

Data Warehouse and Data Mining

Dhruv Gupta

dhruv.gupta@ntnu.no

07-February-2023



NTNU

Norwegian University of
Science and Technology

1 Announcements and References

- Administrative
- References for Today's Lecture

2 Data

- Feature Selection
- Feature Creation
- Discretization and Binarization
- Attribute Transformations
- Similarity and Dissimilarity
- Correlation

3 Association Rules

- Association Rule Mining
- Definitions
- Problem Definition
- Frequent Itemset Generation
- The Apriori Algorithm
 - Efficient Candidate Generation
 - Complexity Factors

4 Summary

1 Announcements and References

- Administrative
- References for Today's Lecture

2 Data

- Feature Selection
- Feature Creation
- Discretization and Binarization
- Attribute Transformations
- Similarity and Dissimilarity
- Correlation

3 Association Rules

- Association Rule Mining
- Definitions
- Problem Definition
- Frequent Itemset Generation
- The Apriori Algorithm
 - Efficient Candidate Generation
 - Complexity Factors

4 Summary

Administrative

1 First Assignment

- Due by 09.February.2023.

2 Volunteers for feedback regarding course

- Interested? Please contact me by email!

1 Announcements and References

- Administrative
- References for Today's Lecture

2 Data

- Feature Selection
- Feature Creation
- Discretization and Binarization
- Attribute Transformations
- Similarity and Dissimilarity
- Correlation

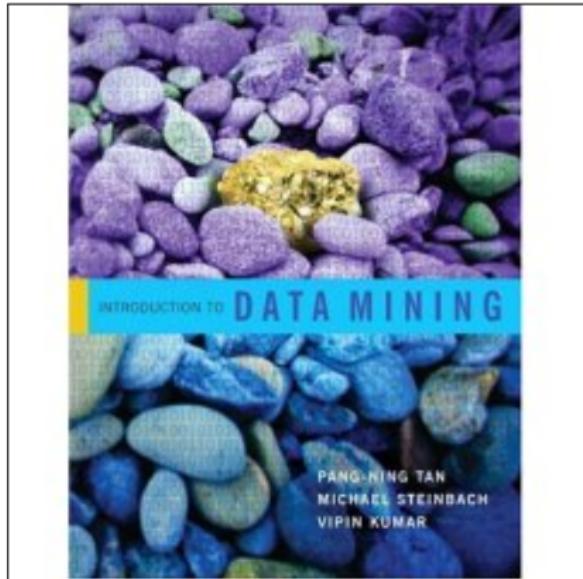
3 Association Rules

- Association Rule Mining
- Definitions
- Problem Definition
- Frequent Itemset Generation
- The Apriori Algorithm
 - Efficient Candidate Generation
 - Complexity Factors

4 Summary

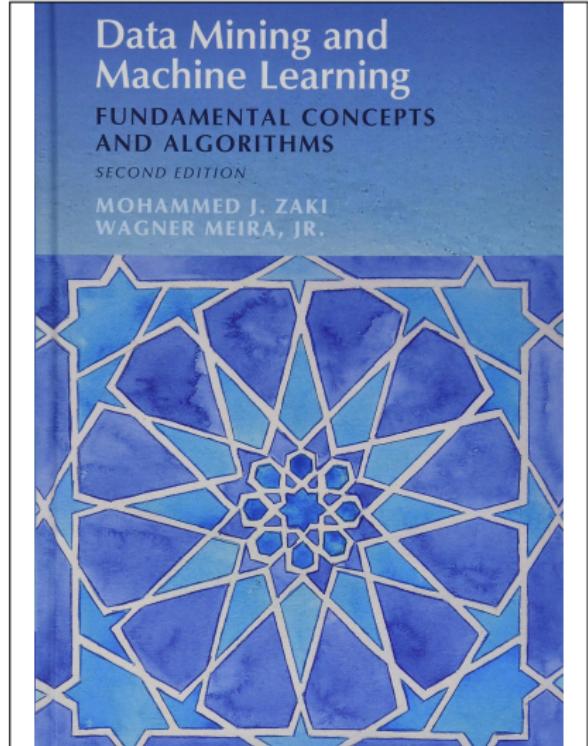
References for "Data" and "Association Rules"

- 1 Book: Tan et al. *"Introduction to Data Mining"*, 1st Edition, 2006, Pearson Education Inc.
- 2 Text and images for majority of slides in "Data" and "Association Rules" subsection are based on the book by Tan et al.



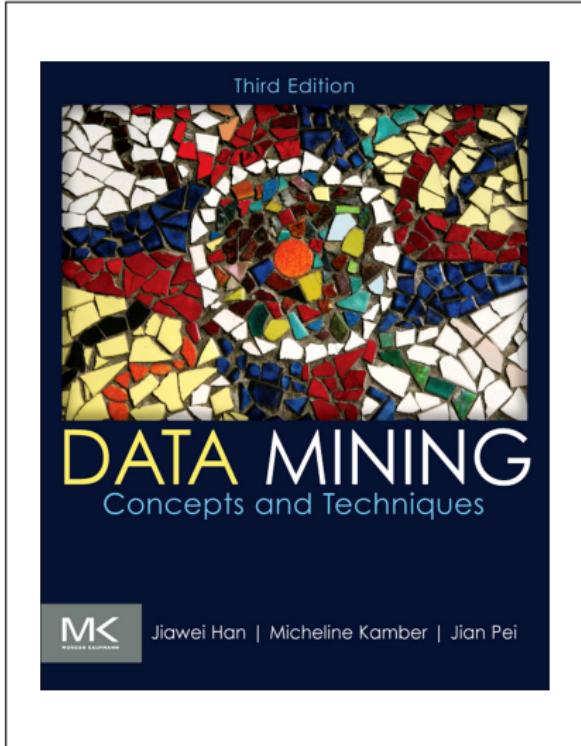
References for "Data" and "Association Rules"

- 1 Book: Zaki and Meira. *"Data Mining and Machine Learning: Fundamental Concepts and Algorithms"*, 2nd Edition, 2020, Cambridge University Press.
- 2 Text and images for some slides in "Data" and "Association Rules" subsection are based on the book by Zaki and Meira.



References for "Data" and "Association Rules"

- 1 Book: Han et al. *"Data Mining Concepts and Techniques"*, 3rd Edition, 2012, Morgan Kaufmann Publishers.
- 2 All text and images for some slides in "Data" and "Association Rules" are based on the book by Han et al.



1 Announcements and References

- Administrative
- References for Today's Lecture

2 Data

- Feature Selection
- Feature Creation
- Discretization and Binarization
- Attribute Transformations
- Similarity and Dissimilarity
- Correlation

3 Association Rules

- Association Rule Mining
- Definitions
- Problem Definition
- Frequent Itemset Generation
- The Apriori Algorithm
 - Efficient Candidate Generation
 - Complexity Factors

4 Summary

- ## 1 Announcements and References
- Administrative
 - References for Today's Lecture

- ## 2 Data
- Feature Selection
 - Feature Creation
 - Discretization and Binarization
 - Attribute Transformations
 - Similarity and Dissimilarity
 - Correlation

- ## 3 Association Rules
- Association Rule Mining
 - Definitions
 - Problem Definition
 - Frequent Itemset Generation
 - The Apriori Algorithm
 - Efficient Candidate Generation
 - Complexity Factors

- ## 4 Summary

Feature Selection

- Another way to **reduce dimensionality of data**.
- **Redundant features:**
 - Duplicate much or all of the information contained in one or more other attributes.
 - Example: purchase price of a product and the amount of sales tax paid.
- **Irrelevant features:**
 - Contain no information that is useful for the data mining task at hand.
 - Example: students' ID is often irrelevant to the task of predicting students' GPA.
 - Many techniques developed, especially for classification.

Feature Selection — Ideal Approach

- Irrelevant and redundant attributes can be eliminated by using common sense domain knowledge.
- A more systematic approach is however required.
- Ideal Approach:
 - Enumerate all possible subsets of features.
 - For each subset obtain performance results from the data mining algorithm being used.
- Pros: reflects the objective and bias of the data mining algorithm that will eventually be used.
- Cons: number of subsets involving n attributes is 2^n .

Feature Selection — Other Approaches

- **Embedded Approaches:** algorithm itself decides which attributes to use and which to ignore (e.g., decision trees).
- **Filter Approaches:** features are selected independently and separately from the algorithm being run (e.g., where pairwise correlation is as low as possible).
- **Wrapper Approaches:** same as Ideal Approach but without enumerating all possible subsets.



Figure 2.11. Flowchart of a feature subset selection process.

