

Big Data Mining

巨量資料探勘

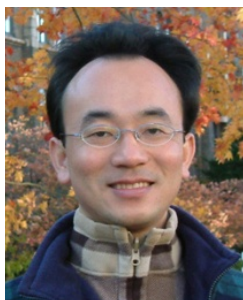
分類與預測

(Classification and Prediction)

1062DM04

MI4 (M2244) (2995)

Wed, 9, 10 (16:10-18:00) (B206)



Min-Yuh Day

戴敏育

Assistant Professor

專任助理教授

Dept. of Information Management, Tamkang University

淡江大學 資訊管理學系

<http://mail.tku.edu.tw/myday/>

2018-03-28



課程大綱 (Syllabus)

週次 (Week)	日期 (Date)	內容 (Subject/Topics)
1	2018/02/28	和平紀念日(放假一天) (Peace Memorial Day) (Day off)
2	2018/03/07	巨量資料探勘課程介紹 (Course Orientation for Big Data Mining)
3	2018/03/14	大數據、AI人工智慧與深度學習 (Big Data, Artificial Intelligence and Deep Learning)
4	2018/03/21	關連分析 (Association Analysis)
5	2018/03/28	分類與預測 (Classification and Prediction)
6	2018/04/04	兒童節(放假一天)(Children's Day) (Day off)
7	2018/04/11	分群分析 (Cluster Analysis)
8	2018/04/18	個案分析與實作一 (SAS EM 分群分析) : Case Study 1 (Cluster Analysis - K-Means using SAS EM)

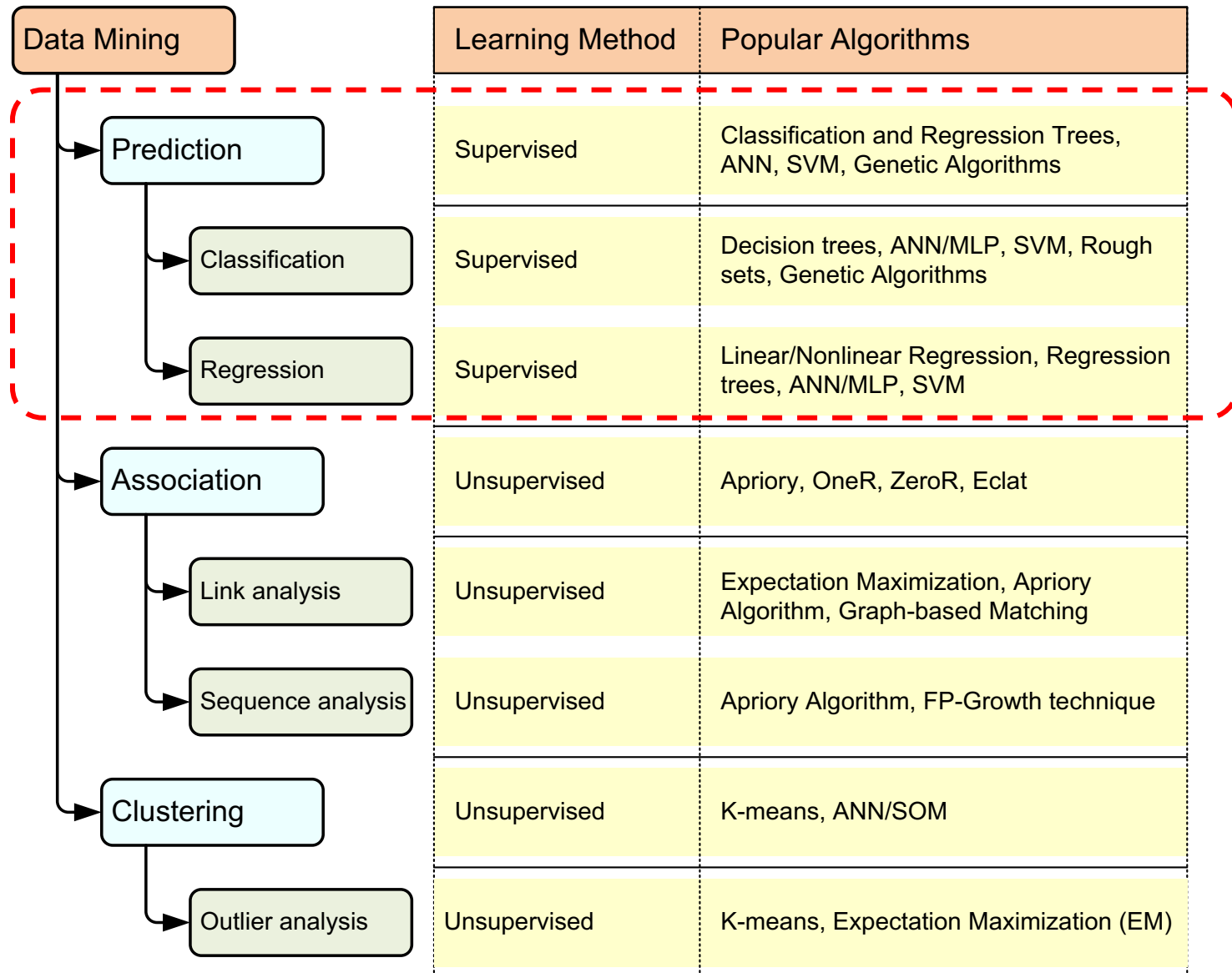
課程大綱 (Syllabus)

週次 (Week)	日期 (Date)	內容 (Subject/Topics)
9	2018/04/25	期中報告 (Midterm Project Presentation)
10	2018/05/02	期中考試週
11	2018/05/09	個案分析與實作二 (SAS EM 關連分析) : Case Study 2 (Association Analysis using SAS EM)
12	2018/05/16	個案分析與實作三 (SAS EM 決策樹、模型評估) : Case Study 3 (Decision Tree, Model Evaluation using SAS EM)
13	2018/05/23	個案分析與實作四 (SAS EM 迴歸分析、類神經網路) : Case Study 4 (Regression Analysis, Artificial Neural Network using SAS EM)
14	2018/05/30	期末報告 (Final Project Presentation)
15	2018/06/06	畢業考試週

Outline

- Classification and Prediction
- Supervised Learning (Classification)
- Decision Tree (DT)
 - Information Gain (IG)
- Support Vector Machine (SVM)
- Data Mining Evaluation
 - Accuracy
 - Precision
 - Recall
 - F1 score (F-measure) (F-score)

A Taxonomy for Data Mining Tasks



Customer database

ID	age	income	student	credit_rating	Class: buys_computer
1	youth	high	no	fair	no
2	middle_aged	high	no	fair	yes
3	youth	high	no	excellent	no
4	senior	medium	no	fair	yes
5	senior	high	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	excellent	yes

What is the class
(buys_computer = “yes” or
buys_computer = “no”)
for a customer
(age=youth, income=medium,
student =yes, credit= fair)?

Customer database

ID	age	income	student	credit_rating	Class: buys_computer
1	youth	high	no	fair	no
2	middle_aged	high	no	fair	yes
3	youth	high	no	excellent	no
4	senior	medium	no	fair	yes
5	senior	high	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	excellent	yes
11	youth	medium	yes	fair	?

Customer database

ID	age	income	student	credit_rating	Class: buys_computer
1	youth	high	no	fair	no
2	middle_aged	high	no	fair	yes
3	youth	high	no	excellent	no
4	senior	medium	no	fair	yes
5	senior	high	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	excellent	yes
11	youth	medium	yes	fair	Yes (0.0889)

What is the **class**

(**buys_computer = "yes"**) or
buys_computer = "no")

for a **customer**

(age=youth, income=medium,
student =yes, credit= fair)?

Yes = 0.0889

No = 0.0167

Classification vs. Prediction

- Classification
 - predicts **categorical class** labels (discrete or nominal)
 - classifies data (constructs a model) based on the training set and the values (**class labels**) in a classifying attribute and uses it in classifying new data
- Prediction
 - models **continuous-valued** functions
 - i.e., predicts unknown or missing values
- Typical applications
 - Credit approval
 - Target marketing
 - Medical diagnosis
 - Fraud detection

Data Mining Methods: Classification

- Most frequently used DM method
- Part of the machine-learning family
- Employ supervised learning
- Learn from past data, classify new data
- The output variable is categorical (nominal or ordinal) in nature
- Classification versus regression?
- Classification versus clustering?

Classification Techniques

- **Decision Tree analysis (DT)**
- Statistical analysis
- **Neural networks (NN)**
- **Deep Learning (DL)**
- **Support Vector Machines (SVM)**
- Case-based reasoning
- Bayesian classifiers
- Genetic algorithms (GA)
- Rough sets

Text Mining

(Text Data Mining)



Example of Opinion: review segment on iPhone



“I bought an iPhone a few days ago.

It was such a nice phone.

The touch screen was really cool.

The voice quality was clear too.

However, my mother was mad with me as I did not tell her before I bought it.

She also thought the phone was too expensive, and wanted me to return it to the shop. ... ”

Example of Opinion: review segment on iPhone

“(1) I bought an iPhone a few days ago.

(2) It was such a **nice** phone.

(3) The touch screen was really **cool**.

(4) The voice quality was **clear** too.

(5) However, my mother was mad with me as I did not tell her before I bought it.

(6) She also thought the phone was too expensive, and wanted me to return it to the shop. ...”



+Positive
Opinion



-Negative
Opinion

Text mining

Text Data Mining

Intelligent Text Analysis

Knowledge-Discovery in Text (KDT)

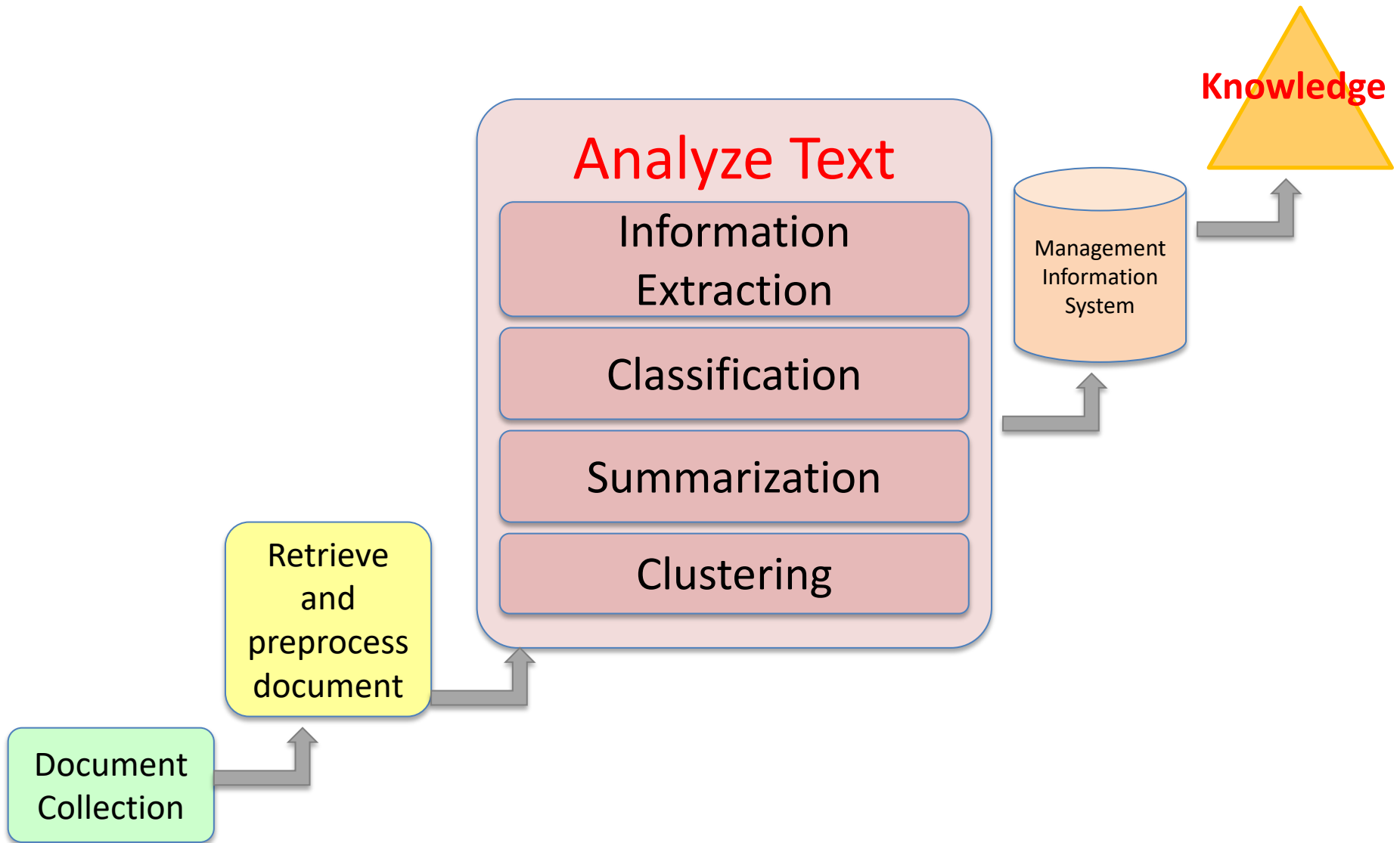
Text Mining:
the process of extracting
interesting and non-trivial
information and knowledge
from unstructured text.

Text Mining:
discovery by computer of
new, previously
unknown information,
by automatically
extracting information
from different written resources.

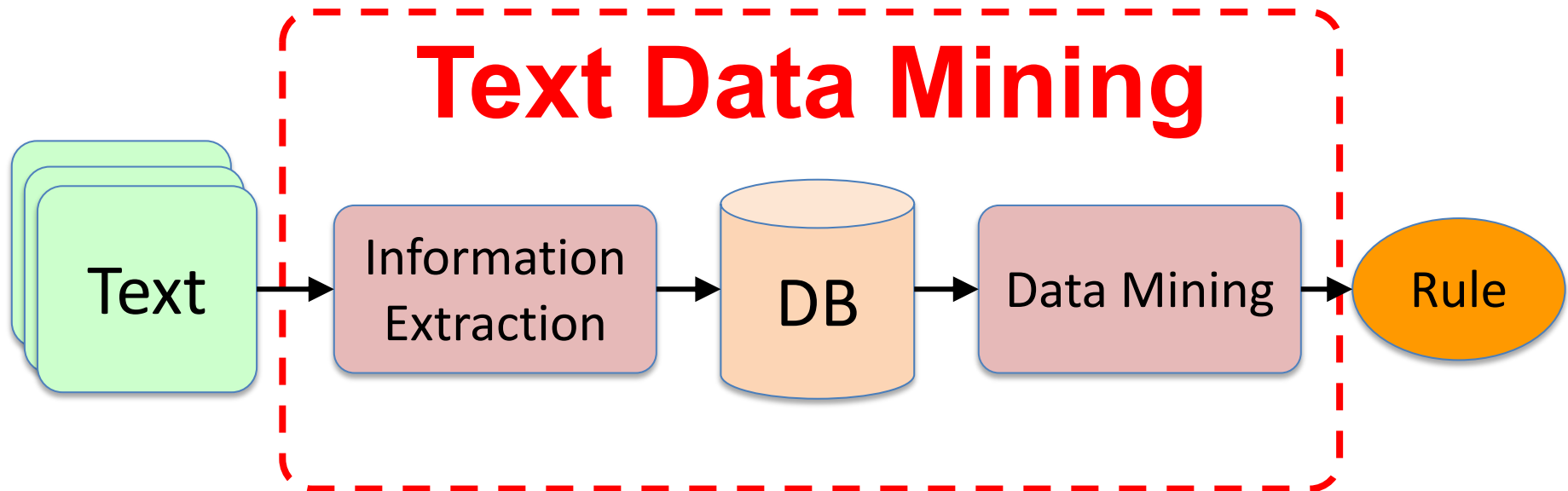
Text Mining (TM)

