# Class 11

Supervised learning

# Learning Objectives

- Algorithms
  - Decision trees
  - Naïve Bayesian
  - Artificial Neural Networks
- Evaluation methods
  - Precision

# Goals and Requirements

- Goals:
  - To produce an accurate classifier/regression function
  - To understand the structure of the problem
- Requirements on the model:
  - High accuracy
  - Understandable by humans, interpretable
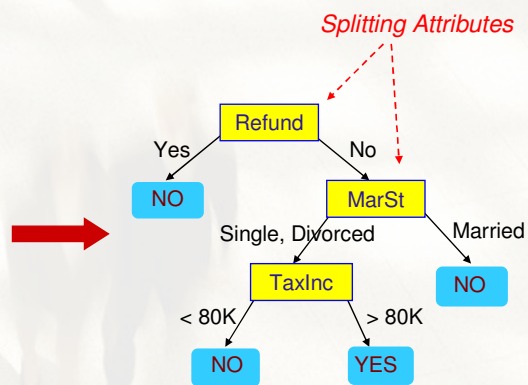  - Fast construction for very large training databases

3

# Example of a Decision Tree

| Tid | Refund | Marital Status | Taxable Income | Cheat |
|-----|--------|----------------|----------------|-------|
| 1 | Yes | Single | 125K | No |
| 2 | No | Married | 100K | No |
| 3 | No | Single | 70K | No |
| 4 | Yes | Married | 120K | No |
| 5 | No | Divorced | 95K | Yes |
| 6 | No | Married | 60K | No |
| 7 | Yes | Divorced | 220K | No |
| 8 | No | Single | 85K | Yes |
| 9 | No | Married | 75K | No |
| 10 | No | Single | 90K | Yes |

*categorical* *categorical* *continuous* *class*

Training Data

Splitting Attributes

Refund
Yes → NO
No → MarSt
Single, Divorced → TaxInc
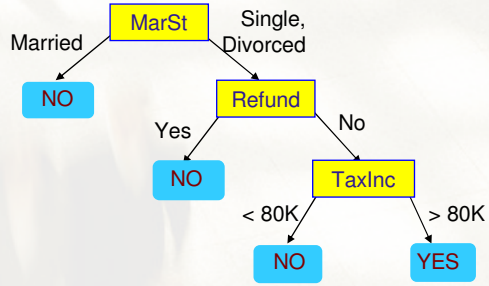Married → NO
TaxInc < 80K → NO
TaxInc > 80K → YES

Model: Decision Tree

4

# Another Example of Decision Tree

categorical categorical continuous class

| Tid | Refund | Marital Status | Taxable Income | Cheat |
|-----|--------|----------------|----------------|-------|
| 1 | Yes | Single | 125K | No |
| 2 | No | Married | 100K | No |
| 3 | No | Single | 70K | No |
| 4 | Yes | Married | 120K | No |
| 5 | No | Divorced | 95K | Yes |
| 6 | No | Married | 60K | No |
| 7 | Yes | Divorced | 220K | No |
| 8 | No | Single | 85K | Yes |
| 9 | No | Married | 75K | No |
| 10 | No | Single | 90K | Yes |

MarSt

Married → NO

Single, Divorced → Refund

Refund: Yes → NO, No → TaxInc
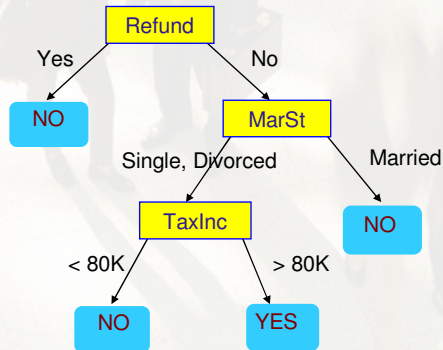
TaxInc: < 80K → NO, > 80K → YES

There could be more than one tree that fits the same data!

---

# Apply Model to Test Data

Start from the root of tree.

