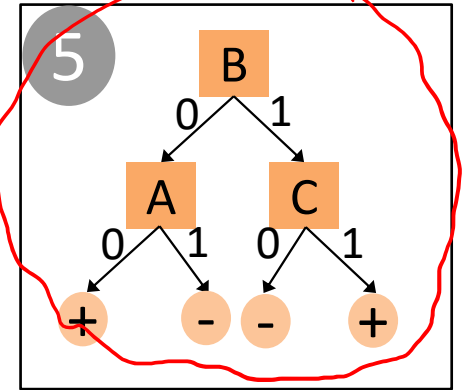
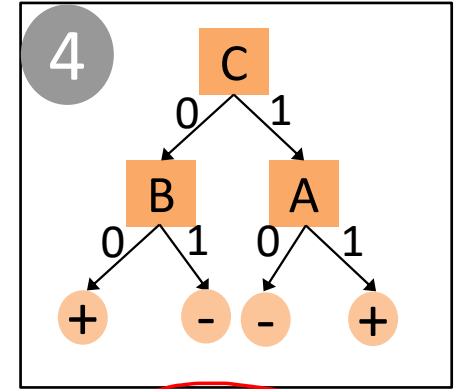
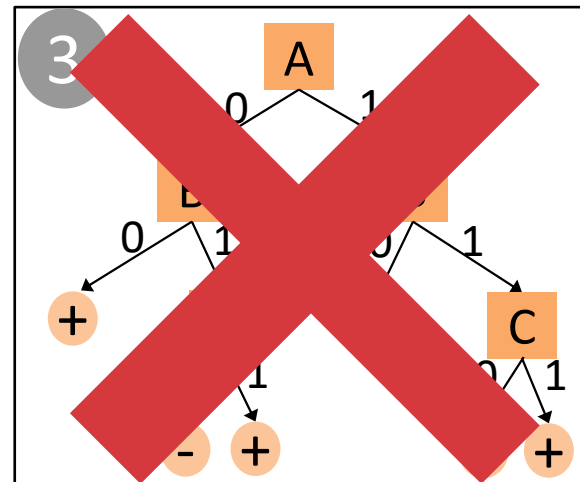
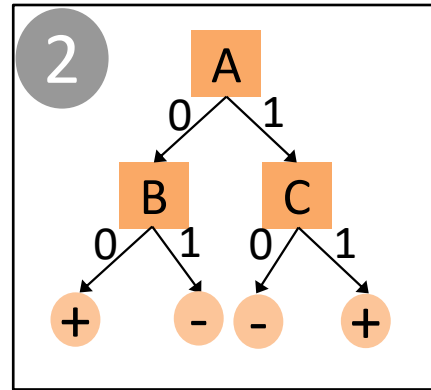
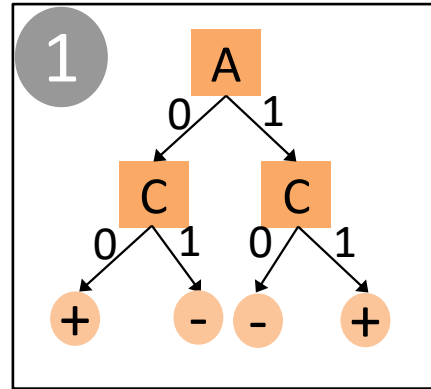


Poll Question 2:

Which decision tree is the smallest decision tree that achieves the lowest possible training error?

Open the poll, either by clicking the [Poll] link on the schedule page of our course website or going to <http://poll.mlcourse.org>

<i>A</i>	<i>B</i>	<i>C</i>	<i>y</i>
0	0	0	+
0	0	1	+
0	1	0	-
0	1	1	+
1	0	0	-
1	0	1	-
1	1	0	-
1	1	1	+



Decision Trees: Inductive Bias

- The **inductive bias** of a machine learning algorithm is the principle by which it generalizes to unseen examples
- What is the inductive bias of the ID3 algorithm i.e., decision tree learning with mutual information maximization as the splitting criterion?
 - Try to find the smallest tree that achieves (purity?) low/zero? training err. with high mutual info. features at the top

Decision Trees: Pros & Cons

- Pros
 - Interpretable
 - Efficient (computational cost and storage)
 - Can be used for classification and regression tasks
 - Compatible with categorical and real-valued features
- Cons

Decision Trees & Real-Valued Features

x Resting Systolic Pressure	y Heart Disease?
135	No
142	Yes
118	No
147	Yes
121	Yes
133	No
140	No
125	Yes
150	Yes
126	No



x Resting Systolic Pressure	y Heart Disease?
118	No
121	Yes
125	Yes
126	No
133	No
135	No
140	No
142	Yes
147	Yes
150	Yes

← $x < 119.5?$

Decision Trees & Real-Valued Features

x Resting Systolic Pressure	y Heart Disease?
135	No
142	Yes
118	No
147	Yes
121	Yes
133	No
140	No
125	Yes
150	Yes
126	No



x Resting Systolic Pressure	y Heart Disease?
118	No
121	Yes
125	Yes
126	No
133	No
135	No
140	No
142	Yes
147	Yes
150	Yes

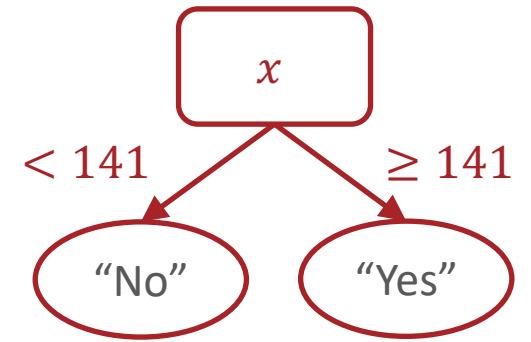
← $x < 123?$

Decision Trees & Real-Valued Features

x Resting Systolic Pressure	y Heart Disease?
135	No
142	Yes
118	No
147	Yes
121	Yes
133	No
140	No
125	Yes
150	Yes
126	No



x Resting Systolic Pressure	y Heart Disease?
118	No
121	Yes
125	Yes
126	No
133	No
135	No
140	No
142	Yes
147	Yes
150	Yes



