

# Introduction to Data Mining

Yücel SAYGIN

[ysaygin@sabanciuniv.edu](mailto:ysaygin@sabanciuniv.edu)

<http://people.sabanciuniv.edu/~ysaygin/>

. Sabancı .  
Universitesi

# A Brief History

- Historically, we had operational databases, ex: for accounts, customers, personnel of a bank
- Data collection is now very easy and storage is very cheap
- Enterprises are in a data collection frenzy hoping that they can use that later on
- Data warehouses:
  - Integrating historical data from multiple sources
  - For high level decision making
- Data Mining

# Overview of Data Mining

- Why do we need data mining?
  - Data collection is easy, and huge amounts of data is collected everyday into flat files, databases and data warehouses
  - We have lots of data but this data needs to be turned into knowledge
  - Data mining technology tries to extract useful knowledge from huge collections of data

# Overview of Data Mining

- Data mining definition: Extraction of interesting information from large data sources
- The extracted information should be
  - Implicit
  - Non-trivial
  - Previously unknown
  - and potentially useful
- Query processing, simple statistics are not data mining
- Databases + Statistics + Machine Learning = Data Mining

# Overview of Data Mining

- Data mining applications
  - Market basket analysis
  - CRM (loyalty detection, churn detection)
  - Fraud detection
  - Stream mining
  - Web mining
  - Mining of bioinformatics data

# Overview of Data Mining

- Retail market, as a case study:
  - What type of data is collected?
  - What type of knowledge do we need about customers?
  - Is it useful to know the customer buying patterns?
  - Is it useful to segment the customers?

# Overview of Data Mining

- Advertisement of a product: A case study
  - Send all the customers a brochure
  - Or send a targeted list of customers a brochure
  - Sending a smaller targeted list aims to guarantee a high percentage of response, cutting the mailing cost

# Overview of Data Mining

- What complicates things in data mining?
  - Incomplete and noisy data
  - Complex data types
  - Heterogeneous data sources
  - Size of data (need to have distributed, parallel scalable algorithms)



# Data Mining Models

- Patterns (Associations, sequences, temporal sequences)
- Clusters (Descriptive)
- Predictive models (Classification)

# Associations (As an example of patterns)

- Remember the case study of retail market, and market basket analysis
- Remember the type of data collected
- Associations are among the most popular patterns that can be extracted from transactional data.
- We will explain the properties of associations and how they could be extracted from large collections of transactions efficiently based on the slide of the book : “Data Mining Concepts and Techniques” by Jiawei Han and Micheline Kamber.

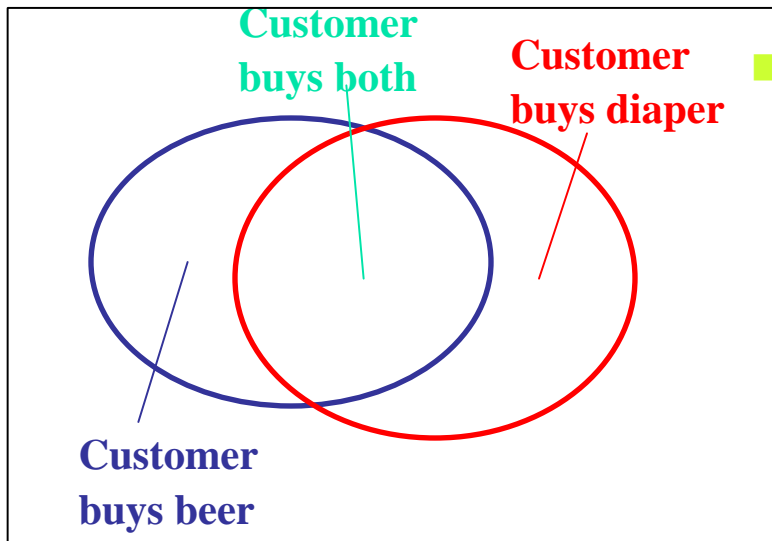
# What Is Association Mining?

- Association rule mining:
  - Finding frequent patterns, associations, correlations, or causal structures among sets of items or objects in transaction databases, relational databases, and other information repositories.
- Applications:
  - Basket data analysis, cross-marketing, catalog design, clustering, classification, etc.
- Examples.
  - Rule form: “Body  $\rightarrow$  Head [support, confidence]”.
  - buys(x, “diapers”)  $\rightarrow$  buys(x, “beers”) [0.5%, 60%]
  - major(x, “CS”)  $\wedge$  takes(x, “DB”)  $\rightarrow$  grade(x, “A”) [1%, 75%]

# Association Rule: Basic Concepts

- **Given:**
  - (1) database of transactions,
  - (2) each transaction is a list of items (purchased by a customer in a visit)
- **Find: all rules that correlate the presence of one set of items with that of another set of items**
  - E.g., *98% of people who purchase tires and auto accessories also get automotive services done*
- **Applications**
  - \*  $\Rightarrow$  *Maintenance Agreement* (What the store should do to boost Maintenance Agreement sales)
  - *Home Electronics*  $\Rightarrow$  \* (What other products should the store stocks up?)
  - Attached mailing in direct marketing
  - Detecting “ping-pong”ing of patients, faulty “collisions”

# Rule Measures: Support and Confidence



Find all the rules  $X \& Y \Rightarrow Z$  with minimum confidence and support

- **support**,  $s$  probability that a transaction contains  $\{X \& Y \& Z\}$
- **confidence**,  $c$  conditional probability that a transaction having  $\{X \& Y\}$  also contains  $Z$

Transaction ID	Items Bought
2000	A,B,C
1000	A,C
4000	A,D
5000	B,E,F

*Let minimum support 50%, and minimum confidence 50%, we have*

- $A \Rightarrow C$  (50%, 66.6%)
- $C \Rightarrow A$  (50%, 100%)

# Association Rule Mining: A Road Map

- Boolean vs. quantitative associations (Based on the types of values handled)
  - $\text{buys}(x, \text{"SQLServer"}) \wedge \text{buys}(x, \text{"DMBook"}) \rightarrow \text{buys}(x, \text{"DBMiner"})$  [0.2%, 60%]
  - $\text{age}(x, \text{"30..39"}) \wedge \text{income}(x, \text{"42..48K"}) \rightarrow \text{buys}(x, \text{"PC"})$  [1%, 75%]
- Single dimension vs. multiple dimensional associations (see ex. Above)
- Single level vs. multiple-level analysis
  - What brands of beers are associated with what brands of diapers?
- Various extensions
  - Correlation, causality analysis
    - Association does not necessarily imply correlation or causality
  - Maxpatterns and closed itemsets
  - Constraints enforced E.g., small sales (sum < 100) trigger big buys (sum > 1,000)?

