

Classification Algo

AVAILABLE AT:

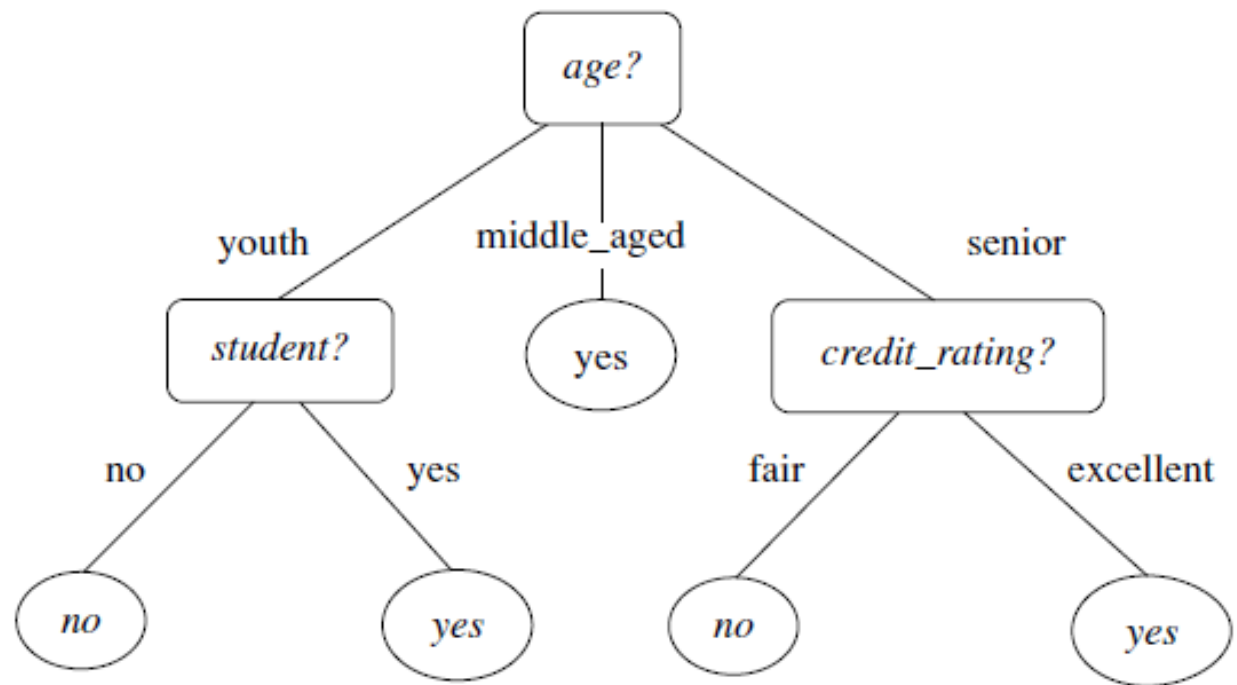
Onebyzero Edu - Organized Learning, Smooth Career

The Comprehensive Academic Study Platform for University Students in Bangladesh (www.onebyzeroedu.com)

Decision Tree Learning

- Decision tree induction is the learning of decision trees from class-labeled training tuples.
- A decision tree is a flowchart-like tree structure, where each internal node (nonleaf node) denotes a test on an attribute, each branch represents an outcome of the test, and each leaf node (or *terminal node*) holds a class label.
- *The topmost node in* a tree is the root node.

A decision tree for the concept *buys_computer*, indicating whether an *AllElectronics* customer is likely to purchase a computer. Each internal (nonleaf) node represents a test on an attribute. Each leaf node represents a class (either *buys_computer* = yes or *buys_computer* = no).



Example: Good versus Evil

- **problem**: identify people as good or bad from their appearance

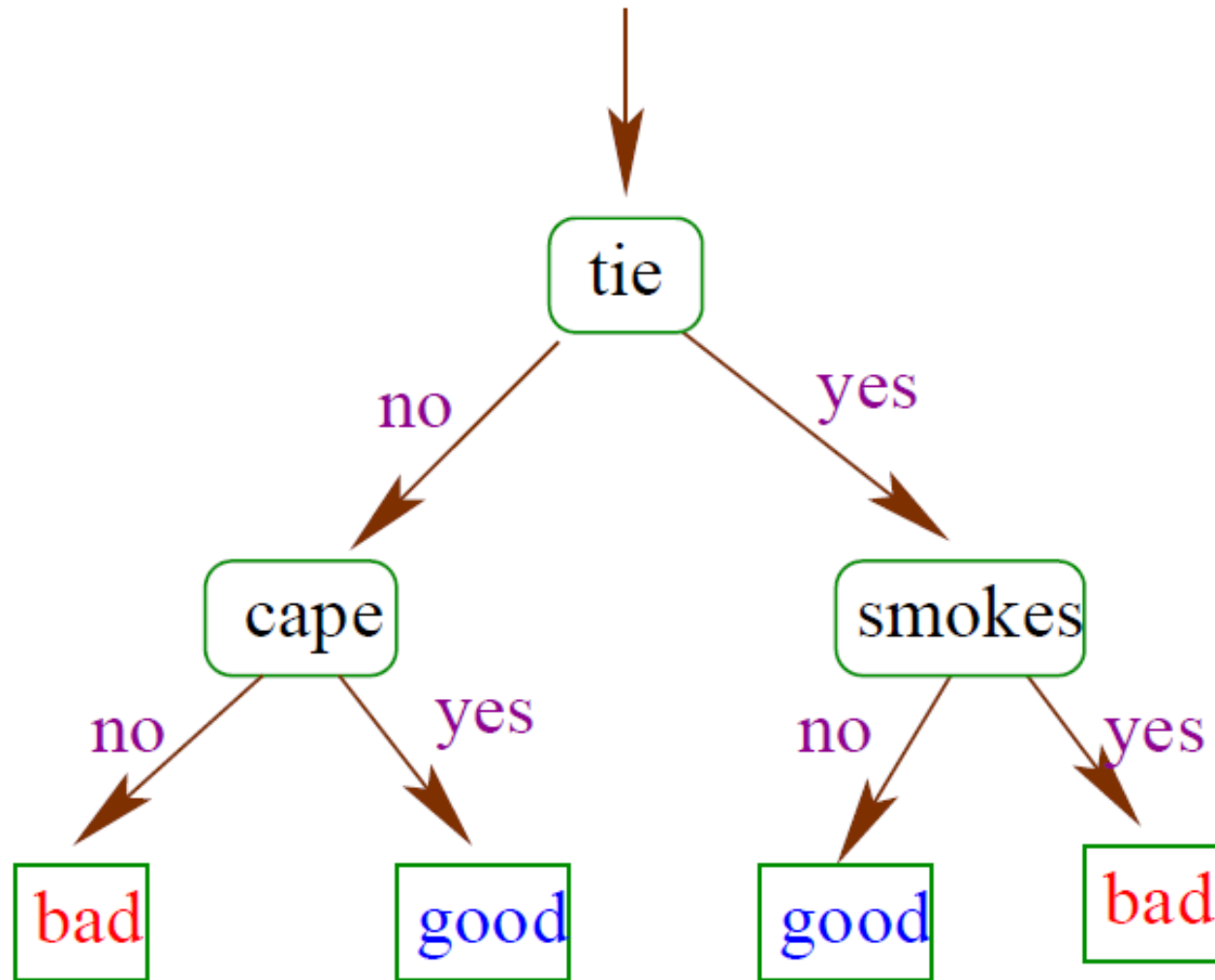
	sex	mask	cape	tie	ears	smokes	class
training data							
batman	male	yes	yes	no	yes	no	Good
robin	male	yes	yes	no	no	no	Good
alfred	male	no	no	yes	no	no	Good
penguin	male	no	no	yes	no	yes	Bad
catwoman	female	yes	no	no	yes	no	Bad
joker	male	no	no	no	no	no	Bad
test data							
batgirl	female	yes	yes	no	yes	no	??
riddler	male	yes	no	no	no	no	??

AVAILABLE AT:

Onebyzero Edu - Organized Learning, Smooth Career

The Comprehensive Academic Study Platform for University Students in Bangladesh (www.onebyzeroedu.com)

A Decision Tree Classifier

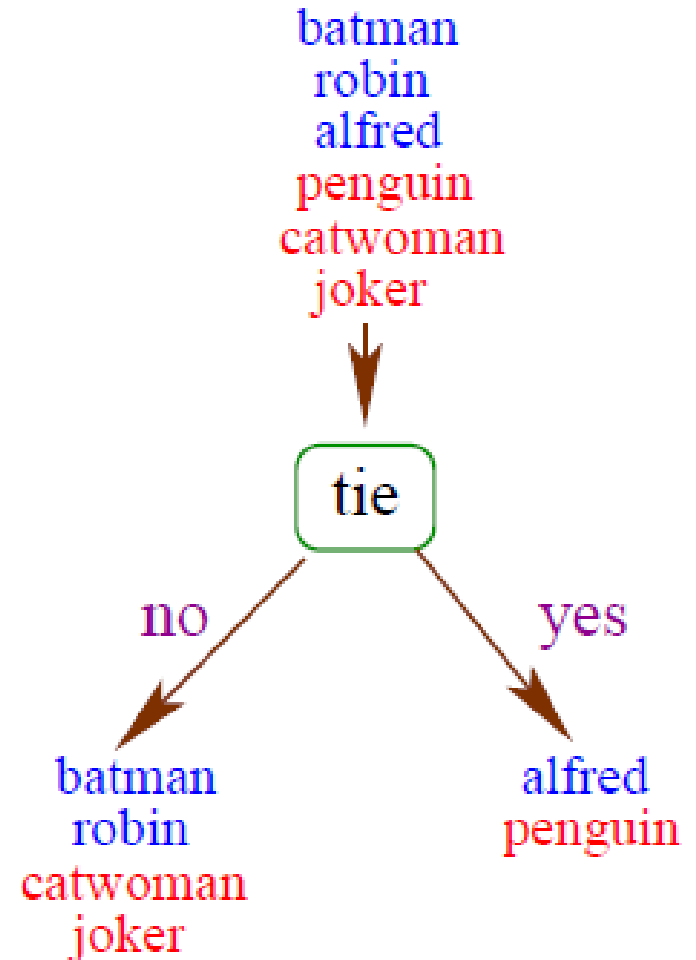


AVAILABLE AT:

Onebyzero Edu - Organized Learning, Smooth Career
The Comprehensive Academic Study Platform for University Students in Bangladesh (www.onebyzeroedu.com)

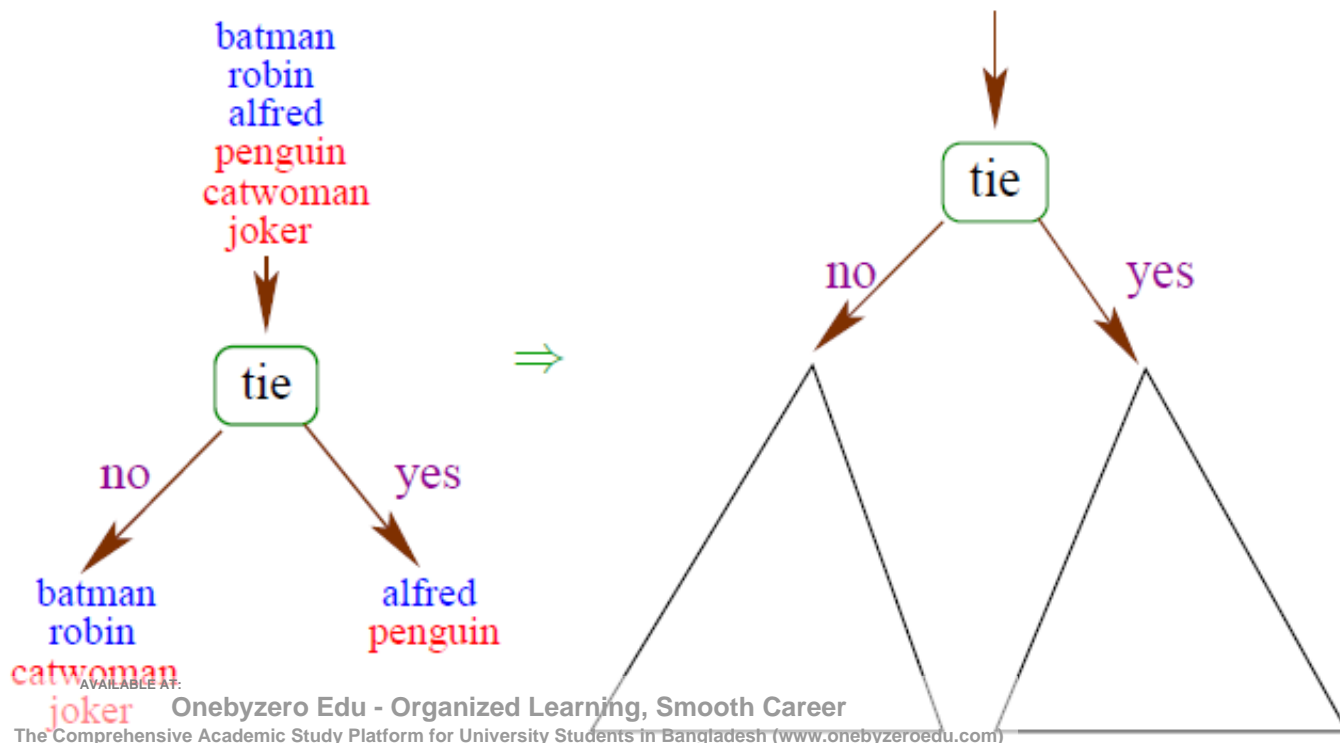
How to Build Decision Trees

- choose rule to split on
- divide data using splitting rule into disjoint subsets



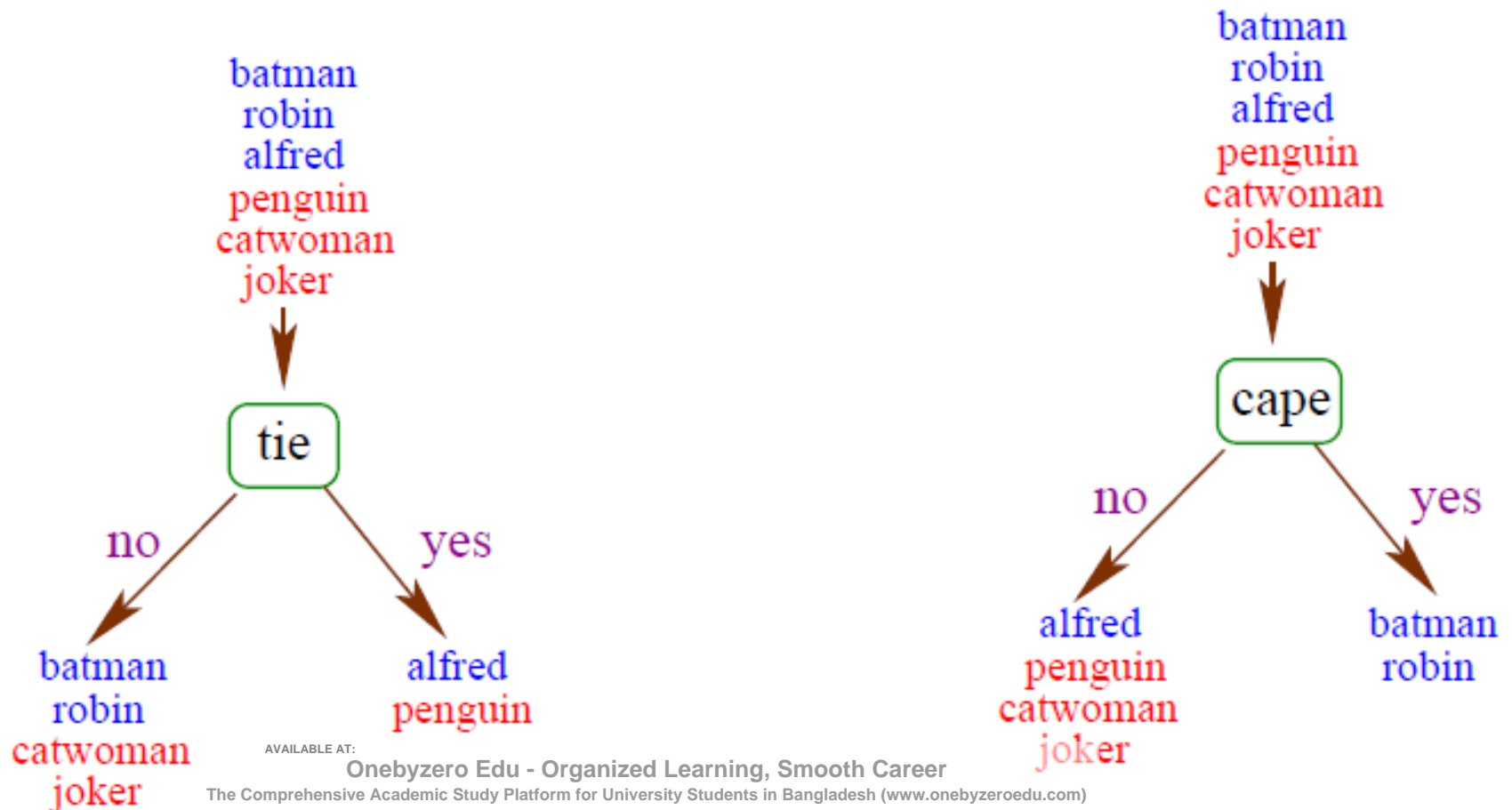
How to Build Decision Trees

- choose rule to split on
- divide data using splitting rule into disjoint subsets
- repeat recursively for each subset
- stop when leaves are (almost) “pure”



