

CS 360: Machine Learning

Prof. Sara Mathieson

Fall 2019



Haverford
COLLEGE

Admin

- Office hours: **TODAY 12:30-1:30pm** (L310)
 - Also make use of TA office hours Sun/Mon
- **Lab 1 due** Tues night (TODAY)
 - Make sure to **push your final code** on git
 - Turn in **math portion on paper** (slide under door)
 - There is one with no name
- Next reading quiz: Thursday on **Duame 1.3-1.6 and chapter 2**

Outline for September 10

- Recap featurization discussion
- Begin: Decision Trees
- ID3 Decision Tree algorithm
- Entropy

Outline for September 10

- Recap featurization discussion
- Begin: Decision Trees
- ID3 Decision Tree algorithm
- Entropy

Back to Handout 2

4. Using your response from the previous question, what would the *feature vector* become for \mathbf{x}_1 ?

Day	Outlook	Temperature	Humidity	Wind	PlayTennis (y)
\mathbf{x}_1	Sunny	Hot	High	Weak	No
\mathbf{x}_2	Sunny	Hot	High	Strong	No
\mathbf{x}_3	Overcast	Hot	High	Weak	Yes
\mathbf{x}_4	Rain	Mild	High	Weak	Yes
\mathbf{x}_5	Rain	Cool	Normal	Weak	Yes
\mathbf{x}_6	Rain	Cool	Normal	Strong	No
\mathbf{x}_7	Overcast	Cool	Normal	Strong	Yes
\mathbf{x}_8	Sunny	Mild	High	Weak	No
\mathbf{x}_9	Sunny	Cool	Normal	Weak	Yes
\mathbf{x}_{10}	Rain	Mild	Normal	Weak	Yes

Data from Machine Learning by Tom Mitchell (Table 3.2)

Back to Handout 2

4. Using your response from the previous question, what would the *feature vector* become for \mathbf{x}_1 ?

Day	Outlook	Temperature	Humidity	Wind	PlayTennis (y)
\mathbf{x}_1	Sunny	Hot	High	Weak	No
\mathbf{x}_2	Sunny	Hot	High	Strong	No
\mathbf{x}_3	Overcast	Hot	High	Weak	Yes
\mathbf{x}_4	Rain	Mild	High	Weak	Yes
\mathbf{x}_5	Rain	Cool	Normal	Weak	Yes
\mathbf{x}_6	Rain	Cool	Normal	Strong	No
\mathbf{x}_7	Overcast	Cool	Normal	Strong	Yes
\mathbf{x}_8	Sunny	Mild	High	Weak	No
\mathbf{x}_9	Sunny	Cool	Normal	Weak	Yes
\mathbf{x}_{10}	Rain	Mild	Normal	Weak	Yes

Sunny:	{0,1}
Overcast:	{0,1}
Rain:	{0,1}
Temperature:	{0, 1, 2} (Cool, Mild, Hot)
Humidity:	{0,1} (Normal, High)
Wind	{0,1} (Weak, Strong)

Data from Machine Learning by Tom Mitchell (Table 3.2)

Back to Handout 2

4. Using your response from the previous question, what would the *feature vector* become for x_1 ?

Day	Outlook	Temperature	Humidity	Wind	PlayTennis (y)
x_1	Sunny	Hot	High	Weak	No
x_2	Sunny	Hot	High	Strong	No
x_3	Overcast	Hot	High	Weak	Yes
x_4	Rain	Mild	High	Weak	Yes
x_5	Rain	Cool	Normal	Weak	Yes
x_6	Rain	Cool	Normal	Strong	No
x_7	Overcast	Cool	Normal	Strong	Yes
x_8	Sunny	Mild	High	Weak	No
x_9	Sunny	Cool	Normal	Weak	Yes
x_{10}	Rain	Mild	Normal	Weak	Yes

Data from Machine Learning by Tom Mitchell (Table 3.2)

Sunny: {0,1}
 Overcast: {0,1}
 Rain: {0,1}
 Temperature: {0, 1, 2} (Cool, Mild, Hot)
 Humidity: {0,1} (Normal, High)
 Wind {0,1} (Weak, Strong)

	Sunny	Overcast	Rain	Temp	Humidity	Wind
x_1	1	0	0	2	1	0

Featurization (rule of thumb)

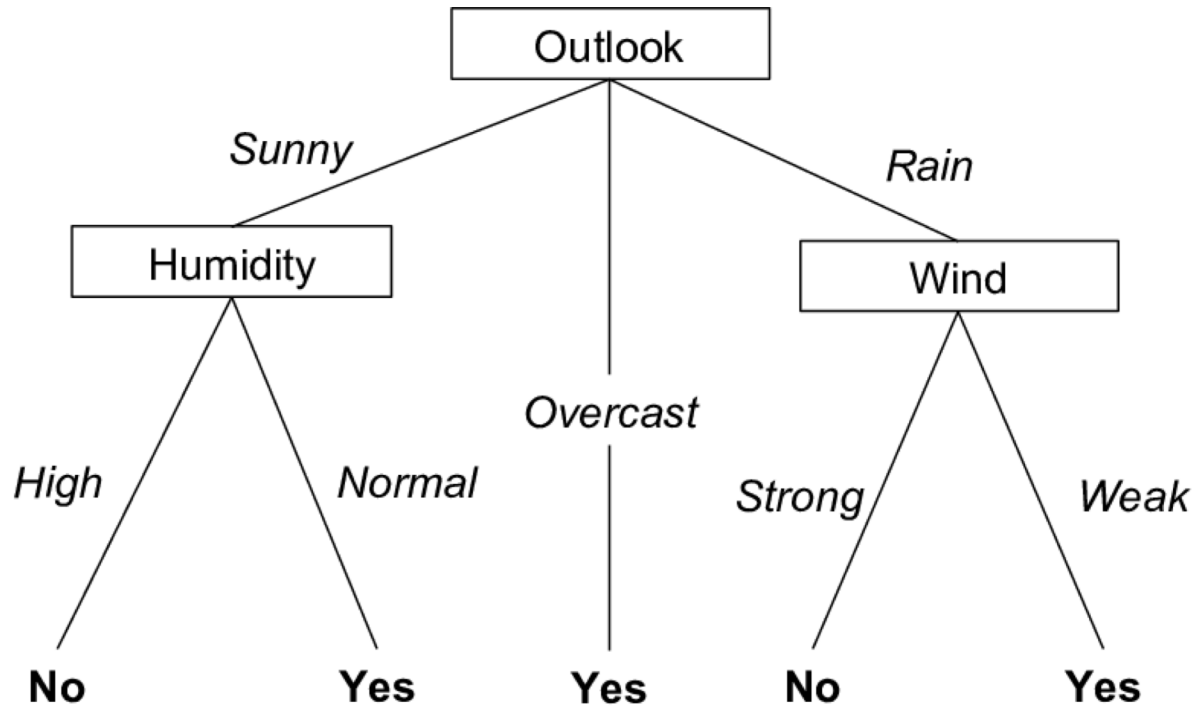
- Real-valued features get copied directly. *Duame, Chap 3*
- Binary features become 0 (for false) or 1 (for true).
- Categorical features with V possible values get mapped to V -many binary indicator features.