- ## 1 Announcements and References
- Administrative
 - References for Today's Lecture

- ## 2 Data
- Feature Selection
 - **Feature Creation**
 - Discretization and Binarization
 - Attribute Transformations
 - Similarity and Dissimilarity
 - Correlation

- ## 3 Association Rules
- Association Rule Mining
 - Definitions
 - Problem Definition
 - Frequent Itemset Generation
 - The Apriori Algorithm
 - Efficient Candidate Generation
 - Complexity Factors

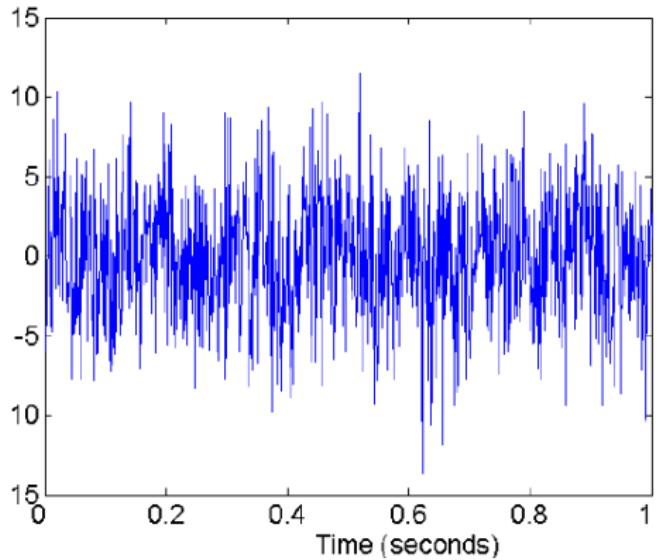
- ## 4 Summary

Feature Creation

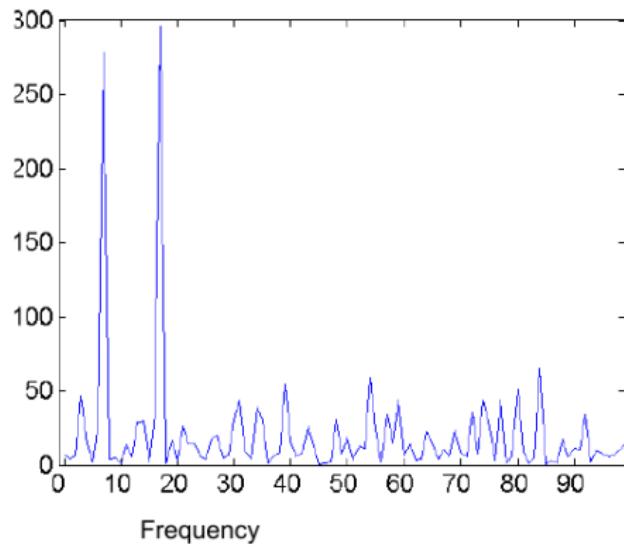
- Create new attributes that can capture the important information in a data set much more efficiently than the original attributes.
- Three general methodologies:
 - Feature extraction.
 - Example: extracting edges from images.
 - Feature construction.
 - Example: dividing mass by volume to get density.
 - Mapping data to new space.
 - Example: Fourier and wavelet analysis.

Feature Creation

- Example of mapping data to new space via Fourier Transform.



Two Sine Waves + Noise



Frequency

1 Announcements and References

- Administrative
- References for Today's Lecture

2 Data

- Feature Selection
- Feature Creation
- **Discretization and Binarization**
- Attribute Transformations
- Similarity and Dissimilarity
- Correlation

3 Association Rules

- Association Rule Mining
- Definitions
- Problem Definition
- Frequent Itemset Generation
- The Apriori Algorithm
 - Efficient Candidate Generation
 - Complexity Factors

4 Summary

Discretization

- Certain classification algorithms, require **data in the form of categorical attributes**.
- Algorithms that find association patterns require that the data be in the form of binary attributes.
- **Discretization**: transform a continuous attribute into a categorical attribute.
- **Binarization**: Both continuous and discrete attributes may need to be transformed into one or more binary attributes.

Binarization

- A simple approach to binarize m categorical values:
 - 1 Assign a unique integer in the interval $[0, m - 1]$. Maintain order during assignment if the attribute is ordinal.
 - 2 Convert each of these m integers to a binary number.
 - 3 $n = \lceil \log_2(m) \rceil$ binary attributes are required to represent these integers.
- Problems: unintended correlations introduced.

Categorical Value	Integer Value	x_1	x_2	x_3
awful	0	0	0	0
poor	1	0	0	1
okay	2	0	1	0
good	3	0	1	1
great	4	1	0	0

Binarization

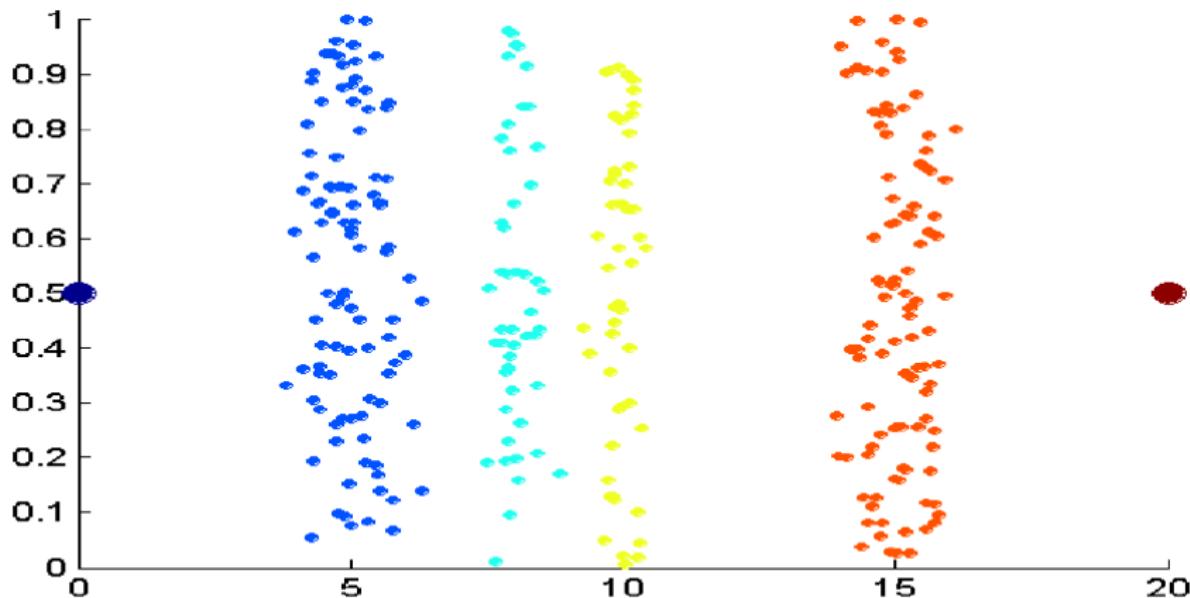
- For association problems **asymmetrical attributes are required.**
- Therefore necessary to introduce one binary attribute for each categorical value.

Categorical Value	Integer Value	x_1	x_2	x_3	x_4	x_5
awful	0	1	0	0	0	0
poor	1	0	1	0	0	0
okay	2	0	0	1	0	0
good	3	0	0	0	1	0
great	4	0	0	0	0	1

Discretization

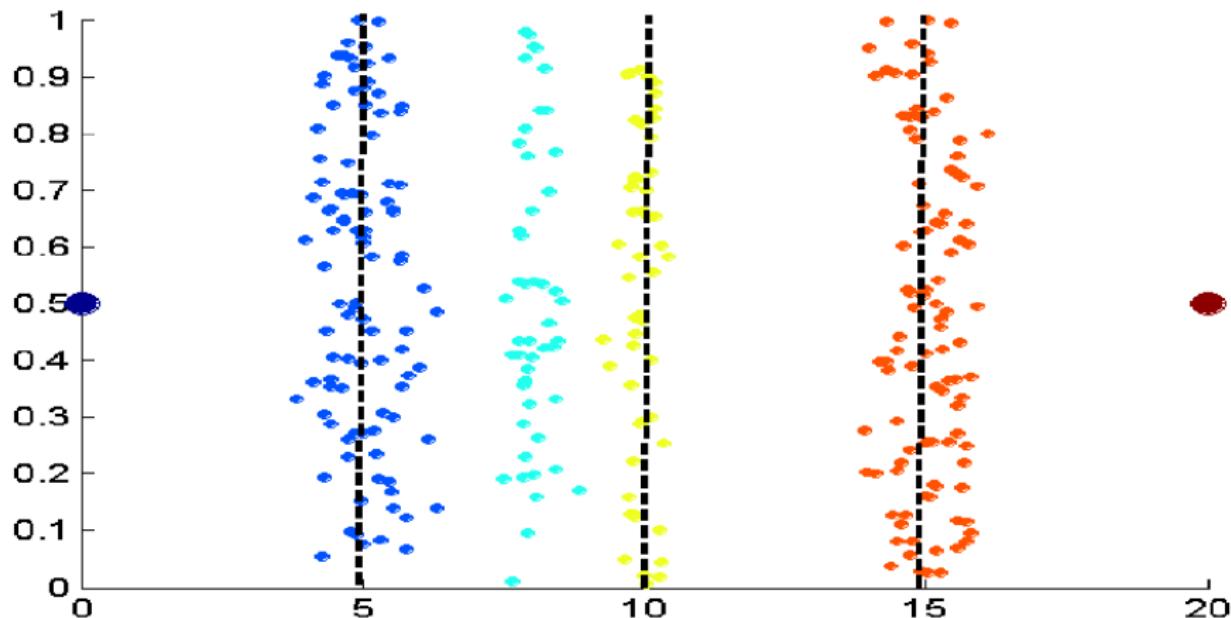
- Transformation of a **continuous attribute** to a **categorical attribute** involves two subtasks:
 - 1 **First step:** continuous attribute are sorted and then are divided into n intervals by specifying $n - 1$ split points.
 - 2 **Second step:** all the values in one interval are mapped to the same categorical value.
- **Unsupervised Approaches:** label associated with the objects are not used.
 - Example: equal width, equal frequency / equal depth, and k-means.
- **Supervised Approaches:** labels associated with the objects are used.
 - Example: a simple approach for partitioning a continuous attribute starts by bisecting the initial values so that the resulting two intervals give minimum entropy.