**Natural Language Processing
(NLP)**

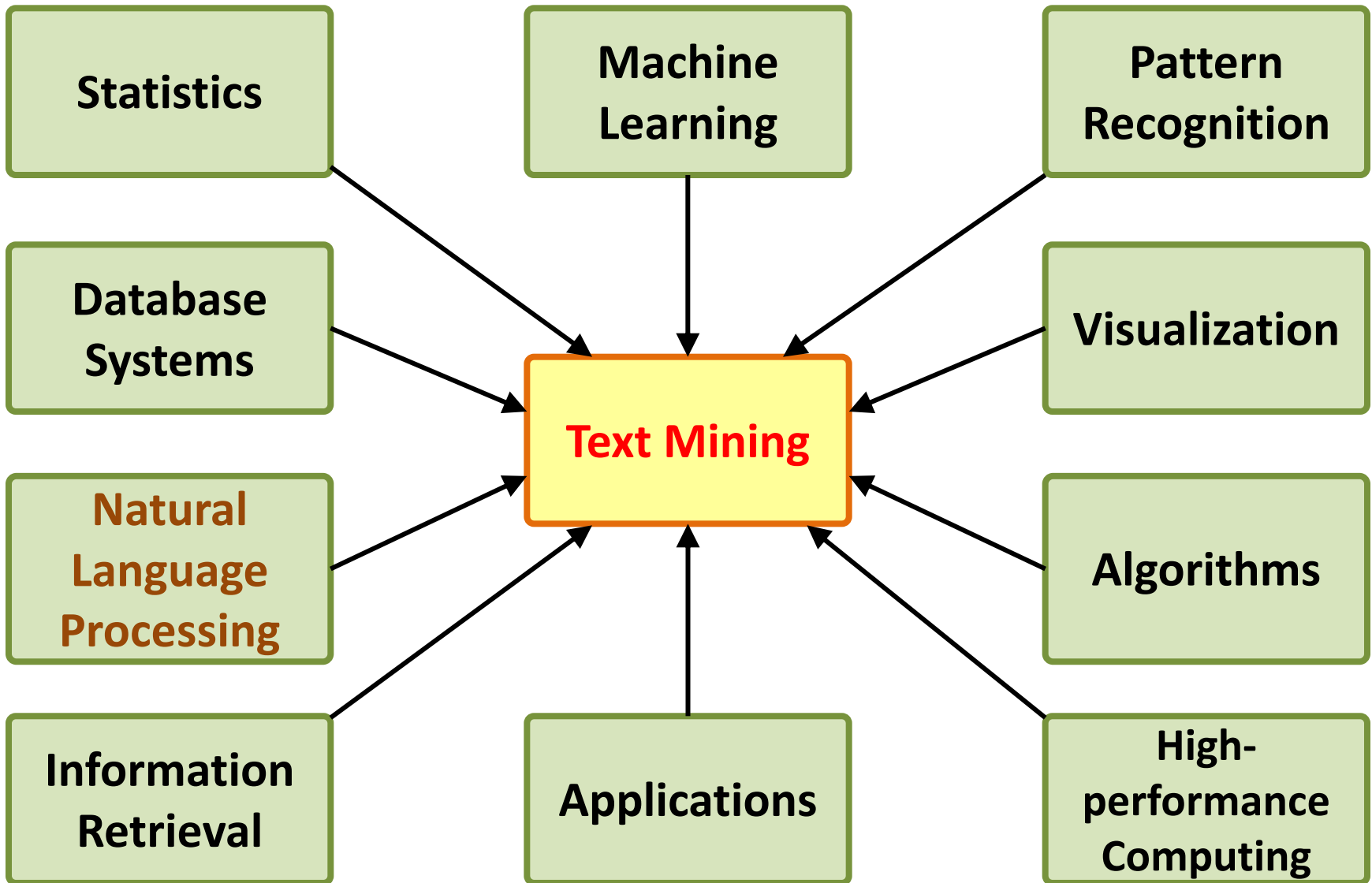
An Example of Text Mining



Overview of Information Extraction based Text Mining Framework



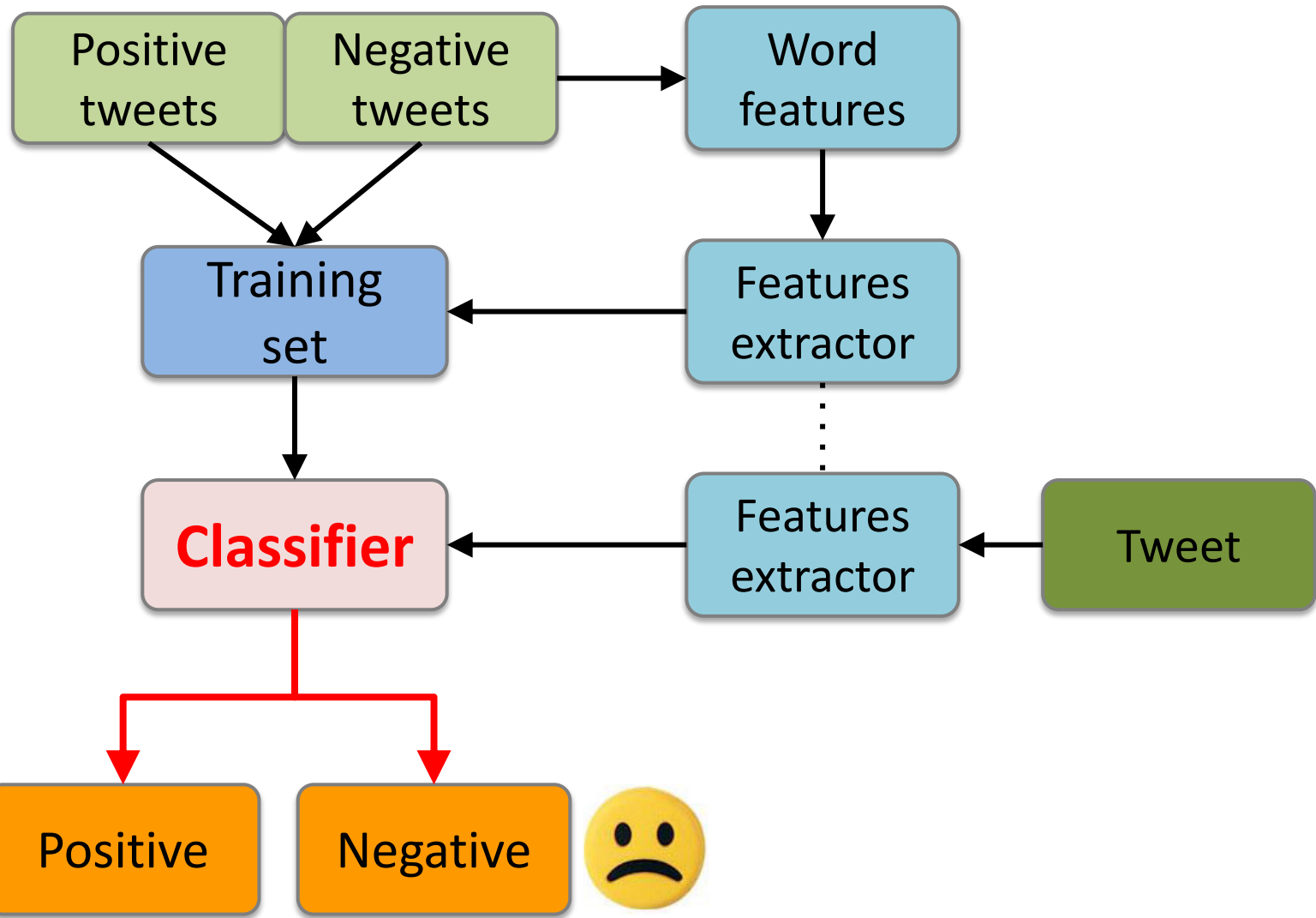
Text Mining Technologies



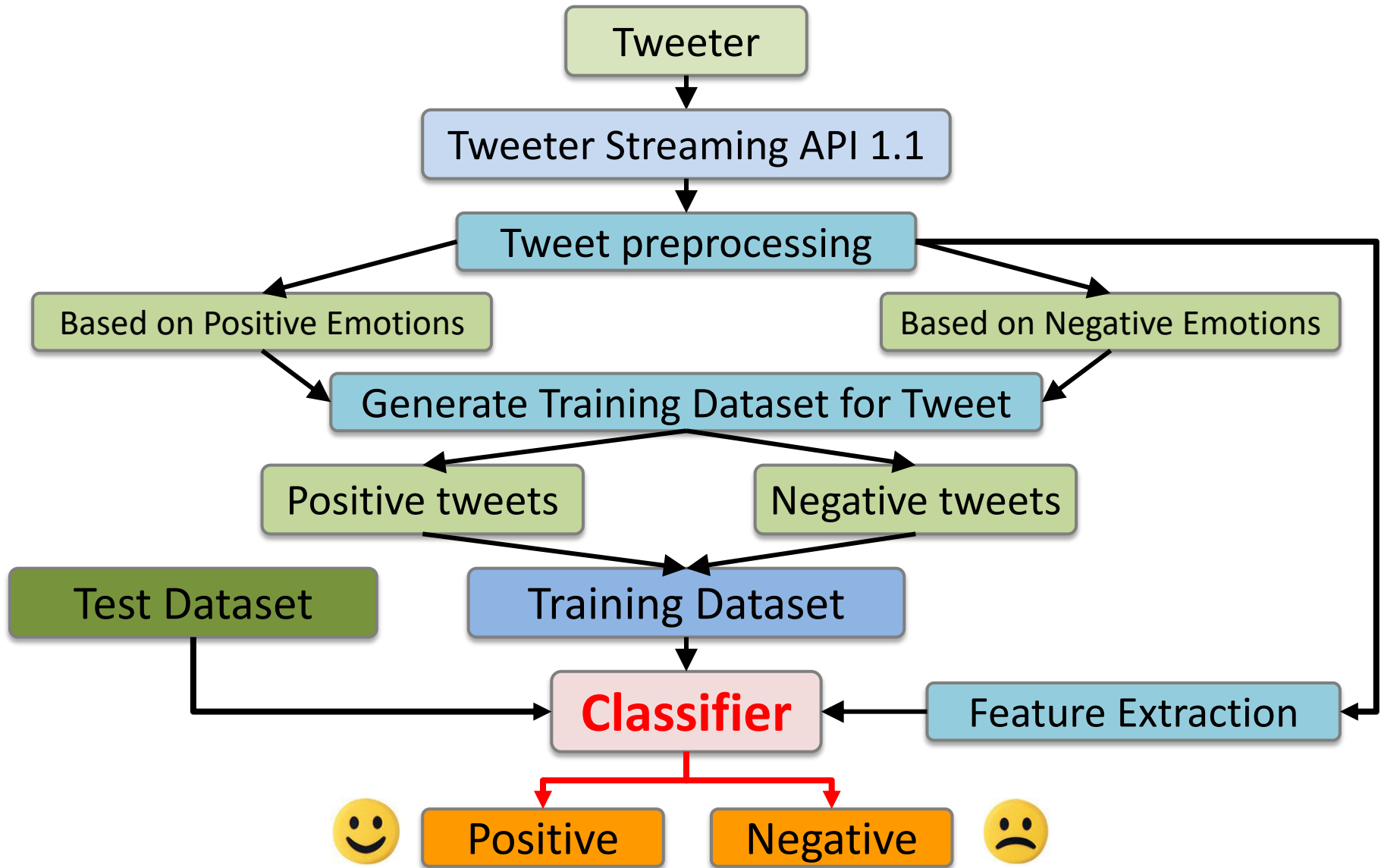
Data Mining versus Text Mining

- Both seek for novel and useful patterns
- Both are semi-automated processes
- Difference is the nature of the data:
 - Structured versus unstructured data
 - **Structured data:** in databases
 - **Unstructured data:** Word documents, PDF files, text excerpts, XML files, and so on
- Text mining – first, impose structure to the data, then mine the structured data

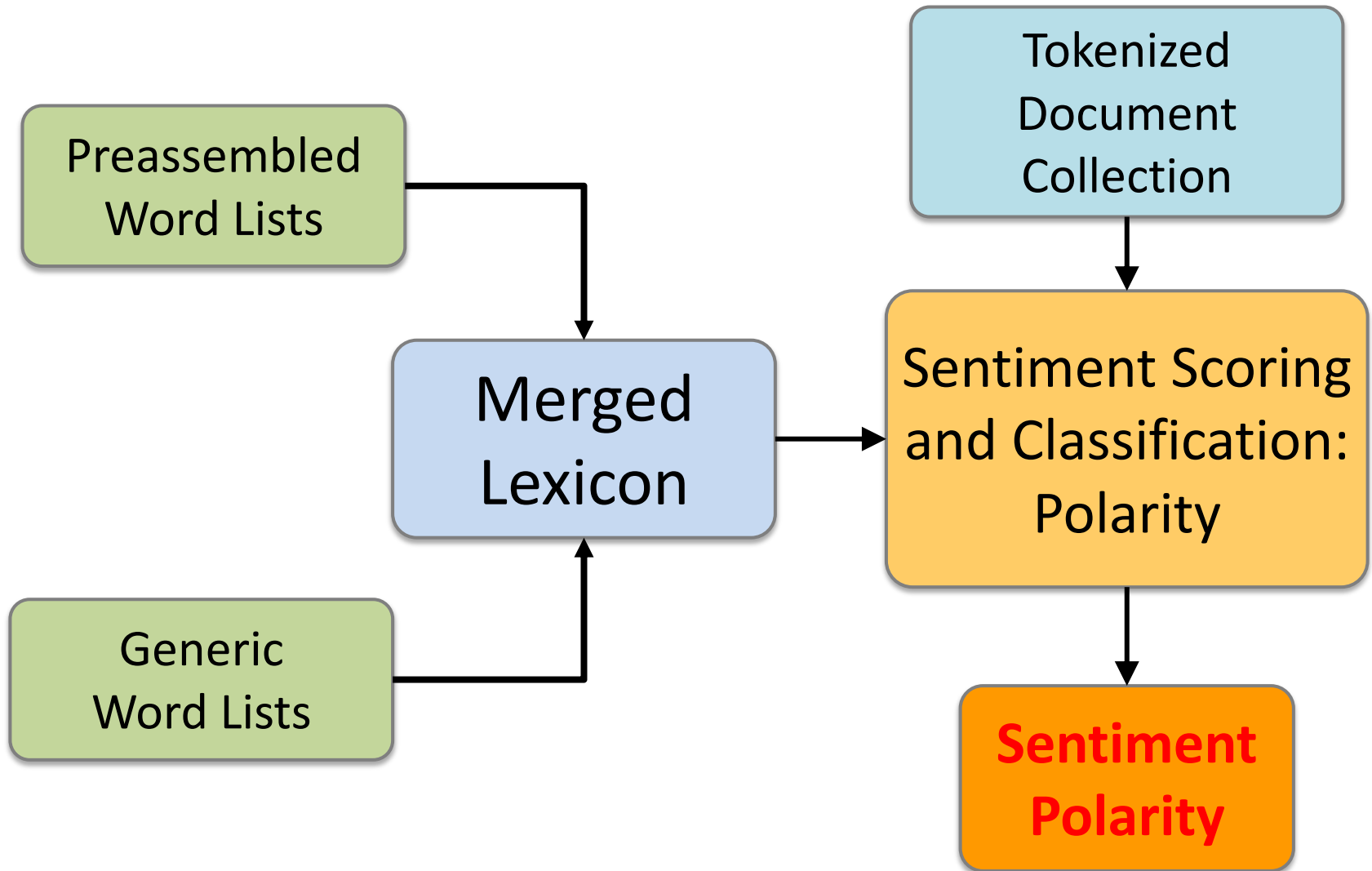
Sentiment Analysis Architecture



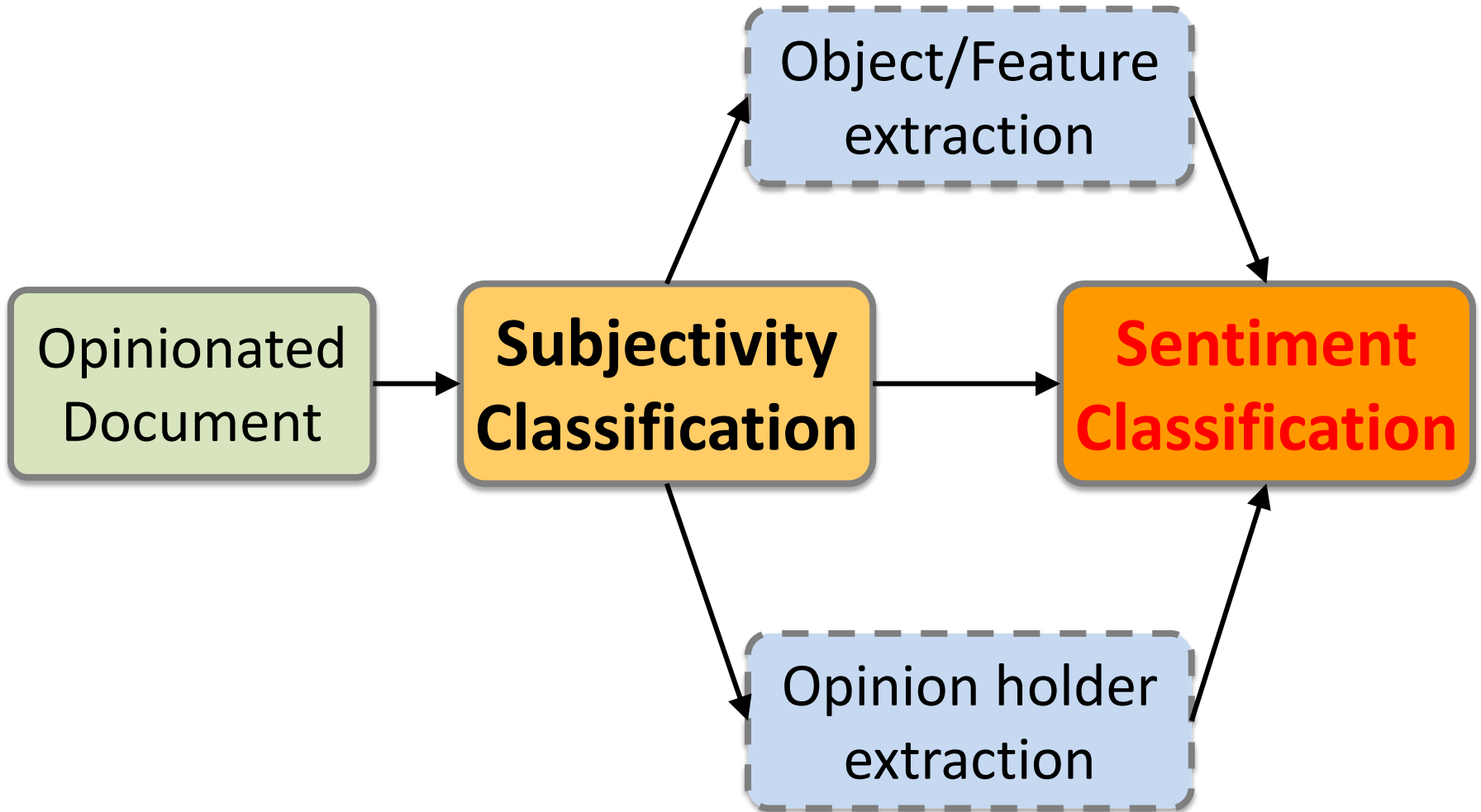
Sentiment Classification Based on Emoticons



Lexicon-Based Model



Sentiment Analysis Tasks



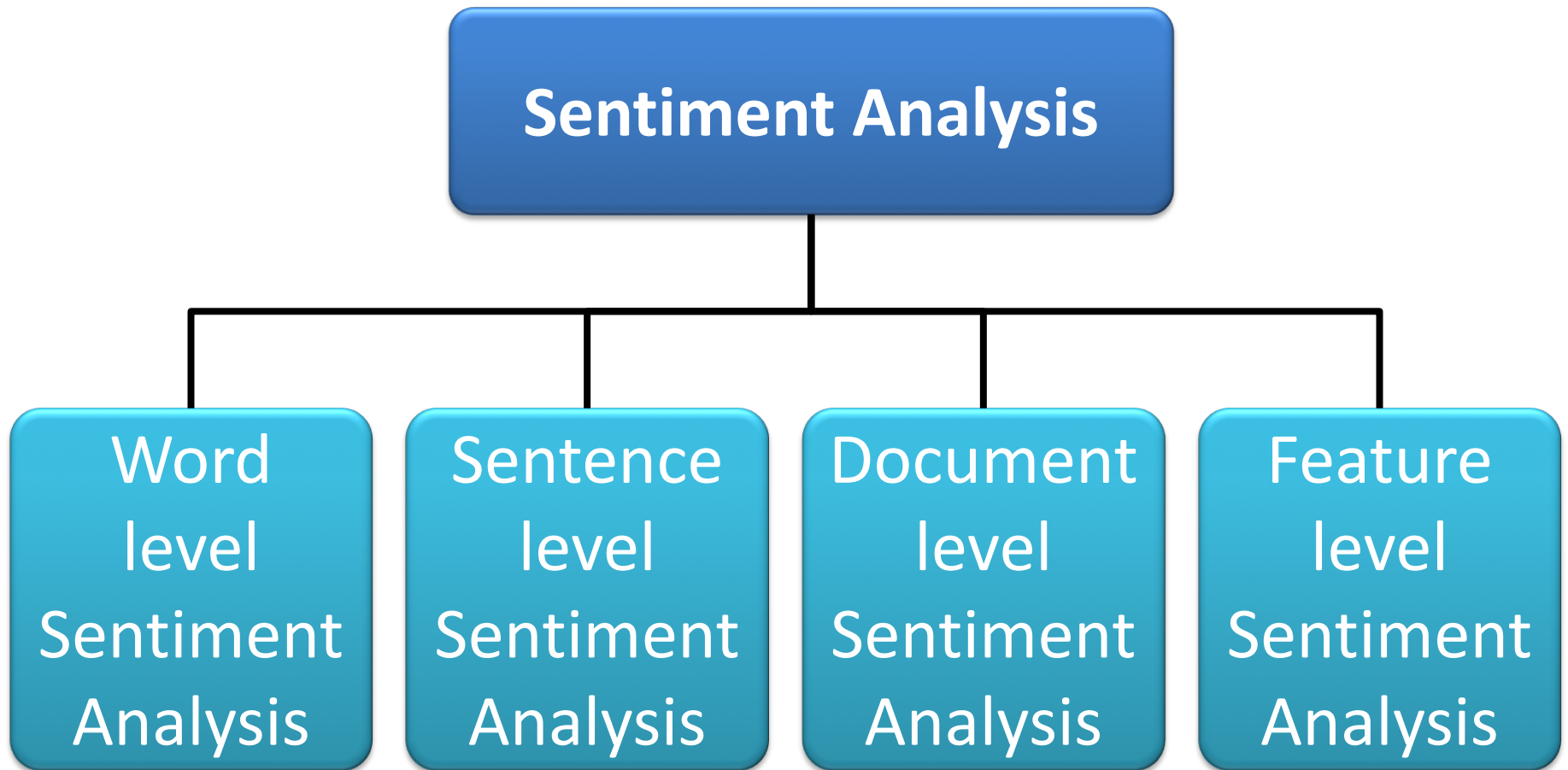
Sentiment Analysis

vs.

