



Week 2

السنة الخامسة – هندسة المعلوماتية / الذكاء الصناعي

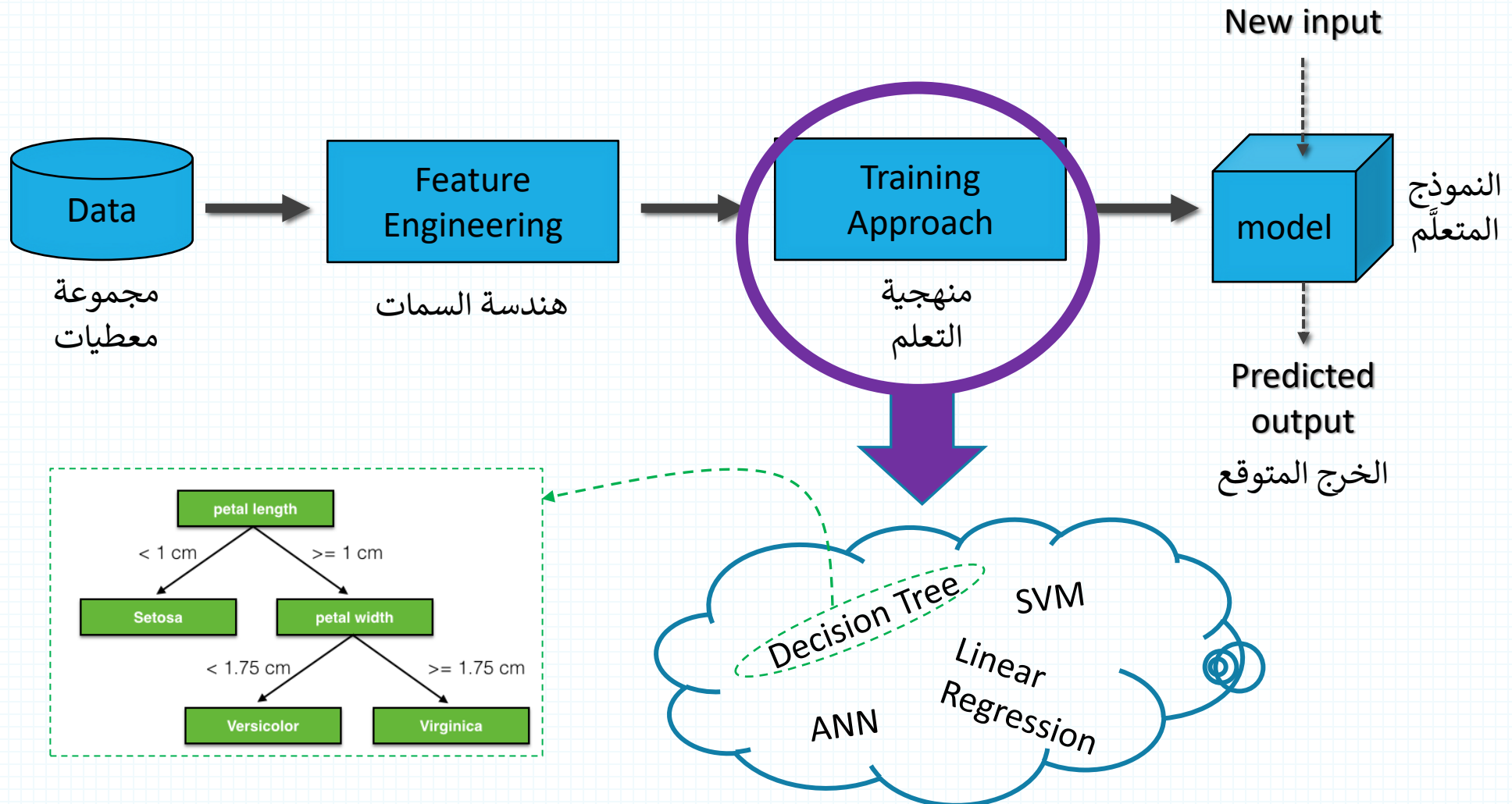
مقرر التعلم التلقائي

أشجار القرار Decision Tree

د. رياض سنبل

[Access Course Materials](#) 

Traditional ML Pipeline

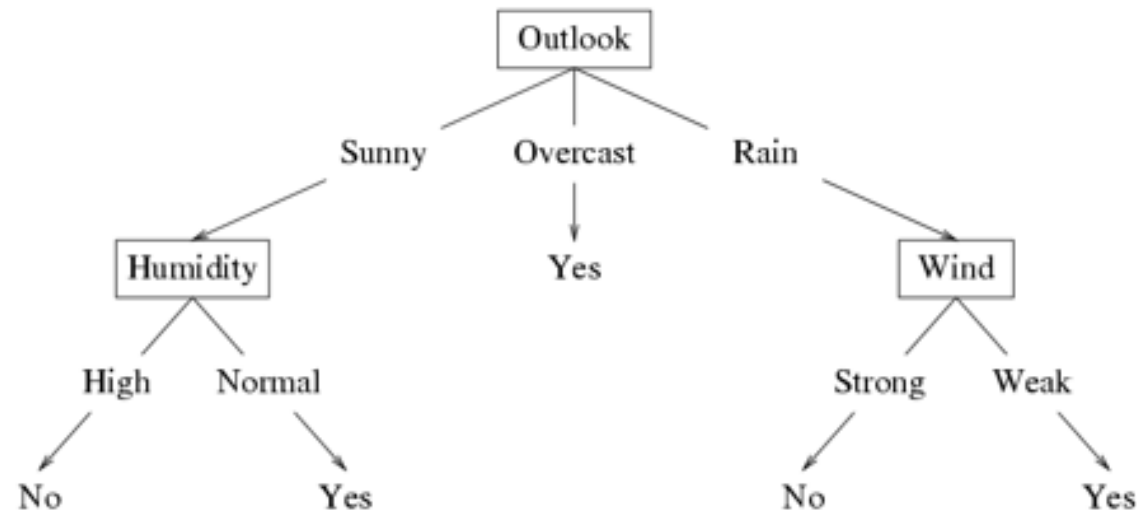


Example: Predict When This Player Will Play Tennis?

Day	Outlook	Temp.	Humidity	Wind	Play Tennis
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	Weak	Yes
D8	Sunny	Mild	High	Weak	No
D9	Sunny	Cool	Normal	Weak	Yes
D10	Rain	Mild	Normal	Strong	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

Decision Tree Hypothesis Space

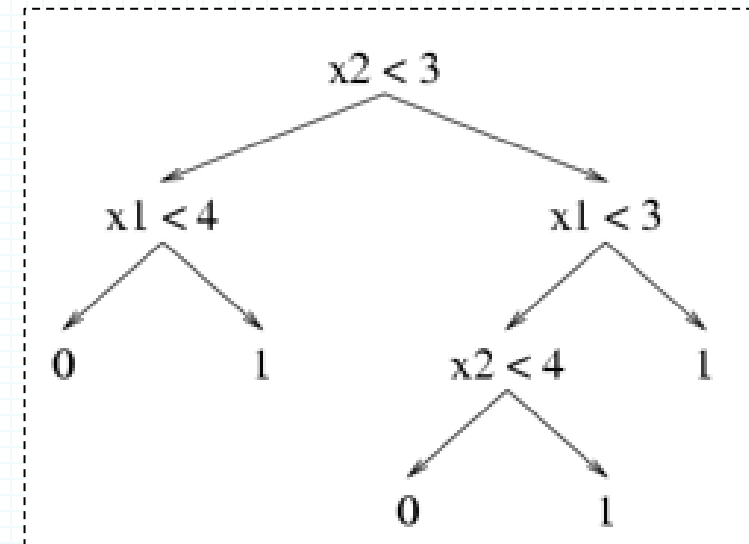
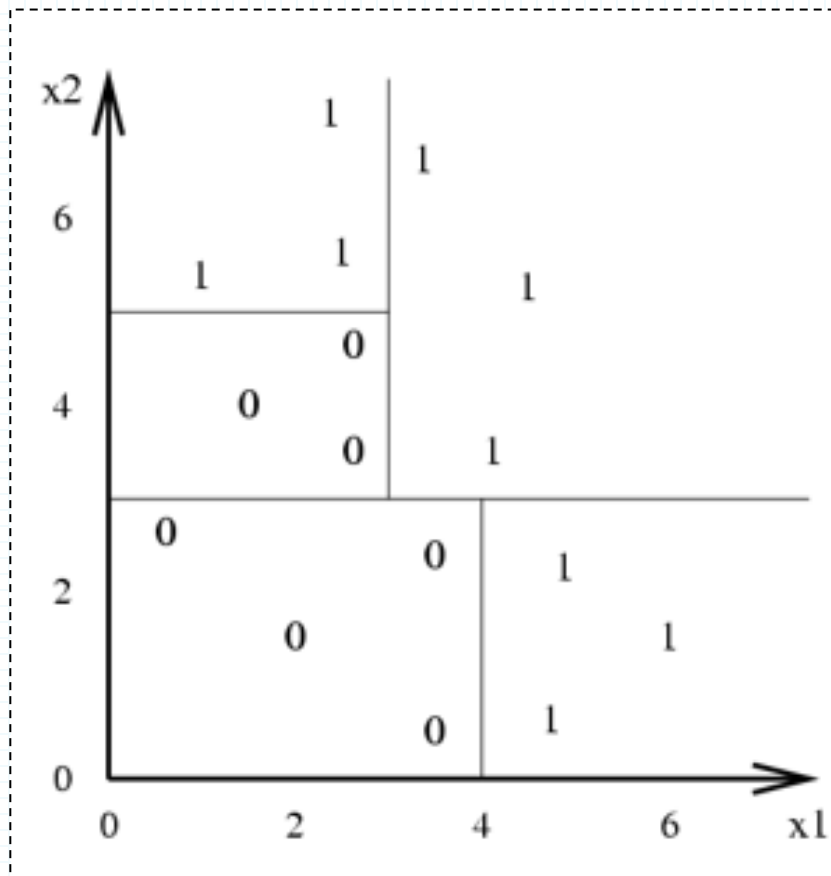
- Internal nodes test the value of particular features x_i and branch according to the results of the test.
- Leaf nodes specify the class.



