

Data Visualization

- Categorization of visualization methods:
 - Pixel-oriented** visualization techniques
 - Geometric projection** visualization techniques
 - Icon-based** visualization techniques
 - Hierarchical** visualization techniques
 - Visualizing complex data and relations**

Geometric: Direct visualization; scatterplot and scatterplot matrices; landscape; parallel coordinates

Icon-based: chernoff faces; stick figures

Hierarchical: dimensional stack; tree-map; cone trees; infocube

Distance on Numeric Data: Minkowski Distance

- Minkowski distance: A popular distance measure

$$d(i, j) = \sqrt[p]{|x_{i1} - x_{j1}|^p + |x_{i2} - x_{j2}|^p + \dots + |x_{ip} - x_{jp}|^p}$$

where $i = (x_{i1}, x_{i2}, \dots, x_{ip})$ and $j = (x_{j1}, x_{j2}, \dots, x_{jp})$ are two p -dimensional data objects, and p is the order (the distance so defined is also called L $_p$ norm)

- $p = 1$: (L₁ norm) Manhattan (or city block) distance

E.g., the Hamming distance: the number of bits that are different between two binary vectors

$$d(i, j) = |x_{i1} - x_{j1}| + |x_{i2} - x_{j2}| + \dots + |x_{ip} - x_{jp}|$$

- $p = 2$: (L₂ norm) Euclidean distance

$$d(i, j) = \sqrt{|x_{i1} - x_{j1}|^2 + |x_{i2} - x_{j2}|^2 + \dots + |x_{ip} - x_{jp}|^2}$$

- $p \rightarrow \infty$: (L_{max} norm, L _{∞} norm) "supremum" distance

The maximum difference between any component (attribute) of the vectors

$$d(i, j) = \lim_{p \rightarrow \infty} \sqrt[p]{|x_{i1} - x_{j1}|^p + |x_{i2} - x_{j2}|^p + \dots + |x_{ip} - x_{jp}|^p} = \max_{f=1}^p |x_{if} - x_{jf}|$$

Proximity Measure for Binary Attributes

- A contingency table for binary data

Object /	Object /			sum
	1	q	r	
0	s	t	s+t	p
sum	q+s	r+t	s+t	p

Distance measure for symmetric binary variables: $d(i, j) = \frac{r+s}{q+r+s+t}$

Distance measure for asymmetric binary variables: $d(i, j) = \frac{r+s}{q+r+s+t}$

Jaccard coefficient (similarity measure for asymmetric binary variables): $sim_{Jaccard}(i, j) = \frac{q}{q+r+s+t}$

Note: Jaccard coefficient is the same as

[a concept discussed in Pattern Discovery]

$$coherence(i, j) = \frac{sup(i, j)}{sup(i) + sup(j) - sup(i, j)} = \frac{q}{(q+r) + (q+s) - q}$$

Proximity Measure for Categorical Attributes

- Categorical data, also called nominal attributes

Example: Color (red, yellow, blue, green), profession, etc.

- 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 attributes

Creating a new binary attribute for each of the M nominal states

Ordinal Variables

- An ordinal variable can be discrete or continuous
- Order is important, e.g., rank (e.g., freshman, sophomore, junior, senior)
- Can be treated like interval-scaled
- Replace an ordinal variable value by its rank: $r_j \in \{1, \dots, M_f\}$
- Map the range of each variable onto [0, 1] by replacing i -th object by the f -th variable by $z_{if} = \frac{r_f - 1}{M_f - 1}$

Example: freshman: 0; sophomore: 1/3; junior: 2/3; senior 1

Then distance: $d(\text{freshman}, \text{senior}) = 1$, $d(\text{junior}, \text{senior}) = 1/3$

Compute the dissimilarity using methods for interval-scaled variables

Cosine Similarity of Two Vectors

- A document can be represented by a bag of terms or a long vector, with each attribute recording the frequency of a particular term (such as word, keyword, or phrase) in the document

Document1	team	coach	hockey	baseball	soccer	penalty	score	win	loss	0
Document1	5	0	3	0	2	1	1	0	1	0
Document2	3	0	2	0	1	1	0	1	0	1
Document3	0	7	0	2	1	0	0	3	0	1
Document4	0	1	0	0	1	2	0	2	3	0

Other vector objects: Gene features in micro-arrays

Applications: Information retrieval, biologic taxonomy, gene feature mapping, etc.

Cosine measure: If d_1 and d_2 are two vectors (e.g., term-frequency vectors), then

$$\cos(d_1, d_2) = \frac{d_1 \cdot d_2}{\|d_1\| \|d_2\|}$$

where \cdot indicates vector dot product, $\|\cdot\|$: the length of vector d

KL Divergence: Comparing Two Probability Distributions

- The Kullback-Leibler (KL) divergence:

Measure the difference between two probability distributions over the same variable x

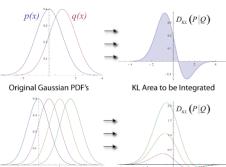
From information theory, closely related to relative entropy, information divergence, and information for discrimination

$D_{KL}(p(x) || q(x))$: divergence of $q(x)$ from $p(x)$, measuring the information lost when $q(x)$ is used to approximate $p(x)$

$$D_{KL}(p(x) || q(x)) = \sum_{x \in X} p(x) \ln \frac{p(x)}{q(x)}$$

Discrete form: $D_{KL}(p(x) || q(x)) = \sum_{x \in X} p(x) \ln \frac{p(x)}{q(x)}$

$$D_{KL}(p(x) || q(x)) = \int_{-\infty}^{\infty} p(x) \ln \frac{p(x)}{q(x)} dx$$



Ack: Wikipedia entry: The Kullback-Leibler (KL) divergence

Chapter 3. Data Preprocessing

- Redundant attributes may be able to be detected by correlation analysis and covariance analysis

Correlation Analysis (for Categorical Data)

- X² (chi-square) test:

$$\chi^2 = \sum_i \frac{(O_i - E_i)^2}{E_i}$$

Null hypothesis: The two distributions are independent

The cells that contribute the most to the X² value are those whose actual count is very different from the expected count

The larger the X² value, the more likely the variables are related

Note: Correlation does not imply causality

of hospitals and # of car-theft in a city are correlated

Both are causally linked to the third variable: population

Variance for Single Variable (Numerical Data)

- The variance of a random variable X provides a measure of how much the value of X deviates from the mean or expected value of X :

$$\sigma^2 = \text{var}(X) = E[(X - \mu)^2] = \sum_x (x - \mu)^2 f(x)$$

$$\sigma^2 = \text{var}(X) = E[(X - \mu)^2] = E[X^2] - \mu^2 = E[X^2] - [E(X)]^2$$

Covariance for Two Variables

- Covariance between two variables X_1 and X_2

$$\sigma_{12} = E[(X_1 - \mu_1)(X_2 - \mu_2)] = E[X_1 X_2] - \mu_1 \mu_2 = E[X_1 X_2] - E[X_1] E[X_2]$$

where $\mu_1 = E[X_1]$ is the respective mean or expected value of X_1 ; similarly for μ_2

- Positive covariance: If $\sigma_{12} > 0$

- Negative covariance: If $\sigma_{12} < 0$

- Independence: If X_1 and X_2 are independent, $\sigma_{12} = 0$ but the reverse is not true

Some pairs of random variables may have a covariance of 0 but are not independent

Only under some additional assumptions (e.g., the data follow multivariate normal distributions) does a covariance of 0 imply independence

Correlation between Two Numerical Variables

- Correlation between two variables X_1 and X_2 is the standard covariance, obtained by normalizing the covariance with the standard deviation of each variable

$$\rho_{12} = \frac{\sigma_{12}}{\sigma_1 \sigma_2} = \frac{\sigma_{12}}{\sqrt{\sigma_1^2 \sigma_2^2}}$$