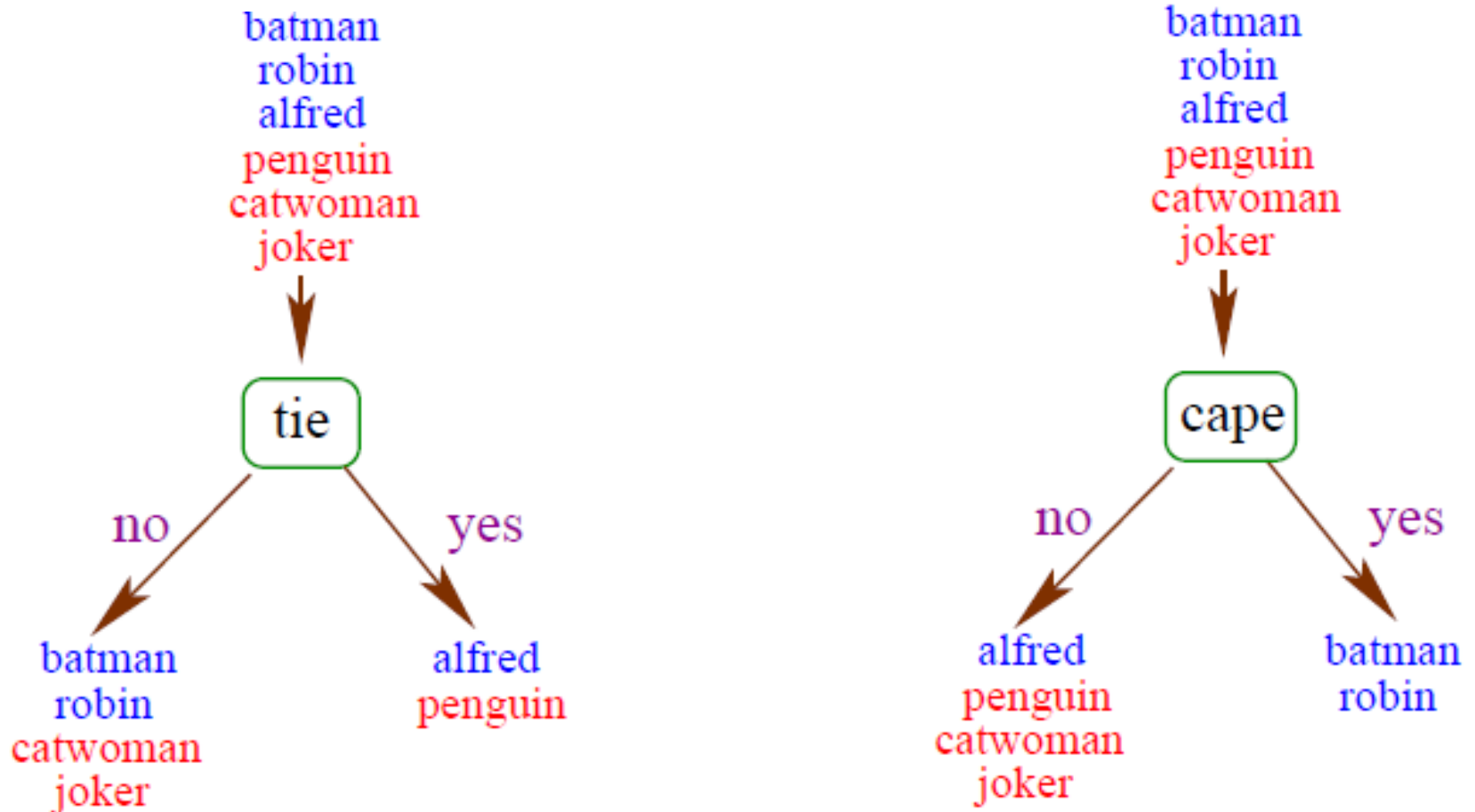
How to Choose the Splitting Rule

- key problem: choosing best rule to split on:



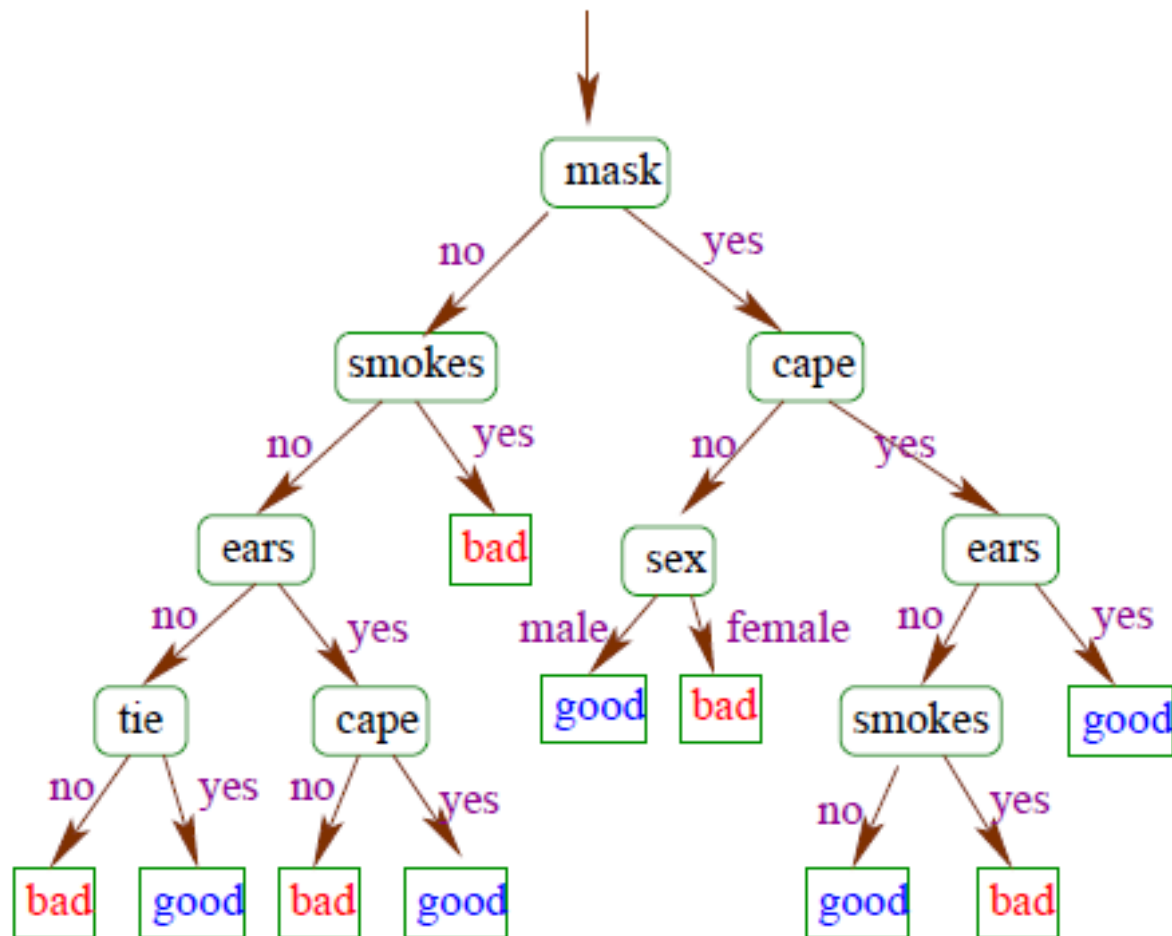
How to Choose the Splitting Rule

key problem: choosing best rule to split on:



idea: choose rule that leads to greatest increase in “purity”

A Possible Classifier



AVAILABLE AT:

Onebyzero Edu - Organized Learning, Smooth Career
The Comprehensive Academic Study Platform for University Students in Bangladesh (www.onebyzeroedu.com)

How to Measure Purity

- Information gain
- Gini Index
- Gain Ratio

Information Gain

- Expected information (entropy) needed to classify a tuple in D

$$Info(D) = -\sum_{i=1}^m p_i \log_2(p_i)$$

- Information needed (after using A to split D into v partitions) to classify D:

$$Info_A(D) = \sum_{j=1}^v \frac{|D_j|}{|D|} \times Info(D_j)$$

- Information gained by branching on attribute A

$$Gain(A) = Info(D) - Info_A(D)$$

$$Gain(A) = Info(D) - Info_A(D)$$

- **Gain(A)** tells us how much would be gained by branching on A
- The attribute A with the highest **information gain, Gain (A)**, is chosen as the splitting attribute at node N.

An Illustrative Example

Class-Labeled Training Tuples from the *AllElectronics* Customer Database

<i>RID</i>	<i>age</i>	<i>income</i>	<i>student</i>	<i>credit_rating</i>	<i>Class: buys_computer</i>
1	youth	high	no	fair	no
2	youth	high	no	excellent	no
3	middle_aged	high	no	fair	yes
4	senior	medium	no	fair	yes
5	senior	low	yes	fair	yes
6	senior	low	yes	excellent	no
7	middle_aged	low	yes	excellent	yes
8	youth	medium	no	fair	no
9	youth	low	yes	fair	yes
10	senior	medium	yes	fair	yes
11	youth	medium	yes	excellent	yes
12	middle_aged	medium	no	excellent	yes
13	middle_aged	high	yes	fair	yes
14	senior	medium	no	excellent	no

Quinlan [Qui86].

SATILABEAT:

Onebyzero Edu - Organized Learning, Smooth Career

The Comprehensive Academic Study Platform for University Students in Bangladesh (www.onebyzeroedu.com)

How to choose best splitting criterion?

- The class label attribute, **buys computer**, has two distinct values (namely, {**yes**, **no**});
- two distinct classes (i.e., $m = 2$).
- class $C_1 = \text{yes}$
- class $C_2 = \text{no}$
- 09 tuples of class = yes
- 05 tuples of class = no
- A (root) node N is created for the tuples in D.
- To find the splitting criterion for these tuples, **compute the information gain** of each attribute.

- Expected information needed to classify a tuple in D :

$$Info(D) = -\sum_{i=1}^m p_i \log_2(p_i)$$

$$Info(D) = -\frac{9}{14} \log_2\left(\frac{9}{14}\right) - \frac{5}{14} \log_2\left(\frac{5}{14}\right) = 0.940 \text{ bits.}$$

- Next, we need to compute the expected information requirement for each attribute.
- Start with attribute: *age*
- age category “youth,”: yes = 02 tuples & no = 03 tuples.
- category “middle aged,”: yes = 04 tuples & no = 0 tuples.
- category “senior,”: *yes = 03 tuples & no = 02 tuples.*