Discretization



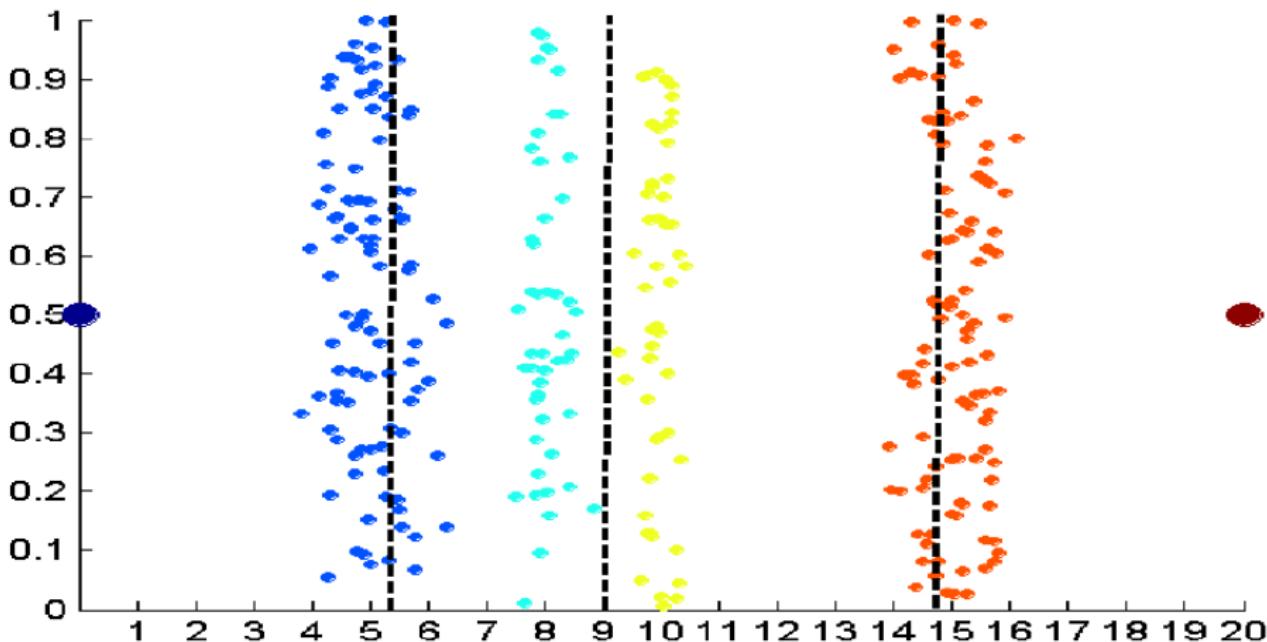
Data consists of four groups of points and two outliers. Data is one-dimensional, but a random y component is added to reduce overlap.

Discretization



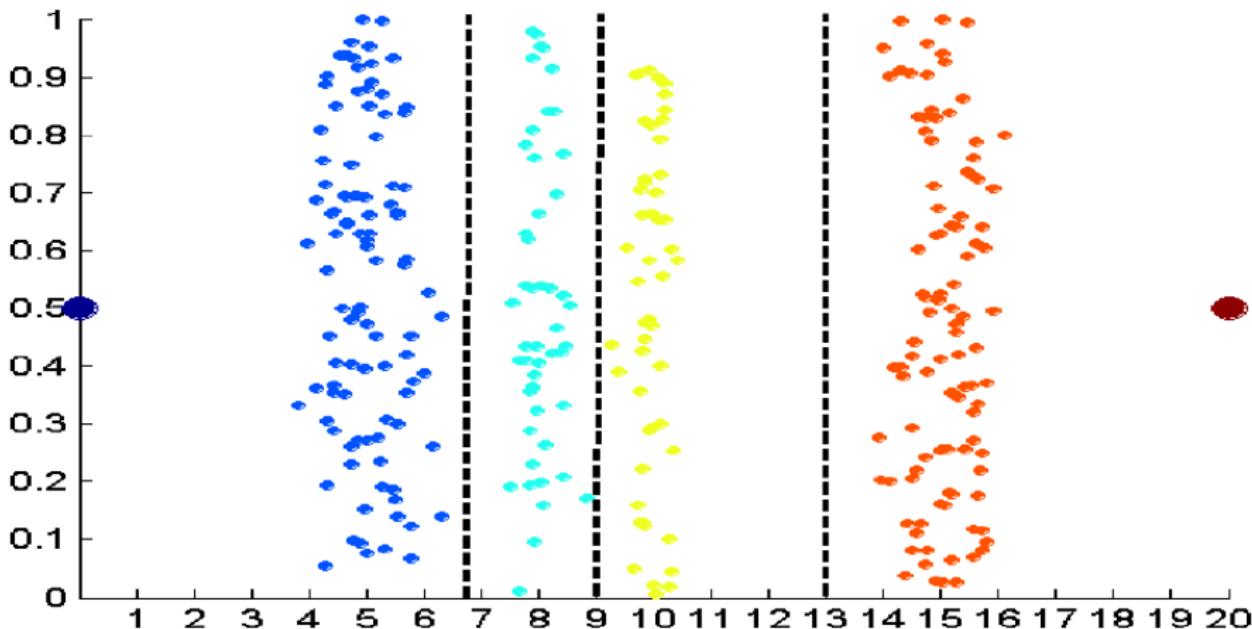
Equal interval width approach used to obtain 4 values.

Discretization



Equal frequency approach used to obtain 4 values.

Discretization



K-means approach to obtain 4 values.

1 Announcements and References

- Administrative
- References for Today's Lecture

2 Data

- Feature Selection
- Feature Creation
- Discretization and Binarization
- **Attribute Transformations**
- Similarity and Dissimilarity
- Correlation

3 Association Rules

- Association Rule Mining
- Definitions
- Problem Definition
- Frequent Itemset Generation
- The Apriori Algorithm
 - Efficient Candidate Generation
 - Complexity Factors

4 Summary

Simple Functions

- **Attribute Transform:** a function that maps the entire set of values of a given attribute to a new set of replacement values such that each old value can be identified with one of the new values.
- **Using Simple Functions:**
 - Using transformations such as x^k , $\log(x)$, e^x , \sqrt{x} , $1/x$, or $|x|$.
 - Variable transformations should be applied with caution since they change the nature of the data.
 - Does the transformation apply to all values?
 - Especially negative values and 0?
 - What is the effect of the transformation on the values between 0 and 1?

Normalization or Standardization

- Normalization refers to various techniques to **adjust to differences among attributes** in terms of **frequency of occurrence, mean, variance, and range.**
- Take out unwanted, common signal, e.g., seasonality.
- In statistics, **standardization** refers to **subtracting off the means and dividing by the standard deviation.**

$$x' = \frac{x - \bar{x}}{s_x} \quad (1)$$

1 Announcements and References

- Administrative
- References for Today's Lecture

2 Data

- Feature Selection
- Feature Creation
- Discretization and Binarization
- Attribute Transformations
- **Similarity and Dissimilarity**
- Correlation

3 Association Rules

- Association Rule Mining
- Definitions
- Problem Definition
- Frequent Itemset Generation
- The Apriori Algorithm
 - Efficient Candidate Generation
 - Complexity Factors

4 Summary

Similarity and Dissimilarity

- **Similarity** between two objects:
 - A **numerical measure** indicating degree to which the **two objects** are alike.
 - Usually **non-negative** and are often **between 0** (no similarity) and **1** (complete similarity).
- **Dissimilarity** between two objects:
 - A **numerical measure** of the degree to which the **two objects** are different.
 - Dissimilarities are lower for more similar pairs of objects.
 - Dissimilarities sometimes fall in the interval $[0,1]$, but it is also common for them to range from 0 to ∞ .
 - **Distances**, which are **dissimilarities** with certain properties.

Similarity and Dissimilarity for Simple Attributes

Attribute Type	Dissimilarity	Similarity
Nominal	$d = \begin{cases} 0 & \text{if } x = y \\ 1 & \text{if } x \neq y \end{cases}$	$s = \begin{cases} 1 & \text{if } x = y \\ 0 & \text{if } x \neq y \end{cases}$
Ordinal	$d = x - y / (n - 1)$ (values mapped to integers 0 to $n - 1$, where n is the number of values)	$s = 1 - d$
Interval or Ratio	$d = x - y $	$s = -d, s = \frac{1}{1+d}, s = e^{-d},$ $s = 1 - \frac{d - \min_d}{\max_d - \min_d}$

Figure: The following table shows the similarity and dissimilarity between two objects, x and y , with respect to a single, simple attribute.

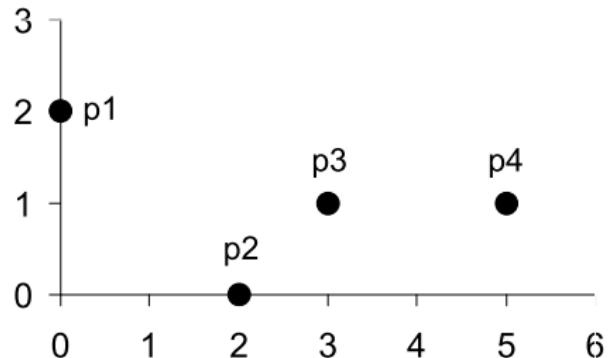
Distances (Dissimilarities with Certain Properties)

- Euclidean Distance:

$$d(x, y) = \sqrt{\sum_{k=1}^n (x_k - y_k)^2} \quad (2)$$

- where, n is the number of dimensions (attributes) and x_k and y_k are respectively the k^{th} attributes (components) or data objects x and y .
- Standardization is necessary if scales differ.

Euclidean Distance



point	x	y
p1	0	2
p2	2	0
p3	3	1
p4	5	1

	p1	p2	p3	p4
p1	0	2.828	3.162	5.099
p2	2.828	0	1.414	3.162
p3	3.162	1.414	0	2
p4	5.099	3.162	2	0

Distance Matrix

Minkowski Distance

- Minkowski Distance — generalization of Euclidean distance:

$$d(x, y) = \left(\sum_{k=1}^n |x_k - y_k|^r \right)^{1/r} \quad (3)$$

- where, r is a parameter, n is the number of dimensions (attributes) and x_k and y_k are, respectively, the k^{th} attributes (components) or data objects x and y .

Minkowski Distance: Examples

- $r = 1$. City block (Manhattan, taxicab, L_1 norm) distance.
 - A common example of this for binary vectors is the Hamming distance, which is just the number of bits that are different between two binary vectors
- $r = 2$. Euclidean distance.
- $r \rightarrow \infty$. “supremum” (L_{\max} norm, L_∞ norm, Uniform norm, and Chebyshev) distance.
 - This is the maximum difference between any component of the vectors.
 - That is,

$$d(x, y) = \lim_{r \rightarrow \infty} \left(\sum_{k=1}^n |x_k - y_k|^r \right)^{1/r} = \max_k^n |x_k - y_k|. \quad (4)$$

- Do not confuse r with n , i.e., all these distances are defined for all numbers of dimensions.

