

Lecture Slides for

INTRODUCTION TO

# Machine Learning

2nd Edition

CHAPTER 9:

## Decision Trees

ETHEM ALPAYDIN

© The MIT Press, 2010

Edited and expanded for CS 4641 by Chris Simpkins

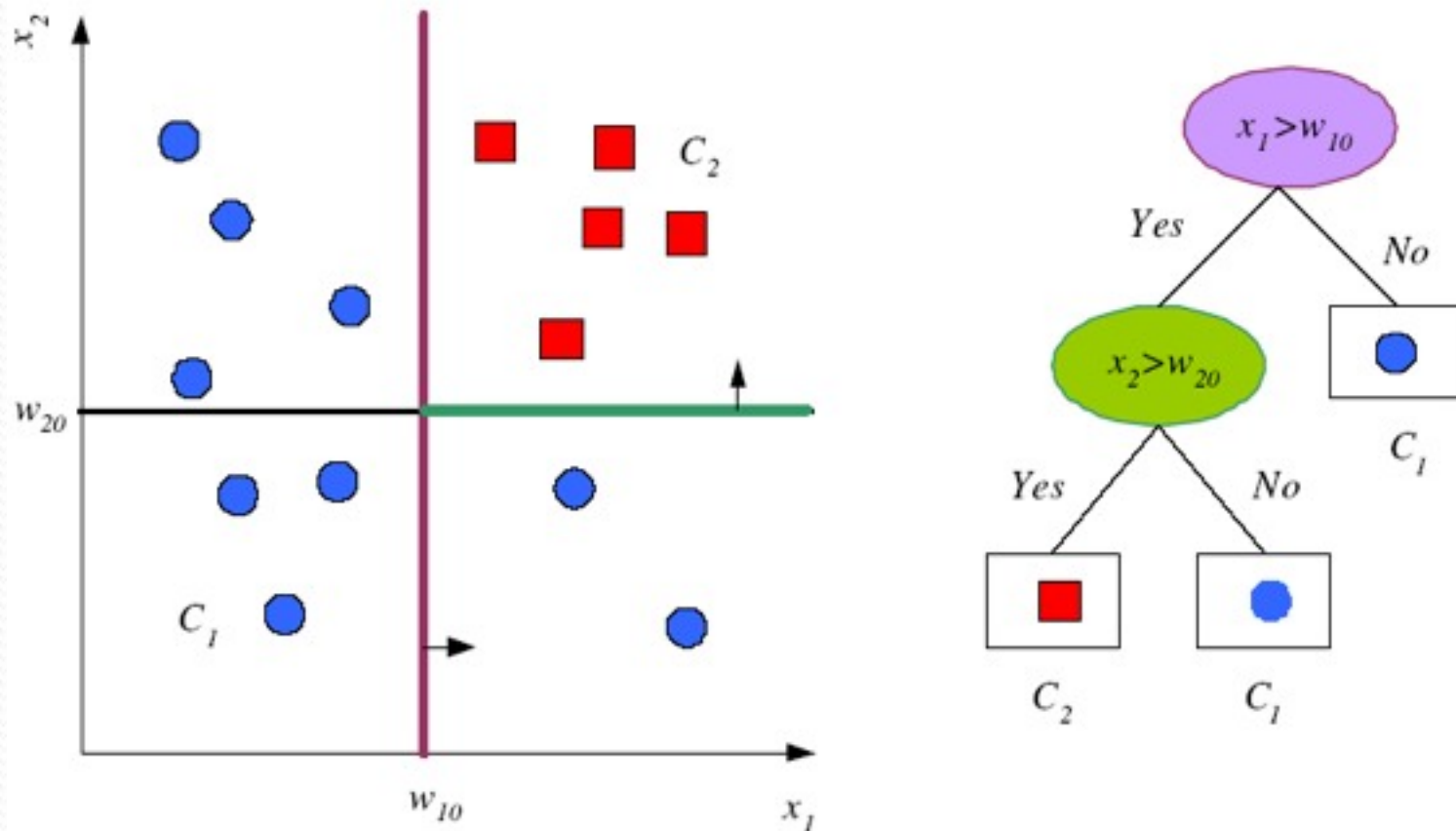
[alpaydin@boun.edu.tr](mailto:alpaydin@boun.edu.tr)

