

Data Mining: Data

Instructor: Dr.Eng. Rahmat Widyanto

Dr.Eng. Rahmat Widyanto

What is Data?

- Collection of data objects and their attributes
- An attribute is a property or characteristic of an object
 - Examples: eye color of a person, temperature, etc.
 - Attribute is also known as variable, field, characteristic, or feature
- A collection of attributes describe an object
 - Object is also known as record, point, case, sample, entity, or instance

Attributes					
Tid	Refund	Marital Status	Taxable Income	Cheat	
1	Yes	Single	125K	No	
2	No	Married	100K	No	
3	No	Single	70K	No	
4	Yes	Married	120K	No	
5	No	Divorced	95K	Yes	
6	No	Married	60K	No	
7	Yes	Divorced	220K	No	
8	No	Single	85K	Yes	
9	No	Married	75K	No	
10	No	Single	90K	Yes	

Dr.Eng. Rahmat Widyanto

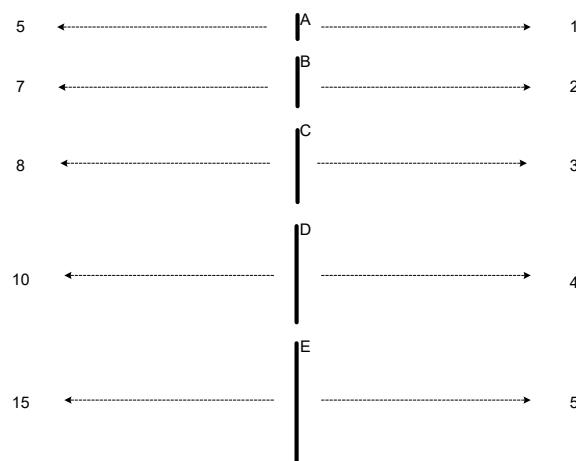
Attribute Values

- Attribute values are numbers or symbols assigned to an attribute
- Distinction between attributes and attribute values
 - Same attribute can be mapped to different attribute values
 - ◆ Example: height can be measured in feet or meters
 - Different attributes can be mapped to the same set of values
 - ◆ Example: Attribute values for ID and age are integers
 - ◆ But properties of attribute values can be different
 - ID has no limit but age has a maximum and minimum value

Dr.Eng. Rahmat Widyanto

Measurement of Length

- The way you measure an attribute is somewhat may not match the attributes properties.

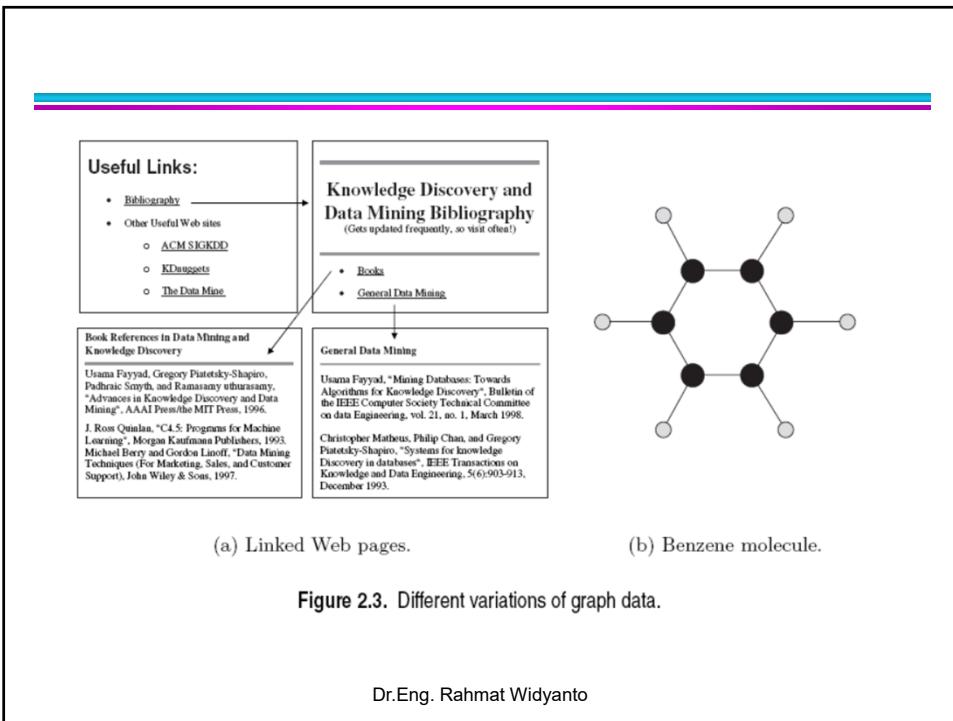


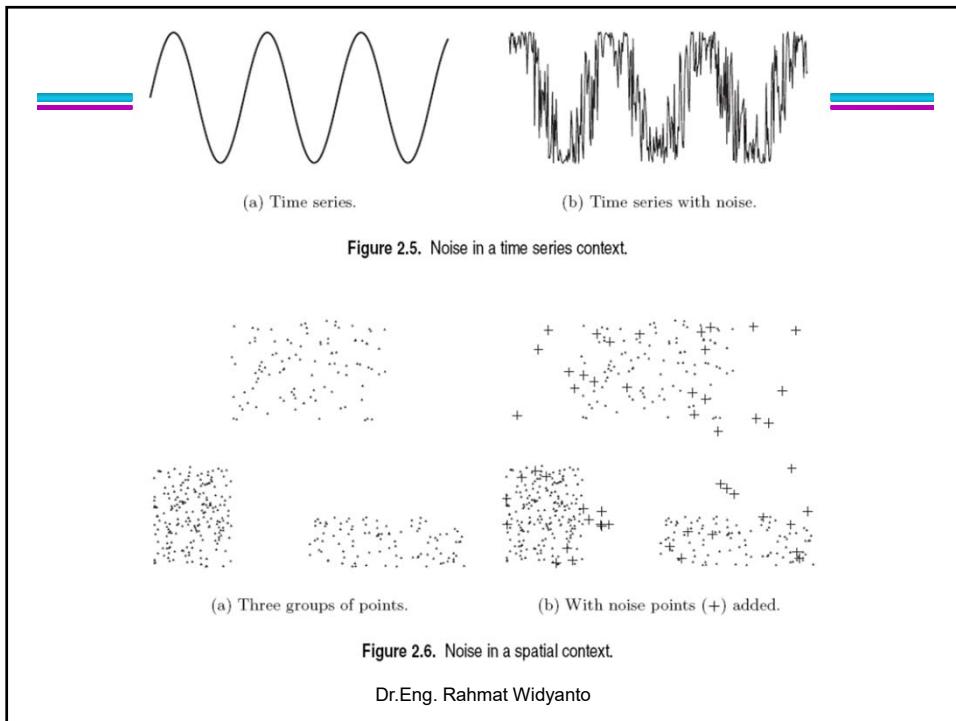
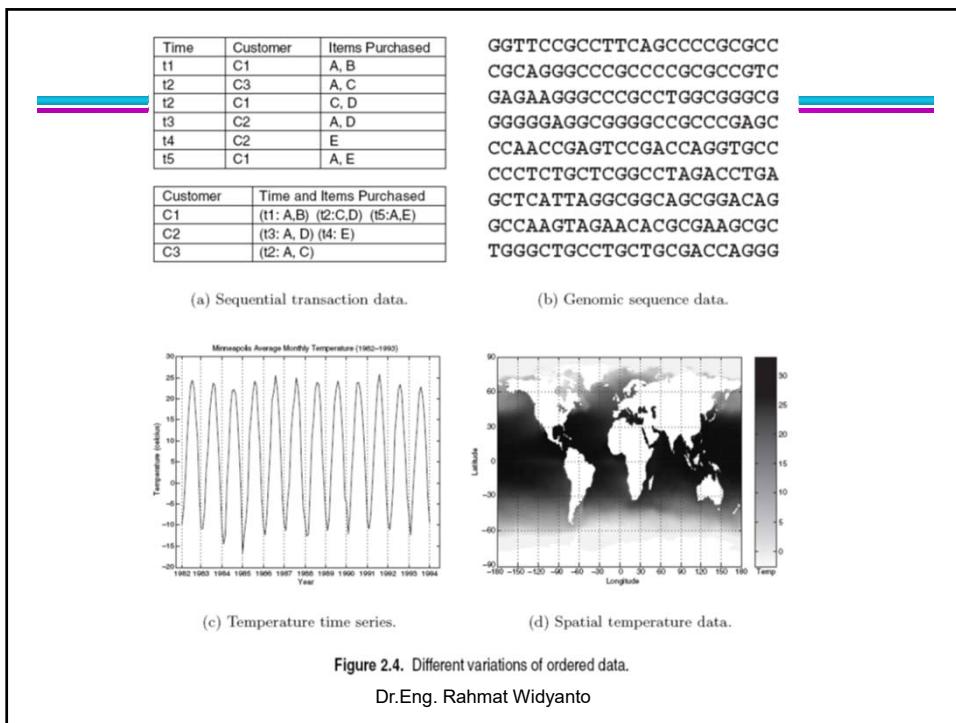
Dr.Eng. Rahmat Widyanto

<table border="1"> <thead> <tr> <th>Tid</th> <th>Refund</th> <th>Marital Status</th> <th>Taxable Income</th> <th>Defaulted Borrower</th> </tr> </thead> <tbody> <tr><td>1</td><td>Yes</td><td>Single</td><td>125K</td><td>No</td></tr> <tr><td>2</td><td>No</td><td>Married</td><td>100K</td><td>No</td></tr> <tr><td>3</td><td>No</td><td>Single</td><td>70K</td><td>No</td></tr> <tr><td>4</td><td>Yes</td><td>Married</td><td>120K</td><td>No</td></tr> <tr><td>5</td><td>No</td><td>Divorced</td><td>95K</td><td>Yes</td></tr> <tr><td>6</td><td>No</td><td>Married</td><td>60K</td><td>No</td></tr> <tr><td>7</td><td>Yes</td><td>Divorced</td><td>220K</td><td>No</td></tr> <tr><td>8</td><td>No</td><td>Single</td><td>85K</td><td>Yes</td></tr> <tr><td>9</td><td>No</td><td>Married</td><td>75K</td><td>No</td></tr> <tr><td>10</td><td>No</td><td>Single</td><td>90K</td><td>Yes</td></tr> </tbody> </table>	Tid	Refund	Marital Status	Taxable Income	Defaulted Borrower	1	Yes	Single	125K	No	2	No	Married	100K	No	3	No	Single	70K	No	4	Yes	Married	120K	No	5	No	Divorced	95K	Yes	6	No	Married	60K	No	7	Yes	Divorced	220K	No	8	No	Single	85K	Yes	9	No	Married	75K	No	10	No	Single	90K	Yes	<table border="1"> <thead> <tr> <th>TID</th> <th>ITEMS</th> </tr> </thead> <tbody> <tr><td>1</td><td>Bread, Soda, Milk</td></tr> <tr><td>2</td><td>Beer, Bread</td></tr> <tr><td>3</td><td>Beer, Soda, Diaper, Milk</td></tr> <tr><td>4</td><td>Beer, Bread, Diaper, Milk</td></tr> <tr><td>5</td><td>Soda, Diaper, Milk</td></tr> </tbody> </table>	TID	ITEMS	1	Bread, Soda, Milk	2	Beer, Bread	3	Beer, Soda, Diaper, Milk	4	Beer, Bread, Diaper, Milk	5	Soda, Diaper, Milk
Tid	Refund	Marital Status	Taxable Income	Defaulted Borrower																																																																
1	Yes	Single	125K	No																																																																
2	No	Married	100K	No																																																																
3	No	Single	70K	No																																																																
4	Yes	Married	120K	No																																																																
5	No	Divorced	95K	Yes																																																																
6	No	Married	60K	No																																																																
7	Yes	Divorced	220K	No																																																																
8	No	Single	85K	Yes																																																																
9	No	Married	75K	No																																																																
10	No	Single	90K	Yes																																																																
TID	ITEMS																																																																			
1	Bread, Soda, Milk																																																																			
2	Beer, Bread																																																																			
3	Beer, Soda, Diaper, Milk																																																																			
4	Beer, Bread, Diaper, Milk																																																																			
5	Soda, Diaper, Milk																																																																			
(a) Record data.	(b) Transaction data.																																																																			
<table border="1"> <thead> <tr> <th>Projection of x Load</th> <th>Projection of y Load</th> <th>Distance</th> <th>Load</th> <th>Thickness</th> </tr> </thead> <tbody> <tr><td>10.23</td><td>5.27</td><td>15.22</td><td>27</td><td>1.2</td></tr> <tr><td>12.65</td><td>6.25</td><td>16.22</td><td>22</td><td>1.1</td></tr> <tr><td>13.54</td><td>7.23</td><td>17.34</td><td>23</td><td>1.2</td></tr> <tr><td>14.27</td><td>8.43</td><td>18.45</td><td>25</td><td>0.9</td></tr> </tbody> </table>	Projection of x Load	Projection of y Load	Distance	Load	Thickness	10.23	5.27	15.22	27	1.2	12.65	6.25	16.22	22	1.1	13.54	7.23	17.34	23	1.2	14.27	8.43	18.45	25	0.9	<table border="1"> <thead> <tr> <th>team</th> <th>coach</th> <th>play</th> <th>goal</th> <th>solve</th> <th>win</th> <th>lost</th> <th>tournament</th> <th>session</th> </tr> </thead> <tbody> <tr><td>Document 1</td><td>3</td><td>0</td><td>5</td><td>0</td><td>2</td><td>0</td><td>2</td><td>0</td></tr> <tr><td>Document 2</td><td>0</td><td>7</td><td>0</td><td>2</td><td>1</td><td>0</td><td>0</td><td>3</td></tr> <tr><td>Document 3</td><td>0</td><td>1</td><td>0</td><td>0</td><td>1</td><td>2</td><td>2</td><td>0</td></tr> </tbody> </table>	team	coach	play	goal	solve	win	lost	tournament	session	Document 1	3	0	5	0	2	0	2	0	Document 2	0	7	0	2	1	0	0	3	Document 3	0	1	0	0	1	2	2	0						
Projection of x Load	Projection of y Load	Distance	Load	Thickness																																																																
10.23	5.27	15.22	27	1.2																																																																
12.65	6.25	16.22	22	1.1																																																																
13.54	7.23	17.34	23	1.2																																																																
14.27	8.43	18.45	25	0.9																																																																
team	coach	play	goal	solve	win	lost	tournament	session																																																												
Document 1	3	0	5	0	2	0	2	0																																																												
Document 2	0	7	0	2	1	0	0	3																																																												
Document 3	0	1	0	0	1	2	2	0																																																												
(c) Data matrix.	(d) Document-term matrix.																																																																			