Minkowski Distance

point	x	y
p1	0	2
p2	2	0
p3	3	1
p4	5	1

L1	p1	p2	p3	p4
p1	0	4	4	6
p2	4	0	2	4
p3	4	2	0	2
p4	6	4	2	0

L2	p1	p2	p3	p4
p1	0	2.828	3.162	5.099
p2	2.828	0	1.414	3.162
p3	3.162	1.414	0	2
p4	5.099	3.162	2	0

L ∞	p1	p2	p3	p4
p1	0	2	3	5
p2	2	0	1	3
p3	3	1	0	2
p4	5	3	2	0

Distance Matrix

Common Properties of Distance

- Distances, such as the Euclidean distance, have some well known properties.
 - Positivity: $d(x, y) \geq 0$ for all x and y and $d(x, y) = 0$ if and only if $x = y$.
 - Symmetry: $d(x, y) = d(y, x)$ for all x and y .
 - Triangle Inequality: $d(x, z) \leq d(x, y) + d(y, z)$ for all points x , y , and z .
- A distance that satisfies these properties is a metric.

Similarity between Binary Vectors

- Common situation is that objects, x and y , have only binary attributes.
- Compute similarities using the following quantities:
 - f_{00} = the number of attributes where x was 0 and y was 0.
 - f_{01} = the number of attributes where x was 0 and y was 1.
 - f_{10} = the number of attributes where x was 1 and y was 0.
 - f_{11} = the number of attributes where x was 1 and y was 1.
- Simple Matching Coefficient:

$$SMC = \frac{\text{number of matches}}{\text{number of attributes}} = \frac{f_{11} + f_{00}}{f_{00} + f_{01} + f_{10} + f_{11}}. \quad (5)$$

- Jaccard Coefficient:

$$J = \frac{\text{number of } 11 \text{ matches}}{\text{number of non-zero attributes}} = \frac{f_{11}}{f_{01} + f_{10} + f_{11}}. \quad (6)$$

Simple Matching and Jaccard Coefficient Example

- $x = \langle 1, 0, 0, 0, 0, 0, 0, 0, 0, 0 \rangle$
- $y = \langle 0, 0, 0, 0, 0, 0, 1, 0, 0, 1 \rangle$
 - $f_{00} = 7$ the number of attributes where x was 0 and y was 0.
 - $f_{01} = 2$ the number of attributes where x was 0 and y was 1.
 - $f_{10} = 1$ the number of attributes where x was 1 and y was 0.
 - $f_{11} = 0$ the number of attributes where x was 1 and y was 1.
- Simple Matching Coefficient:

$$SMC = \frac{f_{11} + f_{00}}{f_{00} + f_{01} + f_{10} + f_{11}} = \frac{0 + 7}{7 + 2 + 1 + 0} = \frac{7}{10}. \quad (7)$$

- Jaccard Coefficient:

$$J = \frac{f_{11}}{f_{01} + f_{10} + f_{11}} = \frac{0}{2 + 1 + 0} = 0. \quad (8)$$

Recall — Document Data

- Each document becomes a term vector.
- Each term is a component (attribute) of the vector.
- The value of each component is the number of times the corresponding term occurs in the document.
- Example of sparse data matrix.

	word_1	word_2	word_3	\dots	$\text{word}_{ \mathcal{V} }$
document_1	10	0	0	\dots	5
document_2	1	2	0	\dots	0
\vdots	\vdots	\vdots	\vdots	\dots	\vdots
document_n	0	1	0	\dots	0

Cosine Similarity

- If x and y are two document vectors, then:

$$\cos(x, y) = \frac{\langle x, y \rangle}{\|x\| \cdot \|y\|} = \frac{\sum_{k=1}^n x_k \cdot y_k}{\sqrt{\sum_{k=1}^n x_k^2} \cdot \sqrt{\sum_{k=1}^n y_k^2}}. \quad (9)$$

- where, $\langle x, y \rangle$ indicates the inner product or vector dot product of vectors, x and y , and $\|x\|$ is the length of vector x .
- Example:

$$d_1 = \langle 3, 2, 0, 5, 0, 0, 0, 2, 0, 0 \rangle,$$

$$d_2 = \langle 1, 0, 0, 0, 0, 0, 0, 1, 0, 2 \rangle.$$

$$\langle d_1, d_2 \rangle = 3 \cdot 1 + 2 \cdot 0 + 0 \cdot 0 + 5 \cdot 0 + 0 \cdot 0 + 0 \cdot 0 + 0 \cdot 0 + 2 \cdot 1 + 0 \cdot 0 + 0 \cdot 2 = 5$$

$$\|d_1\| = \sqrt{(3 \cdot 3 + 2 \cdot 2 + 0 \cdot 0 + 5 \cdot 5 + 0 \cdot 0 + 0 \cdot 0 + 0 \cdot 0 + 2 \cdot 2 + 0 \cdot 0 + 0 \cdot 0)} = \sqrt{42} = 6.481$$

$$\|d_2\| = \sqrt{(1 \cdot 1 + 0 \cdot 0 + 1 \cdot 1 + 0 \cdot 0 + 2 \cdot 2)} = \sqrt{6} = 2.449$$

$$\cos(d_1, d_2) = \frac{5}{6.481 \cdot 2.449} = 0.3150$$

1 Announcements and References

- Administrative
- References for Today's Lecture

2 Data

- Feature Selection
- Feature Creation
- Discretization and Binarization
- Attribute Transformations
- Similarity and Dissimilarity
- Correlation

3 Association Rules

- Association Rule Mining
- Definitions
- Problem Definition
- Frequent Itemset Generation
- The Apriori Algorithm
 - Efficient Candidate Generation
 - Complexity Factors

4 Summary

Correlation

- Correlation between two data objects that have binary or continuous variables is a measure of the linear relationship between the attributes of the objects.
- Perfect Correlation: correlation is always in the range -1 to 1. A correlation of 1 (-1) means that x and y have a perfect positive (negative) linear relationship.
- That is, $x_k = a \cdot y_k + b$, where a , and b are constants.
- Example 1:

$$x = \langle -3, 6, 0, 3, -6 \rangle \text{ and } y = \langle 1, -2, 0, -1, 2 \rangle.$$

- where, $x = -3 \cdot y$ implies correlation is -1.
- Example 2:

$$x = \langle 3, 6, 0, 3, 6 \rangle \text{ and } y = \langle 1, 2, 0, 1, 2 \rangle.$$

- where, $x = 3 \cdot y$ implies correlation is 1.

Pearson's Correlation Coefficient

- Pearson's Correlation Coefficient between two data objects, x and y is defined by the following equation:

$$\text{corr}(x, y) = \frac{\text{covariance}(x, y)}{\text{std_dev}(x) \cdot \text{std_dev}(y)} = \frac{s_{xy}}{s_x \cdot s_y}. \quad (10)$$

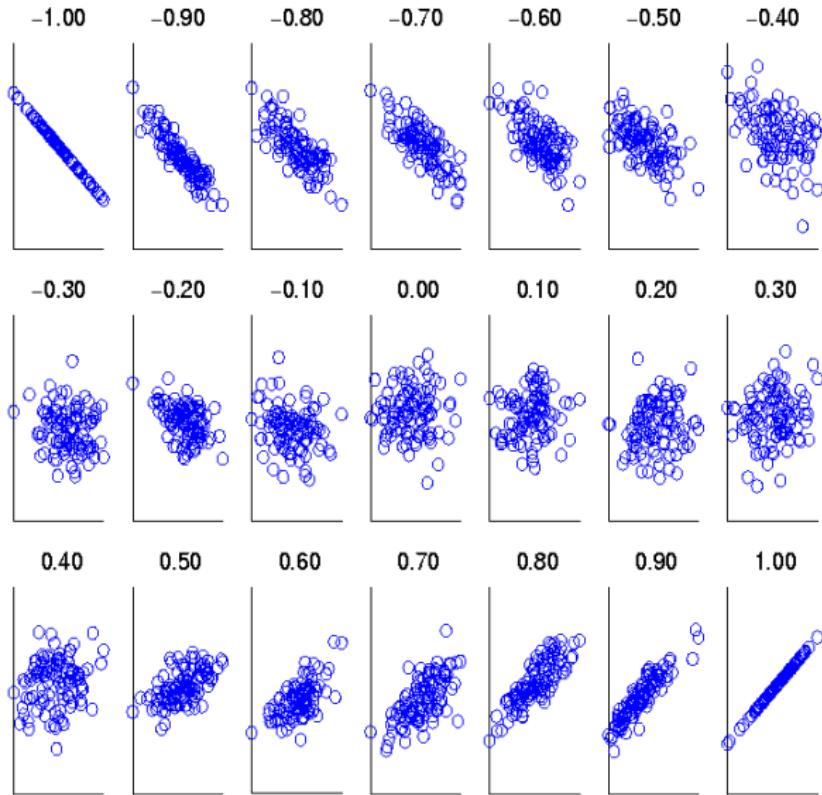
- where, the standard statistical notation and definitions are:

$$\text{covariance}(x, y) = s_{xy} = \frac{1}{n-1} \cdot \sum_{k=1}^n (x_k - \bar{x})(y_k - \bar{y}). \quad (11)$$

$$\text{std_dev}(x) = s_x = \sqrt{\frac{1}{n-1} \cdot \sum_{k=1}^n (x_k - \bar{x})^2}. \quad (12)$$

$$\bar{x} = \frac{1}{n} \sum_{k=1}^n x_k \text{ is the mean of } x. \quad (13)$$

Visualizing Correlations



**Scatter plots
showing the
similarity from
-1 to 1.**

1 Announcements and References

- Administrative
- References for Today's Lecture

2 Data

- Feature Selection
- Feature Creation
- Discretization and Binarization
- Attribute Transformations
- Similarity and Dissimilarity
- Correlation

3 Association Rules

- Association Rule Mining
- Definitions
- Problem Definition
- Frequent Itemset Generation
- The Apriori Algorithm
 - Efficient Candidate Generation
 - Complexity Factors

4 Summary

1 Announcements and References

- Administrative
- References for Today's Lecture

2 Data

- Feature Selection
- Feature Creation
- Discretization and Binarization
- Attribute Transformations
- Similarity and Dissimilarity
- Correlation

3 Association Rules

- **Association Rule Mining**
 - Definitions
 - Problem Definition
 - Frequent Itemset Generation
 - The Apriori Algorithm
 - Efficient Candidate Generation
 - Complexity Factors

4 Summary

Association Rule Mining

- 1 Association Analysis: discovering interesting relationships hidden in large data sets.
- 2 Uncovered relationships can be represented in the form of association rules or sets of frequent items.
- 3 When applied to retail transactions, commonly known as Market Basket Analysis.
- 4 Other application areas: medical diagnosis, Web mining, and scientific data analysis.