<http://www.cmpe.boun.edu.tr/~ethem/i2ml2e>

# Overview

- Univariate decision trees
- Building classification trees
- Dealing with overfitting
- Extracting rules from decision trees

# Tree Uses Nodes, and Leaves



# Divide and Conquer

- Internal decision nodes
  - Univariate: Uses a single attribute,  $x_j$ 
    - Numeric  $x_j$ : Binary split :  $x_j > w_m$
    - Discrete  $x_j$ :  $n$ -way split for  $n$  possible values
  - Multivariate: Uses all attributes,  $\mathbf{x}$
- Leaves
  - Classification: Class labels, or proportions
  - Regression: Numeric;  $r$  average, or local fit
- Learning is greedy; find the best split recursively (Breiman et al, 1984; Quinlan, 1986, 1993)

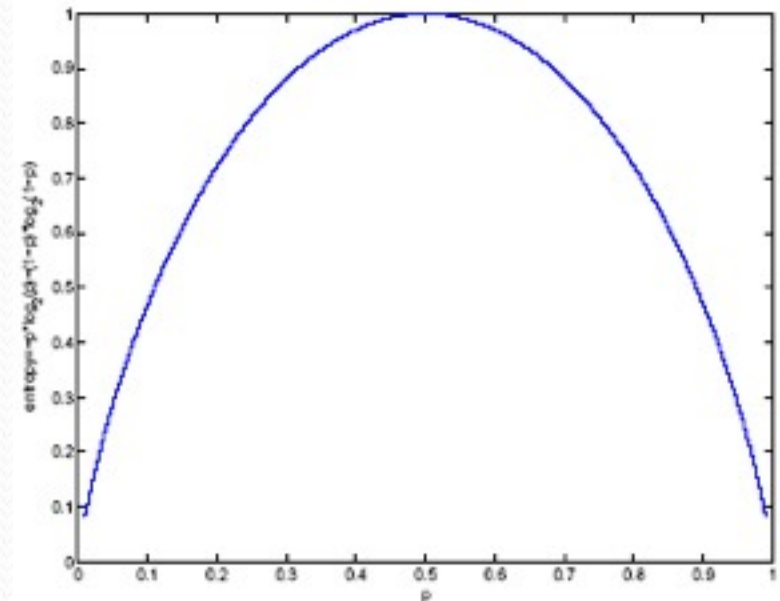
# Classification Trees (ID3, CART, C4.5)

- For node  $m$ ,  $N_m$  instances reach  $m$ ,  $N_m^i$  belong to  $C_i$

$$\hat{P}(C_i | \mathbf{x}, m) \equiv p_m^i = \frac{N_m^i}{N_m}$$

- Node  $m$  is pure if  $p_m^i$  is 0 or 1
- Measure of impurity is entropy

$$I_m = - \sum_{i=1}^K p_m^i \log_2 p_m^i$$



# Best Split

- If node  $m$  is pure, generate a leaf and stop, otherwise split and continue recursively
- Impurity after split:  $N_{mj}$  of  $N_m$  take branch  $j$ .  $N_{mj}^i$  belong to  $C_i$

$$\hat{P}(C_i | \mathbf{x}, m, j) \equiv p_{mj}^i = \frac{N_{mj}^i}{N_{mj}}$$

$$J'_m = - \sum_{j=1}^n \frac{N_{mj}}{N_m} \sum_{i=1}^K p_{mj}^i \log_2 p_{mj}^i$$

- Find the variable and split that min impurity (among all variables -- and split positions for numeric variables)

GenerateTree( $\mathcal{X}$ )

If NodeEntropy( $\mathcal{X}$ ) <  $\theta_I$  /\* eq. 9.3

    Create leaf labelled by majority class in  $\mathcal{X}$

    Return

$i \leftarrow$  SplitAttribute( $\mathcal{X}$ )

    For each branch of  $\mathbf{x}_i$

        Find  $\mathcal{X}_i$  falling in branch

        GenerateTree( $\mathcal{X}_i$ )

SplitAttribute( $\mathcal{X}$ )

    MinEnt  $\leftarrow$  MAX

    For all attributes  $i = 1, \dots, d$

        If  $\mathbf{x}_i$  is discrete with  $n$  values

            Split  $\mathcal{X}$  into  $\mathcal{X}_1, \dots, \mathcal{X}_n$  by  $\mathbf{x}_i$

$e \leftarrow$  SplitEntropy( $\mathcal{X}_1, \dots, \mathcal{X}_n$ ) /\* eq. 9.8 \*/