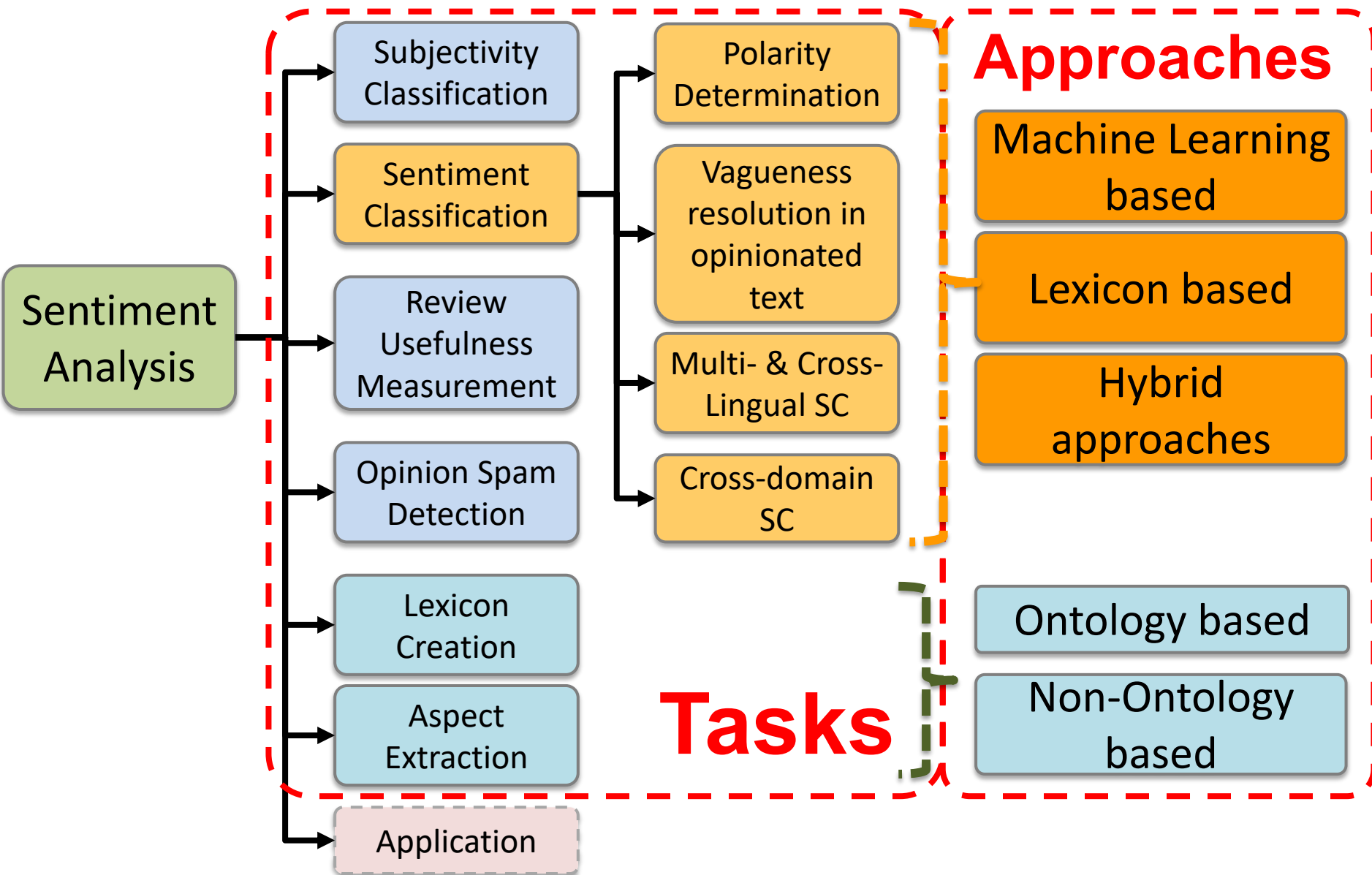
Subjectivity Analysis

Sentiment Analysis	Subjectivity Analysis
Positive	Subjective
Negative	
Neutral	Objective

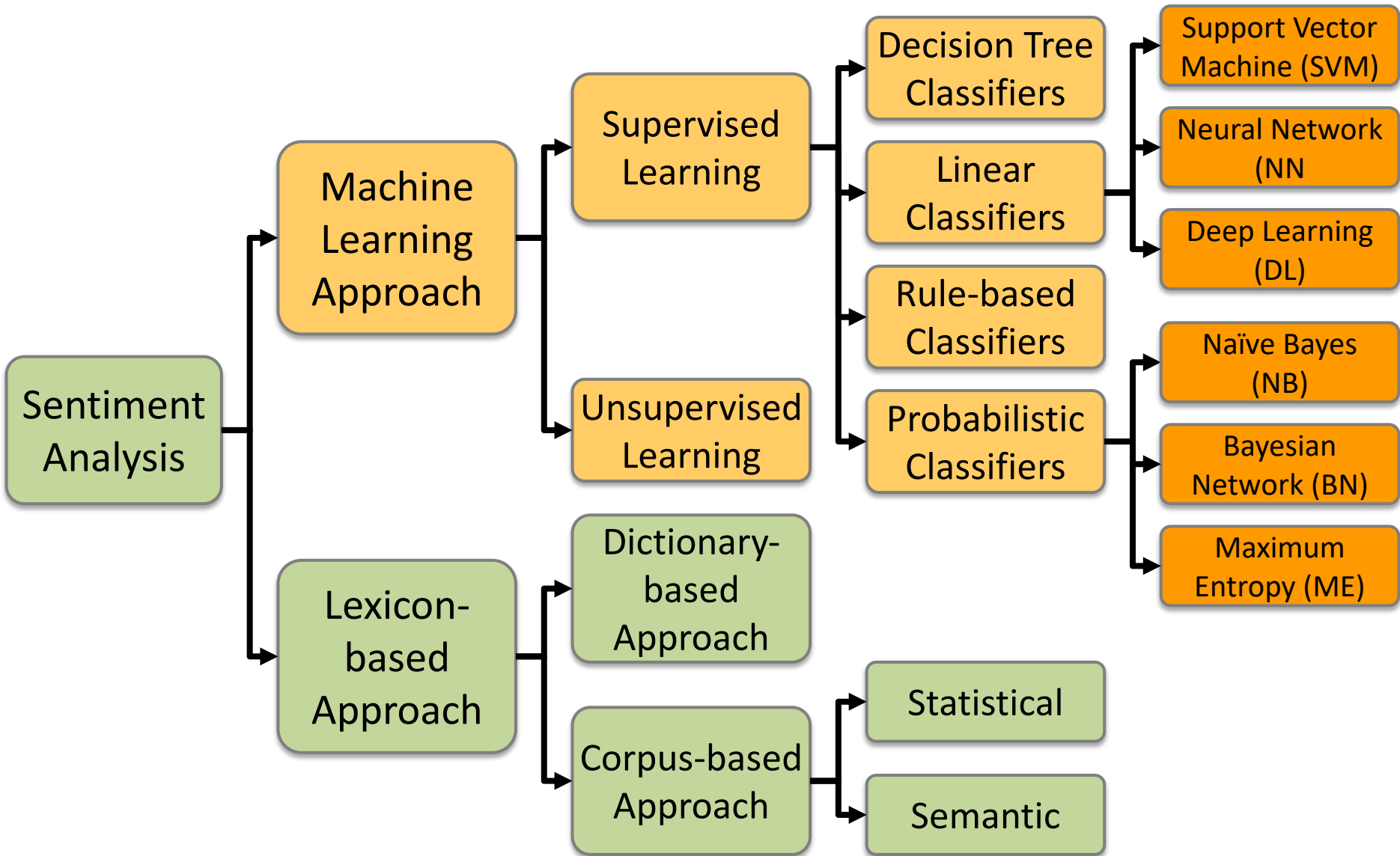
Levels of Sentiment Analysis



Sentiment Analysis



Sentiment Classification Techniques



Machine Learning Models

Deep Learning

Association rules

Decision tree

Clustering

Bayesian

Kernel

Ensemble

Dimensionality reduction

Regression Analysis

Instance based

Example of Classification

- Loan Application Data
 - Which loan applicants are “safe” and which are “risky” for the bank?
 - “Safe” or “risky” for load application data
- Marketing Data
 - Whether a customer with a given profile will buy a new computer?
 - “yes” or “no” for marketing data
- **Classification**
 - Data analysis task
 - A model or **Classifier** is constructed to predict categorical labels
 - Labels: “safe” or “risky”; “yes” or “no”; “treatment A”, “treatment B”, “treatment C”

What Is Prediction?

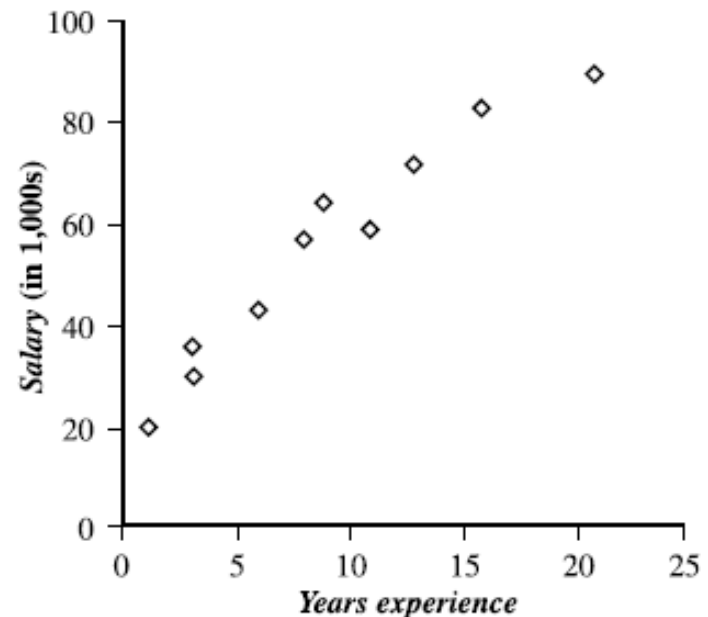
- (Numerical) prediction is similar to classification
 - construct a model
 - use model to predict continuous or ordered value for a given input
- Prediction is different from classification
 - **Classification** refers to predict **categorical class** label
 - **Prediction** models **continuous-valued** functions
- Major method for prediction: **regression**
 - model the relationship between one or more *independent* or **predictor** variables and a *dependent* or **response** variable
- Regression analysis
 - Linear and multiple regression
 - Non-linear regression
 - Other regression methods: generalized linear model, Poisson regression, log-linear models, regression trees

Prediction Methods

- Linear Regression
- Nonlinear Regression
- Other Regression Methods

Salary data.

<i>x</i> years experience	<i>y</i> salary (in \$1000s)
3	30
8	57
9	64
13	72
3	36
6	43
11	59
21	90
1	20
16	83



Classification and Prediction

- **Classification** and **prediction** are two forms of data analysis that can be used to extract **models** describing important data classes or to predict future data trends.
- **Classification**
 - Effective and scalable methods have been developed for **decision trees** induction, **Naive Bayesian classification**, **Bayesian belief network**, **rule-based classifier**, **Backpropagation**, **Support Vector Machine (SVM)**, **associative classification**, **nearest neighbor classifiers**, and **case-based reasoning**, and other classification methods such as **genetic algorithms**, **rough set** and **fuzzy set** approaches.
- **Prediction**
 - **Linear, nonlinear, and generalized linear models of regression** can be used for **prediction**. Many nonlinear problems can be converted to linear problems by performing transformations on the predictor variables. **Regression trees** and **model trees** are also used for prediction.

Classification

—A Two-Step Process

- 1. Model construction:** describing a set of predetermined classes
 - Each tuple/sample is assumed to belong to a predefined class, as determined by the **class label attribute**
 - The set of tuples used for model construction is **training set**
 - The model is represented as classification rules, decision trees, or mathematical formulae
- 2. Model usage:** for classifying future or unknown objects
 - **Estimate accuracy** of the model
 - The known label of test sample is compared with the classified result from the model
 - **Accuracy rate** is the percentage of test set samples that are correctly classified by the model
 - **Test set** is independent of **training set**, otherwise over-fitting will occur
 - If the accuracy is acceptable, use the model to **classify data** tuples whose class labels are not known

Supervised Learning vs. Unsupervised Learning

- **Supervised learning (classification)**
 - Supervision: The training data (observations, measurements, etc.) are accompanied by labels indicating the class of the observations
 - New data is classified based on the training set
- **Unsupervised learning (clustering)**
 - The class labels of training data is unknown
 - Given a set of measurements, observations, etc. with the aim of establishing the existence of classes or clusters in the data

Issues Regarding Classification and Prediction: Data Preparation

- Data cleaning
 - Preprocess data in order to reduce noise and handle missing values
- Relevance analysis (**feature selection**)
 - Remove the irrelevant or redundant attributes
 - Attribute subset selection
 - **Feature Selection** in machine learning
- Data transformation
 - Generalize and/or normalize data
 - Example
 - Income: low, medium, high

Issues:

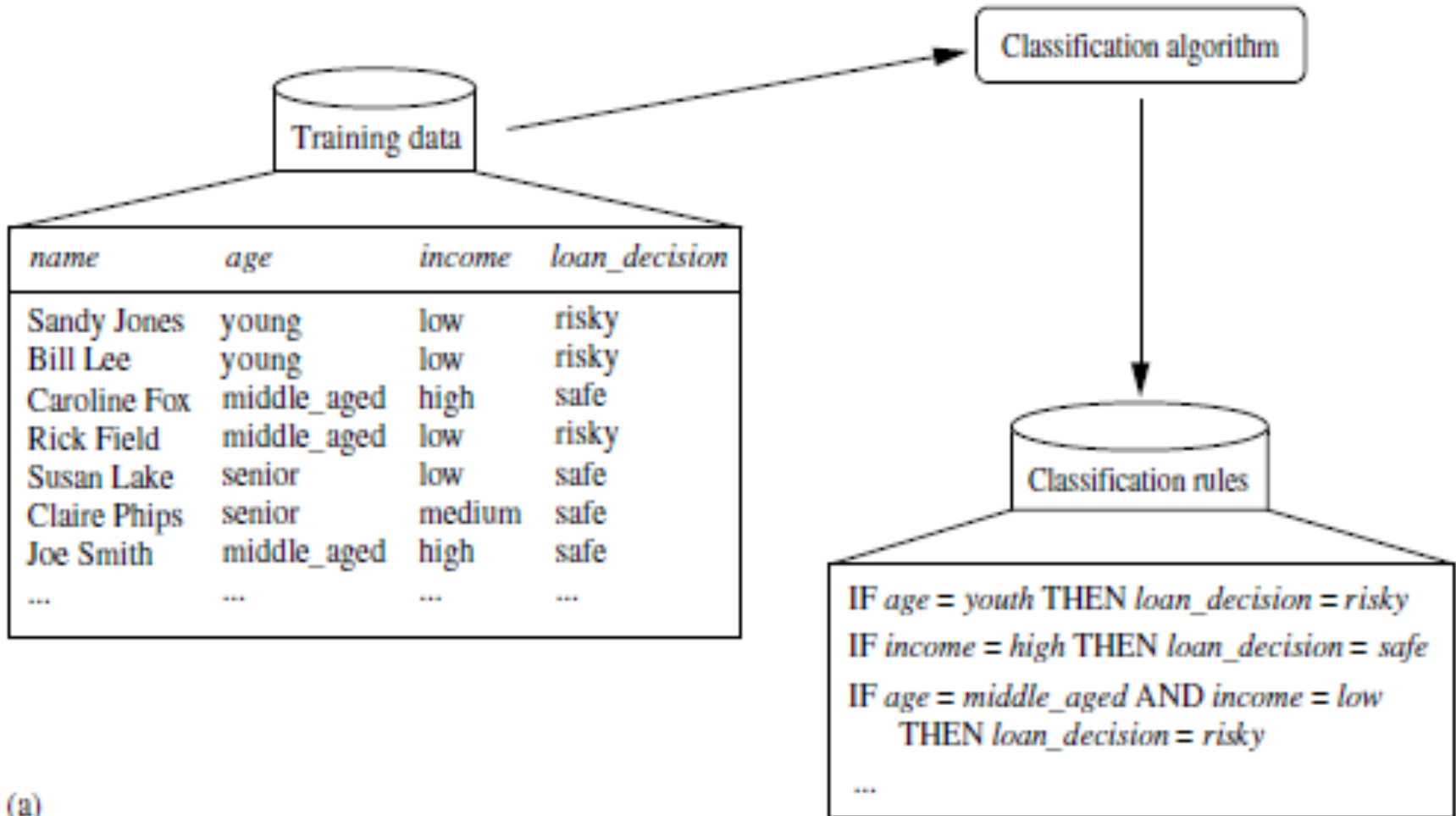
Evaluating Classification and Prediction Methods

- **Accuracy**
 - classifier accuracy: predicting class label
 - predictor accuracy: guessing value of predicted attributes
 - estimation techniques: cross-validation and bootstrapping
- Speed
 - time to construct the model (training time)
 - time to use the model (classification/prediction time)
- Robustness
 - handling noise and missing values
- Scalability
 - ability to construct the classifier or predictor efficiently given large amounts of data
- Interpretability
 - understanding and insight provided by the model

Data Classification Process 1: **Learning (Training)** Step

(a) **Learning**: Training data are analyzed by classification algorithm

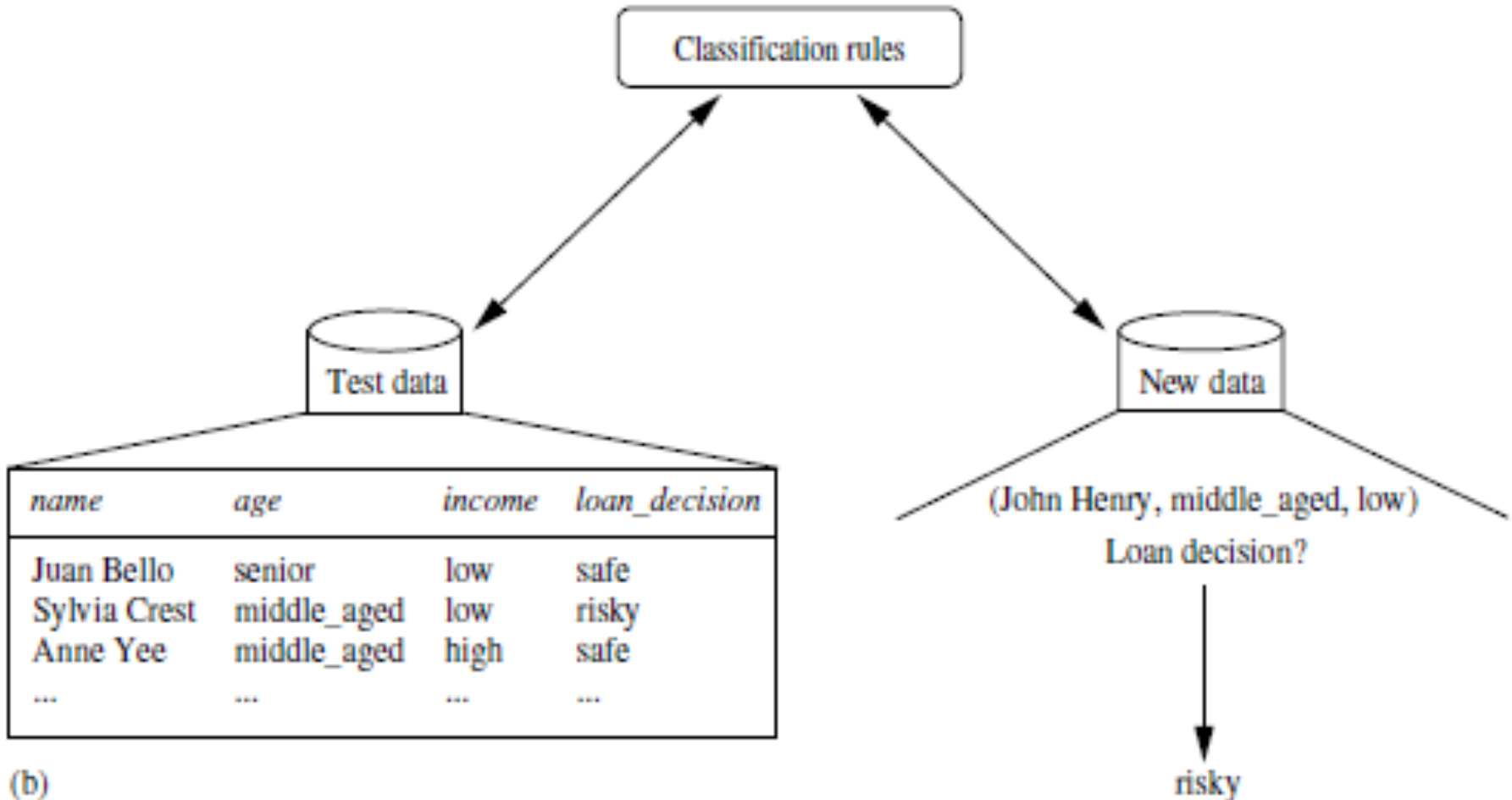
$$y = f(X)$$



(a)

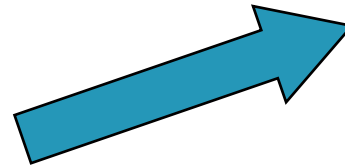
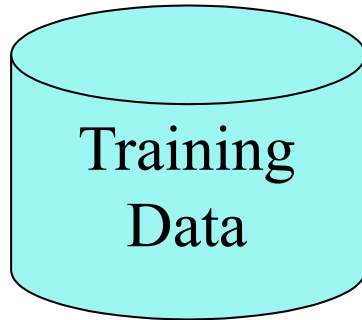
Data Classification Process 2

(b) **Classification:** Test data are used to estimate the accuracy of the classification rules.