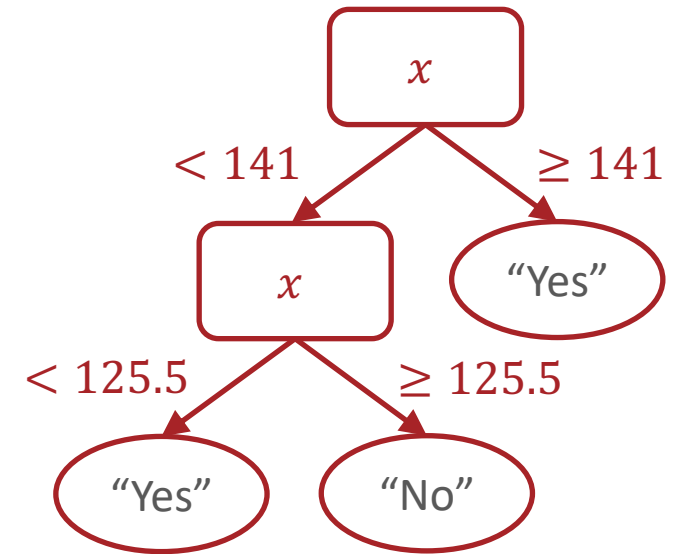
← $x < 141?$

Decision Trees & Real-Valued Features

x Resting Systolic Pressure	y Heart Disease?
135	No
142	Yes
118	No
147	Yes
121	Yes
133	No
140	No
125	Yes
150	Yes
126	No



x Resting Systolic Pressure	y Heart Disease?
118	No
121	Yes
125	Yes
126	No
133	No
135	No
140	No
142	Yes
147	Yes
150	Yes



← $x < 125.5?$

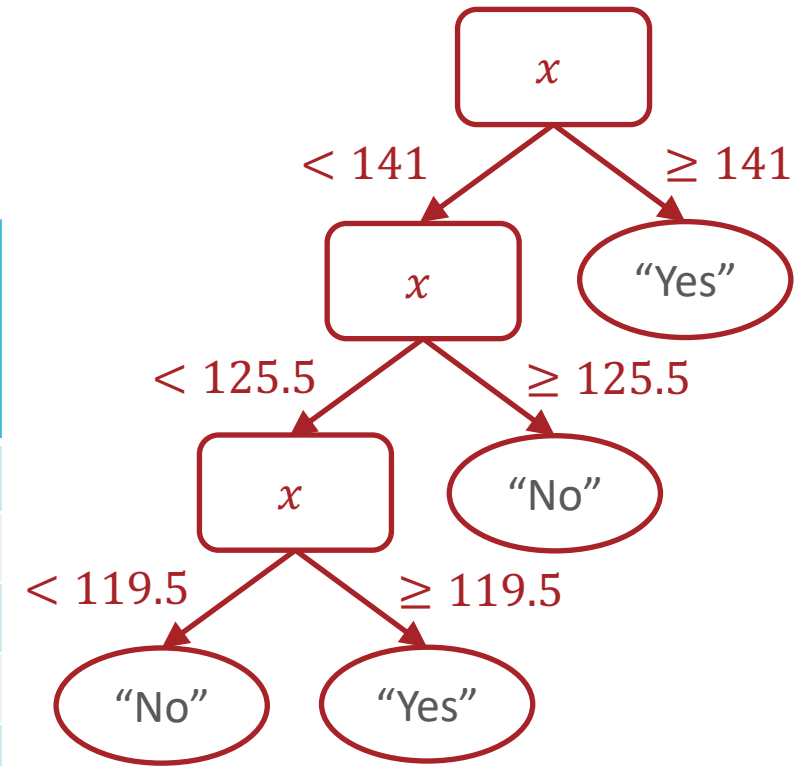
← $x < 141?$

Decision Trees & Real-Valued Features

x Resting Systolic Pressure	y Heart Disease?
135	No
142	Yes
118	No
147	Yes
121	Yes
133	No
140	No
125	Yes
150	Yes
126	No



x Resting Systolic Pressure	y Heart Disease?
118	No
121	Yes
125	Yes
126	No
133	No
135	No
140	No
142	Yes
147	Yes
150	Yes



- Discretize real-valued features using thresholds (effectively creating *new* categorical features)
- Can split on real-valued features multiple times in the same tree

Decision Trees: Pros & Cons

- Pros
 - Interpretable
 - Efficient (computational cost and storage)
 - Can be used for classification and regression tasks
 - Compatible with categorical and real-valued features
- Cons
 - Learned greedily: each split only considers the immediate impact on the splitting criterion
 - Not guaranteed to find the smallest (fewest number of splits) tree that achieves a training error rate of 0.
 - Liable to overfit!