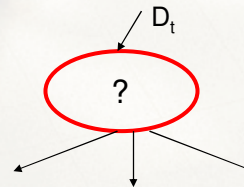### Test Data

| Refund | Marital Status | Taxable Income | Cheat |
|--------|----------------|----------------|-------|
| No | Married | 80K | ? |

Refund

Yes → NO

No → MarSt

MarSt: Single, Divorced → TaxInc, Married → NO

TaxInc: < 80K → NO, > 80K → YES

# General algorithm

l Let $D_t$ be the set of training records that reach a node t

l General Procedure:

- If $D_t$ contains records that belong the same class $y_t$, then t is a leaf node labeled as $y_t$
- If $D_t$ is an empty set, then t is a leaf node labeled by the default class, $y_d$
- If $D_t$ contains records that belong to more than one class, use an attribute test to split the data into smaller subsets. Recursively apply the procedure to each subset.
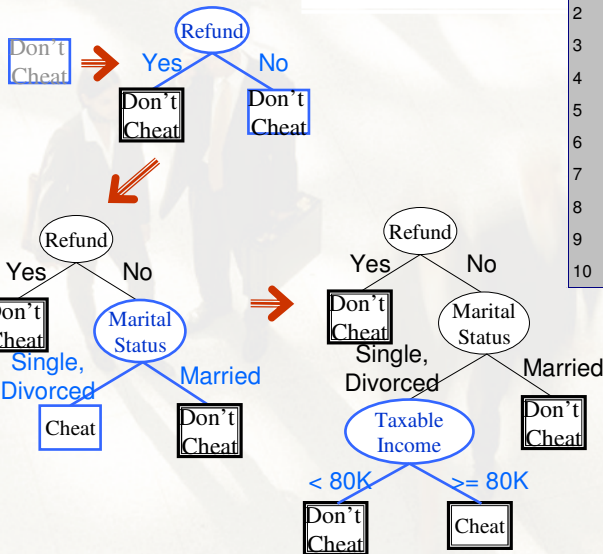
| Tid | Refund | Marital Status | Taxable Income | Cheat |
|-----|--------|----------------|----------------|-------|
| 1 | Yes | Single | 125K | No |
| 2 | No | Married | 100K | No |
| 3 | No | Single | 70K | No |
| 4 | Yes | Married | 120K | No |
| 5 | No | Divorced | 95K | Yes |
| 6 | No | Married | 60K | No |
| 7 | Yes | Divorced | 220K | No |
| 8 | No | Single | 85K | Yes |
| 9 | No | Married | 75K | No |
| 10 | No | Single | 90K | Yes |

$D_t$

?

# Example



| Tid | Refund | Marital Status | Taxable Income | Cheat |
|-----|--------|----------------|----------------|-------|
| 1 | Yes | Single | 125K | No |
| 2 | No | Married | 100K | No |
| 3 | No | Single | 70K | No |
| 4 | Yes | Married | 120K | No |
| 5 | No | Divorced | 95K | Yes |
| 6 | No | Married | 60K | No |
| 7 | Yes | Divorced | 220K | No |
| 8 | No | Single | 85K | Yes |
| 9 | No | Married | 75K | No |
| 10 | No | Single | 90K | Yes |

# Stopping Criteria for Tree Induction

- Stop expanding a node when all the records belong to the same class

- Stop expanding a node when all the records have similar attribute values
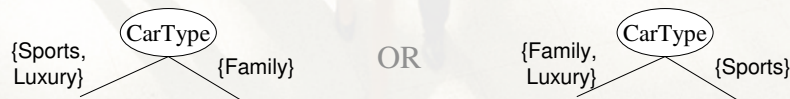
- Early termination (to be discussed later)

9

# Splitting Based on Nominal Attributes

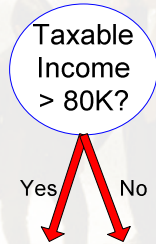- Multi-way split: Use as many partitions as distinct values.

CarType

Family — Sports — Luxury

- Binary split:  Divides values into two subsets. Need to find optimal partitioning.

CarType

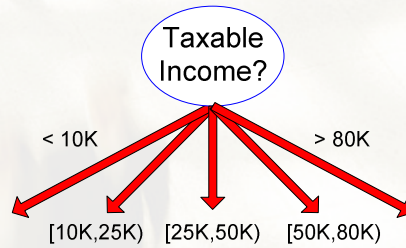{Sports, Luxury} — {Family}

OR

CarType

{Family, Luxury} — {Sports}

10

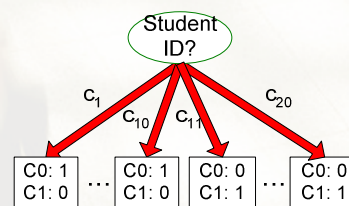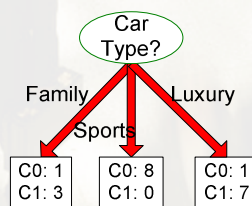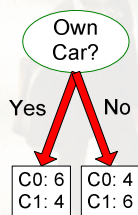# Splitting Based on Continuous Attributes



(i) Binary split

(ii) Multi-way split

# How to determine the Best Split

Before Splitting: 10 records of class 0,
10 records of class 1



Which test condition is the best?