Figure 2.2. Different variations of record data.

Dr.Eng. Rahmat Widyanto





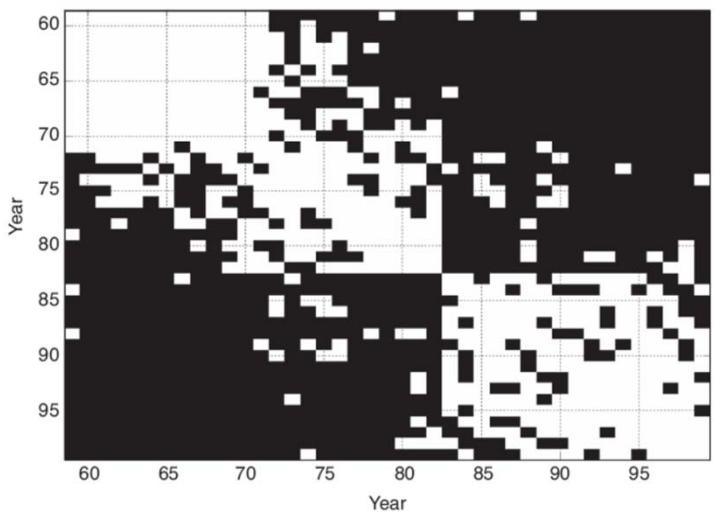
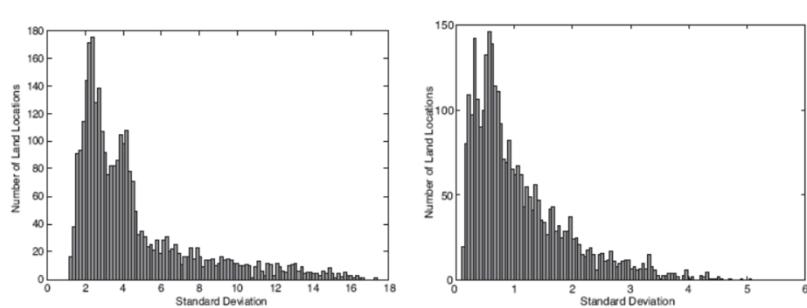


Figure 2.7. Correlation of SST data between pairs of years. White areas indicate positive correlation. Black areas indicate negative correlation.

Dr.Eng. Rahmat Widyanto

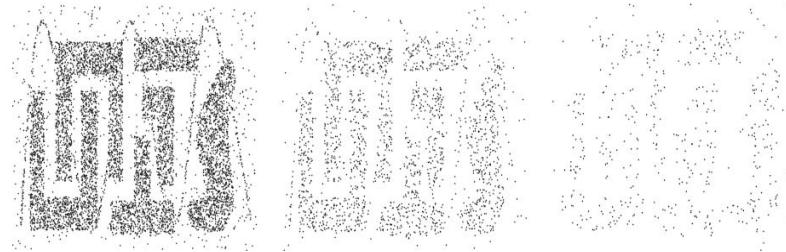


(a) Histogram of standard deviation of average monthly precipitation

(b) Histogram of standard deviation of average yearly precipitation

Figure 2.8. Histograms of standard deviation for monthly and yearly precipitation in Australia for the period 1982 to 1993.

Dr.Eng. Rahmat Widyanto



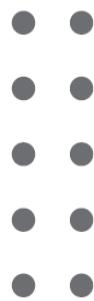
(a) 8000 points

(b) 2000 points

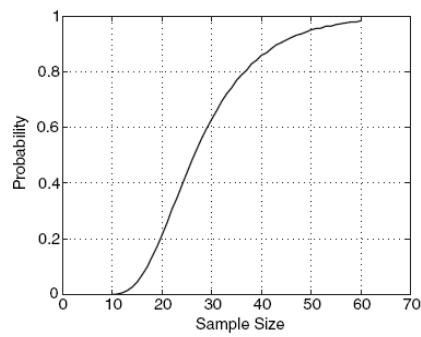
(c) 500 points

Figure 2.9. Example of the loss of structure with sampling.

Dr.Eng. Rahmat Widyanto



(a) Ten groups of points.



(b) Probability a sample contains points from each of 10 groups.

Figure 2.10. Finding representative points from 10 groups.

Dr.Eng. Rahmat Widyanto

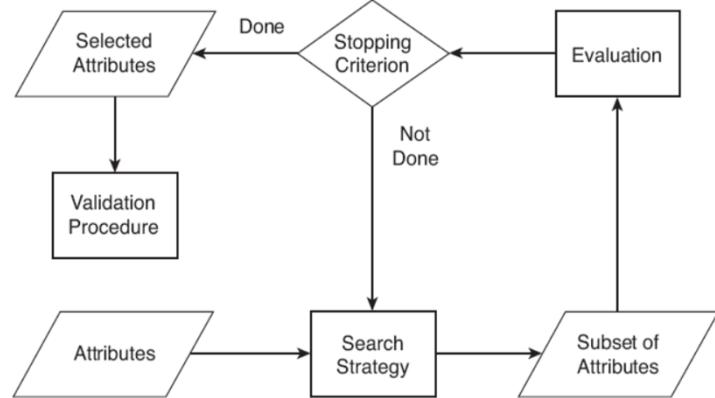
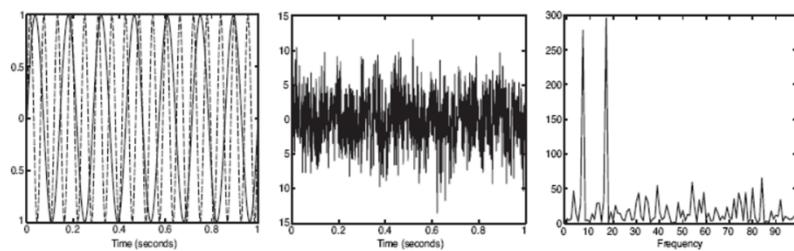


Figure 2.11. Flowchart of a feature subset selection process.

Dr.Eng. Rahmat Widyanto



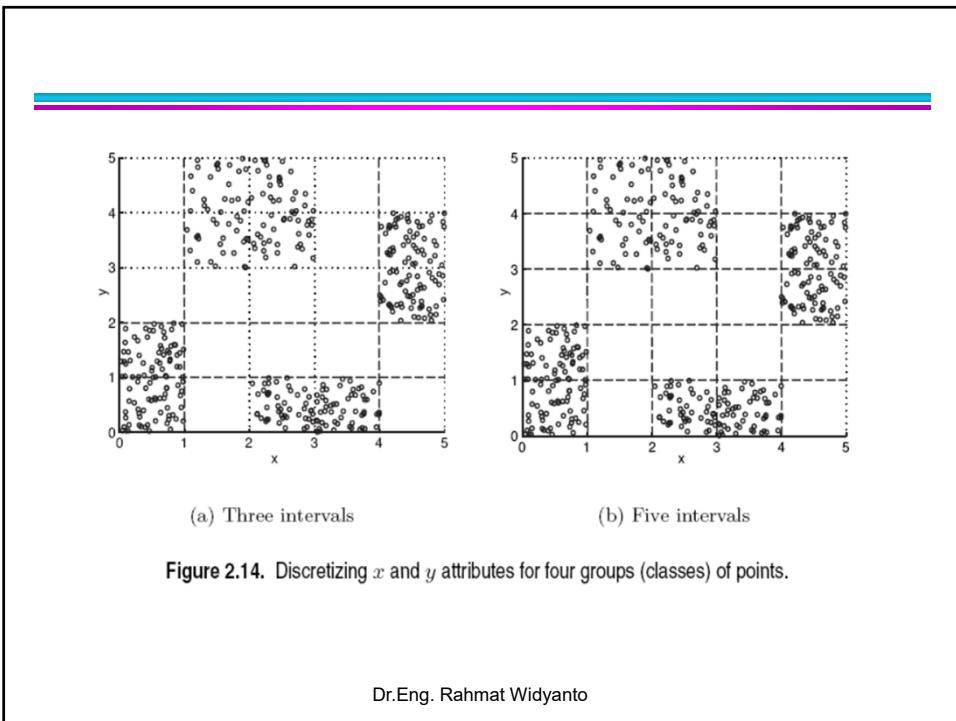
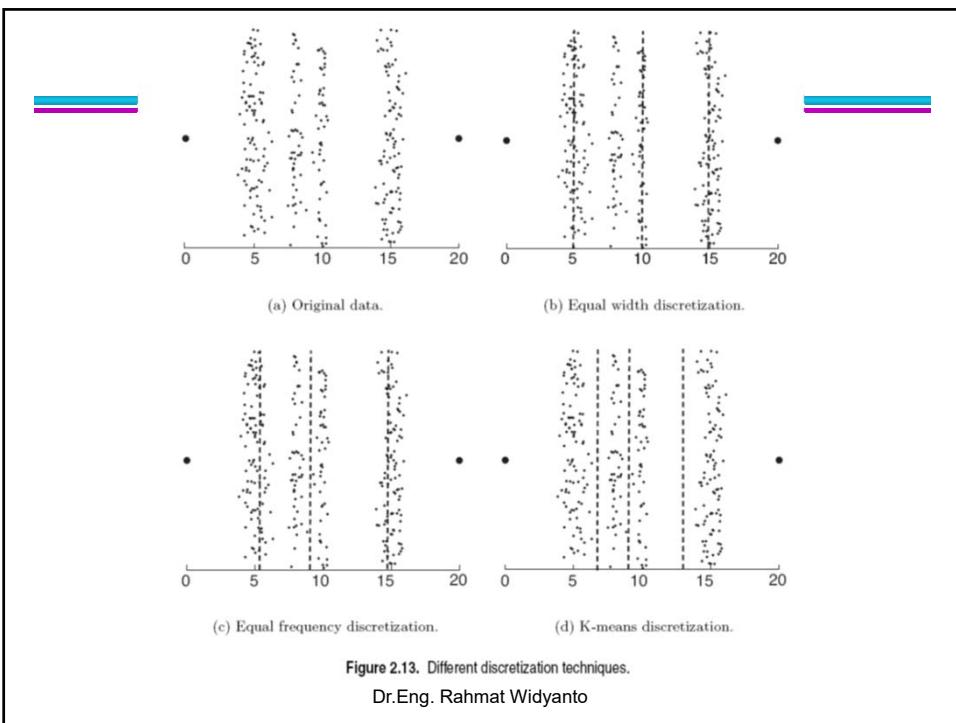
(a) Two time series.

(b) Noisy time series.

(c) Power spectrum

Figure 2.12. Application of the Fourier transform to identify the underlying frequencies in time series data.