Association Rule Mining

- 1 Given a set of transactions, find rules that will predict the occurrence of an item based on the occurrences of other items in the transaction.
- 2 Example of association rules:
 - $\{Diaper\} \rightarrow \{Beer\}$
 - $\{Milk, Bread\} \rightarrow \{Eggs, Cola\}$
 - $\{Beer, Bread\} \rightarrow \{Milk\}$
- 3 Note that implication means co-occurrence, not causality.

T-ID	Items
1	Bread, Milk
2	Bread, Diaper, Beer, Eggs
3	Milk, Diaper, Beer, Cola
4	Bread, Milk, Diaper, Beer
5	Bread, Milk, Diaper, Cola

1 Announcements and References

- Administrative
- References for Today's Lecture

2 Data

- Feature Selection
- Feature Creation
- Discretization and Binarization
- Attribute Transformations
- Similarity and Dissimilarity
- Correlation

3 Association Rules

- Association Rule Mining
- **Definitions**
- Problem Definition
- Frequent Itemset Generation
- The Apriori Algorithm
 - Efficient Candidate Generation
 - Complexity Factors

4 Summary

Definitions

- **Binary Representation:** presence of an item in a transaction is often considered more important than its absence, therefore an item is represented as an asymmetric binary variable.

T-ID	Items
1	Bread, Milk
2	Bread, Diaper, Beer, Eggs
3	Milk, Diaper, Beer, Cola
4	Bread, Milk, Diaper, Beer
5	Bread, Milk, Diaper, Cola

Definitions

- **Binary Representation:** presence of an item in a transaction is often considered more important than its absence, therefore an item is represented as an asymmetric binary variable.

T-ID	Bread	Milk	Diapers	Beers	Eggs	Cola
1	1	1	0	0	0	0
2	1	0	1	1	1	0
3	0	1	1	1	0	1
4	1	1	1	1	0	0
5	1	1	1	0	0	1

Definitions

- Set of all transactions: $T = \{t_1, t_2, \dots, t_N\}$.
- Set of all Items: $I = \{i_1, i_2, \dots, i_d\}$.
- Each transaction t_i contains a subset of items chosen from I .
- Itemset: a collection of zero or more items.
- k -Itemset: itemset containing k items.

Example: $\{Beer, Diapers, Milk\}$ is a 3-itemset.

T-ID	Bread	Milk	Diapers	Beers	Eggs	Cola
1	1	1	0	0	0	0
2	1	0	1	1	1	0
3	0	1	1	1	0	1
4	1	1	1	1	0	0
5	1	1	1	0	0	1

Definitions

- **Support Count ($\sigma(X)$):** number of transactions that contain a particular itemset.

$$\sigma(X) = |\{t_i | X \subseteq t_i, t_i \in T\}| \quad (14)$$

- Example: $\sigma(\{\text{Beer}, \text{Diapers}, \text{Milk}\}) = 2$.

T-ID	Bread	Milk	Diapers	Beers	Eggs	Cola
1	1	1	0	0	0	0
2	1	0	1	1	1	0
3	0	1	1	1	0	1
4	1	1	1	1	0	0
5	1	1	1	0	0	1

Definitions

- **Association Rule:** is an implication expression of the form $X \rightarrow Y$, where, $X \cap Y = \emptyset$.
- **Support** $s(X \rightarrow Y) = \frac{\sigma(X \cup Y)}{N}$.
- Example: $s(\{Milk, Diaper\} \rightarrow \{Beer\}) = \frac{\sigma(\{Milk, Diaper, Beer\})}{|T|} = 2/5 = 0.4$
- **Confidence** $c(X \rightarrow Y) = \frac{\sigma(X \cup Y)}{\sigma(X)}$.
- Example: $c(\{Milk, Diaper\} \rightarrow \{Beer\}) = \frac{\sigma(\{Milk, Diaper, Beer\})}{\sigma(\{Milk, Diaper\})} = 2/3 = 0.6$

T-ID	Bread	Milk	Diapers	Beers	Eggs	Cola
1	1	1	0	0	0	0
2	1	0	1	1	1	0
3	0	1	1	1	0	1
4	1	1	1	1	0	0
5	1	1	1	0	0	1

1 Announcements and References

- Administrative
- References for Today's Lecture

2 Data

- Feature Selection
- Feature Creation
- Discretization and Binarization
- Attribute Transformations
- Similarity and Dissimilarity
- Correlation

3 Association Rules

- Association Rule Mining
- Definitions
- **Problem Definition**
- Frequent Itemset Generation
- The Apriori Algorithm
 - Efficient Candidate Generation
 - Complexity Factors

4 Summary

Association Rule Mining Problem

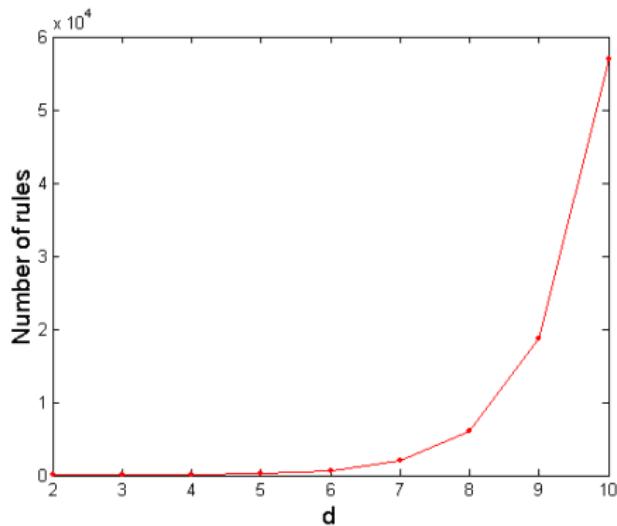
- Given a set of transactions T , the goal of association rule mining is to find all rules having:
 - $\text{support} \geq \text{minimum support}$ threshold,
 - $\text{confidence} \geq \text{minimum confidence}$ threshold.
- Minimum support (minsup) and confidence (minconf) represent thresholds on interesting association rules.
- Brute-Force Approach:
 - List all possible association rules.
 - Compute the support and confidence for each rule.
 - Prune rules that do not satisfy the minsup and minconf thresholds.
- Brute-force is computationally prohibitive!

Brute-Force Association Rule Mining

- Consider, we have d unique items.
- Total number of itemsets = 2^d .
- Total number of possible association rules

$$R = \sum_{k=1}^d \left[{}^d C_k \cdot \sum_{j=1}^{d-k} {}^{d-k} C_j \right], \\ = 3^d - 2^d + 1.$$

- Example: for $d = 6$, then $R = 602$ rules to analyze.



Mining Association Rules

- Consider, all rules that can be generated from the itemset $\{\text{Milk}, \text{Diaper}, \text{Beer}\}$:

- 1 $\{\text{Milk}, \text{Diaper}\} \rightarrow \{\text{Beer}\}$ ($s = 0.4, c = 0.\bar{6}$).
- 2 $\{\text{Milk}, \text{Beer}\} \rightarrow \{\text{Diaper}\}$ ($s = 0.4, c = 1.0$).
- 3 $\{\text{Diaper}, \text{Beer}\} \rightarrow \{\text{Milk}\}$ ($s = 0.4, c = 0.\bar{6}$).
- 4 $\{\text{Beer}\} \rightarrow \{\text{Milk}, \text{Diaper}\}$ ($s = 0.4, c = 0.\bar{6}$).
- 5 $\{\text{Diaper}\} \rightarrow \{\text{Milk}, \text{Beer}\}$ ($s = 0.4, c = 0.5$).
- 6 $\{\text{Milk}\} \rightarrow \{\text{Diaper}, \text{Beer}\}$ ($s = 0.4, c = 0.5$).

T-ID	Bread	Milk	Diapers	Beers	Eggs	Cola
1	1	1	0	0	0	0
2	1	0	1	1	1	0
3	0	1	1	1	0	1
4	1	1	1	1	0	0
5	1	1	1	0	0	1

Mining Association Rules

- Consider, all rules that can be generated from the itemset $\{\text{Milk}, \text{Diaper}, \text{Beer}\}$:
 - 1 $\{\text{Milk}, \text{Diaper}\} \rightarrow \{\text{Beer}\}$ ($s = 0.4, c = 0.\bar{6}$).
 - 2 $\{\text{Milk}, \text{Beer}\} \rightarrow \{\text{Diaper}\}$ ($s = 0.4, c = 1.0$).
 - 3 $\{\text{Diaper}, \text{Beer}\} \rightarrow \{\text{Milk}\}$ ($s = 0.4, c = 0.\bar{6}$).
 - 4 $\{\text{Beer}\} \rightarrow \{\text{Milk}, \text{Diaper}\}$ ($s = 0.4, c = 0.\bar{6}$).
 - 5 $\{\text{Diaper}\} \rightarrow \{\text{Milk}, \text{Beer}\}$ ($s = 0.4, c = 0.5$).
 - 6 $\{\text{Milk}\} \rightarrow \{\text{Diaper}, \text{Beer}\}$ ($s = 0.4, c = 0.5$).
- Rules originating from the same itemset have identical support but can have different confidence.
- Thus, we may decouple the support and confidence requirements.