- Sample correlation for two attributes X_1 and X_2 : $\hat{\rho}_{12} = \frac{\sum_i (x_{i1} - \bar{x}_1)(x_{i2} - \bar{x}_2)}{\sqrt{\sum_i (x_{i1} - \bar{x}_1)^2} \sqrt{\sum_i (x_{i2} - \bar{x}_2)^2}}$

where n is the number of tuples, μ_1 and μ_2 are the respective means of X_1 and X_2 , σ_1 and σ_2 are the respective standard deviation of X_1 and X_2

- If $\rho_{12} > 0$: A and B are positively correlated (X_1 's values increase as X_2 's)

The higher, the stronger correlation

- If $\rho_{12} = 0$: independent (under the same assumption as discussed in co-variance)

- If $\rho_{12} < 0$: negatively correlated

Covariance Matrix

- The variance and covariance information for the two variables X_1 and X_2 can be summarized as 2x2 covariance matrix as

$$\Sigma = E[(X - \mu)(X - \mu)^T] = \begin{pmatrix} \sigma_{11} & \sigma_{12} \\ \sigma_{21} & \sigma_{22} \end{pmatrix} = \frac{1}{n} \sum_{i=1}^n \begin{pmatrix} (x_{i1} - \bar{x}_1)(x_{i1} - \bar{x}_1) & (x_{i1} - \bar{x}_1)(x_{i2} - \bar{x}_2) \\ (x_{i2} - \bar{x}_2)(x_{i1} - \bar{x}_1) & (x_{i2} - \bar{x}_2)(x_{i2} - \bar{x}_2) \end{pmatrix}$$

Generalizing it to d dimensions, we have,

$$D = \begin{pmatrix} x_{11} & x_{12} & \dots & x_{1d} \\ x_{21} & x_{22} & \dots & x_{2d} \\ \vdots & \vdots & \ddots & \vdots \\ x_{d1} & x_{d2} & \dots & x_{dd} \end{pmatrix} \quad \Sigma = E[(X - \mu)(X - \mu)^T] = \begin{pmatrix} \sigma_{11} & \sigma_{12} & \dots & \sigma_{1d} \\ \sigma_{21} & \sigma_{22} & \dots & \sigma_{2d} \\ \vdots & \vdots & \ddots & \vdots \\ \sigma_{d1} & \sigma_{d2} & \dots & \sigma_{dd} \end{pmatrix}$$

Data Reduction

regression, histogram, sampling, clustering

Principal Component Analysis (PCA)

- PCA: A statistical procedure that uses an orthogonal transformation to convert a set of observations of possibly correlated variables into a set of values of linearly uncorrelated variables called principal components

The original data are projected onto a much smaller space, resulting in dimensionality reduction

- Method: Find the eigenvectors of the covariance matrix, and these eigenvectors define the new space

Star Schema: An Example

Fact Table: Sales Fact Table

Dim Table: DimCustomer

Dim Table: DimProduct

Dim Table: DimRegion

Dim Table: DimSupplier

Dim Table: DimTime

Dim Table: DimShipping

Dim Table: DimRegion

Dim Table: DimCountry

Dim Table: DimCity

Dim Table: DimState

Dim Table: DimProductCategory

Dim Table: DimProductSubcategory

Dim Table: DimProductSubcategory

Snowflake Schema: An Example

Fact Table: Sales Fact Table

Dim Table: DimCustomer

Dim Table: DimProduct

Dim Table: DimRegion

Dim Table: DimSupplier

Dim Table: DimTime

Dim Table: DimShipping

Dim Table: DimRegion

Dim Table: DimCountry

Dim Table: DimCity

Dim Table: DimState

Dim Table: DimProductCategory

Dim Table: DimProductSubcategory

Dim Table: DimProductSubcategory

Fact Constellation: An Example

Fact Table: Sales Fact Table

Dim Table: DimCustomer

Dim Table: DimProduct

Dim Table: DimRegion

Dim Table: DimSupplier

Dim Table: DimTime

Dim Table: DimShipping

Dim Table: DimRegion

Dim Table: DimCountry

Dim Table: DimCity

Dim Table: DimState

Dim Table: DimProductCategory

Dim Table: DimProductSubcategory

Dim Table: DimProductSubcategory

Data Cube Measures: Three Categories

- Distributive:** if the result derived by applying the function to n aggregate values is the same as that derived by applying the function on all the data without partitioning

E.g., count(), sum(), min(), max()

- Algebraic:** if it can be computed by an algebraic function with M arguments (where M is a bounded integer), each of which is obtained by applying a distributive aggregate function

avg(x) = sum(x) / count(x)

- Is min_(N) an algebraic measure? How about standard deviation()?

- Holistic:** if there is no constant bound on the storage size needed to describe a subaggregate.

E.g., median(), mode(), rank()

Typical OLAP Operations

- Roll up (drill-up): summarize data

by climbing up hierarchy or by dimension reduction

- Drill down (roll down): reverse of roll-up

from higher level summary to lower level summary or detailed data, or introducing new dimensions

- Slice and dice: project and select

Pivot (rotate):

reorient the cube, visualization, 3D to series of 2D planes

- Other operations

Drill across: involving (across) more than one fact table

- Drill through: through the bottom level of the cube to its back-end relational tables (using SQL)

Efficient Data Cube Computation

- Data cube can be viewed as a lattice of cuboids

The bottom-most cuboid is the base cuboid

Why this formula?

$$T = \prod_{i=1}^n (L_i + 1)$$

- How many cuboids in an n -dimensional cube with L levels?

- Materialization of data cube

Full materialization: Materialize every (cuboid)

No materialization: Materialize none (cuboid)

Partial materialization: Materialize some cuboids

Which cuboids to materialize?

Selection based on size, sharing, access frequency, etc.

Indexing OLAP Data: Bitmap Index

- Index on a particular column

Each value in the column has a bit vector: bit-op is fast

- The length of the bit vector: # of records in the base table

The i -th bit is set if the i -th row

Cube Computation: Computing in Reverse Order

❑ BUC (Beyer & Ramakrishnan, SIGMOD'99)

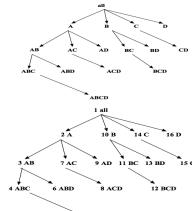
BUC: acronym of Bottom-Up (cube) Computation
(Note: It is "top-down" in our view since we put Apex cuboid on the top!)

❑ Divides dimensions into partitions and facilitates iceberg pruning

❑ If a partition does not satisfy min_sup , its descendants can be pruned

❑ If $\text{minsup} = 1$ compute full CUBE!

❑ No simultaneous aggregation



BUC: Partitioning and Aggregating

❑ Usually, entire data set cannot fit in main memory

❑ Sort distinct values

❑ partition into blocks that fit

❑ Continue processing

❑ Optimizations

❑ Partitioning

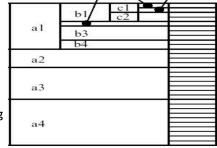
❑ External Sorting, Hashing, Counting Sort

❑ Ordering dimensions to encourage pruning

❑ Cardinality, Skew, Correlation

❑ Collapsing duplicates

❑ Cannot do holistic aggregates anymore!



Shell Fragment Cubes: Size and Design

❑ Given a database of T tuples, D dimensions, and F shell fragment size, the resulting cubes' space requirement is: $O\left(\tau \left[\frac{D}{F}\right] (2^F - 1)\right)$

❑ For $F < 5$, the growth is sub-linear

❑ Shell fragments do not have to be disjoint

❑ Fragment groupings can be arbitrary to allow for maximum online performance

❑ Known common combinations (e.g., city, state) should be grouped together

❑ Shell fragment sizes can be adjusted for optimal balance between offline and online computation

Cell	Intersection	TID List	List Size
a1 b1	1 2 3 n 4 5	1	1
a1 b2	1 2 3 n 4 5	2	2
a2 b1	1 2 3 n 4 5	4	2
a2 b2	1 2 3 n 4 5	φ	0

$(a_1, a_2, a_3, a_4, \dots, a_{10}) : 1, (a_1, b_2, a_3, b_4, \dots, b_{10}) : 1$

i. [3] How many nonempty cuboids are there in this data cube?