Overfitting

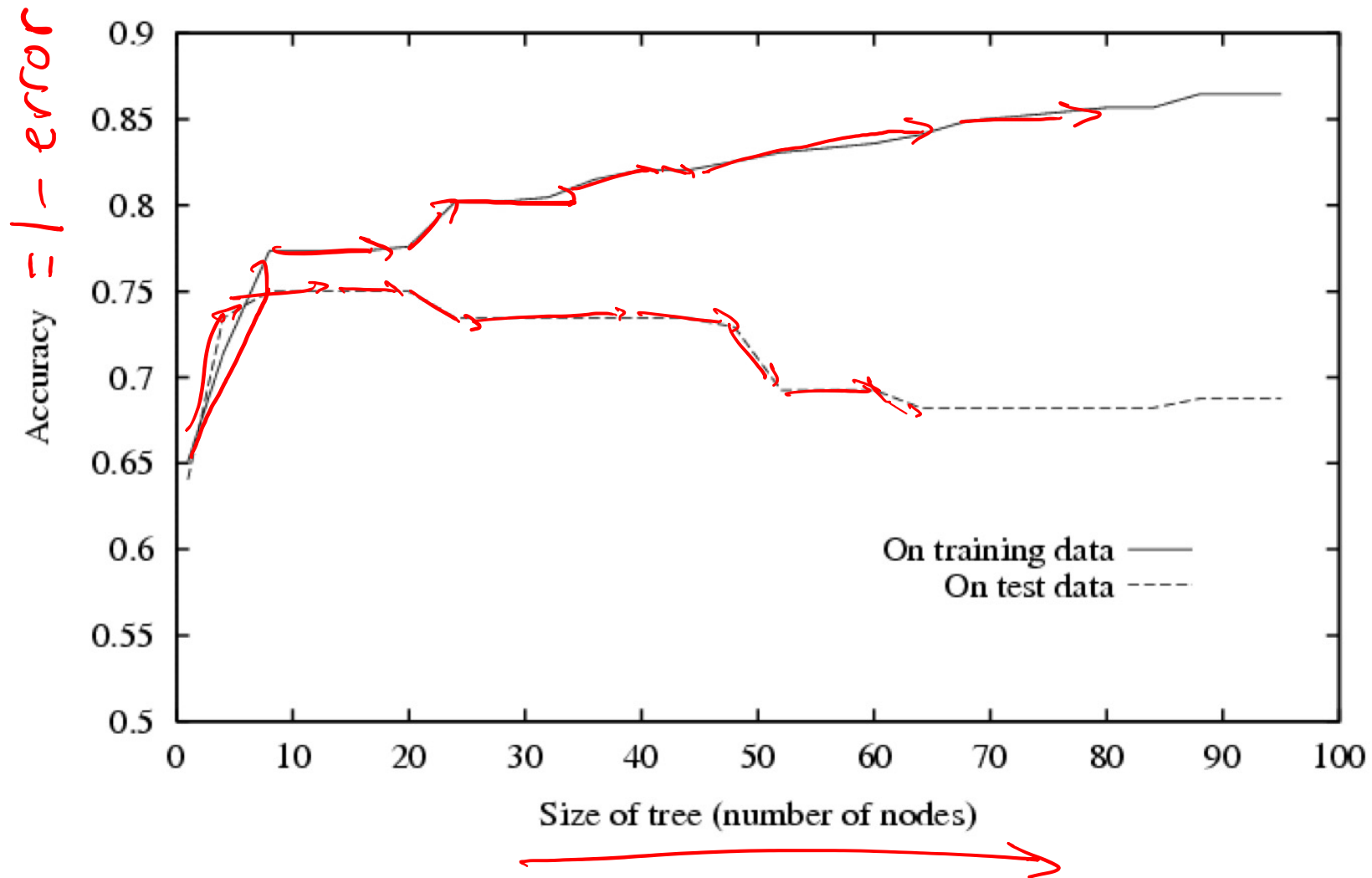
- Overfitting occurs when the classifier (or model)...
 - is too complex
 - fits noise or “outliers” in the training dataset as opposed to the actual pattern of interest
 - doesn’t have enough inductive bias pushing it to generalize (e.g., the memorizer)
- Underfitting occurs when the classifier (or model)...
 - is too simple
 - can’t capture the actual pattern of interest in the training dataset
 - has too much inductive bias (e.g., majority vote)

Recall: Different Kinds of Error

- Training error rate = $err(h, \mathcal{D}_{train})$
- Test error rate = $err(h, \mathcal{D}_{test})$
- True error rate = $err(h)$
 - = the error rate of h on all possible examples
 - In machine learning, this is the quantity that we care about but, in most cases, it is unknowable.
- Overfitting occurs when $err(h) > err(h, \mathcal{D}_{train})$
 - $err(h) - err(h, \mathcal{D}_{train})$ can be thought of as a measure of overfitting

Overfitting in Decision Trees

rate



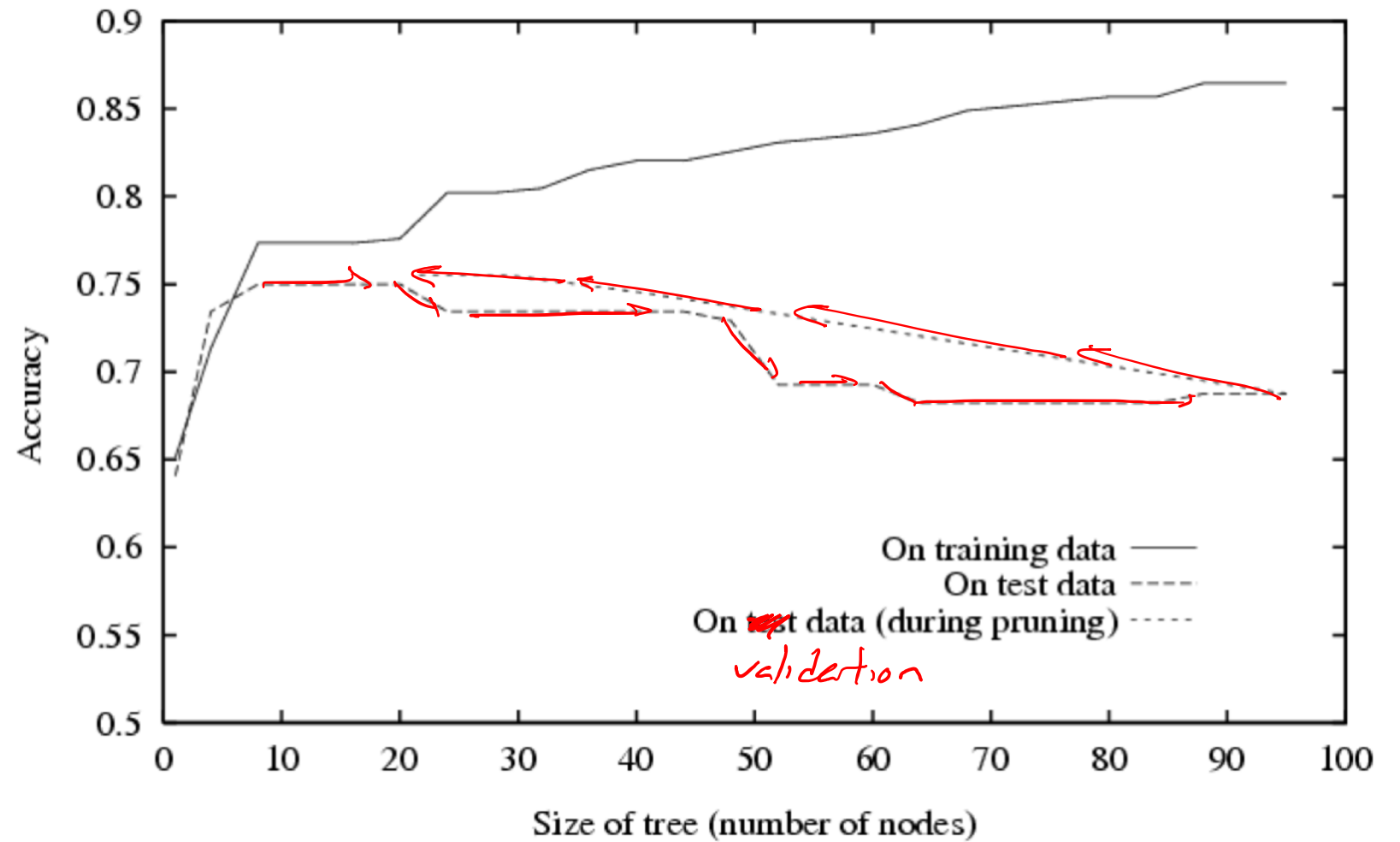
Combatting Overfitting in Decision Trees

- Heuristics:
 - Do not split leaves past a fixed depth, δ
 - Do not split leaves with fewer than c data points
 - Do not split leaves where the maximal information gain is less than τ
- Take a majority vote in impure leaves

Combatting Overfitting in Decision Trees

- Pruning:
 - First, learn a decision tree
 - Then, evaluate each split using a “validation” dataset by comparing the validation error rate with and without that split
 - Greedily remove the split that most decreases the validation error rate
 - Stop if no split is removed

Pruning Decision Trees



Decision Tree Learning Objectives

You should be able to...