Mining Association Rules

- Two-step approach:
 - 1 Frequent Itemset Generation: Generate all itemsets whose support \geq minimum support.
 - 2 Rule Generation: Generate high confidence rules from each frequent itemset, where each rule is a binary partitioning of a frequent itemset.
- Frequent itemset generation is still a computationally expensive process.

Frequent Itemset Generation — Brute-Force Approach

- Given d items, there are 2^d possible candidate itemsets.

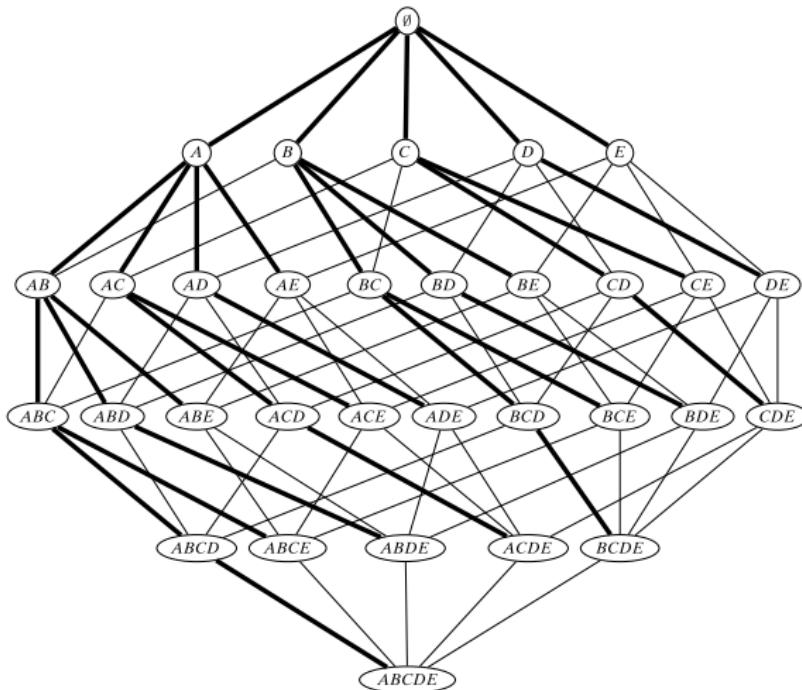


Figure 8.2. Itemset lattice and prefix-based search tree (in bold).

1 Announcements and References

- Administrative
- References for Today's Lecture

2 Data

- Feature Selection
- Feature Creation
- Discretization and Binarization
- Attribute Transformations
- Similarity and Dissimilarity
- Correlation

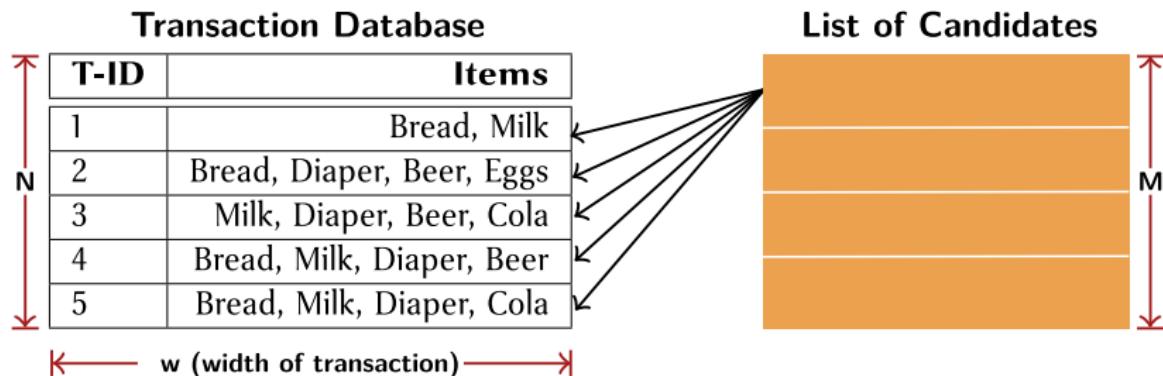
3 Association Rules

- Association Rule Mining
- Definitions
- Problem Definition
- Frequent Itemset Generation**
- The Apriori Algorithm
 - Efficient Candidate Generation
 - Complexity Factors

4 Summary

Brute-Force Frequent Itemset Generation

- Each itemset in the lattice is a candidate frequent itemset.
- Count the support of each candidate by scanning the database.
- Match each transaction against every candidate.
- Complexity — $\mathcal{O}(NMw)$ — very expensive since $M = 2^d$.



Improving Brute-Force Frequent Itemset Generation

- Reduce the number of candidates (M)
 - Complete search: $M = 2^d$.
 - Use pruning techniques to reduce M.
- Reduce the number of comparisons (NM)
 - Use efficient data structures to store the candidates or transactions.
 - No need to match every candidate against every transaction.

1 Announcements and References

- Administrative
- References for Today's Lecture

2 Data

- Feature Selection
- Feature Creation
- Discretization and Binarization
- Attribute Transformations
- Similarity and Dissimilarity
- Correlation

3 Association Rules

- Association Rule Mining
- Definitions
- Problem Definition
- Frequent Itemset Generation
- **The Apriori Algorithm**
 - Efficient Candidate Generation
 - Complexity Factors

4 Summary

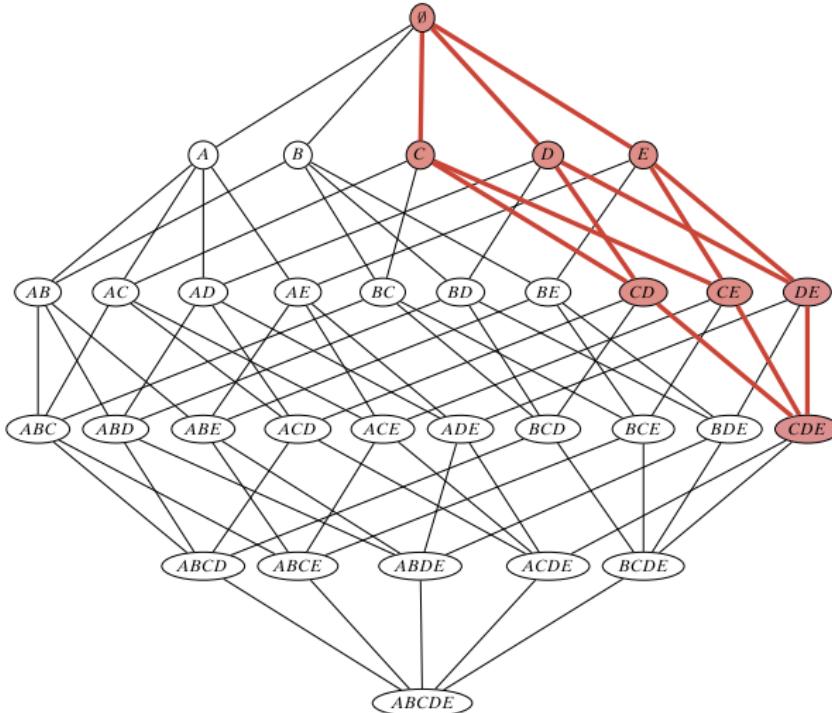
Reducing Number of Candidates — The Apriori Principle

- **Apriori Principle:** If an itemset is frequent, then all of its subsets must also be frequent.

$$\forall X, Y : (X \subseteq Y) \implies s(X) \geq s(Y) \quad (15)$$

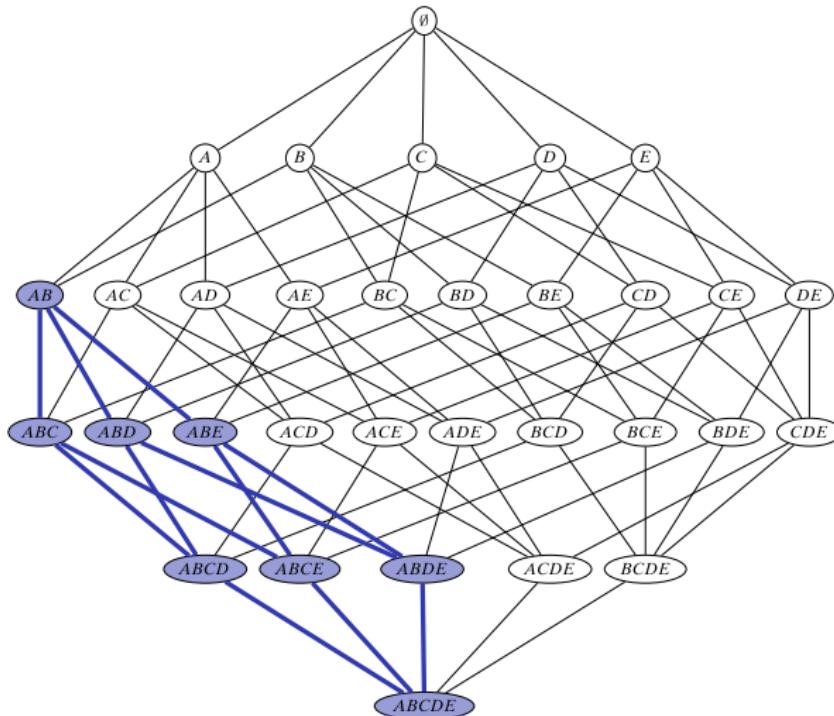
- Support of an itemset never exceeds the support of its subsets.
- This is known as the anti-monotone property of support.
- Conversely, if an itemset is infrequent then all of its supersets must be infrequent too.