# How to determine the Best Split

l Greedy approach:
  – Nodes with homogeneous class distribution are preferred
l Need a measure of node impurity:

| C0: 5 |
| C1: 5 |

| C0: 9 |
| C1: 1 |

Non-homogeneous,

High degree of impurity

Homogeneous,

Low degree of impurity

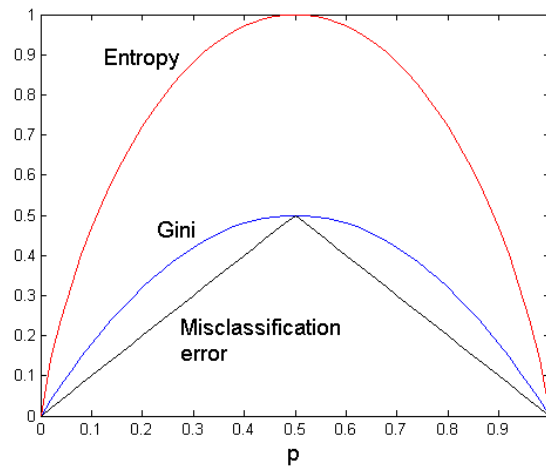# Measure of Node Impurity

l Entropy at a given node t:

$$Entropy(t) = -\sum_{j} p(j\,|\,t)\log p(j\,|\,t)$$

(NOTE: $p(j\,|\,t)$ is the relative frequency of class j at node t).

– Measures homogeneity of a node.
  u Maximum ($\log n_c$) when records are equally distributed among all classes implying least information
  u Minimum (0.0) when all records belong to one class, implying most information

# Entropy function

# Example

$$Entropy(t) = -\sum_j p(j \mid t) \log_2 p(j \mid t)$$

| C1 | 0 |
|----|---|
| C2 | 6 |

P(C1) = 0/6 = 0     P(C2) = 6/6 = 1

Entropy = – 0 log 0 – 1 log 1 = – 0 – 0 = 0

| C1 | 1 |
|----|---|
| C2 | 5 |

P(C1) = 1/6       P(C2) = 5/6

Entropy = – (1/6) $\log_2$ (1/6) – (5/6) $\log_2$ (1/6) = 0.65

| C1 | 2 |
|----|---|
| C2 | 4 |

P(C1) = 2/6       P(C2) = 4/6

Entropy = – (2/6) $\log_2$ (2/6) – (4/6) $\log_2$ (4/6) = 0.92

# Information Gain

- Information Gain:

$$GAIN_{split} = Entropy(p) - \left( \sum_{i=1}^{k} \frac{n_i}{n} Entropy(i) \right)$$

Parent Node, p is split into k partitions;

$n_i$ is number of records in partition i

- – Measures Reduction in Entropy achieved because of the split. Choose the split that achieves most reduction (maximizes GAIN)
- – Used in ID3 and C4.5
- – Disadvantage: Tends to prefer splits that result in large number of partitions, each being small but pure.

# Expressiveness

- Decision tree provides expressive representation for learning discrete-valued function
  - – But they do not generalize well to certain types of Boolean functions
    - u Example: parity function:
      - – Class = 1 if there is an even number of Boolean attributes with truth value = True
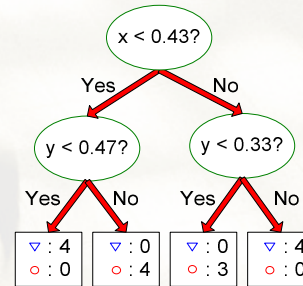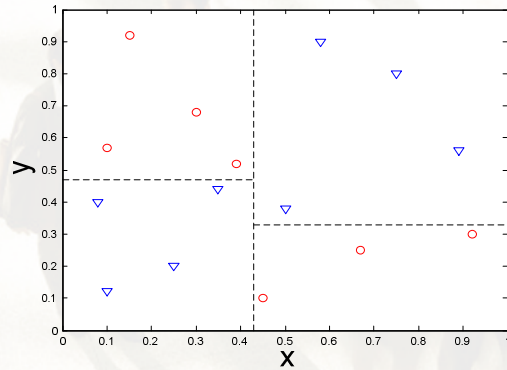      - – Class = 0 if there is an odd number of Boolean attributes with truth value = True
    - u For accurate modeling, must have a complete tree

