

# Data Mining: Data

---

## Lecture Notes for Chapter 2

Introduction to Data Mining , 2<sup>nd</sup> Edition  
by  
Tan, Steinbach, Kumar

# Outline

---

- Attributes and Objects
- Types of Data
- Data Quality
- Data Preprocessing
- Similarity and Distance

# 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, dimension, 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

**Objects**

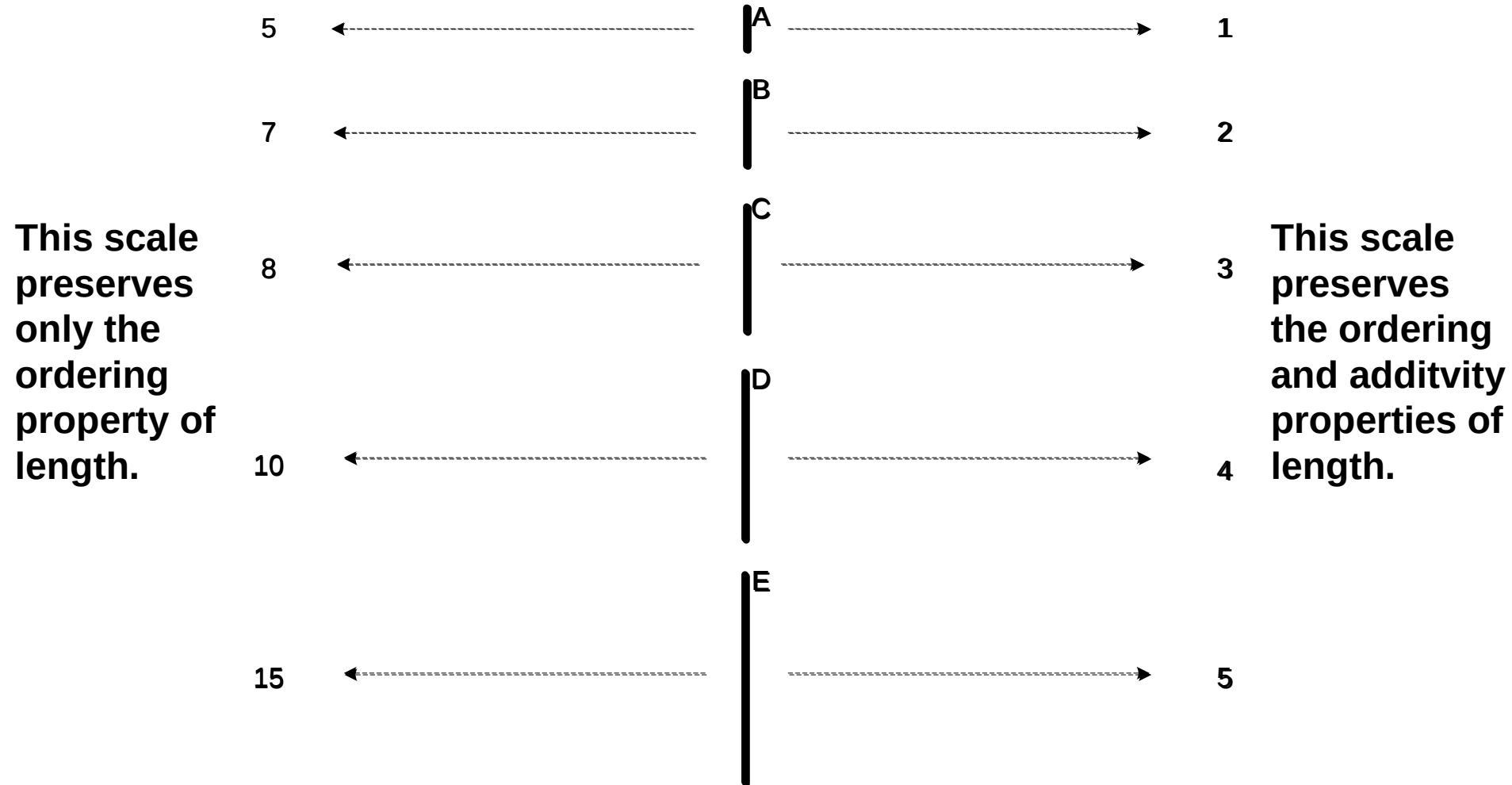
# Attribute Values

---

- *Attribute values* are numbers or symbols assigned to an attribute for a particular object
- 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 can be different than the properties of the values used to represent the attribute

# Measurement of Length

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



# 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 {tall, medium, short}
  - **Interval**
    - ◆ Examples: calendar dates, temperatures in Celsius or Fahrenheit.
  - **Ratio**
    - ◆ Examples: temperature in Kelvin, length, counts, elapsed time (e.g., time to run a race)

# Properties of Attribute Values

---

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

# Types of Attributes - Nominal

---

## Examples of Nominal Variables

- Gender (Male, Female, Transgender).
- Eye color (Blue, Green, Brown, Hazel).
- Type of house (Bungalow, Duplex, Ranch).
- Type of pet (Dog, Cat, Rodent, Fish, Bird).
- Genotype ( AA, Aa, or aa).



# Types of Attributes - Ordinal

---

## Examples:

- **High school class ranking:** 1st, 9th, 87th...
- **Socioeconomic status:** poor, middle class, rich.
- The **Likert Scale**: strongly disagree, disagree, neutral, agree, strongly agree.
- **Level of Agreement:** yes, maybe, no.
- **Time of Day:** dawn, morning, noon, afternoon, evening, night.
- **Political Orientation:** left, center, right.

# Types of Attributes - Interval

---

## Examples:

- Celsius Temperature.
- Fahrenheit Temperature.
- IQ (intelligence scale).
- SAT scores.
- Time on a clock with hands.

# Types of Attributes - Ratio

---

## Examples:

- Age.\*
- Weight.
- Height.
- Sales Figures.
- Ruler measurements.
- Income earned in a week.
- Years of education.
- Number of children.

		Attribute Type	Description	Examples	Operations
Categorical Qualitative		Nominal	Nominal attribute values only distinguish. (=, $\neq$ )	zip codes, employee ID numbers, eye color, sex: { <i>male</i> , <i>female</i> }	mode, entropy, contingency correlation, $\chi^2$ test
		Ordinal	Ordinal attribute values also order objects. (<, >)	hardness of minerals, { <i>good</i> , <i>better</i> , <i>best</i> }, grades, street numbers	median, percentiles, rank correlation, run tests, sign tests
Numeric Quantitative		Interval	For interval attributes, differences between values are meaningful. (+, -)	calendar dates, temperature in Celsius or Fahrenheit	mean, standard deviation, Pearson's correlation, <i>t</i> and <i>F</i> tests
		Ratio	For ratio variables, both differences and ratios are meaningful. (*, /)	temperature in Kelvin, monetary quantities, counts, age, mass, length, current	geometric mean, harmonic mean, percent variation

**This categorization of attributes is due to S. S. Stevens**

		Attribute Type	Transformation	Comments
Categorical Qualitative		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}.
Numeric Quantitative		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.

**This categorization of attributes is due to S. S. Stevens**

# 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.

# Important Characteristics of Data

---

- Dimensionality (number of attributes)
  - ◆ High dimensional data brings a number of challenges
- Sparsity
  - ◆ Only presence counts
- Resolution
  - ◆ Patterns depend on the scale
- Size
  - ◆ Type of analysis may depend on size of data

# 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



# Record Data

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

<i>Tid</i>	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

# 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 a 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

# 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	play	ball	score	game	win	lost	timeout	season
Document 1	3	0	5	0	2	6	0	2	0	2
Document 2	0	7	0	2	1	0	0	3	0	0
Document 3	0	1	0	0	1	2	2	0	3	0

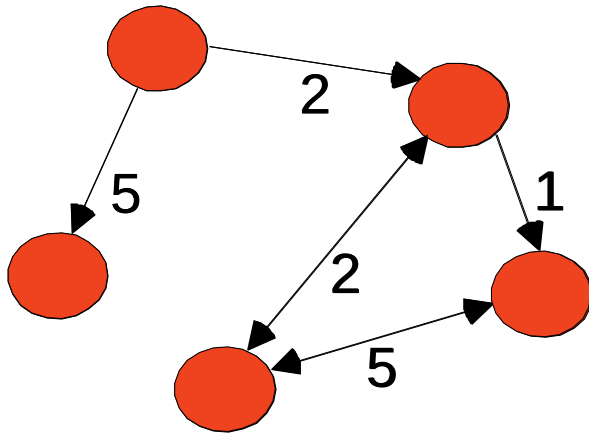
# Transaction Data

- A special type of data, where
  - Each 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.
  - Can represent transaction data as record data

<i><b>TID</b></i>	<i><b>Items</b></i>
<b>1</b>	<b>Bread, Coke, Milk</b>
<b>2</b>	<b>Beer, Bread</b>
<b>3</b>	<b>Beer, Coke, Diaper, Milk</b>
<b>4</b>	<b>Beer, Bread, Diaper, Milk</b>
<b>5</b>	<b>Coke, Diaper, Milk</b>

# Graph Data

- Examples: Generic graph, a molecule, and webpages



## Useful Links:

- [Bibliography](#)
- Other Useful Web sites
  - [ACM SIGKDD](#)
  - [KDnuggets](#)
  - [The Data Mine](#)

## Knowledge Discovery and Data Mining Bibliography

(Gets updated frequently, so visit often!)

