

# Decision Trees

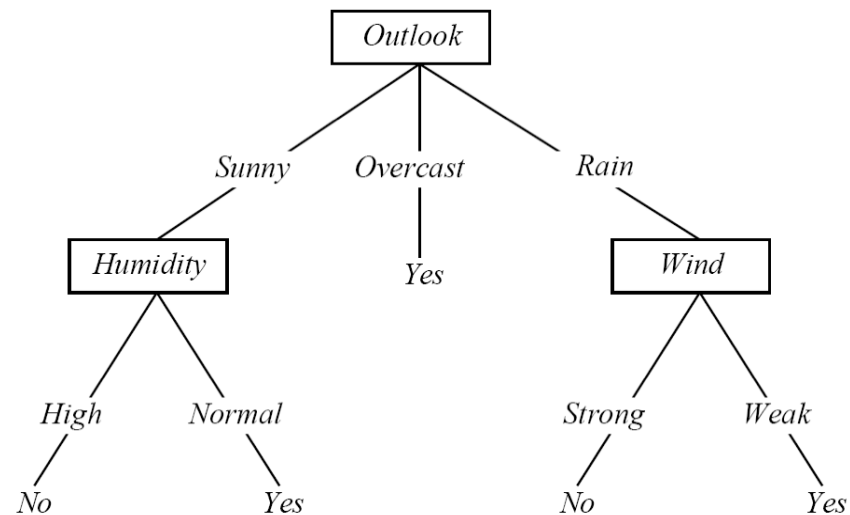
- Learn from labeled observations - supervised learning
- Represent the knowledge learned in form of a tree

Example: learning when to play tennis.

- Examples/observations are days with their observed characteristics and whether we played tennis or not

# Play Tennis Example

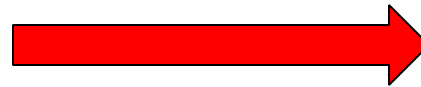
Outlook	Temperature	Humidity	Windy	PlayTennis
Sunny	Hot	High	False	No
Sunny	Hot	High	True	No
Overcast	Hot	High	False	Yes
Rainy	Mild	High	False	Yes
Rainy	Cool	Normal	False	Yes
Rainy	Cool	Normal	True	No
Overcast	Cool	Normal	True	Yes
Sunny	Mild	High	False	No
Sunny	Cool	Normal	False	Yes
Rainy	Mild	Normal	False	Yes
Sunny	Mild	Normal	True	Yes
Overcast	Mild	High	True	Yes
Overcast	Hot	Normal	False	Yes
Rainy	Mild	High	True	No



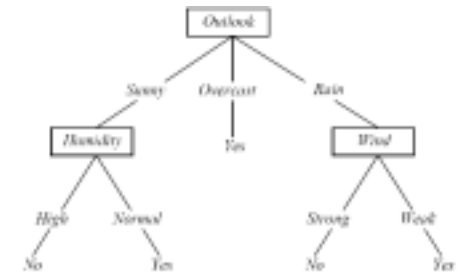
# Decision Tree Learning

Outlook	Temperature	Humidity	Windy	PlayTennis
Sunny	Hot	High	False	No
Sunny	Hot	High	True	No
Overcast	Hot	High	False	Yes
Rainy	Mild	High	False	Yes
Rainy	Cool	Normal	False	Yes
Rainy	Cool	Normal	True	No
Overcast	Cool	Normal	True	Yes
Sunny	Mild	High	False	No
Sunny	Cool	Normal	False	Yes
Rainy	Mild	Normal	False	Yes
Sunny	Mild	Normal	True	Yes
Overcast	Mild	High	True	Yes
Overcast	Hot	Normal	False	Yes
Rainy	Mild	High	True	No

