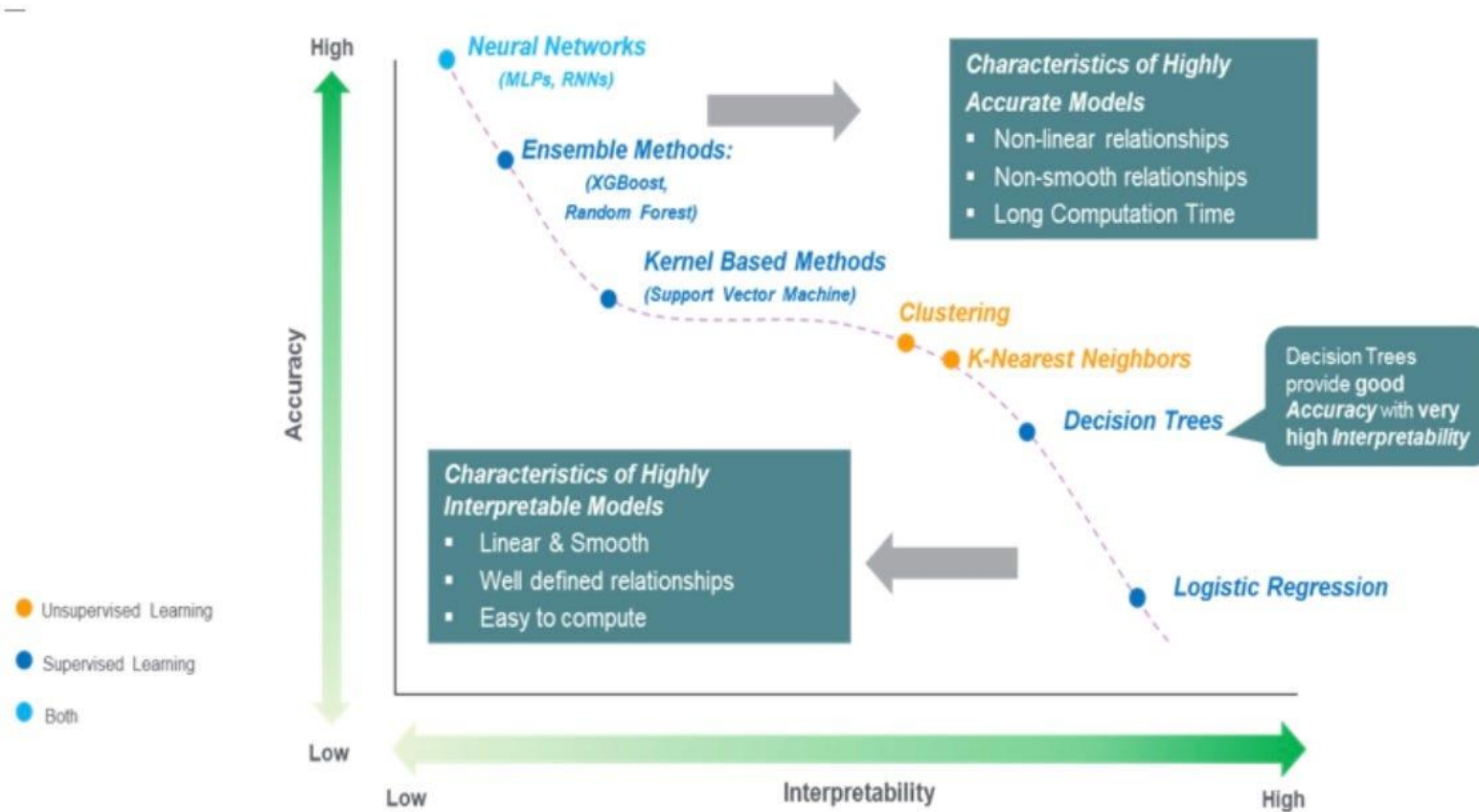


Classification Methods

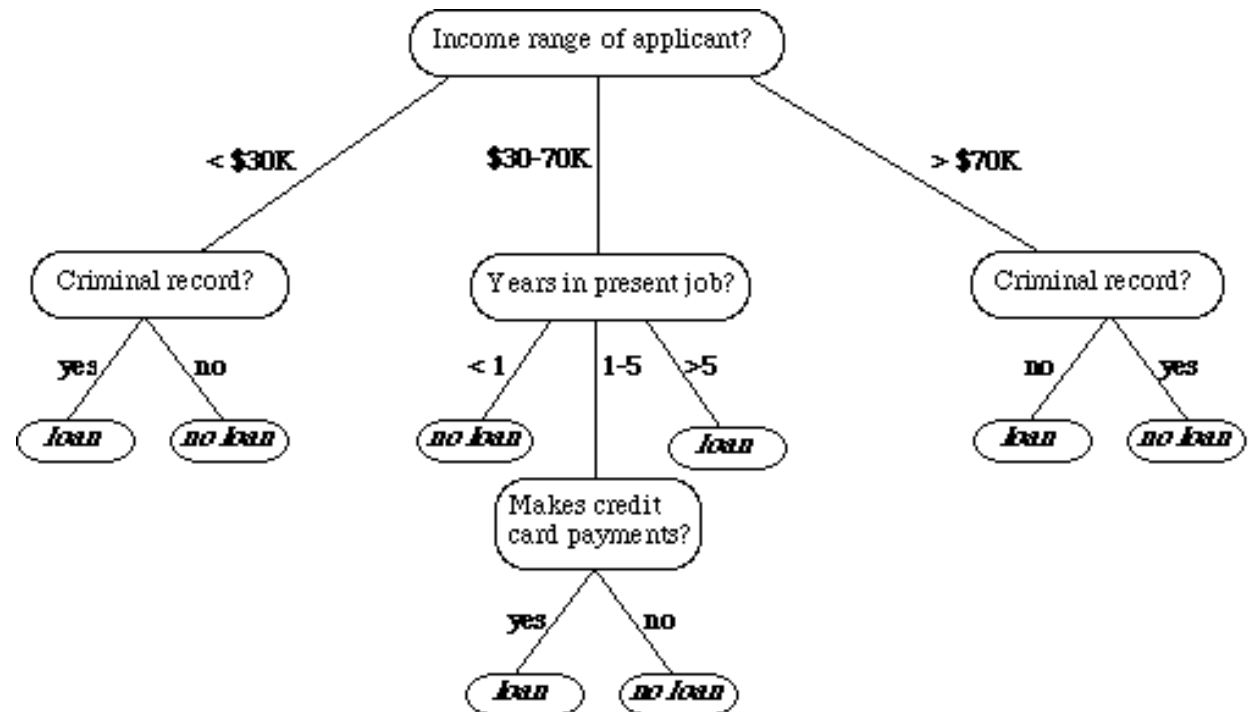
- Logistic Regression
- Naïve Bayes
- k-Nearest Neighbors
- **Decision Trees**
- Ensemble Methods (Boosting, Random Forests)
- Support Vector Machines
- Neural Networks

Generalization capability vs Interpretability



Bank loan

- what kind of income does the person have?
- How long have they held their current job?