- Not expressive enough for modeling continuous variables
  - – Particularly when test condition involves only a single attribute at-a-time

# Decision Boundary



• Border line between two neighboring regions of different classes is known as decision boundary

• Decision boundary is parallel to axes because test condition involves a single attribute at-a-time

# Learning Curve



 l Learning curve shows how accuracy changes with varying sample size

l Requires a sampling schedule for creating learning curve:

 l Arithmetic sampling (Langley, et al)

 l Geometric sampling (Provost et al)

Effect of small sample size:

 - Bias in the estimate

 - Variance of estimate

# Decision Trees: Summary

- Many application of decision trees
- There are many algorithms available for:
  - Split selection
  - Pruning
  - Handling Missing Values
  - Data Access
- Decision tree construction still active research area (after 20+ years!)
- Challenges: Performance, scalability, evolving datasets, new applications

# Bayes Classifier

- A probabilistic framework for solving classification problems
- Conditional Probability:

$$P(C \mid A) = \frac{P(A,C)}{P(A)}$$

$$P(A \mid C) = \frac{P(A,C)}{P(C)}$$

- Bayes theorem:

$$P(C \mid A) = \frac{P(A \mid C)P(C)}{P(A)}$$

# Example of Bayes Theorem

- Given:
  - A doctor knows that meningitis causes stiff neck 50% of the time
  - Prior probability of any patient having meningitis is 1/50,000
  - Prior probability of any patient having stiff neck is 1/20

- If a patient has stiff neck, what's the probability he/she has meningitis?

$$P(M \mid S) = \frac{P(S \mid M)P(M)}{P(S)} = \frac{0.5 \times 1/50000}{1/20} = 0.0002$$

# Bayesian Classifiers

- Consider each attribute and class label as random variables

- Given a record with attributes $(A_1, A_2, \ldots, A_n)$
  - Goal is to predict class C
  - Specifically, we want to find the value of C that maximizes $P(C \mid A_1, A_2, \ldots, A_n)$

- Can we estimate $P(C \mid A_1, A_2, \ldots, A_n)$ directly from data?

# Bayesian Classifiers

- Approach:
  - compute the posterior probability $P(C \mid A_1, A_2, \ldots, A_n)$ for all values of C using the Bayes theorem

$$P(C \mid A_1 A_2 \mathrm{K} \; A_n) = \frac{P(A_1 A_2 \mathrm{K} \; A_n \mid C) P(C)}{P(A_1 A_2 \mathrm{K} \; A_n)}$$

  - Choose value of C that maximizes
    $P(C \mid A_1, A_2, \ldots, A_n)$

  - Equivalent to choosing value of C that maximizes
    $P(A_1, A_2, \ldots, A_n \mid C) \; P(C)$

- How to estimate $P(A_1, A_2, \ldots, A_n \mid C)$?

# Naïve Bayes Classifier

- Assume independence among attributes $A_i$ when class is given:
  - $P(A_1, A_2, \ldots, A_n \mid C) = P(A_1 \mid C_j) \, P(A_2 \mid C_j) \ldots P(A_n \mid C_j)$

  - Can estimate $P(A_i \mid C_j)$ for all $A_i$ and $C_j$.

  - New point is classified to $C_j$ if $P(C_j) \, \Pi \, P(A_i \mid C_j)$ is maximal.