Dr.Eng. Rahmat Widyanto



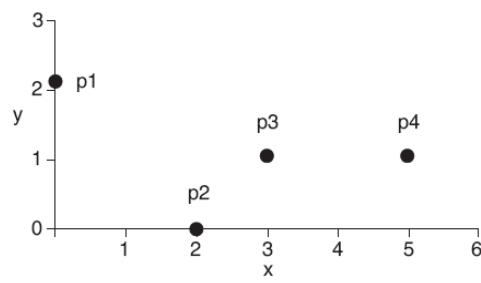


Figure 2.15. Four two-dimensional points.

Dr.Eng. Rahmat Widyanto

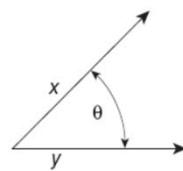


Figure 2.16. Geometric illustration of the cosine measure.

Dr.Eng. Rahmat Widyanto

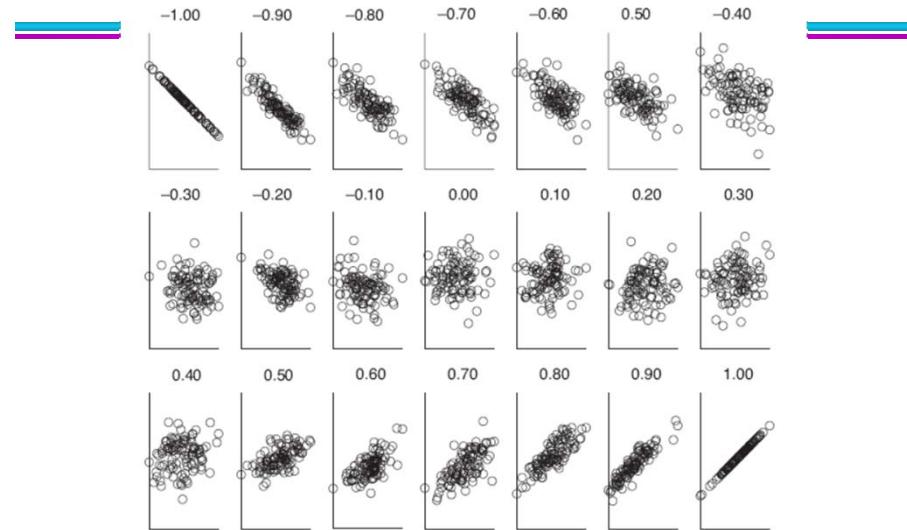


Figure 2.17. Scatter plots illustrating correlations from -1 to 1 .

Dr.Eng. Rahmat Widyanto

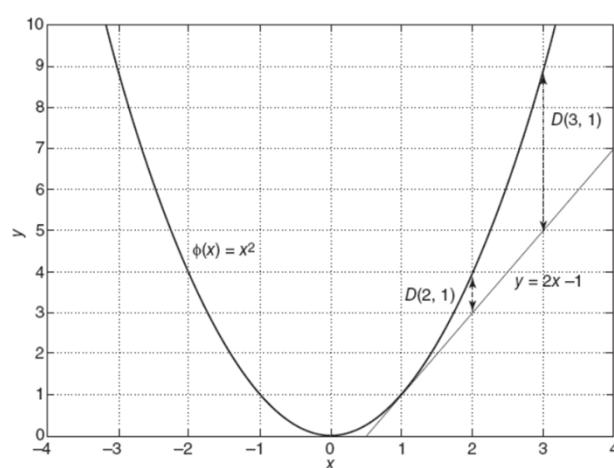


Figure 2.18. Illustration of Bregman divergence.

Dr.Eng. Rahmat Widyanto

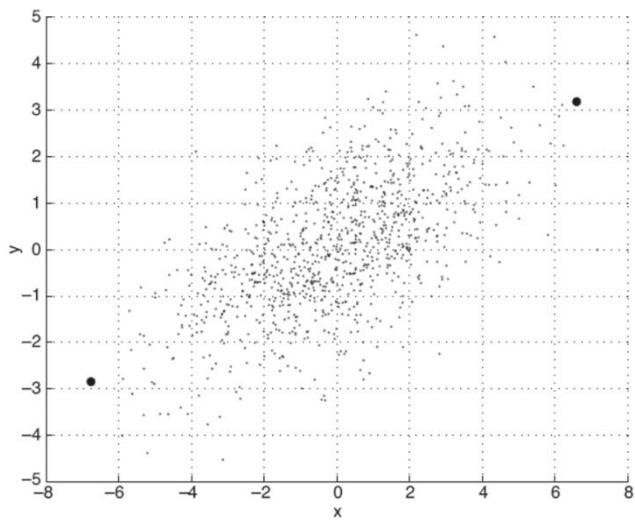
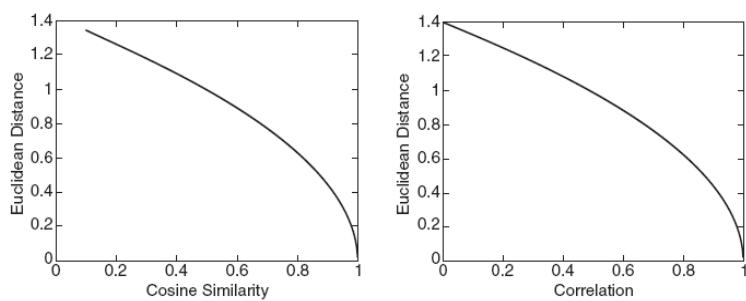


Figure 2.19. Set of two-dimensional points. The Mahalanobis distance between the two points represented by large dots is 6; their Euclidean distance is 14.7.

Dr.Eng. Rahmat Widyanto



(a) Relationship between Euclidean distance and the cosine measure.

(b) Relationship between Euclidean distance and correlation.

Figure 2.20. Graphs for Exercise 20.

Dr.Eng. Rahmat Widyanto

Types of Attributes

- There are different types of attributes
 - Nominal
 - ◆ Examples: ID numbers, eye color, zip codes
 - Ordinal
 - ◆ Examples: rankings (e.g., taste of potato chips on a scale from 1-10), grades, height in {tall, medium, short}
 - Interval
 - ◆ Examples: calendar dates, temperatures in Celsius or Fahrenheit.
 - Ratio
 - ◆ Examples: temperature in Kelvin, length, time, counts

Dr.Eng. Rahmat Widyanto

Properties of Attribute Values

- The type of an attribute depends on which of the following properties it possesses:
 - Distinctness: $= \neq$
 - Order: $< >$
 - Addition: $+ -$
 - Multiplication: $* /$
 - Nominal attribute: distinctness
 - Ordinal attribute: distinctness & order
 - Interval attribute: distinctness, order & addition
 - Ratio attribute: all 4 properties

Dr.Eng. Rahmat Widyanto

Attribute Type	Description	Examples	Operations
Nominal	The values of a nominal attribute are just different names, i.e., nominal attributes provide only enough information to distinguish one object from another. ($=, \neq$)	zip codes, employee ID numbers, eye color, sex: $\{male, female\}$	mode, entropy, contingency correlation, χ^2 test
Ordinal	The values of an ordinal attribute provide enough information to order objects. ($<, >$)	hardness of minerals, $\{good, better, best\}$, grades, street numbers	median, percentiles, rank correlation, run tests, sign tests
Interval	For interval attributes, the differences between values are meaningful, i.e., a unit of measurement exists. ($+, -$)	calendar dates, temperature in Celsius or Fahrenheit	mean, standard deviation, Pearson's correlation, t and F tests
Ratio	For ratio variables, both differences and ratios are meaningful. ($*, /$)	temperature in Kelvin, monetary quantities, counts, age, mass, length, electrical current	geometric mean, harmonic mean, percent variation

Attribute Level	Transformation	Comments
Nominal	Any permutation of values	If all employee ID numbers were reassigned, would it make any difference?
Ordinal	An order preserving change of values, i.e., $new_value = f(old_value)$ where f is a monotonic function.	An attribute encompassing the notion of good, better best can be represented equally well by the values $\{1, 2, 3\}$ or by $\{0.5, 1, 10\}$.
Interval	$new_value = a * old_value + b$ where a and b are constants	Thus, the Fahrenheit and Celsius temperature scales differ in terms of where their zero value is and the size of a unit (degree).
Ratio	$new_value = a * old_value$	Length can be measured in meters or feet.

Discrete and Continuous Attributes

- Discrete Attribute
 - Has only a finite or countably infinite set of values
 - Examples: zip codes, counts, or the set of words in a collection of documents
 - Often represented as integer variables.
 - Note: binary attributes are a special case of discrete attributes
- Continuous Attribute
 - Has real numbers as attribute values
 - Examples: temperature, height, or weight.
 - Practically, real values can only be measured and represented using a finite number of digits.
 - Continuous attributes are typically represented as floating-point variables.

Dr.Eng. Rahmat Widyanto

Types of data sets

- Record
 - Data Matrix
 - Document Data
 - Transaction Data
- Graph
 - World Wide Web
 - Molecular Structures
- Ordered
 - Spatial Data
 - Temporal Data
 - Sequential Data
 - Genetic Sequence Data

Dr.Eng. Rahmat Widyanto

Important Characteristics of Structured Data

- Dimensionality
 - ◆ Curse of Dimensionality
- Sparsity
 - ◆ Only presence counts
- Resolution
 - ◆ Patterns depend on the scale

Dr.Eng. Rahmat Widyanto

Record Data

- Data that consists of a collection of records, each of which consists of a fixed set of attributes

Tid	Refund	Marital Status	Taxable Income	Cheat
1	Yes	Single	125K	No
2	No	Married	100K	No
3	No	Single	70K	No
4	Yes	Married	120K	No
5	No	Divorced	95K	Yes
6	No	Married	60K	No
7	Yes	Divorced	220K	No
8	No	Single	85K	Yes
9	No	Married	75K	No
10	No	Single	90K	Yes

Dr.Eng. Rahmat Widyanto

Data Matrix

- If data objects have the same fixed set of numeric attributes, then the data objects can be thought of as points in a multi-dimensional space, where each dimension represents a distinct attribute
- Such data set can be represented by an m by n matrix, where there are m rows, one for each object, and n columns, one for each attribute

Projection of x Load	Projection of y load	Distance	Load	Thickness
10.23	5.27	15.22	2.7	1.2
12.65	6.25	16.22	2.2	1.1

Dr.Eng. Rahmat Widyanto

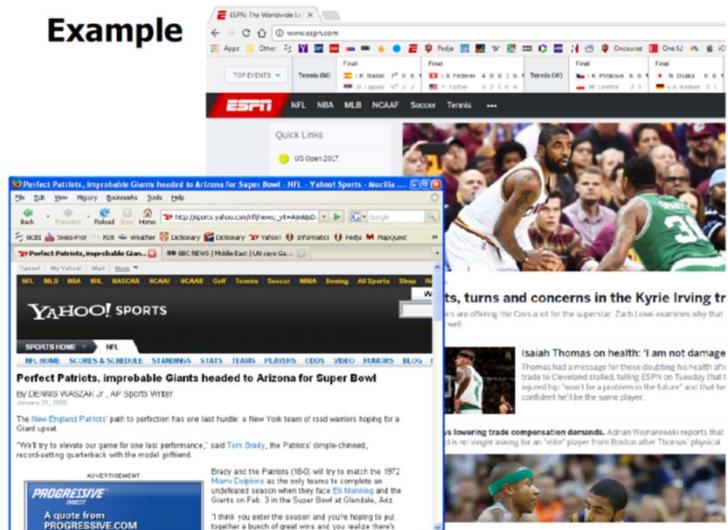
Document Data

- 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.

	team	coach	pla	ball	score	game	u	wi	lost	timeout	season
Document 1	3	0	5	0	2	6	0	2	0	0	2
Document 2	0	7	0	2	1	0	0	0	3	0	0
Document 3	0	1	0	0	1	2	2	0	0	3	0

Dr.Eng. Rahmat Widyanto

Example



Dr.Eng. Rahmat Widyanto

Recommendations Data

● Sparse matrix

- each row is a person
- each column is a movie (book, disease, ...)
- each number is a rating

	Movies							
	Spiderman	Ocean's 11	Matrix	Titanic	JFK	Star wars	Creed	Rocky
Persons								
Person 1		3		4				
Person 2							5	
Person 3							4	5
Person 4	1		3				2	

Dr.Eng. Rahmat Widyanto

Transaction Data

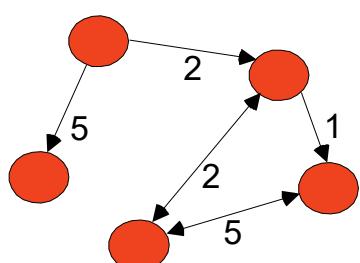
- A special type of record data, where
 - each record (transaction) involves a set of items.
 - For example, consider a grocery store. The set of products purchased by a customer during one shopping trip constitute a transaction, while the individual products that were purchased are the items.