Input

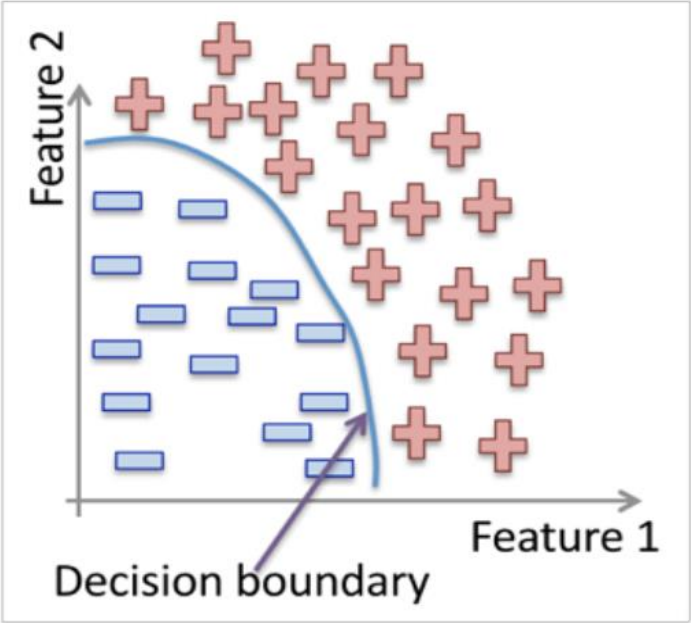
Input $\mathbf{x} \in \mathcal{R}^d$

and output y a label. Learn a function

$$f : \mathbf{x} \rightarrow y$$

Output

example $x_1 \rightarrow$	x_{11}	x_{12}	\dots	x_{1d}	$y_1 \leftarrow \text{label}$
\dots	\dots	\dots	\dots	\dots	\dots
example $x_i \rightarrow$	x_{i1}	x_{i2}	\dots	x_{id}	$y_i \leftarrow \text{label}$
\dots	\dots	\dots	\dots	\dots	\dots
example $x_n \rightarrow$	x_{n1}	x_{n2}	\dots	x_{nd}	$y_n \leftarrow \text{label}$



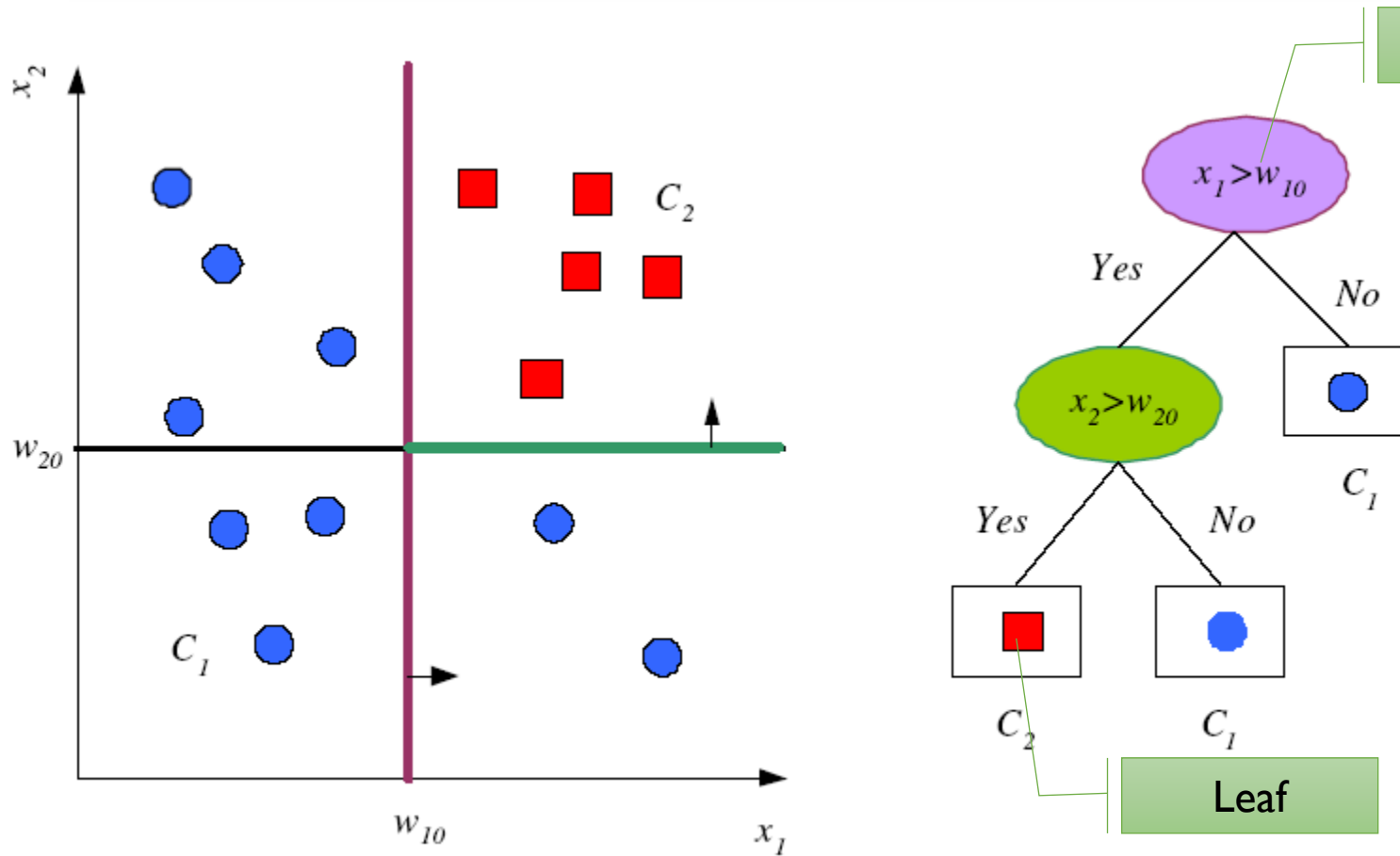
fruit	length	width	weight	label
fruit 1	165	38	172	Banana
fruit 2	218	39	230	Banana
fruit 3	76	80	145	Orange
fruit 4	145	35	150	Banana
fruit 5	90	88	160	Orange
...				
fruit n