- [Books](#)
- [General Data Mining](#)

## Book References in Data Mining and Knowledge Discovery

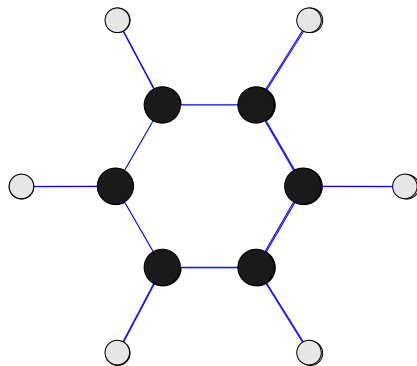
Usama Fayyad, Gregory Piatetsky-Shapiro, Padhraic Smyth, and Ramasamy uthurasamy, "Advances in Knowledge Discovery and Data Mining", AAAI Press/the MIT Press, 1996.

J. Ross Quinlan, "C4.5: Programs for Machine Learning", Morgan Kaufmann Publishers, 1993.  
Michael Berry and Gordon Linoff, "Data Mining Techniques (For Marketing, Sales, and Customer Support)", John Wiley & Sons, 1997.

## General Data Mining

Usama Fayyad, "Mining Databases: Towards Algorithms for Knowledge Discovery", Bulletin of the IEEE Computer Society Technical Committee on data Engineering, vol. 21, no. 1, March 1998.

Christopher Matheus, Philip Chan, and Gregory Piatetsky-Shapiro, "Systems for knowledge Discovery in databases", IEEE Transactions on Knowledge and Data Engineering, 5(6):903-913, December 1993.



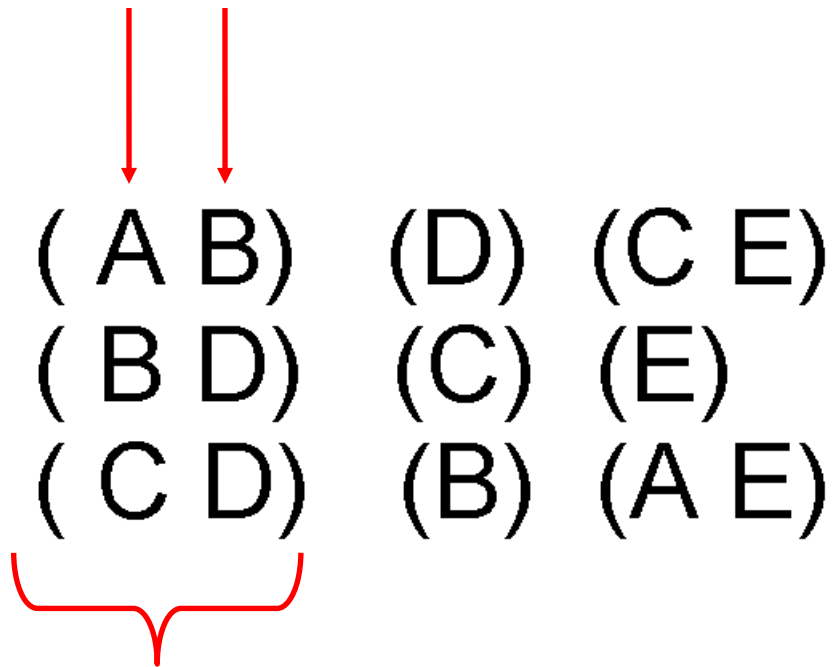
Benzene Molecule: C<sub>6</sub>H<sub>6</sub>

# Ordered Data

---

- Sequences of transactions

Items/Events



**An element of  
the sequence**

# Ordered Data

---

- Genomic sequence data

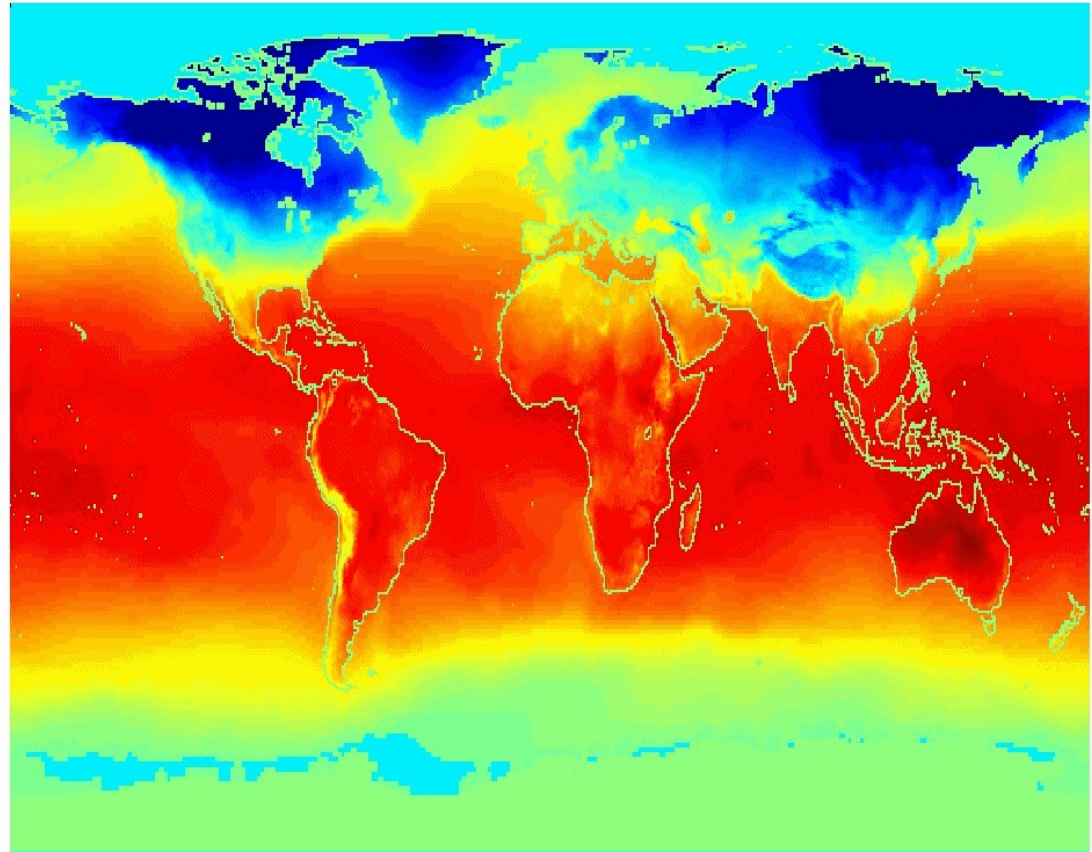
**GGTTCCGCCTTCAGCCCCGCGCC  
CGCAGGGCCCCGCCCCGCGCCGTC  
GAGAAGGGCCCCGCCTGGCGGGCG  
GGGGGAGGCGGGGGCCGCCCGAGC  
CCAACCGAGTCCGACCAGGTGCC  
CCCTCTGCTCGGCCTAGACCTGA  
GCTCATTAGGCGGCAGCGGACAG  
GCCAAGTAGAACACGCGAAGCGC  
TGGGCTGCCTGCTGCGACCAGGG**

# Ordered Data

- Spatio-Temporal Data

**Average Monthly  
Temperature of  
land and ocean**

Jan





# Data Quality

---

- Poor data quality negatively affects many data processing efforts

Data mining algorithms gives results(extracts) only what is there in the data.

If data quality issues are not handled carefully, then Data mining algorithms will produce erroneous or spurious output.