Haven't discussed:
-normalization

Outline for September 10

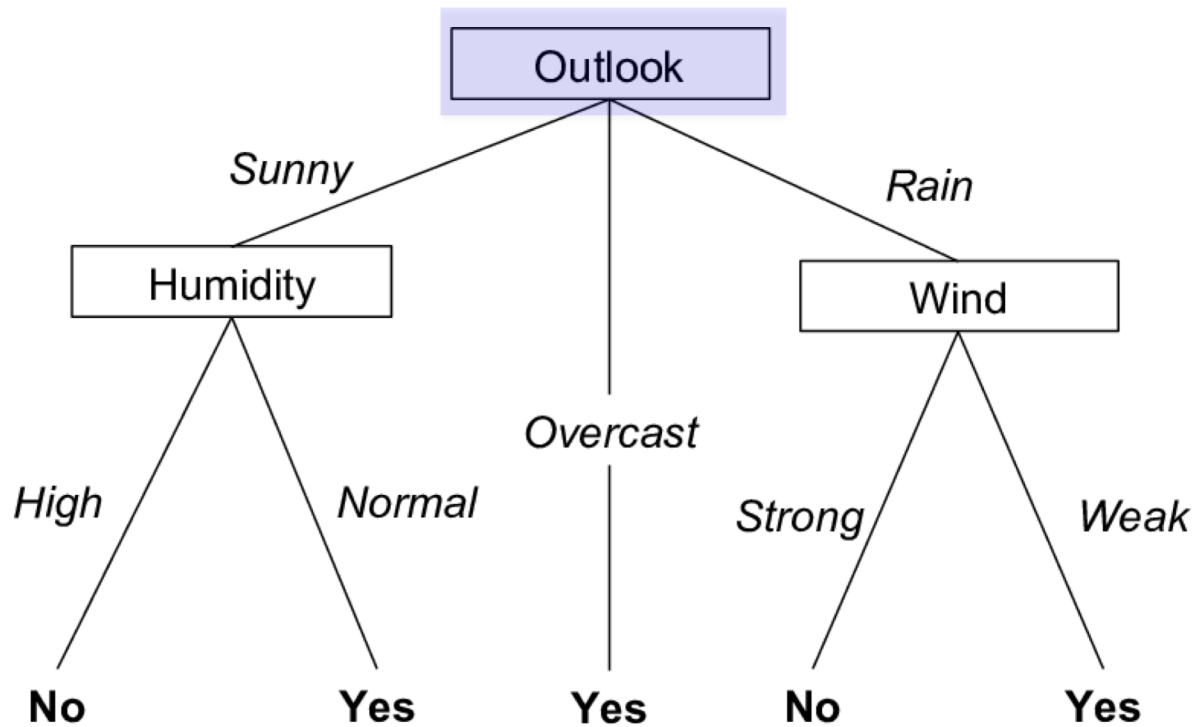
- Recap featurization discussion
- **Begin: Decision Trees**
- ID3 Decision Tree algorithm
- Entropy

Decision Tree example (tennis data)



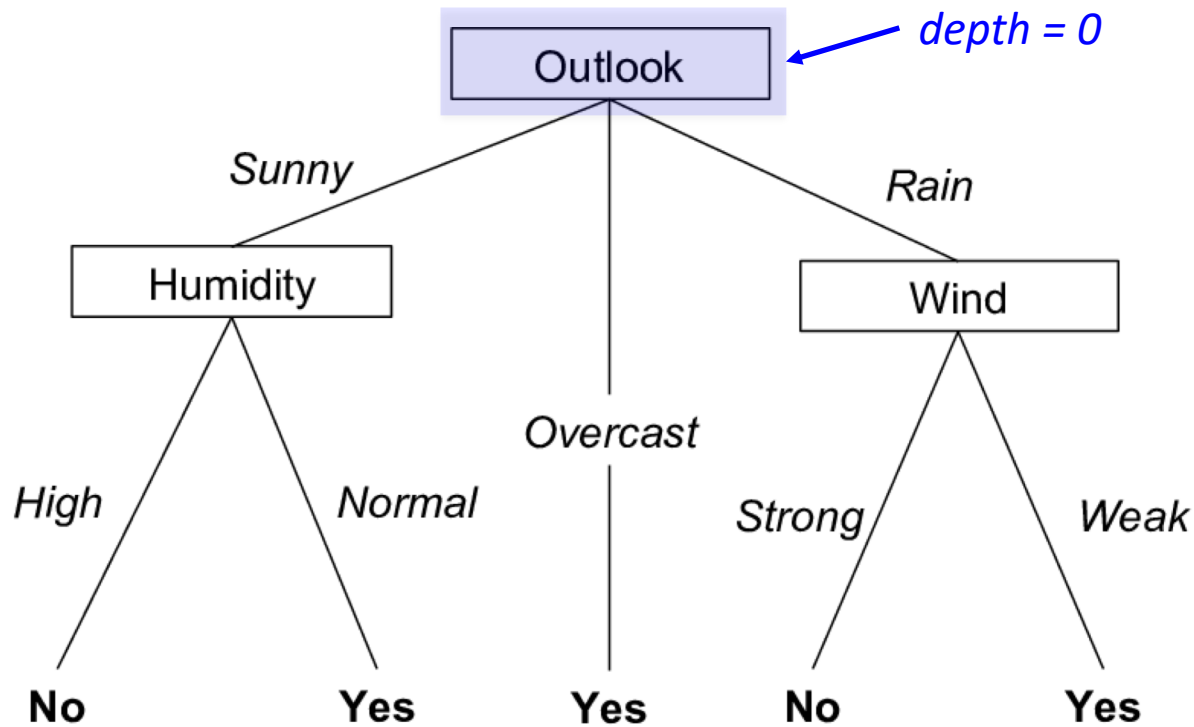
- Each internal node: test one feature
- Each branch from node: selects one value of the feature
- Each leaf node: predict y

Decision Tree example (tennis data)



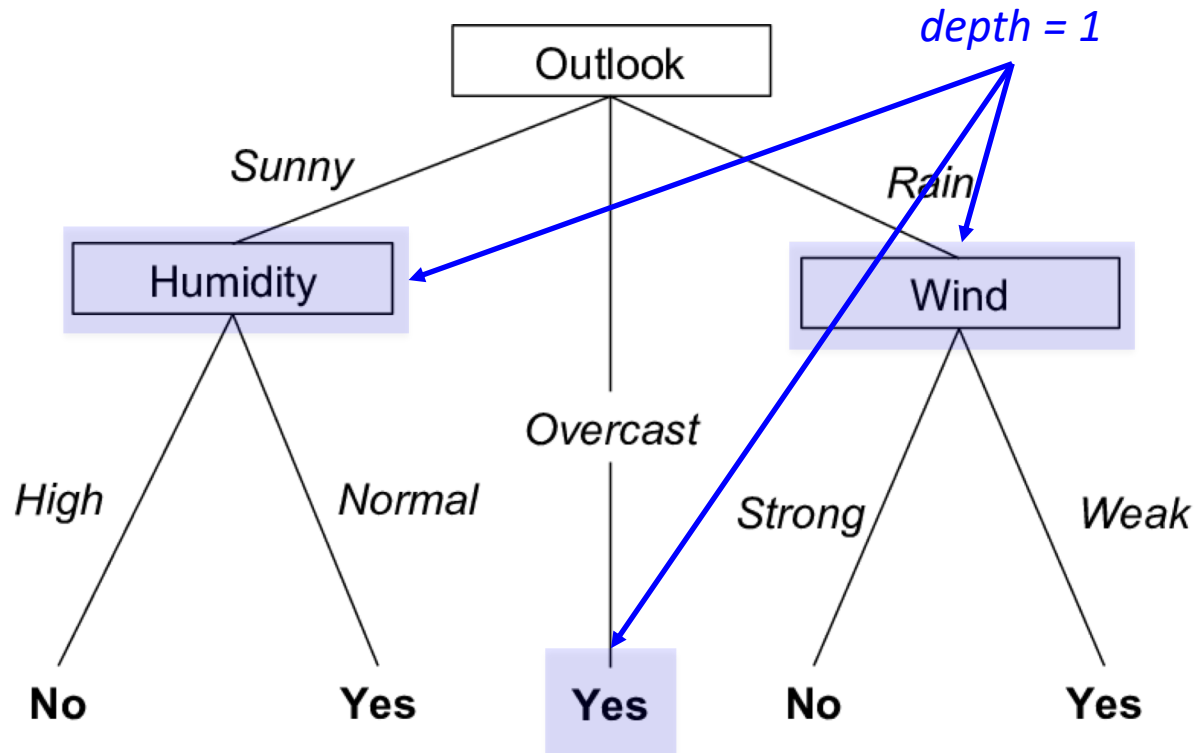
Key term: *depth*

Decision Tree example (tennis data)