Answer: 2^{10} . Since we have 10 dimensions with no concept hierarchy, there are 2^{10} cuboids and all of them should not be empty.

ii. [3] How many (nonempty) aggregate closed cells are there in this data cube?

Answer: 1. There are 3 closed cells, including the two base cells and $(a_1, *, a_2, *, a_3, *, a_4, *, a_5, *, a_6, *, a_7, *, a_8, *)$. But only the latter one is a aggregated closed cell.

iii. [3] How many (nonempty) aggregate cells are there in this data cube?

Answer: 2014. For each base cell, there are 2^{10-1} aggregated cells. However, there are 2^5 cells that are counted twice since there are 5 common dimensions. Therefore, the total number of nonempty aggregate cells is $2 \cdot (2^{10} - 1) - 2^5 = 2014$.

iv. [3] If we set minimum support = 2, how many (nonempty) aggregate cells are there in the corresponding iceberg cube?

Answer: 2⁵. These two base cells have common value in 5 dimensions; therefore, there are 2^5 nonempty cells with support = 2 and all of them are aggregate cells.

From Frequent Itemsets to Association Rules

Tid	Items bought
10	Beer, Nuts, Diaper
20	Beer, Coffee, Diaper
30	Beer, Diaper, Eggs
40	Nuts, Eggs, Milk
50	Nuts, Coffee, Diaper, Eggs, Milk

Containing beer
Containing Diaper
 $\{\text{Beer} \cup \{\text{Beer}\} \cup \{\text{Diaper}\} = \{\text{Beer, Diaper}\}$

Note: Itemset: $X \cup Y$, a subtle notation!

Expressing Patterns in Compressed Form: Closed Patterns

❑ How to handle such a challenge?

❑ Solution 1: **Closed patterns**: A pattern (itemset) X is **closed** if X is frequent, and there exists no super-pattern Y $\supseteq X$, with the same support as X

❑ Let Transaction DB TDB₁: T₁: {a₁, ..., a₅₀}, T₂: {a₁, ..., a₁₀₀}

❑ Suppose $\text{minsup} = 1$. How many closed patterns does TDB₁ contain?

❑ Two: P₁: {"a₁, ..., a₅₀": 2}; P₂: {"a₁, ..., a₁₀₀": 1"}

❑ **Closed pattern** is a lossless compression of frequent patterns

❑ Reduces the # of patterns but does not lose the support information!

❑ You will still be able to say: {"a₁, ..., a₅₀": 2}; {"a₁, ..., a₅₁": 1"}

Expressing Patterns in Compressed Form: Max-Patterns

❑ Solution 2: **Max-patterns**: A pattern X is a **max-pattern** if X is frequent and there exists no frequent super-pattern Y $\supseteq X$

❑ Difference from close-patterns?

❑ Do not care the real support of the sub-patterns of a max-pattern

❑ Let Transaction DB TDB₁: T₁: {a₁, ..., a₅₀}, T₂: {a₁, ..., a₁₀₀}

❑ Suppose $\text{minsup} = 1$. How many max-patterns does TDB₁ contain?

❑ One: P: {"a₁, ..., a₁₀₀": 1"}

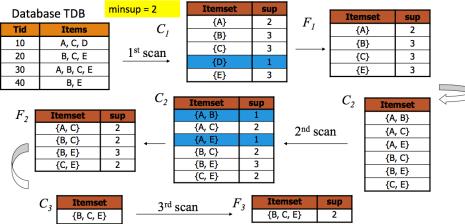
❑ **Max-pattern** is a lossy compression!

❑ We only know {a₁, ..., a₅₀} is frequent

❑ But we do not know the real support of {a₁, ..., a₁₀₀}, any more!

❑ Thus in many applications, mining close-patterns is more desirable than mining max-patterns

The Apriori Algorithm—An Example



Apriori: Implementation Tricks

❑ How to generate candidates?

❑ Step 1: self-joining F_k

❑ Step 2: pruning

❑ Example of candidate-generation

❑ $F_3 = \{abc, abd, acd, ace, bcd\}$

❑ Self-joining: $F_3 * F_3$

❑ abc from abc and abd

❑ acde from acd and ace

❑ Pruning:

❑ acde is removed because ade is not in F_3

❑ $C_4 = \{abcd\}$

Exploring Vertical Data Format: ECLAT

❑ ECLAT (Equivalence Class Transformation): A depth-first search algorithm using set intersection (Zaki et al. @KDD'97)

❑ Tid-List: List of transaction-ids containing an itemset

❑ Vertical format: t(e) = {T₁₀, T₂₀, T₃₀}; t(a) = {T₁₀, T₂₀}; t(ae) = {T₁₀, T₂₀}

❑ Properties of Tid-Lists

❑ t(x) = T(Y): X and Y always happen together (e.g., t(ac) = t(d))

❑ t(X) ⊂ t(Y): transaction having X always has Y (e.g., t(ac) ⊂ t(cce))

❑ Deriving frequent patterns based on vertical intersections

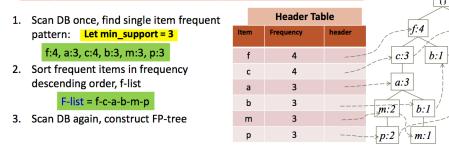
❑ Using diffset to accelerate mining

❑ Only keep track of differences of tids

❑ t(e) = {T₁₀, T₂₀, T₃₀}, t(c) = {T₁₀, T₃₀} → Diffset(t(c), e) = {T₁₀}

Example: Construct FP-tree from a Transactional DB

TID	Items in the Transaction	Ordered, frequent items
100	{f, a, c, i, g, k, m, p}	{f, c, a, m, p}
200	{a, b, c, f, l, m, o}	{f, c, a, b, m}
300	{b, f, h, j, k, o, w}	{f, b}
400	{b, k, l, m, o, p}	{b, k, p}
500	{a, f, c, i, l, m, o, p}	{f, c, a, m, p}



Interestingness Measure: Lift

❑ Measure of dependent/correlated events: lift

$lift(B, C) = \frac{s(B \cap C)}{s(B) \cdot s(C)}$

❑ Lift(B, C) may tell how B and C are correlated

❑ Lift(B, C) = 1: B and C are independent

❑ > 1: positively correlated

❑ < 1: negatively correlated

❑ For our example, $lift(B, C) = \frac{400/1000}{600/1000 \times 750/1000} = 0.89$

$lift(B, -C) = \frac{200/1000}{600/1000 \times 250/1000} = 1.33$

❑ Thus, B and C are negatively correlated since lift(B, C) < 1;

❑ B and -C are positively correlated since lift(B, -C) > 1

❑ Lift is more telling than the support-confidence framework

Lift is more telling than s & c

B	-B	Σ_{all}
C	400	350
-C	200	50
Σ_{all}	600	400

❑ Expected value

❑ Observed value

Interestingness Measures & Null-Invariance

❑ Null Invariance: Value does not change with the # of null-transactions

❑ A few interestingness measures: Some are null invariant

Measure	Definition	Range	Null-Invariant
$\chi^2(A, B)$	$\sum_{i,j=0,1} (e(a_i b_j) - o(a_i b_j))^2 / o(a_i b_j)$	[0, ∞]	No
$Lift(A, B)$	$\frac{s(A) \cdot s(B)}{s(A \cup B)}$	[0, ∞]	No
$AllConf(A, B)$	$\max\{s(A), s(B)\}$	[0, 1]	Yes
$Jaccard(A, B)$	$\frac{s(A \cap B)}{s(A) + s(B)}$	[0, 1]	Yes
$Cosine(A, B)$	$\frac{s(A) \cdot s(B)}{\sqrt{s(A) \cdot s(B)} \cdot \sqrt{s(B) \cdot s(A)}}$	[0, 1]	Yes
$Kulczynski(A, B)$	$\frac{1}{2} \left(\frac{s(A \cap B)}{s(A)} + \frac{s(A \cap B)}{s(B)} \right)$	[0, 1]	Yes
$MaxConf(A, B)$	$\max\left\{ \frac{s(A)}{s(A \cup B)}, \frac{s(B)}{s(A \cup B)} \right\}$	[0, 1]	Yes