- Data mining example: a classification model for detecting people who are loan risks is built using poor data
  - Some credit-worthy candidates are denied loans
  - More loans are given to individuals that default

# Data Quality ..

---

- To overcome the poor data quality problem, Data mining focuses on:
  - 1) The detection and correction of data quality problem ( is often called **data cleaning**)
  - 2) The use of algorithms that can tolerate poor data quality

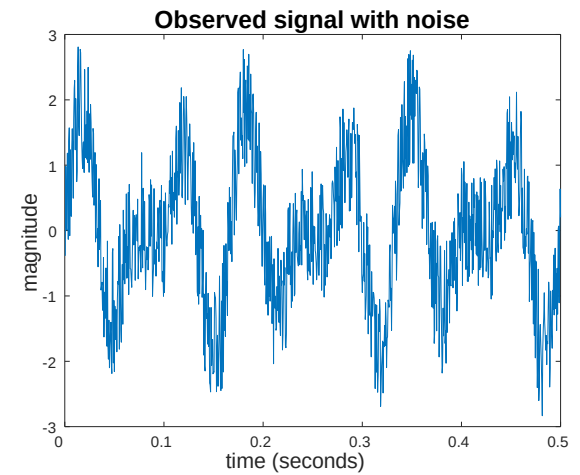
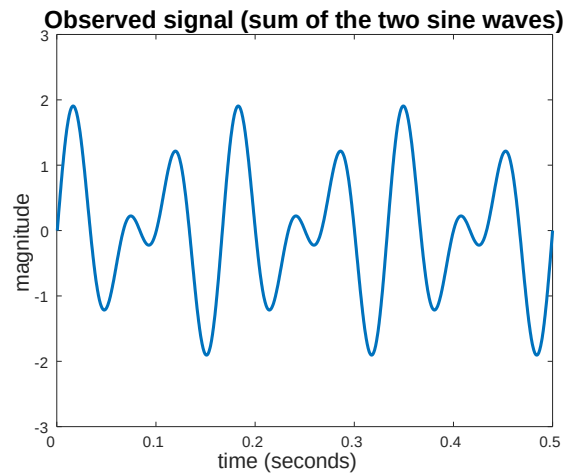
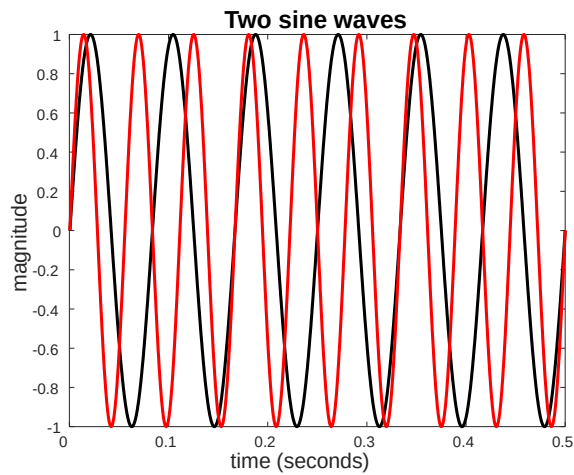
# Data Quality ...

---

- 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**
  - Wrong data
  - Fake data
  - **Missing values**
  - **Duplicate data**

# Noise

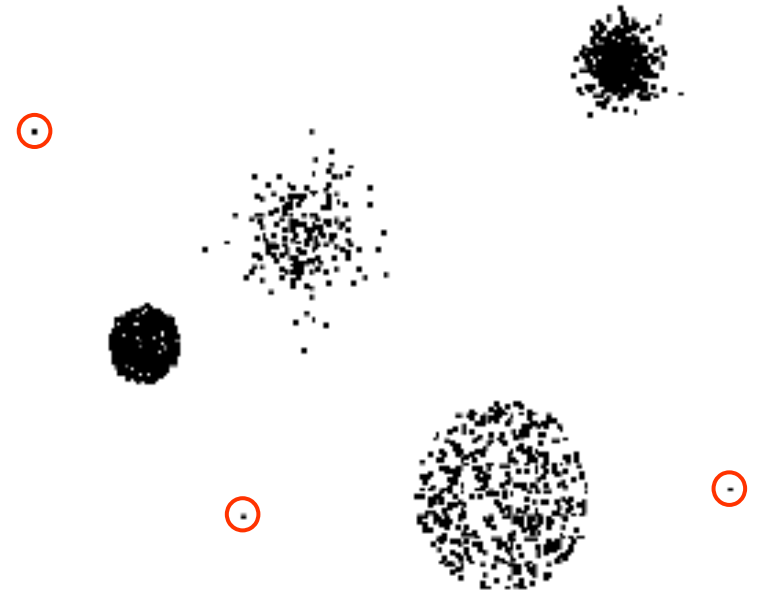
- For objects, noise is an extraneous object
- For attributes, noise refers to modification of original values
  - Examples: distortion of a person's voice when talking on a poor phone and "snow" on television screen
  - The figures below show two sine waves of the same magnitude and different frequencies, the waves combined, and the two sine waves with random noise
    - ♦ The magnitude and shape of the original signal is distorted



# Outliers

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

***For example:** In fraud and network intrusion detection, the goal is to find unusual objects or events from among a large number of normal ones.*



# 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 or variables
  - Estimate missing values
    - ◆ Example: time series of temperature
    - ◆ Example: census results
  - Ignore the missing value during analysis
  - Replace with all possible values(weighted by their probabilities)

# 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

# Data Quality - In a nutshell

---

Data mining algorithms gives results(extracts) only what is there in the data.

If data quality issues are not handled carefully, then Data mining algorithms will produce erroneous or spurious output.

So the **Preprocessing** is indeed a very important step to solve the data quality problems. (next topic)



# Data Preprocessing

---

- Aggregation
- Sampling
- Dimensionality Reduction
- Feature Subset Selection
- Feature Creation
- Discretization and Binarization
- Variable Transformation

# 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.
    - ◆ Days aggregated into weeks, months, or years
  - **More “stable” data** - aggregated data tends to have less variability

**Table 2.4.** Data set containing information about customer purchases.

Transaction ID	Item	Store Location	Date	Price	...
⋮	⋮	⋮	⋮	⋮	
101123	Watch	Chicago	09/06/04	\$25.99	...
101123	Battery	Chicago	09/06/04	\$5.99	...
101124	Shoes	Minneapolis	09/06/04	\$75.00	...
⋮	⋮	⋮	⋮	⋮	

# Aggregation

---

- An obvious issue is how an aggregate transaction is created?
- Quantitative attributes
  - such as **price**, are typically aggregated by taking a **sum** or an **average**
- Qualitative attributes
  - such as **item**, can either be omitted or summarized in terms of a higher level category, e.g., televisions versus electronics
- Disadvantages of aggregation
  - Potential loss of interesting details
  - In store example: aggregation over months loses information about which day of the week has the highest sales.

# Example: Precipitation in Australia

---

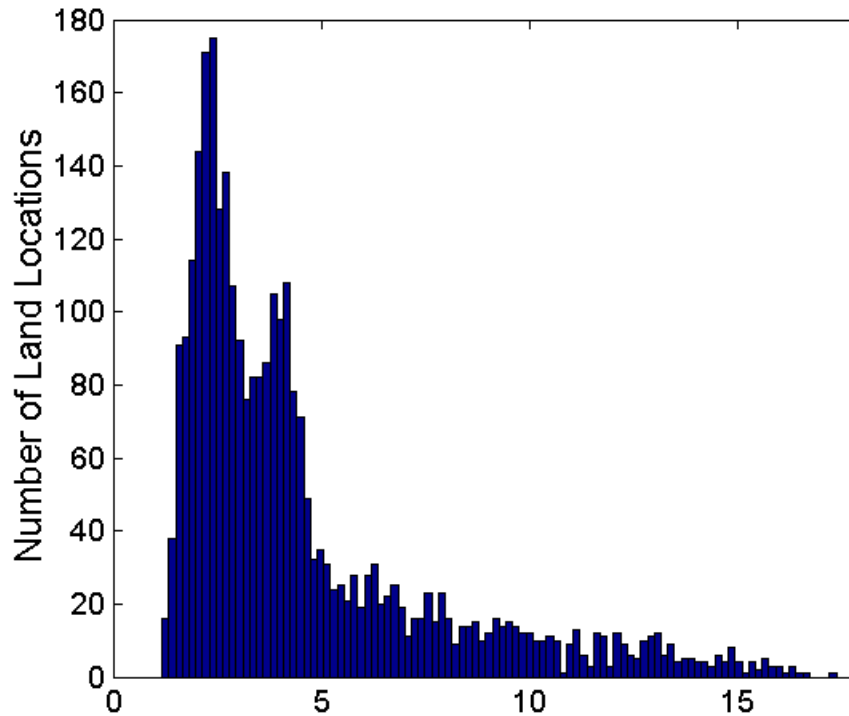
- This example is based on precipitation in Australia from the period 1982 to 1993.