# Mining Association Rules—An Example

Transaction ID	Items Bought
2000	A,B,C
1000	A,C
4000	A,D
5000	B,E,F

Min. support 50%  
Min. confidence 50%

Frequent Itemset	Support
{A}	75%
{B}	50%
{C}	50%
{A,C}	50%

For rule  $A \Rightarrow C$ :

support = support({A C}) = 50%

confidence = support({A C})/support({A}) = 66.6%

The **Apriori** principle: (Agrawal and Srikant)

Any subset of a frequent itemset must be frequent

# Mining Frequent Itemsets: the Key Step

- Find the *frequent itemsets*: the sets of items that have minimum support
  - A subset of a frequent itemset must also be a frequent itemset
    - i.e., if  $\{A\ B\}$  is a frequent itemset, both  $\{A\}$  and  $\{B\}$  should be a frequent itemset
  - Iteratively find frequent itemsets with cardinality from 1 to  $k$  ( $k$ -itemset)
- Use the frequent itemsets to generate association rules.



# The Apriori Algorithm

- **Join Step:**  $C_k$  is generated by joining  $L_{k-1}$  with itself
- **Prune Step:** Any  $(k-1)$ -itemset that is not frequent cannot be a subset of a frequent  $k$ -itemset
- Pseudo-code:

$C_k$ : Candidate itemset of size  $k$   
 $L_k$ : frequent itemset of size  $k$

$L_1 = \{\text{frequent items}\};$

**for** ( $k = 1; L_k \neq \emptyset; k++$ ) **do begin**

$C_{k+1}$  = candidates generated from  $L_k$ ;

**for each** transaction  $t$  in database **do**

        increment the count of all candidates in  $C_{k+1}$   
        that are contained in  $t$

$L_{k+1}$  = candidates in  $C_{k+1}$  with min\_support

**end**

**return**  $\cup_k L_k$ ;

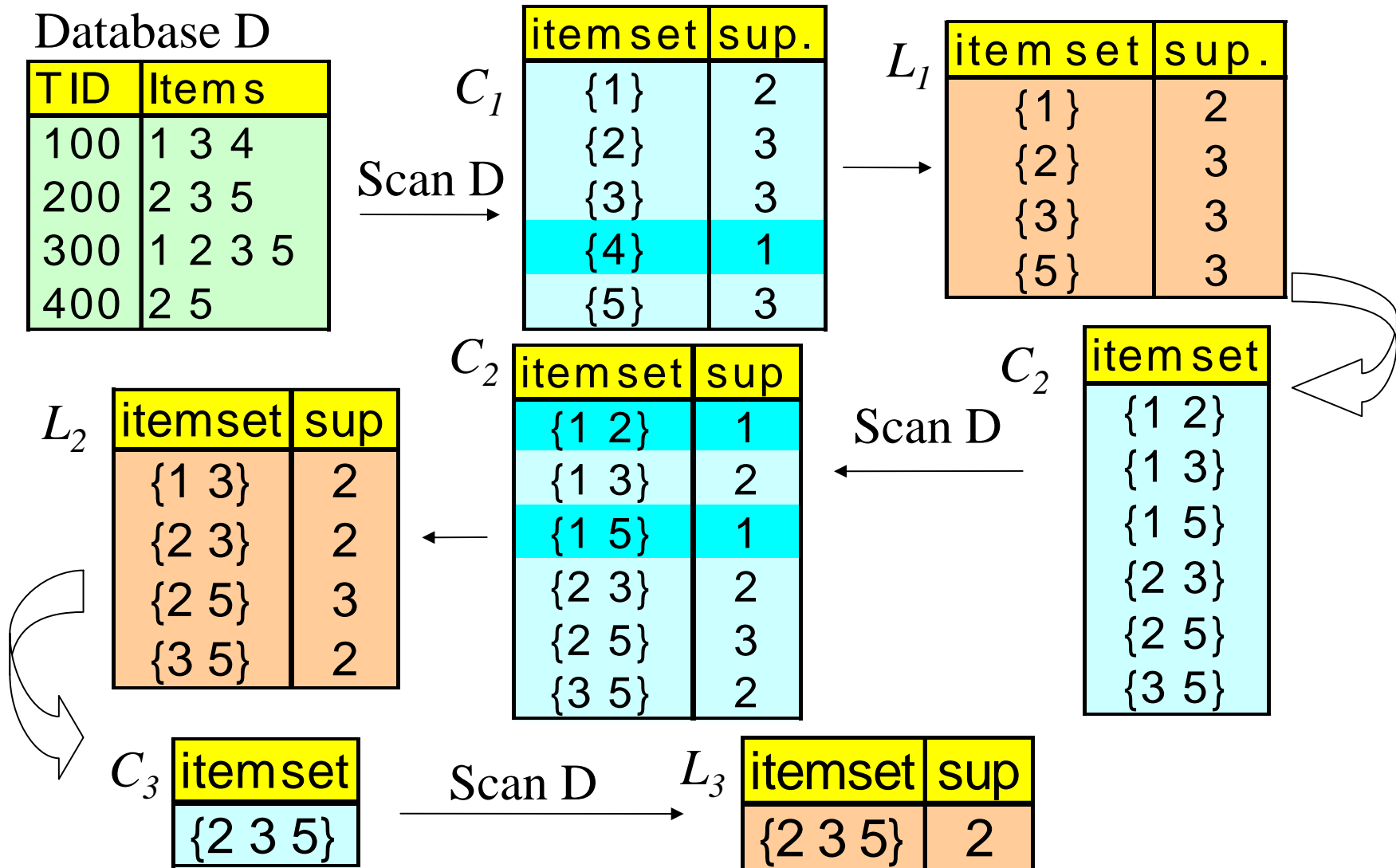
# The Apriori Algorithm — Example

Database D

TID	Items
100	1 3 4
200	2 3 5
300	1 2 3 5
400	2 5

- Find frequent sets of items with support 40% or more

# The Apriori Algorithm — Example



# How to Generate Candidates?

- Suppose the items in  $L_{k-1}$  are listed in an order

- **Step 1:** self-joining  $L_{k-1}$

insert into  $C_k$

select  $p.item_1, p.item_2, \dots, p.item_{k-1}, q.item_{k-1}$

from  $L_{k-1} p, L_{k-1} q$

where  $p.item_1=q.item_1, \dots, p.item_{k-2}=q.item_{k-2}, p.item_{k-1} < q.item_{k-1}$