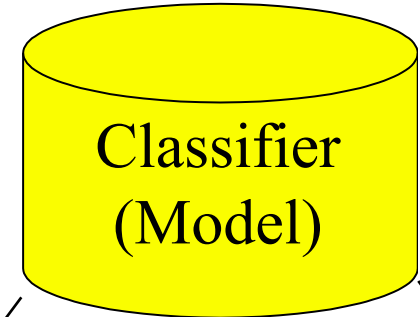


(b)

Process (1): Model Construction



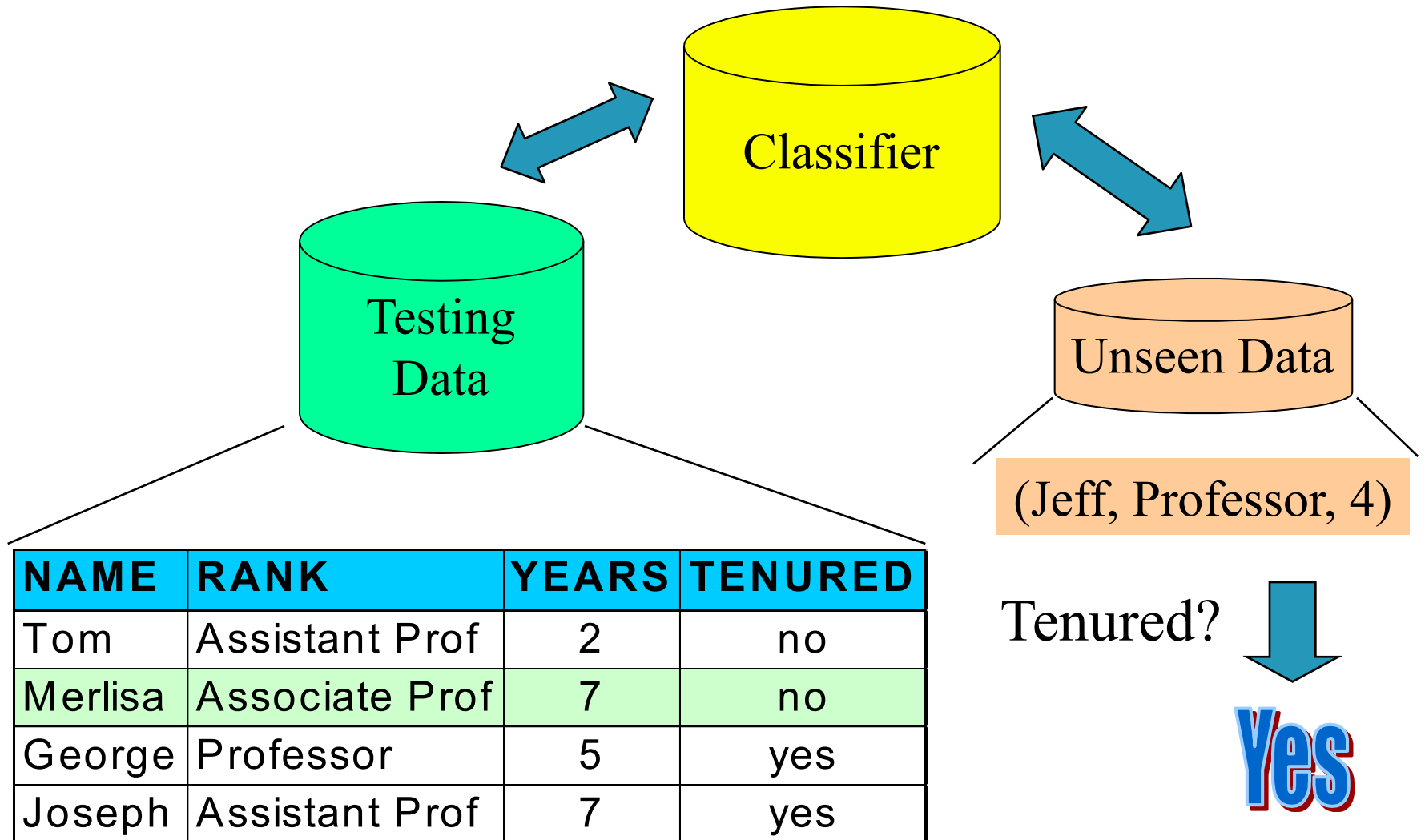
Classification Algorithms



NAME	RANK	YEARS	TENURED
Mike	Assistant Prof	3	no
Mary	Assistant Prof	7	yes
Bill	Professor	2	yes
Jim	Associate Prof	7	yes
Dave	Assistant Prof	6	no
Anne	Associate Prof	3	no

IF rank = 'professor'
OR years > 6
THEN tenured = 'yes'

Process (2): Using the Model in Prediction



Decision Trees

Decision Trees

A general algorithm for decision tree building

- Employs the divide and conquer method
- Recursively divides a training set until each division consists of examples from one class
 1. Create a root node and assign all of the training data to it
 2. Select the best splitting attribute
 3. Add a branch to the root node for each value of the split. Split the data into mutually exclusive subsets along the lines of the specific split
 4. Repeat the steps 2 and 3 for each and every leaf node until the stopping criteria is reached

Decision Trees

- DT algorithms mainly differ on
 - Splitting criteria
 - Which variable to split first?
 - What values to use to split?
 - How many splits to form for each node?
 - Stopping criteria
 - When to stop building the tree
 - Pruning (generalization method)
 - Pre-pruning versus post-pruning
- Most popular DT algorithms include
 - ID3, C4.5, C5; CART; CHAID; M5

Decision Trees

- Alternative splitting criteria
 - **Gini index** determines the purity of a specific class as a result of a decision to branch along a particular attribute/value
 - Used in CART
 - **Information gain** uses entropy to measure the extent of uncertainty or randomness of a particular attribute/value split
 - Used in ID3, C4.5, C5
 - **Chi-square statistics** (used in CHAID)

Classification by Decision Tree Induction

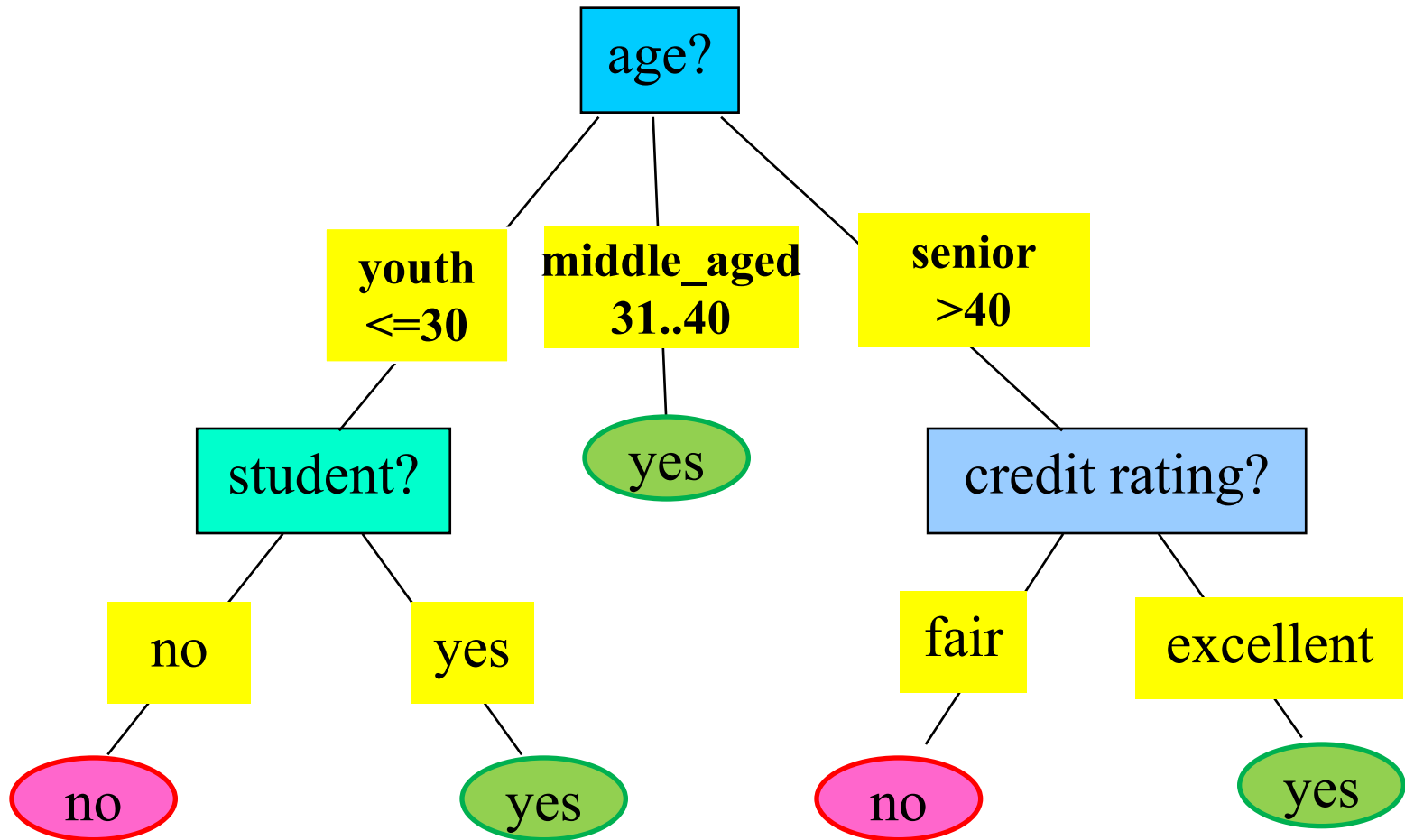
Training Dataset

age	income	student	credit_rating	buys_computer
<=30	high	no	fair	no
<=30	high	no	excellent	no
31...40	high	no	fair	yes
>40	medium	no	fair	yes
>40	low	yes	fair	yes
>40	low	yes	excellent	no
31...40	low	yes	excellent	yes
<=30	medium	no	fair	no
<=30	low	yes	fair	yes
>40	medium	yes	fair	yes
<=30	medium	yes	excellent	yes
31...40	medium	no	excellent	yes
31...40	high	yes	fair	yes
>40	medium	no	excellent	no

This follows an example of Quinlan's ID3 (Playing Tennis)

Classification by Decision Tree Induction

Output: A Decision Tree for “*buys_computer*”

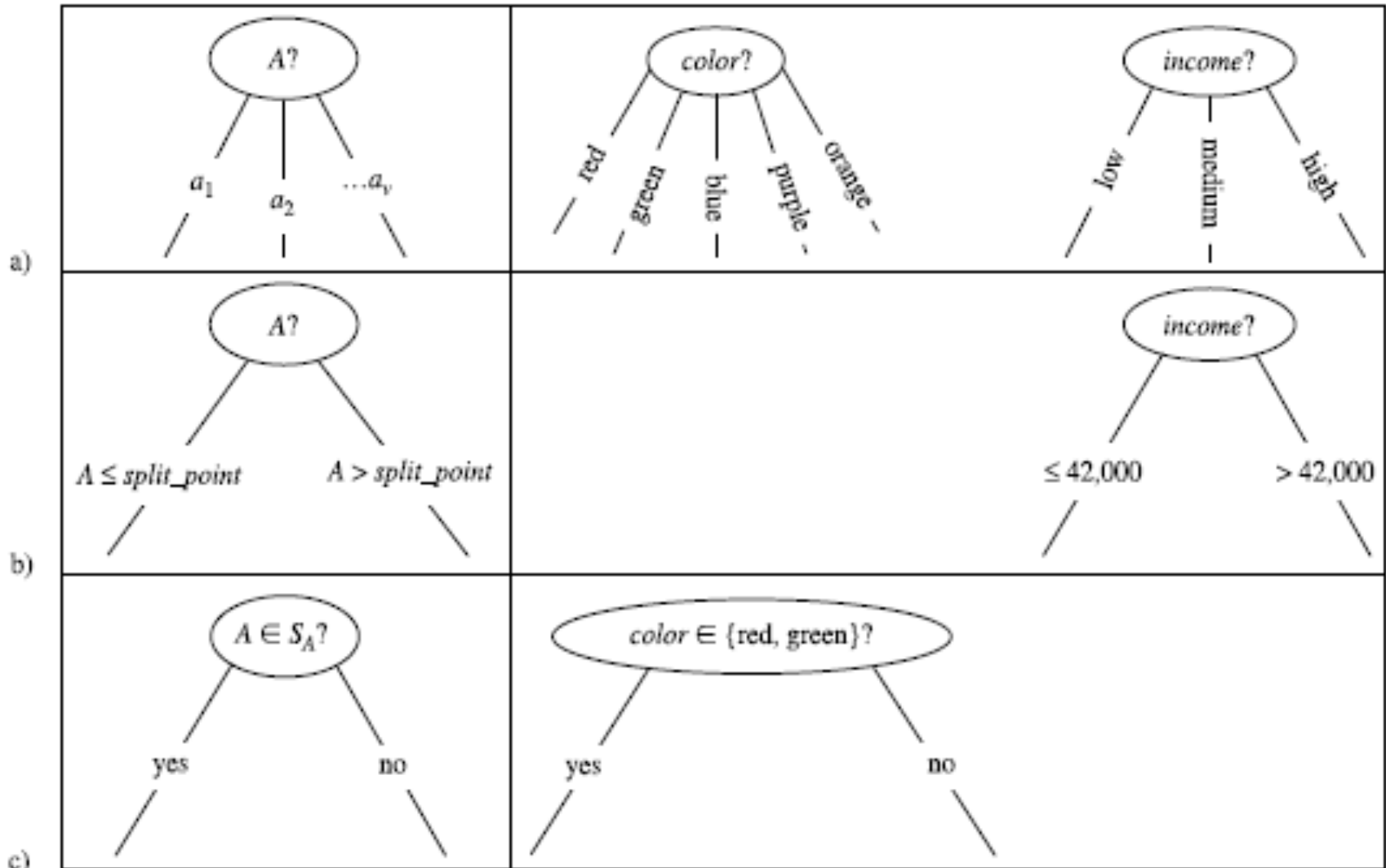


buys_computer="yes" or *buys_computer*="no"

Three possibilities for partitioning tuples based on the splitting Criterion

Partitioning Scenarios

Examples



Algorithm for Decision Tree Induction

- Basic algorithm (a greedy algorithm)
 - Tree is constructed in a **top-down recursive divide-and-conquer manner**
 - At start, all the training examples are at the root
 - Attributes are categorical (if continuous-valued, they are discretized in advance)
 - Examples are partitioned recursively based on selected attributes
 - Test attributes are selected on the basis of a heuristic or statistical measure (e.g., **information gain**)
- Conditions for stopping partitioning
 - All samples for a given node belong to the same class
 - There are no remaining attributes for further partitioning – **majority voting** is employed for classifying the leaf
 - There are no samples left

Attribute Selection Measure

- Notation: Let D , the data partition, be a training set of class-labeled tuples.
Suppose the class label attribute has m distinct values defining m distinct classes, C_i (for $i = 1, \dots, m$).
Let $C_{i,D}$ be the set of tuples of class C_i in D .
Let $|D|$ and $|C_{i,D}|$ denote the number of tuples in D and $C_{i,D}$, respectively.
- Example:
 - Class: `buys_computer` = “yes” or “no”
 - Two distinct classes ($m=2$)
 - Class C_i ($i=1,2$):
 $C_1 = \text{“yes”}$,
 $C_2 = \text{“no”}$

Attribute Selection Measure: Information Gain (ID3/C4.5)

- Select the attribute with the highest information gain
- Let p_i be the probability that an arbitrary tuple in D belongs to class C_i , estimated by $|C_{i,D}|/|D|$

- **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 :

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

- **Information gained** by branching on attribute A

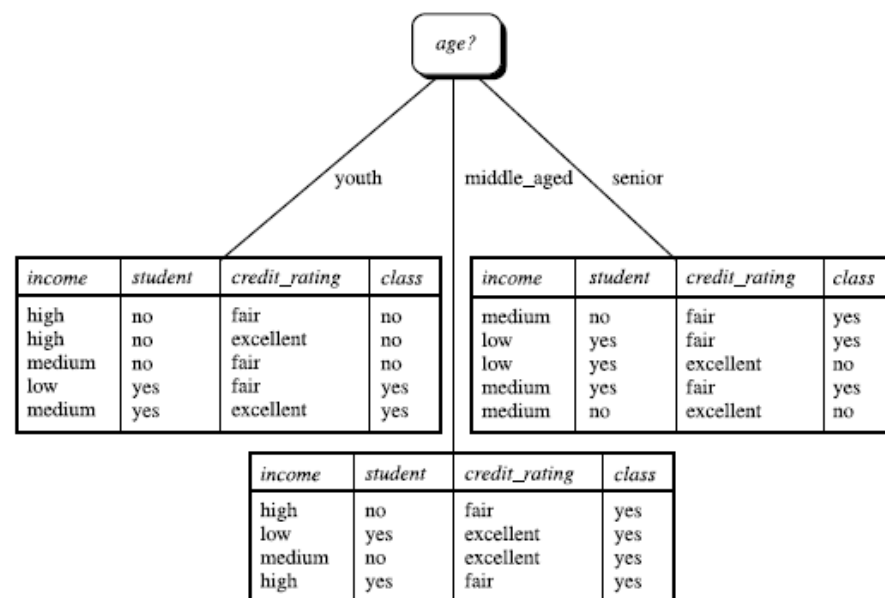
$$Gain(A) = Info(D) - Info_A(D)$$

$$\begin{aligned}\log_2 (1) &= 0 \\ \log_2 (2) &= 1 \\ \log_2 (3) &= 1.5850 \\ \log_2 (4) &= 2 \\ \log_2 (5) &= 2.3219 \\ \log_2 (6) &= 2.5850 \\ \log_2 (7) &= 2.8074 \\ \log_2 (8) &= 3 \\ \log_2 (9) &= 3.1699 \\ \log_2 (10) &= 3.3219\end{aligned}$$

$$\begin{aligned}\log_2 (0.1) &= -3.3219 \\ \log_2 (0.2) &= -2.3219 \\ \log_2 (0.3) &= -1.7370 \\ \log_2 (0.4) &= -1.3219 \\ \log_2 (0.5) &= -1 \\ \log_2 (0.6) &= -0.7370 \\ \log_2 (0.7) &= -0.5146 \\ \log_2 (0.8) &= -0.3219 \\ \log_2 (0.9) &= -0.1520 \\ \log_2 (1) &= 0\end{aligned}$$

Class-labeled training tuples from the *AllElectronics* customer database

RID	age	income	student	credit_rating	Class: buys_computer
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



The attribute age has the highest information gain and therefore becomes the splitting attribute at the root node of the decision tree

Attribute Selection: Information Gain

- Class P: buys_computer = "yes"
- Class N: buys_computer = "no"

$$Info(D) = I(9,5) = -\frac{9}{14} \log_2\left(\frac{9}{14}\right) - \frac{5}{14} \log_2\left(\frac{5}{14}\right) = 0.940$$

$$Info_{age}(D) = \frac{5}{14} I(2,3) + \frac{4}{14} I(4,0) + \frac{5}{14} I(3,2) = 0.694$$

age	p_i	n_i	$I(p_i, n_i)$
≤ 30	2	3	0.971
31...40	4	0	0
> 40	3	2	0.971

$\frac{5}{14} I(2,3)$ means "age ≤ 30 " has 5 out of 14 samples, with 2 yes'es and 3 no's. Hence

$$Gain(age) = Info(D) - Info_{age}(D) = 0.246$$

Similarly,

$$Gain(income) = 0.029$$

$$Gain(student) = 0.151$$

$$Gain(credit_rating) = 0.048$$

age	income	student	credit_rating	buys_computer
≤ 30	high	no	fair	no
≤ 30	high	no	excellent	no
31...40	high	no	fair	yes
> 40	medium	no	fair	yes
> 40	low	yes	fair	yes
> 40	low	yes	excellent	no
31...40	low	yes	excellent	yes
≤ 30	medium	no	fair	no
≤ 30	low	yes	fair	yes
> 40	medium	yes	fair	yes
≤ 30	medium	yes	excellent	yes
31...40	medium	no	excellent	yes
31...40	high	yes	fair	yes
> 40	medium	no	excellent	no

Decision Tree Information Gain

Customer database

ID	age	income	student	credit_rating	Class: buys_computer
1	youth	high	no	fair	no
2	middle_aged	high	no	fair	yes
3	youth	high	no	excellent	no
4	senior	medium	no	fair	yes
5	senior	high	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	excellent	yes

What is the class
(buys_computer = “yes” or
buys_computer = “no”)
for a customer
(age=youth, income=medium,
student =yes, credit= fair)?

Customer database

ID	age	income	student	credit_rating	Class: buys_computer
1	youth	high	no	fair	no
2	middle_aged	high	no	fair	yes
3	youth	high	no	excellent	no
4	senior	medium	no	fair	yes
5	senior	high	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	excellent	yes
11	youth	medium	yes	fair	?

What is the **class**

(**buys_computer = "yes"**) or
buys_computer = "no")

for a **customer**

(age=youth, income=medium,
student =yes, credit= fair)?

Yes = 0.0889

No = 0.0167

Table 1 shows the class-labeled training tuples from customer database. Please calculate and illustrate the final **decision tree** returned by decision tree induction using **information gain**.