- The expected information needed to classify a tuple in D if the tuples are partitioned according to **age**:

$$Info_A(D) = \sum_{j=1}^v \frac{|D_j|}{|D|} \times Info(D_j)$$

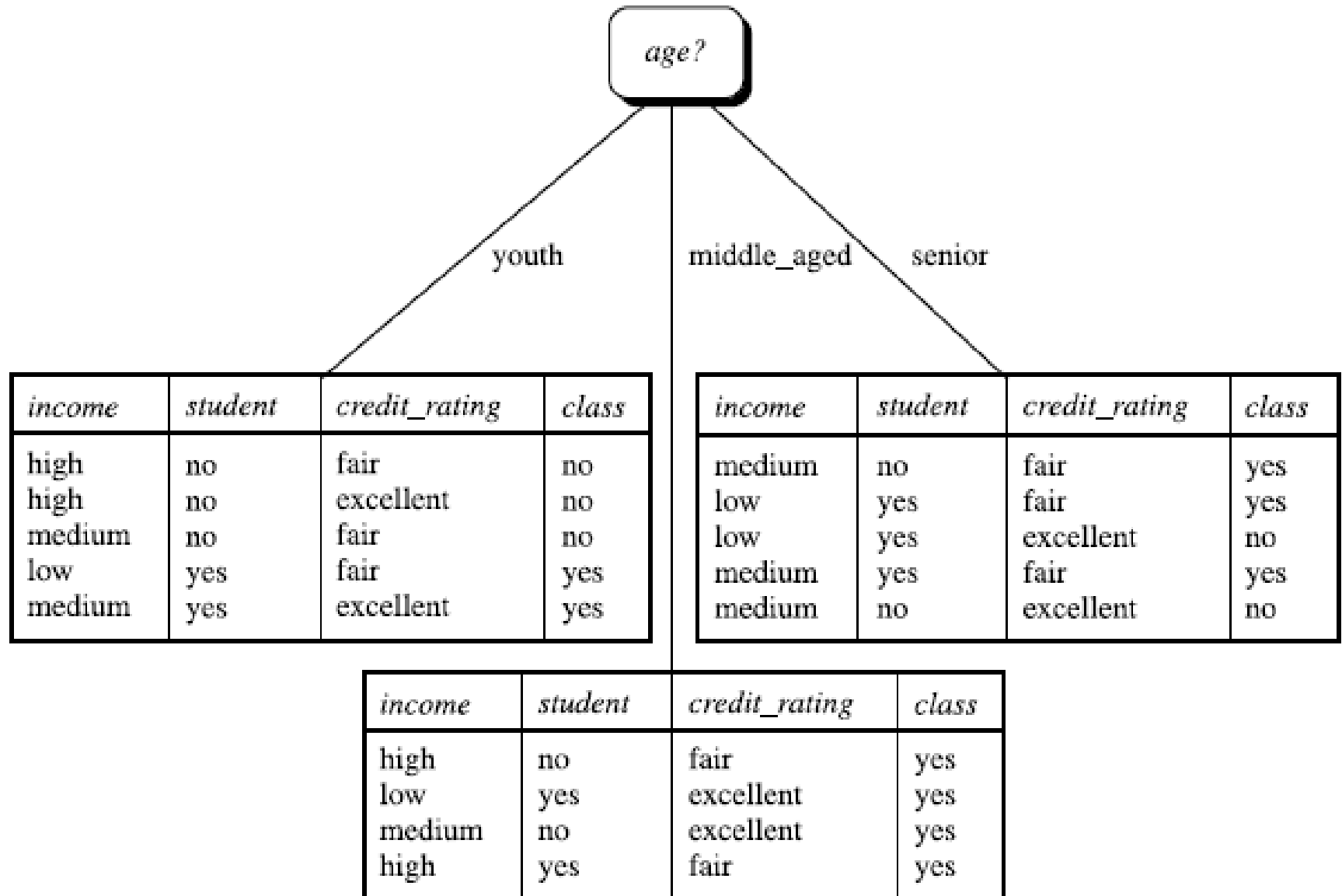
$$\begin{aligned}
 Info_{age}(D) &= \frac{5}{14} \times \left(-\frac{2}{5} \log_2 \frac{2}{5} - \frac{3}{5} \log_2 \frac{3}{5} \right) + \frac{4}{14} \times \left(-\frac{4}{4} \log_2 \frac{4}{4} \right) \\
 &\quad + \frac{5}{14} \times \left(-\frac{3}{5} \log_2 \frac{3}{5} - \frac{2}{5} \log_2 \frac{2}{5} \right) \\
 &= 0.694 \text{ bits.}
 \end{aligned}$$

- Hence, the gain in information from such a partitioning would be

$$\text{Gain}(\text{age}) = \text{Info}(D) - \text{Info}_{\text{age}}(D) = 0.940 - 0.694 = 0.246 \text{ bits.}$$

Similarly, we can compute $\text{Gain}(\text{income}) = 0.029 \text{ bits}$, $\text{Gain}(\text{student}) = 0.151 \text{ bits}$, and $\text{Gain}(\text{credit_rating}) = 0.048 \text{ bits}$.

Because age has the highest information gain among the attributes, it is selected as the splitting attribute.



AVAILABLE AT:

Onebyzero Edu - Organized Learning, Smooth Career

The Comprehensive Academic Study Platform for University Students in Bangladesh (www.onebyzeroedu.com)

Gain Ratio

- The information gain measure is **biased** toward tests with many outcomes.
- That is, it prefers to select attributes having a large number of values.
- For example, consider an attribute that acts as a unique identifier such as *product_ID*.
- *A split on product_ID would result in a large number of partitions (as many as there are values), each one containing just one tuple.*
- Because each partition is pure, the information required to classify data set *D* based on this partitioning would be **$Info_{product_ID}(D) = 0$** .
- *information* gained by partitioning on this attribute is maximal.
- Such a partitioning is **useless** for classification.

- C4.5, a successor of ID3, uses an extension to information gain known as **gain ratio**, which attempts to overcome this bias.
- It applies a kind of normalization to information gain using a “**split information**” value defined analogously with $Info(D)$:

$$SplitInfo_A(D) = - \sum_{j=1}^v \frac{|D_j|}{|D|} \times \log_2 \left(\frac{|D_j|}{|D|} \right)$$

This value represents the potential information generated by splitting the training data set, D , into v partitions, corresponding to the v outcomes of a test on attribute A .

- The gain ratio is defined as

$$\text{GainRatio}(A) = \frac{\text{Gain}(A)}{\text{SplitInfo}_A(D)}.$$

- The attribute with the **maximum gain ratio** is selected as the splitting attribute.
- Note, however, that as the split information approaches 0, the ratio becomes unstable.

Example: Computation of gain ratio for the attribute *income*

- A test on *income* splits the data of into 03 partitions:

✓ *Low = 04 tuples*

✓ *Medium = 06 tuples*

✓ *High = 04 tuples*

Class-Labeled Training Tuples from the *AllElectronics* Custom

RID	age	income	student	credit_rating	Class:
1	youth	high	no	fair	no
2	youth	high	no	excellent	no
3	middle_aged	high	no	fair	yes
4	senior	medium	no	fair	yes
5	senior	low	yes	fair	yes
6	senior	low	yes	excellent	no
7	middle_aged	low	yes	excellent	yes
8	youth	medium	no	fair	no
9	youth	low	yes	fair	yes
10	senior	medium	yes	fair	yes
11	youth	medium	yes	excellent	yes
12	middle_aged	medium	no	excellent	yes
13	middle_aged	high	yes	fair	yes
14	senior	medium	no	excellent	no

$$SplitInfo_A(D) = - \sum_{j=1}^v \frac{|D_j|}{|D|} \times \log_2 \left(\frac{|D_j|}{|D|} \right)$$

$$SplitInfo_{income}(D) = -\frac{4}{14} \times \log_2 \left(\frac{4}{14} \right) - \frac{6}{14} \times \log_2 \left(\frac{6}{14} \right) - \frac{4}{14} \times \log_2 \left(\frac{4}{14} \right)$$

$$= 1.557.$$

$$Gain(age) = Info(D) - Info_{age}(D)$$

$$GainRatio(A) = \frac{Gain(A)}{SplitInfo_A(D)}.$$

Gain(income) = 0.029
[See previous example]

$$= \frac{0.029}{1.557}$$

$$= 0.019$$

AVAILABLE AT:

Onebyzero Edu - Organized Learning, Smooth Career
 The Comprehensive Academic Study Platform for University Students in Bangladesh (www.onebyzeroedu.com)

Gini Index

- The Gini index is used in CART (Classification & Regression Tree).
- Gini index measures the impurity of D , *a data partition or set of training tuples*

$$Gini(D) = 1 - \sum_{i=1}^m p_i^2,$$

- where p_i is the probability that a tuple in D belongs to class C_i and is estimated by $|C_{i,D}|/|D|$.
- The sum is computed over m classes.

- The Gini index considers **a binary split** for each attribute.
- To determine the best binary split on A , we *examine all the possible subsets* that can be formed using known values of A .
- Each subset, S_A , *can be considered as a binary test for attribute A of the form “ $A \in S_A$?”*
- *Given a tuple, this test is satisfied if the value of A for the tuple is among the values listed in S_A .*
- *If A has v possible values, then there are 2^v possible subsets.*

- For example, if *income* has three possible values: *low, medium, high*,
- possible subsets are *{low, medium, high}*, *{low, medium}*, *{low, high}*, *{medium, high}*, *{low}*, *{medium}*, *{high}*, and *{}*.
- Exclude the power set, *{low, medium, high}*, and the empty set from consideration since, conceptually, they do not represent a split.
- Therefore, there are $2^v - 2$ possible ways to form two partitions of the data, *D*, based on a binary split on *A*.

- When considering a binary split, we compute a weighted sum of the impurity of each resulting partition.
- For example, if a binary split on A partitions D into D_1 and D_2 , the Gini index of D given that partitioning is

$$Gini_A(D) = \frac{|D_1|}{|D|} Gini(D_1) + \frac{|D_2|}{|D|} Gini(D_2).$$

- For each attribute, each of the possible binary splits is considered.
- For a discrete-valued attribute, the subset that gives the **minimum Gini index for that attribute is selected as its splitting subset.**

AVAILABLE AT:

Onebyzero Edu - Organized Learning, Smooth Career

The Comprehensive Academic Study Platform for University Students in Bangladesh (www.onebyzeroedu.com)

- The reduction in impurity that would be incurred by a binary split on a discrete- or continuous-valued attribute A

$$\Delta Gini(A) = Gini(D) - Gini_A(D).$$

- The attribute that maximizes the reduction in impurity (or, equivalently, has the minimum Gini index) is selected as the splitting attribute.
- This attribute & either its splitting subset or split-point together form the splitting criterion.