- **Step 2:** pruning

forall *itemsets*  $c$  in  $C_k$  do

    forall *(k-1)-subsets*  $s$  of  $c$  do

        if ( $s$  is not in  $L_{k-1}$ ) then delete  $c$  from  $C_k$

# Example of Generating Candidates

- $L_3 = \{abc, abd, acd, ace, bcd\}$
- Self-joining:  $L_3 * L_3$ 
  - $abcd$  from  $abc$  and  $abd$
  - $acde$  from  $acd$  and  $ace$
- Pruning:
  - $acde$  is removed because  $ade$  is not in  $L_3$
- $C_4 = \{abcd\}$

# How to Count Supports of Candidates?

- Why counting supports of candidates is a problem?
  - The total number of candidates can be very huge
  - One transaction may contain many candidates
- Method:
  - Candidate itemsets are stored in a *hash-tree*
  - *Leaf* node of hash-tree contains a list of itemsets and counts
  - *Interior* node contains a hash table
  - *Subset function*: finds all the candidates contained in a transaction

# How to Generate Candidates?

- Assume that you have a relational database for storing customer transactions
- Design simple queries for counting the frequent itemsets
- Design a little system with database interaction for mining associations

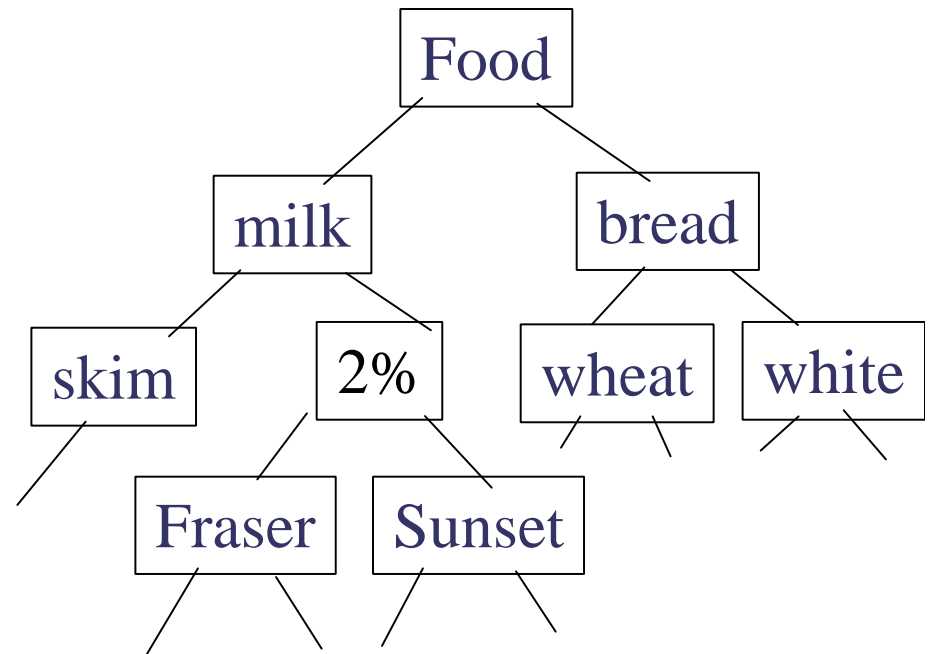
# Methods to Improve Apriori's Efficiency

- **Hash-based itemset counting:** A  $k$ -itemset whose corresponding hashing bucket count is below the threshold cannot be frequent
- **Transaction reduction:** A transaction that does not contain any frequent  $k$ -itemset is useless in subsequent scans
- **Partitioning:** Any itemset that is potentially frequent in DB must be frequent in at least one of the partitions of DB
- **Sampling:** mining on a subset of given data, lower support threshold + a method to determine the completeness
- **Dynamic itemset counting:** add new candidate itemsets only when all of their subsets are estimated to be frequent



# Multiple-Level Association Rules

- Items often form hierarchy.
- Items at the lower level are expected to have lower support.
- Rules regarding itemsets at appropriate levels could be quite useful.
- Transaction database can be encoded based on dimensions and levels
- We can explore shared multi-level mining



TID	Items
T1	{111, 121, 211, 221}
T2	{111, 211, 222, 323}
T3	{112, 122, 221, 411}
T4	{111, 121}
T5	{111, 122, 211, 221, 413}