Example

PlayTennis: training examples

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

Decision Trees

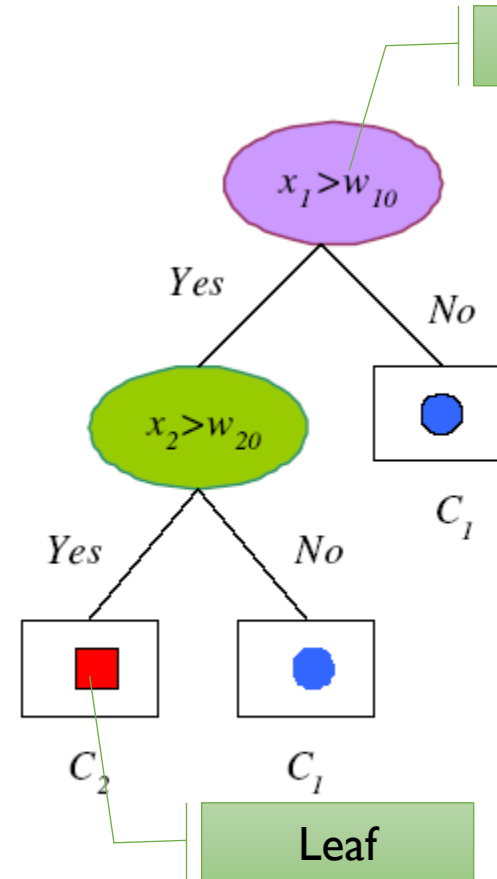
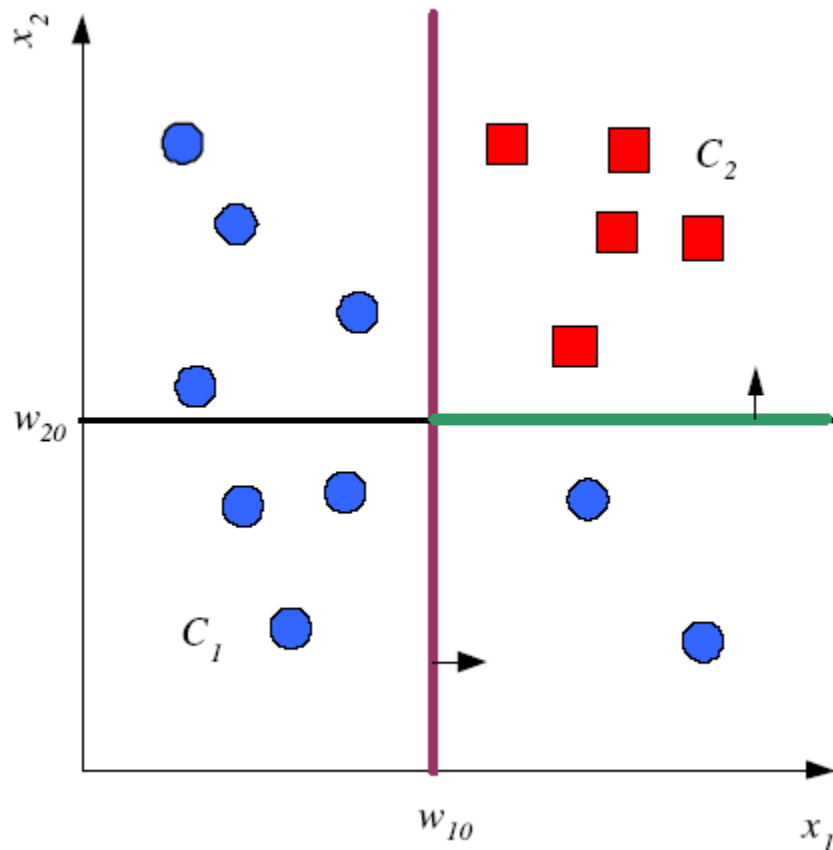


- An efficient nonparametric method
- A hierarchical model
- Divide-and-conquer strategy

Source: Ethem Alpaydin, Introduction to Machine Learning, 3rd Edition (Slides)

Decision Trees

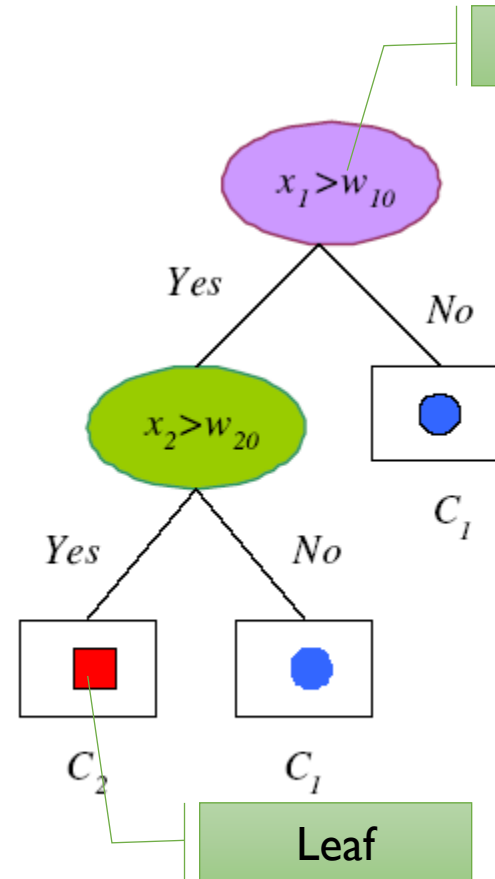
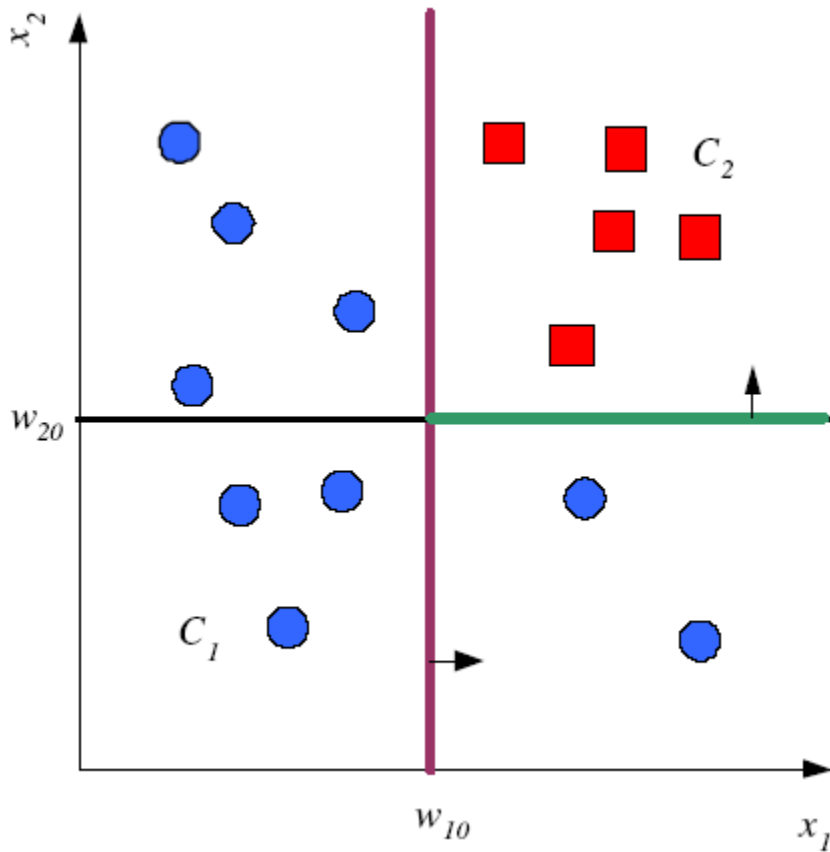
A Non-parametric method means that there are no underlying assumptions about the distribution of the data.