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

Dr.Eng. Rahmat Widyanto

Graph Data

- Examples: Generic graph and HTML Links

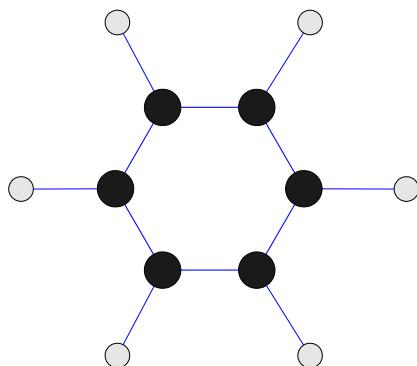


```
<a href="papers/papers.html#bbbb">  
Data Mining </a>  
<li>  
<a href="papers/papers.html#aaaa">  
Graph Partitioning </a>  
<li>  
<a href="papers/papers.html#aaaa">  
Parallel Solution of Sparse Linear System of Equations </a>  
<li>  
<a href="papers/papers.html#ffff">  
N-Body Computation and Dense Linear System Solvers
```

Dr.Eng. Rahmat Widyanto

Chemical Data

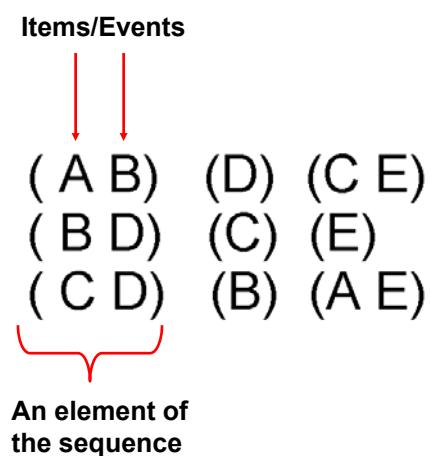
- Benzene Molecule: C₆H₆



Dr.Eng. Rahmat Widyanto

Ordered Data

- Sequences of transactions



Dr.Eng. Rahmat Widyanto

Ordered Data

- Genomic sequence data

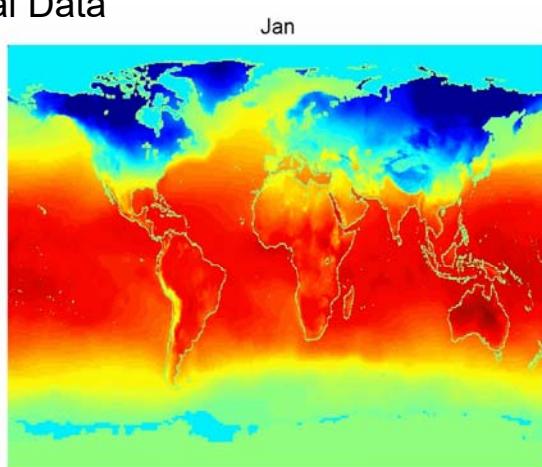
```
GGTTCCGCCTTCAGCCCCGCGCC  
CGCAGGGCCCGCCCCGCGCCGTC  
GAGAAGGGCCCGCCTGGCGGGCG  
GGGGGAGGCAGGGCCGCCGAGC  
CCAACCGAGTCCGACCAGGTGCC  
CCCTCTGCTCGGCCTAGACCTGA  
GCTCATTAGGCAGCGGACAG  
GCCAAGTAGAACACGCGAAGCGC  
TGGGCTGCCTGCTGCGACCAGGG
```

Dr.Eng. Rahmat Widyanto

Ordered Data

- Spatio-Temporal Data

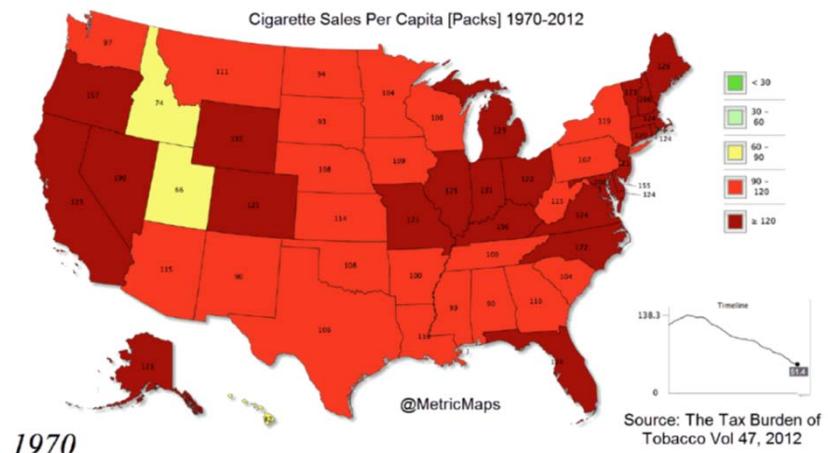
Average Monthly Temperature of land and ocean



Dr.Eng. Rahmat Widyanto

Ordered Data

- Spatio-Temporal Data

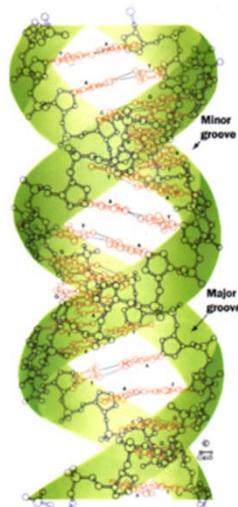


Dr.Eng. Rahmat Widyanto

Ordered/Sequence Data

- Genomic sequence data

1:	...GGTTCCGCCTTCAGCCCCCGGCC...	0
2:	...GGTTCCGCCTTCAGCCCCCGGCC...	1
3:	...GGTTCCGCCTTCAGCCCCCGGCC...	0
4:	...GGTTCCGCCTTCAGCCCCCGGCC...	0
5:	...GGTTCCGCCTTCAGCCCCCTGGCC...	0
6:	...GGTTCCGCCTTCAGCCCCCGGCC...	0
7:	...GGTTCCGCCTTCAGCCCCCTGGCC...	0
8:	...GGTTCCGCATTCAGCCCCCGGCC...	1
9:	...GGTTCCGCCTTCAGCCCCCGGCC...	0



Dr.Eng. Rahmat Widyanto

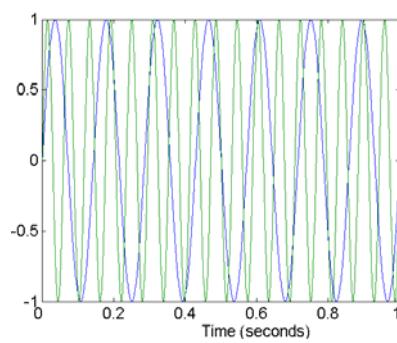
Data Quality

- What kinds of data quality problems?
 - How can we detect problems with the data?
 - What can we do about these problems?
-
- Examples of data quality problems:
 - Noise and outliers
 - missing values
 - duplicate data

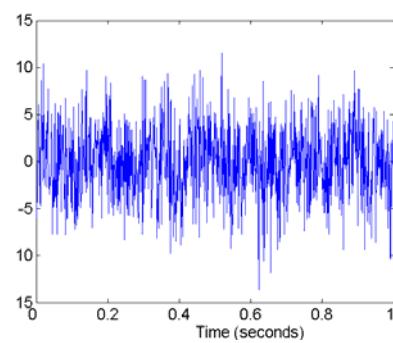
Dr.Eng. Rahmat Widyanto

Noise

- Noise refers to modification of original values
 - Examples: distortion of a person's voice when talking on a poor phone and "snow" on television screen



Two Sine Waves

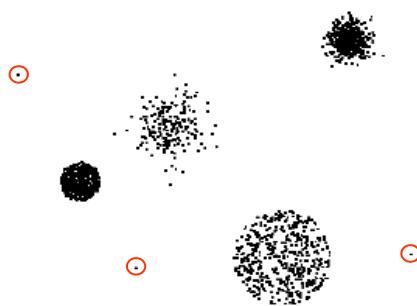


Two Sine Waves + Noise

Dr.Eng. Rahmat Widyanto

Outliers

- Outliers are data objects with characteristics that are considerably different than most of the other data objects in the data set



Dr.Eng. Rahmat Widyanto

Missing Values

- Reasons for missing values
 - Information is not collected
(e.g., people decline to give their age and weight)
 - Attributes may not be applicable to all cases
(e.g., annual income is not applicable to children)
- Handling missing values
 - Eliminate Data Objects
 - Estimate Missing Values
 - Ignore the Missing Value During Analysis
 - Replace with all possible values (weighted by their probabilities)

Dr.Eng. Rahmat Widyanto

Duplicate Data

- Data set may include data objects that are duplicates, or almost duplicates of one another
 - Major issue when merging data from heterogeneous sources
- Examples:
 - Same person with multiple email addresses
- Data cleaning
 - Process of dealing with duplicate data issues

Dr.Eng. Rahmat Widyanto

Data Preprocessing

- Aggregation
- Sampling
- Dimensionality Reduction
- Feature subset selection
- Feature creation
- Discretization and Binarization
- Attribute Transformation

Dr.Eng. Rahmat Widyanto

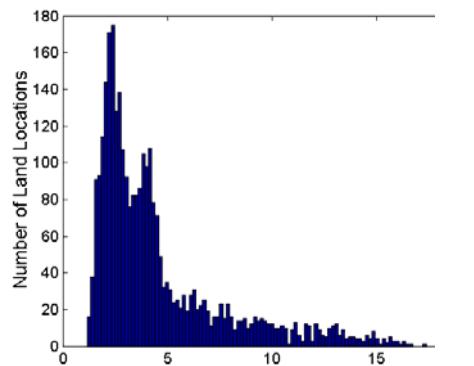
Aggregation

- Combining two or more attributes (or objects) into a single attribute (or object)
- Purpose
 - Data reduction
 - ◆ Reduce the number of attributes or objects
 - Change of scale
 - ◆ Cities aggregated into regions, states, countries, etc
 - More “stable” data
 - ◆ Aggregated data tends to have less variability

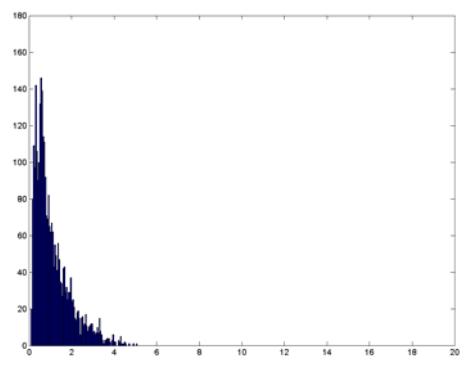
Dr.Eng. Rahmat Widyanto

Aggregation

Variation of Precipitation in Australia



Standard Deviation of Average
Monthly Precipitation



Standard Deviation of Average
Yearly Precipitation

Dr.Eng. Rahmat Widyanto

Sampling

- **Sampling is the main technique employed for data selection.**
 - It is often used for both the preliminary investigation of the data and the final data analysis.
- **Statisticians sample because obtaining the entire set of data of interest is too expensive or time consuming.**
- **Sampling is used in data mining because processing the entire set of data of interest is too expensive or time consuming.**

Dr.Eng. Rahmat Widyanto

Sampling ...

- The key principle for effective sampling is the following:
 - using a sample will work almost as well as using the entire data sets, if the sample is representative
 - A sample is representative if it has approximately the same property (of interest) as the original set of data

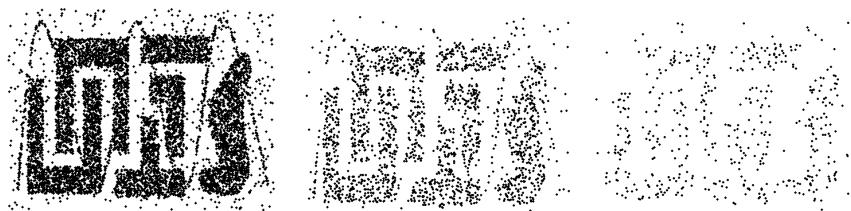
Dr.Eng. Rahmat Widyanto

Types of Sampling

- Simple Random Sampling
 - There is an equal probability of selecting any particular item
- Sampling without replacement
 - As each item is selected, it is removed from the population
- Sampling with replacement
 - Objects are not removed from the population as they are selected for the sample.
 - ◆ In sampling with replacement, the same object can be picked up more than once
- Stratified sampling
 - Split the data into several partitions; then draw random samples from each partition

Dr.Eng. Rahmat Widyanto

Sample Size



8000 points

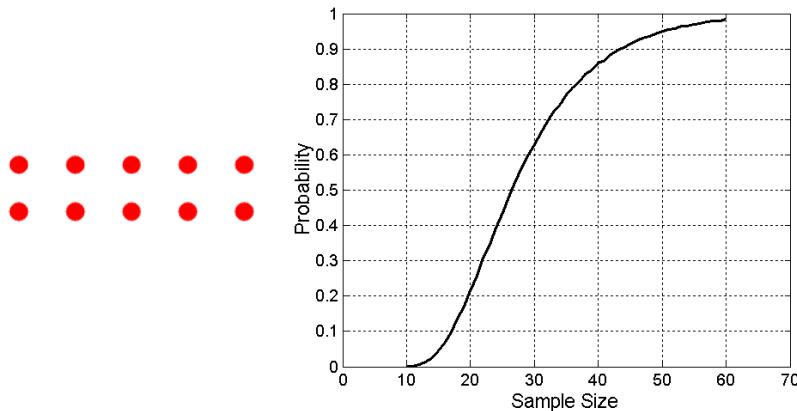
2000 Points

500 Points

Dr.Eng. Rahmat Widyanto

Sample Size

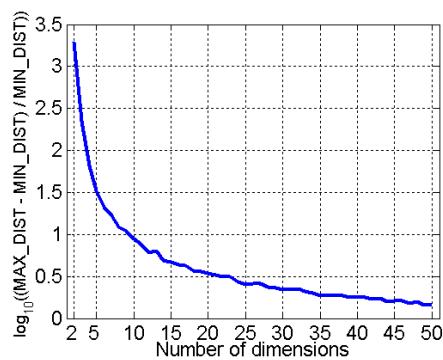
- What sample size is necessary to get at least one object from each of 10 groups.