X² and lift are not null-invariant

Jaccard, cosine, AllConf, MaxConf, and Kulczynski are null-invariant measures

Imbalance Ratio with Kulczynski Measure

❑ IR (Imbalance Ratio): measure the imbalance of two itemsets A and B in rule implications:

$$IR(A, B) = \frac{|s(A) - s(B)|}{s(A) + s(B) - s(A \cup B)}$$

❑ Kulczynski and Imbalance Ratio (IR) together present a clear picture for all the three datasets D₄ through D₆

❑ D₄ is neutral & balanced; D₅ is neutral but imbalanced

❑ D₆ is neutral but very imbalanced

Dataset	minsup	$\neg minsup$	$\neg \neg minsup$	Jaccard	Cosine	Kulc	IR
D ₄	1,000	1,000	1,000	0.78	0.99	0.91	0
D ₅	10,000	1,000	1,000	0.83	0.91	0.92	0
D ₆	1,000	1,000	100,000	0.05	0.69	0.09	0
D ₇	1,000	10	100,000	0.01	0.10	0.05	0.99

§ A database has 5 transactions. Let $\text{minsup} = 0.6$ and $\text{min.conf} = 0.7$.

trans.id	items bought
100	{K, A, D, B, C}
200	{D, A, E, F}
300	{C, D, B, E}
400	{B, A, C, K, D}
500	{B, G, C}

i. [2] List the frequent k-itemset for the largest k, and

Answer: k = 3

Frequent 3-itemset: {B, C, D} with support = 3/5 = 0.6 ≥ min.sup.

ii. [6] all the strong association rules (with support and confidence) for the following shape of rules:

$\forall x \in \text{transaction}, buys(x, item_1) \wedge buys(x, item_2) \Rightarrow buys(x, item_3)$. [s, c]

Answer:

$\forall x \in \text{transaction}, buys(x, B) \wedge buys(x, C) \Rightarrow buys(x, D)$. [0.6, 0.75]

$\forall x \in \text{transaction}, buys(x, B) \wedge buys(x, D) \Rightarrow buys(x, C)$. [0.6, 1]

$\forall x \in \text{transaction}, buys(x, C) \wedge buys(x, D) \Rightarrow buys(x, B)$. [0.6, 1]

[19] Data Warehousing and OLAP for Data Mining

(a) [9] Suppose the base cuboid of a data cube contains only two cells

$(a_1, a_2, a_3, \dots, a_{20}), (b_1, b_2, b_3, \dots, b_{20})$,

where $a_i = b_i$ if i is an odd number; otherwise $a_i \neq b_i$.

i. How many nonempty cuboids are there in this data cube?

Answer: $2^{20} = 1,048,576$

ii. How many nonempty aggregate cells are there in this data cube?

Answer: Each cell generates $2^{20} - 1$ nonempty aggregate cells. Thus in total we should have $2 \cdot 2^{20} - 2$ with overlapped cells not considering. We have 2^{21} overlapped cells (with a_i or b_i , fixed where i is odd number). Therefore, we have $2^{21} - 2 = 2^{20} - 2 = 1,024$.

iii. If we set minimum support = 2, how many nonempty aggregate cells are there in the corresponding iceberg cube?

Answer: Only the overlapped cells in the above question have the support 2, and there are no cells that have the support bigger than 2. Therefore, we have $2^{10} = 1,024$.

Partitioning Methods	<ul style="list-style-type: none"> - Find mutually exclusive clusters of spherical shape, - Distance based - May use mean or mediod (etc) to represent cluster center - Effective for small - to medium-size data sets
Hierarchical Methods	<ul style="list-style-type: none"> - Clustering is a hierarchical decomposition (i.e., multiple levels) - Cannot correct erroneous merges or splits - May incorporate other techniques like microclustering or consider object "linkages"
Density-based methods	<ul style="list-style-type: none"> - Can find arbitrarily shaped clusters - Clusters are dense regions of objects in space that are separated by low-density regions - Cluster density: Each point must have a minimum number of points within its "neighborhood" - May filter out outliers
Grid-based methods	<ul style="list-style-type: none"> - Use a multiresolution grid data structure - Fast processing time (typically independent of the number of data objects, yet dependent on grid size)

Partitioning Algorithms: Basic Concepts

- Partitioning method: Discovering the groupings in the data by optimizing a specific objective function and iteratively improving the quality of partitions
- Partitioning method: Partitioning a dataset D of n objects into a set of K clusters so that an objective function is optimized (e.g., the sum of squared distances is minimized, where c_i is the centroid or medoid of cluster C_i)
- A typical objective function: Sum of Squared Errors (SSE)

$$SSE(C) = \sum_{i=1}^n \|x_i - c_i\|^2$$
- Problem definition: Given K , find a partition of K clusters that optimizes the chosen partitioning criterion
- Global optimal: Needs to exhaustively enumerate all partitions
- Heuristic methods (i.e., greedy algorithms): K-Means, K-Medians, K-Medoids, etc.

The K-Means Clustering Method

- K-Means (MacQueen'67, Lloyd'77/82)
 - Each cluster is represented by the center of the cluster
- Given K , the number of clusters, the K-Means clustering algorithm is outlined as follows
 - Select K points as initial centroids
 - Repeat
 - Form K clusters by assigning each point to its closest centroid
 - Re-compute the centroids (i.e., mean point) of each cluster
 - Until convergence criterion is satisfied
- Different kinds of means are used
 - Manhattan distance (L_1 norm), Euclidean distance (L_2 norm), Cosine similarity

Handling Outliers: From K-Means to K-Medoids

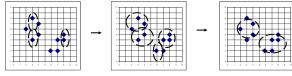
- The K-Means algorithm is sensitive to outliers—since an object with an extremely large value may substantially distort the distribution of the data
- K-Medoids: Instead of taking the mean value of the object in a cluster as a reference point, medoids can be used, which is the most centrally located object in a cluster
- The K-Medoids clustering algorithm:
 - Select K points as the initial representative objects (i.e., as initial K medoids)
 - Repeat
 - Assigning each point to the cluster with the closest medoid
 - Randomly select a non-representative object, o_j
 - Compute the total cost S of swapping the medoid m_i with o_j
 - If $S < 0$, then swap m_i with o_j to form the new set of medoids
 - Until convergence criterion is satisfied

K-Modes: Clustering Categorical Data

- K-Modes: An extension to K-Means by replacing means of clusters with modes
- Dissimilarity measure between object X and the center of a cluster Z
 - $\phi(X, Z) = 1 - n/n_i$, when $n_i = \sum_j x_{ij}$
 - where n_i is the categorical value of attribute i in Z , n_j is the number of objects in cluster i , and n is the number of objects whose attribute value is r
- This dissimilarity measure (distance function) is frequency-based
- Algorithm is still based on iterative object cluster assignment and centroid update
- A fuzzy K-Modes method is proposed to calculate a fuzzy cluster membership value for each object to each cluster
- A mixture of categorical and numerical data: Using a K-Prototype method

Agglomerative Clustering Algorithm

- AGNES (Agglomerative Nesting) (Kaufmann and Rousseeuw, 1990)
 - We use the single-link method and the dissimilarity matrix
 - Continuously merge nodes that have the least dissimilarity
 - Eventually all nodes belong to the same cluster



- Agglomerative clustering varies on different similarity measures among clusters
 - Single link (nearest neighbor)
 - Average link (group average)
 - Complete link (diameter)
 - Centroid link (centroid similarity)