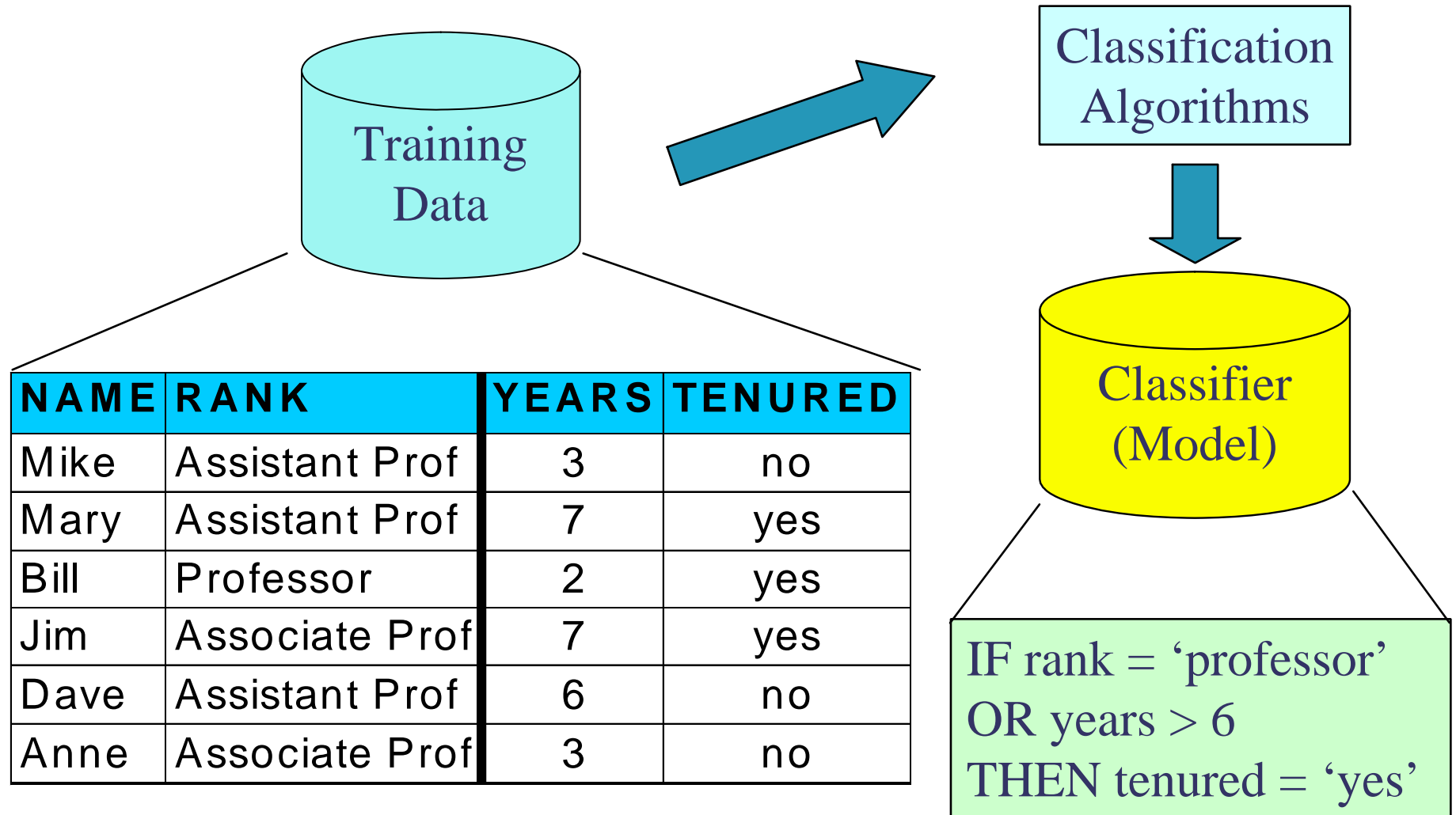
# Classification

- Is an example of predictive modelling
- The basic idea is to build a model using past data to predict the class of a new data sample.
- Lets remember the case of targeted mailing of brochures.
- IF we can work on a small well selected sample to profile the future customers who will respond to the mail ad, then we can save the mailing costs.
- The following slides are based on the slides of the book “Data Mining Concepts and Techniques” by Jiawei Han and Micheline Kamber.

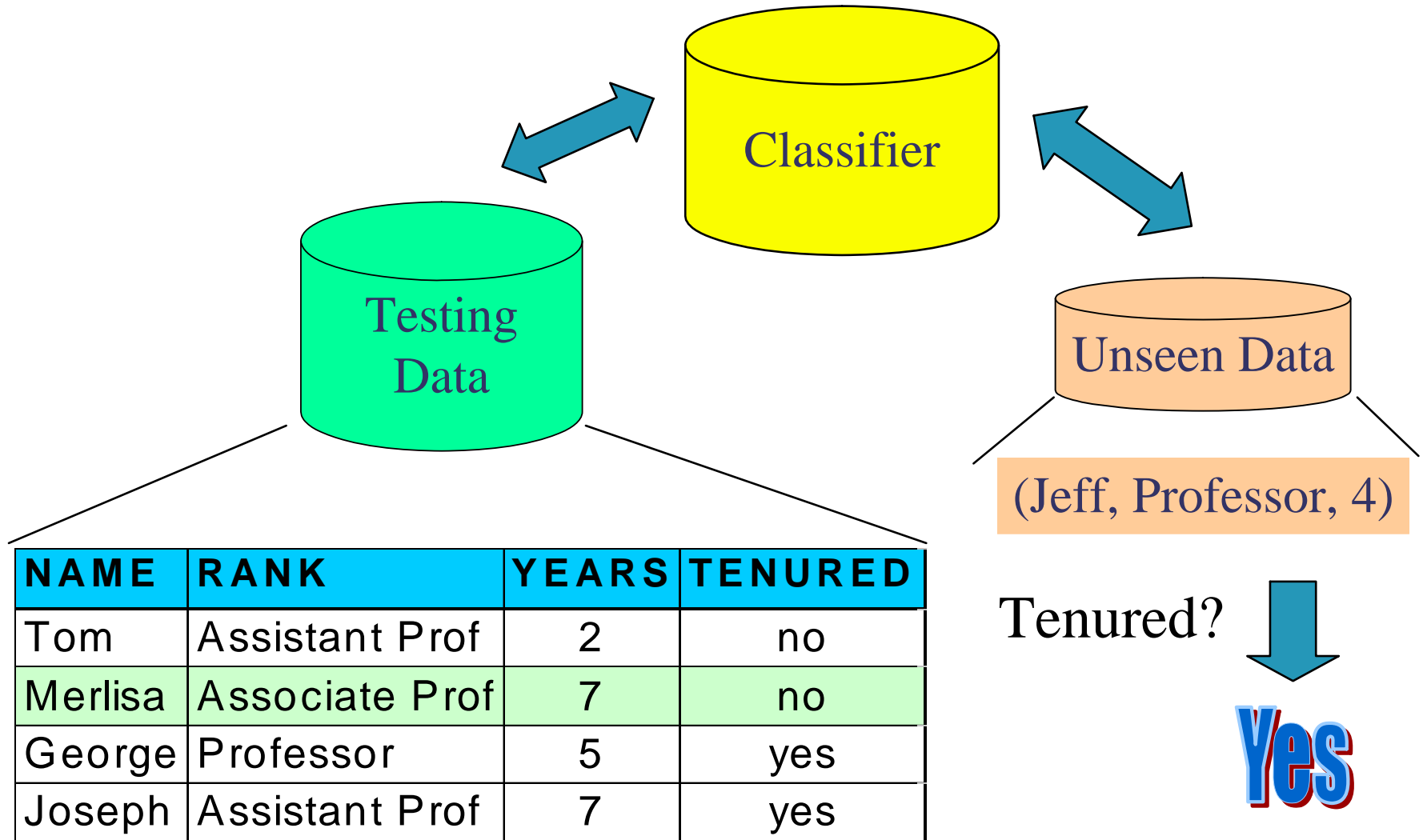
# Classification—A Two-Step Process

- **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
- **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

# Classification Process (1): Model Construction



# Classification Process (2): Use the Model in Prediction



# Supervised 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 (1): Data Preparation

- Data cleaning
  - Pre-process data in order to reduce noise and handle missing values
- Relevance analysis (feature selection)
  - Remove the irrelevant or redundant attributes
- Data transformation
  - Generalize and/or normalize data