The next slide shows

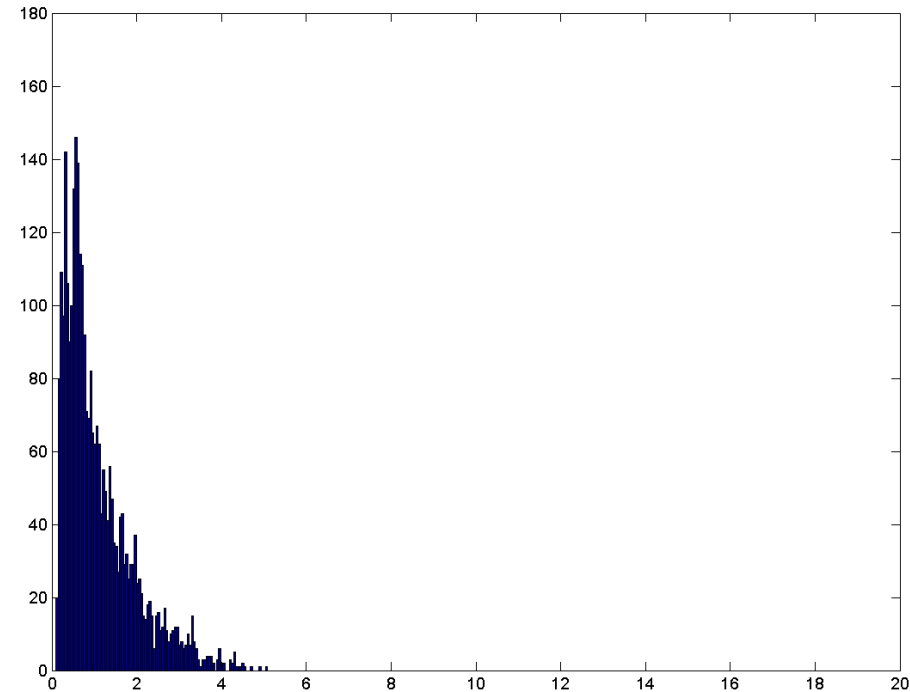
- A histogram for the standard deviation of average monthly precipitation for 3,030  $0.5^\circ$  by  $0.5^\circ$  grid cells in Australia, and
  - A histogram for the standard deviation of the average yearly precipitation for the same locations.
- The average yearly precipitation has **less variability** than the average monthly precipitation.
- All precipitation measurements (and their standard deviations) are in centimeters.

# Example: Precipitation in Australia ...

## Variation of Precipitation in Australia



**Standard Deviation of Average  
Monthly Precipitation**



**Standard Deviation of  
Average Yearly Precipitation**

# Sampling

---

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

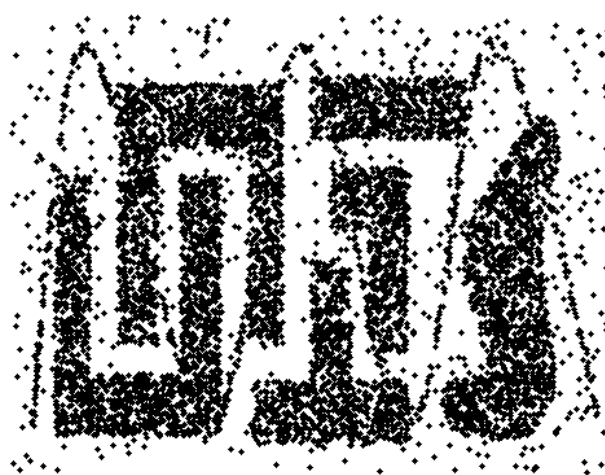
# Sampling ...

---

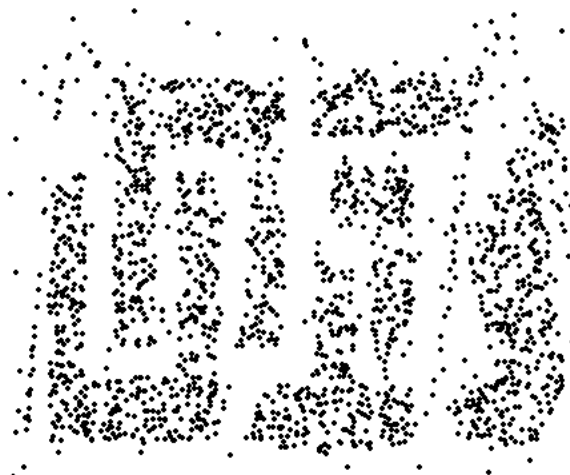
- The key principle for effective sampling is the following:
  - Using a sample will work almost as well as using the entire data set, if the sample is **representative**
  - A sample is **representative** if it has approximately the same properties (of interest) as the original set of data

# Sample Size

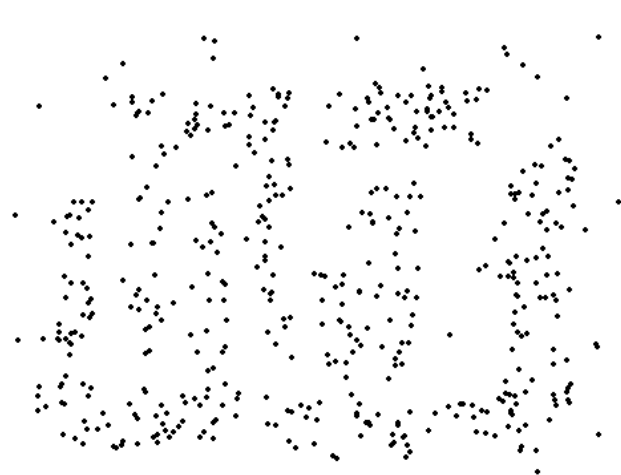
---



8000 points



2000 Points



500 Points

Figure 2.9. Example of the loss of structure with sampling



# 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

# Types of Sampling

- Simple Random Sampling



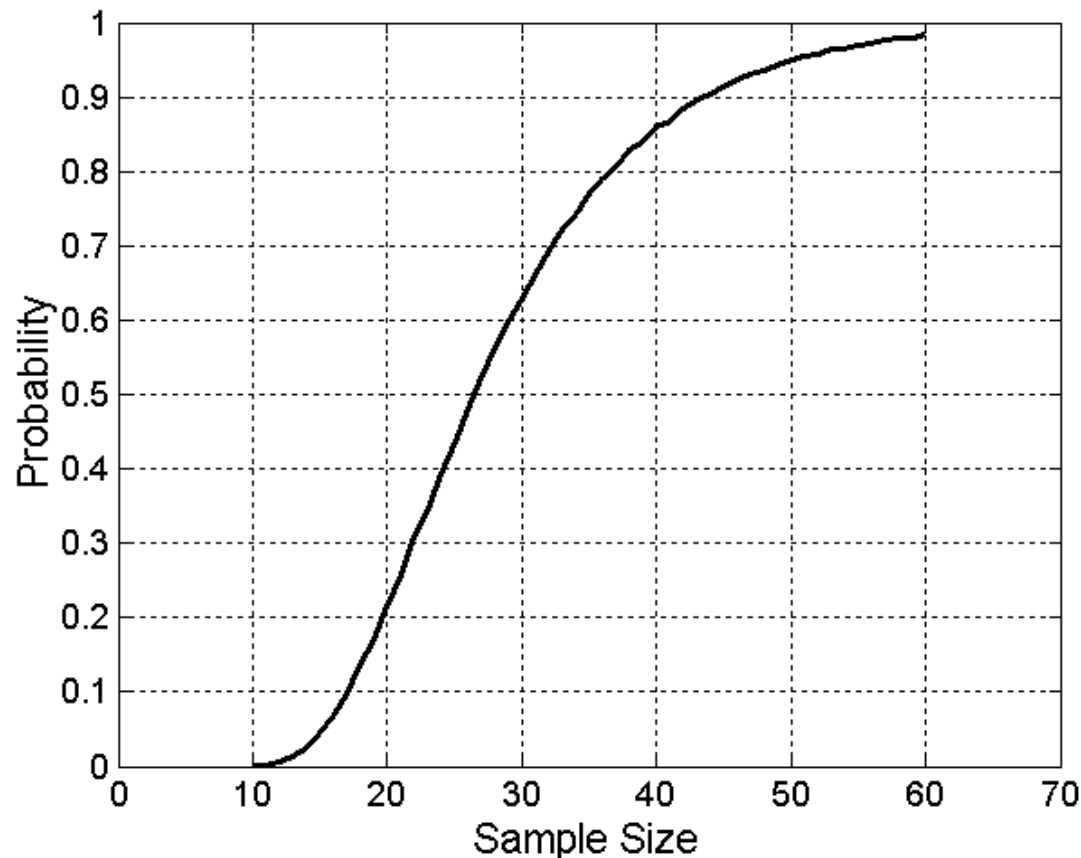
# Types of Sampling

- Stratified Sampling



# Sample Size

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



# Curse of Dimensionality

---

- When dimensionality increases, data becomes increasingly sparse in the space that it occupies
- Definitions of density and distance between points, which are critical for clustering and outlier detection, become less meaningful
- Many clustering and classification algorithms have trouble with high-dimensional data leading to reduced classification accuracy and poor quality clusters.

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

# Dimensionality Reduction: PCA

---

- A data reduction technique that transforms a large number of correlated variables into a smaller set of correlated variables called principal components
  - a method of extracting important variables from a large number of variables available in a dataset
  - it extracts a set of low-dimensional features from a high-dimensional dataset with the goal of capturing as much information as possible(variance) in the data.