- An efficient nonparametric method
- A hierarchical model
- Divide-and-conquer strategy

Source: Ethem Alpaydin, Introduction to Machine Learning, 3rd Edition (Slides)

Decision Trees



Hierarchical means the model is defined by a series of questions that lead to a class label or a value when applied to any observation.

- A hierarchical model
- Divide-and-conquer strategy

Source: Ethem Alpaydin, Introduction to Machine Learning, 3rd Edition (Slides)

Divide and Conquer

- Internal decision nodes
 - **Univariate:** Uses a single attribute, x_i
 - Numeric x_i :
 - Binary split : $x_i > w_m$
 - Discrete x_i :
 - n -way split for n possible values
 - **Multivariate:** Uses more than one attributes, \mathbf{x}
- Leaves
 - Classification: Class labels, or proportions
 - Regression: Numeric; r average, or local fit
- Learning is **greedy**; find the best split recursively

Source: Ethem Alpaydin, Introduction to Machine Learning, 3rd Edition (Slides)

Classification Trees (C4.5, J48)

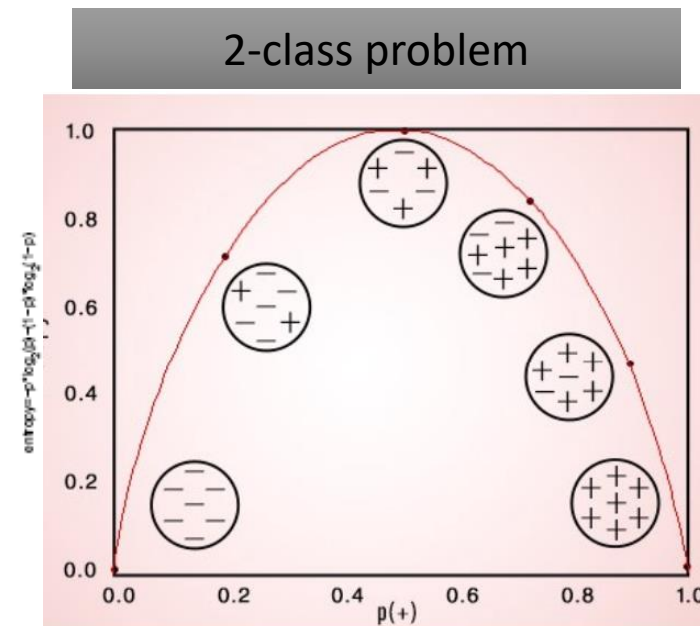
- For node m , N_m instances reach m , N_m^i belong to C_i

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

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

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

- **Gini impurity/index:** $1 - \sum_{j=1}^c p_j^2$



Source: Ethem Alpaydin, Introduction to Machine Learning, 3rd Edition (Slides)

Classification Trees

- 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

$$\hat{P}(C_i | \mathbf{x}, m, j) \equiv p_{mj}^i = \frac{N_{mj}^i}{N_{mj}}$$
$$\mathcal{I}'_m = - \sum_{j=1}^n \frac{N_{mj}}{N_m} \sum_{i=1}^K p_{mj}^i \log_2 p_{mj}^i$$

- **Information Gain:** Expected reduction in impurity measure after split
- Choose the *best* attribute(s) (**with maximum information gain**) to split the remaining instances and make that attribute a decision node
 - You can use same logic to find best splitting value too

Source: Ethem Alpaydin, Introduction to Machine Learning, 3rd Edition (Slides)

Algorithm 1 **DECISIONTREETRAIN**(*data*, *remaining features*)

```
1: guess  $\leftarrow$  most frequent answer in data           // default answer for this data
2: if the labels in data are unambiguous then
3:   return LEAF(guess)                               // base case: no need to split further
4: else if remaining features is empty then
5:   return LEAF(guess)                               // base case: cannot split further
6: else                                                 // we need to query more features
7:   for all  $f \in \textit{remaining features}$  do
8:     NO  $\leftarrow$  the subset of data on which  $f=no$ 
9:     YES  $\leftarrow$  the subset of data on which  $f=yes$ 
10:    score[f]  $\leftarrow$  # of majority vote answers in NO
11:                    + # of majority vote answers in YES
                                // the accuracy we would get if we only queried on f
12:   end for
13:   f  $\leftarrow$  the feature with maximal score(f)
14:   NO  $\leftarrow$  the subset of data on which  $f=no$ 
15:   YES  $\leftarrow$  the subset of data on which  $f=yes$ 
16:   left  $\leftarrow$  DECISIONTREETRAIN(NO, remaining features  $\setminus \{f\}$ )
17:   right  $\leftarrow$  DECISIONTREETRAIN(YES, remaining features  $\setminus \{f\}$ )
18:   return NODE(f, left, right)
19: end if
```

Algorithm 2 `DECISIONTREETEST`(*tree*, *test point*)