1. Implement decision tree training and prediction
2. Use effective splitting criteria for decision trees and be able to define entropy, conditional entropy, and mutual information / information gain
3. Explain the difference between memorization and generalization [CIML]
4. Describe the inductive bias of a decision tree
5. Formalize a learning problem by identifying the input space, output space, hypothesis space, and target function
6. Explain the difference between true error and training error
7. Judge whether a decision tree is "underfitting" or "overfitting"
8. Implement a pruning or early stopping method to combat overfitting in decision tree learning

REAL VALUED ATTRIBUTES



Fisher Iris Dataset

Fisher (1936) used 150 measurements of flowers from 3 different species: Iris setosa (0), Iris virginica (1), Iris versicolor (2) collected by Anderson (1936)

Species	Sepal Length	Sepal Width	Petal Length	Petal Width
0	4.3	3.0	1.1	0.1
0	4.9	3.6	1.4	0.1
0	5.3	3.7	1.5	0.2
1	4.9	2.4	3.3	1.0
1	5.7	2.8	4.1	1.3
1	6.3	3.3	4.7	1.6
1	6.7	3.0	5.0	1.7

Fisher Iris Dataset

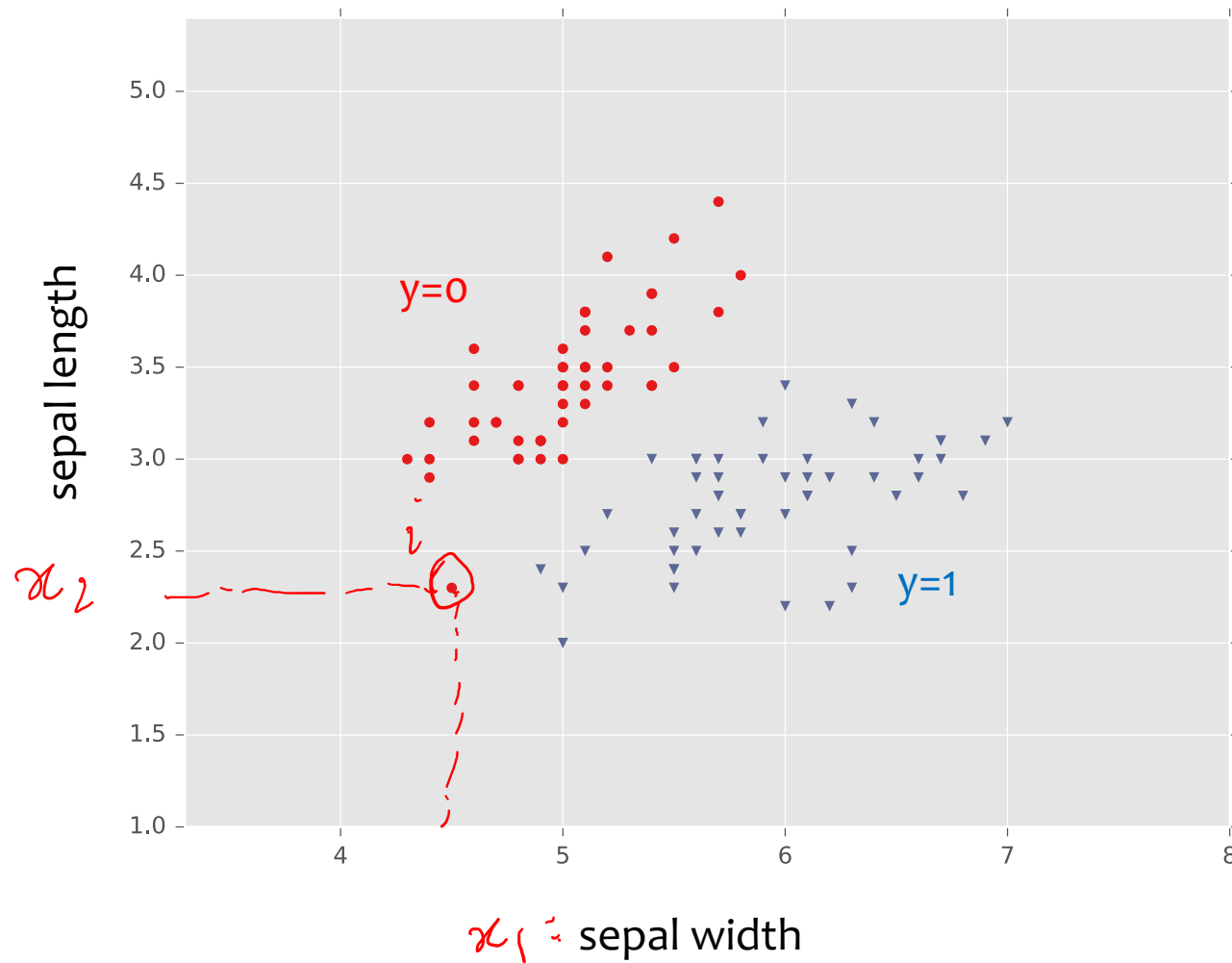
Fisher (1936) used 150 measurements of flowers from 3 different species: Iris setosa (0), Iris virginica (1), Iris versicolor (2) collected by Anderson (1936)

Species	Sepal Length	Sepal Width
0	4.3	3.0
0	4.9	3.6
0	5.3	3.7
1	4.9	2.4
1	5.7	2.8
1	6.3	3.3
1	6.7	3.0

Deleted two of the four features, so that input space is 2D



Fisher Iris Dataset



$$\begin{aligned} & (\vec{x}^{(i)}, y^{(i)}) \\ \vec{x}^{(i)} &= (4.5, 2.2) \\ y^{(i)} &= 0 \end{aligned}$$

