



10-601 Introduction to Machine Learning

Machine Learning Department
School of Computer Science
Carnegie Mellon University

Overfitting + k-Nearest Neighbors

Matt Gormley
Lecture 4
September 10, 2018

(Wrap Up)

DECISION TREE EXAMPLES

Majority Vote Classifier Example

Dataset:

Output Y, Attributes A and B

Y	A	B
-	1	0
-	1	0
+	1	0
+	1	0
+	1	1
+	1	1
+	1	1
+	1	1

In-Class Exercise

What is the **training error** (i.e. *error rate on the training data*) of the **majority vote classifier** on this dataset?

Decision Tree Learning Example

Dataset:

Output Y, Attributes A, B, C

Y	A	B	C
-	1	0	0
-	1	0	1
-	1	0	0
+	0	0	1
+	1	1	0
+	1	1	1
+	1	1	0
+	1	1	1

In-Class Exercise

Using **error rate** as the splitting criterion, what decision tree would be learned by ID3?

Decision Tree Learning Example

Dataset:

Output Y, Attributes A and B

Y	A	B
-	1	0
-	1	0
+	1	0
+	1	0
+	1	1
+	1	1
+	1	1
+	1	1

In-Class Exercise

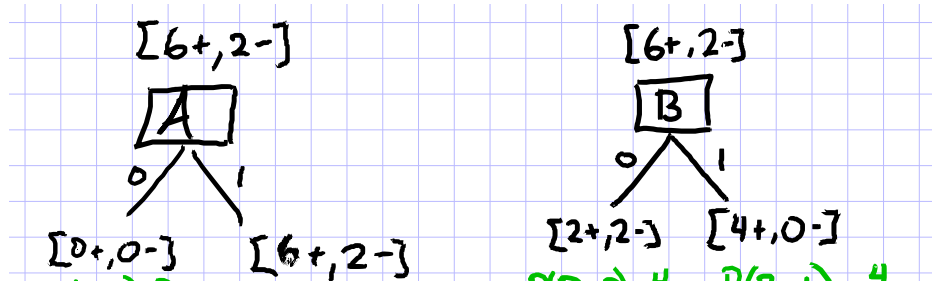
1. Which attribute would **misclassification rate** select for the next split?
2. Which attribute would **information gain** select for the next split?
3. *Justify your answers.*

Decision Tree Learning Example

Dataset:

Output Y, Attributes A and B

Y	A	B
-	1	0
-	1	0
+	1	0
+	1	0
+	1	1
+	1	1
+	1	1
+	1	1



Misclass. Rate

$$r(A) = 2/8$$

$$r(A) = 2/8$$

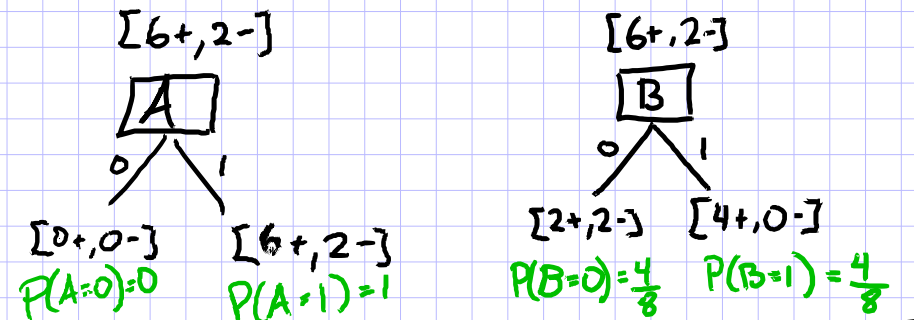
*r(.) treats
attributes as
equally good*

Decision Tree Learning Example

Dataset:

Output Y, Attributes A and B

Y	A	B
-	1	0
-	1	0
+	1	0
+	1	0
+	1	1
+	1	1
+	1	1
+	1	1



Info Gain

$$H(Y) = \frac{2}{8} \log \frac{2}{8} + \frac{6}{8} \log \frac{6}{8}$$

$$\begin{aligned}
 & \left\{ \begin{aligned}
 H(Y|A=0) &= \text{"undefined"} \\
 H(Y|A=1) &= \frac{2}{8} \log \frac{2}{8} + \frac{6}{8} \log \frac{6}{8} = H(Y) \\
 H(Y|A) &= P(A=0)H(Y|A=0) + P(A=1)H(Y|A=1) \\
 &= 0 + H(Y|A=1) = H(Y) \\
 I(Y;A) &= H(Y) - H(Y|A) = 0
 \end{aligned} \right.
 \end{aligned}$$

$$\begin{aligned}
 & \left\{ \begin{aligned}
 H(Y|B=0) &= \frac{2}{4} \log \frac{2}{4} + \frac{2}{4} \log \frac{2}{4} \\
 H(Y|B=1) &= 0 \log 0 + 1 \log 1 = 0 \\
 H(Y|B) &= \frac{4}{8} (0) + \frac{4}{8} (H(Y|B=0)) \\
 I(Y;B) &= H(Y) - \frac{4}{8} H(Y|B=0) > 0
 \end{aligned} \right.
 \end{aligned}$$