Suppose the features are **Outlook** (x_1), **Temperature** (x_2), **Humidity** (x_3), and **Wind** (x_4). Then the feature vector $\mathbf{x} = (\text{Sunny}, \text{Hot}, \text{High}, \text{Strong})$ will be classified as **No**. The **Temperature** feature is irrelevant.

Decision Tree Decision Boundaries

- Decision trees divide the feature space into axis-parallel rectangles, and label each rectangle with one of the K class.



- Advantages:
 - explicit relationship among features
 - human interpretable model
- Limitations?

Learning Algorithm for Decision Trees

$$S = \{(\mathbf{x}_1, y_1), \dots, (\mathbf{x}_N, y_N)\} \quad \begin{array}{l} \mathbf{x} = (x_1, \dots, x_d) \\ x_j, y \in \{0, 1\} \end{array}$$

GROWTREE(S)

if ($y = 0$ for all $\langle \mathbf{x}, y \rangle \in S$) **return** new leaf(0)

else if ($y = 1$ for all $\langle \mathbf{x}, y \rangle \in S$) **return** new leaf(1)

else

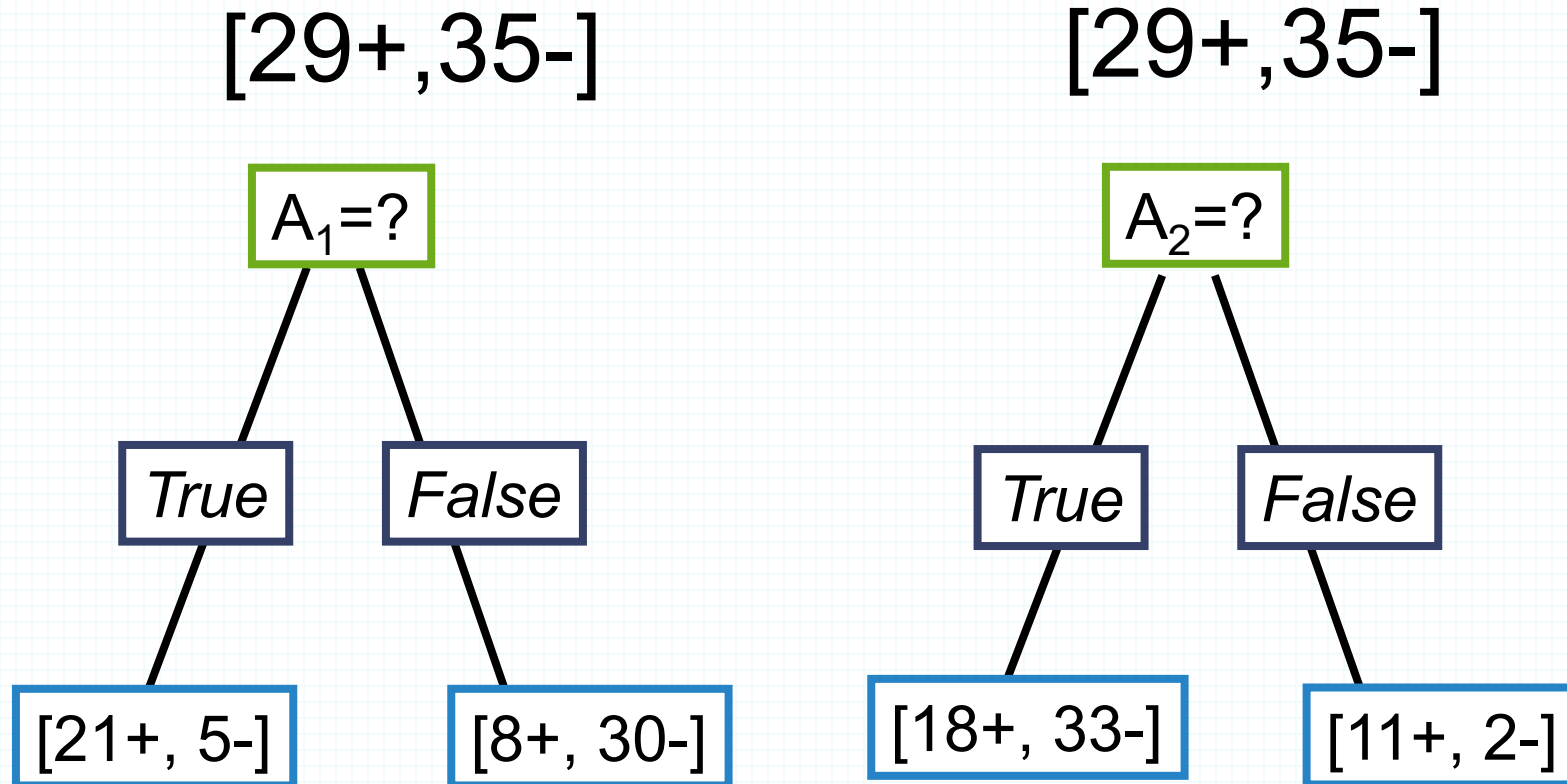
choose best attribute x_j

$S_0 =$ all $\langle \mathbf{x}, y \rangle \in S$ with $x_j = 0$;