# How to Estimate Probabilities from Data?

| Tid | Refund | Marital Status | Taxable Income | Evade |
|-----|--------|----------------|----------------|-------|
| 1 | Yes | Single | 125K | No |
| 2 | No | Married | 100K | No |
| 3 | No | Single | 70K | No |
| 4 | Yes | Married | 120K | No |
| 5 | No | Divorced | 95K | Yes |
| 6 | No | Married | 60K | No |
| 7 | Yes | Divorced | 220K | No |
| 8 | No | Single | 85K | Yes |
| 9 | No | Married | 75K | No |
| 10 | No | Single | 90K | Yes |

- Class: $P(C) = N_c/N$
  - e.g., $P(No) = 7/10$, $P(Yes) = 3/10$

- For discrete attributes:
  $$P(A_i \mid C_k) = |A_{ik}|/ N_{c_k}$$
  - where $|A_{ik}|$ is number of instances having attribute $A_i$ and belongs to class $C_k$
  - Examples:

  $P(Status=Married|No) = 4/7$
  $P(Refund=Yes|Yes)=0$

---

# Example of Naïve Bayes Classifier

| Name | Give Birth | Can Fly | Live in Water | Have Legs | Class |
|------|-----------|---------|---------------|-----------|-------|
| human | yes | no | no | yes | mammals |
| python | no | no | no | no | non-mammals |
| salmon | no | no | yes | no | non-mammals |
| whale | yes | no | yes | no | mammals |
| frog | no | no | sometimes | yes | non-mammals |
| komodo | no | no | no | yes | non-mammals |
| bat | yes | yes | no | yes | mammals |
| pigeon | no | yes | no | yes | non-mammals |
| cat | yes | no | no | yes | mammals |
| leopard shark | yes | no | yes | no | non-mammals |
| turtle | no | no | sometimes | yes | non-mammals |
| penguin | no | no | sometimes | yes | non-mammals |
| porcupine | yes | no | no | yes | mammals |
| eel | no | no | yes | no | non-mammals |
| salamander | no | no | sometimes | yes | non-mammals |
| gila monster | no | no | no | yes | non-mammals |
| platypus | no | no | no | yes | mammals |
| owl | no | yes | no | yes | non-mammals |
| dolphin | yes | no | yes | no | mammals |
| eagle | no | yes | no | yes | non-mammals |

| Give Birth | Can Fly | Live in Water | Have Legs | Class |
|-----------|---------|---------------|-----------|-------|
| yes | no | yes | no | ? |

A: attributes

M: mammals

N: non-mammals

$$P(A \mid M) = \frac{6}{7} \times \frac{6}{7} \times \frac{2}{7} \times \frac{2}{7} = 0.06$$

$$P(A \mid N) = \frac{1}{13} \times \frac{10}{13} \times \frac{3}{13} \times \frac{4}{13} = 0.0042$$

$$P(A \mid M)P(M) = 0.06 \times \frac{7}{20} = 0.021$$

$$P(A \mid N)P(N) = 0.004 \times \frac{13}{20} = 0.0027$$

P(A|M)P(M) >
P(A|N)P(N)

=> Mammals

# Naïve Bayes Classifier

l If one of the conditional probability is zero, then the entire expression becomes zero

l Probability estimation:

$$\text{Original}: P(A_i \mid C) = \frac{N_{ic}}{N_c}$$

$$\text{Laplace}: P(A_i \mid C) = \frac{N_{ic} + 1}{N_c + c}$$

$$\text{m-estimate}: P(A_i \mid C) = \frac{N_{ic} + mp}{N_c + m}$$
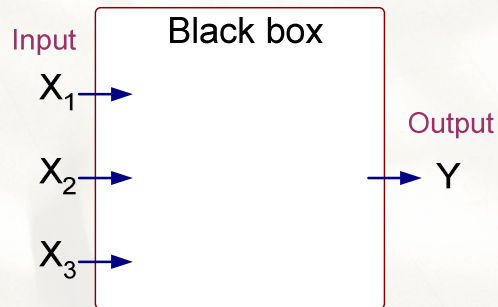
c: number of classes

p: prior probability

m: parameter

---

# Naïve Bayes (Summary)

• Robust to isolated noise points

• Handle missing values by ignoring the instance during probability estimate calculations