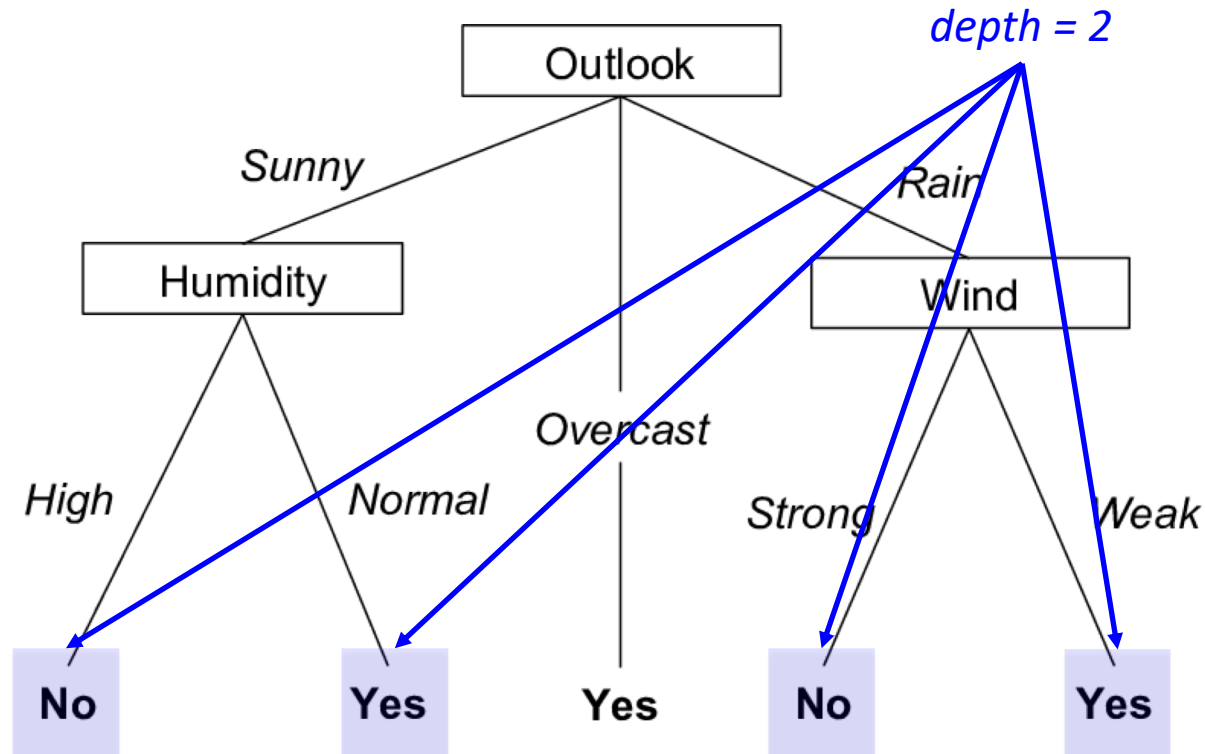
Key term: *depth*

Decision Tree example (tennis data)



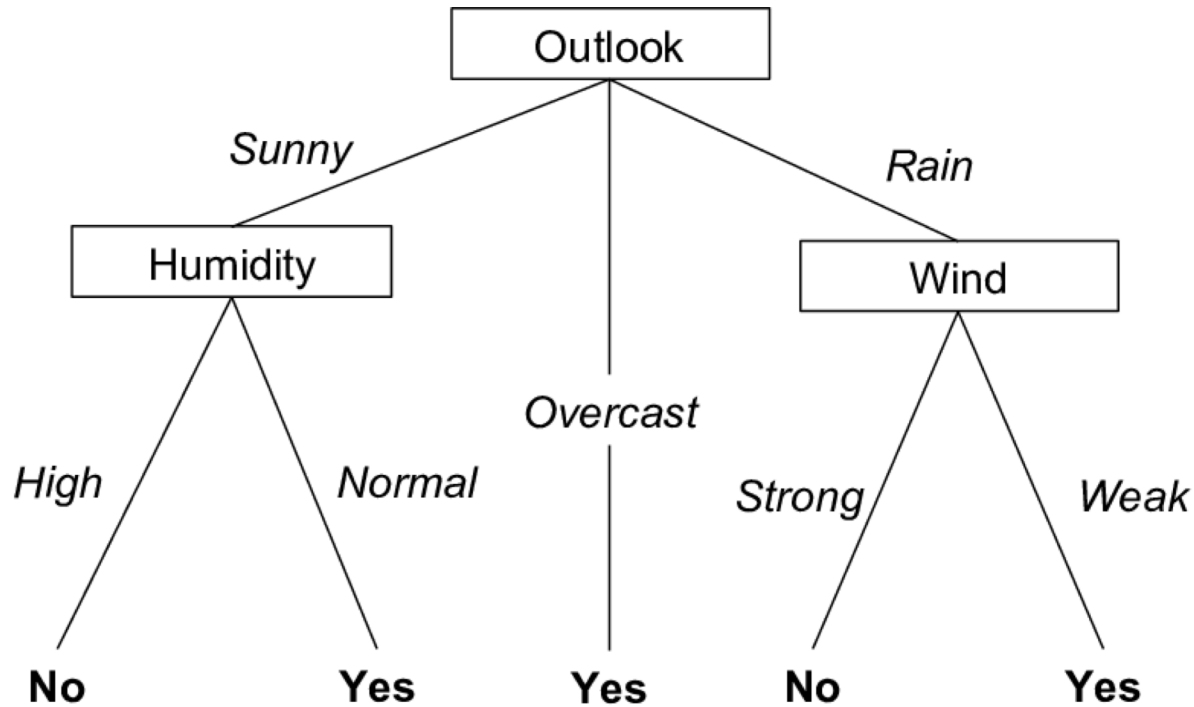
Key term: *depth*

Decision Tree example (tennis data)



Key term: *depth*

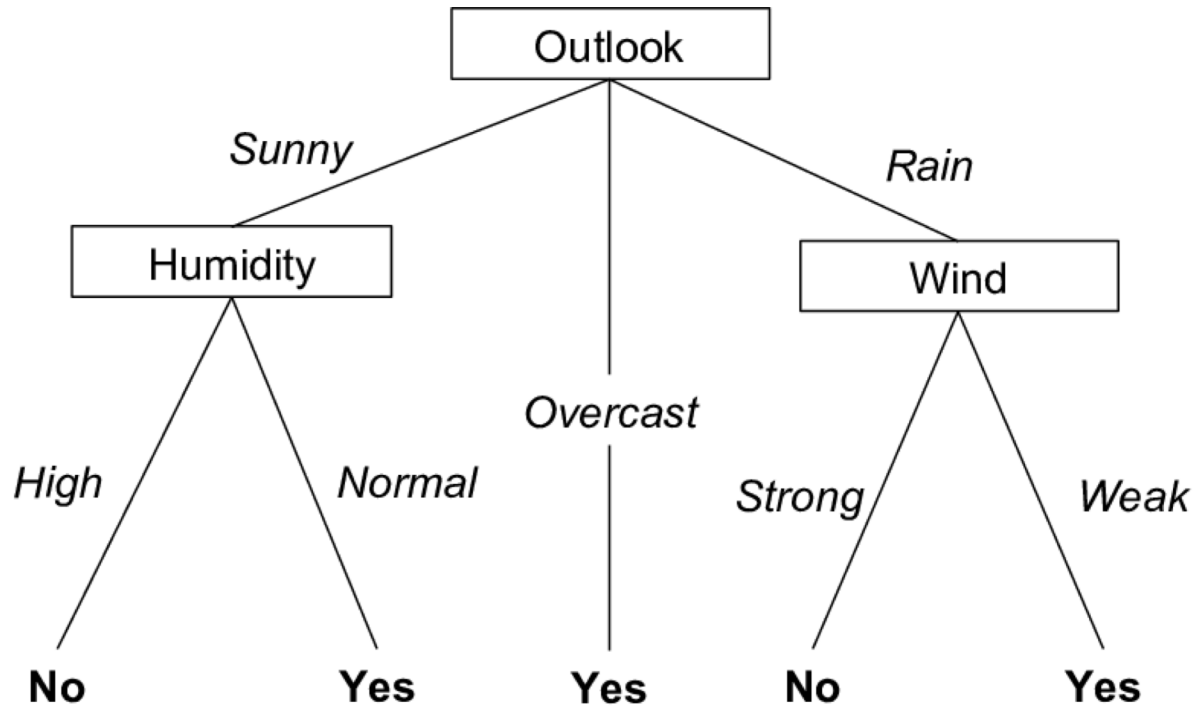
Decision Tree example (tennis data)



(test example) $x =$

Outlook	Temp	Humidity	Wind
Rain	Hot	High	Strong

Decision Tree example (tennis data)