Dr.Eng. Rahmat Widyanto

Curse of Dimensionality

- When dimensionality increases, data becomes increasingly sparse in the space that it occupies
- Definitions of density and distance between points, which is critical for clustering and outlier detection, become less meaningful



- Randomly generate 500 points
- Compute difference between max and min distance between any pair of points

Dr.Eng. Rahmat Widyanto

Dimensionality Reduction

- Purpose:

- Avoid curse of dimensionality
- Reduce amount of time and memory required by data mining algorithms
- Allow data to be more easily visualized
- May help to eliminate irrelevant features or reduce noise

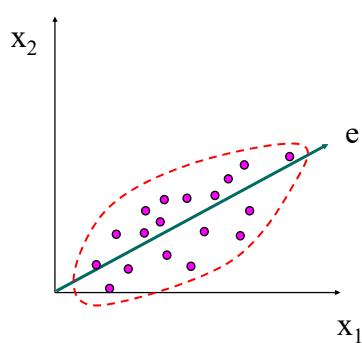
- Techniques

- Principle Component Analysis
- Singular Value Decomposition
- Others: supervised and non-linear techniques

Dr.Eng. Rahmat Widyanto

Dimensionality Reduction: PCA

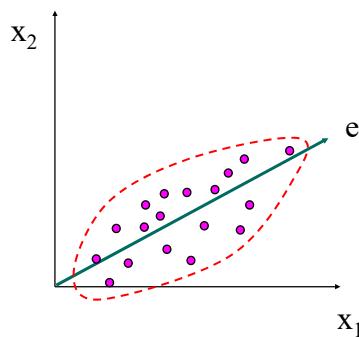
- Goal is to find a projection that captures the largest amount of variation in data



Dr.Eng. Rahmat Widyanto

Dimensionality Reduction: PCA

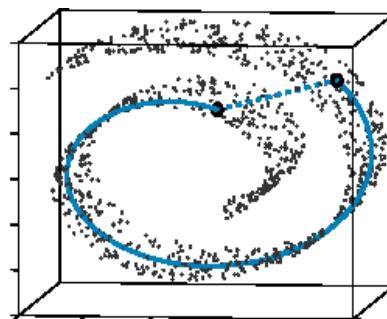
- Find the eigenvectors of the covariance matrix
- The eigenvectors define the new space



Dr.Eng. Rahmat Widyanto

Dimensionality Reduction: ISOMAP

By: Tenenbaum, de Silva,
Langford (2000)



- Construct a neighbourhood graph
- For each pair of points in the graph, compute the shortest path distances – geodesic distances

Dr.Eng. Rahmat Widyanto

Dimensionality Reduction: PCA

Dimensions = 206



Dr.Eng. Rahmat Widyanto

Feature Subset 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

Dr.Eng. Rahmat Widyanto

Feature Subset Selection

- Techniques:
 - Brute-force approach:
 - ◆ Try all possible feature subsets as input to data mining algorithm
 - Embedded approaches:
 - ◆ Feature selection occurs naturally as part of the data mining algorithm
 - Filter approaches:
 - ◆ Features are selected before data mining algorithm is run
 - Wrapper approaches:
 - ◆ Use the data mining algorithm as a black box to find best subset of attributes

Dr.Eng. Rahmat Widyanto

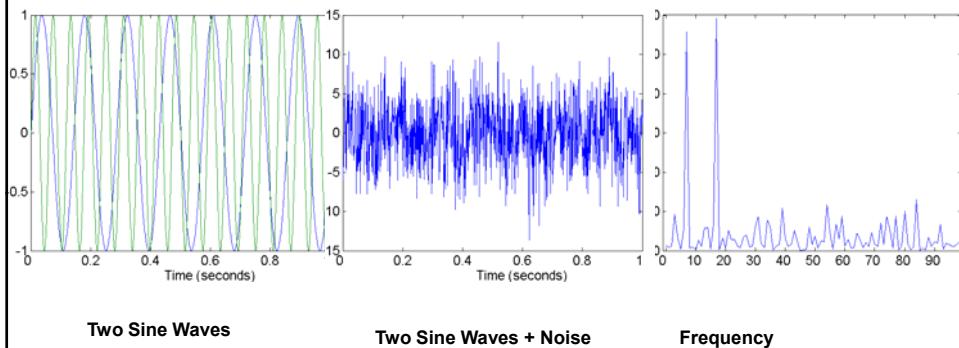
Feature Creation

- 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
 - ◆ domain-specific
 - Mapping Data to New Space
 - Feature Construction
 - ◆ combining features

Dr.Eng. Rahmat Widyanto

Mapping Data to a New Space

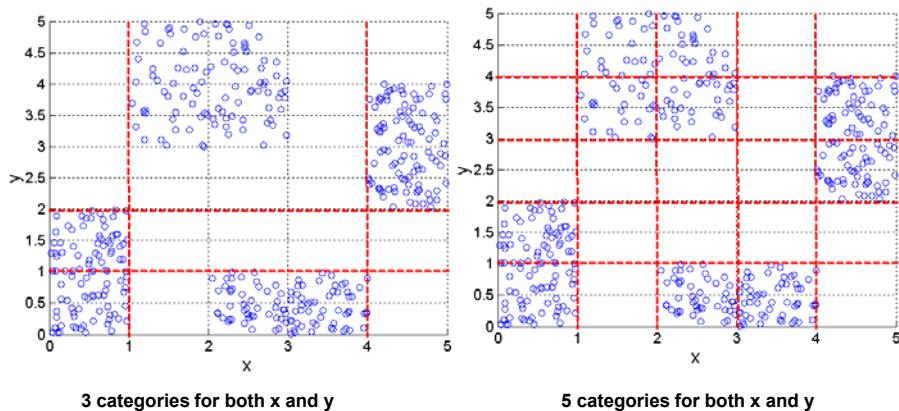
- Fourier transform
- Wavelet transform



Dr.Eng. Rahmat Widyanto

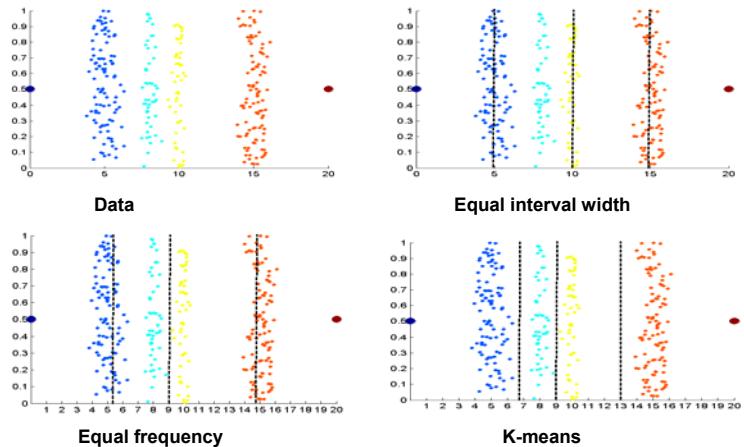
Discretization Using Class Labels

- Entropy based approach



Dr.Eng. Rahmat Widyanto

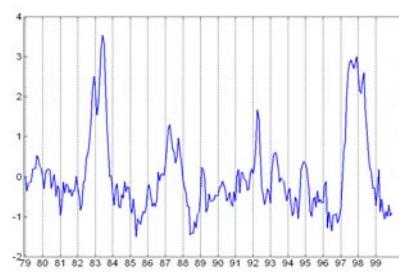
Discretization Without Using Class Labels



Dr.Eng. Rahmat Widyanto

Attribute Transformation

- 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
 - Simple functions: x^k , $\log(x)$, e^x , $|x|$
 - Standardization and Normalization



Dr.Eng. Rahmat Widyanto

Similarity and Dissimilarity

- Similarity

- Numerical measure of how alike two data objects are.
- Is higher when objects are more alike.
- Often falls in the range [0,1]

- Dissimilarity

- Numerical measure of how different are two data objects
- Lower when objects are more alike
- Minimum dissimilarity is often 0
- Upper limit varies

- Proximity refers to a similarity or dissimilarity

Dr.Eng. Rahmat Widyanto

Similarity/Dissimilarity for Simple Attributes

p and q are the attribute values for two data objects.

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

Table 5.1. Similarity and dissimilarity for simple attributes

Dr.Eng. Rahmat Widyanto

Euclidean Distance

- Euclidean Distance

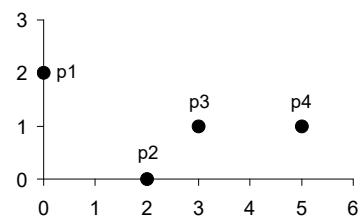
$$dist = \sqrt{\sum_{k=1}^n (p_k - q_k)^2}$$

Where n is the number of dimensions (attributes) and p_k and q_k are, respectively, the k^{th} attributes (components) or data objects p and q .

- Standardization is necessary, if scales differ.

Dr.Eng. Rahmat Widyanto

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

Dr.Eng. Rahmat Widyanto

Minkowski Distance

- Minkowski Distance is a generalization of Euclidean Distance

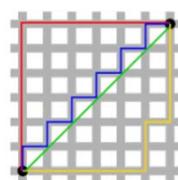
$$dist = \left(\sum_{k=1}^n |p_k - q_k|^r \right)^{\frac{1}{r}}$$

Where r is a parameter, n is the number of dimensions (attributes) and p_k and q_k are, respectively, the k th attributes (components) or data objects p and q .

Dr.Eng. Rahmat Widyanto

Minkowski Distance: Examples

- $r = 1$. City block (Manhattan, taxicab, L_1 norm) distance.
 - A common example of this 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) distance.
 - This is the maximum difference between any component of the vectors
- Do not confuse r with n , i.e., all these distances are defined for all numbers of dimensions.



Dr.Eng. Rahmat Widyanto

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

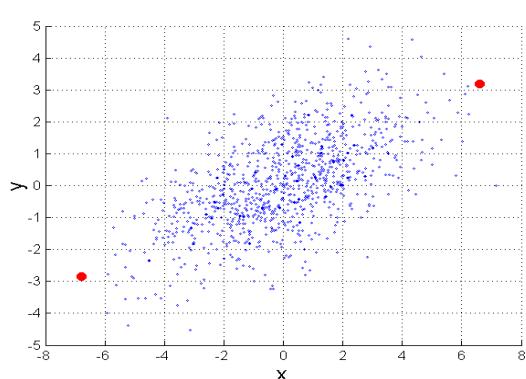
L ∞	p1	p2	p3	p4
p1	0	2	3	5
p2	2	0	1	3
p3	3	1	0	2
p4	5	3	2	0

Distance Matrix

Dr.Eng. Rahmat Widyanto

Mahalanobis Distance

$$mahalanobi s(p, q) = (p - q) \Sigma^{-1} (p - q)^T$$



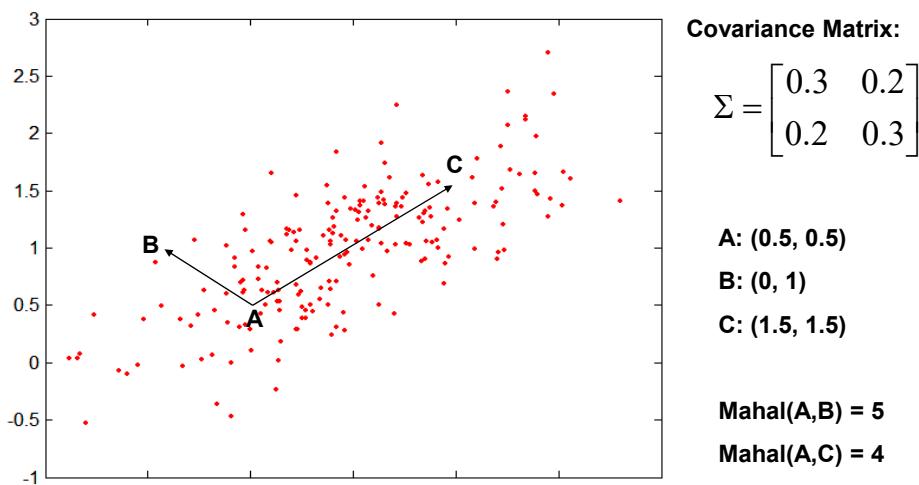
Σ is the covariance matrix of the input data X

$$\Sigma_{j,k} = \frac{1}{n-1} \sum_{i=1}^n (X_{ij} - \bar{X}_j)(X_{ik} - \bar{X}_k)$$

For red points, the Euclidean distance is 14.7, Mahalanobis distance is 6.