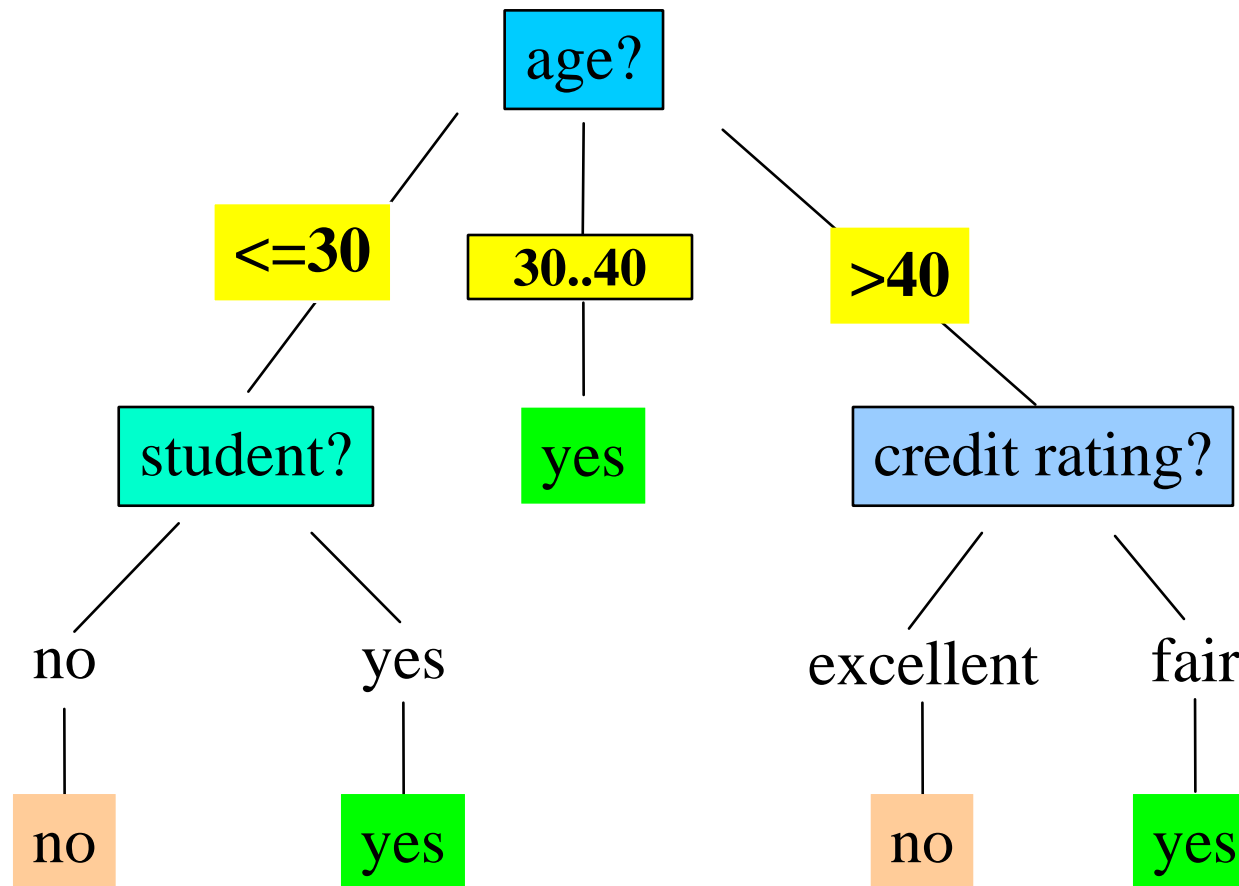
# Training Dataset

This follows an example from Quinlan's ID3

age	income	student	credit_rating	buys_computer
<=30	high	no	fair	no
<=30	high	no	excellent	no
30...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



## Output: A Decision Tree for “*buys\_computer*”



# 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 - Information Gain (ID3/C4.5)

- Select the attribute with the highest information gain
- S contains  $s_i$  tuples of class  $C_i$  for  $i = \{1, \dots, m\}$
- **information** measures info required to classify any arbitrary tuple

$$I(s_1, s_2, \dots, s_m) = - \sum_{i=1}^m \frac{s_i}{S} \log_2 \frac{s_i}{S}$$

- **entropy** of attribute A with values  $\{a_1, a_2, \dots, a_v\}$

$$E(A) = \sum_{j=1}^v \frac{s_{1j} + \dots + s_{mj}}{S} I(s_{1j}, \dots, s_{mj})$$

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

$$Gain(A) = I(s_1, s_2, \dots, s_m) - E(A)$$

# Attribute Selection by Information Gain Computation

- Class P: buys\_computer = “yes”
- Class N: buys\_computer = “no”
- $I(p, n) = I(9, 5) = 0.940$

$$I(s_1, s_2, \dots, s_m) = - \sum_{i=1}^m \frac{s_i}{s} \log_2 \frac{s_i}{s}$$

age	income	student	credit_rating	buys_computer
<=30	high	no	fair	no
<=30	high	no	excellent	no
30...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

# Attribute Selection by Information Gain Computation

- Class P: buys\_computer = “yes”
- Class N: buys\_computer = “no”
- $I(p, n) = I(9, 5) = 0.940$
- Compute the entropy for age:

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

$$E(age) = \frac{5}{14} I(2,3) + \frac{4}{14} I(4,0) + \frac{5}{14} I(3,2) = 0.971$$

Hence

$$Gain(age) = I(p, n) - E(age)$$

Similarly

$$Gain(income) = 0.029$$

$$Gain(student) = 0.151$$

$$Gain(credit\_rating) = 0.048$$

# Other Attribute Selection Measures

- **Gini index** (CART, IBM IntelligentMiner)
  - All attributes are assumed to be continuous-valued
  - Assume there exist several possible split values for each attribute
  - May need other tools, such as clustering, to get the possible split values
  - Can be modified for categorical attributes

# Gini Index (IBM IntelligentMiner)

- If a data set  $T$  contains examples from  $n$  classes, gini index,  $gini(T)$  is defined as 
$$gini(T) = 1 - \sum_{j=1}^n p_j^2$$

where  $p_j$  is the relative frequency of class  $j$  in  $T$ .

- If a data set  $T$  is split into two subsets  $T_1$  and  $T_2$  with sizes  $N_1$  and  $N_2$  respectively, the gini index of the split data contains examples from  $n$  classes, the gini index  $gini(T)$  is defined as

$$gini_{split}(T) = \frac{N_1}{N} gini(T_1) + \frac{N_2}{N} gini(T_2)$$

- The attribute provides the smallest  $gini_{split}(T)$  is chosen to split the node (*need to enumerate all possible splitting points for each attribute*).