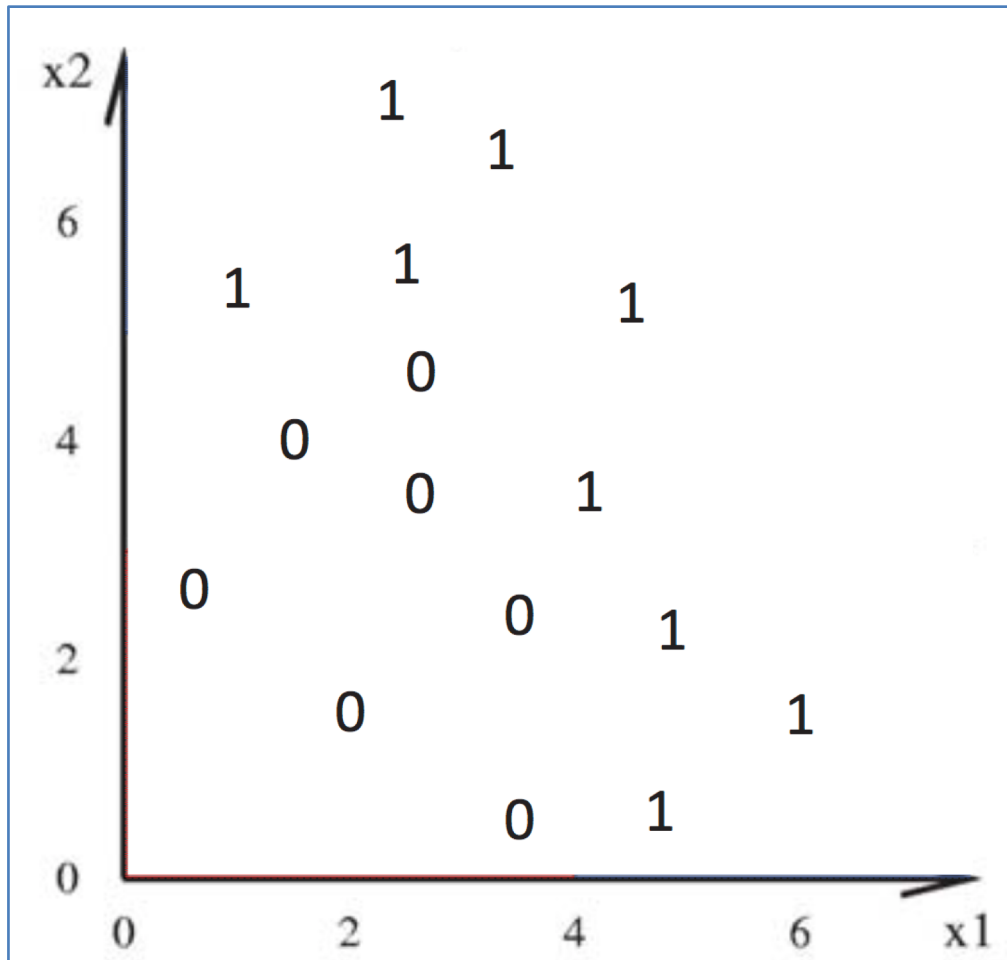


(test example) $x =$

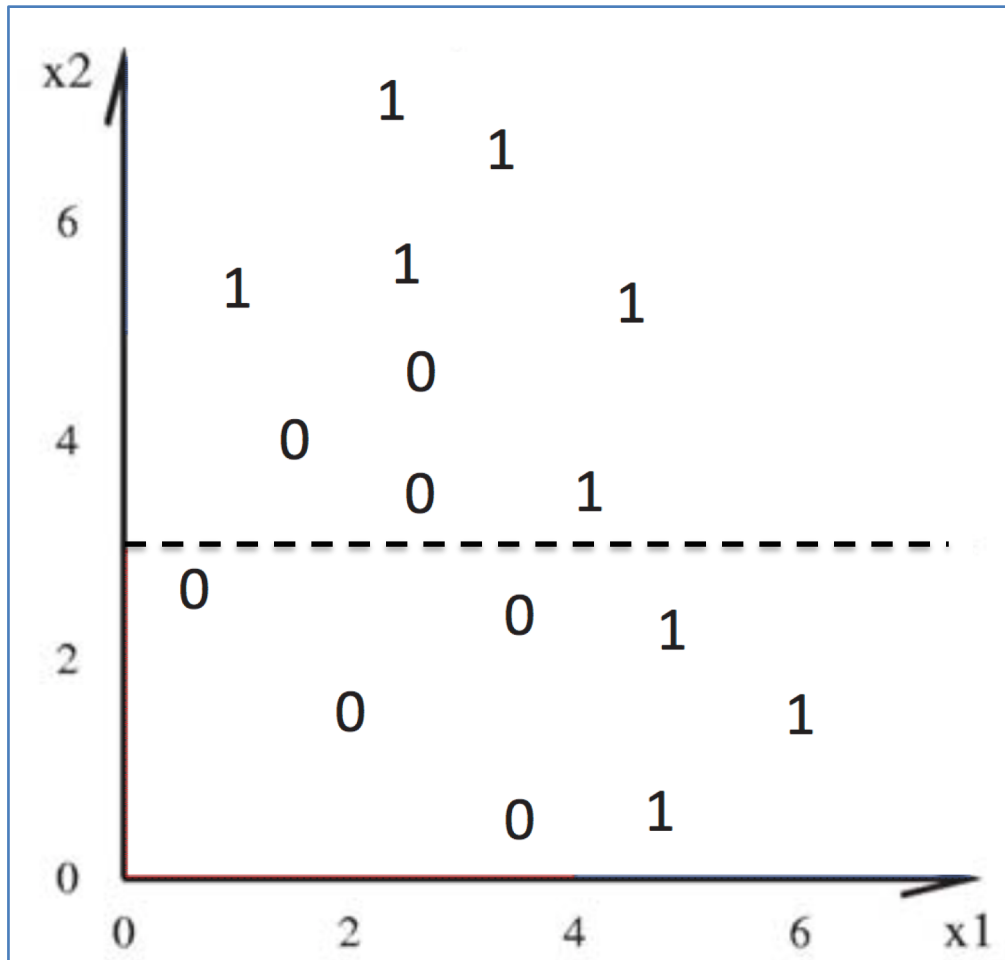
Outlook	Temp	Humidity	Wind
Rain	Hot	High	Strong

$y_{pred} = \text{No}$

Can also consider continuous features

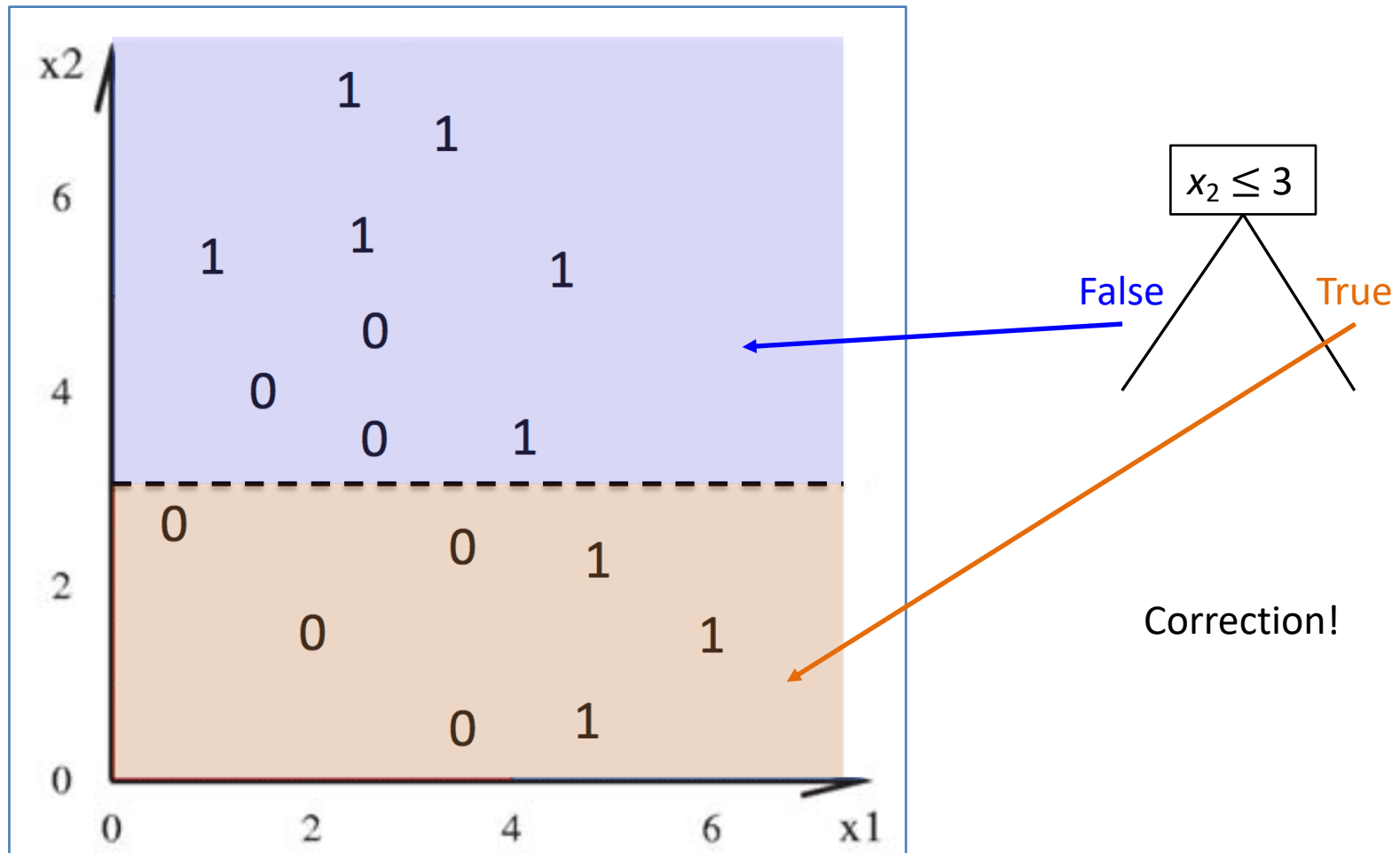


Can also consider continuous features



$$x_2 \leq 3$$

Can also consider continuous features



Decision Tree pros/cons

Discuss with a partner! Think about:

- * training and testing
- * featurization
- * runtime
- * human factors

ICATION

ANIZATION

SYSTEMS

REGEMONY

* faster
testing

* & tree
- depth
is the
test runtime

* Visualize
interpret

* may only
see few
features.

* train could
be slow

Decision Tree pros/cons

- Very interpretable! Easy to say *why* we made a classification (can point to which features)
- Compact representation and fast predictions
- Can be brittle (not looking at each example holistically)
- Featurization and implementation difficulties

Outline for September 10

- Recap featurization discussion
- Begin: Decision Trees
- ID3 Decision Tree algorithm
- Entropy

ID3 Decision Tree algorithm (1986)

- Select feature that “best” informs label prediction (i.e. y)
- **Divide**: partition data into branches based on their value at this feature
- **Conquer**: recurse on each partition

Posted as optional reading

Machine Learning 1: 81–106, 1986
© 1986 Kluwer Academic Publishers, Boston — Manufactured in The Netherlands

Induction of Decision Trees

J.R. QUINLAN (munarni!nswitgould.oz!quinlan@seismo.css.gov)
Centre for Advanced Computing Sciences, New South Wales Institute of Technology, Sydney 2007, Australia

(Received August 1, 1985)

Key words: classification, induction, decision trees, information theory, knowledge acquisition, expert systems

Top-Down decision tree algorithm


MakeSubtree(D, F)

- if stopping criteria met

 - make a leaf node N

 - determine class label/probabilities for N


Top-Down decision tree algorithm

 Dataset (X,y)

MakeSubtree(D, F)

- if stopping criteria met
- make a leaf node N
- determine class label/probabilities for N

Top-Down decision tree algorithm

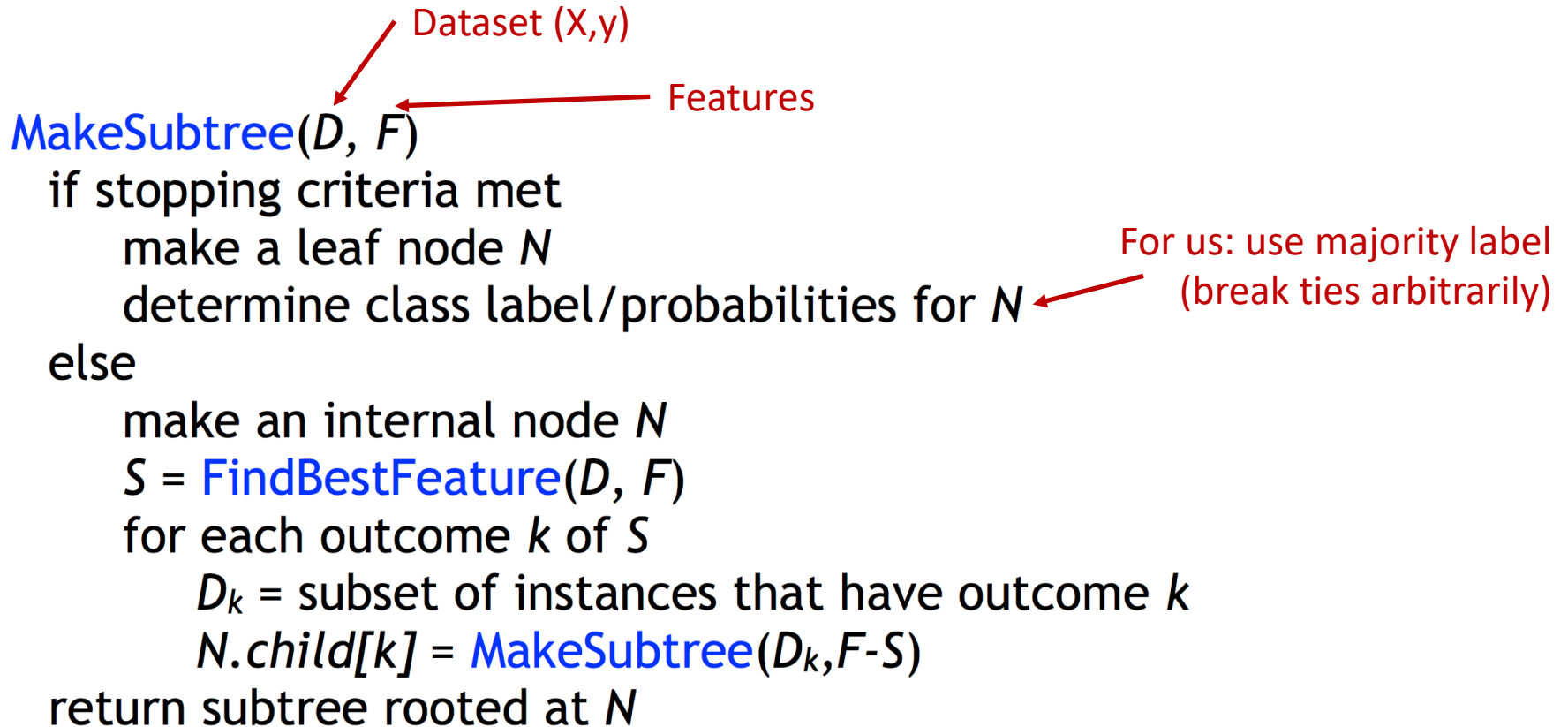
 **MakeSubtree**(D , F)
if stopping criteria met
 make a leaf node N
 determine class label/probabilities for N

Top-Down decision tree algorithm

MakeSubtree(D , F)

- Dataset (X, y)
- Features
- if stopping criteria met
 - make a leaf node N
 - determine class label/probabilities for N For us: use majority label (break ties arbitrarily)