Facts or Observations



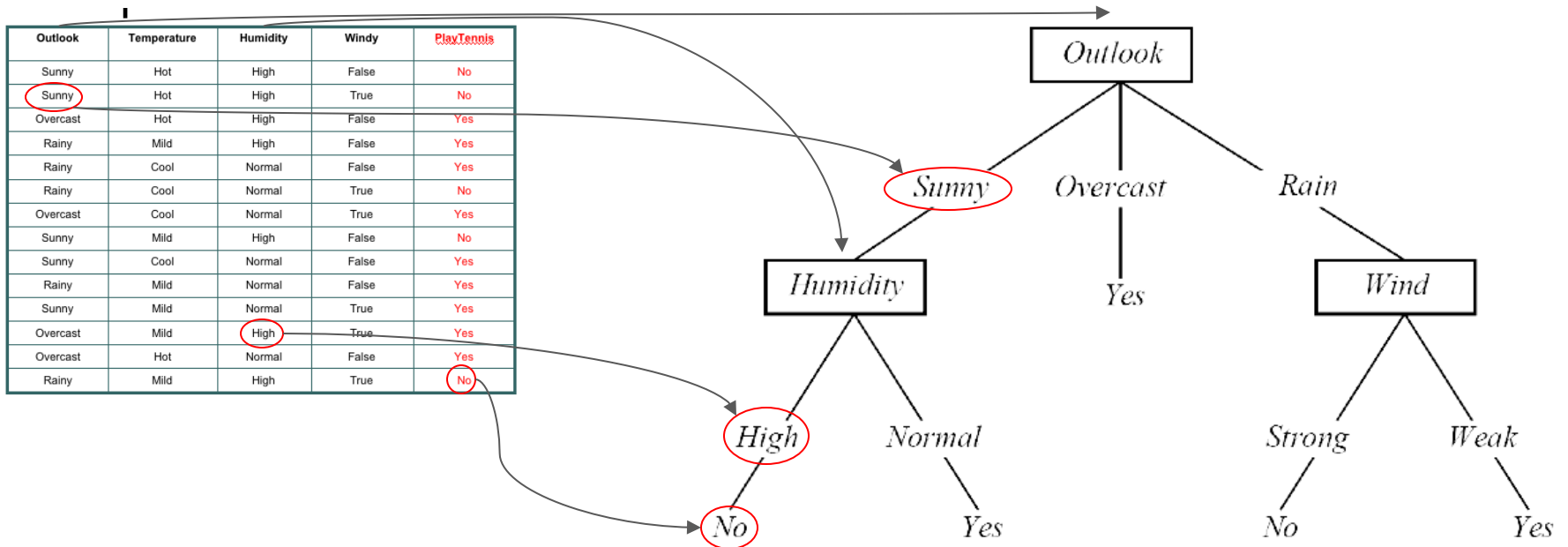
Induction



Theory

# Interpreting a DT

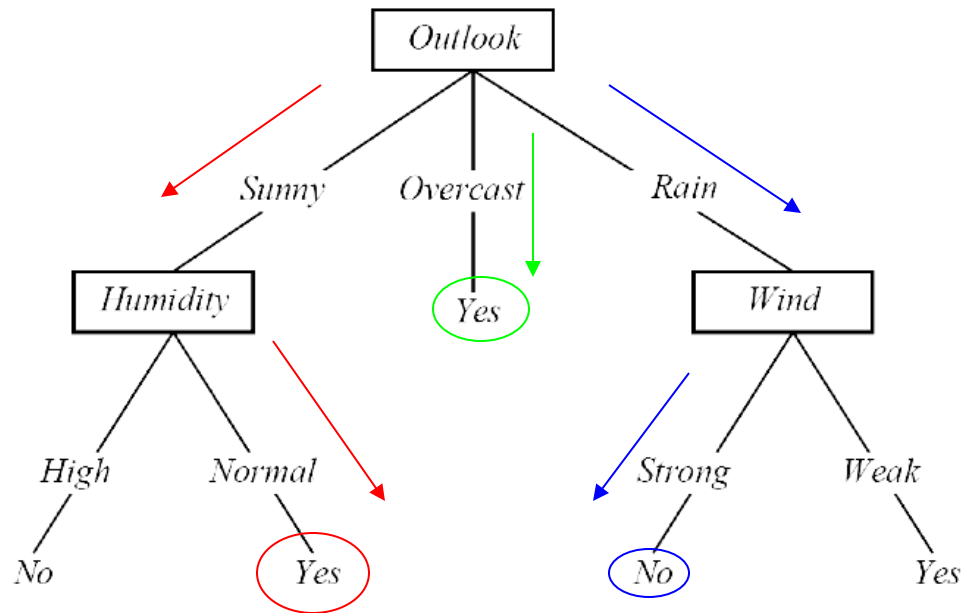
DT  $\equiv$  Decision  
Tree



- A DT uses the features of an observation table as nodes and the feature values as links.
- All feature values of a particular feature need to be represented as links.
- The target feature is special - its values show up as leaf nodes in the DT.

# Interpreting a DT

Each path from the root of the DT to a leaf can be interpreted as a decision rule.