# Dimensionality Reduction: PCA

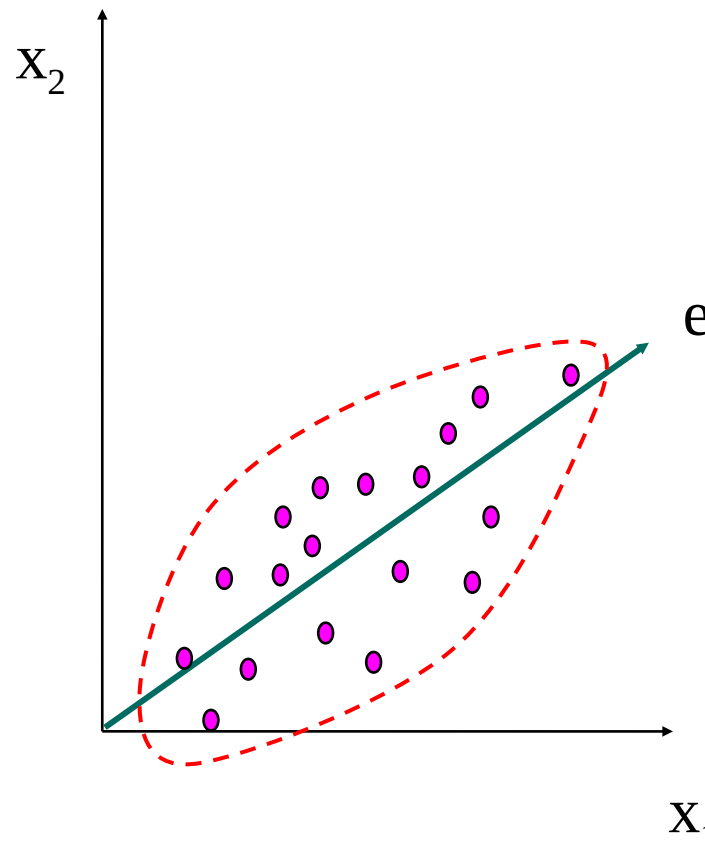
---

- Steps Involved in the Principal Component Analysis:
  - Standardize the dataset
  - Compute the covariance matrix for the features in the dataset
  - Compute the eigenvalues and eigenvectors for the covariance matrix
  - Sort the eigenvalues and their corresponding eigenvectors
  - Choose  $k$  eigenvalues to form an eigenvector matrix
  - Transform the original matrix



# Dimensionality Reduction: PCA

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



# Feature Subset Selection

---

- Another way to reduce dimensionality of data
  - Use only a subset of the features
- 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
- Redundant and irrelevant features can reduce classification accuracy and the quality of the clusters that are found.

# Feature Subset Selection

---

- Techniques

- } Brute-Force approach:

- Try all possible feature subsets as input to data mining algorithm, and then take the subset that produces the best results

- } **Embedded approaches:**

- Feature selection occurs naturally as part of the data mining algorithm
      - } the data mining algorithm itself decides which attributes to use and which to ignore.
      - } For example: Algorithms for building decision tree classifiers

# Feature Subset Selection

---

- Techniques

- } **Filter approaches:**

- Feature are selected before data mining algorithm is run
      - } Using some approach that is independent of the data mining task.
      - } For example: select sets of attributes whose pairwise correlation is as low as possible.

- } **Wrapper approaches:**

- Use the data mining algorithm as a black box to find best subset of attributes

# An Architecture for Feature Subset Selection

- It is possible to encompass both the filter and wrapper approaches within a common architecture
- The feature selection process view as consisting of four parts:
  - } A measure for evaluating a subset
    - Filter methods and Wrapper methods differ only in the way in which they evaluate a subset of features
  - } A search strategy that controls the generation of a new subset of features
  - } A stopping criterion
  - } A validation procedure

# An Architecture for Feature Subset Selection

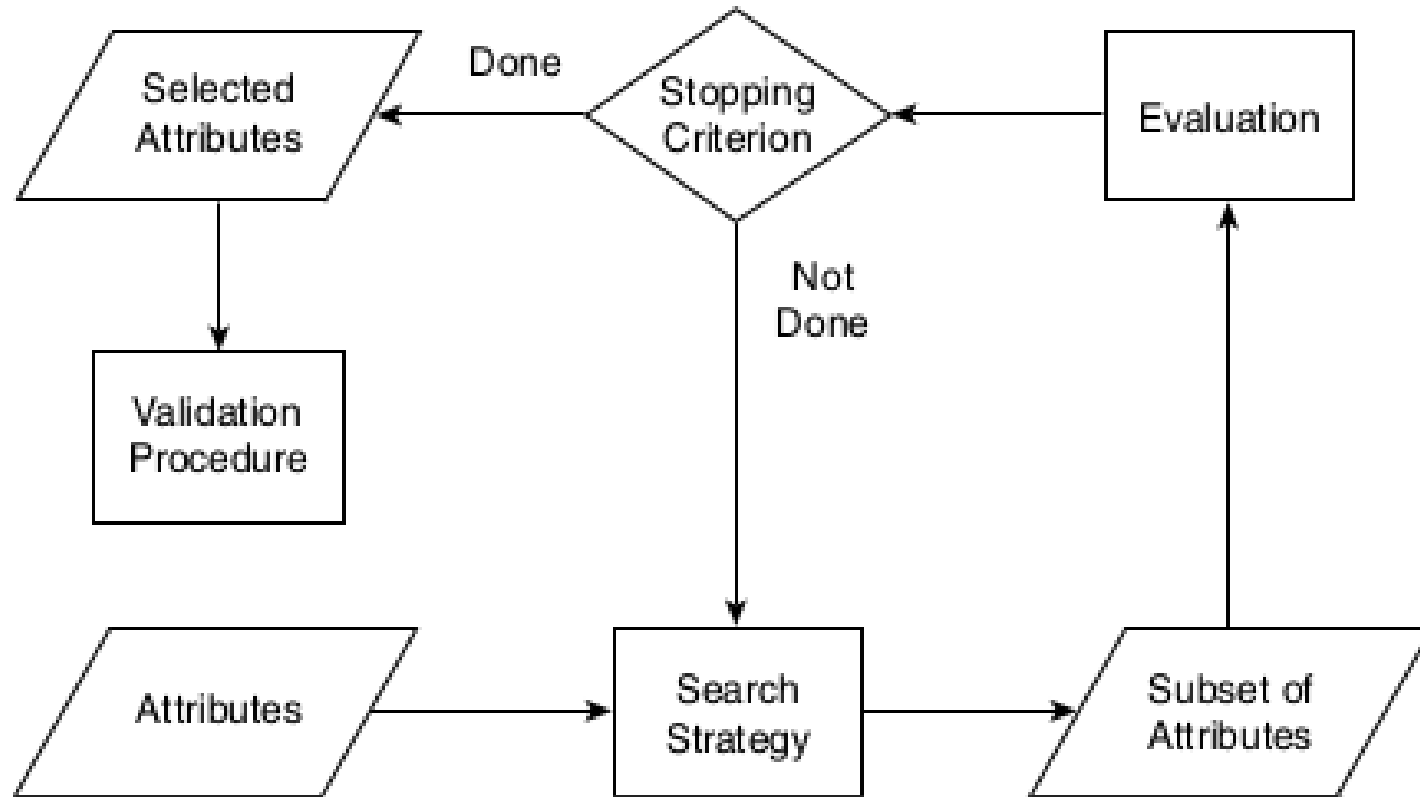


Figure 2.11. Flowchart of a feature subset selection process.