- (1) What is the Information Gain of “age”?
- (2) What is the Information Gain of “income”?
- (3) What is the Information Gain of “student”?
- (4) What is the Information Gain of “credit_rating”?
- (5) What is the class (buys_computer = “yes” or buys_computer = “no”) for a customer (age=youth, income=medium, student =yes, credit= fair) based on the classification result by decision tree induction?

ID	age	income	student	credit_rating	Class: buys_computer
1	youth	high	no	fair	no
2	middle_aged	high	no	fair	yes
3	youth	high	no	excellent	no
4	senior	medium	no	fair	yes
5	senior	high	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	excellent	yes

Attribute Selection Measure: Information Gain (ID3/C4.5)

- Select the attribute with the highest information gain
- Let p_i be the probability that an arbitrary tuple in D belongs to class C_i , estimated by $|C_{i,D}|/|D|$

- **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 :

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

- **Information gained** by branching on attribute A

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

$$\begin{aligned}\log_2 (1) &= 0 \\ \log_2 (2) &= 1 \\ \log_2 (3) &= 1.5850 \\ \log_2 (4) &= 2 \\ \log_2 (5) &= 2.3219 \\ \log_2 (6) &= 2.5850 \\ \log_2 (7) &= 2.8074 \\ \log_2 (8) &= 3 \\ \log_2 (9) &= 3.1699 \\ \log_2 (10) &= 3.3219\end{aligned}$$

$$\begin{aligned}\log_2 (0.1) &= -3.3219 \\ \log_2 (0.2) &= -2.3219 \\ \log_2 (0.3) &= -1.7370 \\ \log_2 (0.4) &= -1.3219 \\ \log_2 (0.5) &= -1 \\ \log_2 (0.6) &= -0.7370 \\ \log_2 (0.7) &= -0.5146 \\ \log_2 (0.8) &= -0.3219 \\ \log_2 (0.9) &= -0.1520 \\ \log_2 (1) &= 0\end{aligned}$$

ID	age	income	student	credit rating	Class: buys_computer
1	youth	high	no	fair	no
2	middle_aged	high	no	fair	yes
3	youth	high	no	excellent	no
4	senior	medium	no	fair	yes
5	senior	high	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	excellent	yes

Class P (Positive): buys_computer = “yes”

Class N (Negative): buys_computer = “no”

$$P(\text{buys} = \text{yes}) = P_{i=1} = P_1 = 6/10 = 0.6$$

$$P(\text{buys} = \text{no}) = P_{i=2} = P_2 = 4/10 = 0.4$$

$$\log_2(0.1) = -3.3219$$

$$\log_2(0.2) = -2.3219$$

$$\log_2(0.3) = -1.7370$$

$$\log_2(0.4) = -1.3219$$

$$\log_2(0.5) = -1$$

$$\log_2(0.6) = -0.7370$$

$$\log_2(0.7) = -0.5146$$

$$\log_2(0.8) = -0.3219$$

$$\log_2(0.9) = -0.1520$$

$$\log_2(1) = 0$$

$$\log_2(1) = 0$$

$$\log_2(2) = 1$$

$$\log_2(3) = 1.5850$$

$$\log_2(4) = 2$$

$$\log_2(5) = 2.3219$$

$$\log_2(6) = 2.5850$$

$$\log_2(7) = 2.8074$$

$$\log_2(8) = 3$$

$$\log_2(9) = 3.1699$$

$$\log_2(10) = 3.3219$$

Step 1: Expected information

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

$$Info(D) = I(6,4) = -\frac{6}{10} \log_2\left(\frac{6}{10}\right) + \left(-\frac{4}{10} \log_2\left(\frac{4}{10}\right)\right)$$

$$= -0.6 \times \log_2(0.6) - 0.4 \times \log_2(0.4)$$

$$= -0.6 \times (-0.737) - 0.4 \times (-1.3219)$$

$$= 0.4422 + 0.5288$$

$$= 0.971$$

$$Info(D) = I(6,4) = 0.971$$

ID	age	income	student	credit_rating	Class: buys_computer
1	youth	high	no	fair	no
2	middle_aged	high	no	fair	yes
3	youth	high	no	excellent	no
4	senior	medium	no	fair	yes
5	senior	high	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	excellent	yes

<i>age</i>	p_i	n_i	<i>total</i>
youth	1	3	4
middle_aged	2	0	2
senior	3	1	4

<i>income</i>	p_i	n_i	<i>total</i>
high	2	2	4
medium	2	1	3
low	2	1	3

<i>student</i>	p_i	n_i	<i>total</i>
yes	4	1	5
no	2	3	5

<i>credit_rating</i>	p_i	n_i	<i>total</i>
excellent	2	2	4
fair	4	2	6

<i>age</i>	p_i	n_i	<i>total</i>	$I(p_i, n_i)$	$I(p_i, n_i)$
youth	1	3	4	$I(1,3)$	0.8112
middle_aged	2	0	2	$I(2,0)$	0
senior	3	1	4	$I(3,1)$	0.8112

Step 2: Information

Step 3: Information Gain

$$\begin{aligned}
 I(1,3) &= -\frac{1}{4} \log_2\left(\frac{1}{4}\right) + \left(-\frac{3}{4} \log_2\left(\frac{3}{4}\right)\right) \\
 &= -0.25 \times [\log_2 1 - \log_2 4] + (-0.75 \times [\log_2 3 - \log_2 4]) \\
 &= -0.25 \times [0 - 2] - 0.75 \times [1.585 - 2] \\
 &= -0.25 \times [-2] - 0.75 \times [-0.415] \\
 &= 0.5 + 0.3112 = 0.8112
 \end{aligned}$$

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

$\text{Info}(D) = I(6,4) = 0.971$

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

$$\begin{aligned}
 I(2,0) &= -\frac{2}{2} \log_2\left(\frac{2}{2}\right) + \left(-\frac{0}{2} \log_2\left(\frac{0}{2}\right)\right) \\
 &= -1 \times \log_2 1 + (-0 \times \log_2 0) \\
 &= -1 \times 0 + (-0 \times -\infty) \\
 &= 0 + 0 = 0
 \end{aligned}$$

$$\begin{aligned}
 \text{Info}_{age}(D) &= \frac{4}{10} I(1,3) + \frac{2}{10} I(2,0) + \frac{4}{10} I(3,1) \\
 &= \frac{4}{10} \times 0.8112 + \frac{2}{10} \times 0 + \frac{4}{10} \times 0.8112 \\
 &= 0.3244 + 0 + 0.3244 = 0.6488
 \end{aligned}$$

$$\begin{aligned}
 I(3,1) &= -\frac{3}{4} \log_2\left(\frac{3}{4}\right) + \left(-\frac{1}{4} \log_2\left(\frac{1}{4}\right)\right) \\
 &= -0.75 \times [\log_2 3 - \log_2 4] + (-0.25 \times [\log_2 1 - \log_2 4]) \\
 &= -0.75 \times [1.585 - 2] - 0.25 \times [0 - 2] \\
 &= -0.75 \times [-0.415] - 0.25 \times [-2] \\
 &= 0.3112 + 0.5 = 0.8112
 \end{aligned}$$

$$\text{Gain}(A) = \text{Info}(D) - \text{Info}_A(D)$$

$$\begin{aligned}
 \text{Gain}(age) &= \text{Info}(D) - \text{Info}_{age}(D) \\
 &= 0.971 - 0.6488 = 0.3221
 \end{aligned}$$

(1) Gain(age) = 0.3221

<i>income</i>	p_i	n_i	<i>total</i>	$I(p_i, n_i)$	$I(p_i, n_i)$
high	2	2	4	$I(2,2)$	1
medium	2	1	3	$I(2,1)$	0.9182
low	2	1	3	$I(2,1)$	0.9182

$$\begin{aligned}
 I(2,2) &= -\frac{2}{4} \log_2\left(\frac{2}{4}\right) + \left(-\frac{2}{4} \log_2\left(\frac{2}{4}\right)\right) \\
 &= -0.5 \times [\log_2 2 - \log_2 4] + (-0.5 \times [\log_2 2 - \log_2 4]) \\
 &= -0.5 \times [1 - 2] - 0.5 \times [1 - 2] \\
 &= -0.5 \times [-1] - 0.5 \times [-1] \\
 &= 0.5 + 0.5 = 1
 \end{aligned}$$

$$\text{Info}(D) = -\sum_{i=1}^m p_i \log_2(p_i) \quad \text{Info}(D) = I(6,4) = 0.971$$

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

$$\begin{aligned}
 I(2,1) &= -\frac{2}{3} \log_2\left(\frac{2}{3}\right) + \left(-\frac{1}{3} \log_2\left(\frac{1}{3}\right)\right) \\
 &= -0.67 \times [\log_2 2 - \log_2 3] + (-0.33 \times [\log_2 1 - \log_2 3]) \\
 &= -0.67 \times [1 - 1.585] - 0.33 \times [0 - 1.585] \\
 &= -0.67 \times [-0.585] - 0.33 \times [-1.585] \\
 &= 0.9182
 \end{aligned}$$

$$\begin{aligned}
 \text{Info}_{\text{income}}(D) &= \frac{4}{10} I(2,2) + \frac{3}{10} I(2,1) + \frac{3}{10} I(2,1) \\
 &= \frac{4}{10} \times 1 + \frac{3}{10} \times 0.9182 + \frac{3}{10} \times 0.9182 \\
 &= 0.4 + 0.2755 + 0.2755 = 0.951
 \end{aligned}$$

$$\text{Gain}(A) = \text{Info}(D) - \text{Info}_A(D)$$

$$\begin{aligned}
 \text{Gain}(\text{income}) &= \text{Info}(D) - \text{Info}_{\text{income}}(D) \\
 &= 0.971 - 0.951 = 0.02
 \end{aligned}$$

(2) Gain(income) = 0.02

<i>student</i>	p_i	n_i	<i>total</i>	$I(p_i, n_i)$	$I(p_i, n_i)$
yes	4	1	5	$I(4,1)$	0.7219
no	2	3	5	$I(2,3)$	0.971

$$\begin{aligned}
 I(4,1) &= -\frac{4}{5} \log_2\left(\frac{4}{5}\right) + \left(-\frac{1}{5} \log_2\left(\frac{1}{5}\right)\right) \\
 &= -0.8 \times [\log_2 4 - \log_2 5] + (-0.2 \times [\log_2 1 - \log_2 5]) \\
 &= -0.8 \times [2 - 2.3219] - 0.2 \times [0 - 2.3219] \\
 &= -0.8 \times [-0.3219] - 0.2 \times [-2.3219] \\
 &= 0.25752 + 0.46438 = 0.7219
 \end{aligned}$$

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

$\text{Info}(D) = I(6,4) = 0.971$

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

$$\begin{aligned}
 I(2,3) &= -\frac{2}{5} \log_2\left(\frac{2}{5}\right) + \left(-\frac{3}{5} \log_2\left(\frac{3}{5}\right)\right) \\
 &= -0.4 \times [\log_2 0.4] + (-0.6 \times [\log_2 0.6]) \\
 &= -0.4 \times [-1.3219] - 0.6 \times [-0.737] \\
 &= 0.5288 + 0.4422 = 0.971
 \end{aligned}$$

$$\begin{aligned}
 \text{Info}_{\text{student}}(D) &= \frac{5}{10} I(4,1) + \frac{5}{10} I(2,3) \\
 &= 0.5 \times 0.7219 + 0.5 \times 0.971 \\
 &= 0.36095 + 0.48545 = 0.8464
 \end{aligned}$$

$$\text{Gain}(A) = \text{Info}(D) - \text{Info}_A(D)$$

$$\begin{aligned}
 \text{Gain}(\text{student}) &= \text{Info}(D) - \text{Info}_{\text{student}}(D) \\
 &= 0.971 - 0.8464 = 0.1245
 \end{aligned}$$

(3) Gain_(student) = 0.1245

<i>credit</i>	p_i	n_i	<i>total</i>	$I(p_i, n_i)$	$I(p_i, n_i)$
excellent	2	2	4	$I(2,2)$	1
fair	4	2	6	$I(4,2)$	0.9183

$$\begin{aligned}
I(2,2) &= -\frac{2}{4} \log_2\left(\frac{2}{4}\right) + \left(-\frac{2}{4} \log_2\left(\frac{2}{4}\right)\right) \\
&= -0.5 \times [\log_2 2 - \log_2 4] + (-0.5 \times [\log_2 2 - \log_2 4]) \\
&= -0.5 \times [1 - 2] - 0.5 \times [1 - 2] \\
&= -0.5 \times [-1] - 0.5 \times [-1] \\
&= 0.5 + 0.5 = 1
\end{aligned}$$

$$\text{Info}(D) = -\sum_{i=1}^m p_i \log_2(p_i) \quad \text{Info}(D) = I(6,4) = 0.971$$

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

$$\begin{aligned}
I(4,2) &= -\frac{4}{6} \log_2\left(\frac{4}{6}\right) + \left(-\frac{2}{6} \log_2\left(\frac{2}{6}\right)\right) \\
&= -0.67 \times [\log_2 2 - \log_2 3] + (-0.33 \times [\log_2 1 - \log_2 3]) \\
&= -0.67 \times [1 - 1.585] - 0.33 \times [0 - 1.585] \\
&= -0.67 \times [-0.585] - 0.33 \times [-1.585] \\
&= 0.9182
\end{aligned}$$

$$\begin{aligned}
\text{Info}_{\text{credit}}(D) &= \frac{4}{10} I(2,2) + \frac{6}{10} I(4,2) \\
&= \frac{4}{10} \times 1 + \frac{6}{10} \times 0.9182 \\
&= 0.4 + 0.5509 = 0.9509
\end{aligned}$$

$$\text{Gain}(A) = \text{Info}(D) - \text{Info}_A(D)$$

$$\begin{aligned}
\text{Gain}(\text{credit}) &= \text{Info}(D) - \text{Info}_{\text{credit}}(D) \\
&= 0.971 - 0.9509 = 0.019
\end{aligned}$$

(4) Gain_(credit) = 0.019

What is the class
(buys_computer = “yes” or
buys_computer = “no”)
for a customer
(age=youth, income=medium,
student =yes, credit= fair)?

<i>age</i>	<i>p_i</i>	<i>n_i</i>	<i>total</i>
youth	1	3	4
middle_aged	2	0	2
senior	3	1	4

<i>student</i>	<i>p_i</i>	<i>n_i</i>	<i>total</i>
yes	4	1	5
no	2	3	5

<i>income</i>	<i>p_i</i>	<i>n_i</i>	<i>total</i>
high	2	2	4
midium	2	1	3
low	2	1	3

<i>credit_rating</i>	<i>p_i</i>	<i>n_i</i>	<i>total</i>
excellent	2	2	4
fair	4	2	6

(5) What is the class (buys_computer = “yes” or buys_computer = “no”) for a customer (age=youth, income=medium, student =yes, credit= fair) based on the classification result by decision three induction?

(5) Yes =0.0889 (No=0.0167)

age (0.3221) > student (0.1245) > income (0.02) > credit (0.019)

buys_computer = “yes”

age:youth (1/4) x student:yes (4/5) x income:medium (2/3) x credit:fair (4/6)

Yes: $1/4 \times 4/5 \times 2/3 \times 4/6 = 4/45 = 0.0889$

buys_computer = “no”

age:youth (3/4) x student:yes (1/5) x income:medium (1/3) x credit:fair (2/6)

No: $3/4 \times 1/5 \times 1/3 \times 2/6 = 0.01667$

What is the **class**

(**buys_computer = "yes"**) or
buys_computer = "no")

for a **customer**

(age=youth, income=medium,
student =yes, credit= fair)?

Yes = 0.0889

No = 0.0167

Customer database

ID	age	income	student	credit_rating	Class: buys_computer
1	youth	high	no	fair	no
2	middle_aged	high	no	fair	yes
3	youth	high	no	excellent	no
4	senior	medium	no	fair	yes
5	senior	high	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	excellent	yes

Customer database

ID	age	income	student	credit_rating	Class: buys_computer
1	youth	high	no	fair	no
2	middle_aged	high	no	fair	yes
3	youth	high	no	excellent	no
4	senior	medium	no	fair	yes
5	senior	high	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	excellent	yes
11	youth	medium	yes	fair	?

Customer database

ID	age	income	student	credit_rating	Class: buys_computer
1	youth	high	no	fair	no
2	middle_aged	high	no	fair	yes
3	youth	high	no	excellent	no
4	senior	medium	no	fair	yes
5	senior	high	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	excellent	yes
11	youth	medium	yes	fair	Yes (0.0889)

Support Vector Machines (SVM)

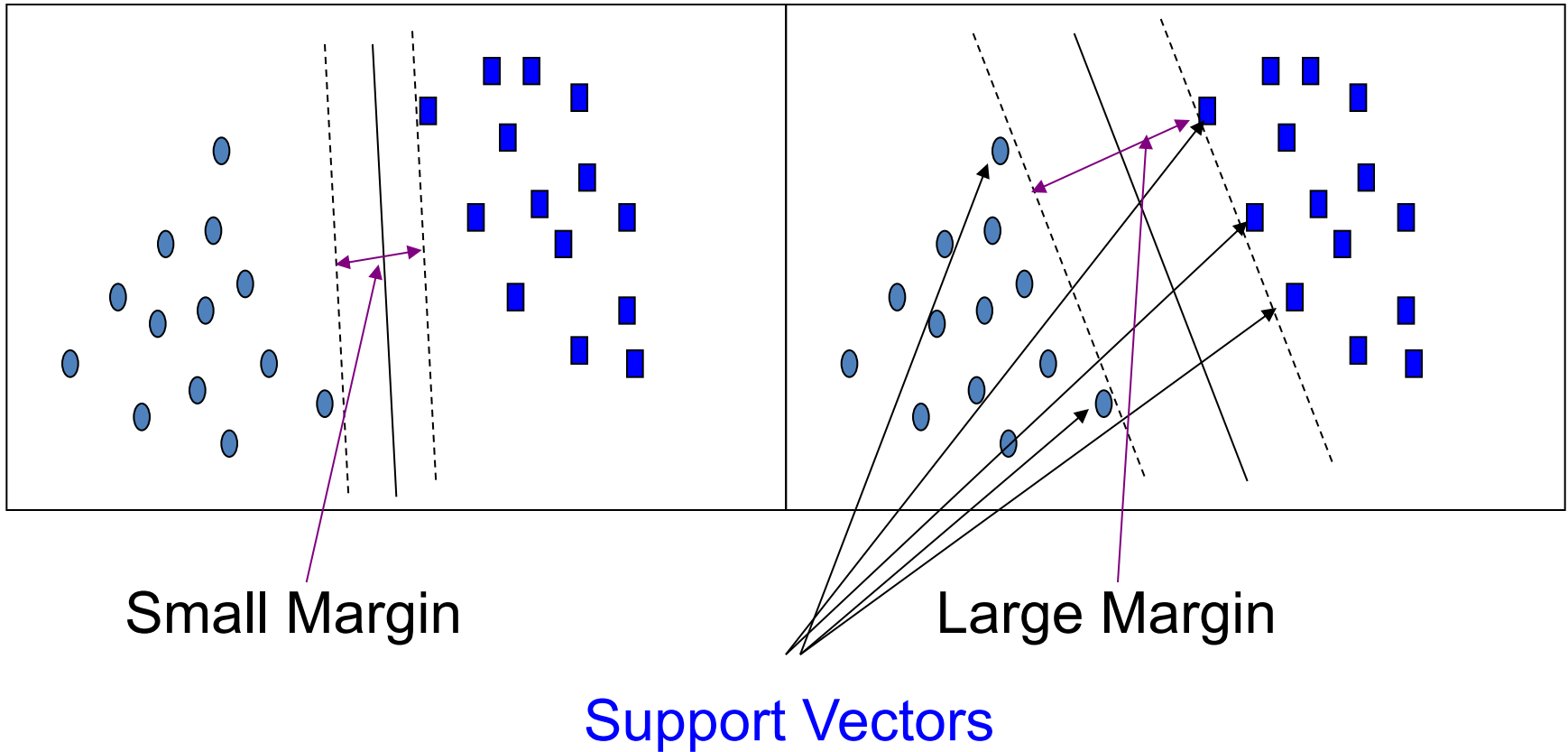
SVM—Support Vector Machines

- A new classification method for both linear and 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) and margins (defined by the support vectors)