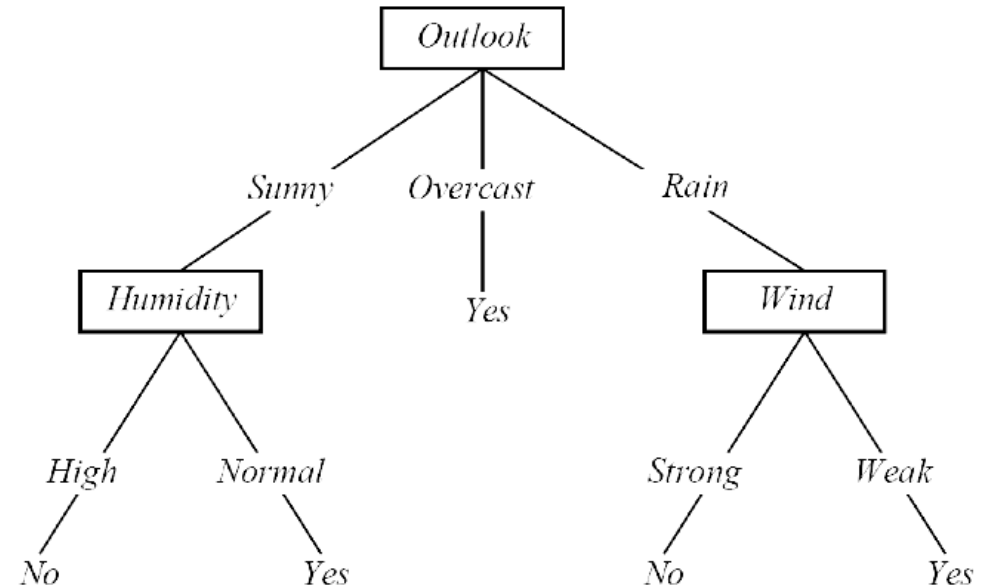
IF *Outlook* = *Sunny* AND *Humidity* = *Normal* THEN *PlayTennis* = *Yes*

IF *Outlook* = *Overcast* THEN *PlayTennis* = *Yes*

IF *Outlook* = *Rain* AND *Wind* = *Strong* THEN *PlayTennis* = *No*

# DT: Explanation & Prediction

Outlook	Temperature	Humidity	Windy	PlayTennis
Sunny	Hot	High	False	No
Sunny	Hot	High	True	No
Overcast	Hot	High	False	Yes
Rainy	Mild	High	False	Yes
Rainy	Cool	Normal	False	Yes
Rainy	Cool	Normal	True	No
Overcast	Cool	Normal	True	Yes
Sunny	Mild	High	False	No
Sunny	Cool	Normal	False	Yes
Rainy	Mild	Normal	False	Yes
Sunny	Mild	Normal	True	Yes
Overcast	Mild	High	True	Yes
Overcast	Hot	Normal	False	Yes
Rainy	Mild	High	True	No



Explanation: the DT summarizes (explains) all the observations in the table perfectly  $\Rightarrow$  100% Accuracy

Prediction: once we have a DT (or model) we can use it to make predictions on observations that are not in the original training table, consider:

Outlook = Sunny, Temperature = Mild, Humidity = Normal, Windy = False, Playtennis = ?

# Constructing DTs

- How do we choose the attributes and the order in which they appear in a DT?
  - Recursive partitioning of the original data table
  - Heuristic - each generated partition has to be “less random” (entropy reduction) than previously generated partitions

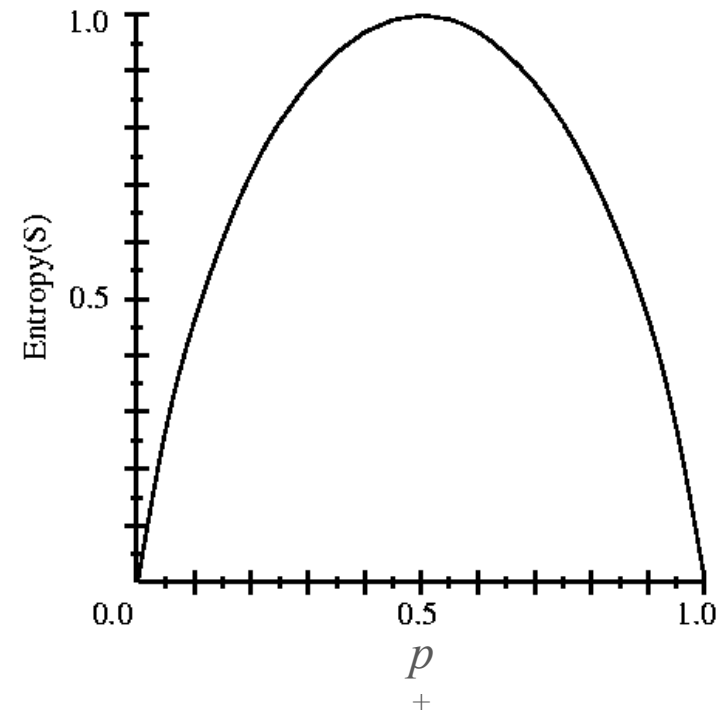
# Entropy

- $S$  is a sample of training examples
- $p^+$  is the proportion of positive examples in  $S$
- $p^-$  is the proportion of negative examples in  $S$
- Entropy measures the impurity (randomness) of  $S$