Single Link vs. Complete Link in Hierarchical Clustering

- Single link (nearest neighbor)
 - The similarity between two clusters is the similarity between their most similar members
 - Local similarity-based: Emphasizing more on close regions, ignoring the overall structure of the clusters
 - Capable of clustering non-elliptical shaped group of objects
 - Sensitive to noise and outliers
- Complete link (diameter)
 - The similarity between two clusters is the similarity between their most dissimilar members
 - Merge two clusters to form one with the smallest diameter
 - Nonlocal in behavior, obtaining compact shaped clusters
 - Sensitive to outliers

Agglomerative Clustering: Average vs. Centroid Links

- Agglomerative clustering with average link
 - Average link: The average distance between an element in one cluster and an element in the other (i.e., all pairs in two clusters)
 - Expensive to compute
- Agglomerative clustering with centroid link
 - Centroid link: The distance between the centroids of two clusters

- Group-Averaged Agglomerative Clustering (GAAC)
 - Let two clusters C_a and C_b be merged into C_{ab} . The new centroid is:

$$c_{ab} = \frac{N_a c_a + N_b c_b}{N_a + N_b}$$
 - The similarity metric for GAAC is the average of their distances

- Agglomerative clustering with k-means
 - World's confusion: The increase in the value of the SSE for the clustering obtained by merging them into C_{ab} is:

$$\Delta_{ab} = \frac{N_a N_b}{N_a + N_b} d(c_{ab}, c_{ab})$$

Divisive Clustering

- DIANA (Divisive Analysis) (Kaufmann and Rousseeuw, 1990)
 - Implemented in some statistical analysis packages, e.g., SPSS
- Inverse order of AGNES: Eventually each node forms a cluster on its own
- Divisive clustering is a top-down approach
 - The process starts at the root with all the points as one cluster
 - It recursively splits the higher level clusters to build the dendrogram
 - Can be considered as a global approach
 - More efficient when compared with agglomerative clustering

Typical Clustering Methodologies (I)

- Distance-based methods
 - Partitioning algorithms: K-Means, K-Medians, K-Medoids
 - Hierarchical algorithms: Agglomerative vs. divisive methods
- Density-based and grid-based methods
 - Density-based: Data space is explored at a high-level of granularity and then post-processing to put together dense regions into an arbitrary shape
 - Grid-based: Individual regions of the data space are formed into a grid-like structure
- Probabilistic and generative models: Modeling data from a generative process
 - Assume a specific form of the generative model (e.g., mixture of Gaussians)
 - Model parameters are estimated with the Expectation-Maximization (EM) algorithm (using the available dataset, for a maximum likelihood fit)
 - Then estimate the generative probability of the underlying data points

Clustering Different Types of Data (I)

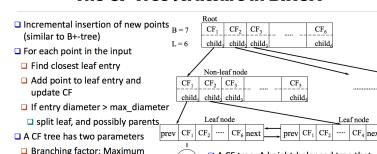
- Numerical data
 - Most earliest clustering algorithms were designed for numerical data
- Categorical data (including binary data)
 - Discrete data, no natural order (e.g., sex, race, zip-code, and market-basket)
- Text data: Popular in social media, Web, and social networks
 - Features: High-dimensional, sparse, value corresponding to word frequencies
 - Methods: Combination of k-means and agglomerative; topic modeling; co-clustering
- Multimedia data: Image, audio, video (e.g., on Flickr, YouTube)
- Multi-modal (often combined with text data)
 - Contextual: Containing both behavioral and contextual attributes
 - Images: Position of a pixel represents its context, value represents its behavior
 - Video and music data: Temporal ordering of records represents its meaning
- Time-series data: Sensor data, stock markets, temporal tracking, forecasting, etc.
 - Data are temporally dependent
 - Time: contextual attribute; value: behavioral attribute
 - Correlation-based online analysis (e.g., online clustering of stock to find stock tickers)
 - Shape-based offline analysis (e.g., cluster ECOS based on overall shapes)
- Sequence data: Weblogs, biological sequences, system command sequences
 - Contextual attribute: Placement (rather than time)
 - Similarity functions: Hamming distance, edit distance, longest common subsequence
 - Sequence clustering: Suffix tree; generative model (e.g., Hidden Markov Model)
- Stream data
 - Real-time, evolution and concept drift, single pass algorithm
 - Create efficient intermediate representation, e.g., micro-clustering

BIRCH (Balanced Iterative Reducing and Clustering Using Hierarchies)

- A multipath clustering algorithm (Zhang, Ramakrishnan, and Livny, SIGMOD'96)
 - Constructs a CF (Clustering Feature) tree, a hierarchical data structure for multipath clustering
- Phase 1: Scan DB to build an initial in-memory CF tree (a multi-level compression of the data that tries to preserve the inherent clustering structure of the data)
- Phase 2: Use an arbitrary clustering algorithm to cluster the leaf nodes of the CF-tree
- Key idea: Multi-level clustering
 - Low-level micro-clustering: Reduce complexity and increase scalability
 - High-level macro-clustering: Leave enough flexibility for high-level clustering
- Complexity: $O(N \log N)$ (If index-based)

- Where N = # of points

The CF Tree Structure in BIRCH



- An integration of agglomerative clustering with other (flexible) clustering methods
 - Low-level micro-clustering
 - Exploring CF-feature and BIRCH tree structure
 - Preserving the inherent clustering structure of the data
 - Higher-level macro-clustering
 - Provide sufficient flexibility for integration with other clustering methods
 - Impact to many other clustering methods and applications
 - Concerns
 - Sensitive to insertion order of data points
 - Due to the fixed size of leaf nodes, clusters may not be so natural
 - Clusters tend to be spherical given the radius and diameter measures

BIRCH: A Scalable and Flexible Clustering Method

- An integration of agglomerative clustering with other (flexible) clustering methods
 - Low-level micro-clustering
 - Exploring CF-feature and BIRCH tree structure
 - Preserving the inherent clustering structure of the data
 - Higher-level macro-clustering
 - Provide sufficient flexibility for integration with other clustering methods
 - Impact to many other clustering methods and applications
 - Concerns
 - Sensitive to insertion order of data points
 - Due to the fixed size of leaf nodes, clusters may not be so natural
 - Clusters tend to be spherical given the radius and diameter measures

CHAMELEON: Hierarchical Clustering Using Dynamic Modeling

- CHAMELEON: A graph partitioning approach (G. Karypis, E. H. Han, and V. Kumar, 1999)
- Measures the similarity based on a dynamic model
 - Two clusters are merged only if the **interconnectivity** and **closeness** (proximity) between two clusters are high relative to the internal interconnectivity of the clusters and closeness of items within the clusters

A graph-based, two-phase algorithm

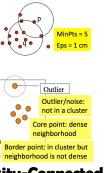
- Use a graph-partitioning algorithm: Cluster objects into a large number of relatively small sub-clusters
- Use an agglomerative hierarchical clustering algorithm: Find the genuine clusters by repeatedly combining these sub-clusters

Density-Based Clustering Methods

- Clustering based on density (a local cluster criterion), such as density-connected points
- Major features:
 - Discover clusters of arbitrary shape
 - Handle noise
 - One scan (only examine the local region to justify density)
 - Need density parameters as termination condition

DBSCAN: A Density-Based Spatial Clustering Algorithm

- DBSCAN (M. Ester, H.-P. Kriegel, J. Sander, and X. Xu, KDD'96)
 - Discover clusters of arbitrary shape
 - Clustering of Applications with Noise
- A density-based notion of cluster
 - A cluster is defined as a maximal set of density-connected points
- Two parameters:
 - ϵ : Maximum radius of the neighborhood
 - MinPts : Minimum number of points in the ϵ -neighborhood of a point
- The ϵ -neighborhood of a point p :
 - $N_{\epsilon}(p) = \{q | p \in N_{\epsilon}(q)\}$
 - $N_{\epsilon}(p) \cap S(p) = \{q | p \in N_{\epsilon}(q) \text{ and } q \in S(p)\}$