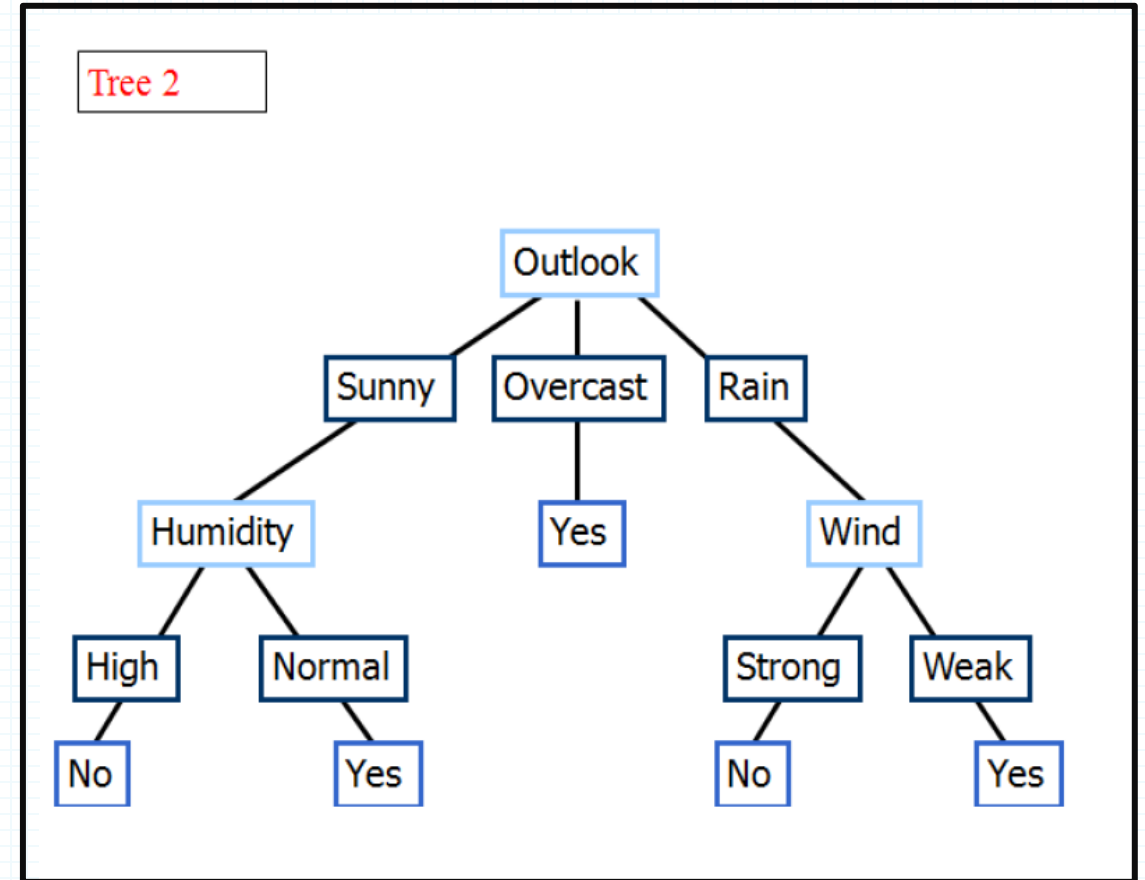
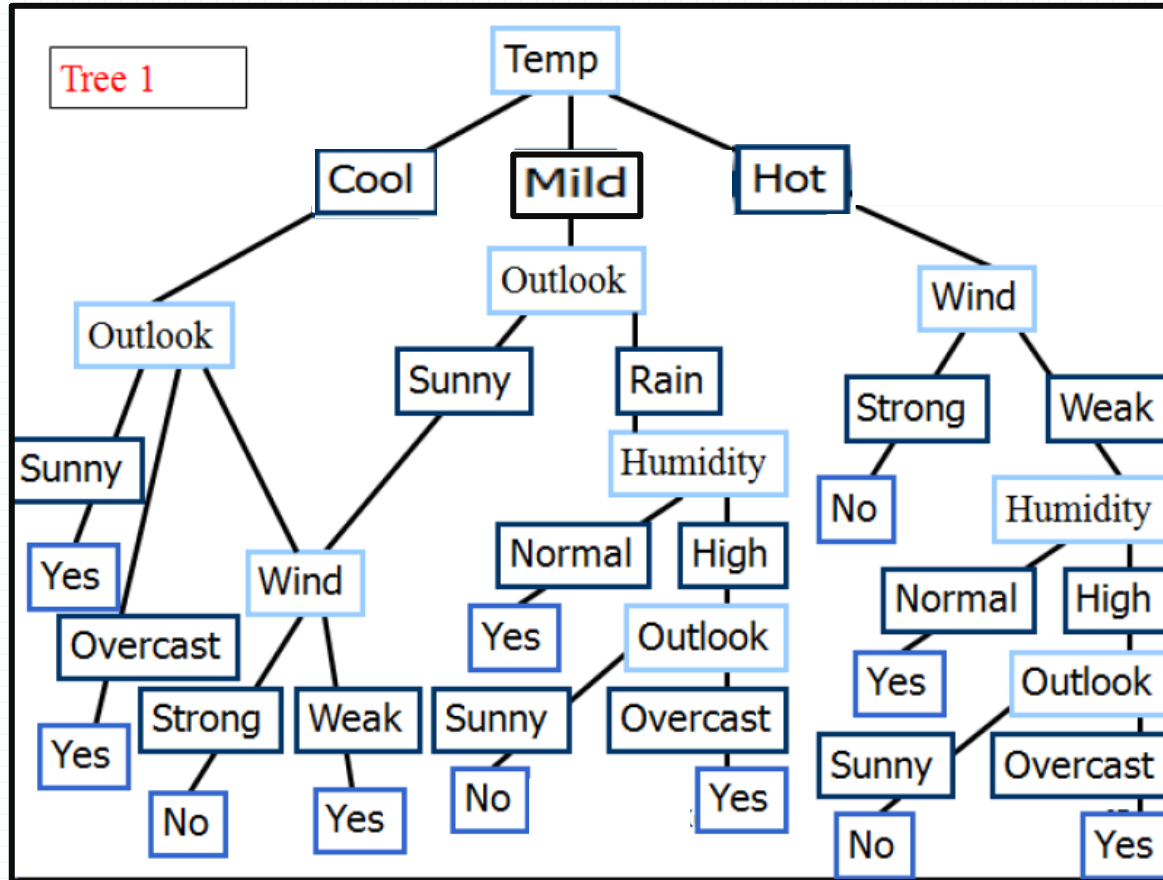
$S_1 =$ all $\langle \mathbf{x}, y \rangle \in S$ with $x_j = 1$;

return new node(x_j , GROWTREE(S_0), GROWTREE(S_1))

Which feature is "best"?



Which feature is "best"?



Using Absolute error

ChooseBestfeature(S)

choose j to minimize J_j , computed as follows:

S_0 = all $\langle \mathbf{x}, \mathbf{y} \rangle \in S$ with $x_j = 0$

S_1 = all $\langle \mathbf{x}, \mathbf{y} \rangle \in S$ with $x_j = 1$

y_0 = the most common value of \mathbf{y} in S_0

y_1 = the most common value of \mathbf{y} in S_1

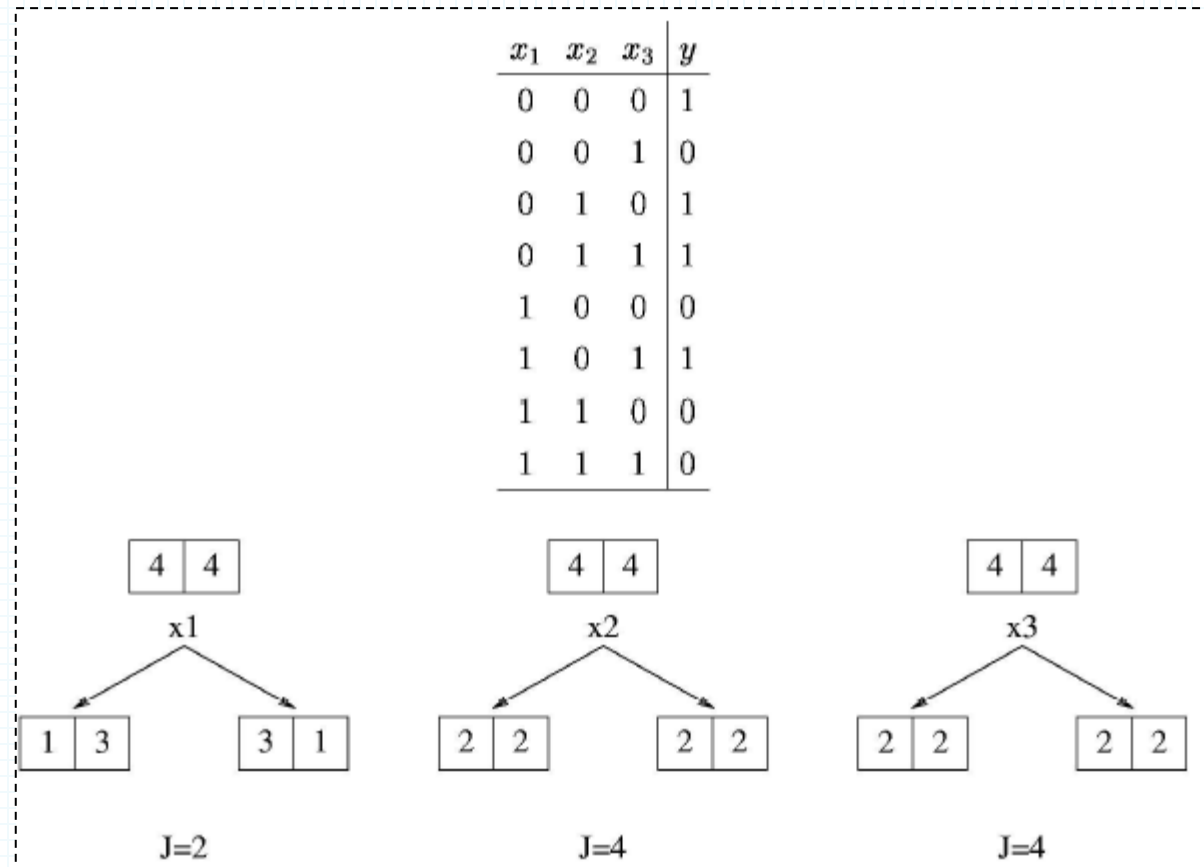
J_0 = number of examples $\langle \mathbf{x}, \mathbf{y} \rangle \in S_0$ with $\mathbf{y} \neq y_0$

J_1 = number of examples $\langle \mathbf{x}, \mathbf{y} \rangle \in S_1$ with $\mathbf{y} \neq y_1$