$S$  {

Outlook	Temperature	Humidity	Windy	PlayTennis
Sunny	Hot	High	False	No
Sunny	Hot	High	True	No
Overcast	Hot	High	False	Yes
Rainy	Mild	High	False	Yes
Rainy	Cool	Normal	False	Yes
Rainy	Cool	Normal	True	No
Overcast	Cool	Normal	True	Yes
Sunny	Mild	High	False	No
Sunny	Cool	Normal	False	Yes
Rainy	Mild	Normal	False	Yes
Sunny	Mild	Normal	True	Yes
Overcast	Mild	High	True	Yes
Overcast	Hot	Normal	False	Yes
Rainy	Mild	High	True	No

$$Entropy(S) = Entropy([9+, 5-]) = .94$$

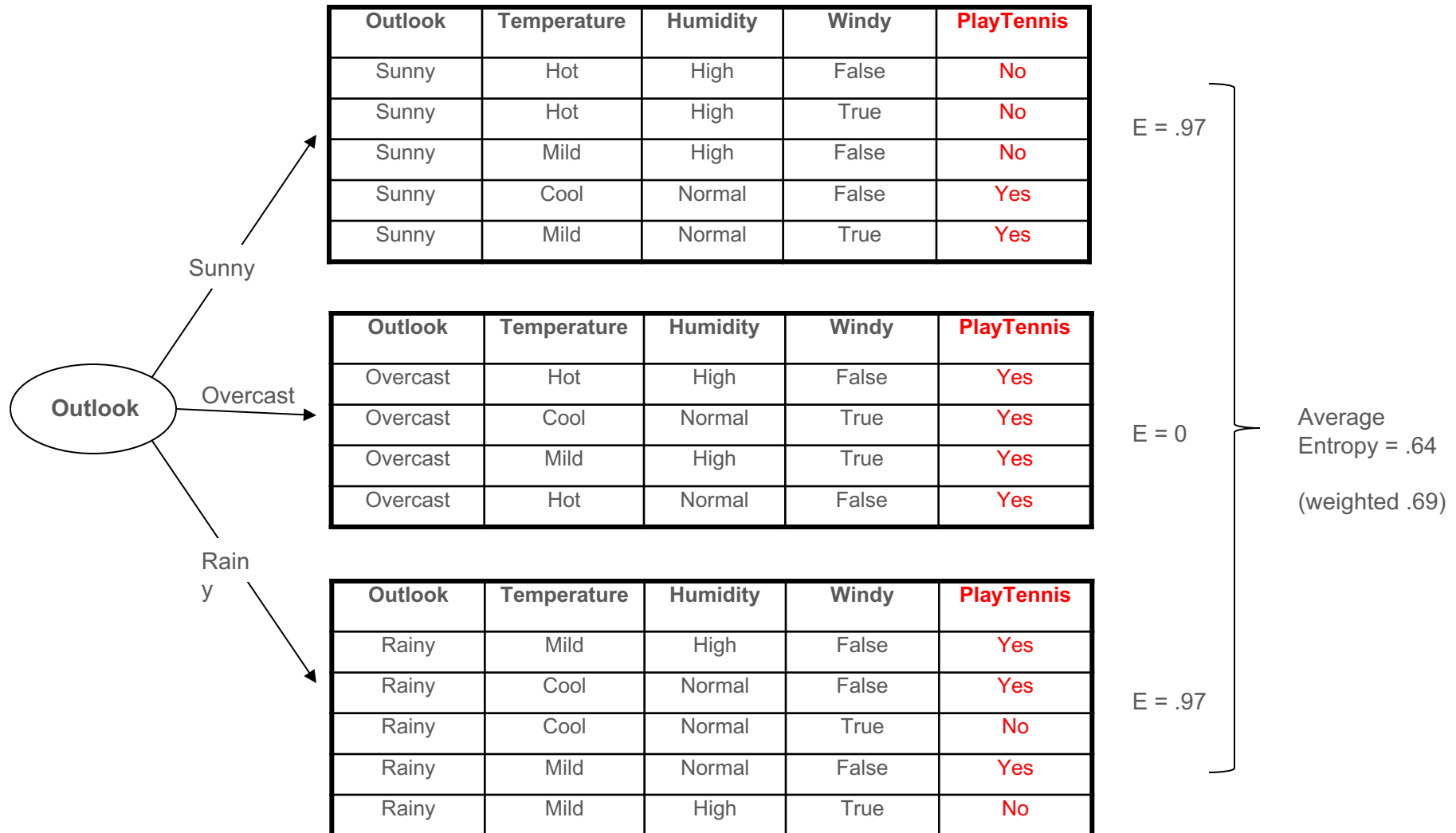


$$\square Entropy(S) \equiv - p^+ \log_2 p^+ - p^- \log_2 p^-$$



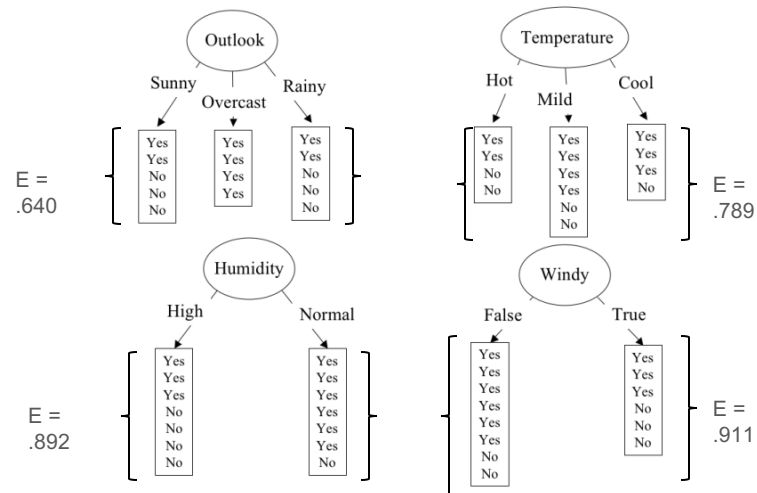
$$\text{AvgEntropy}(S, A) = \sum_{v \in \text{Values}(A)} \frac{|S_v|}{|S|} E(S_v) \quad (\text{weighted average})$$

# Partitioning the Data Set



# Partitioning in Action

Outlook	Temperature	Humidity	Windy	PlayTennis
Sunny	Hot	High	False	No
Sunny	Hot	High	True	No
Overcast	Hot	High	False	Yes
Rainy	Mild	High	False	Yes
Rainy	Cool	Normal	False	Yes
Rainy	Cool	Normal	True	No
Overcast	Cool	Normal	True	Yes
Sunny	Mild	High	False	No
Sunny	Cool	Normal	False	Yes
Rainy	Mild	Normal	False	Yes
Sunny	Mild	Normal	True	Yes
Overcast	Mild	High	True	Yes
Overcast	Hot	Normal	False	Yes
Rainy	Mild	High	True	No



# Recursive Partitioning

Partition(*Examples*, *TargetAttribute*, *Attributes*)

*Examples* are the training examples. *TargetAttribute* is a binary (+/-) categorical dependent variable and *Attributes* is the list of independent variables which are available for testing at this point. This function returns a decision tree.