Tennis Example

Test your understanding

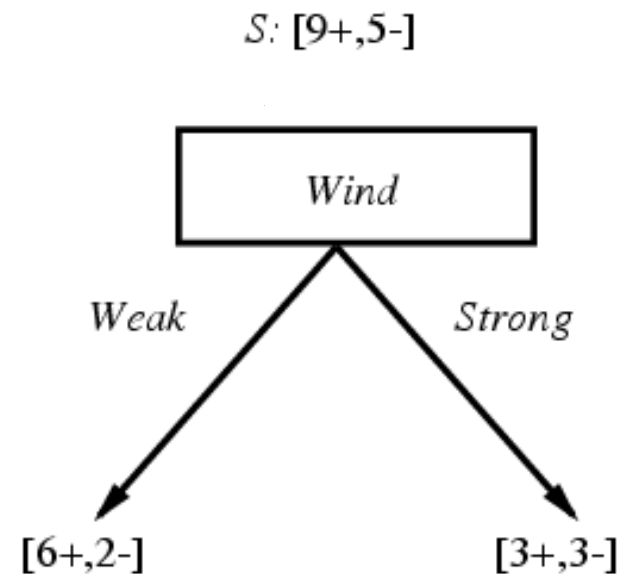
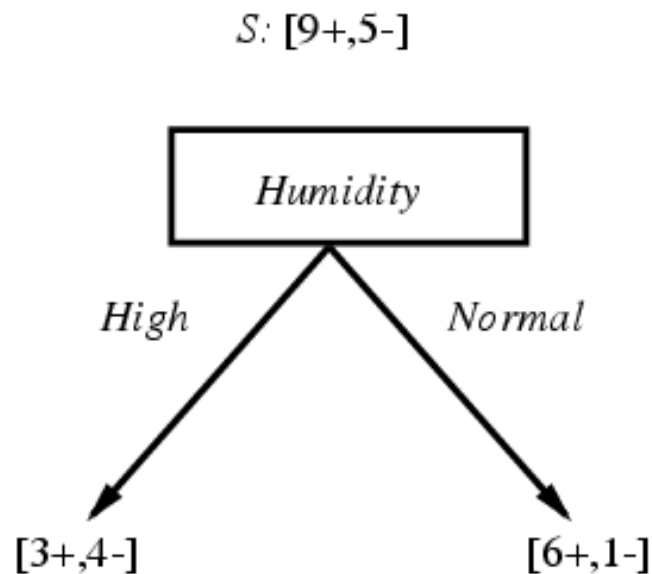
Dataset:

Day	Outlook	Temperature	Humidity	Wind	PlayTennis?
D1	Sunny	Hot	High	Weak	No
D2	Sunny	Hot	High	Strong	No
D3	Overcast	Hot	High	Weak	Yes
D4	Rain	Mild	High	Weak	Yes
D5	Rain	Cool	Normal	Weak	Yes
D6	Rain	Cool	Normal	Strong	No
D7	Overcast	Cool	Normal	Strong	Yes
D8	Sunny	Mild	High	Weak	No
D9	Sunny	Cool	Normal	Weak	Yes
D10	Rain	Mild	Normal	Weak	Yes
D11	Sunny	Mild	Normal	Strong	Yes
D12	Overcast	Mild	High	Strong	Yes
D13	Overcast	Hot	Normal	Weak	Yes
D14	Rain	Mild	High	Strong	No

Tennis Example

Which attribute yields the best classifier?

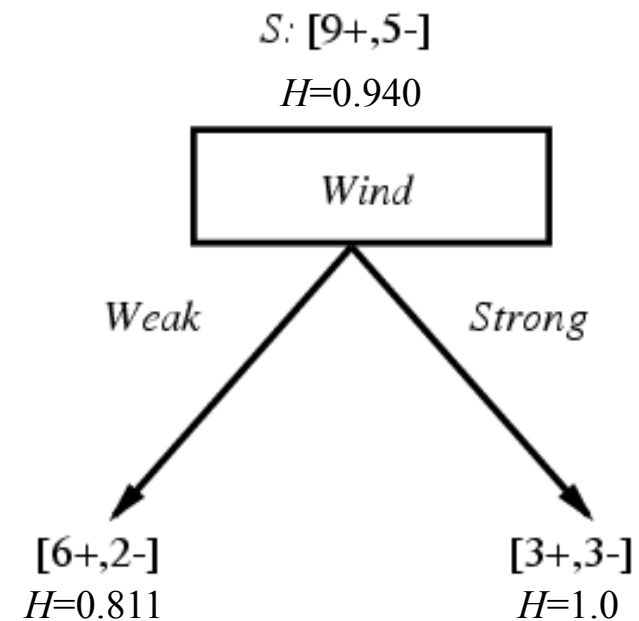
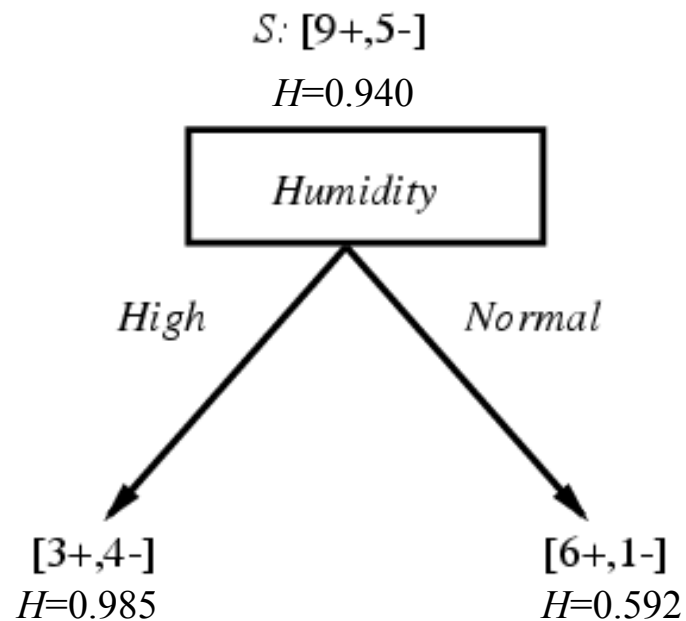
Test your understanding



Tennis Example

Which attribute yields the best classifier?