Dr.Eng. Rahmat Widyanto

Mahalanobis Distance



Dr.Eng. Rahmat Widyanto

Common Properties of a Distance

- Distances, such as the Euclidean distance, have some well known properties.
 - $d(p, q) \geq 0$ for all p and q and $d(p, q) = 0$ only if $p = q$. (Positive definiteness)
 - $d(p, q) = d(q, p)$ for all p and q . (Symmetry)
 - $d(p, r) \leq d(p, q) + d(q, r)$ for all points p , q , and r . (Triangle Inequality)where $d(p, q)$ is the distance (dissimilarity) between points (data objects), p and q .
- A distance that satisfies these properties is a **metric**

Dr.Eng. Rahmat Widyanto

Common Properties of a Similarity

- Similarities, also have some well known properties.

1. $s(p, q) = 1$ (or maximum similarity) only if $p = q$.
2. $s(p, q) = s(q, p)$ for all p and q . (Symmetry)

where $s(p, q)$ is the similarity between points (data objects), p and q .

Dr.Eng. Rahmat Widyanto

Similarity Between Binary Vectors

- Common situation is that objects, p and q , have only binary attributes
- Compute similarities using the following quantities
 - M_{01} = the number of attributes where p was 0 and q was 1
 - M_{10} = the number of attributes where p was 1 and q was 0
 - M_{00} = the number of attributes where p was 0 and q was 0
 - M_{11} = the number of attributes where p was 1 and q was 1
- Simple Matching and Jaccard Coefficients

$$\begin{aligned} \text{SMC} &= \text{number of matches} / \text{number of attributes} \\ &= (M_{11} + M_{00}) / (M_{01} + M_{10} + M_{11} + M_{00}) \end{aligned}$$

$$\begin{aligned} J &= \text{number of 11 matches} / \text{number of not-both-zero attributes values} \\ &= (M_{11}) / (M_{01} + M_{10} + M_{11}) \end{aligned}$$

Dr.Eng. Rahmat Widyanto

SMC versus Jaccard: Example

$p = 1 \ 0 \ 0 \ 0 \ 0 \ 0 \ 0 \ 0 \ 0 \ 0$

$q = 0 \ 0 \ 0 \ 0 \ 0 \ 1 \ 0 \ 0 \ 1$

$M_{01} = 2$ (the number of attributes where p was 0 and q was 1)

$M_{10} = 1$ (the number of attributes where p was 1 and q was 0)

$M_{00} = 7$ (the number of attributes where p was 0 and q was 0)

$M_{11} = 0$ (the number of attributes where p was 1 and q was 1)

$$\text{SMC} = (M_{11} + M_{00}) / (M_{01} + M_{10} + M_{11} + M_{00}) = (0+7) / (2+1+0+7) = 0.7$$

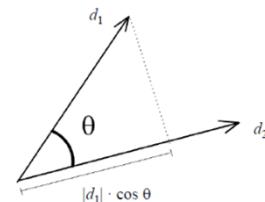
$$J = (M_{11}) / (M_{01} + M_{10} + M_{11}) = 0 / (2 + 1 + 0) = 0$$

Dr.Eng. Rahmat Widyanto

Cosine Similarity

- If d_1 and d_2 are two document vectors, then

$$\cos(d_1, d_2) = (d_1 \bullet d_2) / \|d_1\| \|d_2\|,$$



where \bullet indicates vector dot product and $\|d\|$ is the length of vector d .

- Example:

$d_1 = 3 \ 2 \ 0 \ 5 \ 0 \ 0 \ 0 \ 2 \ 0 \ 0$

$d_2 = 1 \ 0 \ 0 \ 0 \ 0 \ 0 \ 0 \ 1 \ 0 \ 2$

$$d_1 \bullet d_2 = 3*1 + 2*0 + 0*0 + 5*0 + 0*0 + 0*0 + 0*0 + 2*1 + 0*0 + 0*2 = 5$$

$$\|d_1\| = (3^2 + 2^2 + 0^2 + 5^2 + 0^2 + 0^2 + 0^2 + 2^2 + 0^2 + 0^2)^{0.5} = (42)^{0.5} = 6.481$$

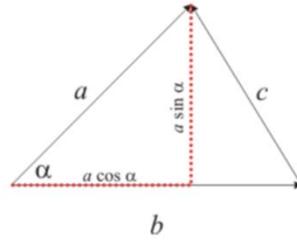
$$\|d_2\| = (1^2 + 0^2 + 0^2 + 0^2 + 0^2 + 0^2 + 0^2 + 1^2 + 0^2 + 2^2)^{0.5} = (6)^{0.5} = 2.245$$

$$\cos(d_1, d_2) = .3150$$

Dr.Eng. Rahmat Widyanto

Proof:

$$\begin{aligned}c^2 &= (b - a \cos \alpha)^2 + (a \sin \alpha)^2 \\&= b^2 - 2ab \cos \alpha + a^2 \cos^2 \alpha + a^2 \sin^2 \alpha \\&= a^2 + b^2 - 2ab \cos \alpha\end{aligned}$$



$$c^2 = \vec{c} \cdot \vec{c}$$

$$\begin{aligned}&= (\vec{a} - \vec{b}) \cdot (\vec{a} - \vec{b}) \\&= \vec{a} \cdot \vec{a} - 2\vec{a} \cdot \vec{b} + \vec{b} \cdot \vec{b} \\&= a^2 - 2\vec{a} \cdot \vec{b} + b^2\end{aligned}$$

$$\cos \alpha = \frac{\vec{a} \cdot \vec{b}}{ab}$$

By combining previous equations we get:

$$a^2 + b^2 - 2ab \cos \alpha = a^2 - 2\vec{a} \cdot \vec{b} + b^2$$

Dr.Eng. Rahmat Widyanto

Extended Jaccard Coefficient (Tanimoto)

- Variation of Jaccard for continuous or count attributes
 - Reduces to Jaccard for binary attributes

$$T(p, q) = \frac{p \bullet q}{\|p\|^2 + \|q\|^2 - p \bullet q}$$

Dr.Eng. Rahmat Widyanto

Correlation

- Correlation measures the linear relationship between objects
- To compute correlation, we standardize data objects, p and q , and then take their dot product

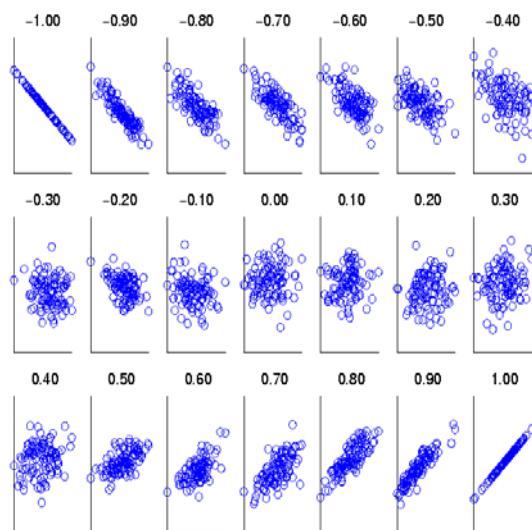
$$p'_k = (p_k - \text{mean}(p)) / \text{std}(p)$$

$$q'_k = (q_k - \text{mean}(q)) / \text{std}(q)$$

$$\text{correlation}(p, q) = p' \bullet q'$$

Dr.Eng. Rahmat Widyanto

Visually Evaluating Correlation



Scatter plots
showing the
similarity from
-1 to 1.

Dr.Eng. Rahmat Widyanto

General Approach for Combining Similarities

- Sometimes attributes are of many different types, but an overall similarity is needed.

1. For the k^{th} attribute, compute a similarity, s_k , in the range $[0, 1]$.

2. Define an indicator variable, δ_k , for the k_{th} attribute as follows:

$$\delta_k = \begin{cases} 0 & \text{if the } k^{th} \text{ attribute is a binary asymmetric attribute and both objects have} \\ & \text{a value of 0, or if one of the objects has a missing values for the } k^{th} \text{ attribute} \\ 1 & \text{otherwise} \end{cases}$$

3. Compute the overall similarity between the two objects using the following formula:

$$similarity(p, q) = \frac{\sum_{k=1}^n \delta_k s_k}{\sum_{k=1}^n \delta_k}$$

Dr.Eng. Rahmat Widyanto

Using Weights to Combine Similarities

- May not want to treat all attributes the same.
 - Use weights w_k which are between 0 and 1 and sum to 1.

$$similarity(p, q) = \frac{\sum_{k=1}^n w_k \delta_k s_k}{\sum_{k=1}^n \delta_k}$$

$$distance(p, q) = \left(\sum_{k=1}^n w_k |p_k - q_k|^r \right)^{1/r}$$

Dr.Eng. Rahmat Widyanto

Data Quality

- What kinds of data quality problems?
- How can we detect problems with the data?
- What can we do about these problems?

- Examples of data quality problems:
 - noise and outliers
 - missing values
 - duplicate data

Dr.Eng. Rahmat Widyanto

Why Is Data Dirty?

- Incomplete data may come from
 - “Not applicable” data value when collected
 - Different considerations between the time when the data was collected and when it is analyzed.
 - Human/hardware/software problems
- Noisy data (incorrect values) may come from
 - Faulty data collection instruments
 - Human or computer error at data entry
 - Errors in data transmission
- Inconsistent data may come from
 - Different data sources
 - Functional dependency violation (e.g., modify some linked data)
- Duplicate records also need data cleaning

Dr.Eng. Rahmat Widyanto

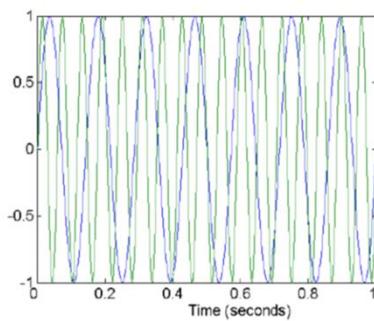
Why Is Data Preprocessing Important?