# Extracting Classification Rules from Trees

- Represent the knowledge in the form of **IF-THEN** rules
- One rule is created for each path from the root to a leaf
- Each attribute-value pair along a path forms a conjunction
- The leaf node holds the class prediction
- Rules are easier for humans to understand
- Example

IF *age* = “<=30” AND *student* = “no” THEN *buys\_computer* = “no”

IF *age* = “<=30” AND *student* = “yes” THEN *buys\_computer* = “yes”

IF *age* = “31...40” THEN *buys\_computer* = “yes”

IF *age* = “>40” AND *credit\_rating* = “excellent” THEN *buys\_computer* = “yes”

IF *age* = “<=30” AND *credit\_rating* = “fair” THEN *buys\_computer* = “no”



# Avoid Overfitting in Classification

- The generated tree may overfit the training data
  - Too many branches, some may reflect anomalies due to noise or outliers
  - Result is in poor accuracy for unseen samples
- Two approaches to avoid overfitting
  - Prepruning: Halt tree construction early—do not split a node if this would result in the goodness measure falling below a threshold
    - Difficult to choose an appropriate threshold
  - Postpruning: Remove branches from a “fully grown” tree—get a sequence of progressively pruned trees
    - Use a set of data different from the training data to decide which is the “best pruned tree”

# Approaches to Determine the Final Tree Size

- Separate training (2/3) and testing (1/3) sets
- Use cross validation, e.g., 10-fold cross validation
- Use all the data for training
  - but apply a **statistical test** (e.g., chi-square) to estimate whether expanding or pruning a node may improve the entire distribution

# Bayesian Classification: Why?

- Probabilistic learning: Calculate explicit probabilities for hypothesis, among the most practical approaches to certain types of learning problems
- Incremental: Each training example can incrementally increase/decrease the probability that a hypothesis is correct. Prior knowledge can be combined with observed data.
- Probabilistic prediction: Predict multiple hypotheses, weighted by their probabilities
- Standard: Even when Bayesian methods are computationally intractable, they can provide a standard of optimal decision making against which other methods can be measured

# Bayesian Theorem: Basics

- Let  $X$  be a data sample whose class label is unknown
- Let  $H$  be a hypothesis that  $X$  belongs to class  $C$
- For classification problems, determine  $P(H/X)$ : the probability that the hypothesis holds given the observed data sample  $X$
- $P(H)$ : prior probability of hypothesis  $H$  (i.e. the initial probability before we observe any data, reflects the background knowledge)
- $P(X)$ : probability that sample data is observed
- $P(X|H)$  : probability of observing the sample  $X$ , given that the hypothesis holds

# Bayesian Theorem

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

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

- Informally, this can be written as  
**posterior = likelihood x prior / evidence**
- MAP (maximum posteriori) hypothesis

$$h_{MAP} \equiv \arg \max_{h \in H} P(h|D) = \arg \max_{h \in H} P(D|h)P(h).$$

- Practical difficulty: require initial knowledge of many probabilities, significant computational cost

# Naïve Bayesian Classifier

- Each data sample  $X$  is represented as a vector  $\{x_1, x_2, \dots, x_n\}$
- There are  $m$  classes  $C_1, C_2, \dots, C_m$
- Given unknown data sample  $X$ , the classifier will predict that  $X$  belongs to class  $C_i$ , iff

$$P(C_i|X) > P(C_j|X) \text{ where } 1 \leq j \leq m, i \neq j$$

$$\text{By Bayes theorem, } P(C_i|X) = P(X|C_i)P(C_i) / P(X)$$

# Naïve Bayes Classifier

- A simplified assumption: attributes are conditionally independent:

$$P(X | C_i) = \prod_{k=1}^n P(x_k | C_i)$$

- The product of occurrence of say 2 elements  $x_1$  and  $x_2$ , given the current class is  $C$ , is the product of the probabilities of each element taken separately, given the same class  $P([y_1, y_2], C) = P(y_1, C) * P(y_2, C)$
- No dependence relation between attributes
- Greatly reduces the computation cost, only count the class distribution.
- Once the probability  $P(X|C_i)$  is known, assign  $X$  to the class with maximum  $P(X|C_i)*P(C_i)$

# Training dataset

Class:

C1: buys\_computer=  
'yes'

C2: buys\_computer=  
'no'

Data sample

X = (age ≤ 30,  
Income = medium,  
Student = yes  
Credit\_rating =  
Fair)

age	income	student	credit_rating	buys_computer
≤30	high	no	fair	no
≤30	high	no	excellent	no
30...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



# Naïve Bayesian Classifier: Example

- Compute  $P(X/C_i)$  for each class

$$P(\text{age}=\text{"<30"} \mid \text{buys\_computer}=\text{"yes"}) = 2/9=0.222$$

$$P(\text{age}=\text{"<30"} \mid \text{buys\_computer}=\text{"no"}) = 3/5 =0.6$$

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

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

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

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

$$P(\text{credit\_rating}=\text{"fair"} \mid \text{buys\_computer}=\text{"yes"})=6/9=0.667$$

$$P(\text{credit\_rating}=\text{"fair"} \mid \text{buys\_computer}=\text{"no"})=2/5=0.4$$

**X=(age<=30 ,income =medium, student=yes,credit\_rating=fair)**

$$\mathbf{P(X|C_i)} : P(X|\text{buys\_computer}=\text{"yes"})= 0.222 \times 0.444 \times 0.667 \times 0.667 =0.044$$

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

$$\mathbf{P(X|C_i)*P(C_i)} : P(X|\text{buys\_computer}=\text{"yes"}) * P(\text{buys\_computer}=\text{"yes"})=0.028$$

$$P(X|\text{buys\_computer}=\text{"no"}) * P(\text{buys\_computer}=\text{"no"})=0.007$$

**X belongs to class "buys\_computer=yes"**

# Naïve Bayesian Classifier: Comments

- Advantages :
  - Easy to implement
  - Good results obtained in most of the cases
- Disadvantages
  - Assumption: class conditional independence , therefore loss of accuracy
  - Practically, dependencies exist among variables
  - E.g., hospitals : patients: Profile : age, family history etc  
Symptoms : fever, cough etc , Disease : lung cancer, diabetes etc  
, Dependencies among these cannot be modeled by Naïve Bayesian Classifier, use a Bayesian network
- How to deal with these dependencies?
  - Bayesian Belief Networks

# k-NN Classifier:

- Learning by analogy,
- Each sample is a point in n-dimensional space
- Given an unknown sample  $u$ ,
  - search for the  $k$  nearest samples
  - Closeness can be defined in Euclidean space
  - Assign the most common class to  $u$ .
- Instance based, (Lazy) learning while decision trees are eager
- K-NN requires the whole sample space for classification therefore indexing is needed for efficient search.