Reducing Number of Candidates — The Apriori Principle



An illustration of the Apriori principle. If $\{c, d, e\}$ is frequent, then all subsets of this itemset are frequent.

Reducing Number of Candidates — The Apriori Principle



An illustration of support-based pruning. If $\{a, b\}$ is infrequent, then all supersets of $\{a, b\}$ are infrequent.

Apriori Algorithm — Example

Example 8.6. Consider the example dataset in Figure 8.1; let $\text{minsup} = 3$. Figure 8.3 shows the itemset search space for the Apriori method, organized as a prefix tree where two itemsets are connected if one is a prefix and immediate subset of the other.

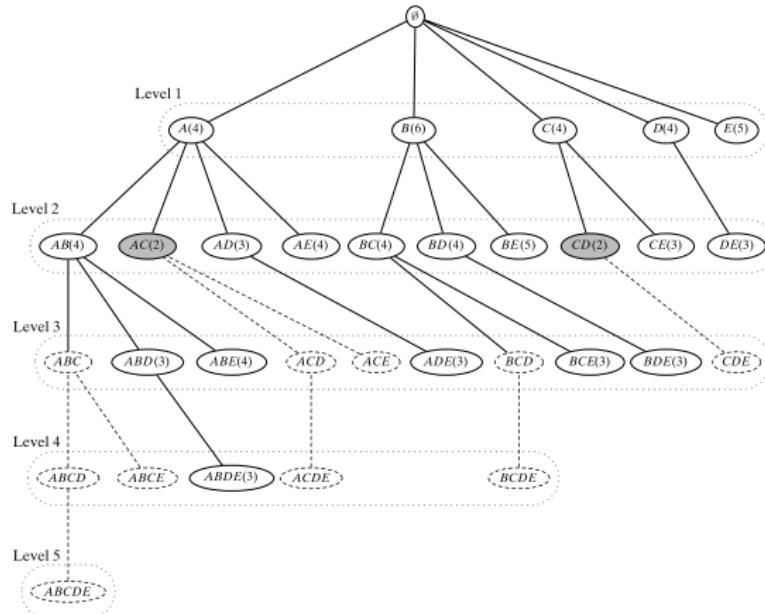


Figure 8.3. Apriori: prefix search tree and effect of pruning. Shaded nodes indicate infrequent itemsets; dashed nodes and lines indicate all of the pruned nodes and branches; solid lines indicate frequent itemsets.

Apriori Algorithm

- F_k : frequent k-itemsets.
- L_k : candidate k-itemsets.
- Algorithm
 - 1 Let $k=1$
 - 2 Generate $F_1 = \{\text{frequent 1-itemsets}\}$
 - 3 Repeat until F_k is empty
 - 1 **Candidate Generation:** Generate L_{k+1} from F_k .
 - 2 **Pruning:** Prune candidate itemsets in L_{k+1} containing subsets of length k that are infrequent.
 - 3 **Support Counting:** Count the support of each candidate in L_{k+1} by scanning the DB.
 - 4 **Candidate Elimination:** Eliminate candidates in L_{k+1} that are infrequent, leaving only those that are frequent $\implies F_{k+1}$.

Illustrating Apriori Principle

TID	Items
1	Bread, Milk
2	Beer, Bread, Diaper, Eggs
3	Beer, Coke, Diaper, Milk
4	Beer, Bread, Diaper, Milk
5	Bread, Coke, Diaper, Milk



Items (1-itemsets)

Item	Count
Bread	4
Coke	2
Milk	4
Beer	3
Diaper	4
Eggs	1

Minimum Support = 3

Illustrating Apriori Principle

Item	Count
Bread	4
Coke	2
Milk	4
Beer	3
Diaper	4
Eggs	1

Items (1-itemsets)



Itemset	Count
{Bread, Milk}	3
{Bread, Beer}	2
{Bread, Diaper}	3
{Milk, Beer}	2
{Milk, Diaper}	3
{Beer, Diaper}	3

Pairs (2-itemsets)

(No need to generate candidates involving Coke or Eggs)

Minimum Support = 3

If every subset is considered,

$${}^6C_1 + {}^6C_2 + {}^6C_3$$

$$6 + 15 + 20 = 41$$

With support-based pruning,

$$6 + 6 + 1 = 13$$



Triplets (3-itemsets)

Itemset	Count
{ Beer, Diaper, Milk }	2
{ Beer, Bread, Diaper }	2
{ Bread, Diaper, Milk }	2
{ Beer, Bread, Milk }	1

1 Announcements and References

- Administrative
- References for Today's Lecture

2 Data

- Feature Selection
- Feature Creation
- Discretization and Binarization
- Attribute Transformations
- Similarity and Dissimilarity
- Correlation

3 Association Rules

- Association Rule Mining
- Definitions
- Problem Definition
- Frequent Itemset Generation
- **The Apriori Algorithm**
 - Efficient Candidate Generation
 - Complexity Factors

4 Summary

Candidate Generation Step — Brute-Force Method

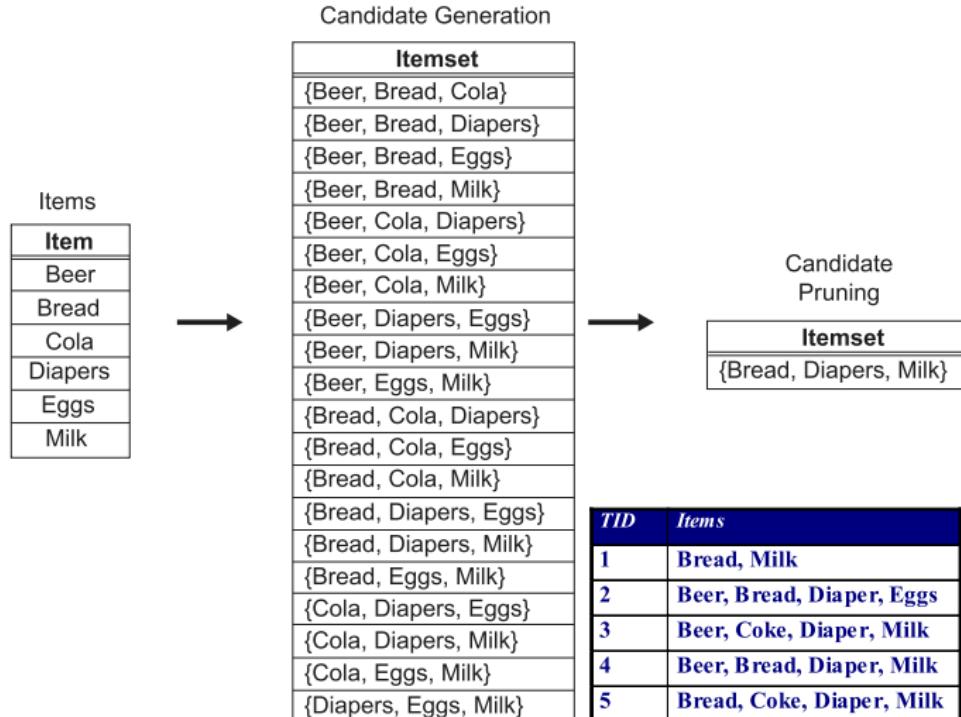


Figure 5.6. A brute-force method for generating candidate 3-itemsets.

Candidate Generation Step — Merge $F_{k-1} \times F_1$ Itemsets

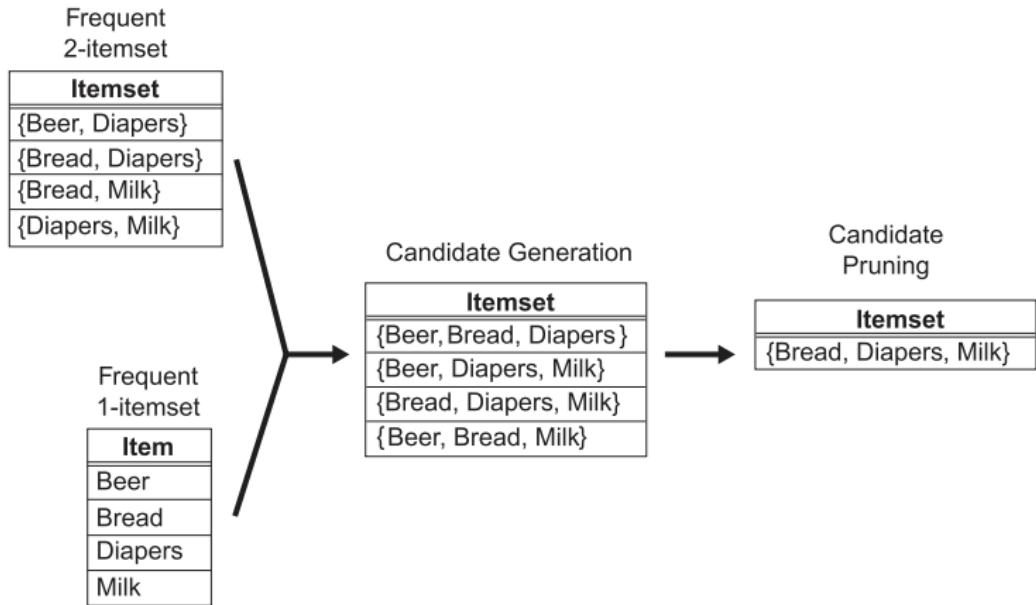


Figure 5.7. Generating and pruning candidate k -itemsets by merging a frequent $(k - 1)$ -itemset with a frequent item. Note that some of the candidates are unnecessary because their subsets are infrequent.

Candidate Generation Step — Merge $F_{k-1} \times F_{k-1}$ Itemsets

- Merge two frequent $(k - 1)$ -itemsets if their first $(k - 2)$ items are identical.
- $F_3 = \{ABC, ABD, ABE, ACD, BCD, BDE, CDE\}$.
 - $\text{Merge}(\underline{ABC}, \underline{ABD}) = \underline{ABCD}$.
 - $\text{Merge}(\underline{ABC}, \underline{ABE}) = \underline{ABCE}$.
 - $\text{Merge}(\underline{ABD}, \underline{ABE}) = \underline{ABDE}$.
- Do not merge $(\underline{ABD}, \underline{ACD})$ because they share only prefix of length 1 instead of length 2.