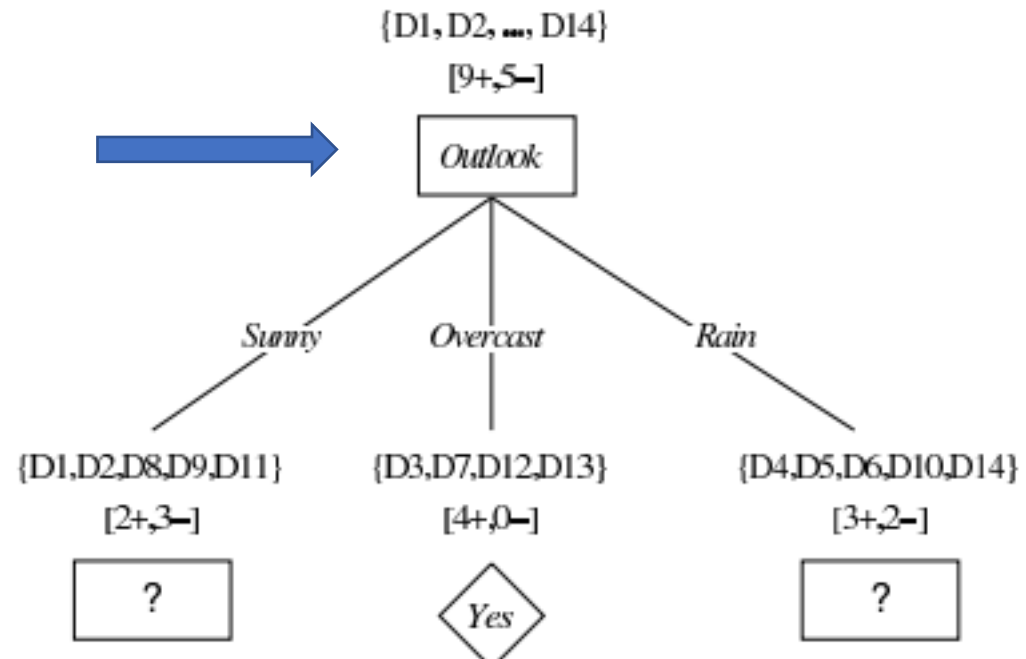
```
1: if tree is of the form LEAF(guess) then  
2:   return guess  
3: else if tree is of the form NODE(f, left, right) then  
4:   if f = no in test point then  
5:     return DECISIONTREETEST(left, test point)  
6:   else  
7:     return DECISIONTREETEST(right, test point)  
8:   end if  
9: end if
```