Top-Down decision tree algorithm

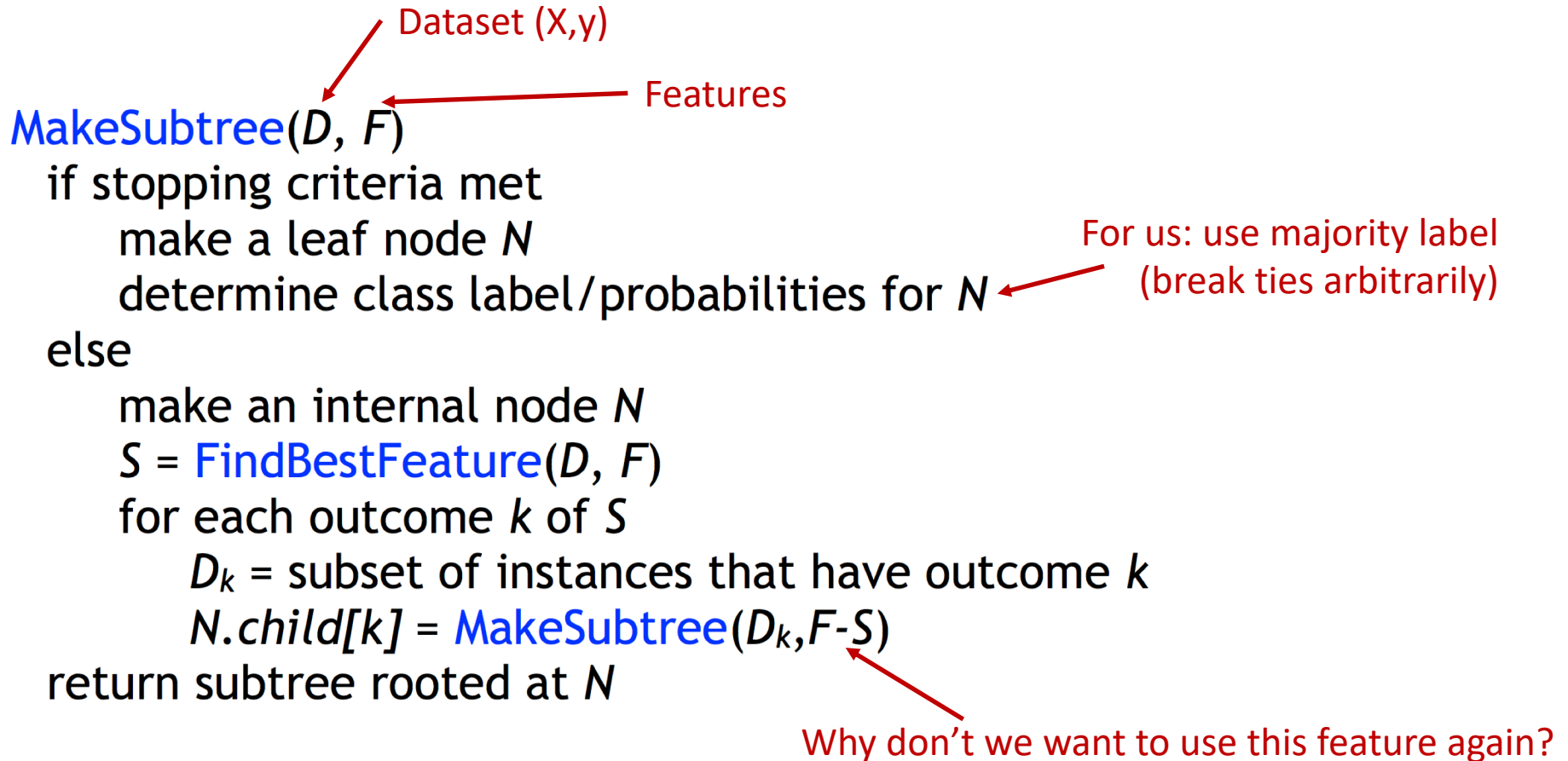


MakeSubtree(D, F)

- if stopping criteria met
 - make a leaf node N
 - determine class label/probabilities for N For us: use majority label (break ties arbitrarily)
- else
 - make an internal node N
 - $S = \text{FindBestFeature}(D, F)$
 - for each outcome k of S
 - D_k = subset of instances that have outcome k
 - $N.\text{child}[k] = \text{MakeSubtree}(D_k, F-S)$

return subtree rooted at N

Top-Down decision tree algorithm

The diagram illustrates the Top-Down decision tree algorithm. It starts with a function call `MakeSubtree(D, F)` in blue. A red arrow points from the text 'Dataset (X,y)' to the parameter `D`. Another red arrow points from the text 'Features' to the parameter `F`. The algorithm proceeds with an if-statement: 'if stopping criteria met', followed by 'make a leaf node N ', and 'determine class label/probabilities for N '. A red arrow points from the text 'For us: use majority label (break ties arbitrarily)' to the 'determine class label/probabilities for N ' step. If the stopping criteria are not met, it goes to the 'else' branch: 'make an internal node N ', ' $S = \text{FindBestFeature}(D, F)$ ', and 'for each outcome k of S '. Inside this loop, it defines ' $D_k = \text{subset of instances that have outcome } k$ ' and ' $N.\text{child}[k] = \text{MakeSubtree}(D_k, F-S)$ '. A red arrow points from the text 'Why don't we want to use this feature again?' to the recursive call `MakeSubtree(D_k, F-S)`. Finally, it 'return subtree rooted at N '.

`MakeSubtree(D, F)`

Dataset (X, y)

Features

if stopping criteria met

make a leaf node N

determine class label/probabilities for N

For us: use majority label
(break ties arbitrarily)

else

make an internal node N

$S = \text{FindBestFeature}(D, F)$

for each outcome k of S

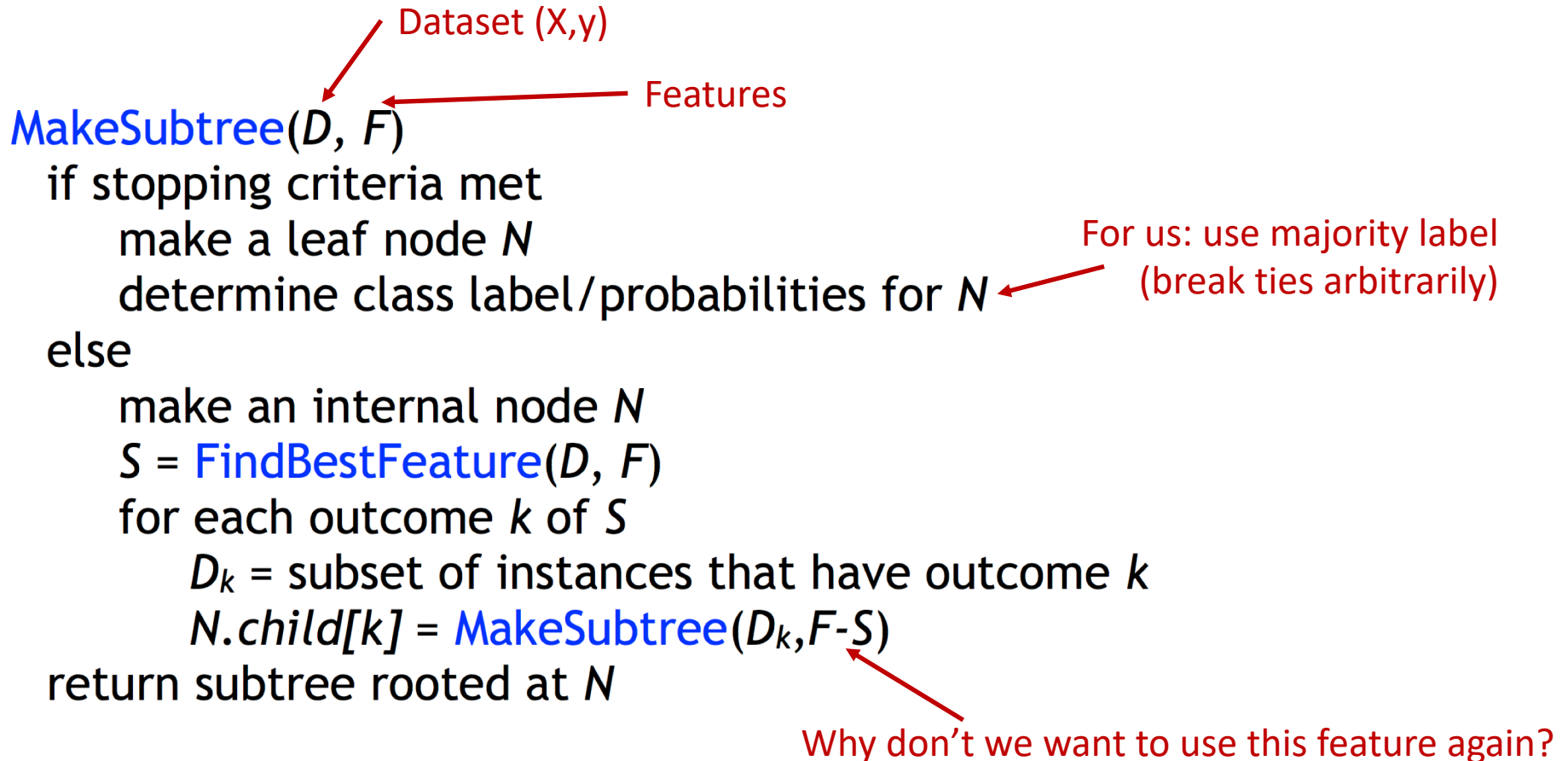
$D_k = \text{subset of instances that have outcome } k$

$N.\text{child}[k] = \text{MakeSubtree}(D_k, F-S)$

Why don't we want to use this feature again?