- **No quality data, no quality mining results!**
 - Quality decisions must be based on quality data
 - ◆ e.g., duplicate or missing data may cause incorrect or even misleading statistics.
 - Data warehouse needs consistent integration of quality data
- Data extraction, cleaning, and transformation comprises the majority of the work of building a data warehouse

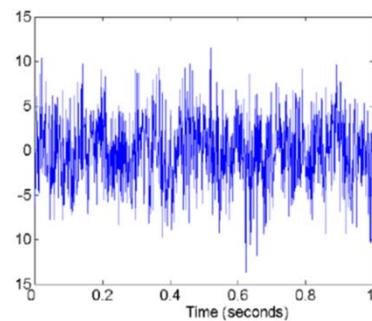
Dr.Eng. Rahmat Widyanto

Noise

- Noise refers to modification of original values
 - Examples: distortion of a person's voice when talking on a poor phone and "snow" on television screen



Two Sine Waves

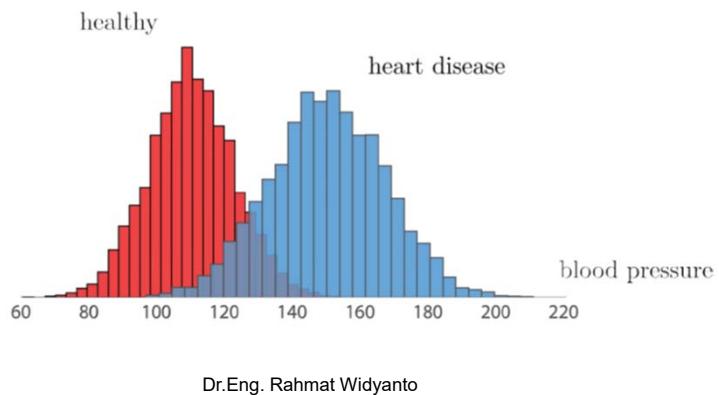


Two Sine Waves + Noise

Dr.Eng. Rahmat Widyanto

Noise vs. uncertainty

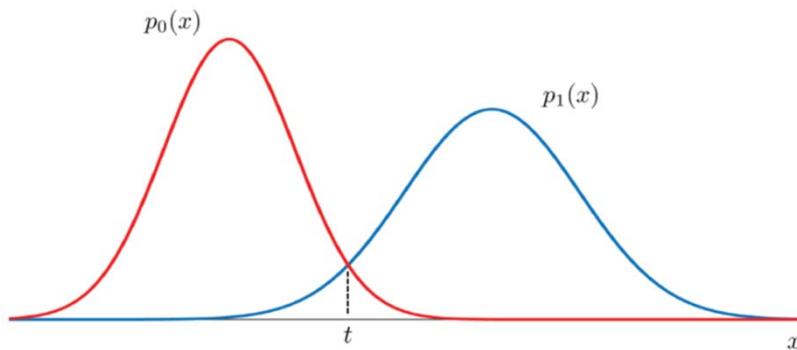
- Distribution overlap is commonly confused with noise
 - noise implies the true value is modified



Dr.Eng. Rahmat Widyatno

Noise vs. uncertainty

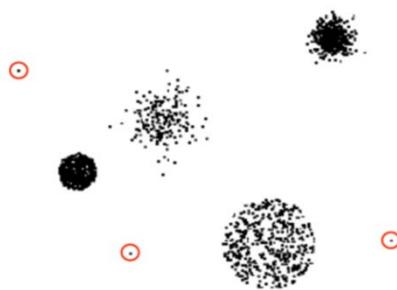
- Distribution overlap is commonly confused with noise
 - noise implies the true value is modified



Dr.Eng. Rahmat Widyatno

Outliers

- Outliers are data objects with characteristics that are considerably different than most of the other data objects in the data set



Dr.Eng. Rahmat Widyanto

Missing Values

- Reasons for missing values
 - Information is not collected (e.g., people decline to give their age and weight)
 - Attributes may not be applicable to all cases (e.g., annual income is not applicable to children)

Dr.Eng. Rahmat Widyanto

Missing Values

- Reasons for missing values
 - Information is not collected
(e.g., people decline to give their age and weight)
 - Attributes may not be applicable to all cases
(e.g., annual income is not applicable to children)
- Handling missing values
 - Eliminate data objects
 - Estimate missing values
 - Ignore the missing value during analysis
 - Replace with all possible values (weighted by their probabilities)

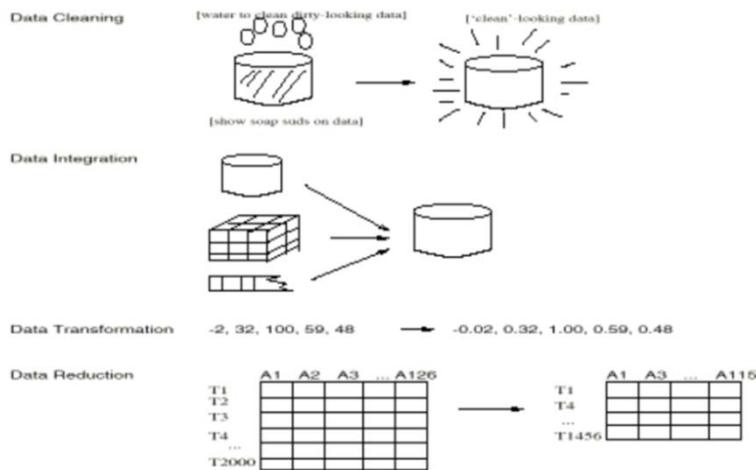
Dr.Eng. Rahmat Widyanto

Duplicate Data

- Data set may include data objects that are duplicates, or almost duplicates of one another
 - Major issue when merging data from heterogeneous sources
- Examples:
 - Same person with multiple email addresses
- Data cleaning
 - Process of dealing with duplicate data issues

Dr.Eng. Rahmat Widyanto

Forms of Data Preprocessing



Dr.Eng. Rahmat Widyanto

Data Preprocessing

- Integration
- Data cleaning
- Aggregation
- Sampling
- Dimensionality reduction
- Feature subset selection
- Feature creation
- Discretization and binarization
- Data transformation

Dr.Eng. Rahmat Widyanto

Data Cleaning

- Importance
 - “Data cleaning is one of the three biggest problems in data warehousing”—Ralph Kimball
 - “Data cleaning is the number one problem in data warehousing”—DCI survey

- Data cleaning tasks
 - Fill in missing values
 - Identify outliers and smooth out noisy data
 - Correct inconsistent data
 - Resolve redundancy caused by data integration

Dr.Eng. Rahmat Widyanto

How to Handle Missing Data?

- Ignore the tuple: usually done when class label is missing (assuming the tasks in classification—not effective when the percentage of missing values per attribute varies considerably).
- Fill in the missing value manually: tedious + infeasible?
- Fill in it automatically with
 - a global constant : e.g., “unknown”, a new class?!
 - the attribute mean
 - the attribute mean for all samples belonging to the same class: smarter
 - the most probable value: inference-based such as Bayesian formula or decision tree

Dr.Eng. Rahmat Widyanto

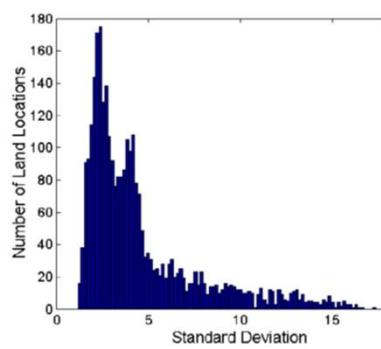
Aggregation

- Combining two or more attributes (or objects) into a single attribute (or object)
- Purpose
 - Data reduction
 - ◆ Reduce the number of attributes or objects
 - Change of scale
 - ◆ Cities aggregated into regions, states, countries, etc
 - More “stable” data
 - ◆ Aggregated data tends to have less variability

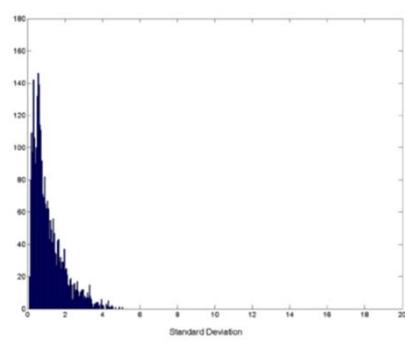
Dr.Eng. Rahmat Widyanto

Aggregation

Variation of Precipitation in Australia



Standard Deviation of Average
Monthly Precipitation



Standard Deviation of Average
Yearly Precipitation

Dr.Eng. Rahmat Widyanto

Data Integration

- Data integration:
 - Combines data from multiple sources into a coherent store
- Schema integration: e.g., A.cust-id ≡ B.cust-#
 - Integrate metadata from different sources
- Entity identification problem:
 - Identify real world entities from multiple data sources, e.g., Bill Clinton = William Clinton
- Detecting and resolving data value conflicts
 - For the same real world entity, attribute values from different sources are different
 - Possible reasons: different representations, different scales, e.g., metric vs. British units

Dr.Eng. Rahmat Widyanto

Handling Redundancy in Data Integration

- Redundant data occur often when integration of multiple databases
 - *Object identification:* The same attribute or object may have different names in different databases
 - *Derivable data:* One attribute may be a “derived” attribute in another table, e.g., annual revenue
- Redundant attributes may be able to be detected by *correlation analysis*
- Careful integration of the data from multiple sources may help reduce/avoid redundancies and inconsistencies and improve mining speed and quality

Dr.Eng. Rahmat Widyanto

Sampling

- Sampling is the main technique employed for data selection.
 - It is often used for both the preliminary investigation of the data and the final data analysis.
- Statisticians sample because **obtaining** the entire set of data of interest is too expensive or time consuming.
- Sampling is used in data mining because **processing** the entire set of data of interest is too expensive or time consuming.

Dr.Eng. Rahmat Widyanto

Sampling ...

- The key principle for effective sampling is the following:
 - using a sample will work almost as well as using the entire data sets, if the sample is representative
 - a sample is representative if it has approximately the same property (of interest) as the original set of data

Dr.Eng. Rahmat Widyanto

Types of Sampling

- Simple random sampling
 - There is an equal probability of selecting any particular item
 - Stratified sampling
 - Split the data into several partitions; then draw random samples from each partition
-
- Sampling without replacement
 - As each item is selected, it is removed from the population
 - Sampling with replacement
 - Objects are not removed from the population as they are selected for the sample. The same object can be picked up more than once.

Dr.Eng. Rahmat Widyanto

Sample Size



8000 points

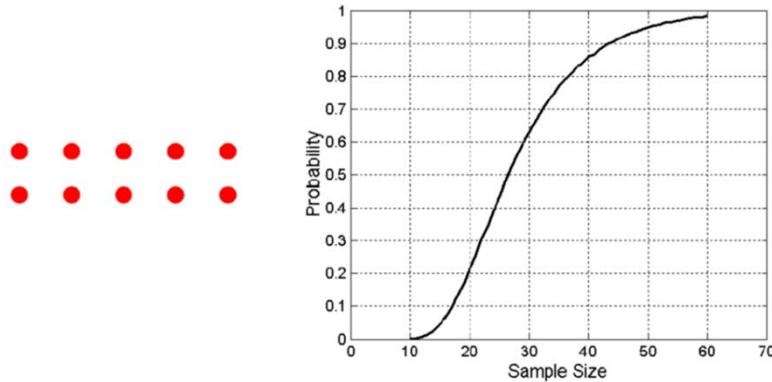
2000 Points

500 Points

Dr.Eng. Rahmat Widyanto

Sample Size

- What sample size is necessary to get at least one object from each of 10 groups.



Dr.Eng. Rahmat Widyanto

Dimensionality Reduction

- Purpose:

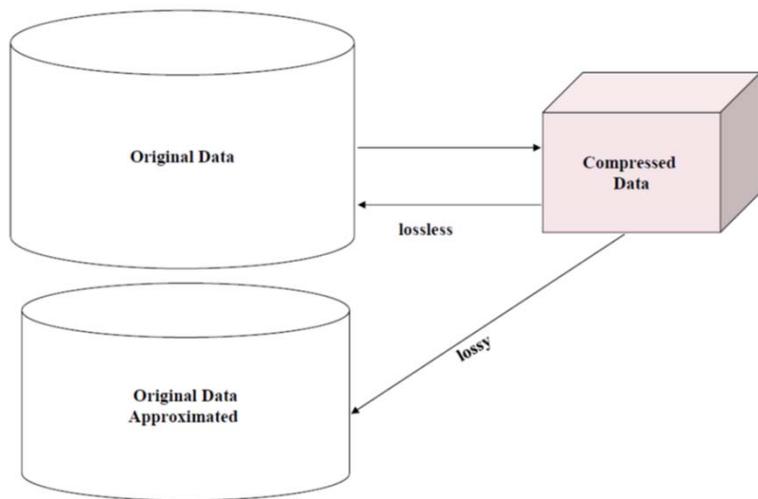
- avoid curse of dimensionality
- reduce amount of time and memory required by data mining algorithms
- allow data to be more easily visualized
- may help to eliminate irrelevant features or reduce noise
- may help to avoid stability problems

- Techniques

- Principal Component Analysis (PCA)
- Singular Value Decomposition (SVD)
- Others: supervised and non-linear techniques

Dr.Eng. Rahmat Widyanto

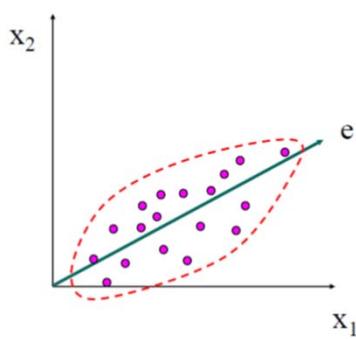
Relationship to Data Compression



Dr.Eng. Rahmat Widyatno

Dimensionality Reduction: PCA

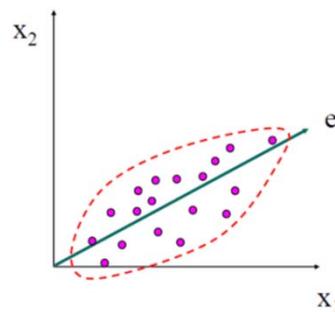
- Goal is to find a projection that captures the largest amount of variation in data