Decision Trees: Example

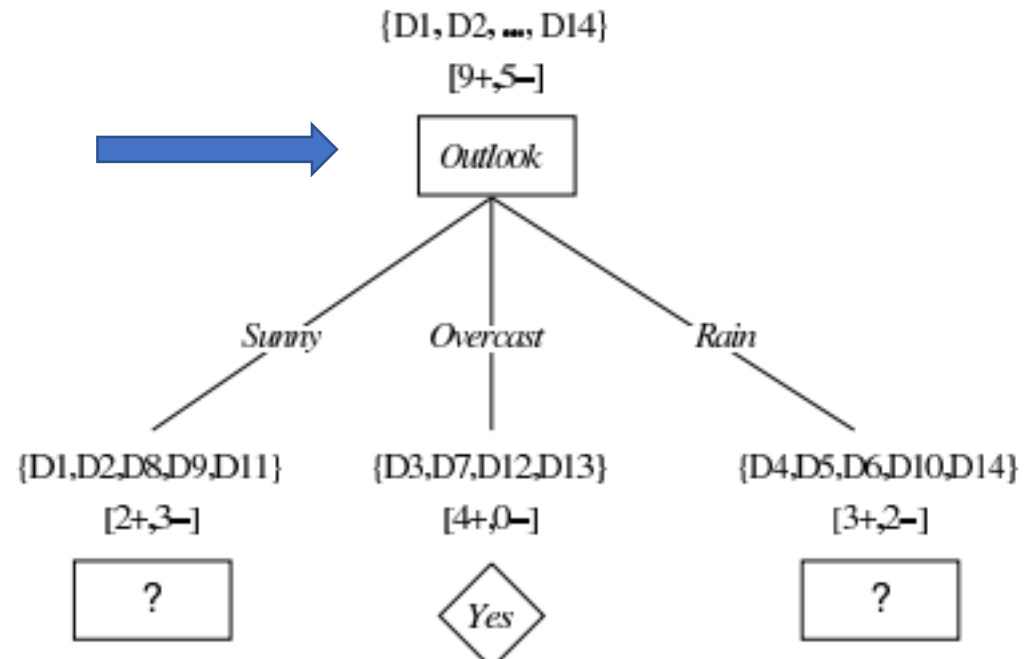
PlayTennis: training examples

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

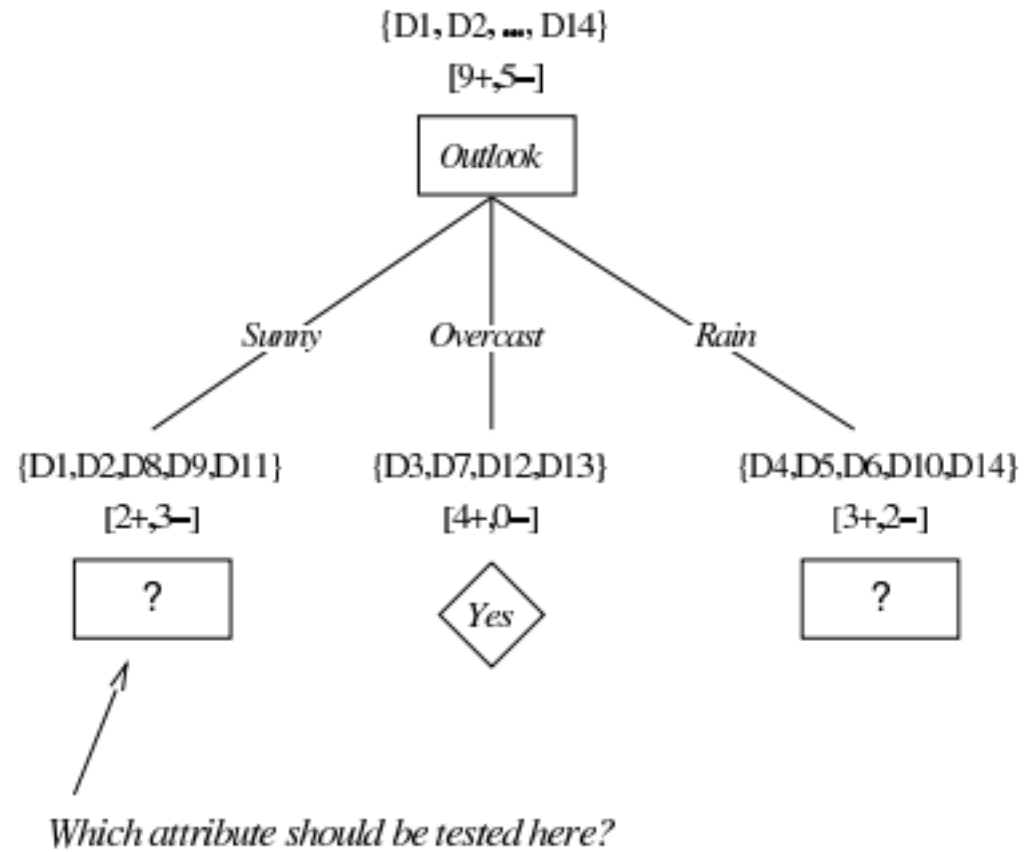
Decision Trees: Example



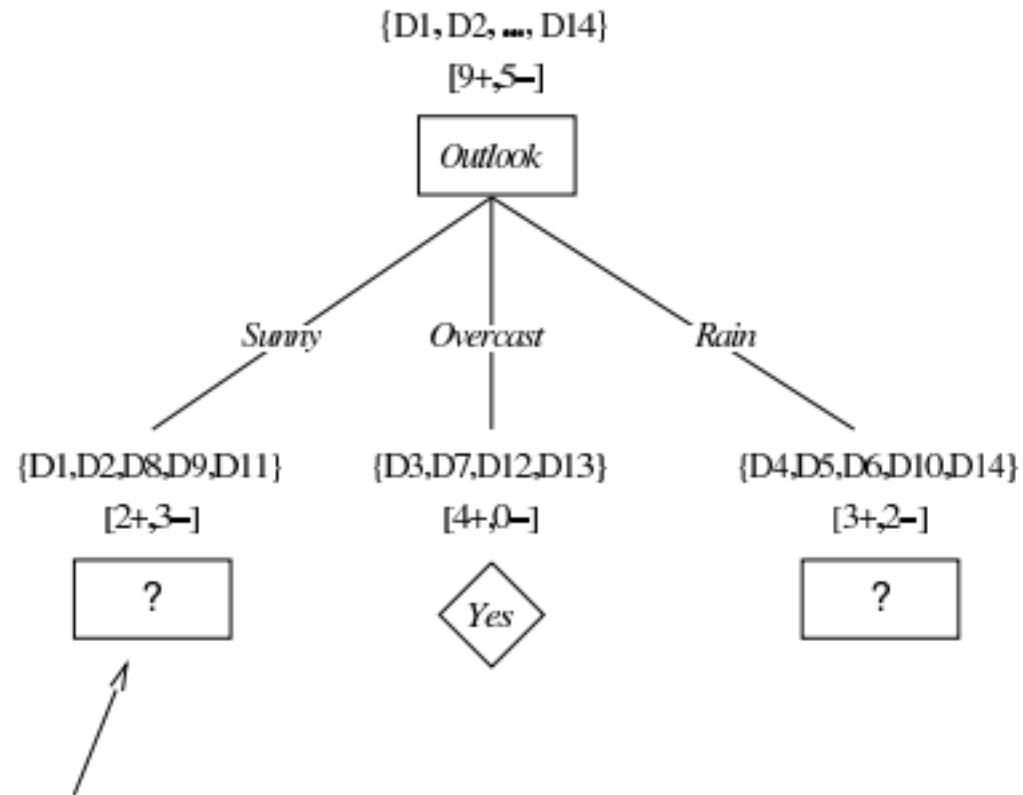
Decision Trees: Example



Decision Trees: Example



Decision Trees: Example



Which attribute should be tested here?

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

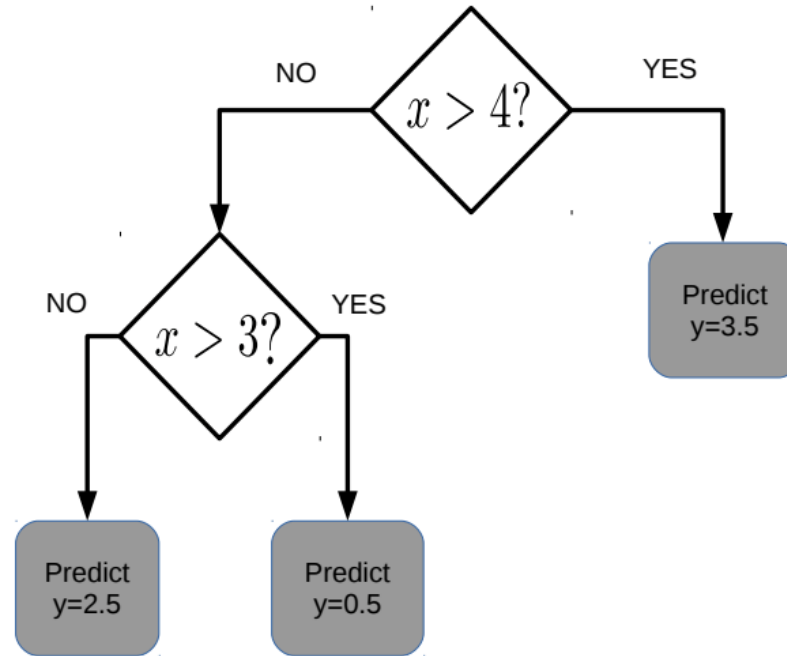
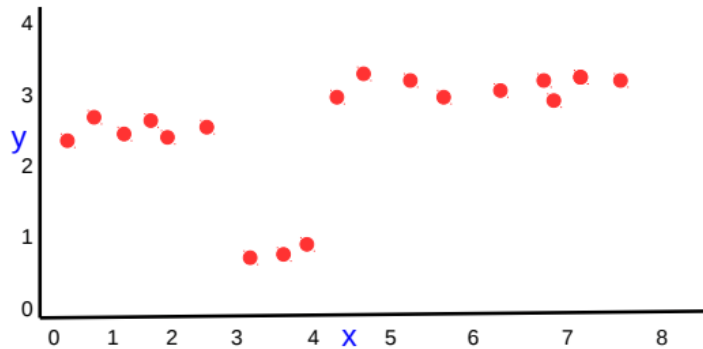
$$\text{Gain}(S_{\text{Sunny}}, \text{Humidity}) = .970 - (3/5) 0.0 - (2/5) 0.0 = .970$$