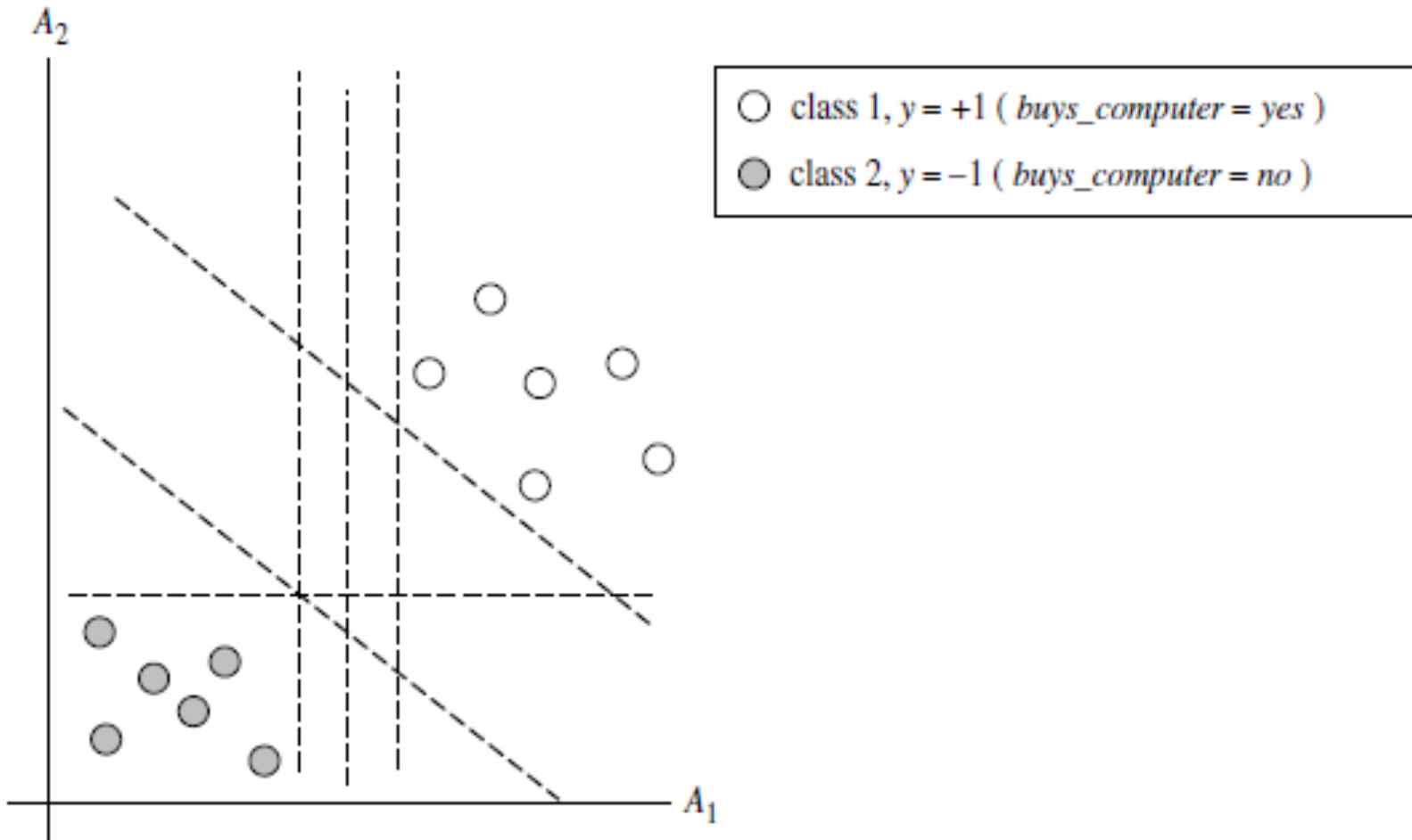
SVM—History and Applications

- 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 both for classification and prediction
- Applications:
 - handwritten digit recognition, object recognition, speaker identification, benchmarking time-series prediction tests, document classification

SVM—General Philosophy

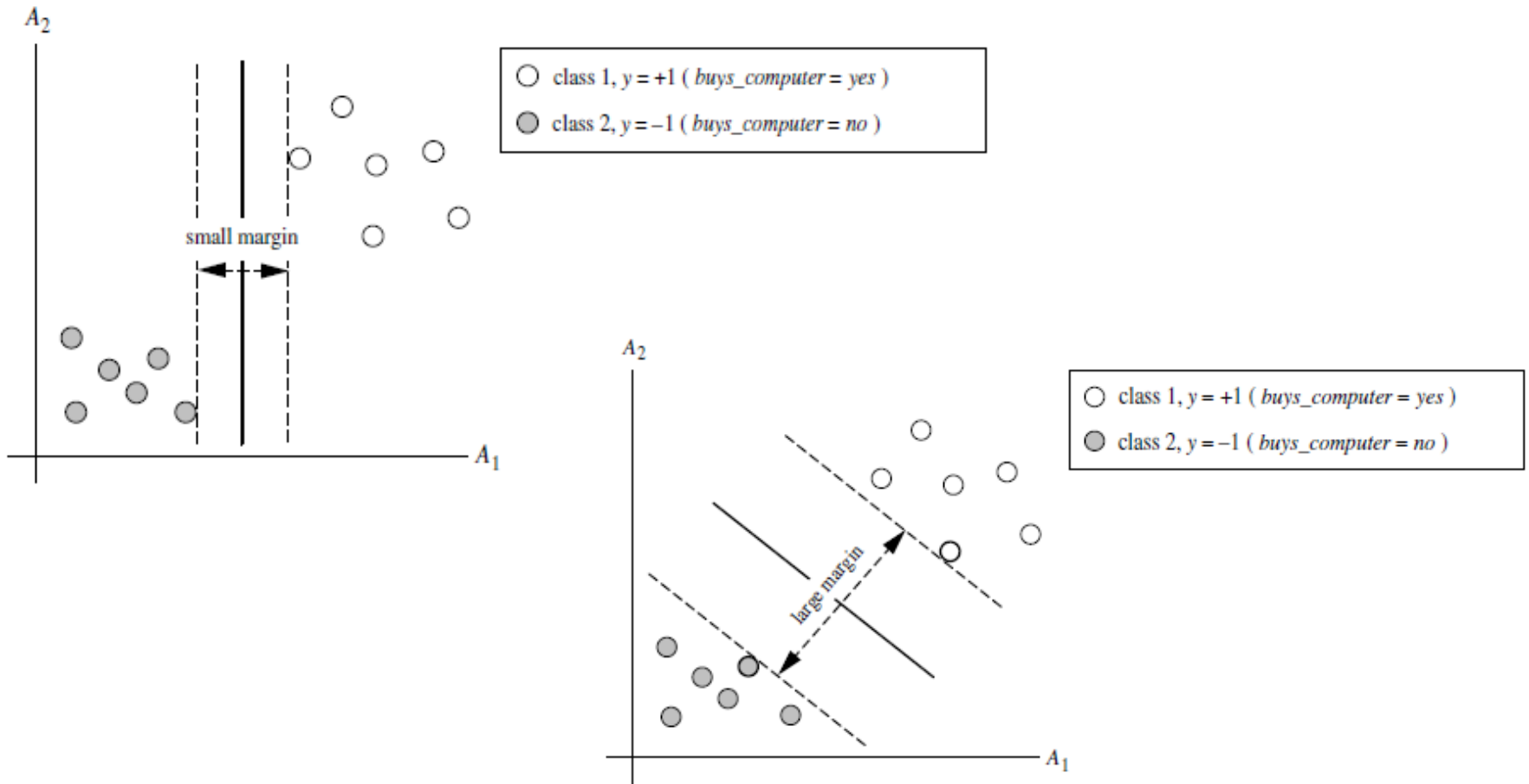


Classification (SVM)



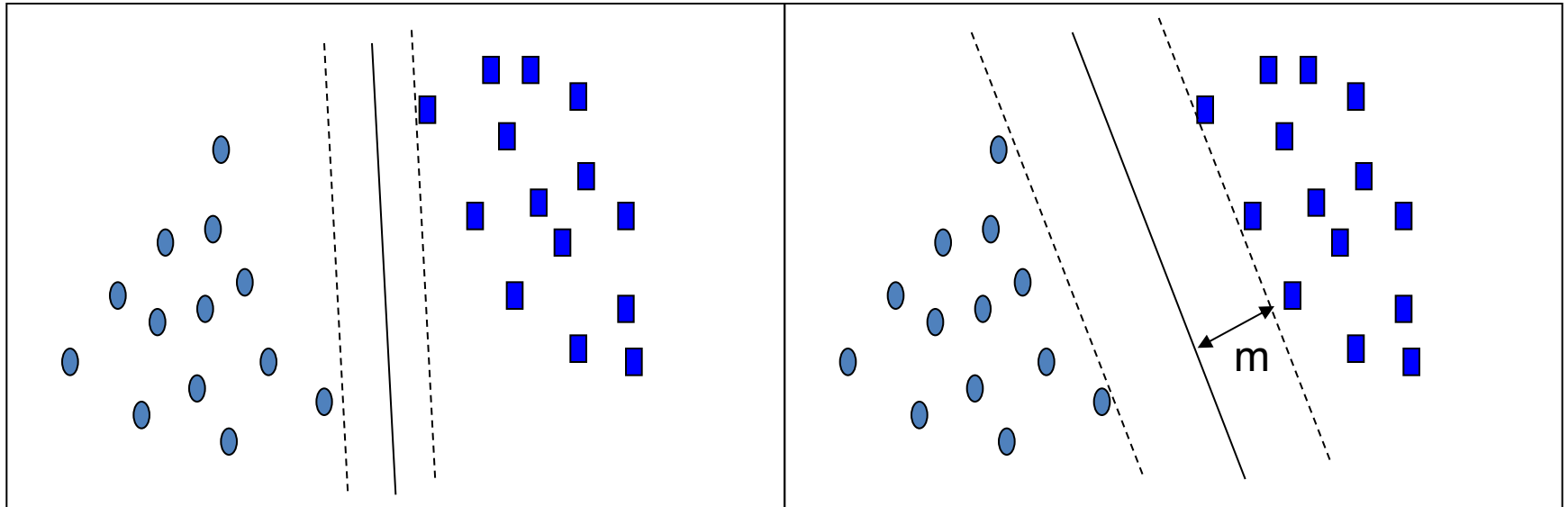
The 2-D training data are linearly separable. There are an infinite number of (possible) separating hyperplanes or “decision boundaries.” Which one is best?

Classification (SVM)



Which one is better? The one with the larger margin should have greater generalization accuracy.

SVM—When Data Is Linearly Separable



Let data D be $(\mathbf{X}_1, y_1), \dots, (\mathbf{X}_{|D|}, y_{|D|})$, where \mathbf{X}_i is the set of training tuples associated with the class labels y_i

There are infinite lines (hyperplanes) separating the two classes but we want to find the best one (the one that minimizes classification error on unseen data)

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

SVM—Linearly Separable

- A separating hyperplane can be written as

$$\mathbf{W} \bullet \mathbf{X} + b = 0$$

where $\mathbf{W} = \{w_1, w_2, \dots, w_n\}$ is a weight vector and b a scalar (bias)

- For 2-D it can be written as

$$w_0 + w_1 x_1 + w_2 x_2 = 0$$

- The hyperplane defining the sides of the margin:

$$H_1: w_0 + w_1 x_1 + w_2 x_2 \geq 1 \quad \text{for } y_i = +1, \text{ and}$$

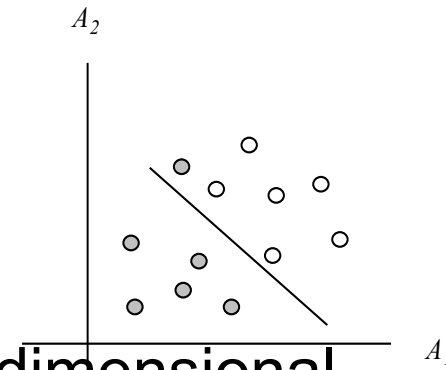
$$H_2: w_0 + w_1 x_1 + w_2 x_2 \leq -1 \quad \text{for } y_i = -1$$

- Any training tuples that fall on hyperplanes H_1 or H_2 (i.e., the sides defining the margin) are **support vectors**
- This becomes a **constrained (convex) quadratic optimization** problem: Quadratic objective function and linear constraints \rightarrow *Quadratic Programming (QP)* \rightarrow Lagrangian multipliers

Why Is SVM Effective on High Dimensional Data?

- The complexity of trained classifier is characterized by the # of support vectors rather than the dimensionality of the data
- The support vectors are the essential or critical training examples — they lie closest to the decision boundary (MMH)
- If all other training examples are removed and the training is repeated, the same separating hyperplane would be found
- The number of support vectors found can be used to compute an (upper) bound on the expected error rate of the SVM classifier, which is independent of the data dimensionality
- Thus, an SVM with a small number of support vectors can have good generalization, even when the dimensionality of the data is high

SVM—Linearly Inseparable



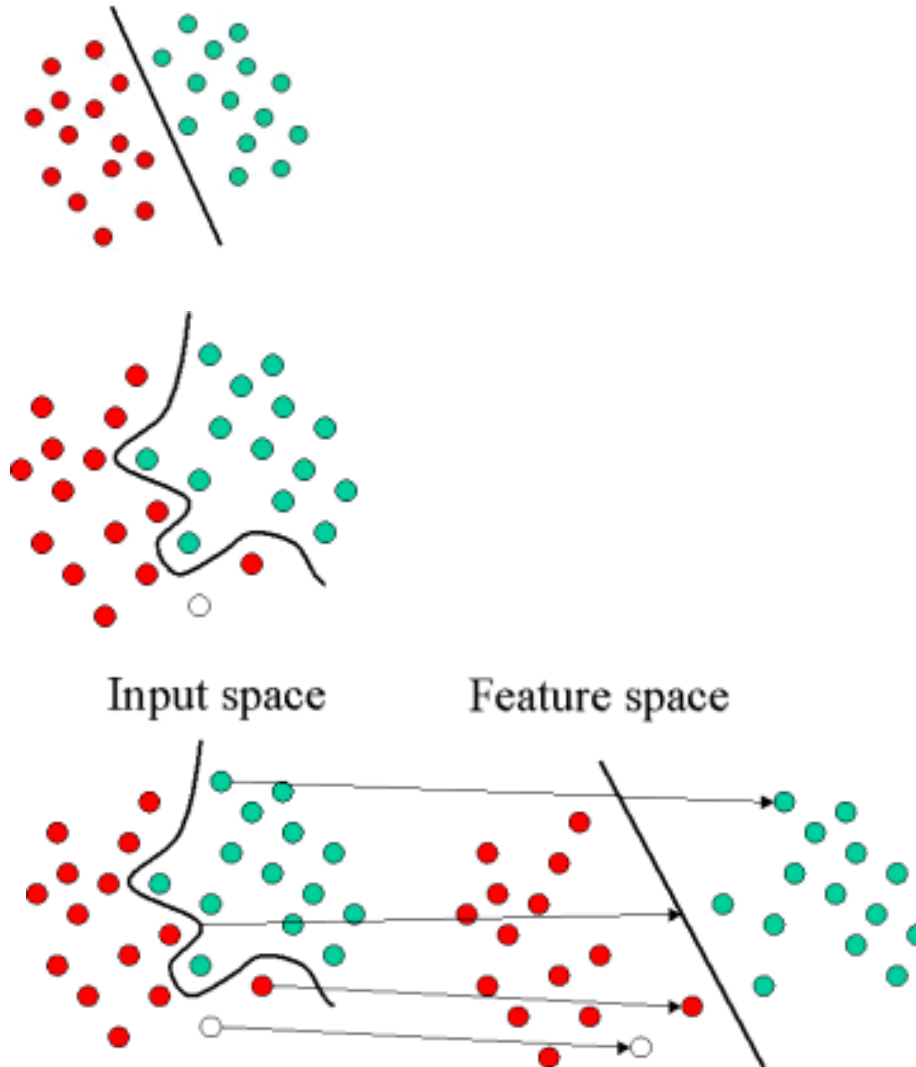
- Transform the original input data into a higher dimensional space

Example 6.8 Nonlinear transformation of original input data into a higher dimensional space. Consider the following example. A 3D input vector $\mathbf{X} = (x_1, x_2, x_3)$ is mapped into a 6D space Z using the mappings $\phi_1(\mathbf{X}) = x_1, \phi_2(\mathbf{X}) = x_2, \phi_3(\mathbf{X}) = x_3, \phi_4(\mathbf{X}) = (x_1)^2, \phi_5(\mathbf{X}) = x_1x_2$, and $\phi_6(\mathbf{X}) = x_1x_3$. A decision hyperplane in the new space is $d(\mathbf{Z}) = \mathbf{WZ} + b$, where \mathbf{W} and \mathbf{Z} are vectors. This is linear. We solve for \mathbf{W} and b and then substitute back so that we see that the linear decision hyperplane in the new (\mathbf{Z}) space corresponds to a nonlinear second order polynomial in the original 3-D input space,

$$\begin{aligned} d(\mathbf{Z}) &= w_1x_1 + w_2x_2 + w_3x_3 + w_4(x_1)^2 + w_5x_1x_2 + w_6x_1x_3 + b \\ &= w_1z_1 + w_2z_2 + w_3z_3 + w_4z_4 + w_5z_5 + w_6z_6 + b \quad \blacksquare \end{aligned}$$

- Search for a linear separating hyperplane in the new space

Mapping Input Space to Feature Space



SVM—Kernel functions

- Instead of computing the dot product on the transformed data tuples, it is mathematically equivalent to instead applying a kernel function $K(\mathbf{X}_i, \mathbf{X}_j)$ to the original data, i.e., $K(\mathbf{X}_i, \mathbf{X}_j) = \Phi(\mathbf{X}_i) \cdot \Phi(\mathbf{X}_j)$
- Typical Kernel Functions

Polynomial kernel of degree h : $K(\mathbf{X}_i, \mathbf{X}_j) = (\mathbf{X}_i \cdot \mathbf{X}_j + 1)^h$

Gaussian radial basis function kernel : $K(\mathbf{X}_i, \mathbf{X}_j) = e^{-\|\mathbf{X}_i - \mathbf{X}_j\|^2 / 2\sigma^2}$

Sigmoid kernel : $K(\mathbf{X}_i, \mathbf{X}_j) = \tanh(\kappa \mathbf{X}_i \cdot \mathbf{X}_j - \delta)$

- SVM can also be used for classifying multiple (> 2) classes and for regression analysis (with additional user parameters)

SVM Related Links

- SVM Website
 - <http://www.kernel-machines.org/>
- Representative implementations
 - **LIBSVM**
 - an efficient implementation of SVM, multi-class classifications, nu-SVM, one-class SVM, including also various interfaces with java, python, etc.
 - SVM-light
 - simpler but performance is not better than LIBSVM, support only binary classification and only C language
 - SVM-torch
 - another recent implementation also written in C.