# Case-based reasoning

- Similar to K-NN,
- When a new case arrives, an identical case is searched
- If not, most similar case is searched
- Depending on the representation, different search techniques are needed, for example graph/subgraph search

# Genetic Algorithms

- Incorporate ideas from natural evolution
- Rules are represented as a sequence of bits
  - IF A1 and NOT A2 THEN C2 : 1 0 0
  - Initially generate a sequence of random rules
  - Choose the fittest rules
  - Create offspring by using genetic operations such as
    - crossover (by swapping substrings from pairs of rules)
    - and mutation (inverting randomly selected bits)

# Data Mining Tools

- WEKA (Univ of Waikato, NZ)
- Open source implementation of data mining algorithms
- Implemented in Java
- Nice API
- Link : [Google WEKA, first entry](#)

# Benchmarks

- Important for testing the performance of various algorithms against various datasets
- UCI Machine Learning Repository is a very good data source
- Try a data sample from there and see how you can apply the classification algorithms

# Clustering

- A descriptive data mining method
- Groups a given dataset into smaller clusters, where the data inside the clusters are similar to each other while the data belonging to different clusters are dissimilar
- Similar to classification in a sense but this time we do not know the labels of clusters. Therefore it is an unsupervised method.
- Lets go back to the retail market example. How can we segment our customers with respect to their profiles and shopping behaviour.
- The following slides are based on the slides of the book “Data Mining Concepts and Techniques” by **Jiawei Han and Micheline Kamber**.



# What Is Good Clustering?

- A good clustering method will produce high quality clusters with
  - high intra-class similarity
  - low inter-class similarity
- The quality of a clustering result depends on both the similarity measure used by the method and its implementation.
- The quality of a clustering method is also measured by its ability to discover some or all of the hidden patterns.

# Requirements of Clustering in Data Mining

- Scalability
- Ability to deal with different types of attributes
- Discovery of clusters with arbitrary shape
- Able to deal with noise and outliers
- Insensitive to order of input records
- High dimensionality
- Interpretability and usability



# Measure the Quality of Clustering

- Dissimilarity/Similarity metric: Similarity is expressed in terms of a distance function, which is typically metric:  
 $d(i, j)$
- There is a separate “quality” function that measures the “goodness” of a cluster.
- The definitions of distance functions are usually very different for interval-scaled, boolean, categorical, ordinal and ratio variables.
- Weights should be associated with different variables based on applications and data semantics.
- It is hard to define “similar enough” or “good enough”
  - the answer is typically highly subjective.

# Type of data in clustering analysis

- Interval-scaled variables:
- Binary variables:
- Nominal, ordinal, and ratio variables:
- Variables of mixed types:

# Interval-scaled variables

- Standardize data

- Calculate the mean absolute deviation:  $s_f = \frac{1}{n}(|x_{1f} - m_f| + |x_{2f} - m_f| + \dots + |x_{nf} - m_f|)$

where  $m_f = \frac{1}{n}(x_{1f} + x_{2f} + \dots + x_{nf})$ .

- Calculate the standardized measurement (z-score)  $z_{if} = \frac{x_{if} - m_f}{s_f}$

- Using mean absolute deviation is more robust than using standard deviation

# Similarity and Dissimilarity Between Objects

- Distances are normally used to measure the similarity or dissimilarity between two data objects
- Some popular ones include: *Minkowski distance*:

$$d(i, j) = \sqrt[q]{(|x_{i_1} - x_{j_1}|^q + |x_{i_2} - x_{j_2}|^q + \dots + |x_{i_p} - x_{j_p}|^q)}$$

where  $i = (x_{i_1}, x_{i_2}, \dots, x_{i_p})$  and  $j = (x_{j_1}, x_{j_2}, \dots, x_{j_p})$  are two  $p$ -dimensional data objects, and  $q$  is a positive integer

- If  $q = 1$ ,  $d$  is Manhattan distance

$$d(i, j) = |x_{i_1} - x_{j_1}| + |x_{i_2} - x_{j_2}| + \dots + |x_{i_p} - x_{j_p}|$$

# Similarity and Dissimilarity Between Objects (Cont.)

- If  $q = 2$ ,  $d$  is Euclidean distance:

$$d(i, j) = \sqrt{(|x_{i_1} - x_{j_1}|^2 + |x_{i_2} - x_{j_2}|^2 + \dots + |x_{i_p} - x_{j_p}|^2)}$$

- Properties

- $d(i, j) \geq 0$
- $d(i, i) = 0$
- $d(i, j) = d(j, i)$
- $d(i, j) \leq d(i, k) + d(k, j)$

- Also one can use weighted distance.



# Binary Variables

- A contingency table for binary data

		Object $j$		
		1	0	$sum$
Object $i$	1	$a$	$b$	$a+b$
	0	$c$	$d$	$c+d$
	$sum$	$a+c$	$b+d$	$p$

- Simple matching coefficient (invariant, if the binary variable is symmetric):

$$d(i, j) = \frac{b + c}{a + b + c + d}$$

- Jaccard coefficient (noninvariant if the binary variable is asymmetric):

$$d(i, j) = \frac{b + c}{a + b + c}$$

# Dissimilarity between Binary Variables

- Example

Name	Gender	Fever	Cough	Test-1	Test-2	Test-3	Test-4
Jack	M	Y	N	P	N	N	N
Mary	F	Y	N	P	N	P	N
Jim	M	Y	P	N	N	N	N

- gender is a symmetric attribute
- the remaining attributes are asymmetric binary
- let the values Y and P be set to 1, and the value N be set to 0

$$d(\text{jack}, \text{mary}) = \frac{0 + 1}{2 + 0 + 1} = 0.33$$

$$d(\text{jack}, \text{jim}) = \frac{1 + 1}{1 + 1 + 1} = 0.67$$

$$d(\text{jim}, \text{mary}) = \frac{1 + 2}{1 + 1 + 2} = 0.75$$

# Nominal Variables

- A generalization of the binary variable in that it can take more than 2 states, e.g., red, yellow, blue, green
- Method 1: Simple matching
  - $m$ : # of matches,  $p$ : total # of variables

$$d(i, j) = \frac{p - m}{p}$$

- Method 2: use a large number of binary variables
  - creating a new binary variable for each of the  $M$  nominal states

# Ordinal Variables

- An ordinal variable can be discrete or continuous
- order is important, e.g., rank
- Can be treated like interval-scaled
  - replacing  $x_{if}$  by their rank  $r_{if} \in \{1, \dots, M_f\}$
  - map the range of each variable onto  $[0, 1]$  by replacing  $i$ -th object in the  $f$ -th variable

by

$$z_{if} = \frac{r_{if} - 1}{M_f - 1}$$

- compute the dissimilarity using methods for interval-scaled variables

# Similarity and Dissimilarity Between Objects (Cont.)

- If  $q = 2$ ,  $d$  is Euclidean distance:

$$d(i,j) = \sqrt{(|x_{i_1} - x_{j_1}|^2 + |x_{i_2} - x_{j_2}|^2 + \dots + |x_{i_p} - x_{j_p}|^2)}$$

- Properties

- $d(i,j) \geq 0$
  - $d(i,i) = 0$
  - $d(i,j) = d(j,i)$
  - $d(i,j) \leq d(i,k) + d(k,j)$
- Also, one can use weighted distance, parametric

# Ratio-Scaled Variables

- Ratio-scaled variable: a positive measurement on a nonlinear scale, approximately at exponential scale, such as  $Ae^{Bt}$  or  $Ae^{-Bt}$
- Methods:
  - treat them like interval-scaled variables — *not a good choice! (why?)*
  - apply logarithmic transformation

$$y_{if} = \log(x_{if})$$

- treat them as continuous ordinal data treat their rank as interval-scaled.