K-NEAREST NEIGHBORS

Nearest Neighbor: Algorithm

def train(\mathcal{D}):

 Store \mathcal{D}

def h(x'):

 Let $x^{(i)}$ = the point in \mathcal{D} that is nearest to x'

return $y^{(i)}$

Classification & Real-Valued Features

Classification

$$\rightarrow \mathcal{D} = \left\{ \left(\vec{x}^{(1)}, y^{(1)} \right), \dots, \left(\vec{x}^{(N)}, y^{(N)} \right) \right\}$$

$$\begin{cases} \vec{x}^{(i)} \in \mathcal{X} = \mathbb{R}^M \\ y^{(i)} \in \mathcal{Y} = \{1, 2, \dots, L\} \end{cases}$$

Binary Classification

$$|\mathcal{Y}| = 2$$

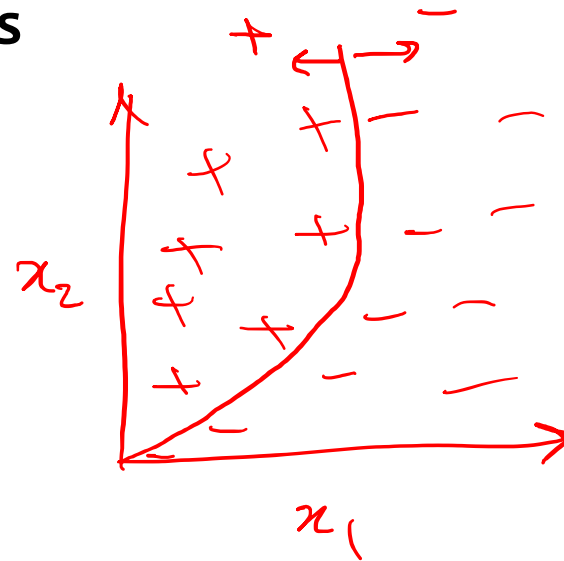
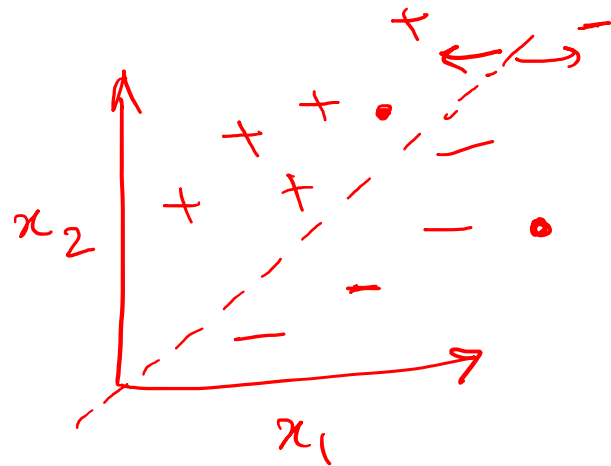
$$\mathcal{Y} = \{-, +\}$$

$$\{0, 1\}$$

$$\{\text{red}, \text{white}\}$$

Classification & Real-Valued Features

Decision Rules / Decision Boundaries



$$h \in \mathcal{H}$$

↓
Class of all linear fu's over x_1, x_2

{ train time: Use D to come up with h
test time: Use h to predict $\hat{y} = h(\vec{x})$

\mathcal{H}'
↓
quadratic --

Nearest Neighbor: Algorithm

def train(\mathcal{D}):

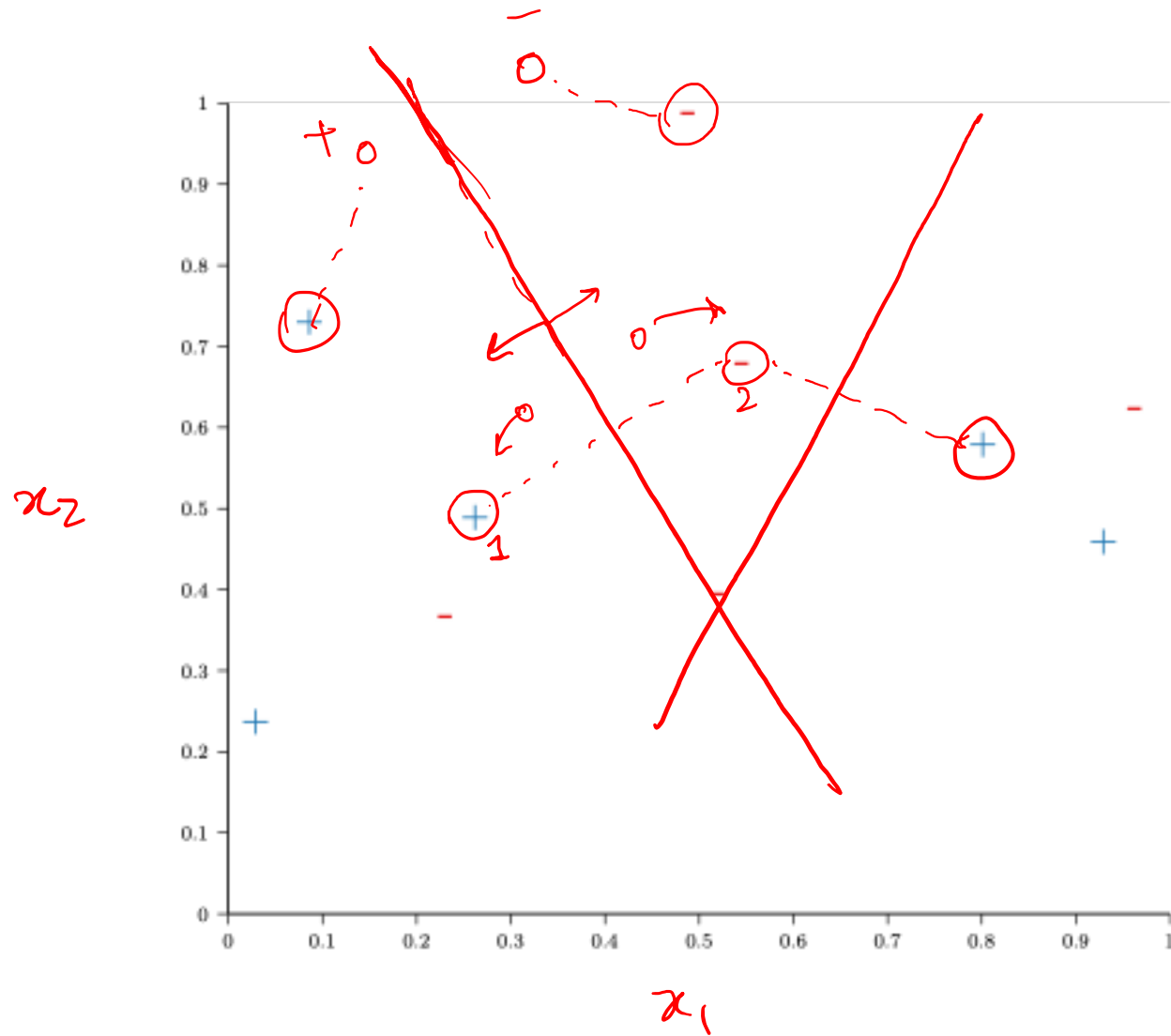
 Store \mathcal{D}

def h(\mathbf{x}'):