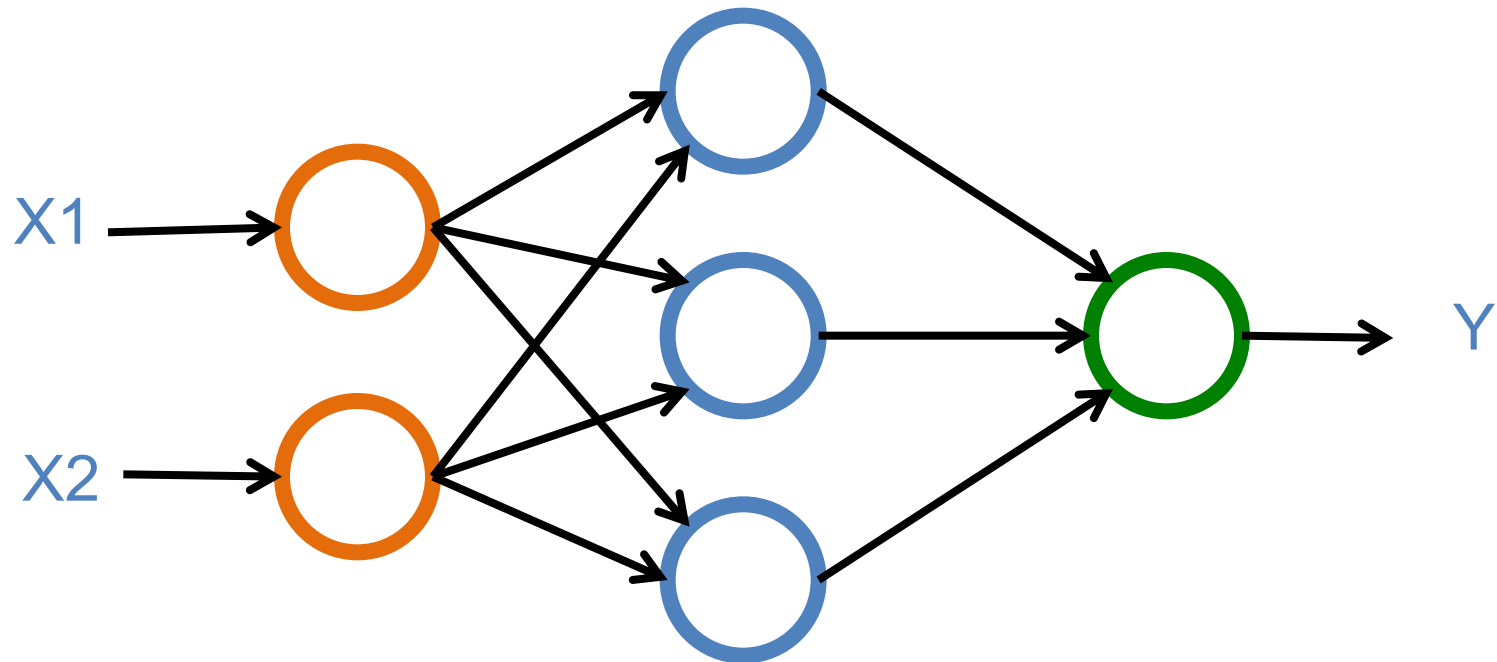
Deep Learning and Neural Networks

Deep Learning and Neural Networks

Input Layer
(X)

Hidden Layer
(H)

Output Layer
(Y)

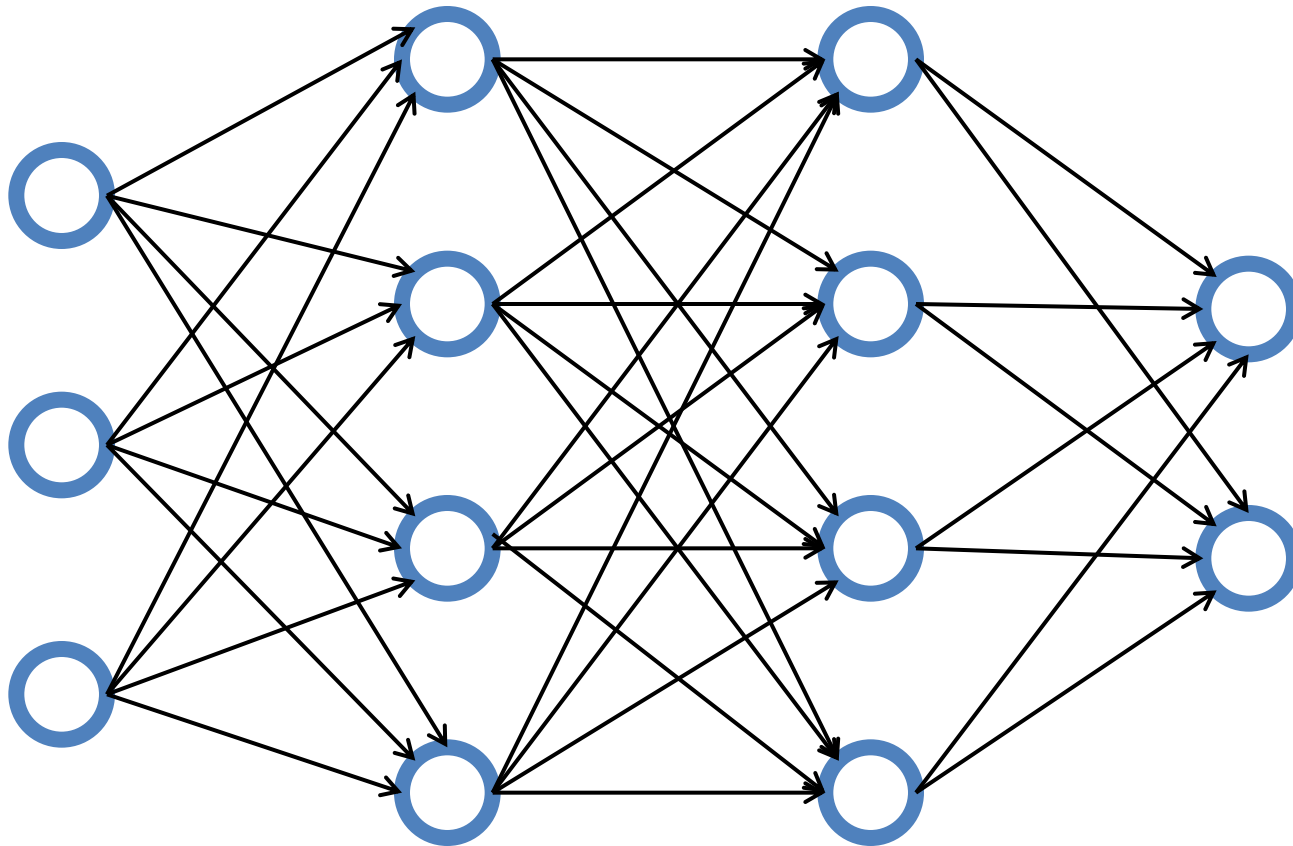


Deep Learning and Neural Networks

Input Layer
(X)

Hidden Layer
(H)

Output Layer
(Y)



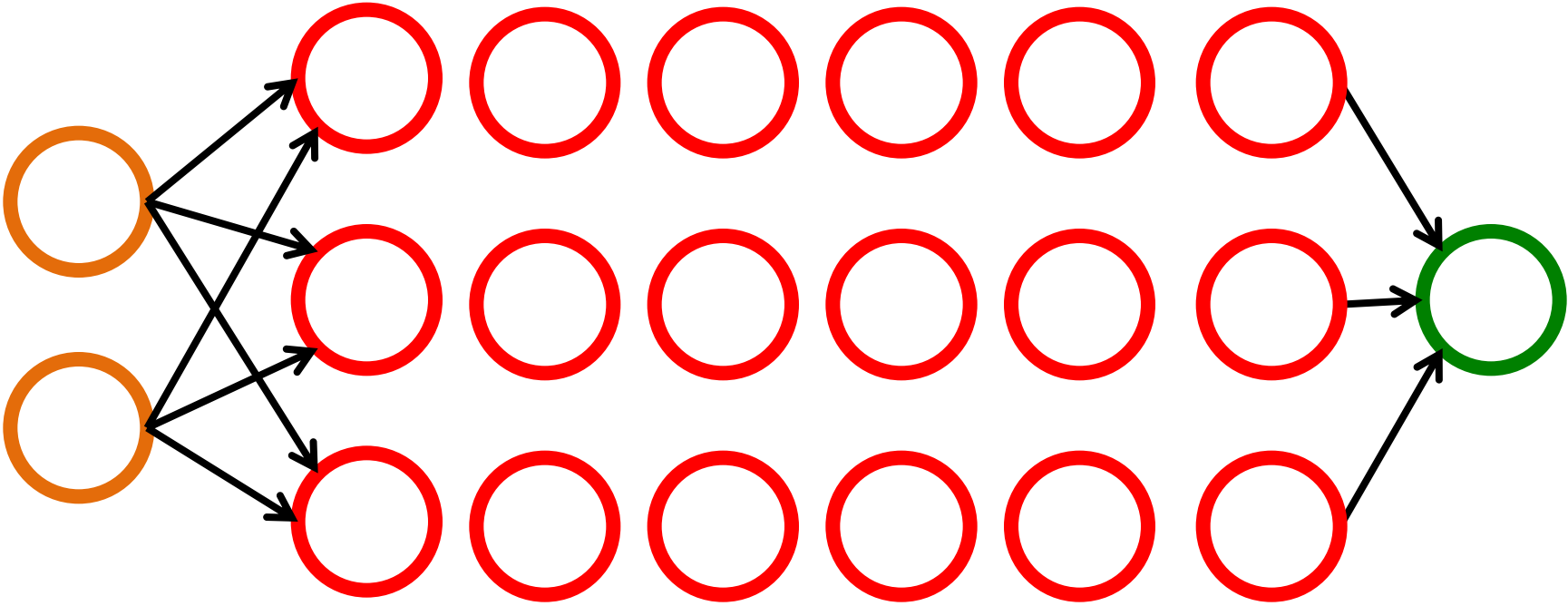
Deep Learning and Neural Networks

Input Layer
(X)

Hidden Layers
(H)

Output Layer
(Y)

Deep Neural Networks
Deep Learning



Data Mining Evaluation

Evaluation

(Accuracy of Classification Model)

Assessment Methods for Classification

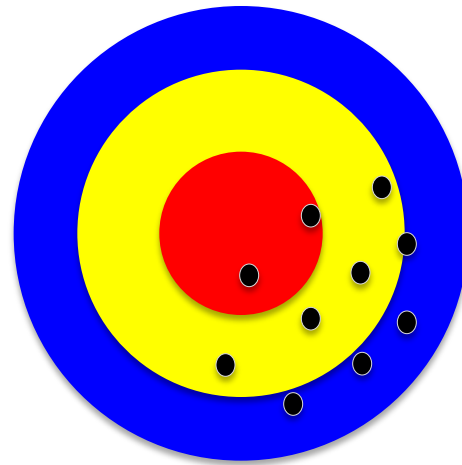
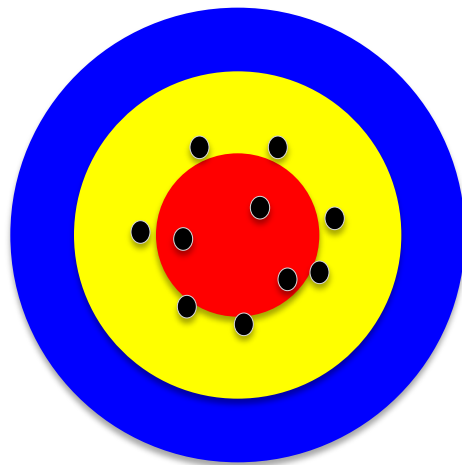
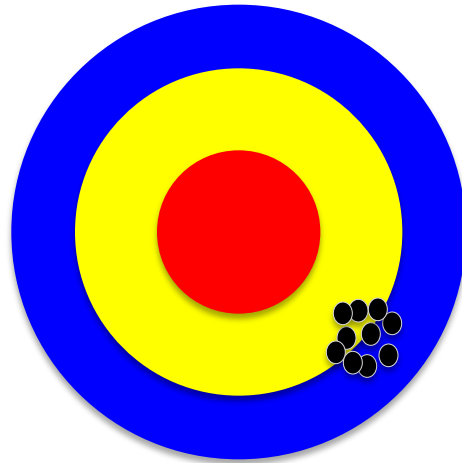
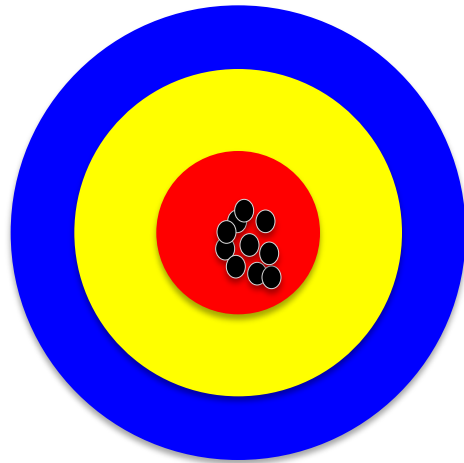
- Predictive accuracy
 - Hit rate
- Speed
 - Model building; predicting
- Robustness
- Scalability
- Interpretability
 - Transparency, explainability

Accuracy

Validity

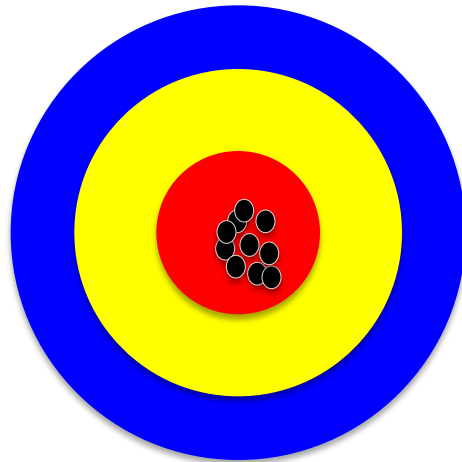
Precision

Reliability



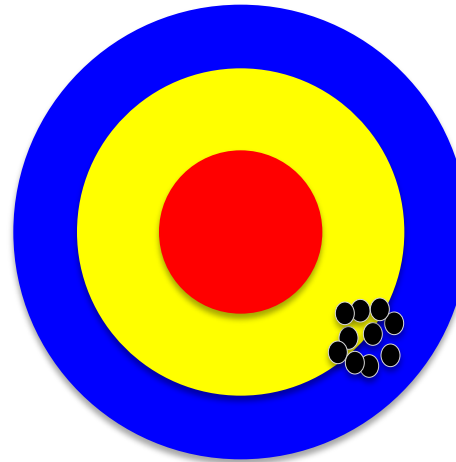
Accuracy vs. Precision

A



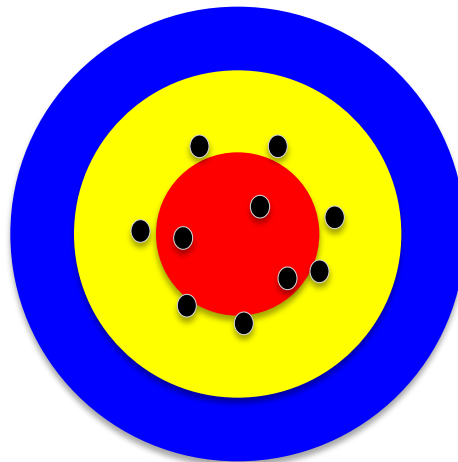
**High Accuracy
High Precision**

B



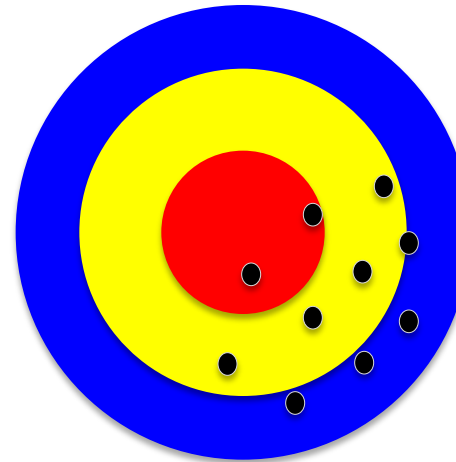
**Low Accuracy
High Precision**

C



**High Accuracy
Low Precision**

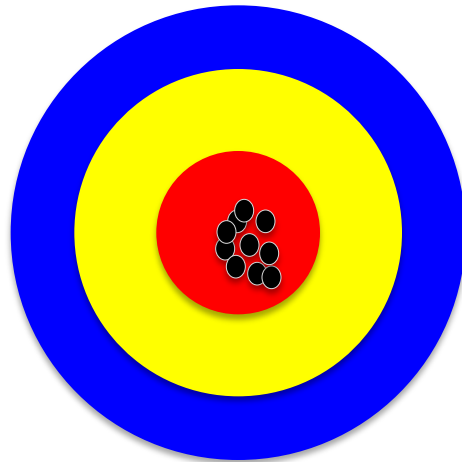
D



**Low Accuracy
Low Precision**

Accuracy vs. Precision

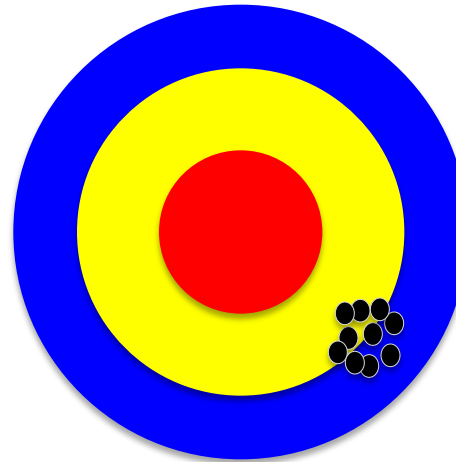
A



**High Accuracy
High Precision**

**High Validity
High Reliability**

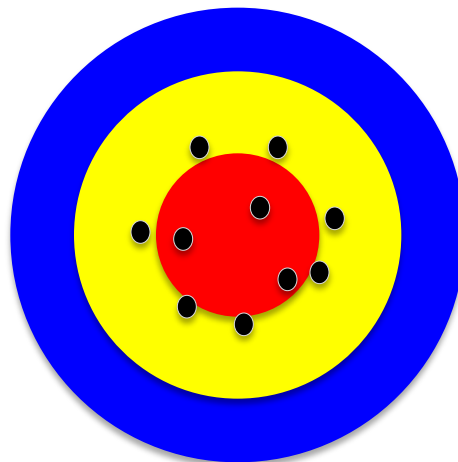
B



**Low Accuracy
High Precision**

**Low Validity
High Reliability**

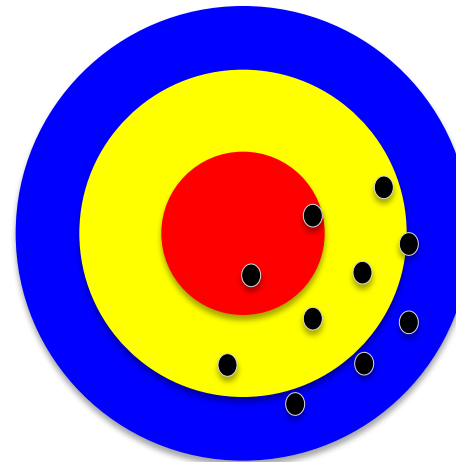
C



**High Accuracy
Low Precision**

**High Validity
Low Reliability**

D

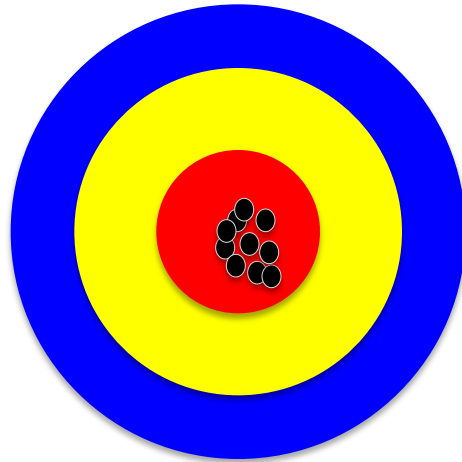


**Low Accuracy
Low Precision**

**Low Validity
Low Reliability**

Accuracy vs. Precision

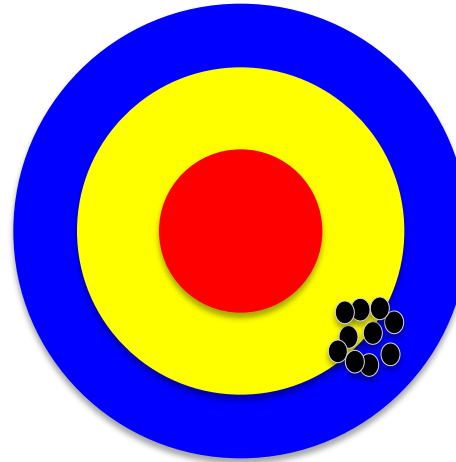
A



High Accuracy
High Precision

High Validity
High Reliability

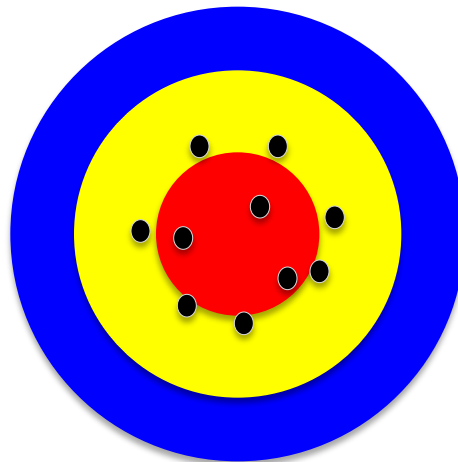
B



Low Accuracy
High Precision

Low Validity
High Reliability

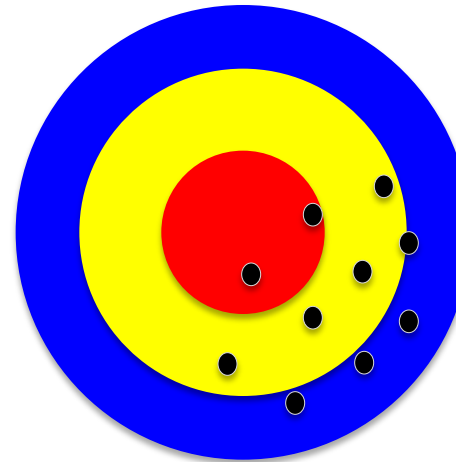
C



High Accuracy
Low Precision

High Validity
Low Reliability

D



Low Accuracy
Low Precision

Low Validity
Low Reliability

Accuracy of Classification Models

- In classification problems, the primary source for accuracy estimation is the **confusion matrix**