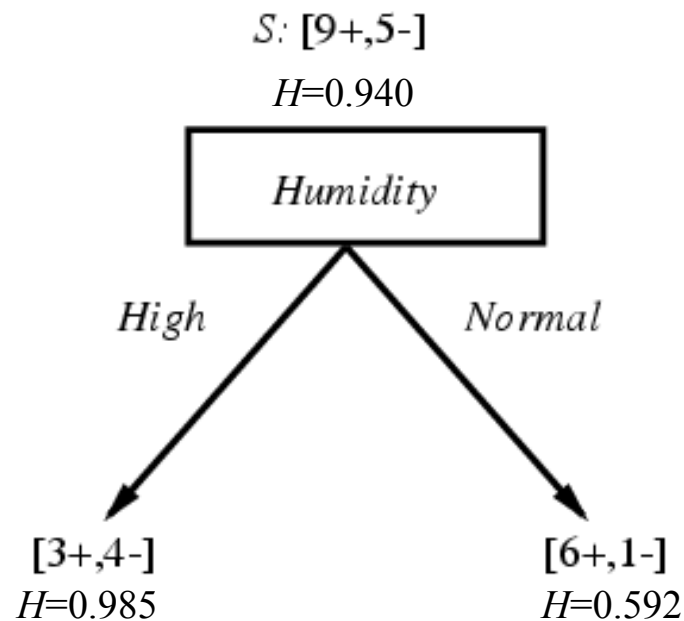
Test your understanding



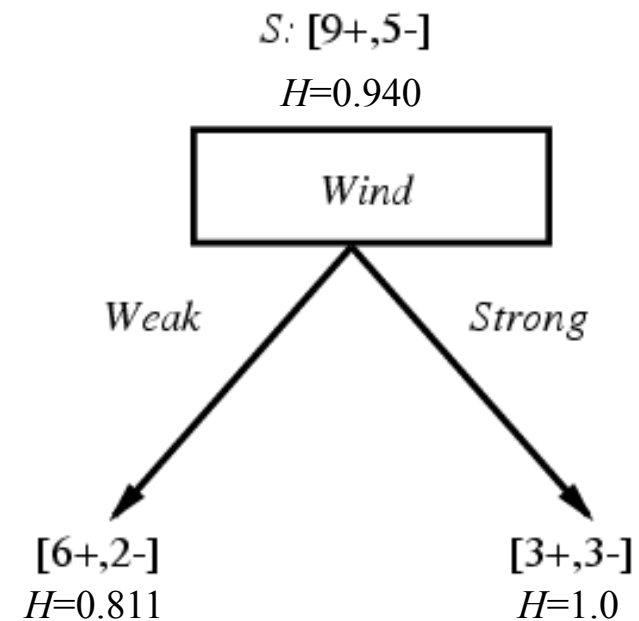
Tennis Example

Which attribute yields the best classifier?

Test your understanding



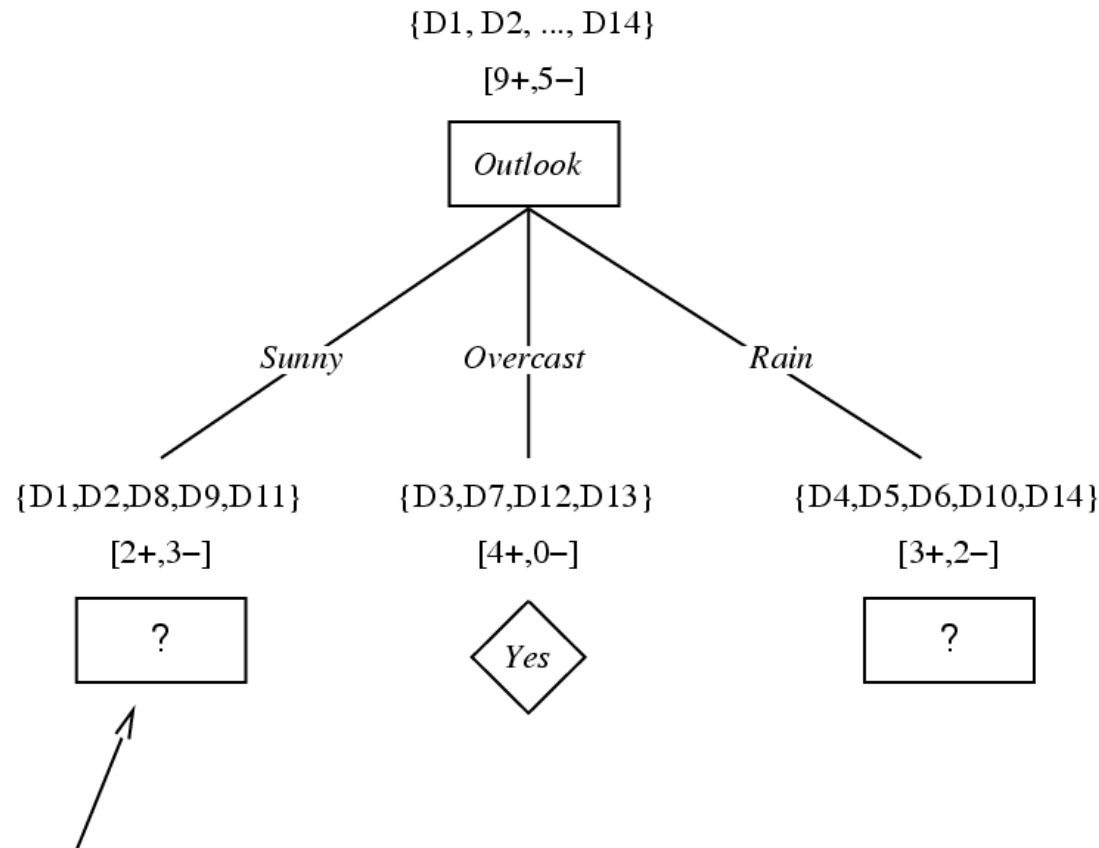
$$\begin{aligned} \text{Gain}(S, \text{Humidity}) &= .940 - (7/14).985 - (7/14).592 \\ &= .151 \end{aligned}$$



$$\begin{aligned} \text{Gain}(S, \text{Wind}) &= .940 - (8/14).811 - (6/14)1.0 \\ &= .048 \end{aligned}$$

Tennis Example

Test your understanding



Which attribute should be tested here?

$$S_{\text{sunny}} = \{D1, D2, D8, D9, D11\}$$

$$\text{Gain}(S_{\text{sunny}}, \text{Humidity}) = .970 - (3/5) 0.0 - (2/5) 0.0 = .970$$

$$\text{Gain}(S_{\text{sunny}}, \text{Temperature}) = .970 - (2/5) 0.0 - (2/5) 1.0 - (1/5) 0.0 = .570$$

$$\text{Gain}(S_{\text{sunny}}, \text{Wind}) = .970 - (2/5) 1.0 - (3/5) .918 = .019$$

Q&A

Q: Why don't my entropy calculations match those on the slides?

A: $H(Y)$ is conventionally reported in “bits” and computed using log base 2.
e.g., $H(Y) = -P(Y=0) \log_2 P(Y=0) - P(Y=1) \log_2 P(Y=1)$

Q: When and how do we decide to stop growing trees? What if the set of values an attribute could take was really large or even infinite?