$J_j = J_0 + J_1$ // total classification errors if we split on feature j

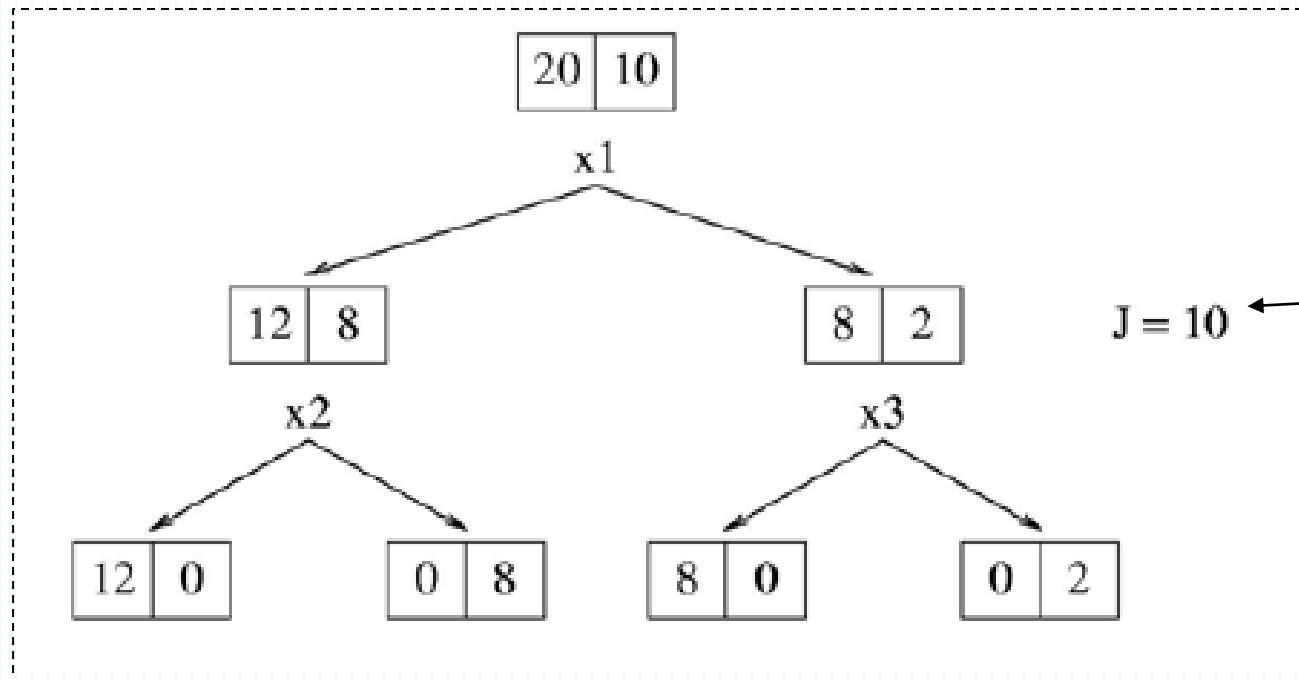
return j

Using Absolute error - Example



Using Absolute error - Example

BUT...

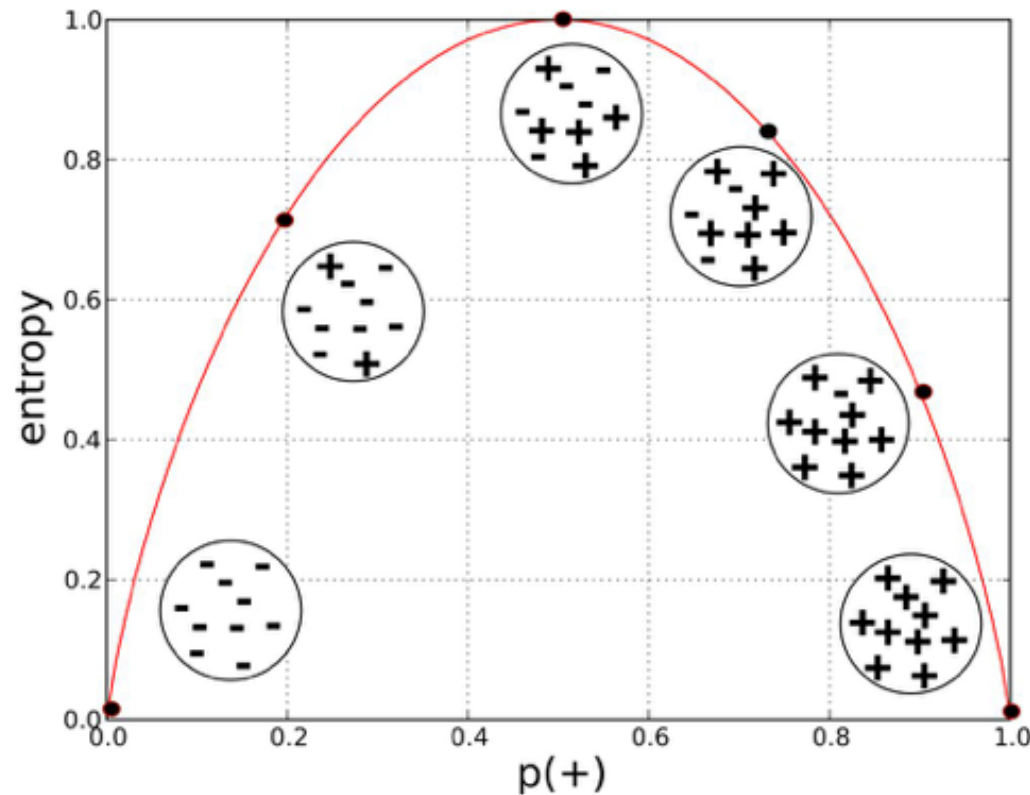


$J = 10$

Classification error is **not as sensitive to changes in data distribution**, so it might not provide enough signal for deciding where to split next.

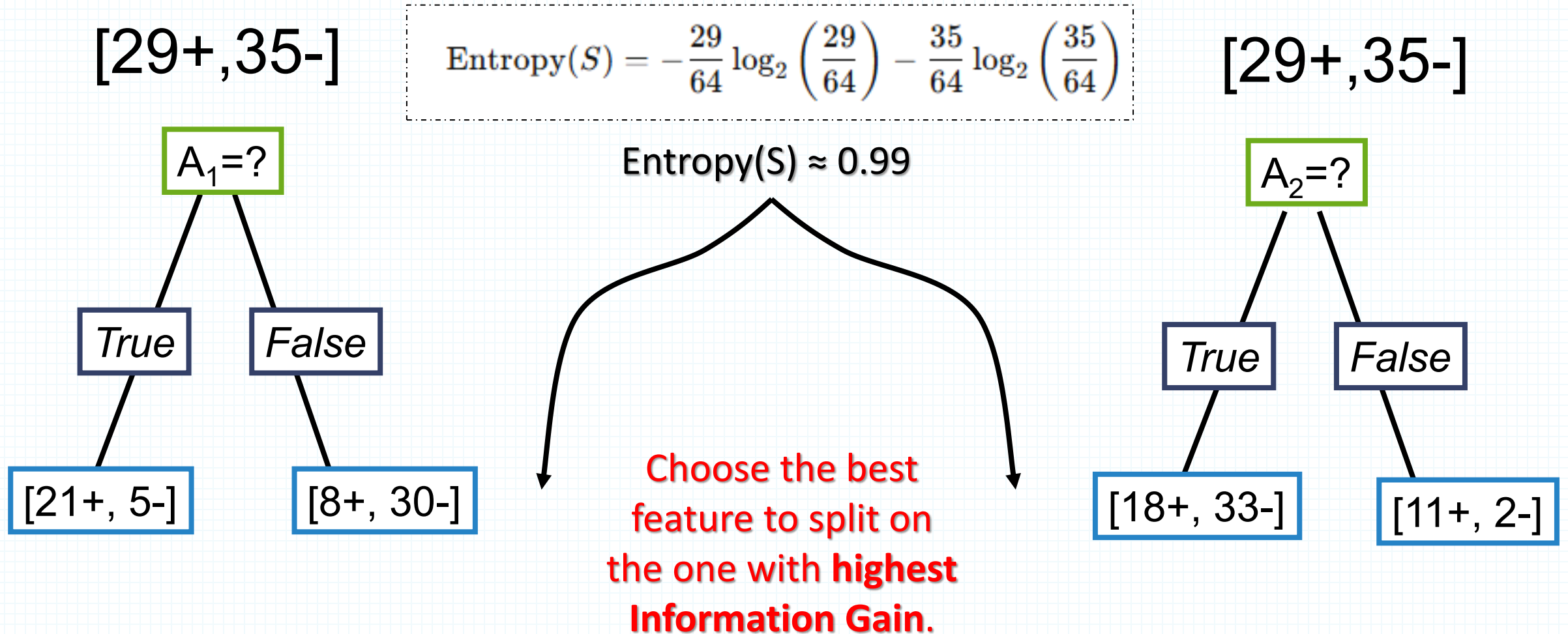
Better Criteria?
Entropy!