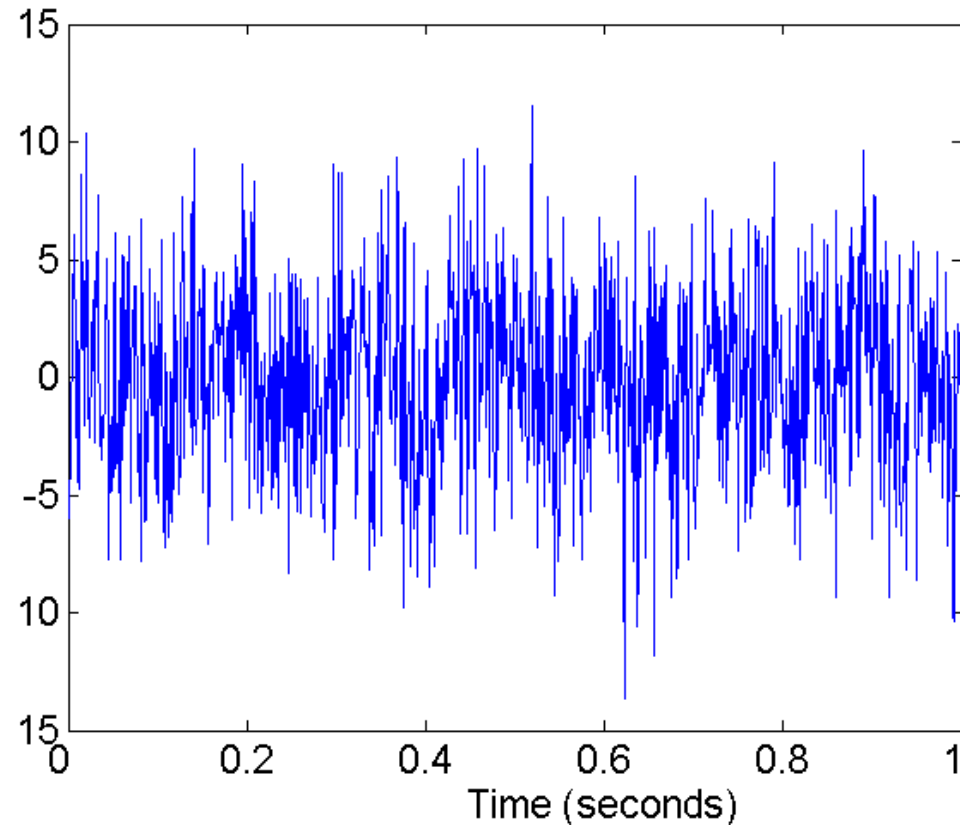
# Feature Creation

---

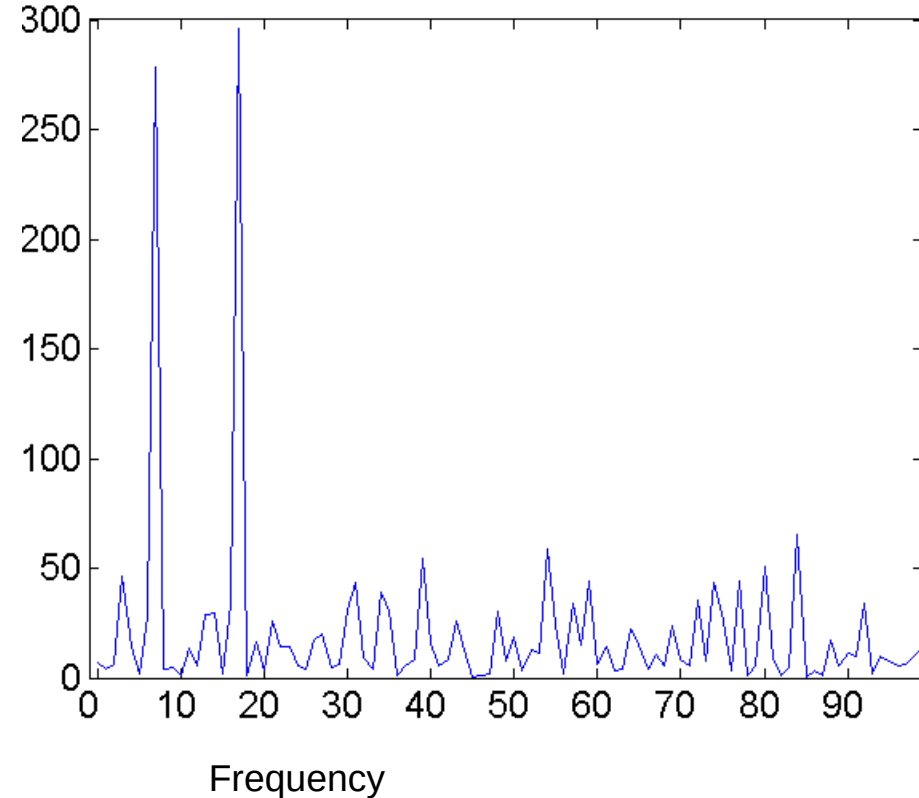
- Create new attributes that can capture the important information in a data set much more efficiently than the original attributes
- Three general methodologies:
  - Feature extraction
    - ◆ Example: extracting edges from images
  - Feature construction
    - ◆ Example: dividing mass by volume to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

# Mapping Data to a New Space

- **Fourier and wavelet transform**



**Two Sine Waves + Noise**



**Frequency**



# Discretization

---

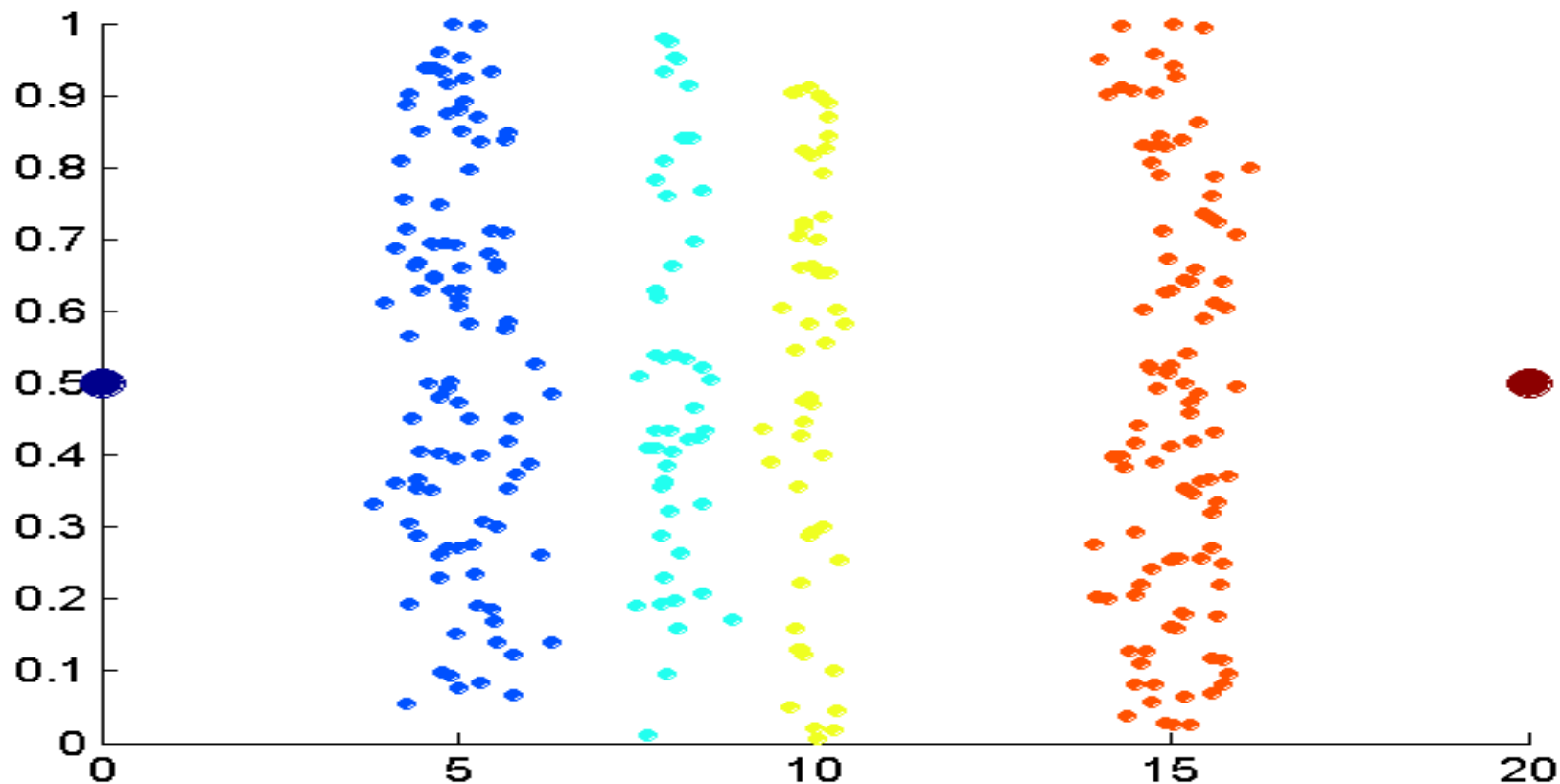
- **Discretization** is the process of converting a continuous attribute into a categorical attribute
  - A potentially infinite number of values are mapped into a small number of categories
  - Discretization is used in both unsupervised and supervised settings
- Discretization is typically applied to attributes that are used in *classification* or *association analysis*