# Major Clustering Approaches

- Partitioning algorithms: Construct various partitions and then evaluate them by some criterion
- Hierarchical algorithms: Create a hierarchical decomposition of the set of data (or objects) using some criterion
- Density-based: based on connectivity and density functions

# Partitioning Algorithms: Basic Concept

- Partitioning method: Construct a partition of a database  $D$  of  $n$  objects into a set of  $k$  clusters
- Given a  $k$ , find a partition of  $k$  clusters that optimizes the chosen partitioning criterion
  - Global optimal: exhaustively enumerate all partitions
  - Heuristic methods: *k-means* and *k-medoids* algorithms
  - *k-means* (MacQueen'67): Each cluster is represented by the center of the cluster
  - *k-medoids* or PAM (Partition around medoids) (Kaufman & Rousseeuw'87): Each cluster is represented by one of the objects in the cluster

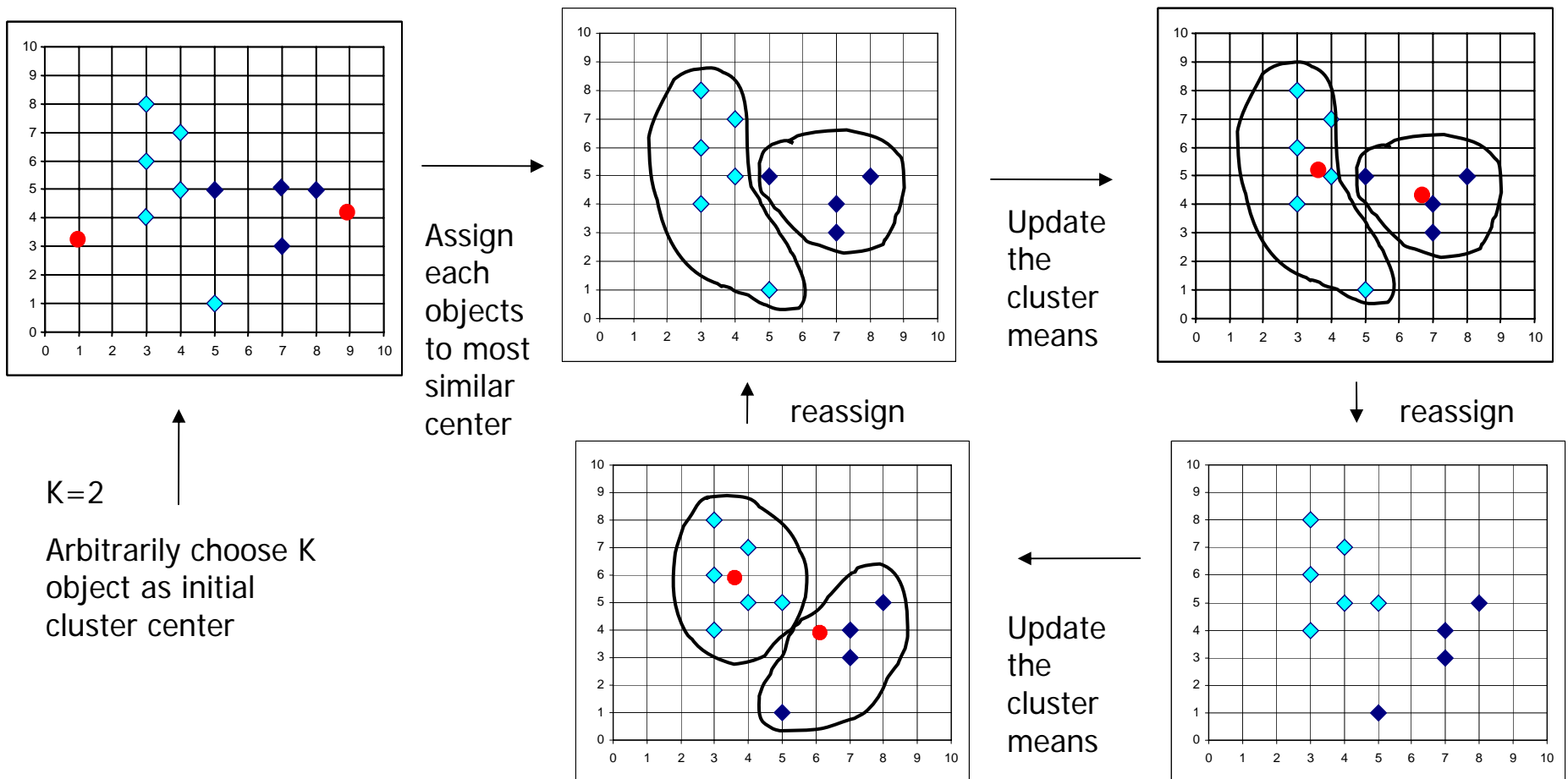


# The *K-Means* Clustering Method

- Given  $k$ , the *k-means* algorithm is implemented in four steps:
  - Partition objects into  $k$  nonempty subsets
  - Compute seed points as the centroids of the clusters of the current partition (the centroid is the center, i.e., *mean point*, of the cluster)
  - Assign each object to the cluster with the nearest seed point
  - Go back to Step 2, stop when no more new assignment

# The *K-Means* Clustering Method

## ■ Example



# Comments on the *K-Means* Method

- Strength: *Relatively efficient*.  $O(tkn)$ , where  $n$  is # objects,  $k$  is # clusters, and  $t$  is # iterations. Normally,  $k, t \ll n$ .
  - Comparing: PAM:  $O(k(n-k)^2)$ , CLARA:  $O(ks^2 + k(n-k))$
- Comment: Often terminates at a *local optimum*. The *global optimum* may be found using techniques such as: *deterministic annealing* and *genetic algorithms*
- Weakness
  - Applicable only when *mean* is defined, then what about categorical data?
  - Need to specify  $k$ , the *number* of clusters, in advance
  - Unable to handle noisy data and *outliers*
  - Not suitable to discover clusters with *non-convex shapes*