 Let $\mathbf{x}^{(i)}$ = the point in \mathcal{D} that is nearest to \mathbf{x}'

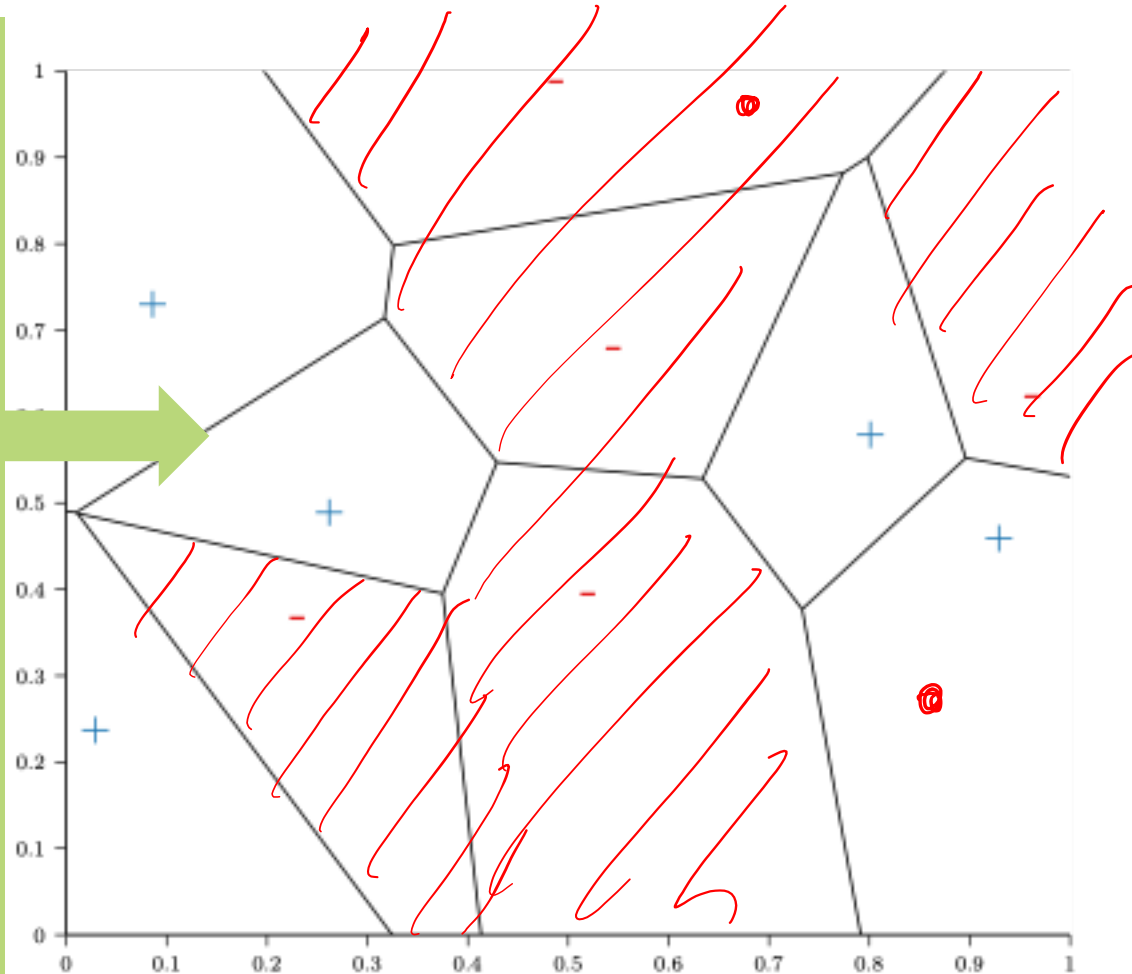
return $y^{(i)}$

Nearest Neighbor: Example

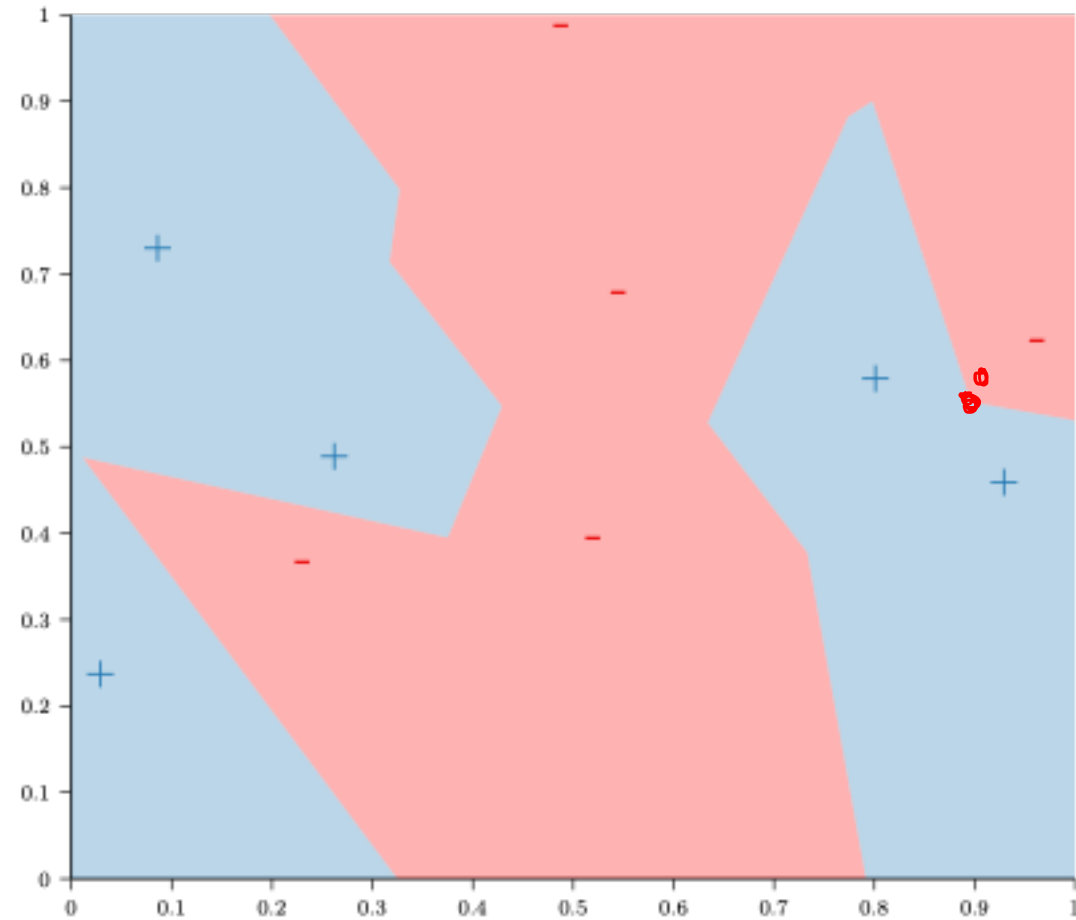


Nearest Neighbor: Example

- This is a **Voronoi diagram**
- Each **cell** contain one of our training examples
- **All points within a cell** are **closer** to that training example, than to any other training example
- **Points on the Voronoi line segments** are **equidistant** to one or more training examples