A: We'll address this question for discrete attributes today. If an attribute is real-valued, there's a clever trick that only considers $O(L)$ splits where $L = \#$ of values the attribute takes in the training set. Can you guess what it does?

Q: Why is entropy based on a sum of $p(.) \log p(.)$ terms?

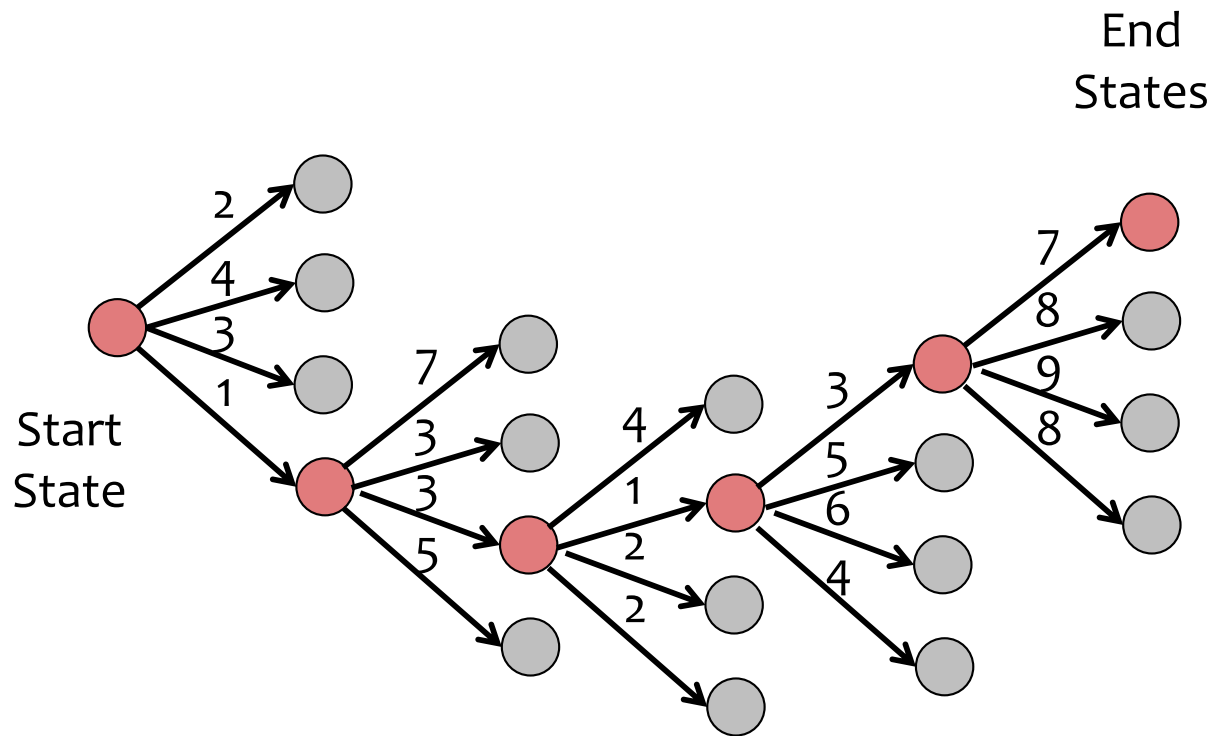
A: We don't have time for a full treatment of why it *has* to be this, but we can develop the right intuition with a few examples...

Reminders

- **Homework 2: Decision Trees**
 - **Out: Wed, Sep 05**
 - **Due: Wed, Sep 19 at 11:59pm**
- **10601 Notation Crib Sheet**

INDUCTIVE BIAS (FOR DECISION TREES)

Background: Greedy Search



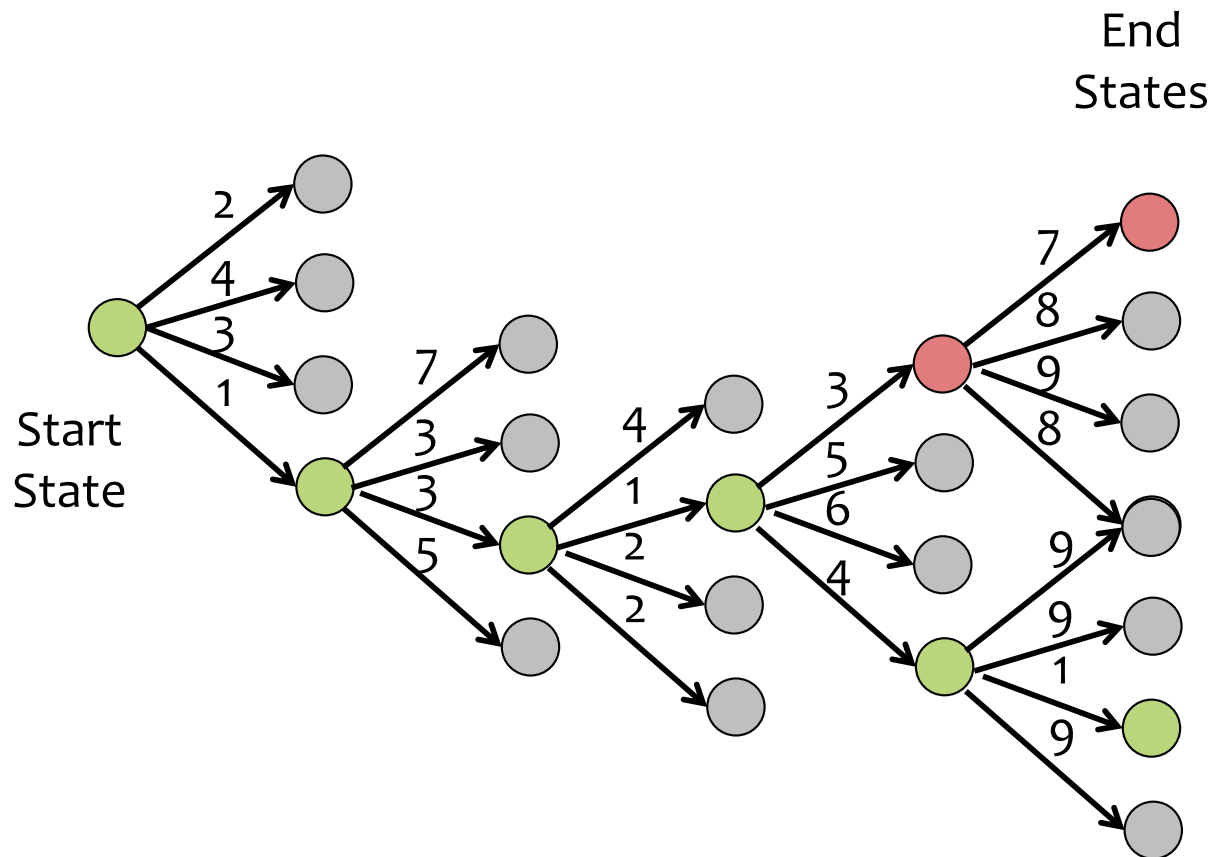
Goal:

- Search space consists of nodes and weighted edges
- Goal is to find the lowest (total) weight path from root to a leaf

Greedy Search:

- At each node, selects the edge with lowest (immediate) weight
- Heuristic method of search (i.e. does not necessarily find the best path)

Background: Greedy Search



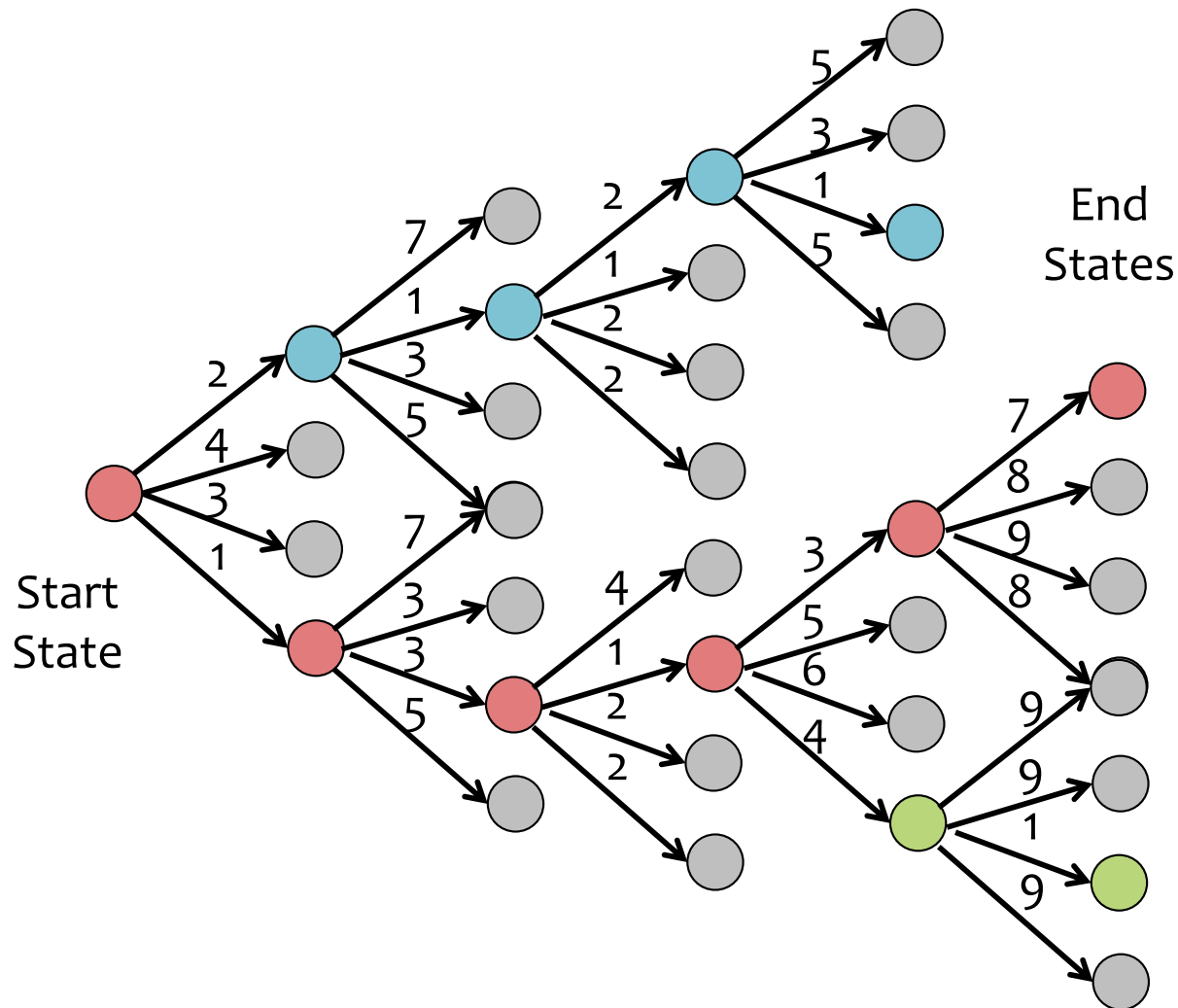
Goal:

- Search space consists of nodes and weighted edges
- Goal is to find the lowest (total) weight path from root to a leaf

Greedy Search:

- At each node, selects the edge with lowest (immediate) weight
- Heuristic method of search (i.e. does not necessarily find the best path)

Background: Greedy Search



Goal:

- Search space consists of nodes and weighted edges
- Goal is to find the lowest (total) weight path from root to a leaf

Greedy Search:

- At each node, selects the edge with lowest (immediate) weight
- Heuristic method of search (i.e. does not necessarily find the best path)

Decision Trees

Chalkboard

- ID3 as Search
- Inductive Bias of Decision Trees
- Occam's Razor

OVERFITTING (FOR DECISION TREES)

Overfitting and Underfitting

Underfitting

- The model...
 - is too simple
 - is unable captures the trends in the data
 - exhibits too much bias
- *Example:* majority-vote classifier (i.e. depth-zero decision tree)
- *Example:* a toddler (that has **not** attended medical school) attempting to carry out medical diagnosis