• Robust to irrelevant attributes

• Independence assumption may not hold for some attributes
  – Use other techniques such as Bayesian Belief Networks (BBN)
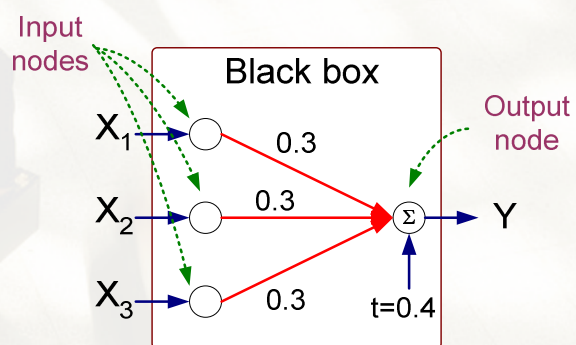
# Artificial Neural Networks (ANN)

| X₁ | X₂ | X₃ | Y |
|----|----|----|---|
| 1  | 0  | 0  | 0 |
| 1  | 0  | 1  | 1 |
| 1  | 1  | 0  | 1 |
| 1  | 1  | 1  | 1 |
| 0  | 0  | 1  | 0 |
| 0  | 1  | 0  | 0 |
| 0  | 1  | 1  | 1 |
| 0  | 0  | 0  | 0 |

Input

Black box

$X_1$

Output

$X_2$ → Y

$X_3$

Output Y is 1 if at least two of the three inputs are equal to 1.

---

# Artificial Neural Networks (ANN)

| X₁ | X₂ | X₃ | Y |
|----|----|----|---|
| 1  | 0  | 0  | 0 |
| 1  | 0  | 1  | 1 |
| 1  | 1  | 0  | 1 |
| 1  | 1  | 1  | 1 |
| 0  | 0  | 1  | 0 |
| 0  | 1  | 0  | 0 |
| 0  | 1  | 1  | 1 |
| 0  | 0  | 0  | 0 |

Input nodes

Black box

Output node

$X_1$ ◯  0.3

$X_2$ ◯  0.3  Σ → Y
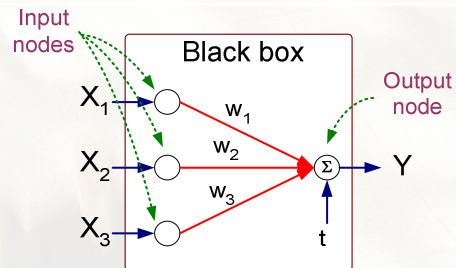
$X_3$ ◯  0.3   t=0.4

$$Y = I(0.3X_1 + 0.3X_2 + 0.3X_3 - 0.4 > 0)$$

$$\text{where } I(z) = \begin{cases} 1 & \text{if } z \text{ is true} \\ 0 & \text{otherwise} \end{cases}$$

# Artificial Neural Networks (ANN)

- Model is an assembly of inter-connected nodes and weighted links

- Output node sums up each of its input value according to the weights of its links
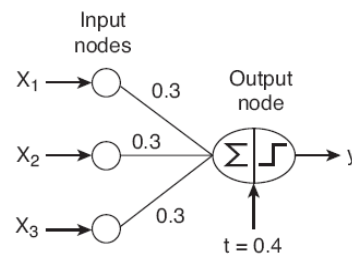
- Compare output node against some threshold t

Input nodes

Black box

Output node

$X_1$ — $w_1$

$X_2$ — $w_2$ — $\Sigma$ — Y

$X_3$ — $w_3$

t

Perceptron Model

$$Y = I(\sum_i w_i X_i - t) \quad \text{or}$$

$$Y = sign(\sum_i w_i X_i - t)$$

---

# Example of perceptron

| $X_1$ | $X_2$ | $X_3$ | y |
|-------|-------|-------|-----|
| 1 | 0 | 0 | −1 |
| 1 | 0 | 1 | 1 |
| 1 | 1 | 0 | 1 |
| 1 | 1 | 1 | 1 |
| 0 | 0 | 1 | −1 |
| 0 | 1 | 0 | −1 |
| 0 | 1 | 1 | 1 |
| 0 | 0 | 0 | −1 |

(a) Data set.

Input nodes

$X_1$ — 0.3

Output node

$X_2$ — 0.3 — $\Sigma$ — y

$X_3$ — 0.3

t = 0.4

(b) Perceptron.