            If  $e <$  MinEnt MinEnt  $\leftarrow$   $e$ ; bestf  $\leftarrow$   $i$

        Else /\*  $\mathbf{x}_i$  is numeric \*/

            For all possible splits

                Split  $\mathcal{X}$  into  $\mathcal{X}_1, \mathcal{X}_2$  on  $\mathbf{x}_i$

$e \leftarrow$  SplitEntropy( $\mathcal{X}_1, \mathcal{X}_2$ )

                If  $e <$  MinEnt MinEnt  $\leftarrow$   $e$ ; bestf  $\leftarrow$   $i$

    Return bestf

# Summary of Main Loop

- Pick  $A$ , the “best” decision attribute for examples at current node using impurity measure
- For each value of  $A$ , create new descendant leaf nodes of current node
- Sort training examples to leaf nodes according to their values for attribute  $A$
- If training examples perfectly classified, Then STOP, Else iterate over new leaf nodes (turning some of them into internal nodes)



# Example: Play tennis today?<sup>1</sup>

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

<sup>1</sup>Example from Mitchell 1997

# Choosing the first split attribute

- Outlook:

$$\begin{aligned} I'_{root} &= - \sum_{j \in \{Sunny, Overcast, Rain\}} \frac{N_j}{N} \sum_{i \in \{Yes, No\}} p_j^i \log_2 p_j^i \\ &= - \left[ \frac{5}{14} \left( \frac{2}{5} \log_2 \frac{2}{5} + \frac{3}{5} \log_2 \frac{3}{5} \right) + \frac{4}{14} \left( \frac{4}{4} \log_2 \frac{4}{4} + \frac{0}{4} \log_2 \frac{0}{4} \right) + \frac{5}{14} \left( \frac{3}{5} \log_2 \frac{3}{5} + \frac{2}{5} \log_2 \frac{2}{5} \right) \right] \\ &= 0.693 \end{aligned}$$

- Temperature:

$$\begin{aligned} I'_{root} &= - \sum_{j \in \{Hot, Mild, Cool\}} \frac{N_j}{N} \sum_{i \in \{Yes, No\}} p_j^i \log_2 p_j^i \\ &= - \left[ \frac{4}{14} \left( \frac{2}{4} \log_2 \frac{2}{4} + \frac{2}{4} \log_2 \frac{2}{4} \right) + \frac{6}{14} \left( \frac{4}{6} \log_2 \frac{4}{6} + \frac{2}{6} \log_2 \frac{2}{6} \right) + \frac{4}{14} \left( \frac{3}{4} \log_2 \frac{3}{4} + \frac{1}{4} \log_2 \frac{1}{4} \right) \right] \\ &= 0.915 \end{aligned}$$

- Humidity: 
$$\begin{aligned} I'_{root} &= - \sum_{j \in \{High, Normal\}} \frac{N_j}{N} \sum_{i \in \{Yes, No\}} p_j^i \log_2 p_j^i \\ &= - \left[ \frac{7}{14} \left( \frac{3}{7} \log_2 \frac{3}{7} + \frac{4}{7} \log_2 \frac{4}{7} \right) + \frac{7}{14} \left( \frac{6}{7} \log_2 \frac{6}{7} + \frac{1}{7} \log_2 \frac{1}{7} \right) \right] \\ &= 0.789 \end{aligned}$$

- Wind: 
$$\begin{aligned} I'_{root} &= - \sum_{j \in \{Strong, Weak\}} \frac{N_j}{N} \sum_{i \in \{Yes, No\}} p_j^i \log_2 p_j^i \\ &= - \left[ \frac{6}{14} \left( \frac{3}{6} \log_2 \frac{3}{6} + \frac{3}{6} \log_2 \frac{3}{6} \right) + \frac{8}{14} \left( \frac{6}{8} \log_2 \frac{6}{8} + \frac{2}{8} \log_2 \frac{2}{8} \right) \right] \\ &= 0.892 \end{aligned}$$

{D1, D2, ..., D14}

[9+,5-]

Outlook

Sunny

Overcast

Rain

{D1,D2,D8,D9,D11}

{D3,D7,D12,D13}

{D4,D5,D6,D10,D14}

[2+,3-]

[4+,0-]

[3+,2-]

?

Yes

?

Which attribute should be tested here?