Overfitting

- The model...
 - is too complex
 - is fitting the noise in the data
 - or fitting random statistical fluctuations inherent in the “sample” of training data
 - does not have enough bias
- *Example:* our “memorizer” algorithm responding to an “orange shirt” attribute
- *Example:* medical student who simply memorizes patient case studies, but does not understand how to apply knowledge to new patients

Overfitting

Consider a hypothesis h and its

- Error rate over training data: $error_{train}(h)$
- True error rate over all data: $error_{true}(h)$

We say h overfits the training data if

$$error_{true}(h) > error_{train}(h)$$

Amount of overfitting =

$$error_{true}(h) - error_{train}(h)$$

Overfitting in Decision Tree Learning

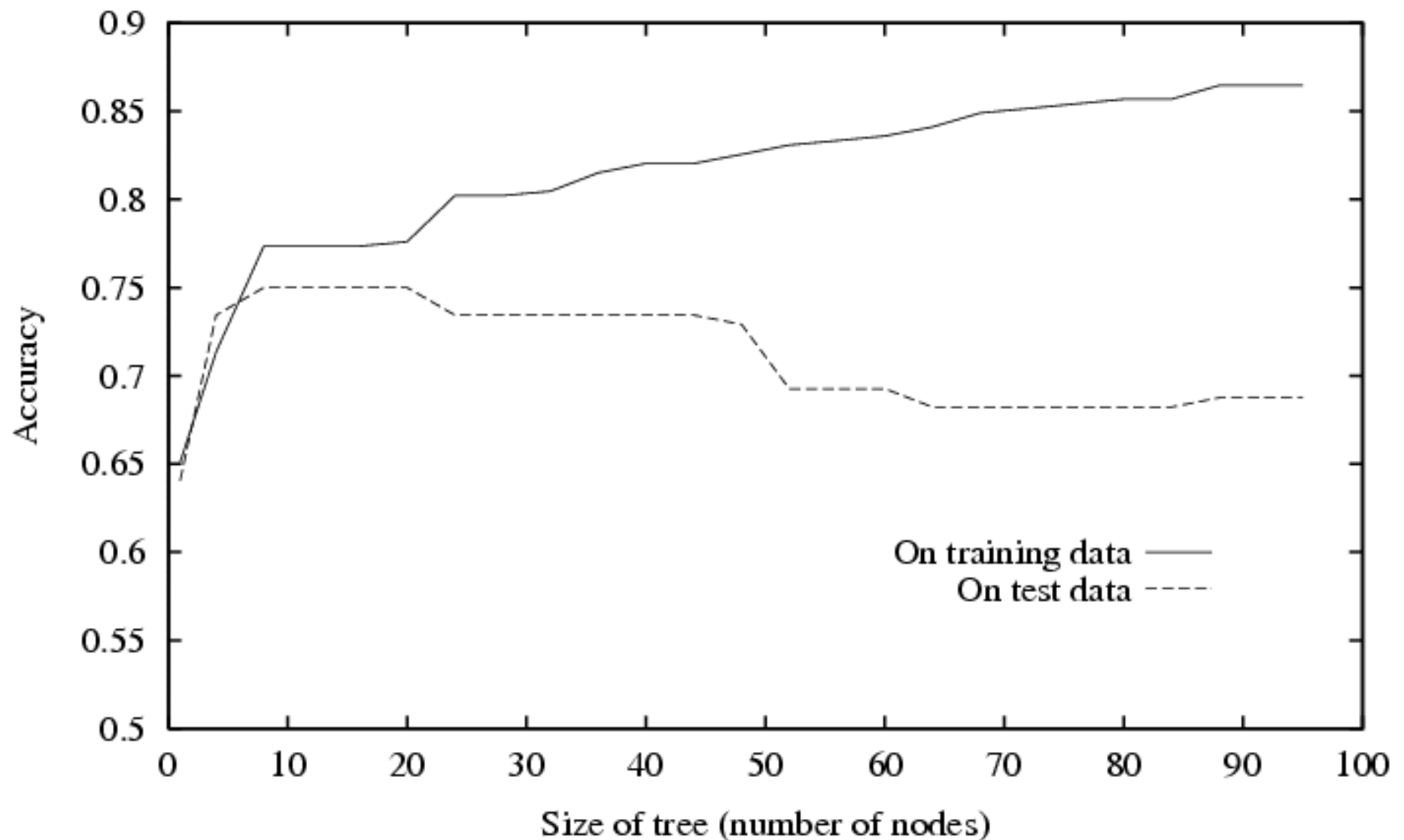


Figure from Tom Mitchell

How to Avoid Overfitting?

For Decision Trees...

1. Do not grow tree beyond some **maximum depth**
2. Do not split if splitting criterion (e.g. Info. Gain) is **below some threshold**
3. Stop growing when the split is **not statistically significant**
4. Grow the entire tree, then **prune**