# Discretization of continuous attributes

---

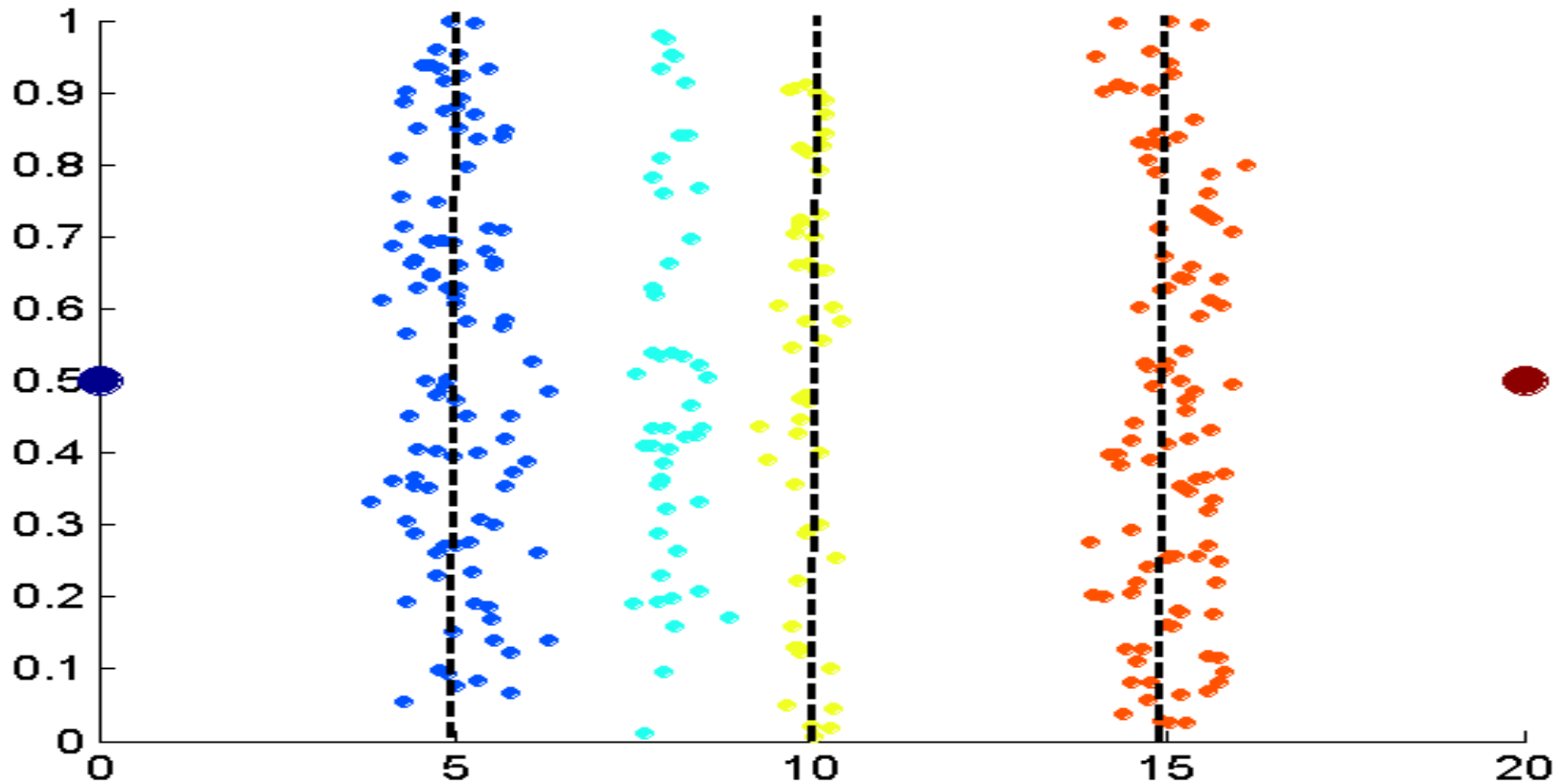
- Transformation of a continuous attribute to a categorical attribute involves two subtasks:
  - deciding how many categories,  $n$ , to have
  - determining how to map the values of the continuous attribute to these categories.
- In the **first step**, after the values of the continuous attribute are sorted, they are then divided into  $n$  intervals by specifying  $n - 1$  **split points**.
- In the **second step**, all the values in one interval are mapped to the same categorical value.

# Unsupervised Discretization



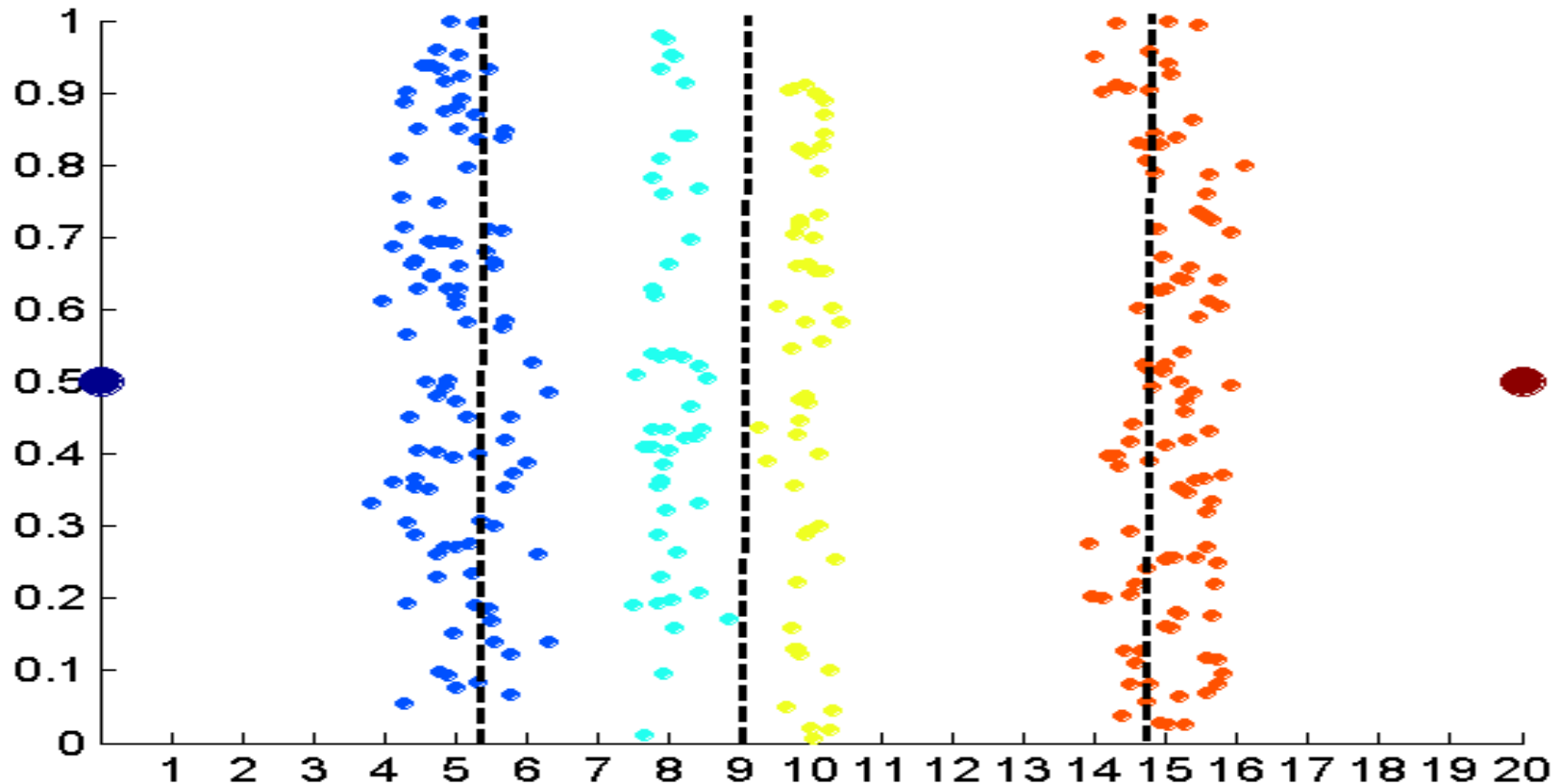
Data consists of four groups of points and two outliers. Data is one-dimensional, but a random y component is added to reduce overlap.