# Choosing an attribute for Sunny node

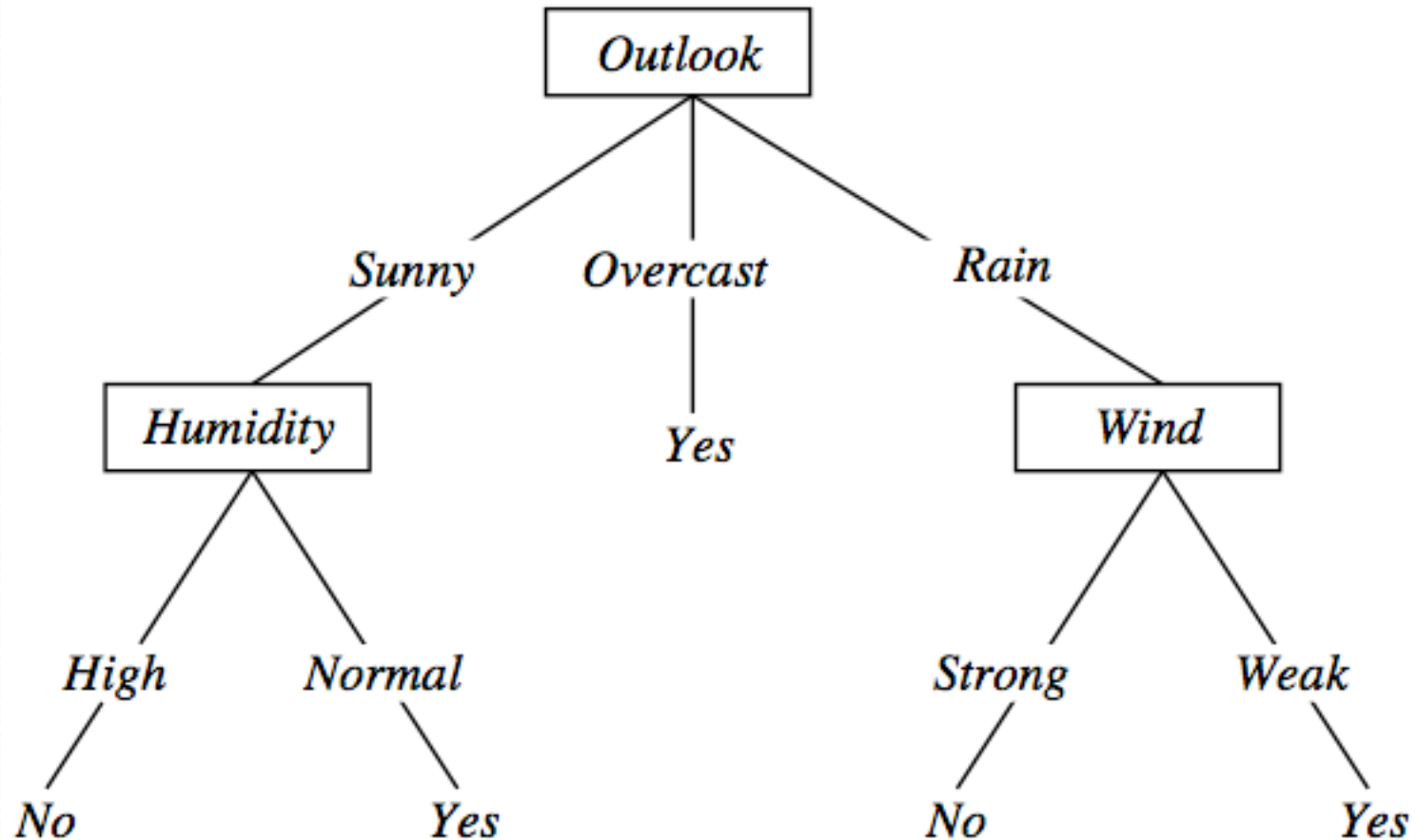
- $X_{\text{Sunny}} = \{D_1, D_2, D_8, D_9, D_{11}\}$

Day	Outlook	Temperature	Humidity	Wind	PlayTennis?
D1	Sunny	Hot	High	Weak	No
D2	Sunny	Hot	High	Strong	No
D8	Sunny	Mild	High	Weak	No
D9	Sunny	Cool	Normal	Weak	Yes
D11	Sunny	Mild	Normal	Strong	Yes