DBSCAN: Density-Reachable and Density-Connected

- Directly density-reachable:
 - A point p is directly density-reachable from a point q w.r.t. ϵ and MinPts if
 - p belongs to $N_{\epsilon}(q)$
 - $core$ point condition: $|N_{\epsilon}(q)| \geq \text{MinPts}$
- Density-reachable:
 - A point p is density-reachable from a point q w.r.t. ϵ and MinPts if there is a chain of points p_0, \dots, p_n such that $p_0 = q$ and p_n is a core point, and p_i is directly density-reachable from p_{i+1} for $i = 0, \dots, n-1$
- Density-connected:
 - A point p is density-connected to a point q w.r.t. ϵ and MinPts if there is a point r such that p and q are density-reachable from r w.r.t. ϵ and MinPts

DBSCAN: The Algorithm

- Algorithm
 - Arbitrarily select a point p
 - Retrieve all points density-reachable from p w.r.t. ϵ and MinPts
 - If p is a core point, a cluster is formed
 - If p is a border point, no points are density-reachable from p , and DBSCAN visits the next point of the database
 - Continue the process until all of the points have been processed

- Computational complexity
 - If a spatial index is used, the computational complexity of DBSCAN is $O(n \log n)$, where n is the number of database objects
 - Otherwise, the complexity is $O(n^2)$

OPTICS: Ordering Points to Identify Clustering Structure

- OPTICS (Ankerst, Breunig, Kriegel, and Sander, SIGMOD'99)
 - DBSCAN is sensitive to parameter setting
 - An extension: Finding clustering structure
 - Observation: Given a MinPts, density-based clusters w.r.t. a higher density are completely contained in clusters w.r.t. a lower density
 - Idea: Higher density points should be processed first—find high-density clusters first
 - OPTICS stores such a clustering order using two pieces of information:
 - Core distance and reachability distance

Reachability plot for a dataset

Figure 10.10: OPTICS ordering based on DBSCAN

Reachability distance undefined

Cluster-order of the objects

Core-distance undefined

Border point: c is not in a cluster

Border point: c is in a cluster but neighborhood is not dense

Border point: c is in a cluster and neighborhood is not dense

Border point: c is in a cluster and neighborhood is not dense

Border point: c is in a cluster and neighborhood is not dense

Border point: c is in a cluster and neighborhood is not dense

Border point: c is in a cluster and neighborhood is not dense

Border point: c is in a cluster and neighborhood is not dense

Border point: c is in a cluster and neighborhood is not dense

Border point: c is in a cluster and neighborhood is not dense

Border point: c is in a cluster and neighborhood is not dense

Border point: c is in a cluster and neighborhood is not dense

Border point: c is in a cluster and neighborhood is not dense

Border point: c is in a cluster and neighborhood is not dense

Border point: c is in a cluster and neighborhood is not dense

Border point: c is in a cluster and neighborhood is not dense

Border point: c is in a cluster and neighborhood is not dense

Border point: c is in a cluster and neighborhood is not dense

Border point: c is in a cluster and neighborhood is not dense

Border point: c is in a cluster and neighborhood is not dense

Border point: c is in a cluster and neighborhood is not dense

Border point: c is in a cluster and neighborhood is not dense

Border point: c is in a cluster and neighborhood is not dense

Border point: c is in a cluster and neighborhood is not dense

Border point: c is in a cluster and neighborhood is not dense

Border point: c is in a cluster and neighborhood is not dense

Border point: c is in a cluster and neighborhood is not dense

Border point: c is in a cluster and neighborhood is not dense

Border point: c is in a cluster and neighborhood is not dense

Border point: c is in a cluster and neighborhood is not dense

Border point: c is in a cluster and neighborhood is not dense

Border point: c is in a cluster and neighborhood is not dense

Border point: c is in a cluster and neighborhood is not dense

Border point: c is in a cluster and neighborhood is not dense

Border point: c is in a cluster and neighborhood is not dense

Border point: c is in a cluster and neighborhood is not dense

Border point: c is in a cluster and neighborhood is not dense

Border point: c is in a cluster and neighborhood is not dense

Border point: c is in a cluster and neighborhood is not dense

Border point: c is in a cluster and neighborhood is not dense

Border point: c is in a cluster and neighborhood is not dense

Border point: c is in a cluster and neighborhood is not dense

Border point: c is in a cluster and neighborhood is not dense

Border point: c is in a cluster and neighborhood is not dense

Border point: c is in a cluster and neighborhood is not dense

Border point: c is in a cluster and neighborhood is not dense

Border point: c is in a cluster and neighborhood is not dense

Border point: c is in a cluster and neighborhood is not dense

Border point: c is in a cluster and neighborhood is not dense

Border point: c is in a cluster and neighborhood is not dense

Border point: c is in a cluster and neighborhood is not dense

Border point: c is in a cluster and neighborhood is not dense

Border point: c is in a cluster and neighborhood is not dense

Border point: c is in a cluster and neighborhood is not dense

Border point: c is in a cluster and neighborhood is not dense

Border point: c is in a cluster and neighborhood is not dense

Border point: c is in a cluster and neighborhood is not dense

Border point: c is in a cluster and neighborhood is not dense

Border point: c is in a cluster and neighborhood is not dense

Border point: c is in a cluster and neighborhood is not dense

Border point: c is in a cluster and neighborhood is not dense

Border point: c is in a cluster and neighborhood is not dense

Border point: c is in a cluster and neighborhood is not dense

Border point: c is in a cluster and neighborhood is not dense

Border point: c is in a cluster and neighborhood is not dense

Border point: c is in a cluster and neighborhood is not dense

Border point: c is in a cluster and neighborhood is not dense

Border point: c is in a cluster and neighborhood is not dense

Border point: c is in a cluster and neighborhood is not dense

Border point: c is in a cluster and neighborhood is not dense

Border point: c is in a cluster and neighborhood is not dense

Border point: c is in a cluster and neighborhood is not dense

Border point: c is in a cluster and neighborhood is not dense

Border point: c is in a cluster and neighborhood is not dense

Border point: c is in a cluster and neighborhood is not dense

Border point: c is in a cluster and neighborhood is not dense

Border point: c is in a cluster and neighborhood is not dense

Border point: c is in a cluster and neighborhood is not dense

</div

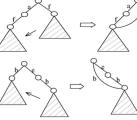
PrefixSpan: A Pattern-Growth Approach

SD Sequence	Prefix Suffix (Projection)
10 <abcde>	<>, <abcde>
20 <abcde>	<>, <bcde>, <cde>, <de>, <e>
30 <abcde>	<>, <abcde>, <bcde>, <cde>, <de>, <e>
40 <efghij>	<>, <efghij>, <fghij>, <ghij>, <hij>, <j>

- PrefixSpan Mining: Pre-Projections
- Step 1: Find length-1 sequential patterns
- <>, <abc>, <bc>, <cde>, <de>, <e>
- Step 2: search space and mine each projected DB

CloSpan: Mining Closed Sequential Patterns

- A **closed sequential pattern** S : There exists no superpattern S' such that $S' \supset S$, and S' and S have the same support.
- Which ones are closed? <abc>, <bc>, <abcde>, <bcde>, <de>, <e>
- Why directly mine closed sequential patterns?
- Reduce # of redundant patterns
- Attain the same expressive power
- Property P_i : if $S \supset S_i$, S is closed iff two project DBs have the same size
- Explore Backward Subspan and Backward Superpattern pruning to prune redundant search space
- Greatly enhances efficiency (Yan, et al., SDM'03)

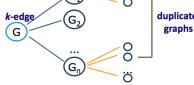


Candidate Generation: Vertex Growing vs. Edge Growing

- Methodology: breadth-search, Apriori joining two-size-k graphs
- Many possibilities at generating size-(k+1) candidate graphs
- Generating new graphs with one more vertex
 - AGM (Inokuchi, et al., POD'00)
 - Generating new graphs with one more edge
 - FSG (Kuramochi and Karypis, ICDM'03)
 - Performance shows via edge growing is more efficient

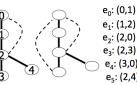
Pattern-Growth Approach

- Depth-first growth of subgraphs from k -edge to $(k+1)$ -edge, then $(k+2)$ -edge subgraphs
- Major challenge
- Generating many duplicate subgraphs
- Major idea to solve the problem
 - Define an order to generate subgraphs
 - DFS spanning tree: Flatten a graph into a sequence with depth-first search
- gSpan (Yan & Han; ICDM'02)



gSPAN: Graph Pattern Growth in Order

- Right-most path extension in subgraph pattern growth
- Right-most path: The path from root to the right-most leaf (choose the vertex w. the smallest index at each step)
- Reduce generation of duplicate subgraphs
- Completeness: The Enumeration of graphs using right-most path extension is complete
- DFS Code: Flatten a graph into a sequence using depth-first search



Why Mining Closed Graph Patterns?