Induction of a decision tree using the Gini index

- $C_1 = \text{buys computer} = \text{yes} = 09$
- $C_2 = \text{buys computer} = \text{no} = 05$

Class-Labeled Training Tuples from the *AllElectronics* Custom

<i>RID</i>	<i>age</i>	<i>income</i>	<i>student</i>	<i>credit_rating</i>	<i>Class:</i>
1	youth	high	no	fair	no
2	youth	high	no	excellent	no
3	middle_aged	high	no	fair	yes
4	senior	medium	no	fair	yes
5	senior	low	yes	fair	yes
6	senior	low	yes	excellent	no
7	middle_aged	low	yes	excellent	yes
8	youth	medium	no	fair	no
9	youth	low	yes	fair	yes
10	senior	medium	yes	fair	yes
11	youth	medium	yes	excellent	yes
12	middle_aged	medium	no	excellent	yes
13	middle_aged	high	yes	fair	yes
14	senior	medium	no	excellent	no

- Gini index to compute the impurity of D

$$Gini(D) = 1 - \sum_{i=1}^m p_i^2,$$

A (root) node N is created for the tuples in D .

$$Gini(D) = 1 - \left(\frac{9}{14}\right)^2 - \left(\frac{5}{14}\right)^2 = 0.459.$$

To find the splitting criterion for the tuples in D , we need to compute the *Gini index* for each attribute

- Let's start with the attribute *income* & consider each of the possible splitting subsets.
- Consider the subset $\{low, medium\}$.
- This would result in 10 tuples in partition D_1 satisfying the condition " $income \in \{low, medium\}$."
- The remaining 04 tuples of D would be assigned to partition D_2 .

<i>RID</i>	<i>age</i>	<i>income</i>	<i>student</i>	<i>credit_rating</i>	<i>Class:</i>
1	youth	high	no	fair	no
2	youth	high	no	excellent	no
3	middle_aged	high	no	fair	yes
4	senior	medium	no	fair	yes
5	senior	low	yes	fair	yes
6	senior	low	yes	excellent	no
7	middle_aged	low	yes	excellent	yes
8	youth	medium	no	fair	no
9	youth	low	yes	fair	yes
10	senior	medium	yes	fair	yes
11	youth	medium	yes	excellent	yes
12	middle_aged	medium	no	excellent	yes
13	middle_aged	high	yes	fair	yes
14	senior	medium	no	excellent	no

D1:**Yes = 07****No = 03****D2:****Yes=2****No=2**

$$Gini_{income \in \{low, medium\}}(D)$$

$$= \frac{10}{14} Gini(D_1) + \frac{4}{14} Gini(D_2)$$

$$= \frac{10}{14} \left(1 - \left(\frac{7}{10} \right)^2 - \left(\frac{3}{10} \right)^2 \right) + \frac{4}{14} \left(1 - \left(\frac{2}{4} \right)^2 - \left(\frac{2}{4} \right)^2 \right)$$

$$= 0.443$$

$$= Gini_{income \in \{high\}}(D).$$

AVAILABLE AT:

- Similarly, the Gini index values for splits on the remaining subsets
 - ✓ $0.458 = \{low, high\} + \{medium\}$
 - ✓ $0.450 = \{medium, high\} + \{low\}$
- Therefore, the best binary split for attribute **income** is on $\{low, medium\}$ or $\{high\} = 0.443$ because it minimizes the Gini index.
- Evaluating *age*, we obtain $\{youth, senior\}$ (or $\{middle_aged\}$) as the best split for *age* with a Gini index of 0.375

- the attributes *student* and *credit_rating* are both binary, with Gini index values of 0.367 & 0.429, respectively.
- The attribute *age* and splitting subset {youth, senior} therefore give the minimum Gini index overall, with a reduction in impurity of $0.459 - 0.357 = 0.102$.
- The binary split “ $age \in \{youth, senior\}$ ” results in the **maximum reduction in impurity** of the tuples in D and is returned as the splitting criterion.
- Node N is labeled with the criterion, two branches are grown from it, and the tuples are partitioned accordingly.

Decision Trees

best known:

- C4.5 (Quinlan)
- CART (Breiman, Friedman, Olshen & Stone)
- very fast to train and evaluate
- relatively easy to interpret
- but: accuracy often not state-of-the-art

Bayes Classification

AVAILABLE AT:

Onebyzero Edu - Organized Learning, Smooth Career

The Comprehensive Academic Study Platform for University Students in Bangladesh (www.onebyzeroedu.com)

What are Bayesian classifiers?"

- *Bayesian classifiers are statistical classifiers.*
- *They can predict class membership probabilities such as the probability that a given tuple belongs to a particular class*

Bayesian Classification: Why?

- A statistical classifier: performs *probabilistic prediction*, i.e., predicts class membership probabilities
- Foundation: Based on Bayes' Theorem.
- Performance: A simple Bayesian classifier, *naïve Bayesian classifier*, has comparable performance with decision tree & selected neural network classifiers
- Incremental: Each training example can incrementally increase/decrease the probability that a hypothesis is correct — prior knowledge can be combined with observed data
- Standard: Even when Bayesian methods are computationally intractable, they can provide a standard of optimal decision making against which other methods can be measured

Bayes' Theorem: Basics

- Total Probability Theorem:

$$P(B) = \sum_{i=1}^M P(B|A_i)P(A_i)$$

Bayes' Theorem: Basics

- Bayes' Theorem:
$$P(H | \mathbf{X}) = \frac{P(\mathbf{X} | H)P(H)}{P(\mathbf{X})} = P(\mathbf{X} | H) \times P(H) / P(\mathbf{X})$$
- Let \mathbf{X} be a data sample (“evidence”): class label is unknown
- Let H be a *hypothesis* that X belongs to class C
- Classification is to determine $P(H|\mathbf{X})$, (i.e., *posteriori probability*): the probability that the hypothesis holds given the observed data sample \mathbf{X}
- $P(H)$ (*prior probability*): the initial probability
 - E.g., \mathbf{X} will buy computer, regardless of age, income, ...
- $P(\mathbf{X})$: probability that sample data is observed
- $P(\mathbf{X}|H)$ (*likelihood*): the probability of observing the sample \mathbf{X} , given that the hypothesis holds
 - E.g., Given that \mathbf{X} will buy computer, the prob. that X is 31..40, medium income

Prediction Based on Bayes' Theorem

- Given training data \mathbf{X} , *posteriori probability of a hypothesis H* , $P(H|\mathbf{X})$, follows the Bayes' theorem

$$P(H|\mathbf{X}) = \frac{P(\mathbf{X}|H)P(H)}{P(\mathbf{X})} = P(\mathbf{X}|H) \times P(H) / P(\mathbf{X})$$

- Informally, this can be viewed as
posteriori = likelihood x prior/evidence
- Predicts \mathbf{X} belongs to C_i iff the probability $P(C_i|\mathbf{X})$ is the highest among all the $P(C_k|\mathbf{X})$ for all the k classes
- Practical difficulty:** It requires initial knowledge of many probabilities, involving significant computational cost

Classification Is to Derive the Maximum Posteriori

- Let D be a training set of tuples & their associated class labels, & each tuple is represented by an n -D attribute vector $\mathbf{X} = (x_1, x_2, \dots, x_n)$
- Suppose there are m classes C_1, C_2, \dots, C_m .
- Classification is to derive the maximum posteriori, i.e., the **maximal $P(C_i | \mathbf{X})$**
- This can be derived from Bayes' theorem

$$P(C_i | \mathbf{X}) = \frac{P(\mathbf{X} | C_i)P(C_i)}{P(\mathbf{X})}$$

- Since $P(\mathbf{X})$ is constant for all classes, only

$$P(C_i | \mathbf{X}) = P(\mathbf{X} | C_i)P(C_i)$$

needs to be maximized

Naïve Bayes Classifier

- **When to use**

- Moderate or large training set available
- Attributes that describes instances are conditionally independent given classifier

- **Applications**

- Diagnosis
- Classifying text documents

Naïve Bayes Classifier

- **A simplified assumption:** attributes are conditionally independent (i.e., no dependence relation between attributes):

$$P(\mathbf{X} | C_i) = \prod_{k=1}^n P(x_k | C_i) = P(x_1 | C_i) \times P(x_2 | C_i) \times \dots \times P(x_n | C_i)$$

- This greatly reduces the computation cost: Only counts the class distribution
- If A_k is categorical, $P(x_k | C_i)$ is the # of tuples in C_i having value x_k for A_k divided by $|C_{i,D}|$ (# of tuples of C_i in D)
- If A_k is continuous-valued, $P(x_k | C_i)$ is usually computed based on Gaussian distribution with a mean μ & standard deviation σ

& $P(x_k | C_i)$ is

$$g(x, \mu, \sigma) = \frac{1}{\sqrt{2\pi}\sigma} e^{-\frac{(x-\mu)^2}{2\sigma^2}}$$

$$P(\mathbf{X} | C_i) = g(x_k, \mu_{C_i}, \sigma_{C_i})$$

AVAILABLE AT:

Predicting a class label using naïve Bayesian classification

Class-Labeled Training Tuples from the *AllElectronics* Customer Database

<i>RID</i>	<i>age</i>	<i>income</i>	<i>student</i>	<i>credit_rating</i>	<i>Class: buys_computer</i>
1	youth	high	no	fair	no
2	youth	high	no	excellent	no
3	middle_aged	high	no	fair	yes
4	senior	medium	no	fair	yes
5	senior	low	yes	fair	yes
6	senior	low	yes	excellent	no
7	middle_aged	low	yes	excellent	yes
8	youth	medium	no	fair	no
9	youth	low	yes	fair	yes
10	senior	medium	yes	fair	yes
11	youth	medium	yes	excellent	yes
12	middle_aged	medium	no	excellent	yes
13	middle_aged	high	yes	fair	yes
14	senior	medium	no	excellent	no

SOURCE: KDD-96

- The data tuples are described by the attributes:
- *age, income, student, & credit_rating*

Class:

C1: *buys_computer* = 'yes'

C2: *buys_computer* = 'no'

Class-Labeled Training Tuples from the *AllElectronics* Customer

<i>RID</i>	<i>age</i>	<i>income</i>	<i>student</i>	<i>credit_rating</i>	<i>Class: buys_computer</i>
1	youth	high	no	fair	no
2	youth	high	no	excellent	no
3	middle_aged	high	no	fair	yes
4	senior	medium	no	fair	yes
5	senior	low	yes	fair	yes
6	senior	low	yes	excellent	no
7	middle_aged	low	yes	excellent	yes
8	youth	medium	no	fair	no
9	youth	low	yes	fair	yes
10	senior	medium	yes	fair	yes
11	youth	medium	yes	excellent	yes
12	middle_aged	medium	no	excellent	yes
13	middle_aged	high	yes	fair	yes
14	senior	medium	no	excellent	no

Data to be classified:

X = (age = youth, income = medium, student = yes, credit_rating = fair)

□ We need to maximize $P(X|C_i)P(C_i)$, for $i = 1, 2$.