Reduced-Error Pruning

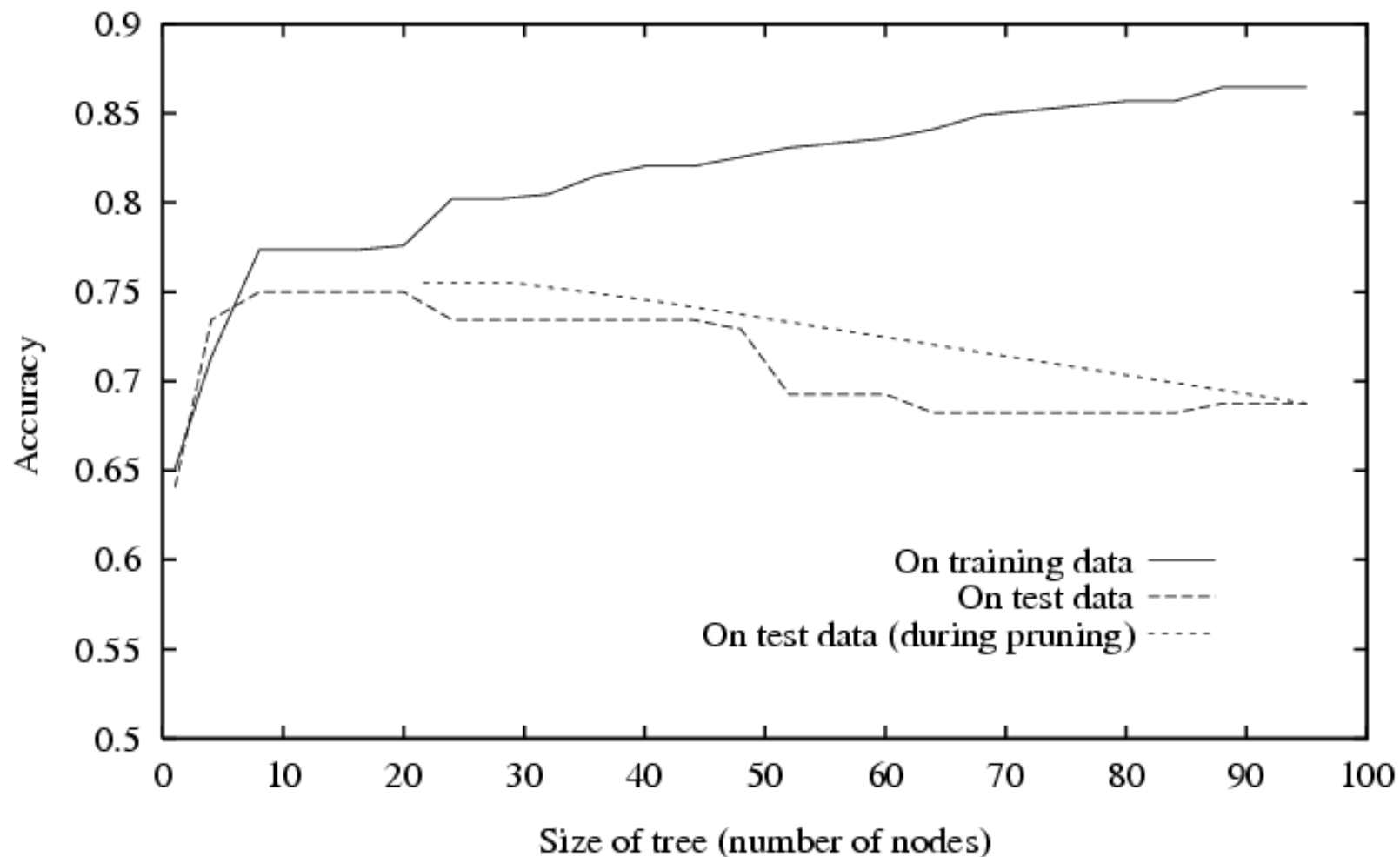
Split data into *training* and *validation* set

Create tree that classifies *training* set correctly

Do until further pruning is harmful:

1. Evaluate impact on *validation* set of pruning each possible node (plus those below it)
 2. Greedily remove the one that most improves *validation* set accuracy
- produces smallest version of most accurate subtree
 - What if data is limited?

Effect of Reduced-Error Pruning



Questions

- Will ID3 always include all the attributes in the tree?
- What if some attributes are real-valued? Can learning still be done efficiently?
- What if some attributes are missing?

Decision Trees (DTs) in the Wild

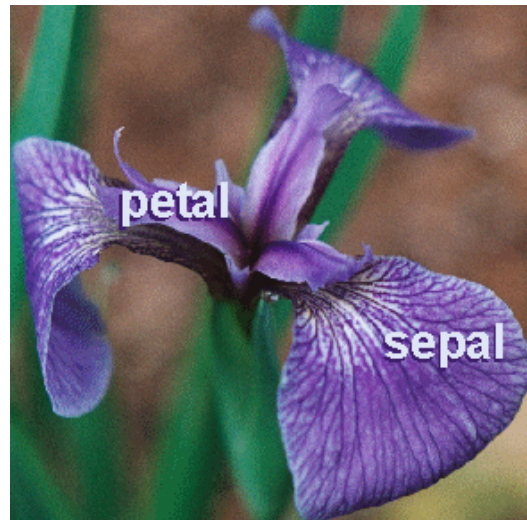
- DTs are one of the most popular classification methods for practical applications
 - Reason #1: The learned representation is **easy to explain** a non-ML person
 - Reason #2: They are **efficient** in both computation and memory
- DTs can be applied to a wide variety of problems including **classification, regression, density estimation**, etc.
- **Applications of DTs** include...
 - medicine, molecular biology, text classification, manufacturing, astronomy, agriculture, and many others
- **Decision Forests** learn many DTs from random subsets of features; the result is a very powerful example of an **ensemble method** (discussed later in the course)