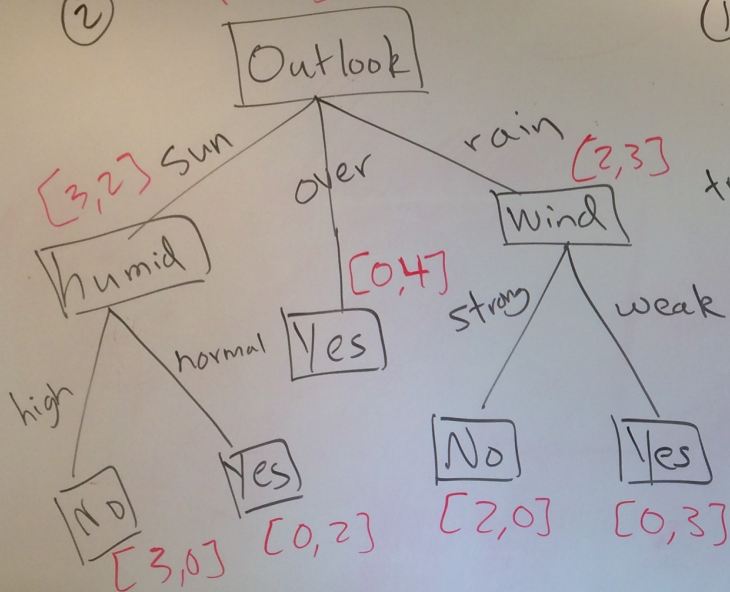
return subtree rooted at N

Top-Down decision tree algorithm


MakeSubtree(D, F)
if stopping criteria met
 make a leaf node N
 determine class label/probabilities for N
else
 make an internal node N
 $S = \text{FindBestFeature}(D, F)$
 for each outcome k of S
 $D_k =$ subset of instances that have outcome k
 $N.\text{child}[k] = \text{MakeSubtree}(D_k, F-S)$
return subtree rooted at N

**Now: Handout 3 + think about:
what design choices do we need to make?**

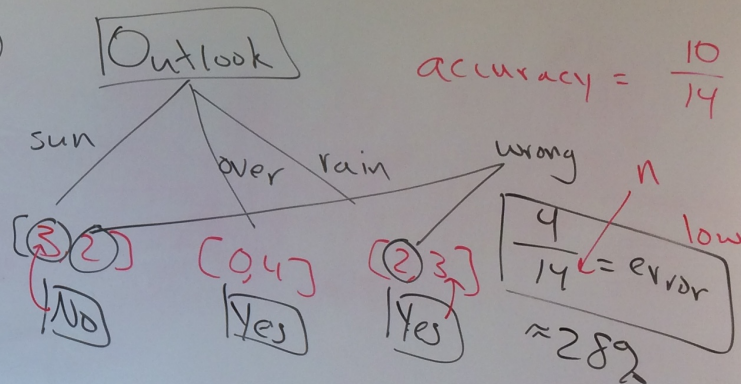
② $no \rightarrow [5, 9] \leftarrow yes$



① $n = 14$
 $p = 4$
 $error = 0\%$
 train \nearrow

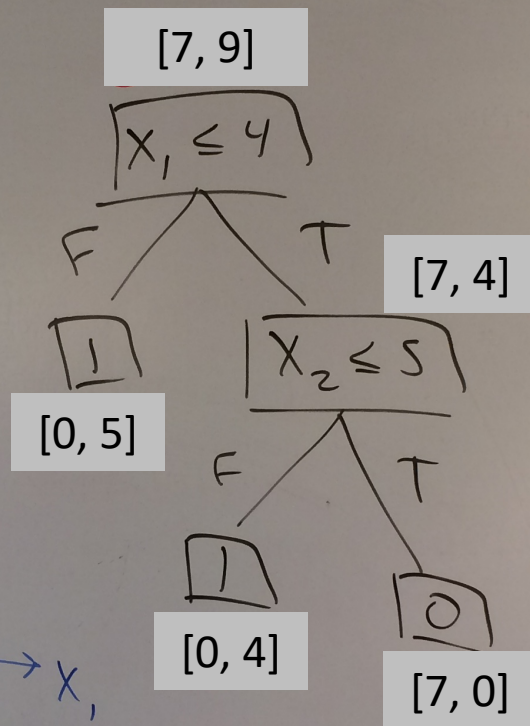
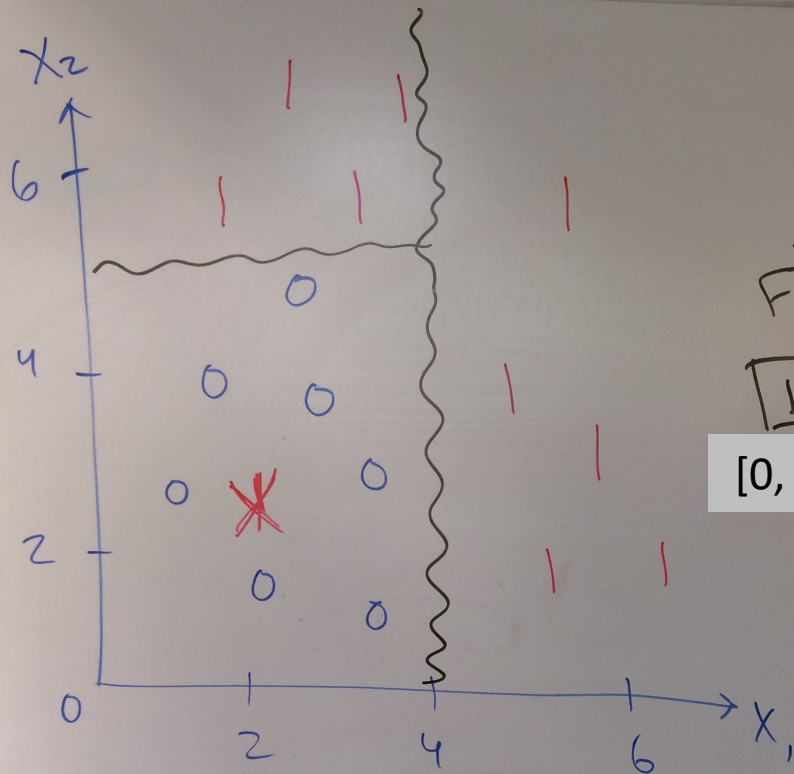
Big Idea
 choose the best feature

③



accuracy = $\frac{10}{14}$

$\frac{4}{14} = error$
 $\approx 28\%$
 low



Design choice: stopping criteria

1. All the data points in our partition have the same label

Design choice: stopping criteria

1. All the data points in our partition have the same label
2. No more features remain to split on

Design choice: stopping criteria

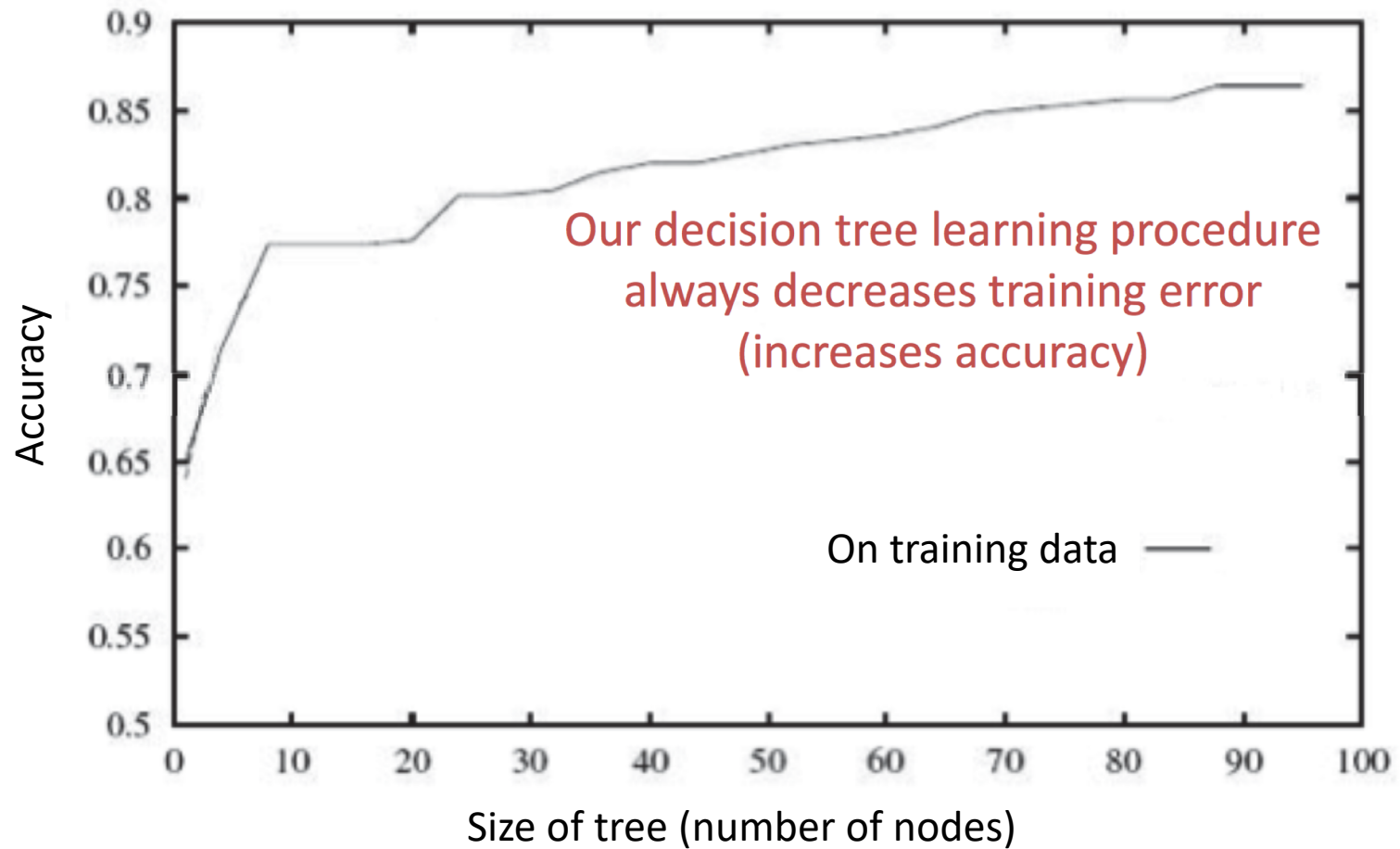
1. All the data points in our partition have the same label
2. No more features remain to split on
3. No features are informative about the label