		True Class	
		Positive	Negative
Predicted Class	Positive	True Positive Count (TP)	False Positive Count (FP)
	Negative	False Negative Count (FN)	True Negative Count (TN)

$$Accuracy = \frac{TP + TN}{TP + TN + FP + FN}$$

$$True\ Positive\ Rate = \frac{TP}{TP + FN}$$

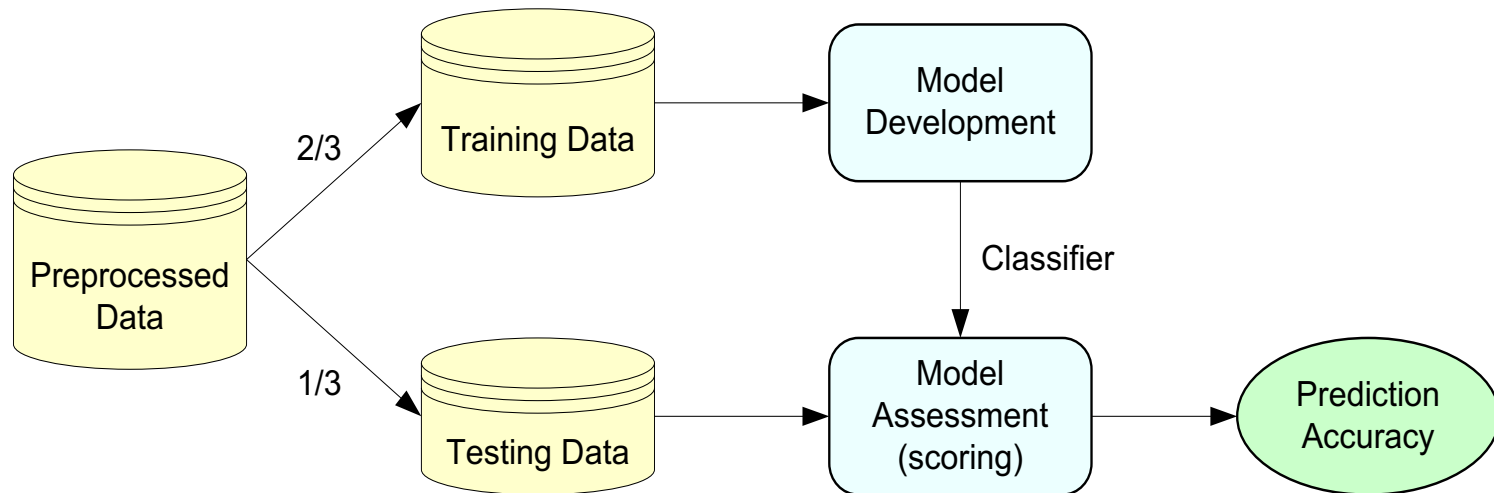
$$True\ Negative\ Rate = \frac{TN}{TN + FP}$$

$$Precision = \frac{TP}{TP + FP}$$

$$Recall = \frac{TP}{TP + FN}$$

Estimation Methodologies for Classification

- **Simple split** (or holdout or test sample estimation)
 - Split the data into 2 mutually exclusive sets training (~70%) and testing (30%)

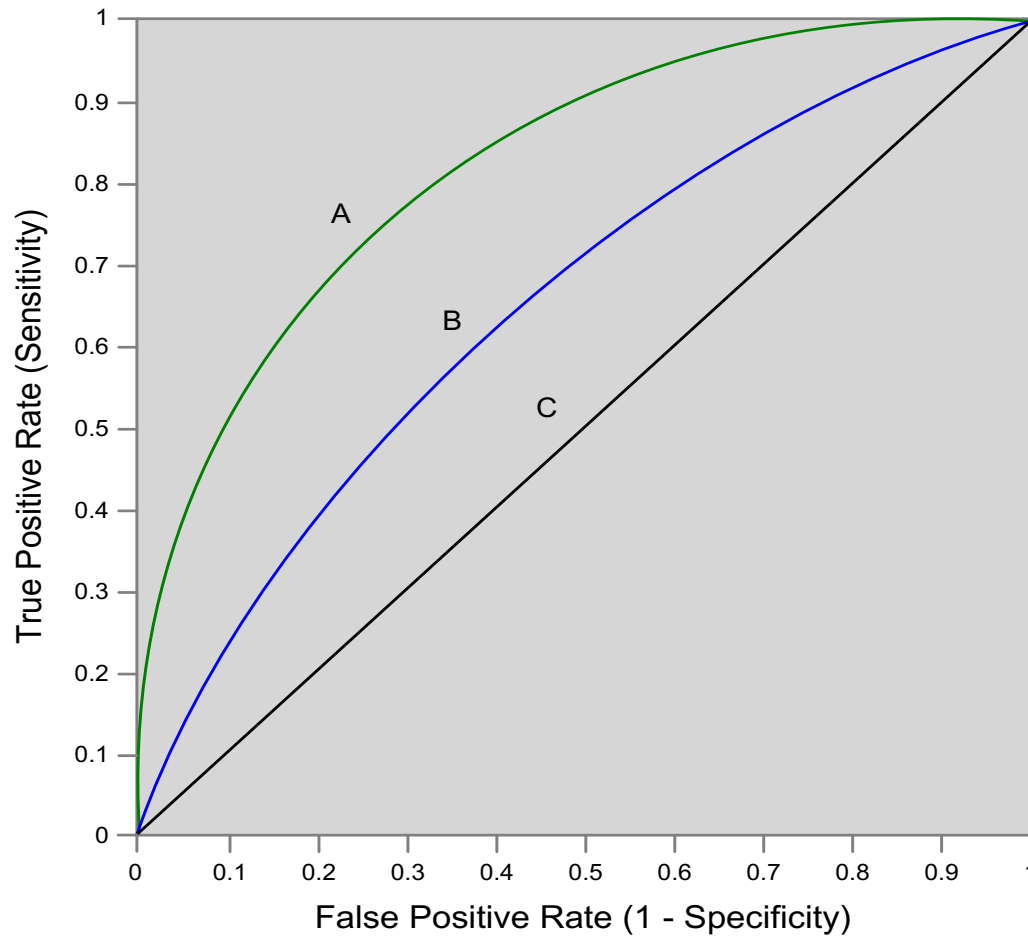


- For ANN, the data is split into three sub-sets (training [~60%], validation [~20%], testing [~20%])

Estimation Methodologies for Classification

- ***k*-Fold Cross Validation** (rotation estimation)
 - Split the data into k mutually exclusive subsets
 - Use each subset as testing while using the rest of the subsets as training
 - Repeat the experimentation for k times
 - Aggregate the test results for true estimation of prediction accuracy training
- Other estimation methodologies
 - Leave-one-out, bootstrapping, jackknifing
 - Area under the ROC curve

Estimation Methodologies for Classification – ROC Curve



Sensitivity = True Positive Rate

Specificity = True Negative Rate

		True Class (actual value)		total
		Positive	Negative	
Predictive Class (prediction outcome)	Positive	True Positive (TP)	False Positive (FP)	P'
	Negative	False Negative (FN)	True Negative (TN)	N'
total		P	N	

$$Accuracy = \frac{TP + TN}{TP + TN + FP + FN}$$

$$True\ Positive\ Rate = \frac{TP}{TP + FN}$$

$$True\ Negative\ Rate = \frac{TN}{TN + FP}$$

$$Precision = \frac{TP}{TP + FP}$$

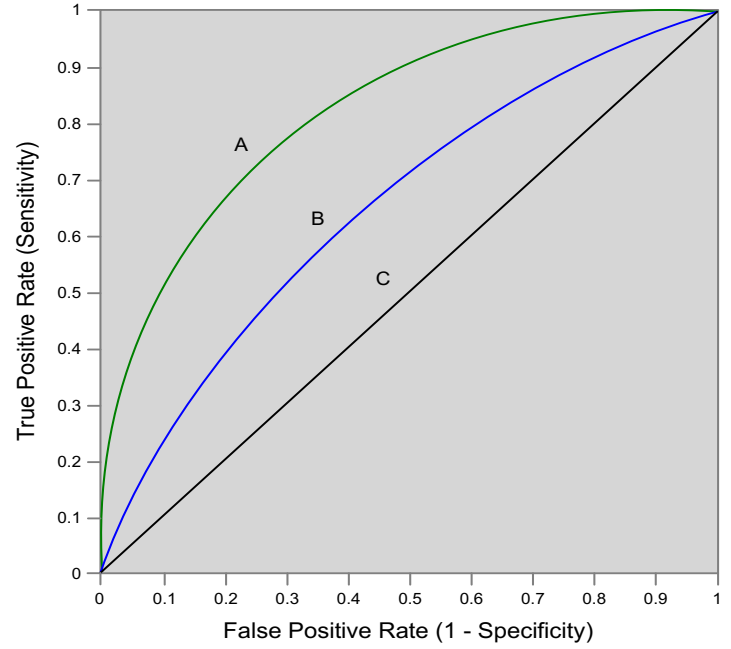
$$Recall = \frac{TP}{TP + FN}$$

$$True\ Positive\ Rate\ (Sensitivity) = \frac{TP}{TP + FN}$$

$$True\ Negative\ Rate\ (Specificity) = \frac{TN}{TN + FP}$$

$$False\ Positive\ Rate = \frac{FP}{FP + TN}$$

$$False\ Positive\ Rate\ (1 - Specificity) = \frac{FP}{FP + TN}$$



Source: http://en.wikipedia.org/wiki/Receiver_operating_characteristic

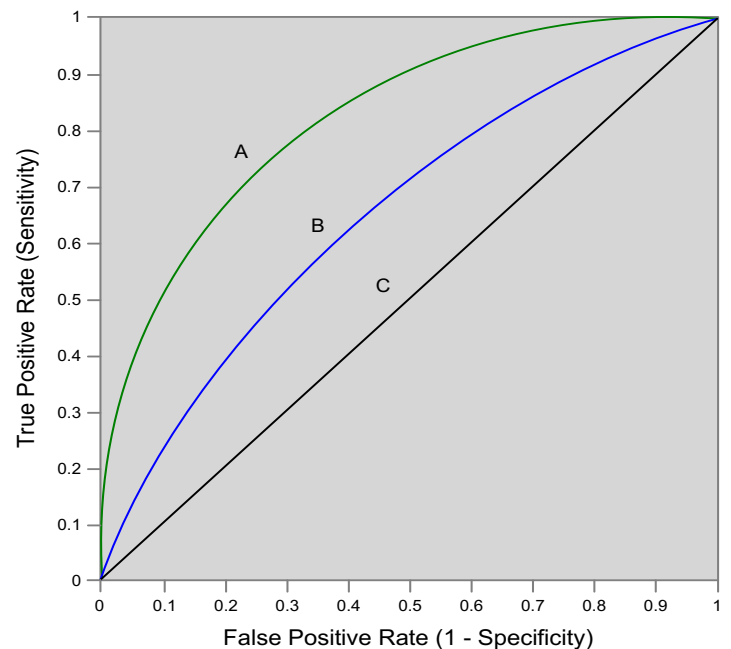
		True Class (actual value)		total
		Positive	Negative	
Predictive Class (prediction outcome)	Positive	True Positive (TP)	False Positive (FP)	P'
	Negative	False Negative (FN)	True Negative (TN)	N'
total		P	N	

$$\text{True Positive Rate} = \frac{TP}{TP + FN}$$

$$\text{Recall} = \frac{TP}{TP + FN}$$

$$\text{True Positive Rate (Sensitivity)} = \frac{TP}{TP + FN}$$

- Sensitivity**
- = True Positive Rate
- = Recall
- = Hit rate
- = $TP / (TP + FN)$



		True Class (actual value)		total
		Positive	Negative	
Predictive Class (prediction outcome)	Positive	True Positive (TP)	False Positive (FP)	P'
	Negative	False Negative (FN)	True Negative (TN)	
total		P	N	

$$\text{True Negative Rate} = \frac{TN}{TN + FP}$$

Specificity

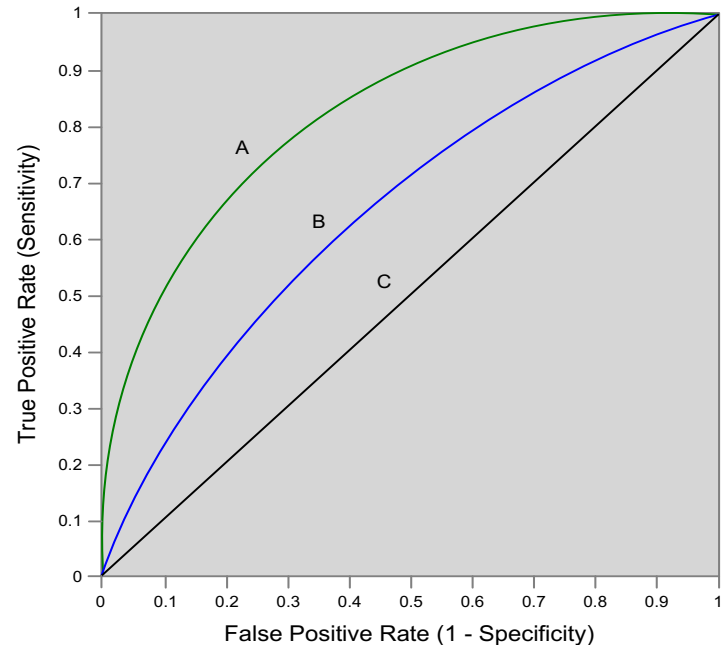
= True Negative Rate

= TN / N

= $TN / (TN + FP)$

$$\text{True Negative Rate (Specificity)} = \frac{TN}{TN + FP}$$

$$\text{False Positive Rate (1-Specificity)} = \frac{FP}{FP + TN}$$



		True Class (actual value)		total
		Positive	Negative	
Predictive Class (prediction outcome)	Positive	True Positive (TP)	False Positive (FP)	P'
	Negative	False Negative (FN)	True Negative (TN)	N'
total		P	N	

Precision

= Positive Predictive Value (PPV)

$$Precision = \frac{TP}{TP + FP}$$

Recall

= True Positive Rate (TPR)

= Sensitivity

= Hit Rate

$$Recall = \frac{TP}{TP + FN}$$

F1 score (F-score)(F-measure)

is the harmonic mean of precision and recall

$$= 2TP / (P + P')$$

$$= 2TP / (2TP + FP + FN)$$

$$F = 2 * \frac{precision * recall}{precision + recall}$$

A		
63 (TP)	28 (FP)	91
37 (FN)	72 (TN)	109
100	100	200

Recall

= True Positive Rate (TPR)
 = Sensitivity
 = Hit Rate
 = $TP / (TP + FN)$

Specificity

= True Negative Rate
 = TN / N
 = $TN / (TN + FP)$

$$TPR = 0.63$$

$$Recall = \frac{TP}{TP + FN}$$

$$True\ Negative\ Rate\ (Specificity) = \frac{TN}{TN + FP}$$

$$FPR = 0.28$$

$$False\ Positive\ Rate\ (1 - Specificity) = \frac{FP}{FP + TN}$$

$$PPV = 0.69$$

$$= 63 / (63 + 28)$$

$$= 63 / 91$$

$$Precision = \frac{TP}{TP + FP}$$

Precision

= Positive Predictive Value (PPV)

$$F1 = 0.66$$

$$= 2 * (0.63 * 0.69) / (0.63 + 0.69)$$

$$= (2 * 63) / (100 + 91)$$

$$= (0.63 + 0.69) / 2 = 1.32 / 2 = 0.66$$

$$F = 2 * \frac{precision * recall}{precision + recall}$$

F1 score (F-score) (F-measure)

is the harmonic mean of precision and recall

$$= 2TP / (P + P')$$

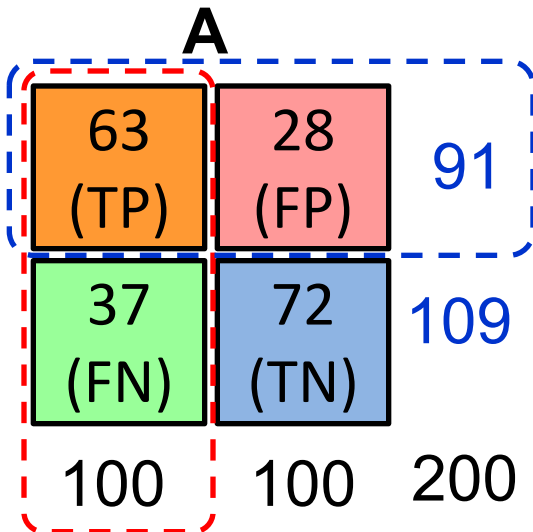
$$= 2TP / (2TP + FP + FN)$$

$$ACC = 0.68$$

$$= (63 + 72) / 200$$

$$= 135 / 200 = 67.5$$

$$Accuracy = \frac{TP + TN}{TP + TN + FP + FN}$$



$$\text{TPR} = 0.63$$

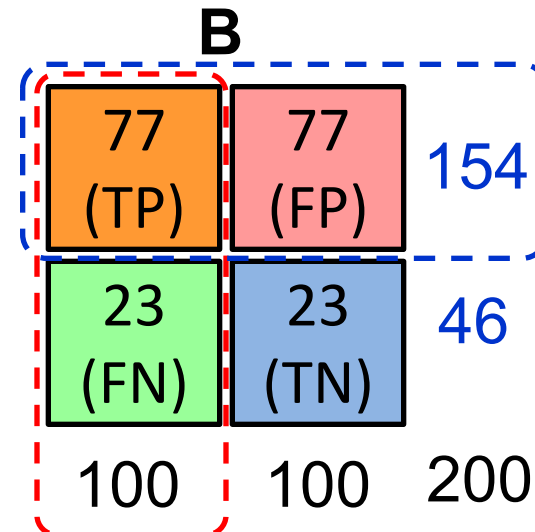
$$\text{FPR} = 0.28$$

$$\begin{aligned} \text{PPV} &= 0.69 \\ &= 63 / (63 + 28) \\ &= 63 / 91 \end{aligned}$$

$$\text{F1} = 0.66$$

$$\begin{aligned} &= 2 * (0.63 * 0.69) / (0.63 + 0.69) \\ &= (2 * 63) / (100 + 91) \\ &= (0.63 + 0.69) / 2 = 1.32 / 2 = 0.66 \end{aligned}$$

$$\begin{aligned} \text{ACC} &= 0.68 \\ &= (63 + 72) / 200 \\ &= 135 / 200 = 67.5 \end{aligned}$$



$$\text{TPR} = 0.77$$

$$\text{FPR} = 0.77$$

$$\text{PPV} = 0.50$$

$$\text{F1} = 0.61$$

$$\text{ACC} = 0.50$$

Recall

= True Positive Rate (TPR)

= Sensitivity

= Hit Rate

$$\text{Recall} = \frac{TP}{TP + FN}$$

Precision

= Positive Predictive Value (PPV)

$$\text{Precision} = \frac{TP}{TP + FP}$$

C

24 (TP)	88 (FP)	112
76 (FN)	12 (TN)	88
100	100	200

$$\text{TPR} = 0.24$$

$$\text{FPR} = 0.88$$

$$\text{PPV} = 0.21$$

$$\text{F1} = 0.22$$

$$\text{ACC} = 0.18$$

C'

76 (TP)	12 (FP)	88
24 (FN)	88 (TN)	112
100	100	200

$$\text{TPR} = 0.76$$

$$\text{FPR} = 0.12$$

$$\text{PPV} = 0.86$$

$$\text{F1} = 0.81$$

$$\text{ACC} = 0.82$$

Recall
 = True Positive Rate (TPR) $\text{Recall} = \frac{TP}{TP + FN}$
 = Sensitivity
 = Hit Rate

Precision
 = Positive Predictive Value (PPV) $\text{Precision} = \frac{TP}{TP + FP}$

Summary

- Classification and Prediction
- Supervised Learning (Classification)
- Decision Tree (DT)
 - Information Gain (IG)
- Support Vector Machine (SVM)
- Data Mining Evaluation
 - Accuracy
 - Precision
 - Recall
 - F1 score (F-measure) (F-score)

References

- Jiawei Han and Micheline Kamber, Data Mining: Concepts and Techniques, Second Edition, Elsevier, 2006.
- Jiawei Han, Micheline Kamber and Jian Pei, Data Mining: Concepts and Techniques, Third Edition, Morgan Kaufmann 2011.
- Efraim Turban, Ramesh Sharda, Dursun Delen, Decision Support and Business Intelligence Systems, Ninth Edition, Pearson, 2011.