Compute: $P(C_i)$, the *prior probability of each* class,

$P(\text{buys_computer} =$
“yes”) = $9/14 = 0.643$

$P(\text{buys_computer} = \text{“no”})$
= $5/14 = 0.357$

Class-Labeled Training Tuples from the *AllElectronics* Customer

RID	age	income	student	credit_rating	Class: bu
1	youth	high	no	fair	no
2	youth	high	no	excellent	no
3	middle_aged	high	no	fair	yes
4	senior	medium	no	fair	yes
5	senior	low	yes	fair	yes
6	senior	low	yes	excellent	no
7	middle_aged	low	yes	excellent	yes
8	youth	medium	no	fair	no
9	youth	low	yes	fair	yes
10	senior	medium	yes	fair	yes
11	youth	medium	yes	excellent	yes
12	middle_aged	medium	no	excellent	yes
13	middle_aged	high	yes	fair	yes
14	senior	medium	no	excellent	no

Compute $P(X|C_i)$ for each class

$$P(\text{age} = \text{youth} \mid \text{buys_computer} = \text{yes}) = 2/9 = 0.222$$

$$P(\text{age} = \text{youth} \mid \text{buys_computer} = \text{no}) = 3/5 = 0.6$$

$$P(\text{income} = \text{medium} \mid \text{buys_computer} = \text{yes}) = 4/9 = 0.444$$

$$P(\text{income} = \text{medium} \mid \text{buys_computer} = \text{no}) = 2/5 = 0.4$$

$$P(\text{student} = \text{yes} \mid \text{buys_computer} = \text{yes}) = 6/9 = 0.667$$

$$P(\text{student} = \text{yes} \mid \text{buys_computer} = \text{no}) = 1/5 = 0.2$$

$$P(\text{credit_rating} = \text{fair} \mid \text{buys_computer} = \text{yes}) = 6/9 = 0.667$$

$$P(\text{credit_rating} = \text{fair} \mid \text{buys_computer} = \text{no}) = 2/5 = 0.4$$

X = (age = youth , income = medium, student = yes, credit_rating = fair)

$P(X|C_i)$:

■ $P(X|\text{buys_computer} = \text{"yes"}) = 0.222 \times 0.444 \times 0.667 \times 0.667 = 0.044$

■ $P(X|\text{buys_computer} = \text{"no"}) = 0.6 \times 0.4 \times 0.2 \times 0.4 = 0.019$

To find the class, C_i , that maximizes $P(X|C_i)P(C_i)$,

Compute:

- $P(X|\text{buys_computer} = \text{"yes"}) * P(\text{buys_computer} = \text{"yes"}) = 0.028$

- $P(X|\text{buys_computer} = \text{"no"}) * P(\text{buys_computer} = \text{"no"}) = 0.007$

Therefore, X belongs to class ("buys_computer = yes")

SVM-Support Vector Machine

AVAILABLE AT:

Onebyzero Edu - Organized Learning, Smooth Career

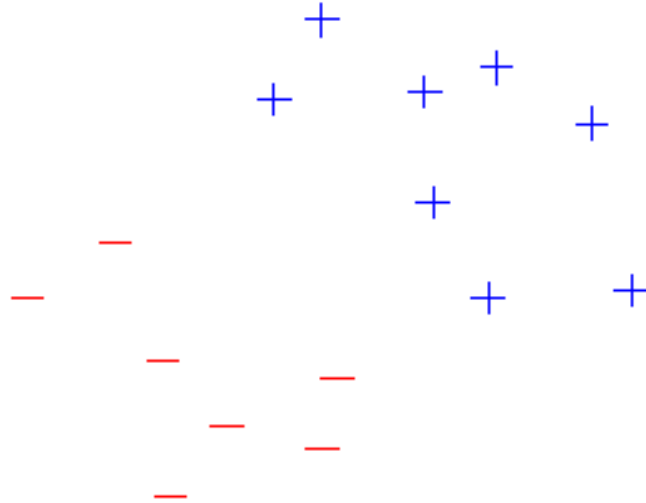
The Comprehensive Academic Study Platform for University Students in Bangladesh (www.onebyzeroedu.com)

SVM-Basics

- A relatively new classification method for both linear & nonlinear data
- It uses a nonlinear mapping to transform the original training data into a higher dimension
- With the new dimension, it searches for the linear optimal separating **hyperplane** (i.e., “decision boundary”)
- With an appropriate nonlinear mapping to a sufficiently high dimension, data from two classes can always be separated by a hyperplane
- SVM finds this hyperplane using **support vectors** (“essential” training tuples) & **margins** (defined by the support vectors)

AVAILABLE AT:

Onebyzero Edu - Organized Learning, Smooth Career
The Comprehensive Academic Study Platform for University Students in Bangladesh (www.onebyzeroedu.com)



given linearly separable data

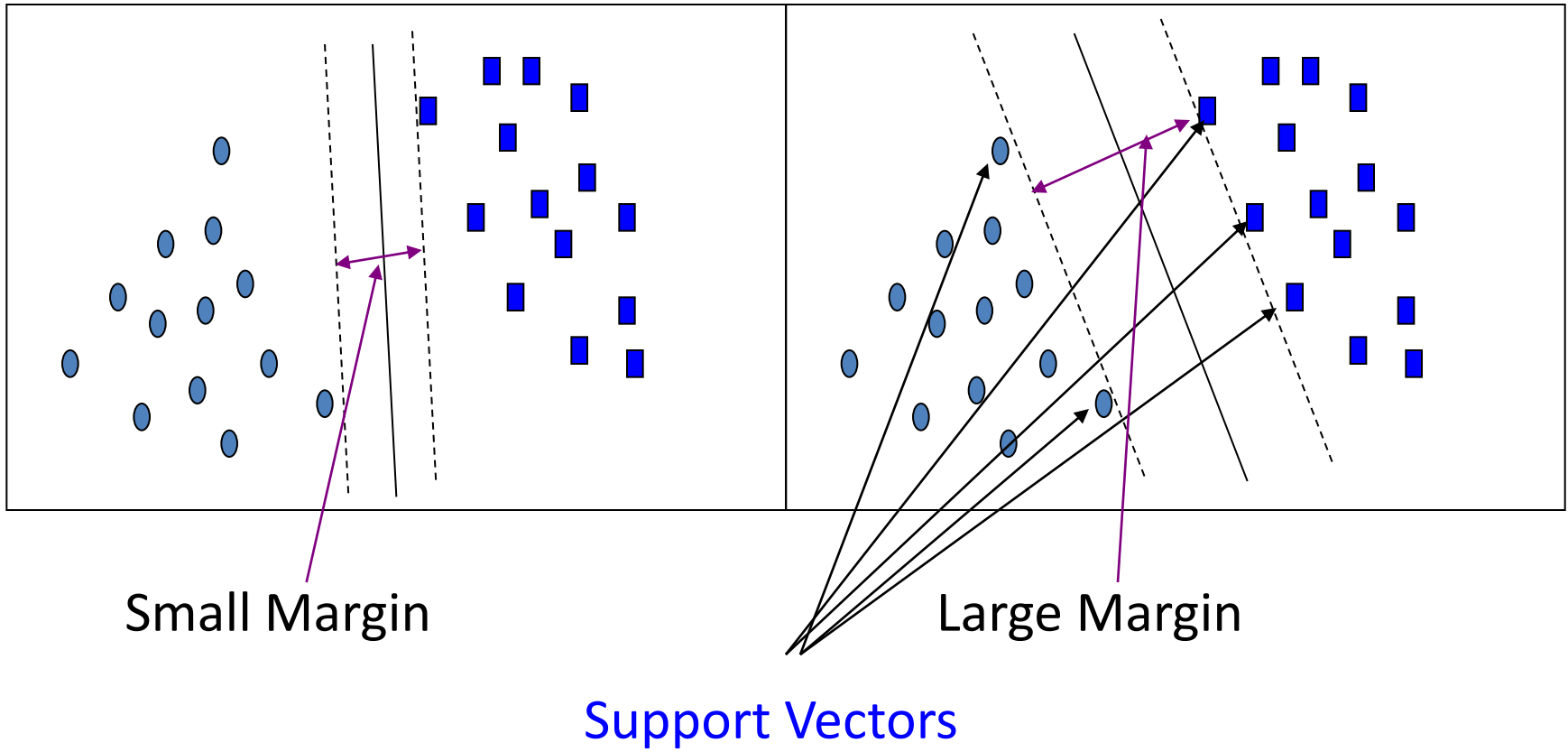
AVAILABLE AT:

Onebyzero Edu - Organized Learning, Smooth Career
The Comprehensive Academic Study Platform for University Students in Bangladesh (www.onebyzeroedu.com)

SVM—History and Applications

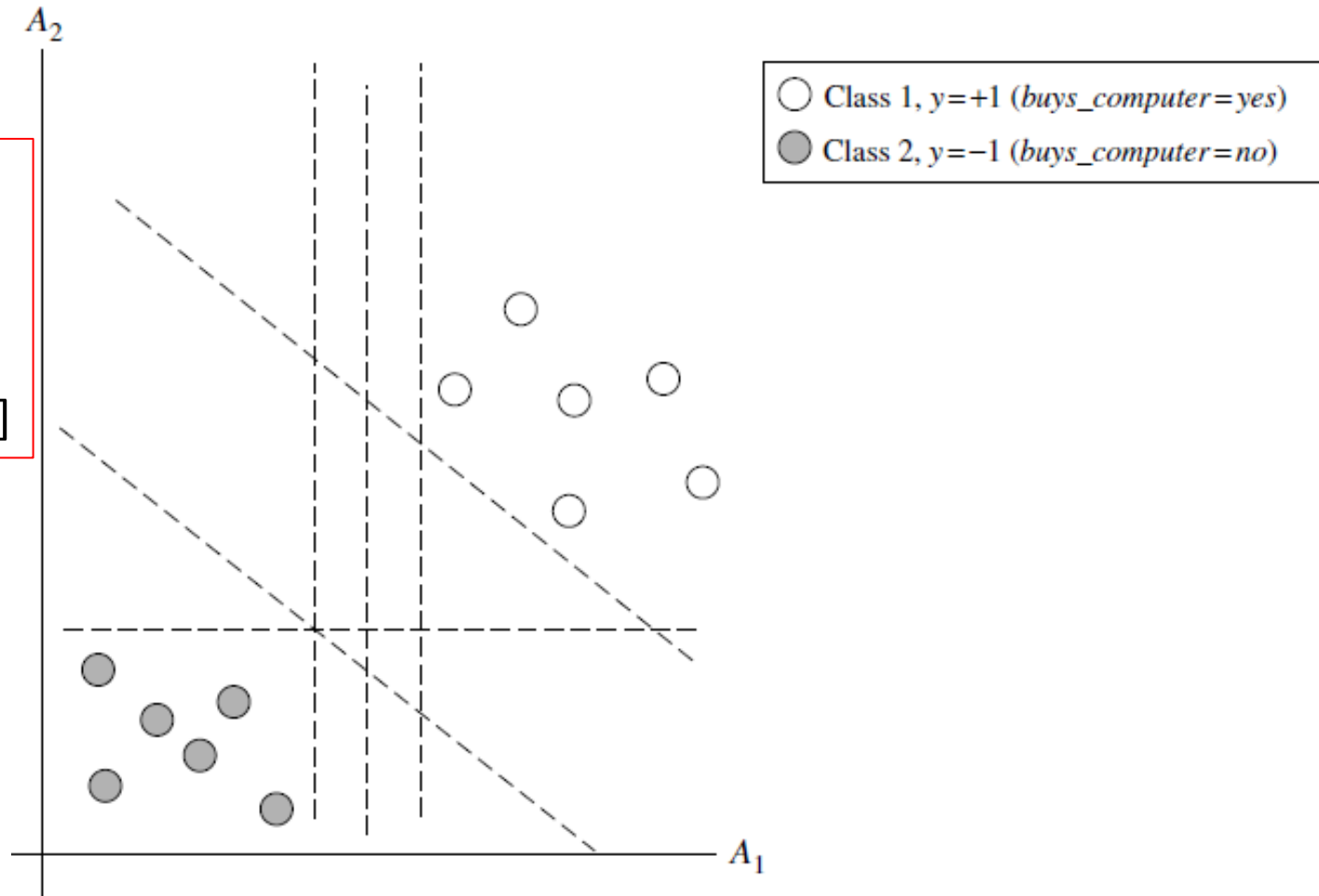
- Vapnik and colleagues (1992)—groundwork from Vapnik & Chervonenkis' statistical learning theory in 1960s
- **Features:** training can be slow but accuracy is high owing to their ability to model complex nonlinear decision boundaries (margin maximization)
- **Used for:** classification and numeric prediction
- **Applications:**
 - handwritten character recognition, object/image recognition, speaker identification, benchmarking time-series prediction tests

SVM—General Philosophy



SVM—When Data Is Linearly Separable

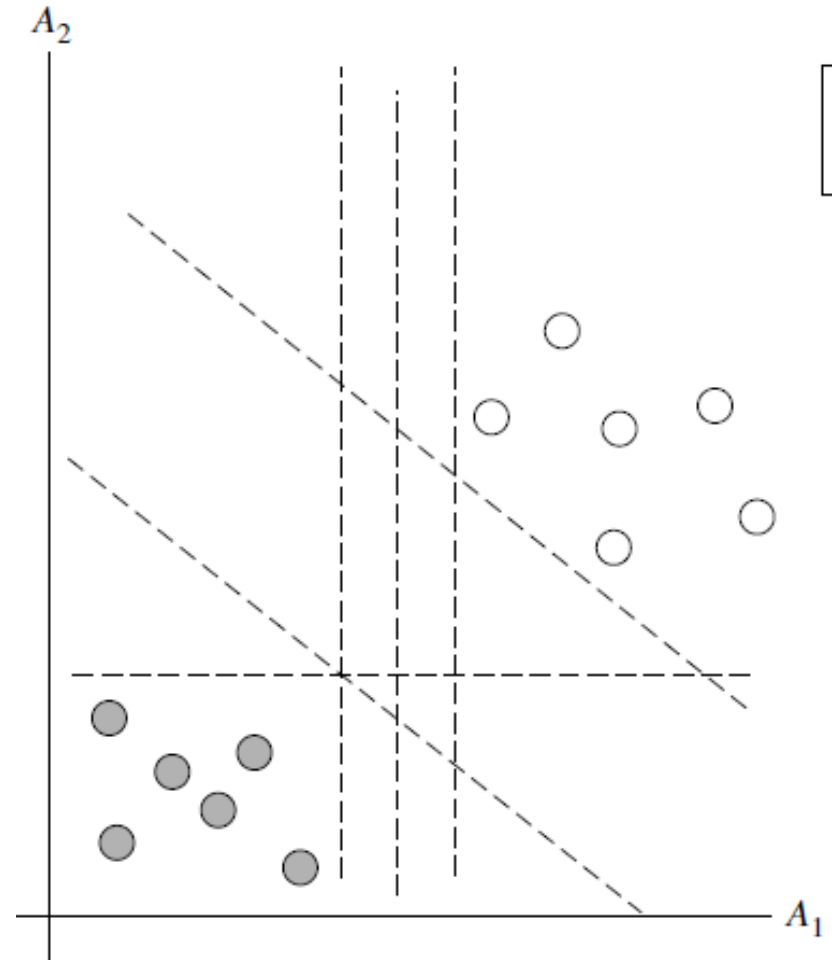
2 class problem:
Database D
 $(X_1, y_1), (X_2, y_2)$ ---
 X_i : training tuples
 y_i : class level [+1/-1]



Data are **linearly separable**, because a straight line can be drawn to separate all the tuples of class +1 from all the tuples of class -1.

Which one is Best?

- infinite number of possible separating hyperplanes or “decision boundaries,”
- **Which one is best?**

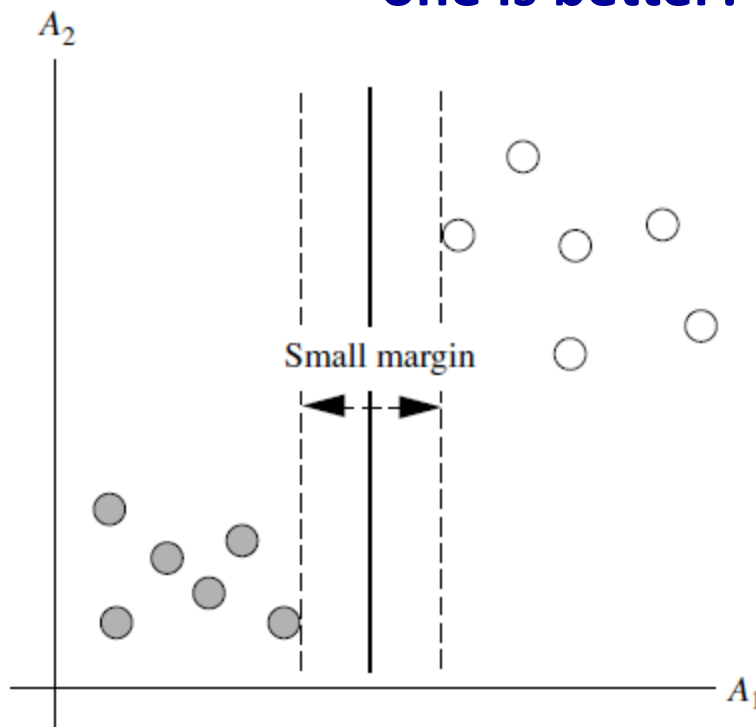


➤ We want to find the “best” one, that is, one that will have the **minimum classification error** on previously unseen tuples.