Dr.Eng. Rahmat Widyatno

Dimensionality Reduction: PCA

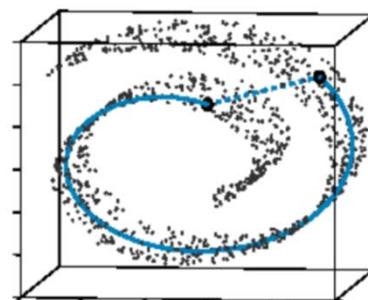
- Find the eigenvectors of the covariance matrix
- The eigenvectors define the new space



Dr.Eng. Rahmat Widyanto

Dimensionality Reduction: ISOMAP

By: Tenenbaum, de Silva, Langford (2000)



- Construct a neighbourhood graph
- For each pair of points in the graph, compute the shortest path distances – geodesic distances

Dr.Eng. Rahmat Widyanto

Feature Subset 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

Dr.Eng. Rahmat Widyanto

Feature Subset Selection

- **Brute-force approach:**
 - Try all possible feature subsets as input to data mining algorithm
- **Embedded approaches:**
 - Feature selection occurs naturally as part of the data mining algorithm
- **Filter approaches (usually one pass through data):**
 - Features are selected before data mining algorithm is run
- **Wrapper approaches (usually many passes through data):**
 - Use the data mining algorithm as a black box to find best subset of attributes

	The	Game	Play	Football	Baseball	Brady	Deflate	Gate
Document 1	12	2	3	14		4	4	6
Document 2	18	5	5		3			5
Document 3	24					4		5
Document 4	56	15					24	

Dr.Eng. Rahmat Widyanto

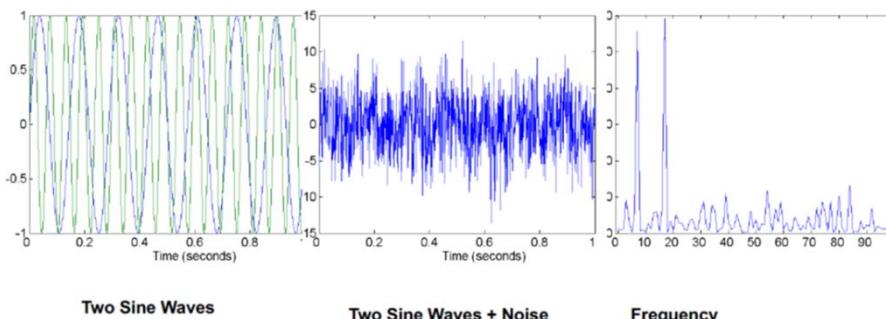
Feature Creation

- 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
 - ◆ domain-specific
 - Mapping Data to New Space
 - Feature Construction
 - ◆ combining features

Dr.Eng. Rahmat Widyanto

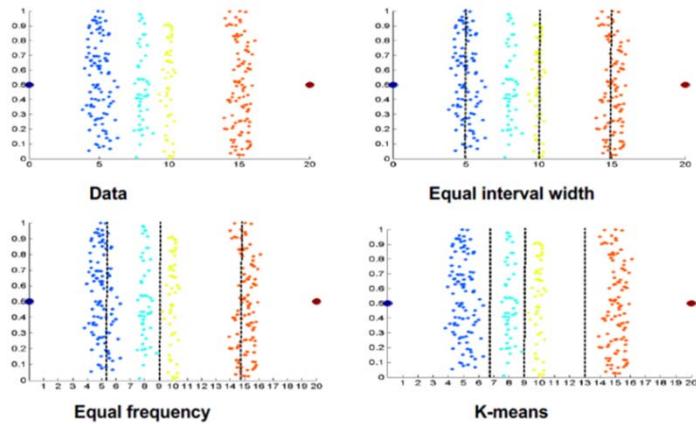
Mapping Data to a New Space

- Fourier transform
- Wavelet transform



Dr.Eng. Rahmat Widyanto

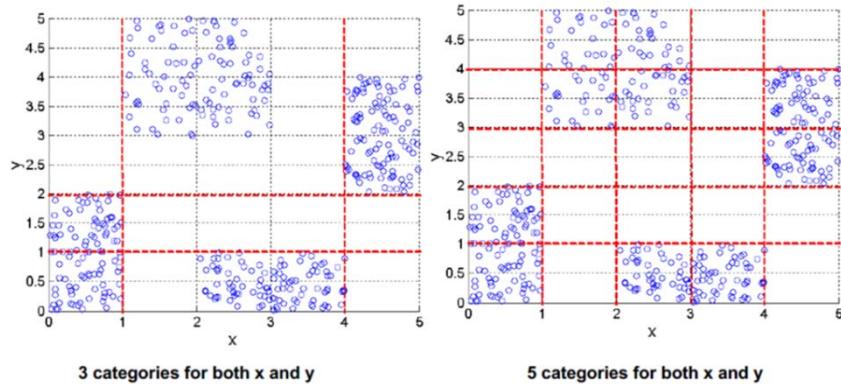
Discretization Without Using Class Labels



Dr.Eng. Rahmat Widyanto

Discretization Using Class Labels

- Entropy based approach



Dr.Eng. Rahmat Widyanto

How to Handle Noisy Data?

- Binning
 - first sort data and partition into (equal-frequency) bins
 - then one can **smooth by bin means, smooth by bin median, smooth by bin boundaries**, etc.
- Regression
 - smooth by fitting the data into regression functions
- Clustering
 - detect and remove outliers
- Combined computer and human inspection
 - detect suspicious values and check by human (e.g., deal with possible outliers)

Dr.Eng. Rahmat Widyanto

Simple Discretization Methods: Binning

- **Equal-width** (distance) partitioning
 - Divides the range into N intervals of equal size: uniform grid
 - if A and B are the lowest and highest values of the attribute, the width of intervals will be: $W = (B - A)/N$.
 - The most straightforward, but outliers may dominate presentation
 - Skewed data is not handled well
- **Equal-depth** (frequency) partitioning
 - Divides the range into N intervals, each containing approximately same number of samples
 - Good data scaling
 - Managing categorical attributes can be tricky

Dr.Eng. Rahmat Widyanto

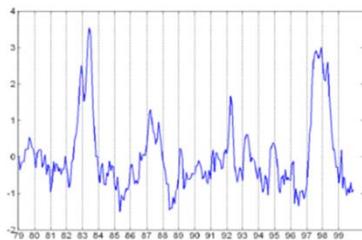
Binning Methods for Data Smoothing

- ❑ Sorted data for price (in dollars): 4, 8, 9, 15, 21, 21, 24, 25, 26, 28, 29, 34
 - * Partition into equal-frequency (equi-depth) bins:
 - **Bin 1:** 4, 8, 9, 15
 - **Bin 2:** 21, 21, 24, 25
 - **Bin 3:** 26, 28, 29, 34
 - * Smoothing by bin means:
 - **Bin 1:** 9, 9, 9, 9
 - **Bin 2:** 23, 23, 23, 23
 - **Bin 3:** 29, 29, 29, 29
 - * Smoothing by bin boundaries:
 - **Bin 1:** 4, 4, 4, 15
 - **Bin 2:** 21, 21, 25, 25
 - **Bin 3:** 26, 26, 26, 34

Dr.Eng. Rahmat Widyanto

Attribute Transformation

- 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
 - Simple functions: x^k , $\log(x)$, e^x , $|x|$
 - Standardization and normalization



Dr.Eng. Rahmat Widyanto

Data Transformation: Normalization

- Min-max normalization: to $[new_min_A, new_max_A]$

$$v' = \frac{v - min_A}{max_A - min_A} (new_max_A - new_min_A) + new_min_A$$

Ex. Let income range \$12,000 to \$98,000 normalized to [0.0, 1.0]. Then \$73,000 is mapped to $\frac{73,600 - 12,000}{98,000 - 12,000} (1.0 - 0) + 0 = 0.716$

- Z-score normalization (μ : mean, σ : standard deviation):

$$v' = \frac{v - \mu}{\sigma}$$

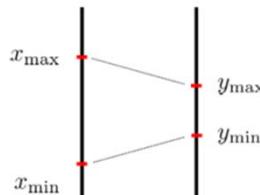
Ex. Let $\mu = 54,000$, $\sigma = 16,000$. Then $\frac{73,600 - 54,000}{16,000} = 1.225$

- Normalization by decimal scaling

$$v' = \frac{v}{10^j} \quad \text{Where } j \text{ is the smallest integer such that } \text{Max}(|v'|) < 1$$

Dr.Eng. Rahmat Widyatno

Deriving the Min-Max Normalization



Find linear transform:

$$y = ax + b$$

$$x_{\min} \rightarrow y_{\min}$$

$$x_{\max} \rightarrow y_{\max}$$

Equation of a line given two points;

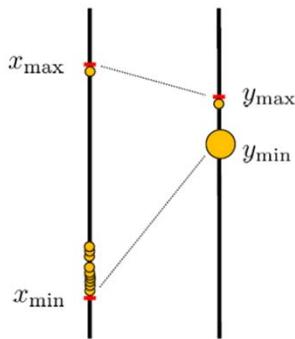
i.e., (x_1, y_1) and (x_2, y_2)

$$y - y_1 = \frac{y_2 - y_1}{x_2 - x_1} (x - x_1)$$

You take it from here...

Dr.Eng. Rahmat Widyatno

Min-max Normalization: Problems?



Dr.Eng. Rahmat Widyanto

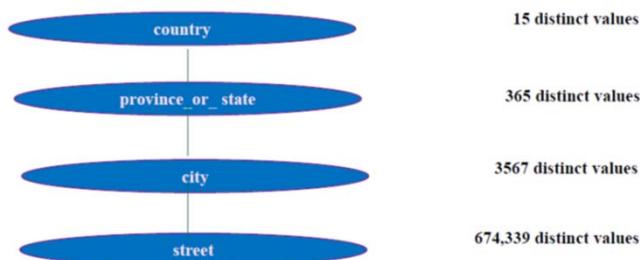
Concept Hierarchy Generation for Categorical Data

- Specification of a partial/total ordering of attributes explicitly at the schema level by users or experts
 - street < city < state < country
- Specification of a hierarchy for a set of values by explicit data grouping
 - {Bloomington, Indianapolis, South Bend} < Indiana
- Specification of only a partial set of attributes
 - E.g., only street < city, not others
- Automatic generation of hierarchies (or attribute levels) by the analysis of the number of distinct values
 - E.g., for a set of attributes: {street, city, state, country}

Dr.Eng. Rahmat Widyanto

Automatic Concept Hierarchy Generation

- Some hierarchies can be automatically generated based on the analysis of the number of distinct values per attribute in the data set
 - The attribute with the most distinct values is placed at the lowest level of the hierarchy
 - Exceptions, e.g., weekday, month, quarter, year



Dr.Eng. Rahmat Widyanto

Density

- Density-based clustering require a notion of density
- Examples:
 - Euclidean density
 - Euclidean density = number of points per unit volume
 - Probability density
 - Graph-based density

Dr.Eng. Rahmat Widyanto

Euclidean Density – Cell-based

- Simplest approach is to divide region into a number of rectangular cells of equal volume and define density as # of points the cell contains

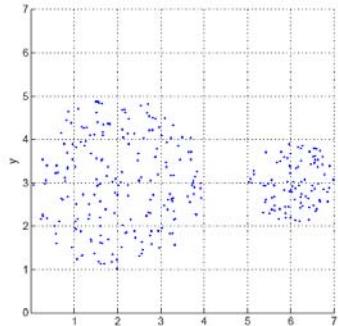


Figure 7.13. Cell-based density.

0	0	0	0	0	0	0
0	0	0	0	0	0	0
4	17	18	6	0	0	0
14	14	13	13	0	18	27
11	18	10	21	0	24	31
3	20	14	4	0	0	0
0	0	0	0	0	0	0

Table 7.6. Point counts for each grid cell.

Dr.Eng. Rahmat Widyanto

Euclidean Density – Center-based

- Euclidean density is the number of points within a specified radius of the point

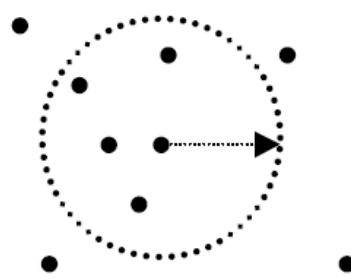


Figure 7.14. Illustration of center-based density.

Dr.Eng. Rahmat Widyanto

Summary

- Data preparation or preprocessing is a big issue for both data warehousing and data mining
- Descriptive data summarization is need for quality data preprocessing
- Data preparation includes
 - Data cleaning and data integration
 - Data reduction and feature selection
 - Discretization
- A lot a methods have been developed but data preprocessing still an active area of research

Dr.Eng. Rahmat Widyanto