Candidate Generation Step — Merge $F_{k-1} \times F_{k-1}$ Itemsets

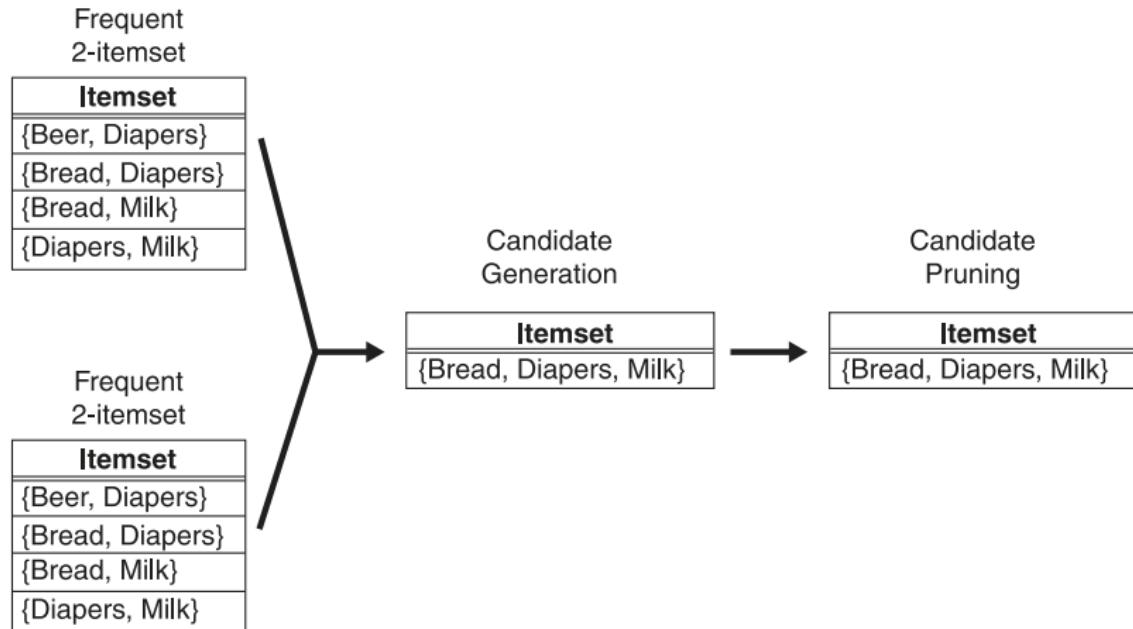


Figure 5.8. Generating and pruning candidate k -itemsets by merging pairs of frequent $(k-1)$ -itemsets.

1 Announcements and References

- Administrative
- References for Today's Lecture

2 Data

- Feature Selection
- Feature Creation
- Discretization and Binarization
- Attribute Transformations
- Similarity and Dissimilarity
- Correlation

3 Association Rules

- Association Rule Mining
- Definitions
- Problem Definition
- Frequent Itemset Generation
- **The Apriori Algorithm**
 - Efficient Candidate Generation
 - Complexity Factors

4 Summary

Factors Affecting Complexity of Apriori Algorithm

- Choice of minimum support threshold:
 - Lowering support threshold results in more frequent itemsets.
 - This may increase number of candidates and max length of frequent itemsets.
- Dimensionality (number of items) of the data set.
 - More space is needed to store support count of itemsets.
 - If number of frequent itemsets also increases, both computation and I/O costs may also increase.
- Size of database.
 - Run time of algorithm increases with number of transactions.
- Average transaction width.
 - Transaction width increases the max length of frequent itemsets.
 - Number of subsets in a transaction increases with its width, increasing computation time for support counting.

1 Announcements and References

- Administrative
- References for Today's Lecture

2 Data

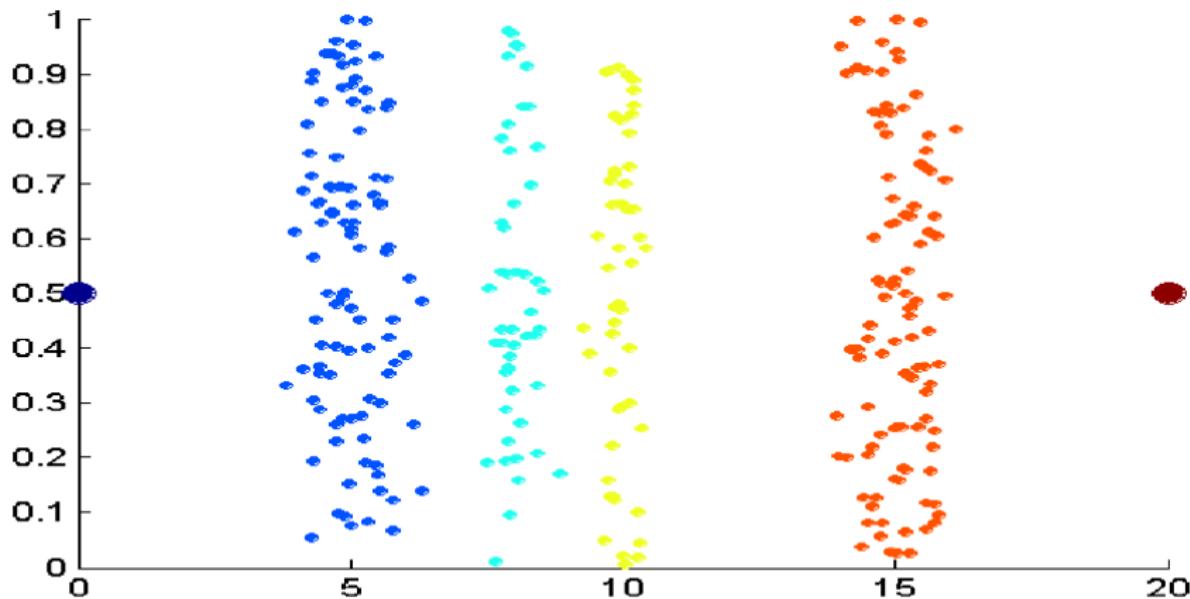
- Feature Selection
- Feature Creation
- Discretization and Binarization
- Attribute Transformations
- Similarity and Dissimilarity
- Correlation

3 Association Rules

- Association Rule Mining
- Definitions
- Problem Definition
- Frequent Itemset Generation
- The Apriori Algorithm
 - Efficient Candidate Generation
 - Complexity Factors

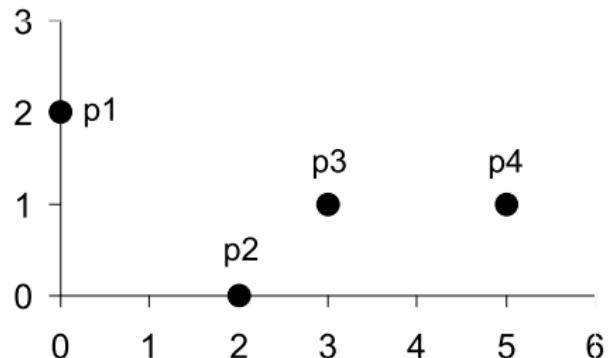
4 Summary

Summary — Discretization



Data consists of four groups of points and two outliers. Data is one-dimensional, but a random y component is added to reduce overlap.

Summary — Euclidean Distance



point	x	y
p1	0	2
p2	2	0
p3	3	1
p4	5	1

	p1	p2	p3	p4
p1	0	2.828	3.162	5.099
p2	2.828	0	1.414	3.162
p3	3.162	1.414	0	2
p4	5.099	3.162	2	0

Distance Matrix

Summary — Frequent Itemset Generation

- Given d items, there are 2^d possible candidate itemsets.

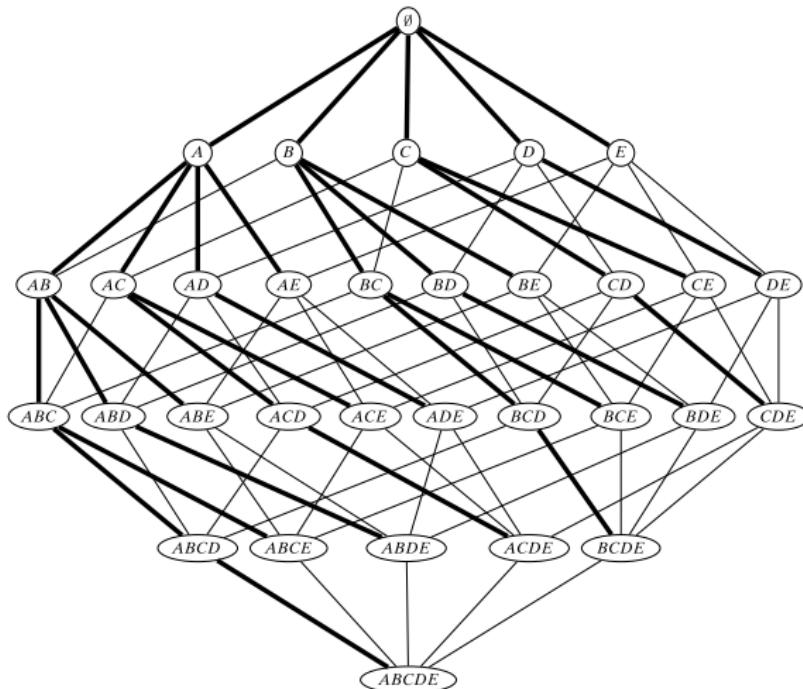


Figure 8.2. Itemset lattice and prefix-based search tree (in bold).

Summary — Administrative

1 First Assignment

- Due by 09.February.2023.

2 Volunteers for feedback regarding course

- Interested? Please contact me by email!