DT 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

CLASSIFICATION

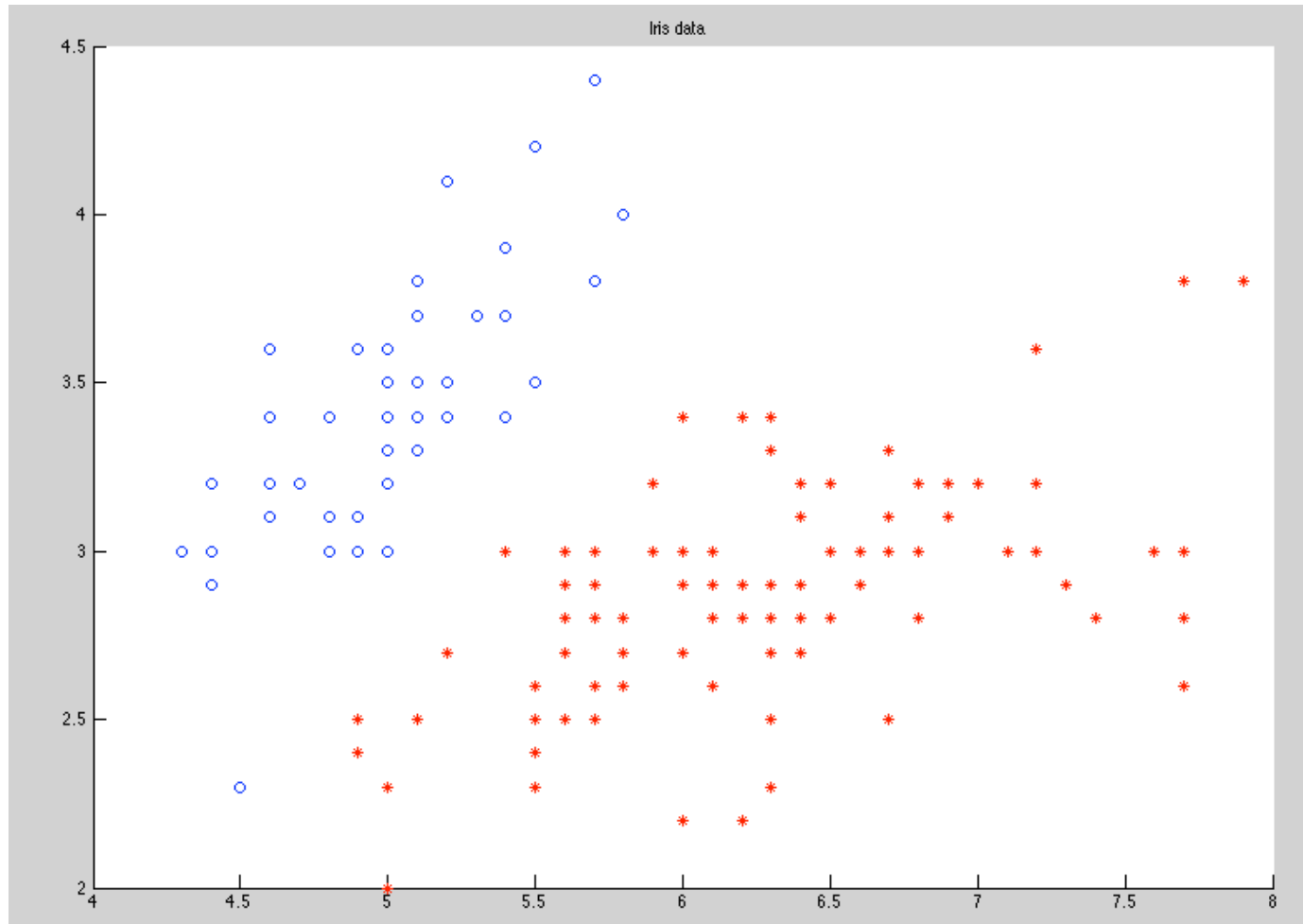


Fisher Iris Dataset

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

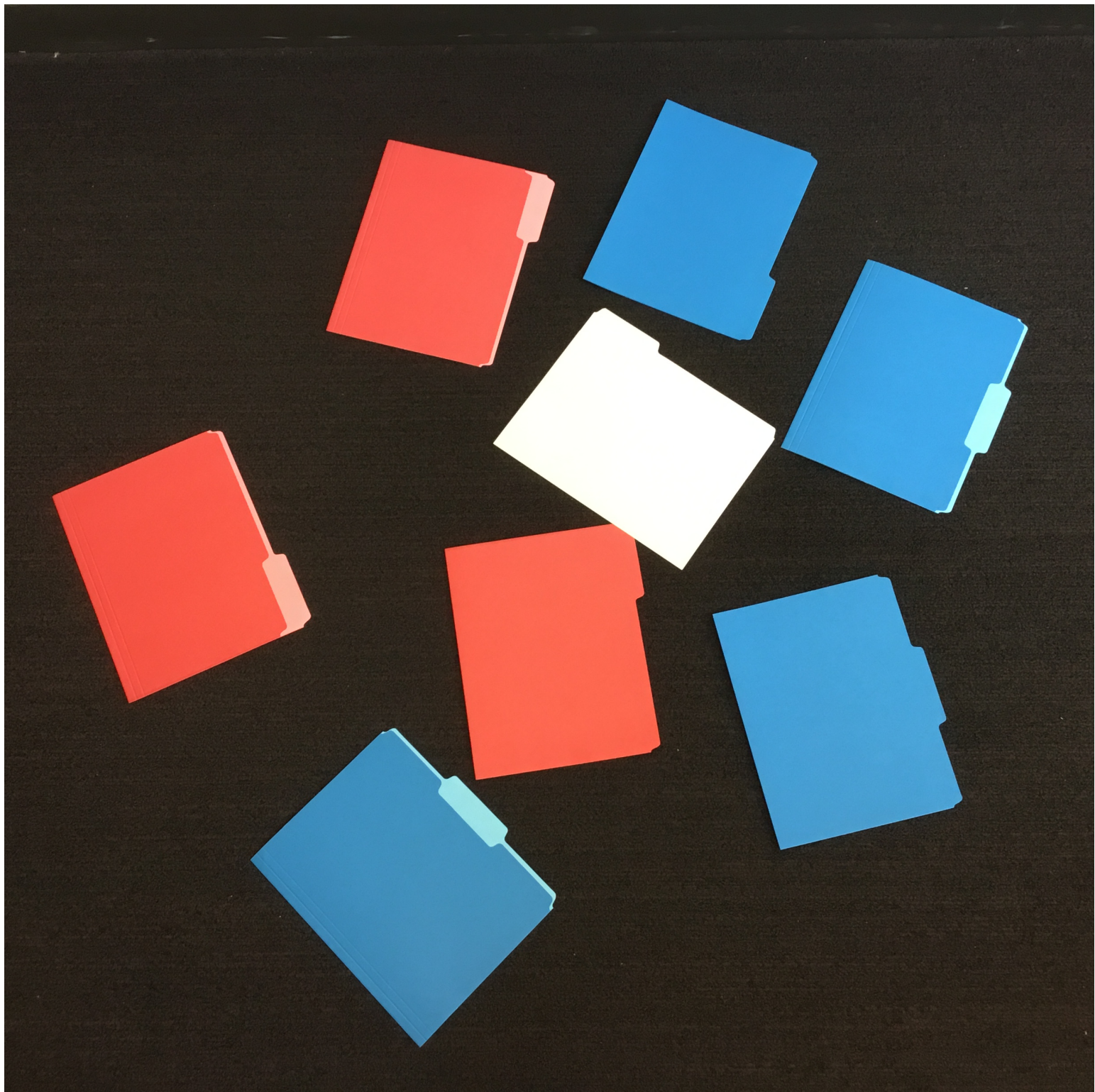


Classification

Chalkboard:

- Binary classification
- 2D examples
- Decision rules / hypotheses

K-NEAREST NEIGHBORS



k-Nearest Neighbors

Chalkboard:

- Nearest Neighbor classifier
- KNN for binary classification