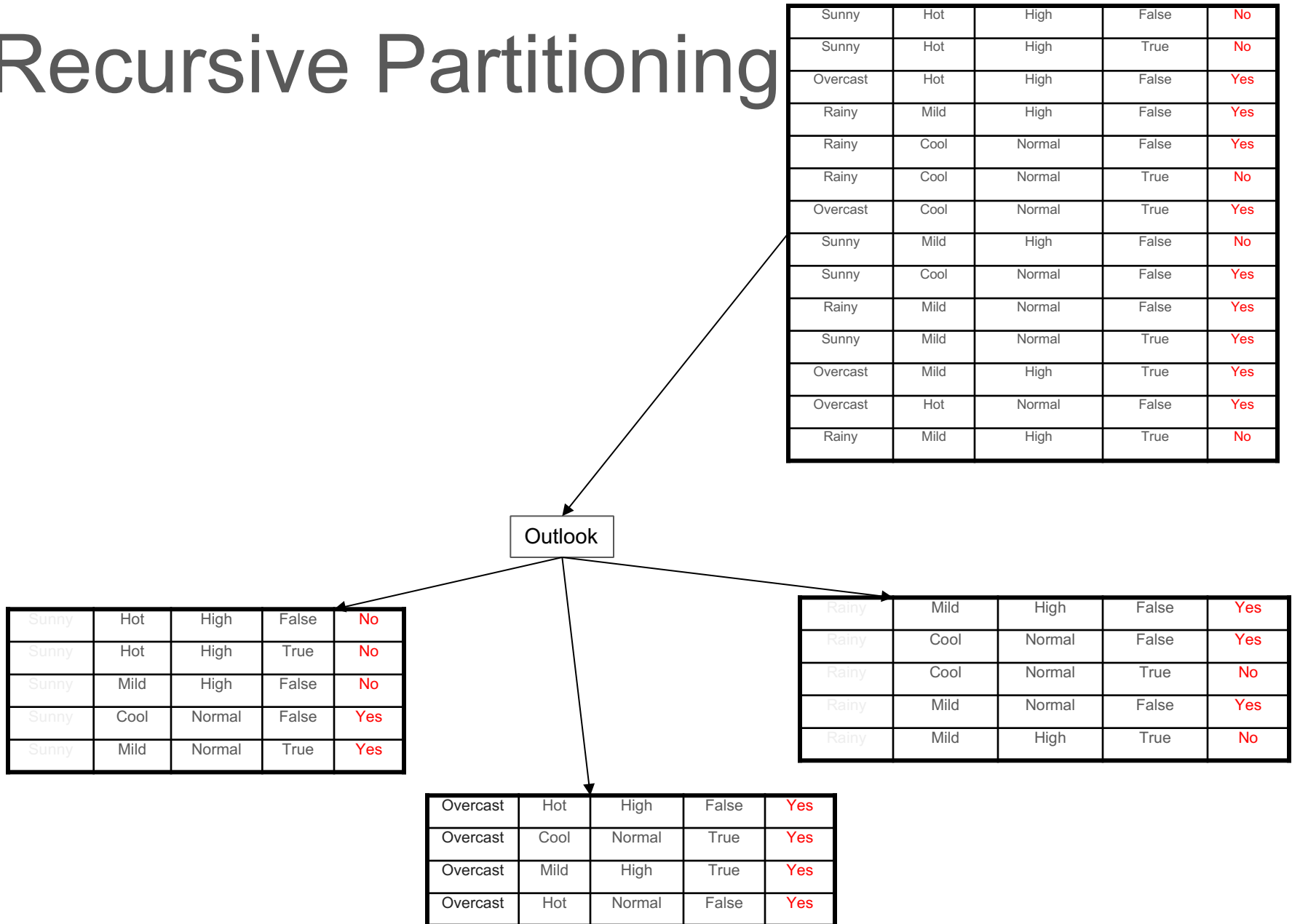
- Create a *Root* node for the tree.
- If all *Examples* are positive then return *Root* as a leaf node with label = +.
- Else if all *Examples* are negative then return *Root* as a leaf node with label = -.
- Else if *Attributes* is empty then return *Root* as a leaf node with label = most common value of *TargetAttribute* in *Examples*.
- Otherwise
  - $A :=$  the attribute from *Attributes* that reduces entropy the most on the *Examples*.
  - $Root := A$
  - For each  $v \in \text{values}(A)$ 
    - Add a new branch below the *Root* node with value  $A = v$
    - Let  $Examples_v$  be the subset of *Examples* where  $A = v$
    - If  $Examples_v$  is empty then add new leaf node to branch with label = most common value of *TargetAttribute* in *Examples*.
    - Else add new subtree to branch  
Partition( $Examples_v$ , *TargetAttribute*,  $Attributes - \{A\}$ )
- Return *Root*

# Recursive Partitioning

Our data set:

Outlook	Temperature	Humidity	Windy	PlayTennis
Sunny	Hot	High	False	No
Sunny	Hot	High	True	No
Overcast	Hot	High	False	Yes
Rainy	Mild	High	False	Yes
Rainy	Cool	Normal	False	Yes
Rainy	Cool	Normal	True	No
Overcast	Cool	Normal	True	Yes
Sunny	Mild	High	False	No
Sunny	Cool	Normal	False	Yes
Rainy	Mild	Normal	False	Yes
Sunny	Mild	Normal	True	Yes
Overcast	Mild	High	True	Yes
Overcast	Hot	Normal	False	Yes
Rainy	Mild	High	True	No

# Recursive Partitioning



# Recursive Partitioning

Outlook

Sunny	Hot	High	False	No
Sunny	Hot	High	True	No
Sunny	Mild	High	False	No
Sunny	Cool	Normal	False	Yes
Sunny	Mild	Normal	True	Yes

Rainy	Mild	High	False	Yes
Rainy	Cool	Normal	False	Yes
Rainy	Cool	Normal	True	No
Rainy	Mild	Normal	False	Yes
Rainy	Mild	High	True	No

Overcast	Hot	High	False	Yes
Overcast	Cool	Normal	True	Yes
Overcast	Mild	High	True	Yes
Overcast	Hot	Normal	False	Yes

# Recursive Partitioning

Outlook

Sunny	Hot	High	False	No
Sunny	Hot	High	True	No
Sunny	Mild	High	False	No
Sunny	Cool	Normal	False	Yes
Sunny	Mild	Normal	True	Yes