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

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

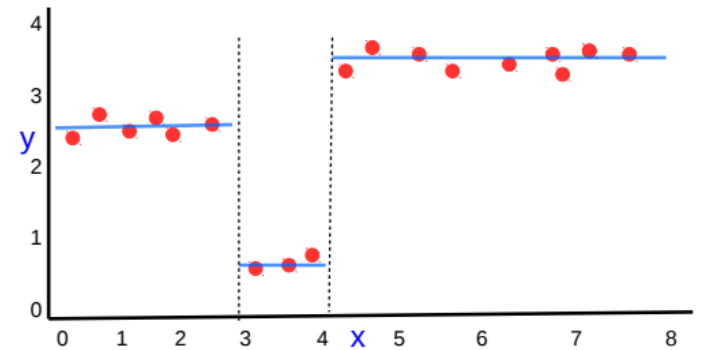
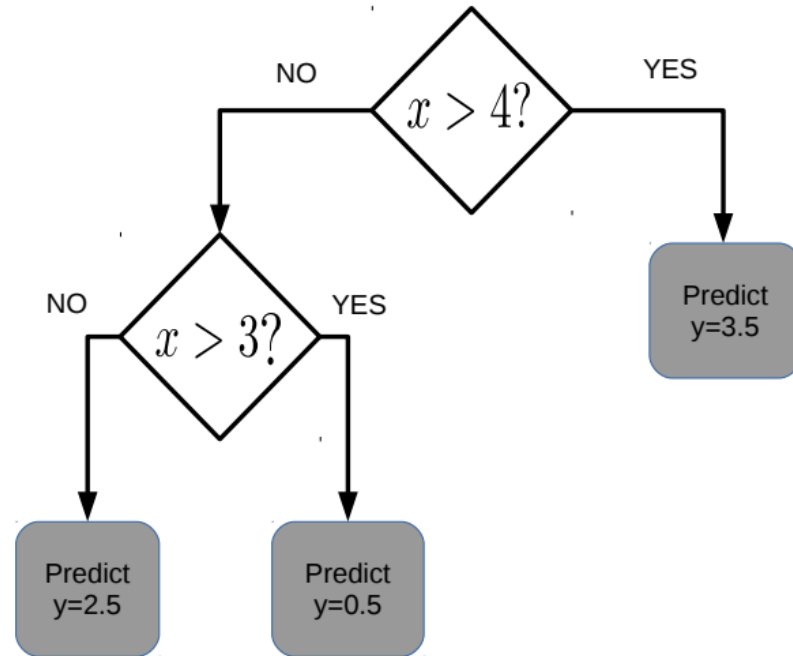
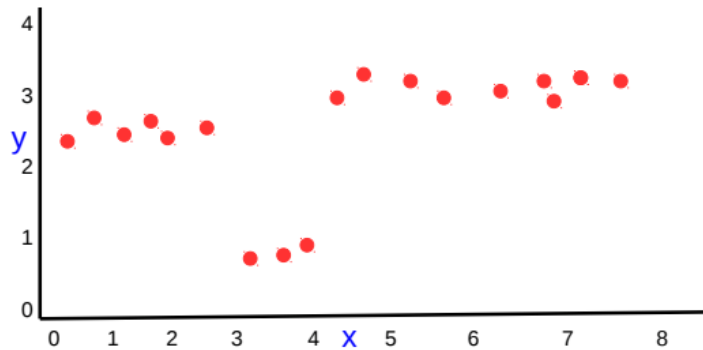
Decision Trees for Regression

Decision Trees can also be used for regression problems



Decision Trees for Regression

Decision Trees can also be used for regression problems



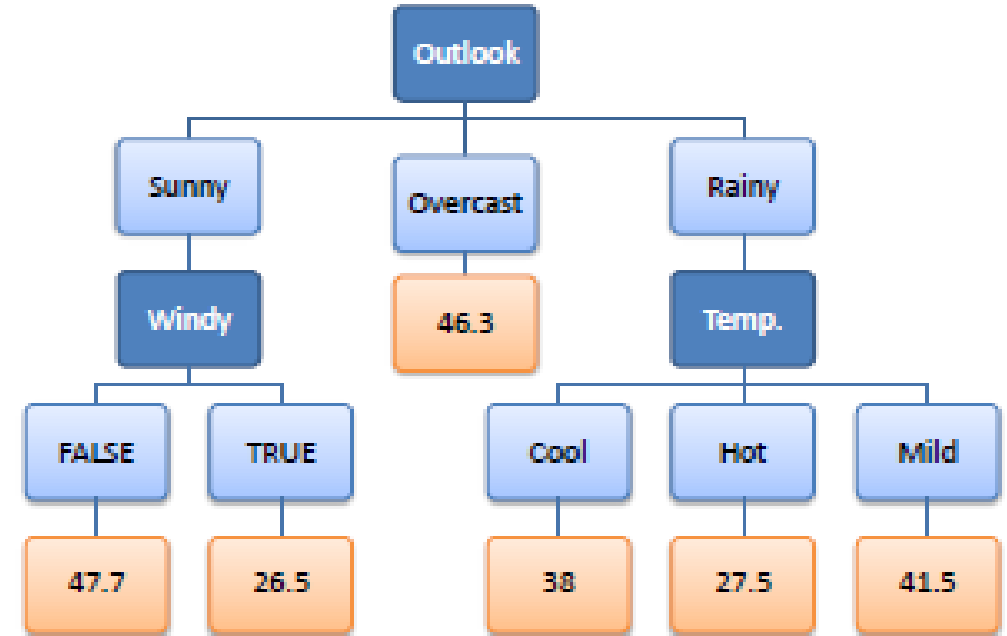
Decision Trees for Regression

- Entropy as a measure of impurity is a useful criteria for classification.
- impurity measure using the **weighted mean squared error (MSE)** of the children nodes or variance

$$\text{MSE}(t) = \frac{1}{N_t} \sum_{i \in D_t} (y^{(i)} - \hat{y}_t)^2 \qquad \hat{y}_t = \frac{1}{N_t} \sum_{i \in D_t} y^{(i)}$$