Nearest Neighbor: Example

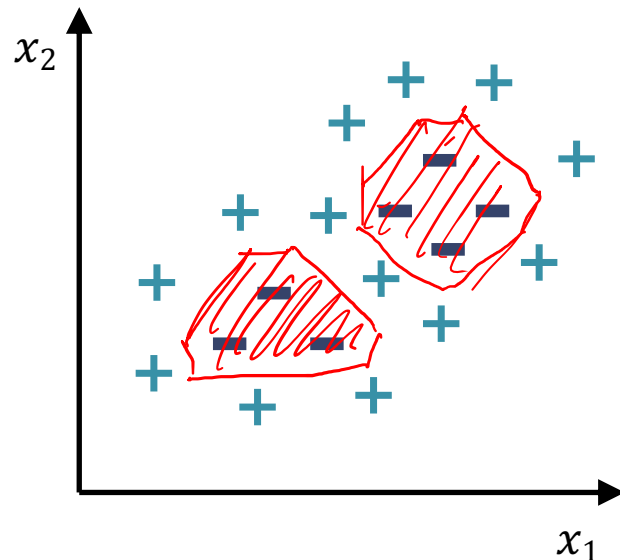


Decision Boundary Exercise

Poll Question 1:

Can a **1-Nearest Neighbor classifier** achieve **zero training error** on this dataset? If so, draw the decision boundary and if not, briefly explain why.

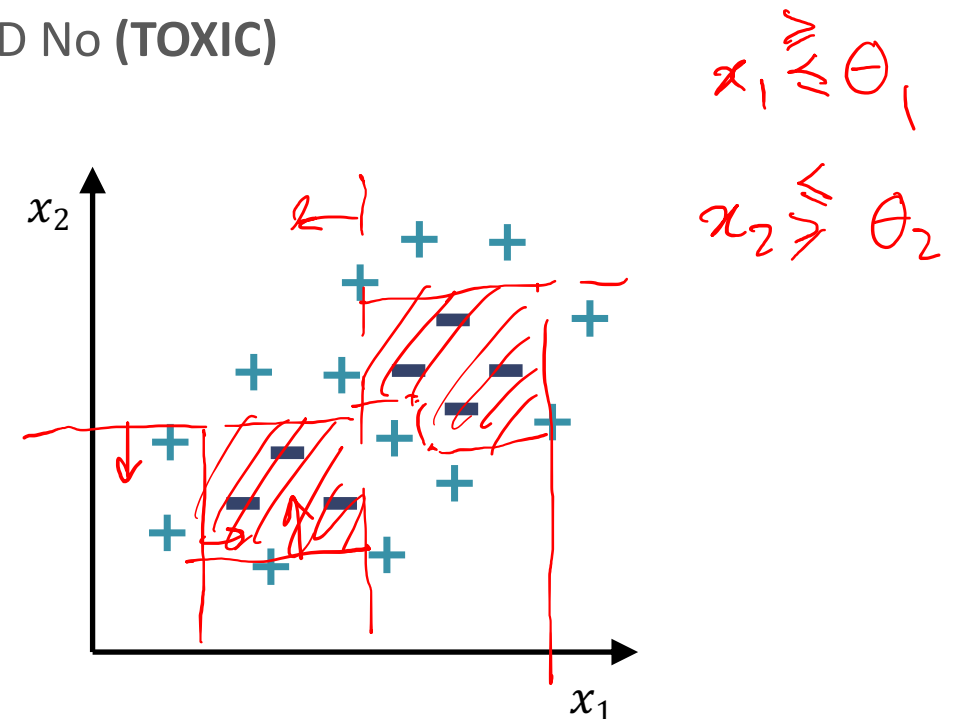
- A. Yes
- B. No
- C. Yes AND No (**TOXIC**)



Poll Question 2:

Can a **Decision Tree classifier** achieve **zero training error** on this dataset? If so, draw the decision boundary and if not, briefly explain why.