- Humidity: 
$$I'_{\text{sunny,humidity}} = - \sum_{j \in \{High, Normal\}} \frac{N_j}{N} \sum_{i \in \{Yes, No\}} p_j^i \log_2 p_j^i$$
$$= - \left[ \frac{3}{5} \left( \frac{0}{3} \log_2 \frac{0}{3} + \frac{3}{3} \log_2 \frac{3}{3} \right) + \frac{2}{5} \left( \frac{2}{2} \log_2 \frac{2}{2} + \frac{0}{2} \log_2 \frac{0}{2} \right) \right]$$
$$= 0$$

- Do same calculation for Temperature and Wind
- Humidity has lowest entropy, so Humidity is split attribute at Sunny node

# Final Induced Tree



# Observations

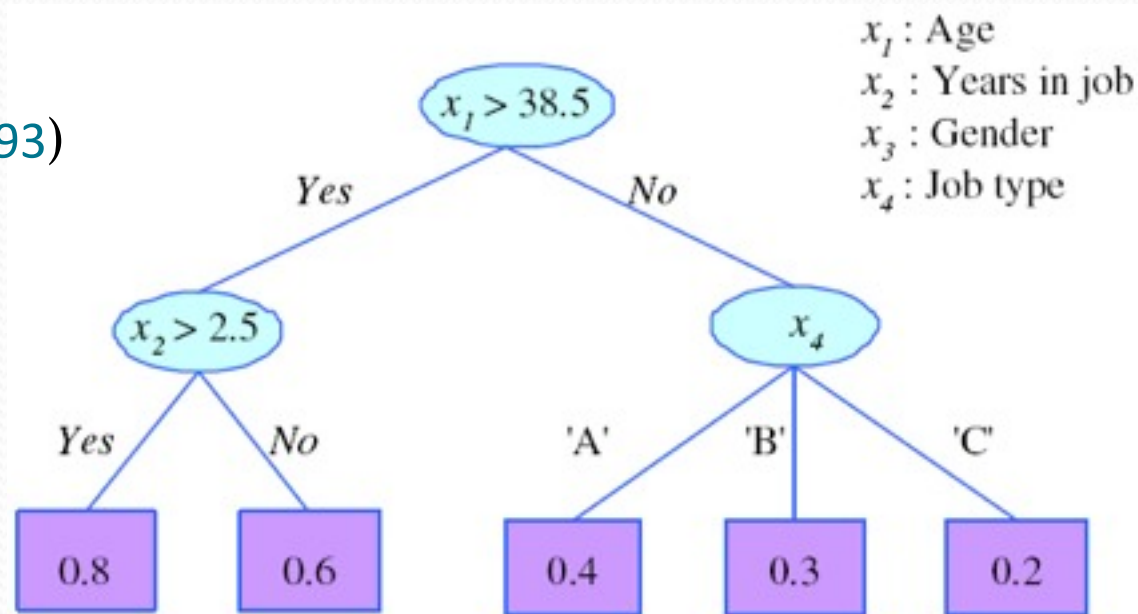
- Temperature was not needed
  - Decision tree can be used for feature extraction (a.k.a. feature selection)
- Tree is short
  - ID<sub>3</sub> family of algorithms have preference bias for short trees
  - Occam's Razor: shorter trees are better

# Pruning Trees

- Remove subtrees for better generalization (decrease variance)
  - Prepruning: Early stopping
  - Postpruning: Grow the whole tree then prune subtrees which overfit on the pruning set
- Prepruning is faster, postpruning is more accurate (requires a separate pruning set)

# Rule Extraction from Trees

C4.5Rules  
(Quinlan, 1993)



- R1: IF (age > 38.5) AND (years-in-job > 2.5) THEN  $y = 0.8$
- R2: IF (age > 38.5) AND (years-in-job  $\leq$  2.5) THEN  $y = 0.6$
- R3: IF (age  $\leq$  38.5) AND (job-type = 'A') THEN  $y = 0.4$
- R4: IF (age  $\leq$  38.5) AND (job-type = 'B') THEN  $y = 0.3$
- R5: IF (age  $\leq$  38.5) AND (job-type = 'C') THEN  $y = 0.2$



# Conclusion

- Decision trees good when:
  - Instances described by attribute-value pairs
  - Target function is discrete
  - Interpretability of learned hypothesis (e.g., as a rule set) is desired
  - Training data may be noisy
- Decision trees are a good first algorithm to try