**Figure 5.14.** Modeling a boolean function using a perceptron.

# Example of perceptron



**Figure 5.15.** Perceptron decision boundary for the data given in Figure 5.14.

# Example of perceptron

| X₁ | X₂ | y |
|----|----|----|
| 0 | 0 | −1 |
| 1 | 0 | 1 |
| 0 | 1 | 1 |
| 1 | 1 | −1 |



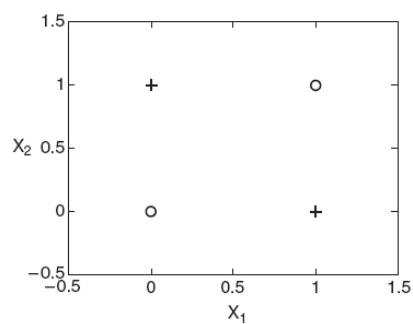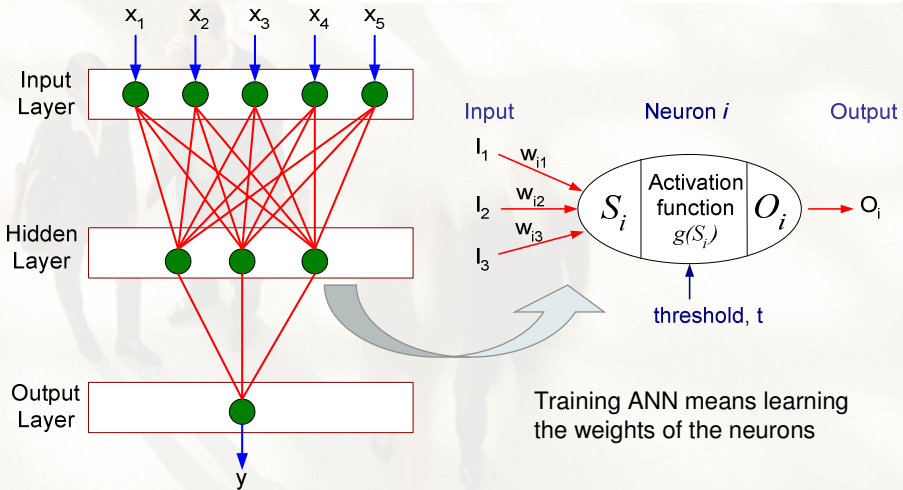**Figure 5.16.** XOR classification problem. No linear hyperplane can separate the two classes.

# General Structure of ANN



$x_1$ $x_2$ $x_3$ $x_4$ $x_5$

Input Layer

Hidden Layer

Output Layer

y

Input    Neuron $i$    Output

$I_1$ $w_{i1}$

$I_2$ $w_{i2}$   $S_i$   Activation function $g(S_i)$   $O_i$   $O_i$

$I_3$ $w_{i3}$

threshold, t

Training ANN means learning the weights of the neurons

# Example of multi-layered ANN



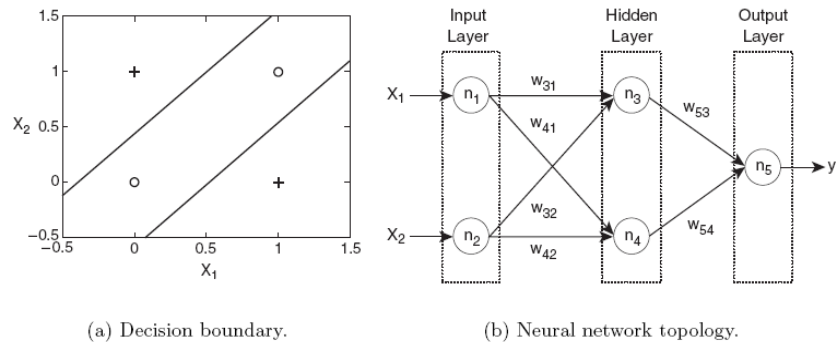(a) Decision boundary.

(b) Neural network topology.

**Figure 5.19.** A two-layer, feed-forward neural network for the XOR problem.

# Algorithm for learning ANN

l Initialize the weights ($w_0$, $w_1$, …, $w_k$)

l Adjust the weights in such a way that the output of ANN is consistent with class labels of training examples
  - Objective function: $E = \sum_i \left[ Y_i - f(w_i, X_i) \right]^2$

  - Find the weights $w_i$'s that minimize the above objective function
    u e.g., backpropagation algorithm

# Neural Networks: Summary

- Pros
  - Accurate
  - Wide range of applications
- Cons
  - Difficult interpretation
  - Tends to 'overfit' the data
  - Extensive amount of training time
  - A lot of data preparation

# Collective comparison

|  | Train time | Run Time | Noise Tolerance | Can Use Prior Knowledge | Accuracy on Customer Modelling | Under-standable |
|---|---|---|---|---|---|---|
| Decision Trees | fast | fast | poor | no | medium | medium |
| Bayesian | slow | fast | good | yes | good | good |
| Neural Networks | slow | fast | good | no | good | poor |