Design choice: stopping criteria

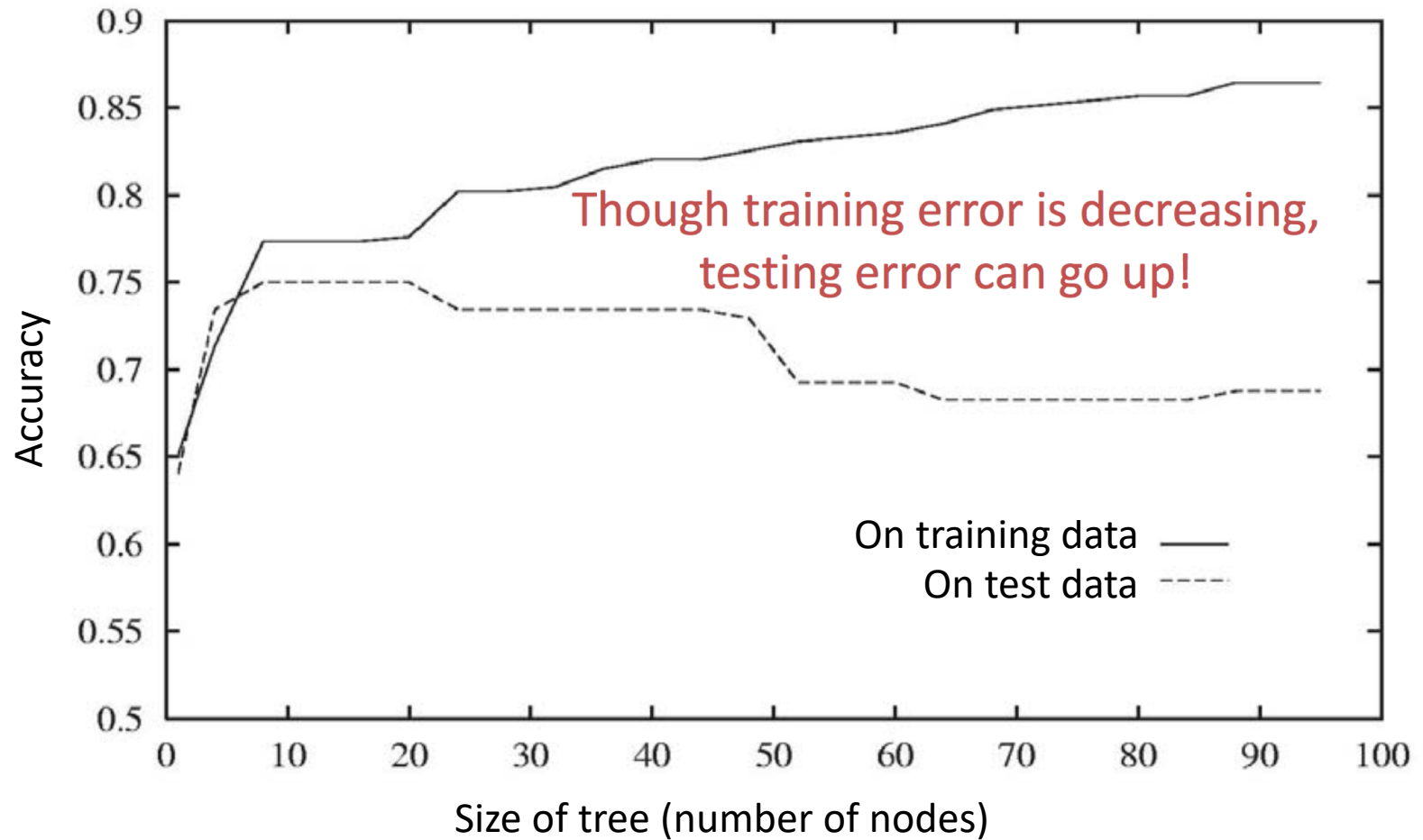
1. All the data points in our partition have the same label
2. No more features remain to split on
3. No features are informative about the label
4. Reached (user specified) max depth in the tree

Detour to overfitting

Overfitting



Overfitting



Overfitting definition

- Consider a hypothesis (tree): h
 - Training error: $error_{train}(h)$
 - Error over all possible data: $error_D(h)$

Overfitting definition

- Consider a hypothesis (tree): h
 - Training error: $error_{train}(h)$
 - Error over all possible data: $error_D(h)$
- A hypothesis h **overfits** training data if there exists another hypothesis h' s.t.
 - $error_{train}(h) < error_{train}(h')$ AND
 - $error_D(h) > error_D(h')$

Avoiding overfitting for us

- Stop when leaf label reaches a certain fraction (i.e. 95% “yes”, 5% “no”)

For our Lab 2 implementation

- Set a maximum depth for the tree
- Set a minimum number of examples in leaf (i.e. if we have a 2-1 split, stop)

Outline for September 10

- Recap featurization discussion
- Begin: Decision Trees
- ID3 Decision Tree algorithm
- Entropy

How to select “best” features?

X

Color	Shape	Size
red	square	big
blue	square	big
red	circle	small
blue	square	small
red	circle	big

Y

Likes toy?
+
+
-
-
+

How to select “best” features?

X

Y

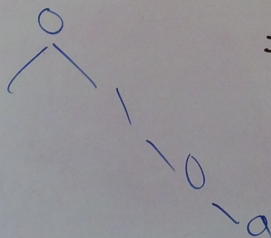
Color	Shape	Size
red	square	big
blue	square	big
red	circle	big
blue	square	big
red	circle	big

Likes toy?
+
+
-
-
+

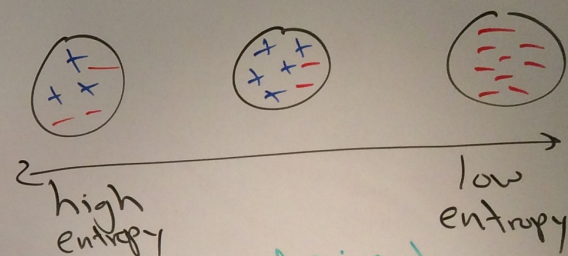
Entropy: avg # bits
needed to transmit info

year	prob (p)	idea	(sort) cumulative prob	binary
senior	0.5	0	0	0.0000...
junior	0.25	1	0.5	0.1000...
soph	0.125	01	0.75	0.1100...
first	0.125	10	0.875	0.1110...

11011010001110



intuition



$$\begin{aligned}
 & 0.5 \cdot 2^0 + 0.25 \cdot 2^1 + 0.125 \cdot 2^2 + 0.125 \cdot 2^3 + \dots \\
 & S = 1.4 + 0.2 + 1.1 \\
 & \Rightarrow 101
 \end{aligned}$$

decimal point

$\frac{1}{2}$ $\frac{1}{4}$