Rainy	Mild	High	False	Yes
Rainy	Cool	Normal	False	Yes
Rainy	Cool	Normal	True	No
Rainy	Mild	Normal	False	Yes
Rainy	Mild	High	True	No

Overcast	Hot	High	False	Yes
Overcast	Cool	Normal	True	Yes
Overcast	Mild	High	True	Yes
Overcast	Hot	Normal	False	Yes

Humidity

Sunny	Cool	Normal	False	Yes
Sunny	Mild	Normal	True	Yes

Sunny	Hot	High	False	No
Sunny	Hot	High	True	No
Sunny	Mild	High	False	No

# Recursive Partitioning

Outlook

Sunny	Hot	High	False	No
Sunny	Hot	High	True	No
Sunny	Mild	High	False	No
Sunny	Cool	Normal	False	Yes
Sunny	Mild	Normal	True	Yes

Humidity

Overcast	Hot	High	False	Yes
Overcast	Cool	Normal	True	Yes
Overcast	Mild	High	True	Yes
Overcast	Hot	Normal	False	Yes

Sunny	Cool	Normal	False	Yes
Sunny	Mild	Normal	True	Yes

Rainy	Mild	High	False	Yes
Rainy	Cool	Normal	False	Yes
Rainy	Cool	Normal	True	No
Rainy	Mild	Normal	False	Yes
Rainy	Mild	High	True	No

Windy

Rainy	Mild	High	False	Yes
Rainy	Cool	Normal	False	Yes
Rainy	Mild	Normal	False	Yes

Rainy	Cool	Normal	True	No
Rainy	Mild	High	True	No

Sunny	Hot	High	False	No
Sunny	Hot	High	True	No
Sunny	Mild	High	False	No



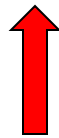
# Continuous-Valued Attributes

Consider:

<i>Temperature:</i>	40	48	60	72	80	90
<i>PlayTennis:</i>	No	No	Yes	Yes	Yes	No



$$(48+60)/2 \\ = 54$$



$$(80+90)/2 \\ = 85$$

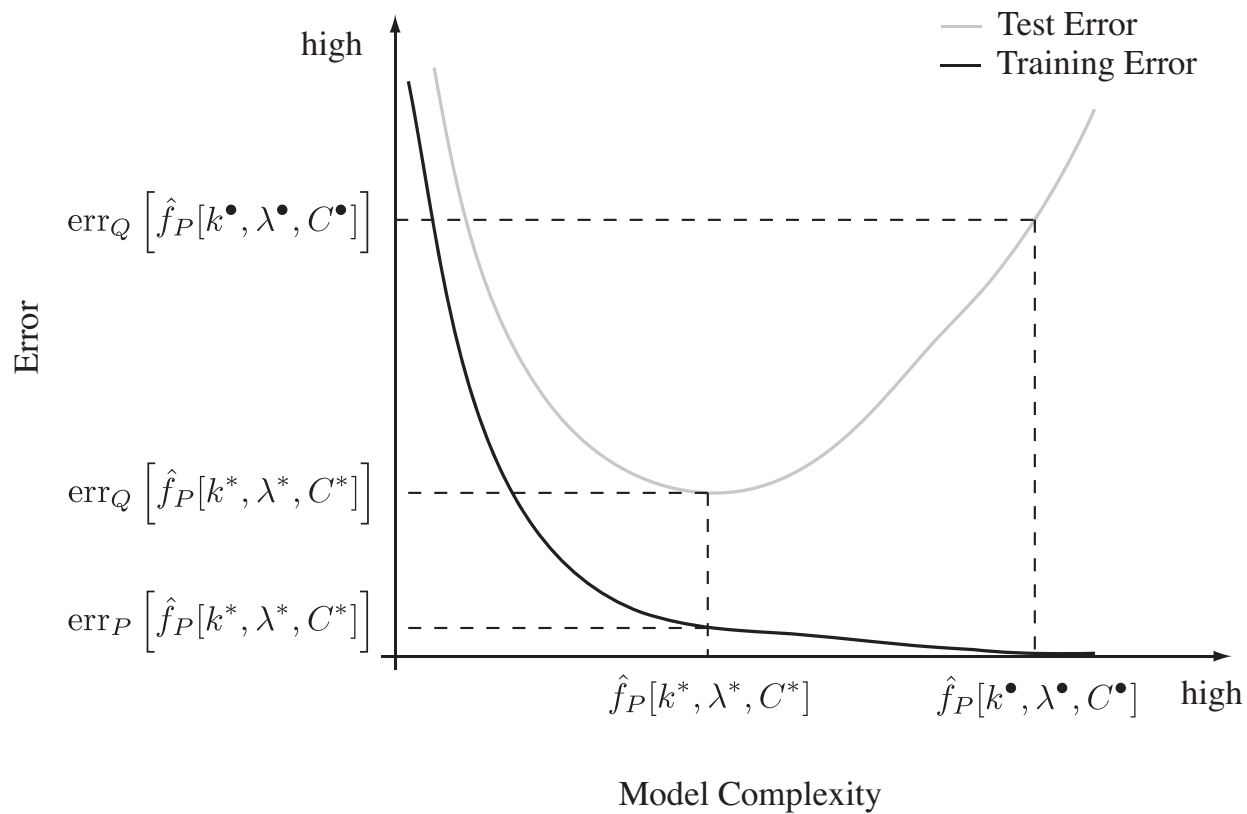
Highest Gain: Temperature > 54

- Sort instances according to the attribute values
- Find “Splits” where the classes change
- Select the split that gives you the highest gain

# Decision Trees & Patterns in Data

- True pattern in domain
  - present in large amounts of data
  - generalizes to unseen instances
- Spurious pattern in training set
  - “noise in the data”
  - present in small amounts of data
  - does not generalize

# Overfitting – Also True for Trees!



Tree Depth!

# Tree Learning Process

- Beginning
  - lots of data
  - discovers true patterns in data
- Later
  - small amount of data
  - likely to learn spurious patterns

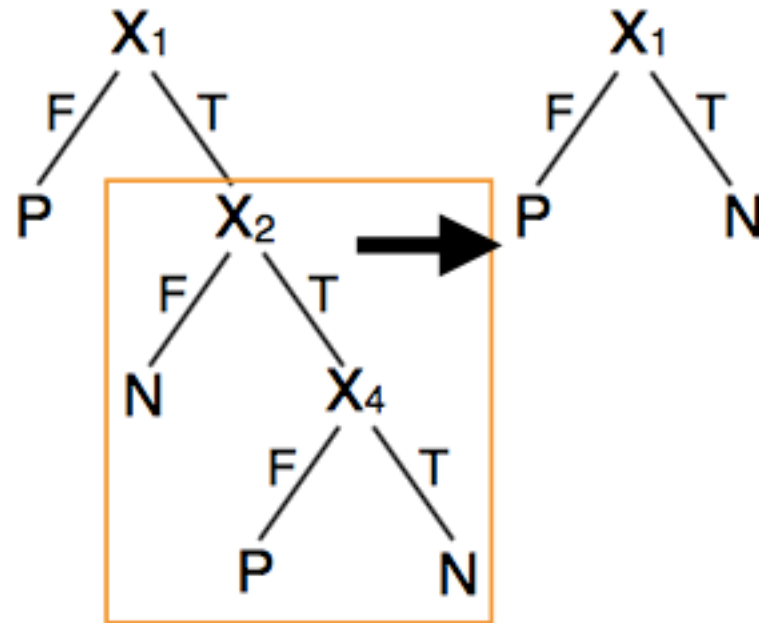
Control the Tree Complexity - Pruning

# Pruning

- One of two ways:
  1. Prevent the tree from overfitting – limit the tree depth.
  2. Build the whole tree and then remove subtrees and replaces with suitable leaves.

# Pruning Example

$X_1$	$X_2$	$X_3$	$X_4$	C
F	F	F	F	P
F	F	T	T	P
F	T	F	T	P
T	T	T	F	P