- Attribute chosen so that it reduces variance or standard deviation

Predictors				Target
Outlook	Temp	Humidity	Windy	Hours Played
Rainy	Hot	High	FALSE	25
Rainy	Hot	High	TRUE	30
Overcast	Hot	High	FALSE	46
Sunny	Mild	High	FALSE	45
Sunny	Cool	Normal	FALSE	52
Sunny	Cool	Normal	TRUE	23
Overcast	Cool	Normal	TRUE	43
Rainy	Mild	High	FALSE	35
Rainy	Cool	Normal	FALSE	38
Sunny	Mild	Normal	FALSE	46
Rainy	Mild	Normal	TRUE	48
Overcast	Mild	High	TRUE	52
Overcast	Hot	Normal	FALSE	44
Sunny	Mild	High	TRUE	30



Hours Played
25
30
46
45
52
23
43
35
38
46
48
52
44
30

$$\text{Count} = n = 14$$

$$\text{Average} = \bar{x} = \frac{\sum x}{n} = 39.8$$

$$\Rightarrow \text{Standard Deviation} = S = \sqrt{\frac{\sum (x - \bar{x})^2}{n}} = 9.32$$

$$\text{Coefficient of Variation} = CV = \frac{S}{\bar{x}} * 100\% = 23\%$$

$$S(T, X) = \sum_{c \in X} P(c) S(c)$$

		Hours Played (StDev)	Count
Outlook	Overcast	3.49	4
	Rainy	7.78	5
	Sunny	10.87	5
			14



$$\begin{aligned}
 S(\text{Hours}, \text{Outlook}) &= P(\text{Sunny}) * S(\text{Sunny}) + P(\text{Overcast}) * S(\text{Overcast}) + P(\text{Rainy}) * S(\text{Rainy}) \\
 &= (4/14) * 3.49 + (5/14) * 7.78 + (5/14) * 10.87 \\
 &= 7.66
 \end{aligned}$$

$$SDR(T, X) = S(T) - S(T, X)$$

$$\begin{aligned}
 SDR(\text{Hours}, \text{Outlook}) &= S(\text{Hours}) - S(\text{Hours}, \text{Outlook}) \\
 &= 9.32 - 7.66 = 1.66
 \end{aligned}$$

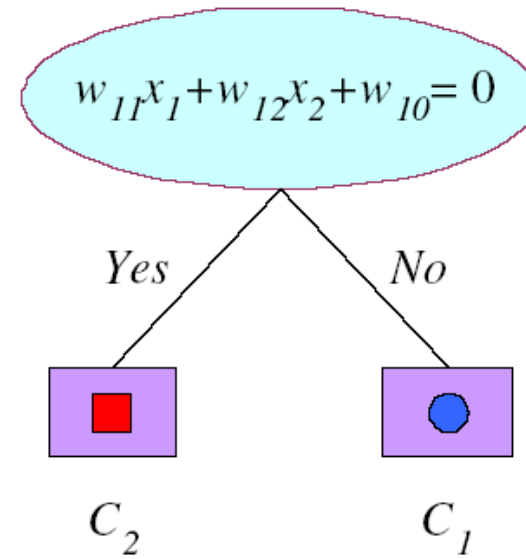
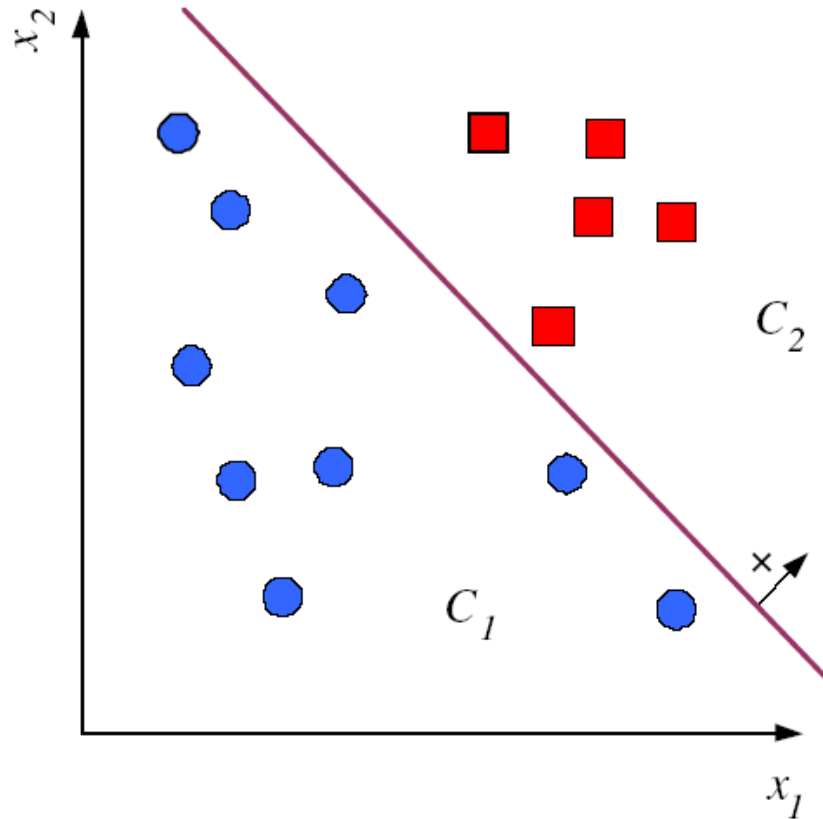
		Hours Played (StDev)
Outlook	Overcast	3.49
	Rainy	7.78
	Sunny	10.87
		SDR=1.66

		Hours Played (StDev)
Humidity	High	9.36
	Normal	8.37
		SDR=0.28

		Hours Played (StDev)
Temp.	Cool	10.51
	Hot	8.95
	Mild	7.65
		SDR=0.17

		Hours Played (StDev)
Windy	False	7.87
	True	10.59
		SDR=0.29

Multivariate Trees



Overfitting and Generalization

- Overfitting can occur with noisy training examples, also when small numbers of examples are associated with leaf nodes. How to handle?
- **Pruning:** 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

Overfitting and Generalization

- **Occam's Razor principle:** when multiple hypotheses can solve a problem, choose the simplest one
- How to select “best” tree:
 - Measure performance over separate validation data set
 - **Minimum Description Length:** Minimize $size(tree) + size(misclassifications(tree))$

Pruning based on validation data

- Reduced-error pruning, is to consider each of the decision nodes in the tree to be candidates for pruning
- Pruning a decision node consists of removing the subtree rooted at that node, making it a leaf node, and assigning it the most common classification of the training examples affiliated with that node
- Nodes are removed only if the resulting pruned tree performs no worse than-the original over the validation set.
- Reduced error pruning has the effect that any leaf node added due to coincidental regularities in the training set is likely to be pruned because these same coincidences are unlikely to occur in the validation set