- A. Yes
- B. No
- C. Yes AND No (**TOXIC**)



The Nearest Neighbor Model

- Requires no training!
- Always has zero training error!
 - *A data point is always its own nearest neighbor*

k-Nearest Neighbors: Algorithm

```
def set_hyperparameters(k, d):
```

```
    Store k
```

```
    Store  $d(\cdot, \cdot)$ 
```

```
def train( $\mathcal{D}$ ):
```

```
    Store  $\mathcal{D}$ 
```

```
def h( $x'$ ):
```

```
    Let  $S$  = the set of  $k$  points in  $\mathcal{D}$  nearest to  $x'$   
    according to distance function
```

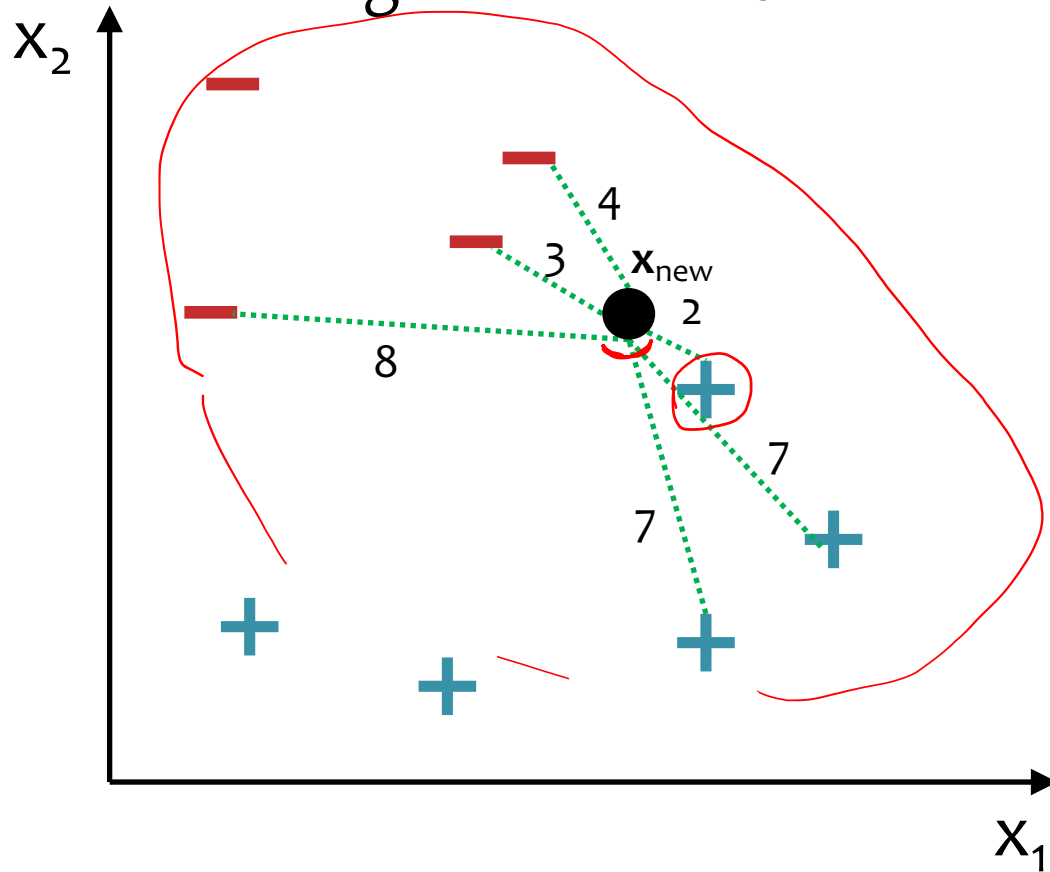
```
     $d(\mathbf{u}, \mathbf{v})$ 
```

```
    Let  $v$  = majority_vote( $S$ )
```

```
    return  $v$ 
```

k-Nearest Neighbors

Suppose we have the training dataset below.



How should we label the new point?

$k=1 \rightarrow +$
It depends on k :
if $k=1$, $h(x_{\text{new}}) = +1$
if $k=3$, $h(x_{\text{new}}) = -1$
if $k=5$, $h(x_{\text{new}}) = +1$

+



-



KNN: Remarks

Distance Functions:

- KNN requires a **distance function**

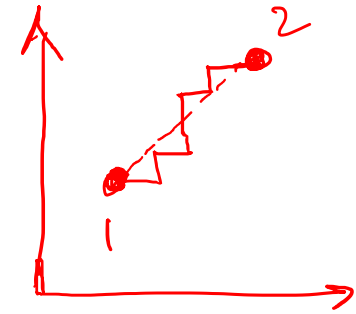
$$d : \mathbb{R}^M \times \mathbb{R}^M \rightarrow \mathbb{R}$$

- The most common choice is **Euclidean distance**

$$\rightarrow d(\mathbf{u}, \mathbf{v}) = \sqrt{\sum_{m=1}^M (u_m - v_m)^2}$$

- But there are other choices (e.g. **Manhattan distance**)

$$\rightarrow d(\mathbf{u}, \mathbf{v}) = \sum_{m=1}^M |u_m - v_m|$$



KNN: Computational Efficiency

- Suppose we have N training examples and each one has M features
- Computational complexity when $k=1$:

Task	Naive	k-d Tree
Train	$O(1)$	$\sim O(M N \log N)$
Predict (one test example)	$O(MN)$	$\sim O(2^M \log N)$ on average

Problem: Very fast for small M , but very slow for large M

In practice: use stochastic approximations (very fast, and empirically often as good)

\vec{x} For $i=1 \dots N$
 $d(x, x^{(i)})$ For $j=1 \dots M$

KNN: Theoretical Guarantees

Cover & Hart (1967)

Let $h(x)$ be a Nearest Neighbor ($k=1$) binary classifier. As the number of training examples N goes to infinity...

$$\text{error}_{\text{true}}(h) < 2 \times \text{Bayes Error Rate}$$

“In this sense, it may be said that half the classification information in an infinite sample set is contained in the nearest neighbor.”

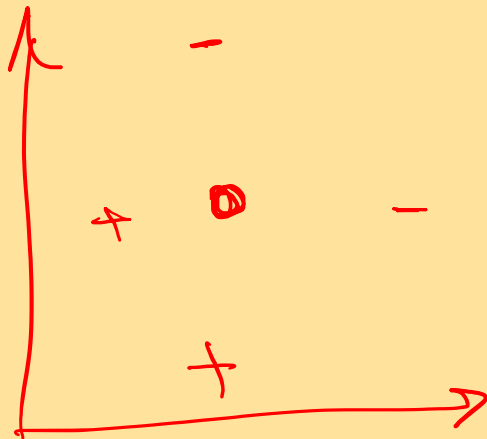
very informally,
Bayes Error Rate can be thought of as:
‘the best you could possibly do’

KNN: Remarks

In-Class Exercises

How can we handle ties for even values of k ?

$k=4$



Answer(s) Here:

- distance-weighted vote

KNN: Remarks

In-Class Exercises

How can we handle ties for even values of k ?

Answer(s) Here:

- Consider another point
- Remove farthest of k points
- Weight votes by distance
- Consider another distance metric

KNN: Inductive Bias

In-Class Exercise

What is the inductive bias of KNN?