Entropy



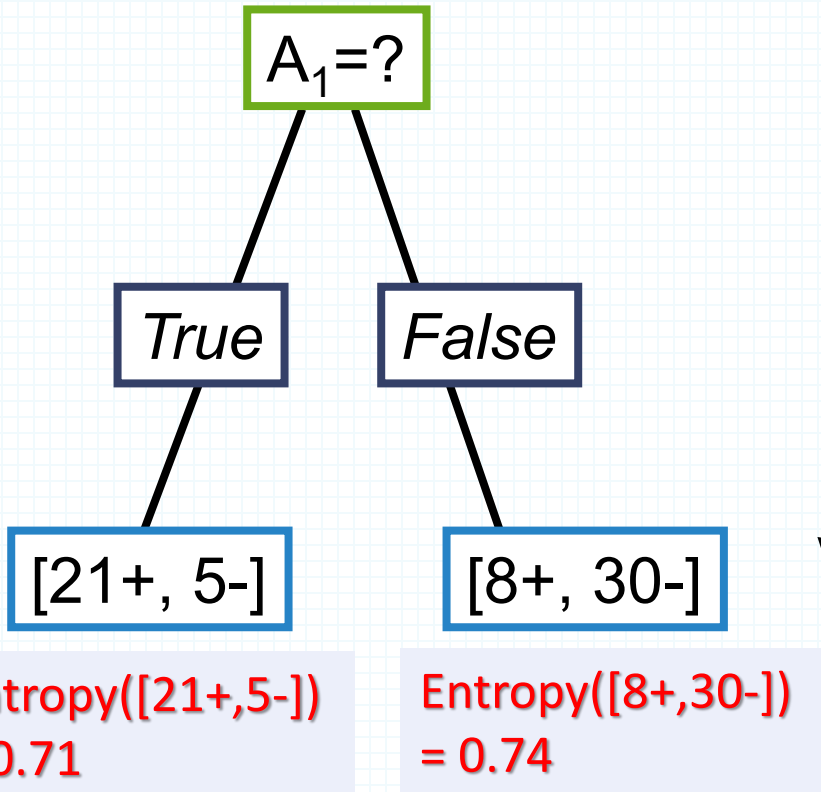
- S is a sample of training examples
- p_+ is the proportion of positive examples
- p_- is the proportion of negative examples
- Entropy measures the impurity of S
$$\text{Entropy}(S) = -p_+ \log_2 p_+ - p_- \log_2 p_-$$

Information Gain = Reduction in Entropy



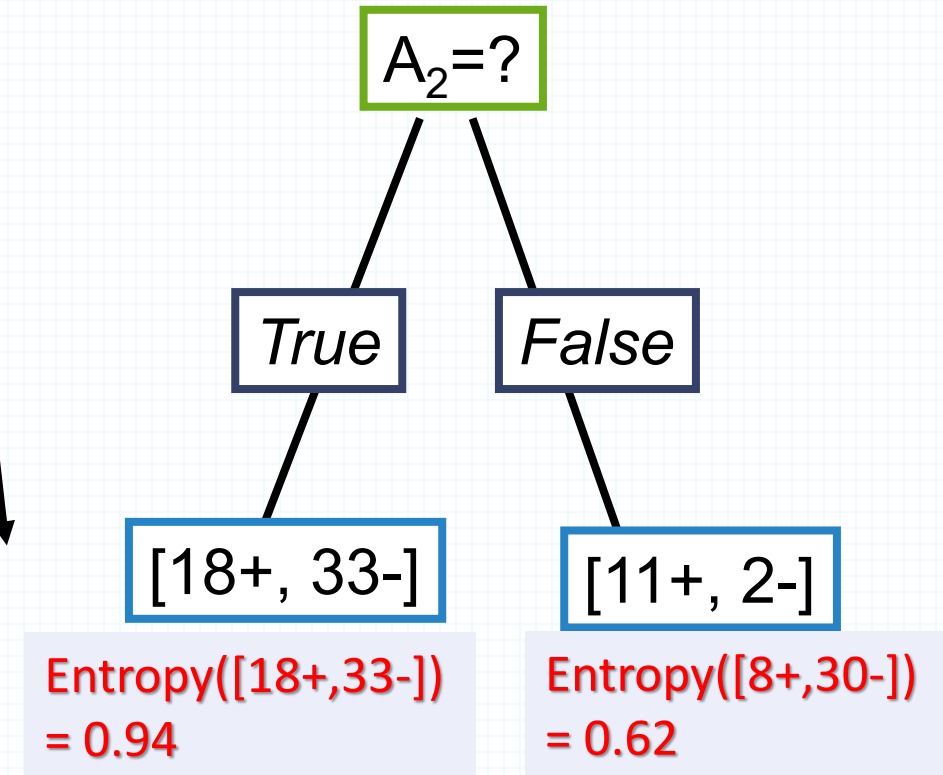
Information Gain = Reduction in Entropy

[29+, 35-]



Entropy(S) \approx 0.99

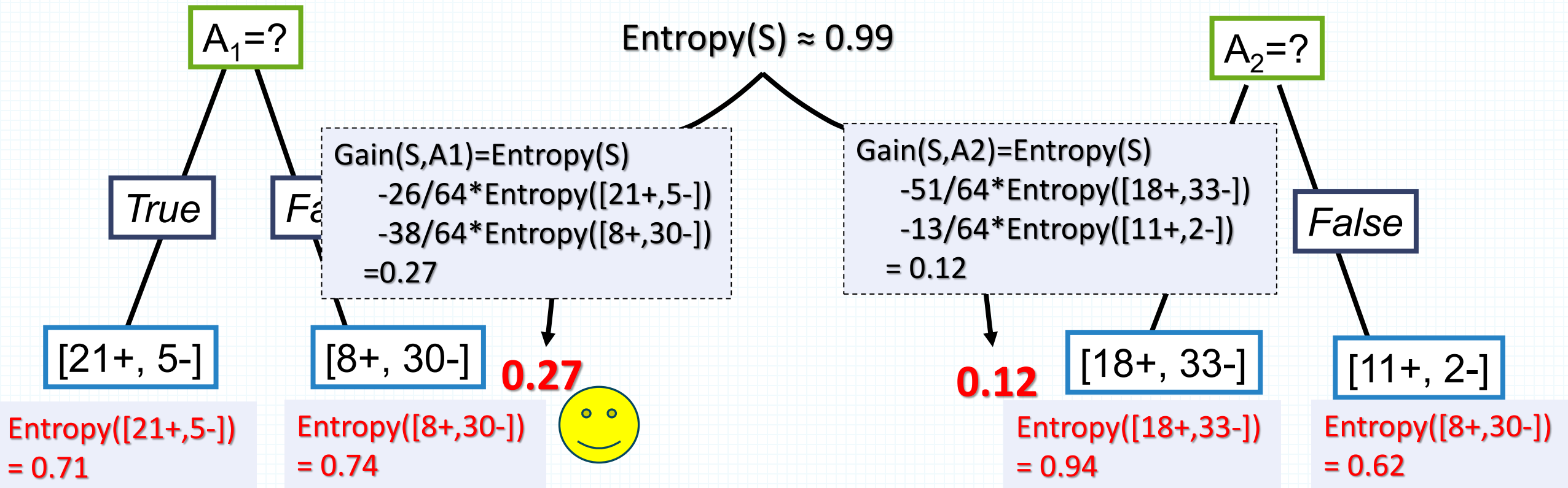
[29+, 35-]



Information Gain = Reduction in Entropy

[29+, 35-]