- Challenge: An n -edge frequent graph may have 2^n subgraphs
- Motivation: Explore closed frequent subgraphs to handle graph pattern explosion problem
- A frequent graph G is closed if there exists no supergraph of G that carries the same support as G
- If this subgraph is closed in the graph dataset, it implies none of its frequent supergraphs carries the same support
- Lossless compression: Does not contain non-closed graphs, but still ensures that the mining result is complete
- Algorithm CloseGraph: Mines closed graph patterns directly

Chapter 8

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 $\frac{|C_i|}{|D|}/|D|$
- Expected information (entropy) needed to classify a tuple in D :
$$Info(D) = -\sum_{i=1}^n p_i \log_2(p_i)$$
- Information needed (after using A to split D into V partitions) to classify D :
$$Info_A(D) = \sum_{j=1}^V |D_j| / |D| \times Info(D_j)$$
- Information gained by branching on attribute A

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

From: Data mining and machine learning

Attribute Selection: Information Gain

- Class P: buys_computer = "yes"
- Class N: buys_computer = "no"
- Info(D) = $I(D) = \frac{2}{14} \log_2(\frac{2}{14}) + \frac{12}{14} \log_2(\frac{12}{14}) = 0.940$
- Age: 14 (0.971), 2 (0.971), 3 (0.971), 4 (0.971), 5 (0.971), 6 (0.971), 7 (0.971), 8 (0.971), 9 (0.971), 10 (0.971), 11 (0.971), 12 (0.971), 13 (0.971), 14 (0.971), 15 (0.971), 16 (0.971), 17 (0.971), 18 (0.971), 19 (0.971), 20 (0.971), 21 (0.971), 22 (0.971), 23 (0.971), 24 (0.971), 25 (0.971), 26 (0.971), 27 (0.971), 28 (0.971), 29 (0.971), 30 (0.971), 31 (0.971), 32 (0.971), 33 (0.971), 34 (0.971), 35 (0.971), 36 (0.971), 37 (0.971), 38 (0.971), 39 (0.971), 40 (0.971), 41 (0.971), 42 (0.971), 43 (0.971), 44 (0.971), 45 (0.971), 46 (0.971), 47 (0.971), 48 (0.971), 49 (0.971), 50 (0.971), 51 (0.971), 52 (0.971), 53 (0.971), 54 (0.971), 55 (0.971), 56 (0.971), 57 (0.971), 58 (0.971), 59 (0.971), 60 (0.971), 61 (0.971), 62 (0.971), 63 (0.971), 64 (0.971), 65 (0.971), 66 (0.971), 67 (0.971), 68 (0.971), 69 (0.971), 70 (0.971), 71 (0.971), 72 (0.971), 73 (0.971), 74 (0.971), 75 (0.971), 76 (0.971), 77 (0.971), 78 (0.971), 79 (0.971), 80 (0.971), 81 (0.971), 82 (0.971), 83 (0.971), 84 (0.971), 85 (0.971), 86 (0.971), 87 (0.971), 88 (0.971), 89 (0.971), 90 (0.971), 91 (0.971), 92 (0.971), 93 (0.971), 94 (0.971), 95 (0.971), 96 (0.971), 97 (0.971), 98 (0.971), 99 (0.971), 100 (0.971)
- Age <= 30 means "age <= 30" has 5 out of 14 samples, with 2 yes's and 3 no's. Hence, $Gain(age) = Info(D) - Info_A(D) = 0.246$
- Similarly,

$$Gain(income) = 0.029$$

$$Gain(student) = -0.151$$

$$Gain(credit_rating) = 0.048$$

- Information gain measure is biased towards attributes with a large number of values
- C4.5 (a successor of ID3) uses **gain ratio** to overcome the problem (normalization to information gain)

$$GainRatio(A) = Gain(A)/SplitInfo(A)$$

- Ex. $SplitInfo_{income}(D) = -\frac{4}{14} \log_2(\frac{4}{14}) - \frac{6}{14} \log_2(\frac{6}{14}) - \frac{4}{14} \times \log_2(\frac{4}{14}) = 1.557$
- gain_ratio(income) = $0.029/1.557 = 0.019$
- The attribute with the maximum gain ratio is selected as the splitting attribute

Still difficult for class imbalance problem on multiclass tasks

Gini Index (CART, IBM IntelligentMiner)

- If a data set D contains examples from n classes, gini index, $gini(D)$ is defined as
$$gini(D) = 1 - \sum_{i=1}^n p_i^2$$
- where p_i is the relative frequency of class i in D
- If a data set D is split on A into two subsets D_1 and D_2 , the gini index $gini(D)$ is defined as
$$gini_A(D) = |D_1|/|D| \cdot gini(D_1) + |D_2|/|D| \cdot gini(D_2)$$

$$Gini(A) = gini(D) - gini_A(D)$$

The attribute provides the **smallest gini(D)** (or the **largest reduction in impurity**) is chosen to split the node (**need to enumerate all the possible splitting points for each attribute**)

Comparing Attribute Selection Measures

The three measures, in general, return good results but

- Information gain:**
 - biased towards multivalued attributes
- Gain ratio:**
 - tends to prefer unbalanced splits in which one partition is much smaller than the others
- Gini index:**
 - biased to multivalued attributes
 - has difficulty when # of classes is large
 - tends to favor tests that result in equal-sized partitions and purity in both partitions

Overfitting and Tree Pruning

- Overfitting:** An induced tree may overfit the training data
 - Too many branches, some may reflect anomalies due to noise or outliers
 - Poor accuracy for unseen samples
- Two approaches to avoid overfitting
 - Pruning:** *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
- Post-pruning:** 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

using testing and prior knowledge on other data

Classifier Evaluation Metrics: Accuracy, Error Rate, Sensitivity and Specificity

Classifier Accuracy, Error Rate, Sensitivity and Specificity

TP	FP	FN	TN	P
True Positives	False Positives	False Negatives	True Negatives	Total Population
True Positives	False Positives	False Negatives	True Negatives	Total Population
True Positives	False Positives	False Negatives	True Negatives	Total Population

- Classifier Accuracy:** percentage of test set tuples that are correctly classified
- Accuracy = $TP + TN / P$
- Error rate: $1 - accuracy$, or
- Error rate = $(FP + FN) / P$

Classifier Evaluation Metrics: Precision and Recall, and F-measures

- Precision: exactness: what % of tuples that the classifier labeled as positive are actually positive

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

- Recall: completeness – what % of positive tuples did the classifier label as positive

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

- Perfect score is 1.0
- Inverse relationship between precision & recall

- F measure (or f-score):** harmonic mean of precision and recall

In general, it is the weighted measure of precision & recall

$$F = \frac{1}{\frac{1}{precision} + \frac{1}{recall}} = \frac{(beta^2 + 1) * PR}{beta^2 * PR + R}$$

Assigning β times as much weight to recall as to precision

F1-measure (balanced F-measure)

- That is, when $\beta = 1$,

$$F_1 = \frac{2PR}{P+R}$$

Holdout method

- Given data is randomly partitioned into two independent sets
 - Training set (e.g., 2/3) for model construction
 - Test set (e.g., 1/3) for accuracy estimation
- Random sampling: a variation of holdout
- Repeat holdout k times, accuracy = avg. of the accuracies obtained

Cross-validation (k fold, where $k = 10$ is most popular)