T	F	F	F	N
T	T	T	T	N
T	T	T	F	N



# Subtree Pruning with Deviation

- At each split ask:
  - Is the pattern found in the data after splitting statistically significant?
- Prune if deviation is small – that is, prune if no significant information gain.

# Given Split

- Given instances  $\mathbf{D}$
- Propose split on  $X_i$
- Notation
  - $N_c$  = number of instances with class  $c$
  - $\mathbf{D}_x$  = data set with value  $x$  for attribute  $X_i$
  - $N_x$  = number of instances in  $\mathbf{D}_x$
  - $N_{xc}$  = number of instances in  $\mathbf{D}_x$  with class  $c$



# Absence of Pattern

- Null hypothesis:  $X_i$  is irrelevant
- In  $\mathbf{D}$ , proportion with class  $c$ :  $N_c/N$
- If null hypothesis is true, we expect on average the number of instances in  $\mathbf{D}_x$  with class  $c$  to be:

$$\hat{N}_{xc} = \frac{|D_x|}{N} N_c$$

# Deviation

- We don't expect to see *exactly* that many, even if null hypothesis is true
- We expect some deviation due to random chance
- Measure the deviation from total absence of pattern:

$$\text{Dev} = \sum_x \sum_c \frac{(N_{xc} - \hat{N}_{xc})^2}{\hat{N}_{xc}}$$

→ Delete split if Dev is small