➤ **How can we find this best line/hyperplane?**

SVM searches for the hyperplane with the largest margin, i.e., maximum marginal hyperplane (MMH)

**Which
one is better?**

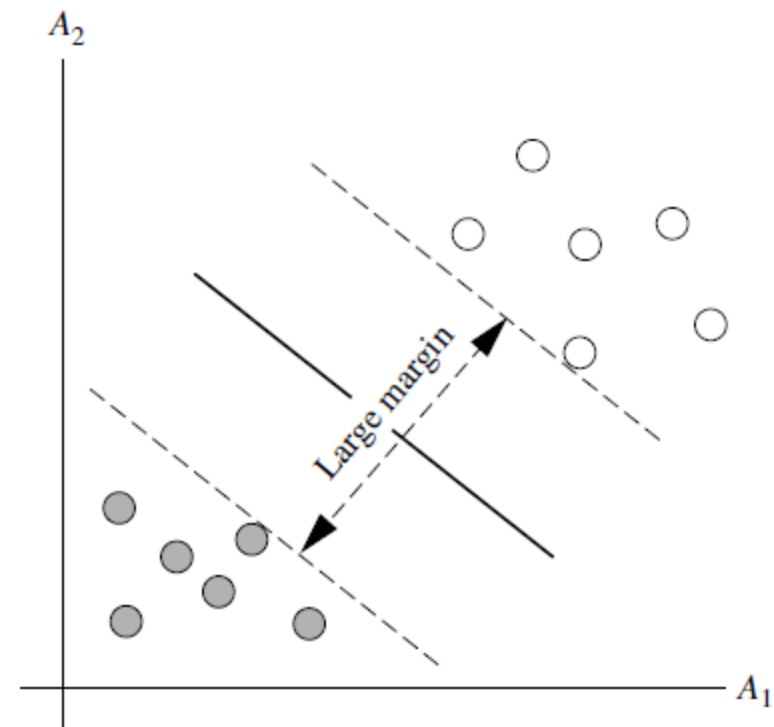


- Class 1, $y=+1$ (*buys_computer=yes*)
- Class 2, $y=-1$ (*buys_computer=no*)

AVAILABLE AT:

Onebyzero Edu - Organized Learning, Smooth Career

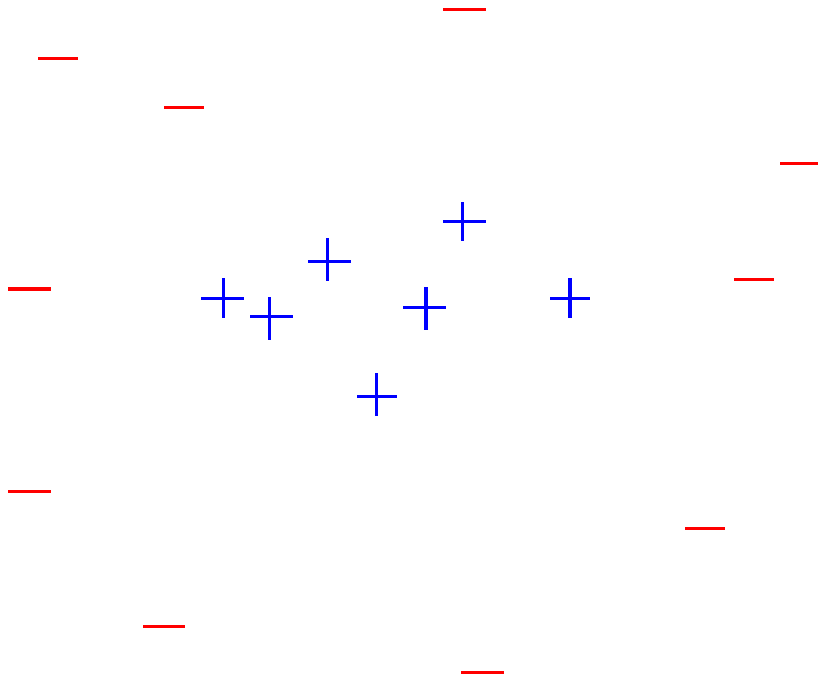
The Comprehensive Academic Study Platform for University Students in Bangladesh (www.onebyzeroedu.com)



- Class 1, $y=+1$ (*buys_computer=yes*)
- Class 2, $y=-1$ (*buys_computer=no*)

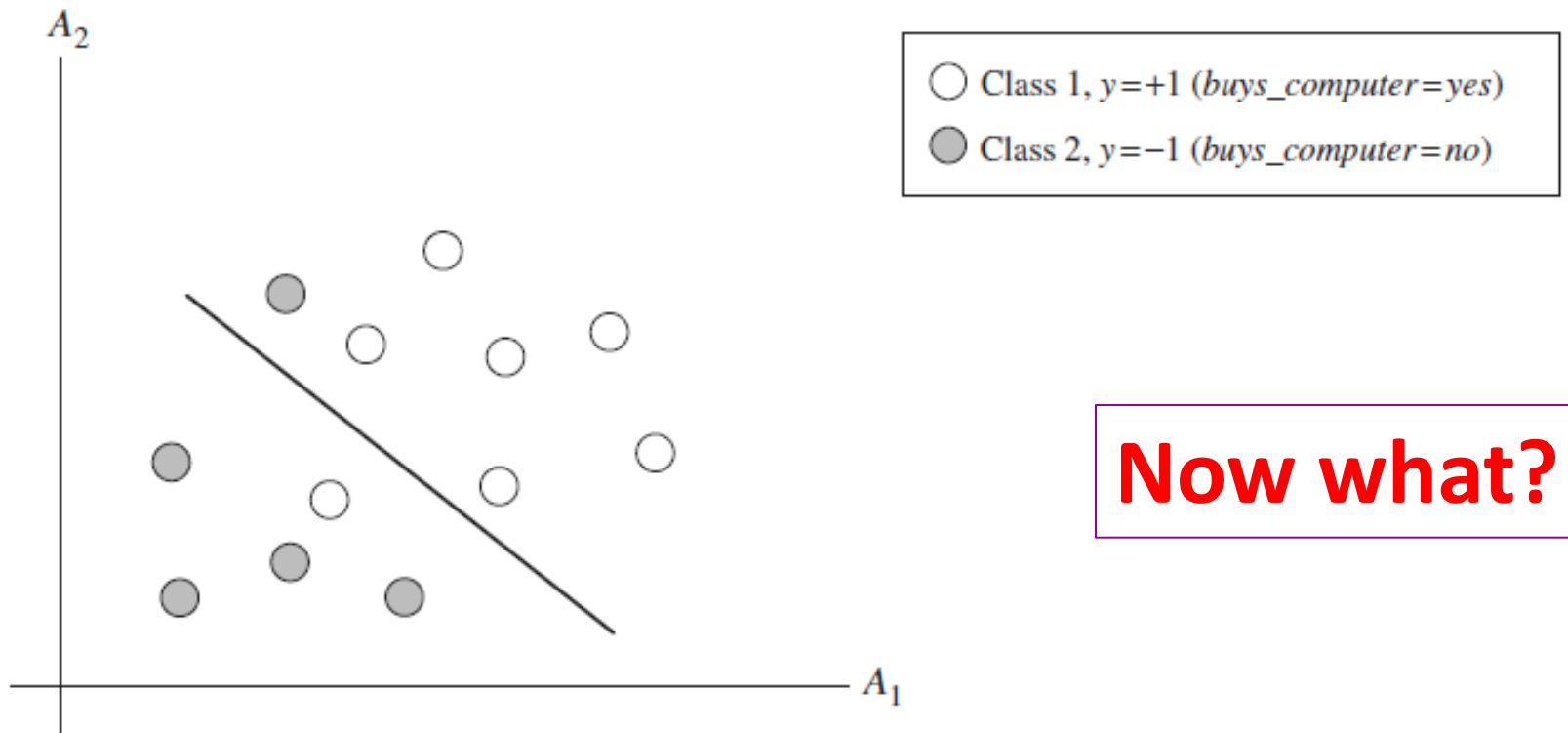
(b)

SVM—Linearly Inseparable



not linearly separable

SVM—Linearly Inseparable



It is **not possible to draw a straight line** to separate the classes. Instead, the decision boundary is nonlinear

AVAILABLE AT:

Onebyzero Edu - Organized Learning, Smooth Career

The Comprehensive Academic Study Platform for University Students in Bangladesh (www.onebyzeroedu.com)

Extend linear approach...

- **How can we extend the linear approach?**
- Two main steps:
 - Transform the original input data into a higher dimensional space using a nonlinear mapping.
 - Searches for a linear separating hyperplane in the new space.
- **Outcomes:** a quadratic optimization problem that can be solved using the linear SVM formulation.
- The maximal marginal hyperplane found in the new space corresponds to a nonlinear separating hypersurface in the original space.

Tips! If Writing Your Own Code

- Matlab are great for easy coding, but for speed, may need C or java
- **debugging** machine learning algorithms is very tricky!
 - hard to tell if working, since **don't know** what to expect
 - run on **small** cases where can figure out answer by hand
 - test each module/subroutine **separately**
 - compare to other **implementations** (written by others, or written in different language)
 - compare to **theory** or published results

Conclusion

- ML: How can we program systems to automatically learn & to improve with experience?
- In order to build up an intelligent or autonomous agents, ML should be used.
- **Still wide gaps:** need to perform plenty of social & cognitive capabilities
- Not so far away from our dreams!
- Artificial Humans co-play in the human world very soon.

Conclusion

- AI dream of someday building machines as intelligent as you or I - Andrew Ng.

A Puzzle.....

Who is the real human?



AVAILABLE AT:

Onebyzero Edu - Organized Learning, Smooth Career

The Comprehensive Academic Study Platform for University Students in Bangladesh (www.onebyzeroedu.com)

Thanks

- [V3](#)

AVAILABLE AT:

Onebyzero Edu - Organized Learning, Smooth Career
The Comprehensive Academic Study Platform for University Students in Bangladesh (www.onebyzeroedu.com)

AVAILABLE AT:

Onebyzero Edu - Organized Learning, Smooth Career

The Comprehensive Academic Study Platform for University Students in Bangladesh (www.onebyzeroedu.com)