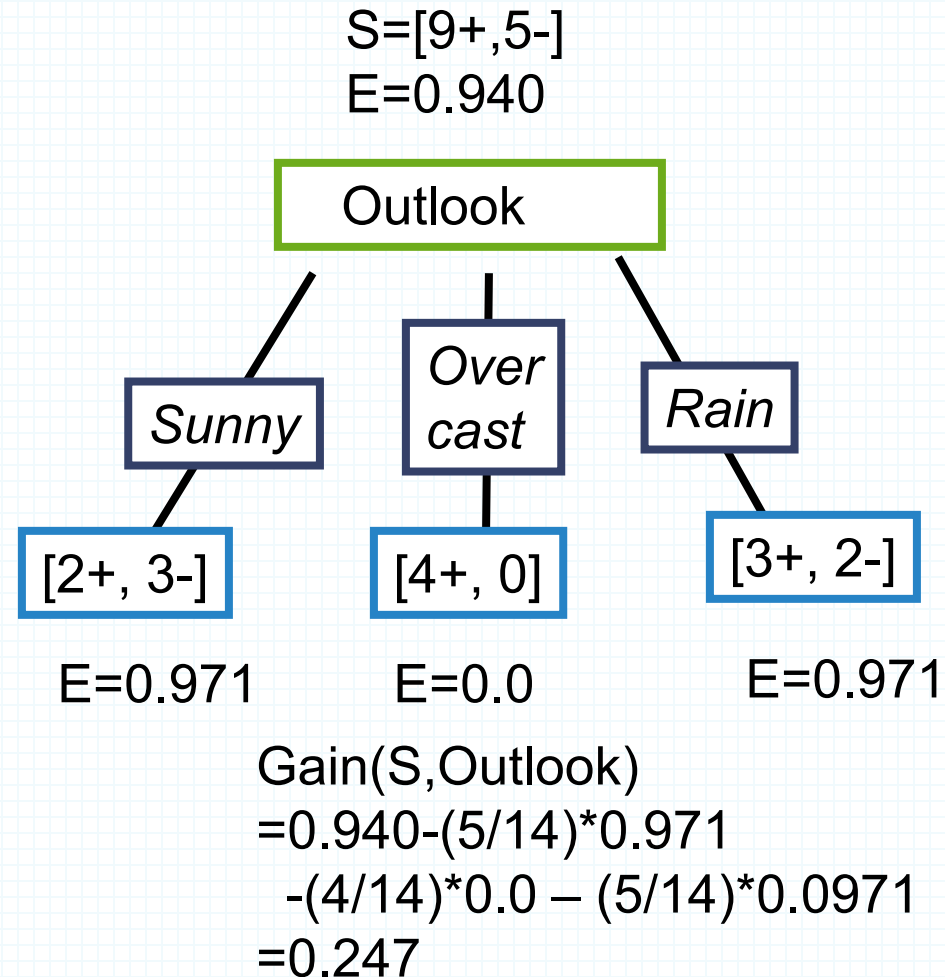
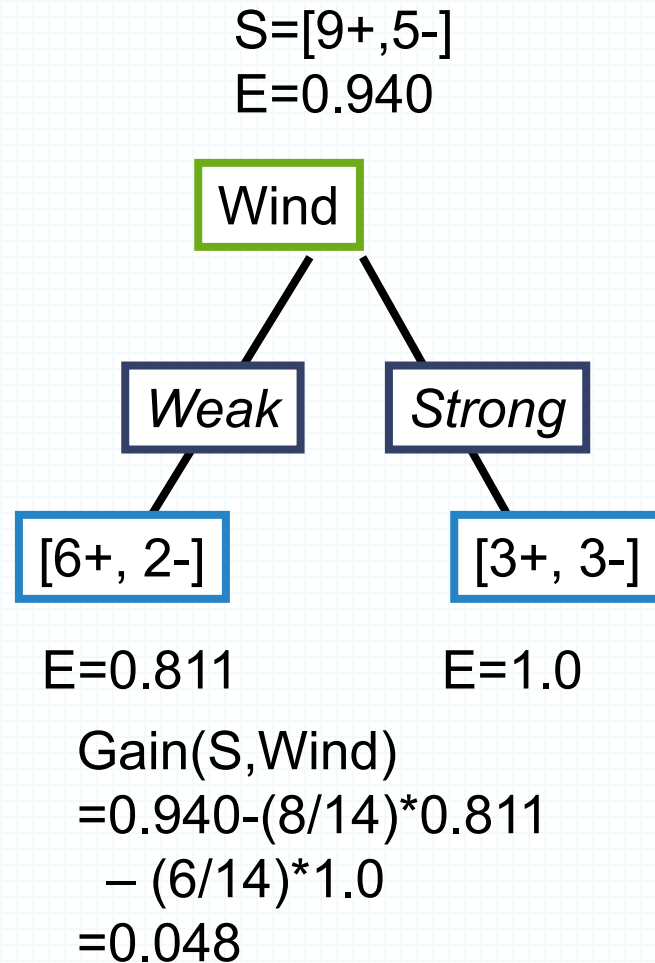
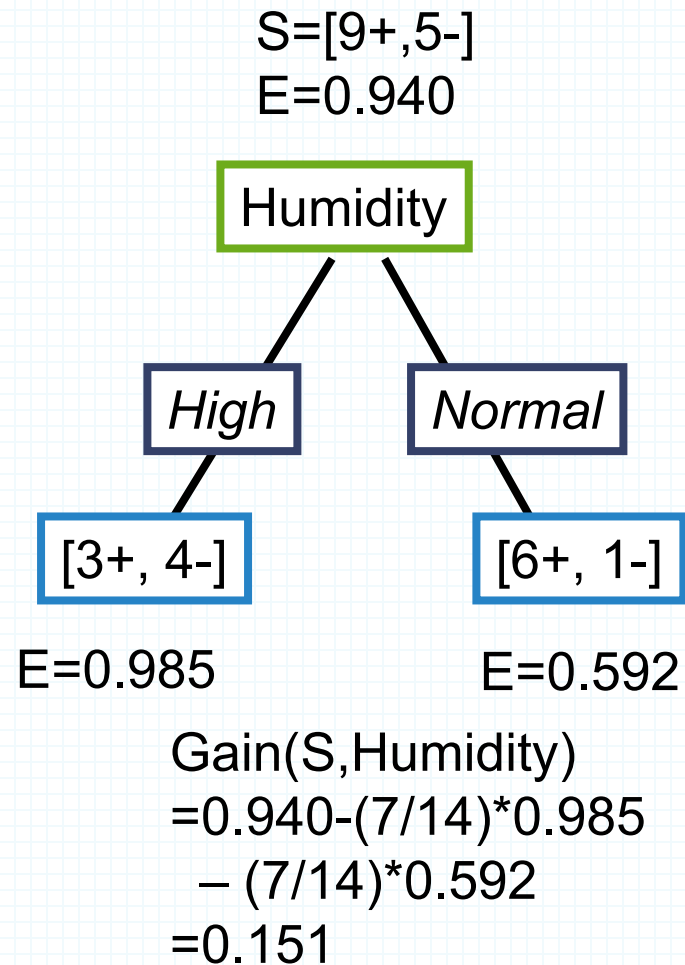
[29+, 35-]



Example

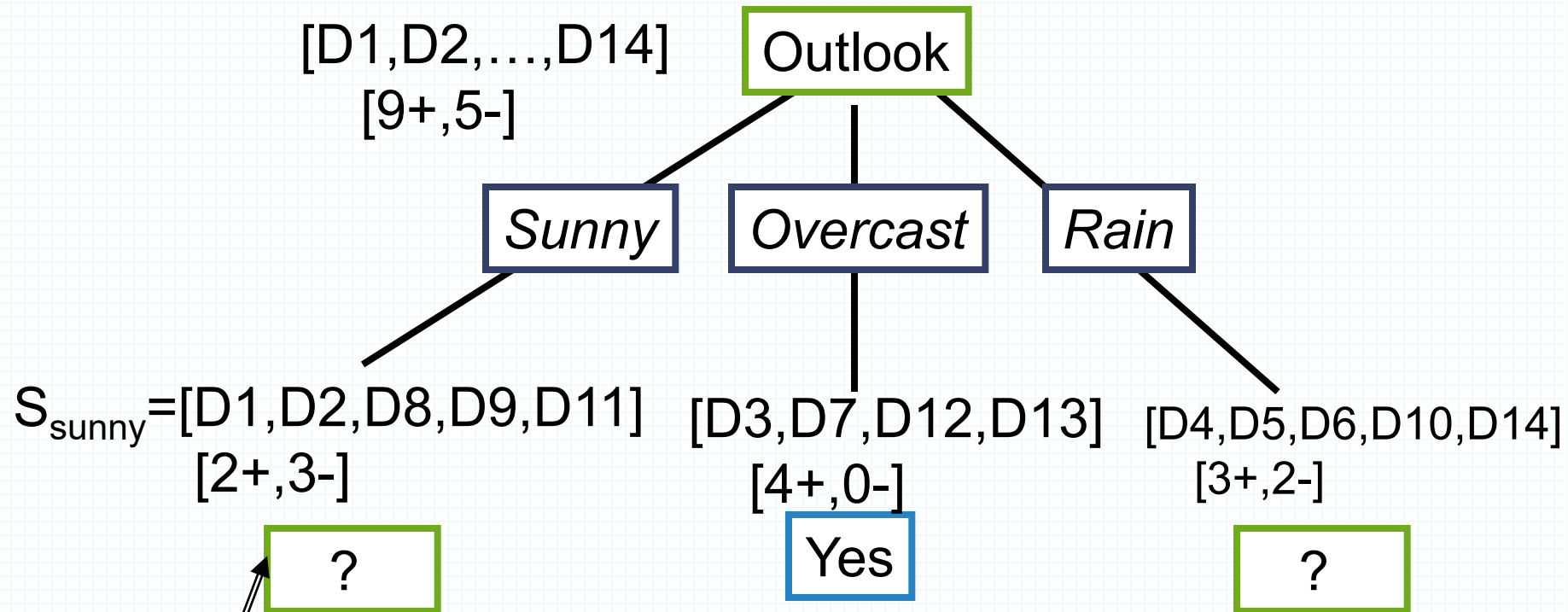
Day	Outlook	Temp.	Humidity	Wind	Play Tennis
D1	Sunny	Hot	High	Weak	No
D2	Sunny	Hot	High	Strong	No
D3	Overcast	Hot	High	Weak	Yes
D4	Rain	Mild	High	Weak	Yes
D5	Rain	Cool	Normal	Weak	Yes
D6	Rain	Cool	Normal	Strong	No
D7	Overcast	Cool	Normal	Weak	Yes
D8	Sunny	Mild	High	Weak	No
D9	Sunny	Cool	Normal	Weak	Yes
D10	Rain	Mild	Normal	Strong	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

Example: Selecting the Best Feature (ID3 Algorithm)



Example: (ID3 Algorithm)

- $\text{Gain}(S, \text{Outlook}) = 0.247$ 😊
- $\text{Gain}(S, \text{Humidity}) = 0.151$
- $\text{Gain}(S, \text{Wind}) = 0.048$
- $\text{Gain}(S, \text{Temperature}) = 0.029$