# Evaluation



| Tid | Attrib1 | Attrib2 | Attrib3 | Class |
|---|---|---|---|---|
| 1 | Yes | Large | 125K | No |
| 2 | No | Medium | 100K | No |
| 3 | No | Small | 70K | No |
| 4 | Yes | Medium | 120K | No |
| 5 | No | Large | 95K | Yes |
| 6 | No | Medium | 60K | No |
| 7 | Yes | Large | 220K | No |
| 8 | No | Small | 85K | Yes |
| 9 | No | Medium | 75K | No |
| 10 | No | Small | 90K | Yes |

Training Set

| Tid | Attrib1 | Attrib2 | Attrib3 | Class |
|---|---|---|---|---|
| 11 | No | Small | 55K | ? |
| 12 | Yes | Medium | 80K | ? |
| 13 | Yes | Large | 110K | ? |
| 14 | No | Small | 95K | ? |
| 15 | No | Large | 67K | ? |

Test Set

Tree Induction algorithm

Induction

Learn Model

Apply Model

Deduction

Model

Decision Tree

# Test Sample Estimate

- Divide D into $D_1$ and $D_2$
- Use $D_1$ to construct the classifier d
- Then use resubstitution estimate $R(d, D_2)$ to calculate the estimated misclassification error of d
- Unbiased and efficient, but removes $D_2$ from training dataset D

# Cross-Validation



—Break up data into subsets of the same size

—

—

—Hold aside one subsets for testing and use the rest for training

Test

—Repeat

# Metrics for Performance Evaluation

l Focus on the predictive capability of a model
  – Rather than how fast it takes to classify or build models, scalability, etc.
l Confusion Matrix:

| | PREDICTED CLASS | | |
|---|---|---|---|
| | | Class=Yes | Class=No |
| **ACTUAL CLASS** | Class=Yes | a | b |
| | Class=No | c | d |

a: TP (true positive)

b: FN (false negative)

c: FP (false positive)

d: TN (true negative)

# Metrics for Performance Evaluation...

| | PREDICTED CLASS | | |
|---|---|---|---|
| | | Class=Yes | Class=No |
| **ACTUAL CLASS** | Class=Yes | a (TP) | b (FN) |
| | Class=No | c (FP) | d (TN) |

l Most widely-used metric:

$$\text{Accuracy} = \frac{a+d}{a+b+c+d} = \frac{TP+TN}{TP+TN+FP+FN}$$

# Limitation of Accuracy

- Consider a 2-class problem
  - Number of Class 0 examples = 9990
  - Number of Class 1 examples = 10

- If model predicts everything to be class 0, accuracy is 9990/10000 = 99.9 %
  - Accuracy is misleading because model does not detect any class 1 example

# Cost Matrix

| | PREDICTED CLASS | | |
|---|---|---|---|
| | C(i\|j) | **Class=Yes** | **Class=No** |
| **ACTUAL CLASS** | **Class=Yes** | C(Yes\|Yes) | C(No\|Yes) |
| | **Class=No** | C(Yes\|No) | C(No\|No) |

C(i|j): Cost of misclassifying class j example as class i

# Example

| Cost Matrix | PREDICTED CLASS | | |
|---|---|---|---|
| | C(i\|j) | + | - |
| ACTUAL CLASS | + | -1 | 100 |
| | - | 1 | 0 |

| Model $M_1$ | PREDICTED CLASS | | |
|---|---|---|---|
| | | + | - |
| ACTUAL CLASS | + | 150 | 40 |
| | - | 60 | 250 |

| Model $M_2$ | PREDICTED CLASS | | |
|---|---|---|---|
| | | + | - |
| ACTUAL CLASS | + | 250 | 45 |
| | - | 5 | 200 |

Accuracy = 80%

Cost = 3910

Accuracy = 90%

Cost = 4255

49

# Underfitting and Overfitting



Underfitting: when model is too simple, both training and test errors are large

# How to Address Overfitting...

l **Post-pruning**

- Grow decision tree to its entirety
- Trim the nodes of the decision tree in a bottom-up fashion
- If generalization error improves after trimming, replace sub-tree by a leaf node.
- Class label of leaf node is determined from majority class of instances in the sub-tree
- Can use MDL for post-pruning