- Randomly partition the data into k mutually exclusive subsets, each approximately equal size
- At i -th iteration, use D_i as test set and others as training set
- Leave-one-out: k folds where $k = #$ of tuples, for small sized data

Stratified cross-validation*: folds are stratified so that class dist. in each fold is approx. the same as that in the initial data

Bootstrap

- Works well with small data sets
- Samples the given training tuples uniformly with replacement
- Each time a tuple is selected, it is equally likely to be selected again and re-added to the training set

Popular ensemble methods

- Bagging:** averaging the prediction over a collection of classifiers
- Boosting:** weighted vote with a collection of classifiers
- Ensemble:** combining a set of heterogeneous classifiers

Adaboost (Freund and Schapire, 1997)

- Given a set of d class-labeled tuples, $(X_1, Y_1), \dots, (X_d, Y_d)$
- Initially, all the weights of tuples are set the same ($1/d$)

- Generate K classifiers in K rounds. At round i ,

- Tuple from D are sampled (with replacement) to form a training set D_i of the same size
- Each tuple's chance of being selected is based on its weight
- A classification model M_i is derived from D_i .
- Its error rate is calculated using D_i as a test set
- If a tuple is misclassified, its weight is increased, o.w. it is decreased

- Error rate(X_j) is the misclassification error of tuple X_j . Classifier M_i , error rate is the sum of the weights of the misclassified tuples:

$$error(M_i) = \sum_j w_j \times err(X_j)$$

- The weight of classifier M_i 's vote is

$$\log \frac{1 - error(M_i)}{error(M_i)}$$

- Two Methods to construct Random Forest:
 - Forest-RI (random selection): Randomly select, at each node, F attributes as candidates for the split at the node. The CART methodology is used to grow the trees to maximum size
 - Forest-RC (random linear combinations): Creates new attributes (or features) that are a linear combination of the existing attributes (reduces the correlation between individual classifiers)

Classification of Class-Imbalanced Data Sets

- Class-imbalance problem: Rare positive example but numerous negative ones, e.g., medical diagnosis, fraud, oil-spill, fault, etc.

- Traditional methods assume a balanced distribution of classes and equal error costs: not suitable for class-imbalanced data

- Typical methods in two-class classification:
 - Oversampling: re-sampling of data from positive class
 - Under-sampling: randomly eliminate tuples from negative class
 - Threshold-moving: move the decision threshold, t, so that the rare class tuples are easier to classify, and hence, less chance of costly false negative errors

- Ensemble techniques:** Ensemble multiple classifiers introduced above

- Still difficult for class imbalance problem on multiclass tasks

Chapter 9 : Advanced Frequent Pattern Mining

- If more than one rule are triggered, need conflict resolution

Size ordering: assign the highest priority to the triggering rules that has the "toughest" requirement (i.e., with the *most attribute tests*)

Class-based ordering: decreasing order of prevalence or misclassification cost per class

Rule-based ordering (list): rules are organized into one long priority list, according to some measure of rule quality or expertise

Rule Extraction from a Decision Tree

- Rules are easier to understand than large trees

- 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 holds the class prediction

- Rules are mutually exclusive and exhaustive



Bayesian Network => based on conditional independence

Neural Network => skip Discriminative Classifier => Pro:

prediction acc higher than bayesian; robust when training contains error; fast evaluation of target function => Cons: long training time; difficult to understand learning function(bayesian can be used for pattern discovery); not easy to incorporate domain knowledge(prior, distribution)

SVM=> linear and nonlinear data => $wx+b=0$ (weight vector) => $w=0 \cdot w_1 x_1 + w_2 x_2 + \dots = 0$ why good? Complexity is # of support vector rather than dimension of data; support vector are critical training samples; good generalization even when high dimensionality.

CF-tree: Hierarchical Micro clustering => SVM not scalable, therefore => Clustering based SVM Algorithm: 1.

Construct two CF-trees. 2. Train SVM from centroid of root entries. 3. De-cluster the entries near the boundary into next level. (child entries declustered from parent are accumulated into training set with non-declustered parent entries) 4. Repeat until nothing is accumulated.

CF-tree: Hierarchical Micro clustering => SVM not scalable, therefore => Clustering based SVM Algorithm: 1.

Construct two CF-trees. 2. Train SVM from centroid of root entries. 3. De-cluster the entries near the boundary into next level. (child entries declustered from parent are accumulated into training set with non-declustered parent entries) 4. Repeat until nothing is accumulated.

Micro Cluster construction=> One scan of the data set: Construct two CF-trees (i.e., statistical summary of the data) from positive and negative classes

Micro Cluster construction=> One scan of the data set: Construct two CF-trees (i.e., statistical summary of the data) from positive and negative classes

Micro Cluster construction=> One scan of the data set: Construct two CF-trees (i.e., statistical summary of the data) from positive and negative classes

Micro Cluster construction=> One scan of the data set: Construct two CF-trees (i.e., statistical summary of the data) from positive and negative classes

Micro Cluster construction=> One scan of the data set: Construct two CF-trees (i.e., statistical summary of the data) from positive and negative classes

Micro Cluster construction=> One scan of the data set: Construct two CF-trees (i.e., statistical summary of the data) from positive and negative classes

Micro Cluster construction=> One scan of the data set: Construct two CF-trees (i.e., statistical summary of the data) from positive and negative classes

Micro Cluster construction=> One scan of the data set: Construct two CF-trees (i.e., statistical summary of the data) from positive and negative classes

Micro Cluster construction=> One scan of the data set: Construct two CF-trees (i.e., statistical summary of the data) from positive and negative classes

Micro Cluster construction=> One scan of the data set: Construct two CF-trees (i.e., statistical summary of the data) from positive and negative classes

Micro Cluster construction=> One scan of the data set: Construct two CF-trees (i.e., statistical summary of the data) from positive and negative classes

Micro Cluster construction=> One scan of the data set: Construct two CF-trees (i.e., statistical summary of the data) from positive and negative classes

Micro Cluster construction=> One scan of the data set: Construct two CF-trees (i.e., statistical summary of the data) from positive and negative classes

Micro Cluster construction=> One scan of the data set: Construct two CF-trees (i.e., statistical summary of the data) from positive and negative classes

Micro Cluster construction=> One scan of the data set: Construct two CF-trees (i.e., statistical summary of the data) from positive and negative classes

Micro Cluster construction=> One scan of the data set: Construct two CF-trees (i.e., statistical summary of the data) from positive and negative classes

Micro Cluster construction=> One scan of the data set: Construct two CF-trees (i.e., statistical summary of the data) from positive and negative classes

Micro Cluster construction=> One scan of the data set: Construct two CF-trees (i.e., statistical summary of the data) from positive and negative classes

Micro Cluster construction=> One scan of the data set: Construct two CF-trees (i.e., statistical summary of the data) from positive and negative classes

Micro Cluster construction=> One scan of the data set: Construct two CF-trees (i.e., statistical summary of the data) from positive and negative classes

Micro Cluster construction=> One scan of the data set: Construct two CF-trees (i.e., statistical summary of the data) from positive and negative classes

Micro Cluster construction=> One scan of the data set: Construct two CF-trees (i.e., statistical summary of the data) from positive and negative classes

Micro Cluster construction=> One scan of the data set: Construct two CF-trees (i.e., statistical summary of the data) from positive and negative classes

Micro Cluster construction=> One scan of the data set: Construct two CF-trees (i.e., statistical summary of the data) from positive and negative classes

Micro Cluster construction=> One scan of the data set: Construct two CF-trees (i.e., statistical summary of the data) from positive and negative classes

Micro Cluster construction=> One scan of the data set: Construct two CF-trees (i.e., statistical summary of the data) from positive and negative classes

Micro Cluster construction=> One scan of the data set: Construct two CF-trees (i.e., statistical summary of the data) from positive and negative classes

Micro Cluster construction=> One scan of the data set: Construct two CF-trees (i.e., statistical summary of the data) from positive