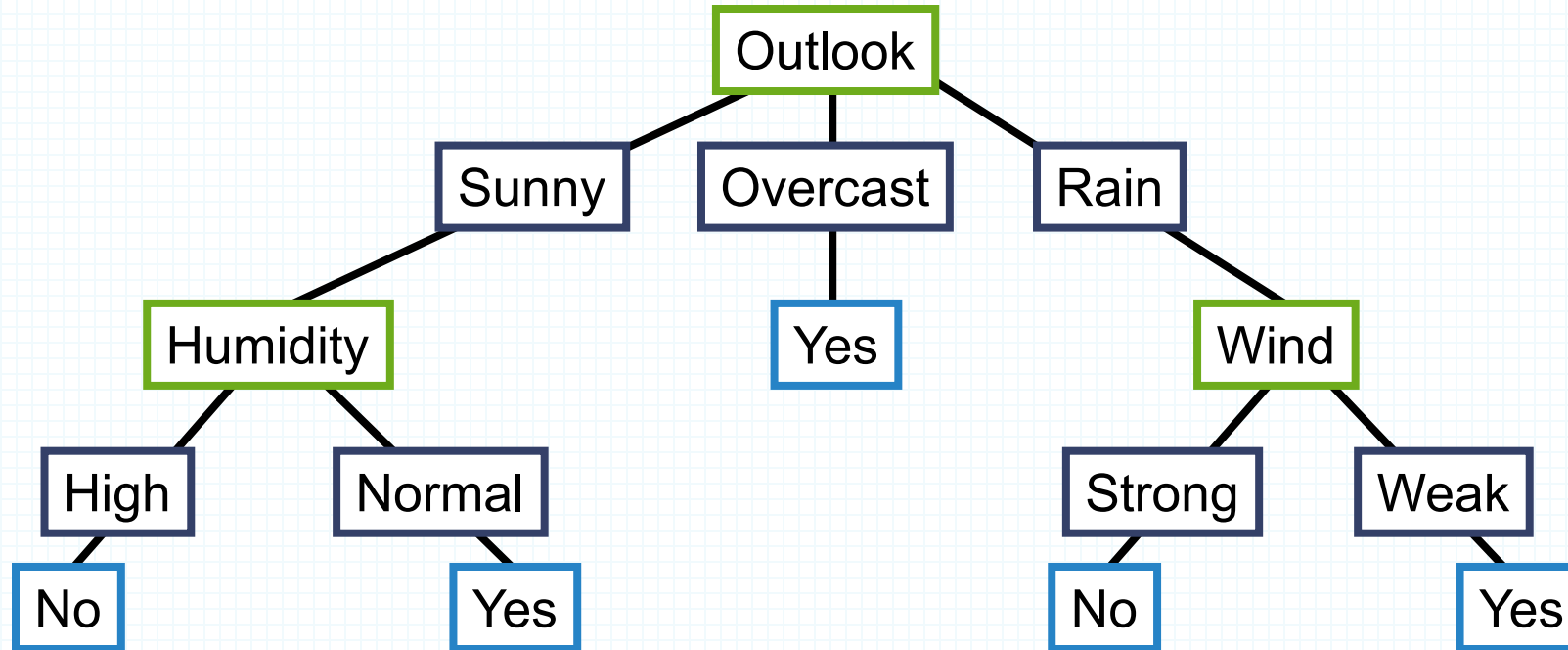


$$\begin{aligned}
 \text{Gain}(S_{\text{sunny}}, \text{Humidity}) &= 0.970 - (3/5)0.0 - 2/5(0.0) = 0.970 \\
 \text{Gain}(S_{\text{sunny}}, \text{Temp.}) &= 0.970 - (2/5)0.0 - 2/5(1.0) - (1/5)0.0 = 0.570 \\
 \text{Gain}(S_{\text{sunny}}, \text{Wind}) &= 0.970 - (2/5)1.0 - 3/5(0.918) = 0.019
 \end{aligned}$$

When should we stop

- All samples for a given node belong to **the same class**.
 - If all instances are "Yes", there's no need to ask more questions.
- There are **no remaining features** for further partitioning.
 - You ran out of features, but still have both "Yes" and "No" samples → choose the majority class (say, "Yes").
- There are **no samples left**
 - If a split condition results in an empty group, fill it with the majority class of the parent node.

Converting a Tree to Rules



- R_1 : If (Outlook=Sunny) \wedge (Humidity=High) Then PlayTennis=No
 R_2 : If (Outlook=Sunny) \wedge (Humidity=Normal) Then PlayTennis=Yes
 R_3 : If (Outlook=Overcast) Then PlayTennis=Yes
 R_4 : If (Outlook=Rain) \wedge (Wind=Strong) Then PlayTennis=No
 R_5 : If (Outlook=Rain) \wedge (Wind=Weak) Then PlayTennis=Yes

Challenges: (1) Continuous features

Convert continuous features to discrete intervals

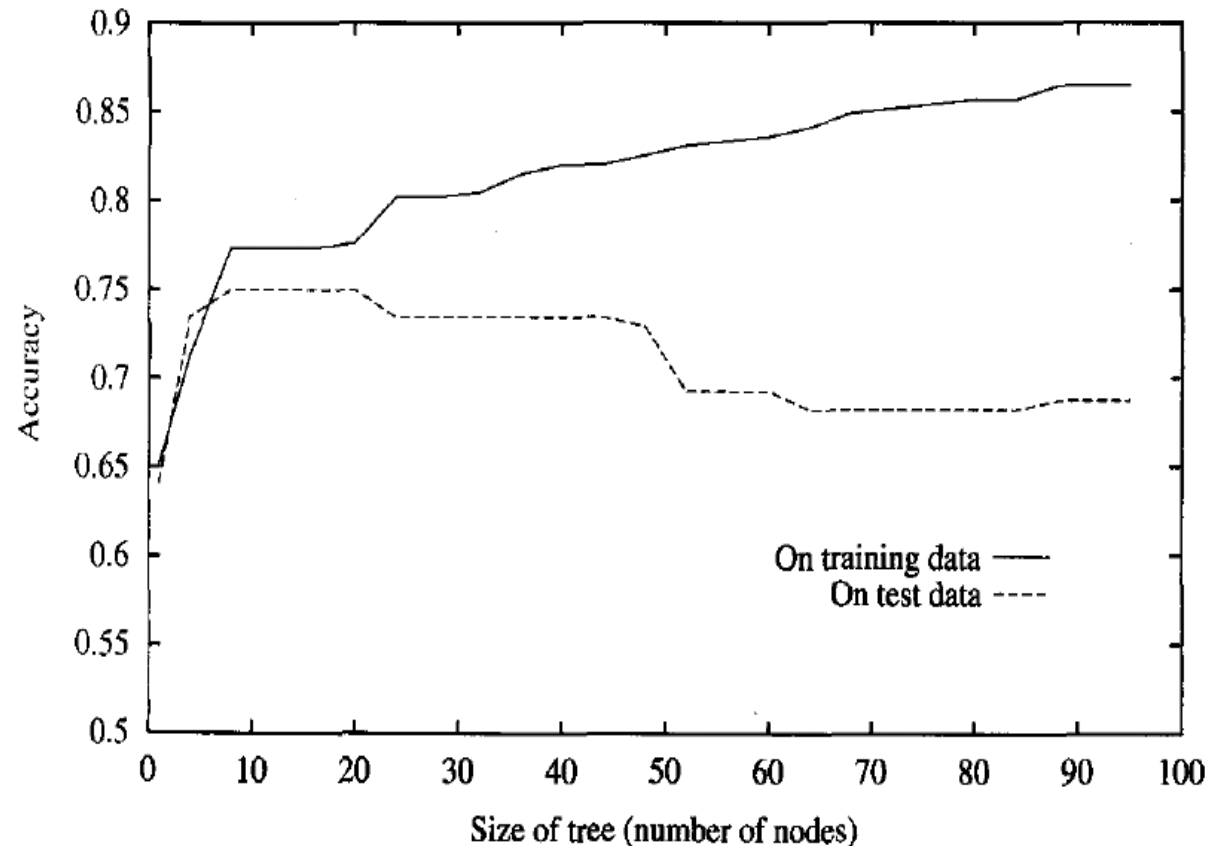
- Temperature=24°C, Temperature =27°C
- (Temperature > 20.0°C) = {true, false}

How to recognise a good threshold?

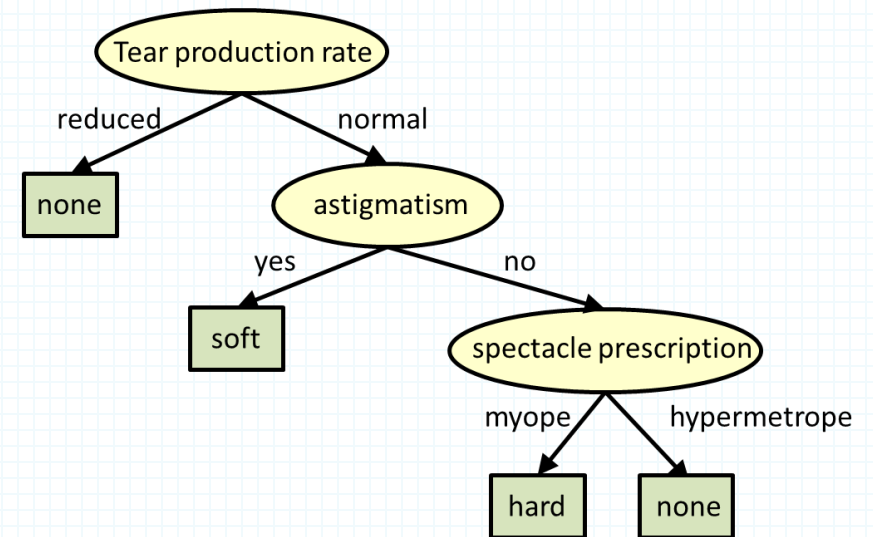
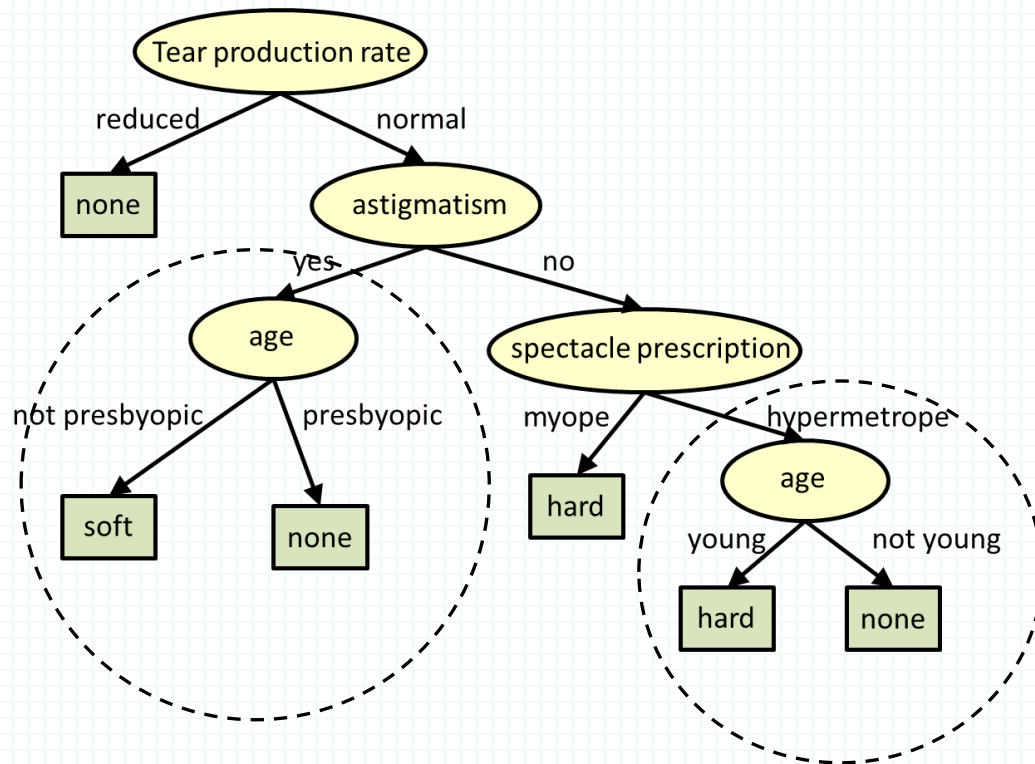
Temperature	15°C	18°C	19°C	22°C	24°C	27°C
Tennis	No	No	Yes	Yes	Yes	No

Challenges: (2) Overfitting

- The idea: Stop growing the tree if the improvement is minor
- Pruning Decision Trees involves a set of techniques that can be used to simplify a Decision Tree, and enable it to generalize better (to avoid overtraining).
- Overfitting occurs when the model trains too well and starts learning noise other than the required characteristics



Challenges: (2) Overfitting



Challenges: (2) Overfitting

- Some Possible methods to apply pruning:
 - **Max depth:** controls the maximum depth of the tree.
 - **Min samples split:** controls the minimum number of samples that are required to split an internal node (prevent the tree from creating overfitting by creating branches with very few samples)
 - **Threshold for Class-Related Samples:** Consider setting a threshold for the percentage of samples related to a specific class.
 - Experimentation and evaluation are key to determining the most effective threshold for your particular scenario.

Challenges: (3) Attributes with many values

- Information Gain tends to prefer Attributes with many values.
- It would prefer to split on say the 6-valued feature, over the current 2- or 3-valued attributes.

[10 10 10 10 10 10]

[30 30]

- Solution?

$$\text{GainRatio}(S, A) = \frac{\text{Gain}(S, A)}{\text{Split}(S, A)}$$

"Why don't we just divide Information Gain by the number of values the attribute has?"

[10 10 10 10 10 10]

[25 25 2 2 3 3]

[30 30]

example

Challenges: (3) Attributes with many values

- Solution?

$$\text{GainRatio}(S, A) = \frac{\text{Gain}(S, A)}{\text{Split}(S, A)}$$

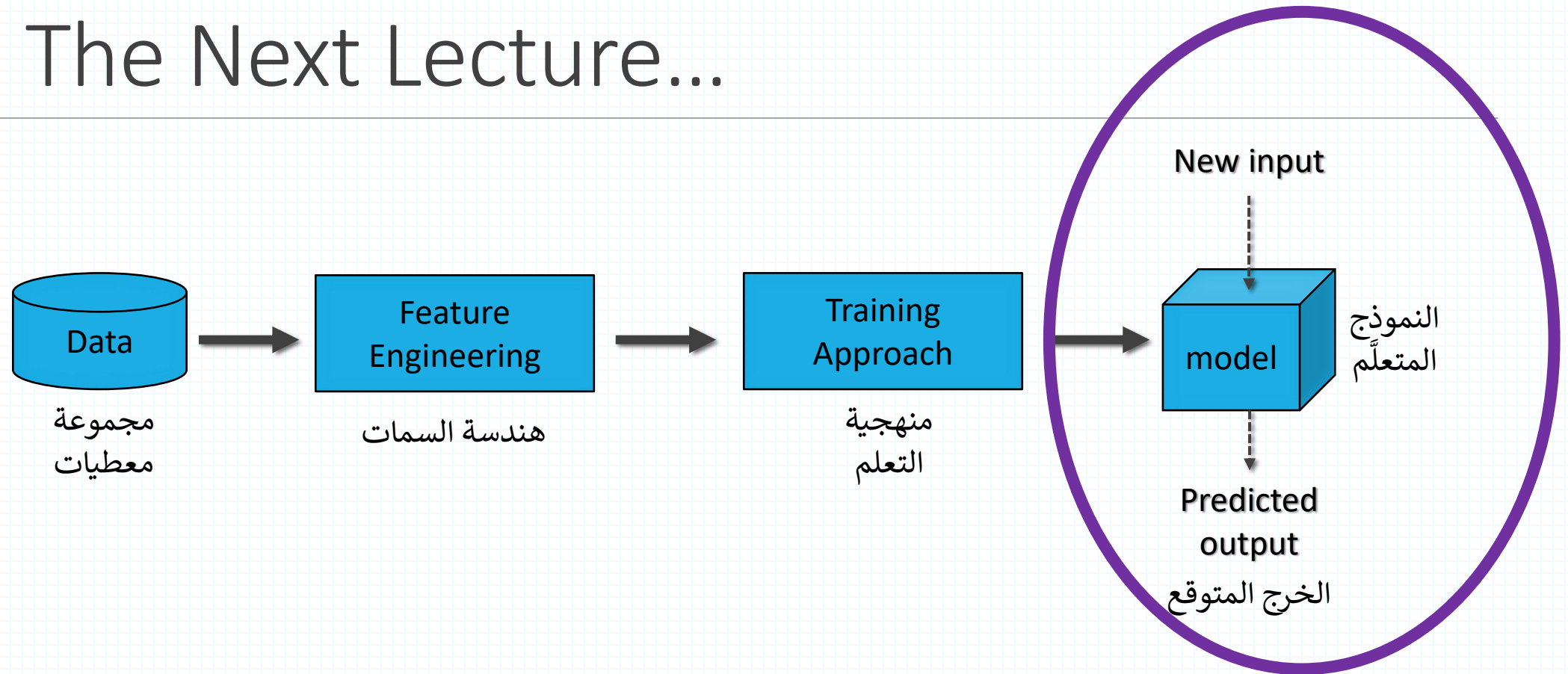
- Split Information ($\text{Split}(S, A)$) measures **how broadly** and **evenly** the data is split when you use feature **A**.

$$\text{Split}(S, A) = - \sum_{i=1}^c \frac{|S_i|}{|S|} \cdot \log_2 \left(\frac{|S_i|}{|S|} \right)$$

Where:

- c : number of distinct values of feature A
- S_i : subset of S where feature $A=v_i$
- $|S_i|$: number of instances in subset S_i
- $|S|$: total number of instances

In The Next Lecture...



**Estimation Strategies
and
Evaluation Criteria**