# Unsupervised Discretization



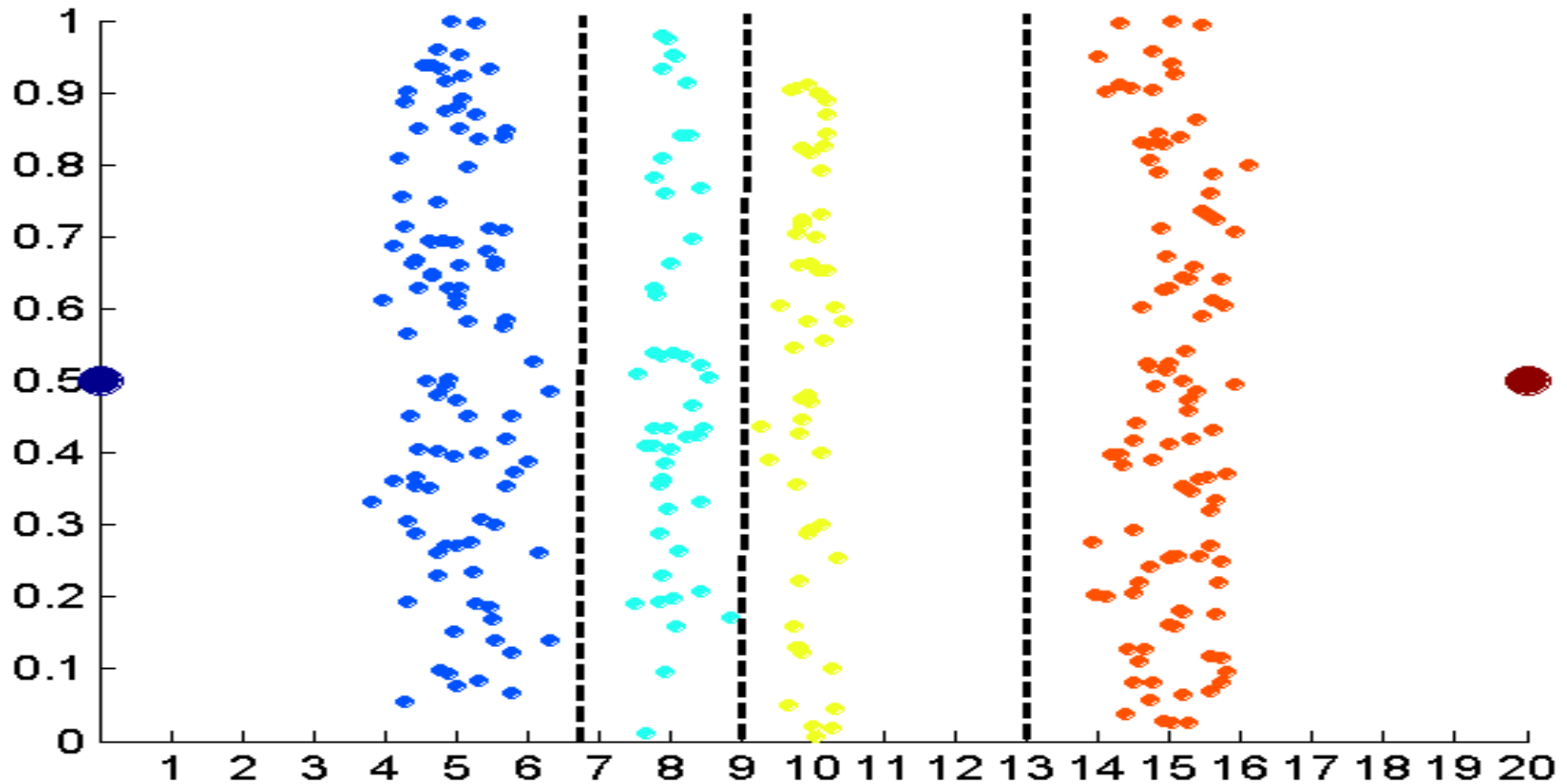
**Equal interval width** approach used to obtain 4 values.

# Unsupervised Discretization



**Equal frequency(equal depth)** approach used to obtain 4 values.

# Unsupervised Discretization



**K-means** approach to obtain 4 values.

# Binarization

- Binarization maps a continuous or categorical attribute into one or more binary variables

**Table 2.5.** Conversion of a categorical attribute to three binary attributes.

Categorical Value	Integer Value	$x_1$	$x_2$	$x_3$
<i>awful</i>	0	0	0	0
<i>poor</i>	1	0	0	1
<i>OK</i>	2	0	1	0
<i>good</i>	3	0	1	1
<i>great</i>	4	1	0	0

**Table 2.6.** Conversion of a categorical attribute to five asymmetric binary attributes.

Categorical Value	Integer Value	$x_1$	$x_2$	$x_3$	$x_4$	$x_5$
<i>awful</i>	0	1	0	0	0	0
<i>poor</i>	1	0	1	0	0	0
<i>OK</i>	2	0	0	1	0	0
<i>good</i>	3	0	0	0	1	0
<i>great</i>	4	0	0	0	0	1

# Attribute Transformation

---

- An **attribute transform** is 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**



# Similarity and Dissimilarity Measures

---

- Similarity measure
  - Numerical measure of how alike two data objects are.
  - Is higher when objects are more alike.
  - Often falls in the range  $[0,1]$
- Dissimilarity measure
  - Numerical measure of how different two data objects are
  - Lower when objects are more alike
  - Minimum dissimilarity is often 0
  - Upper limit varies
- Proximity refers to a similarity or dissimilarity

# Similarity/Dissimilarity for Simple Attributes

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

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 $	$s = -d, s = \frac{1}{1+d}, s = e^{-d},$ $s = 1 - \frac{d - \min\_d}{\max\_d - \min\_d}$

# Euclidean Distance

---

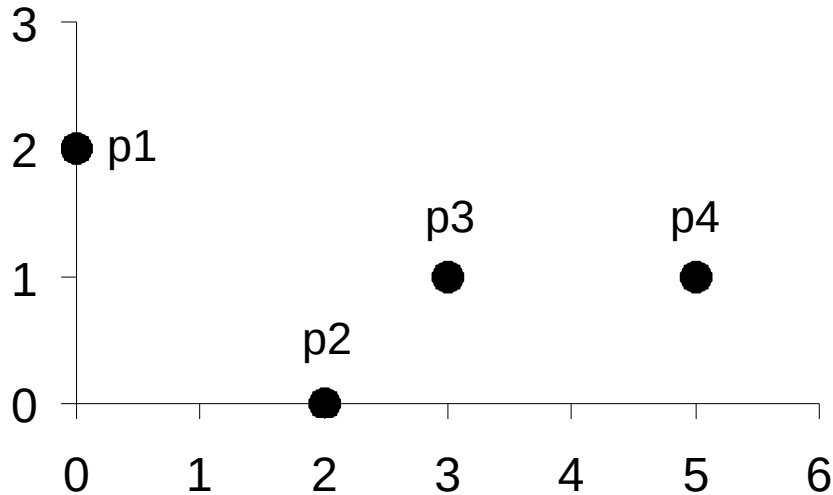
- Euclidean Distance

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

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

- 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 is a generalization of Euclidean Distance

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

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

# 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
- $r = \infty$ . “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.

# 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

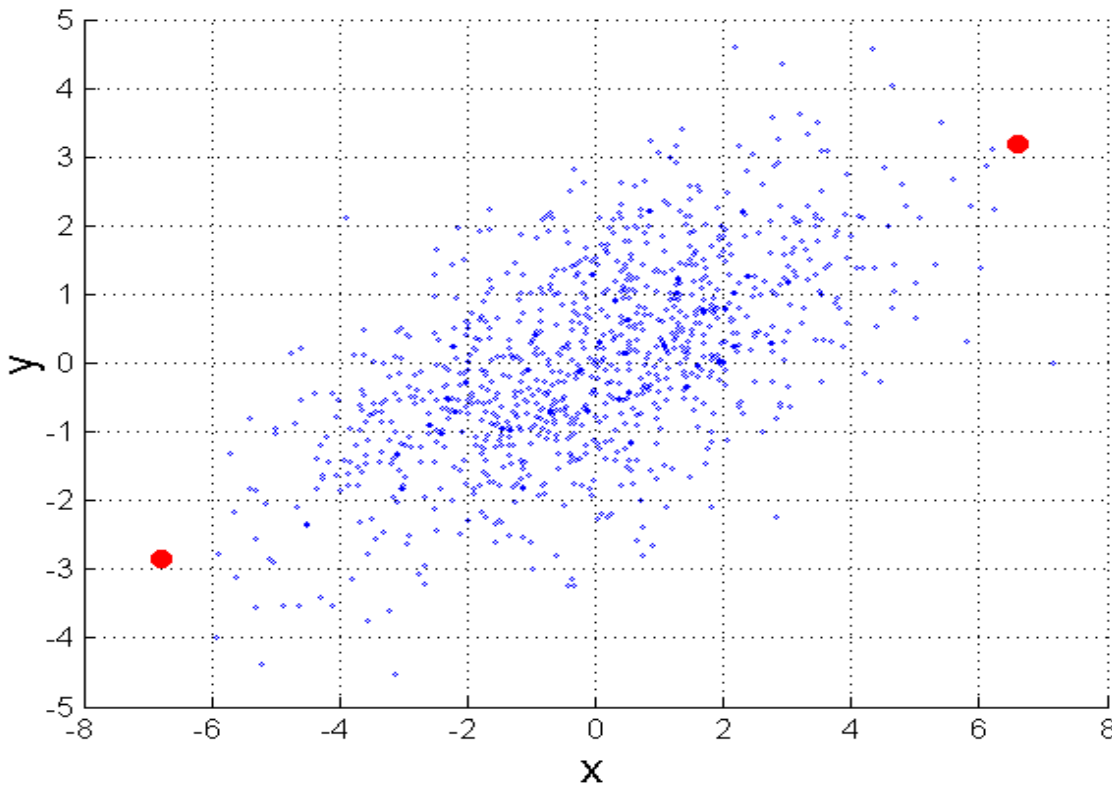
L2	p1	p2	p3	p4
p1	0	2.828	3.162	5.099
p2	2.828	0	1.414	3.162
p3	3.162	1.414	0	2
p4	5.099	3.162	2	0

$L_{\infty}$	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

# Mahalanobis Distance

-0.5

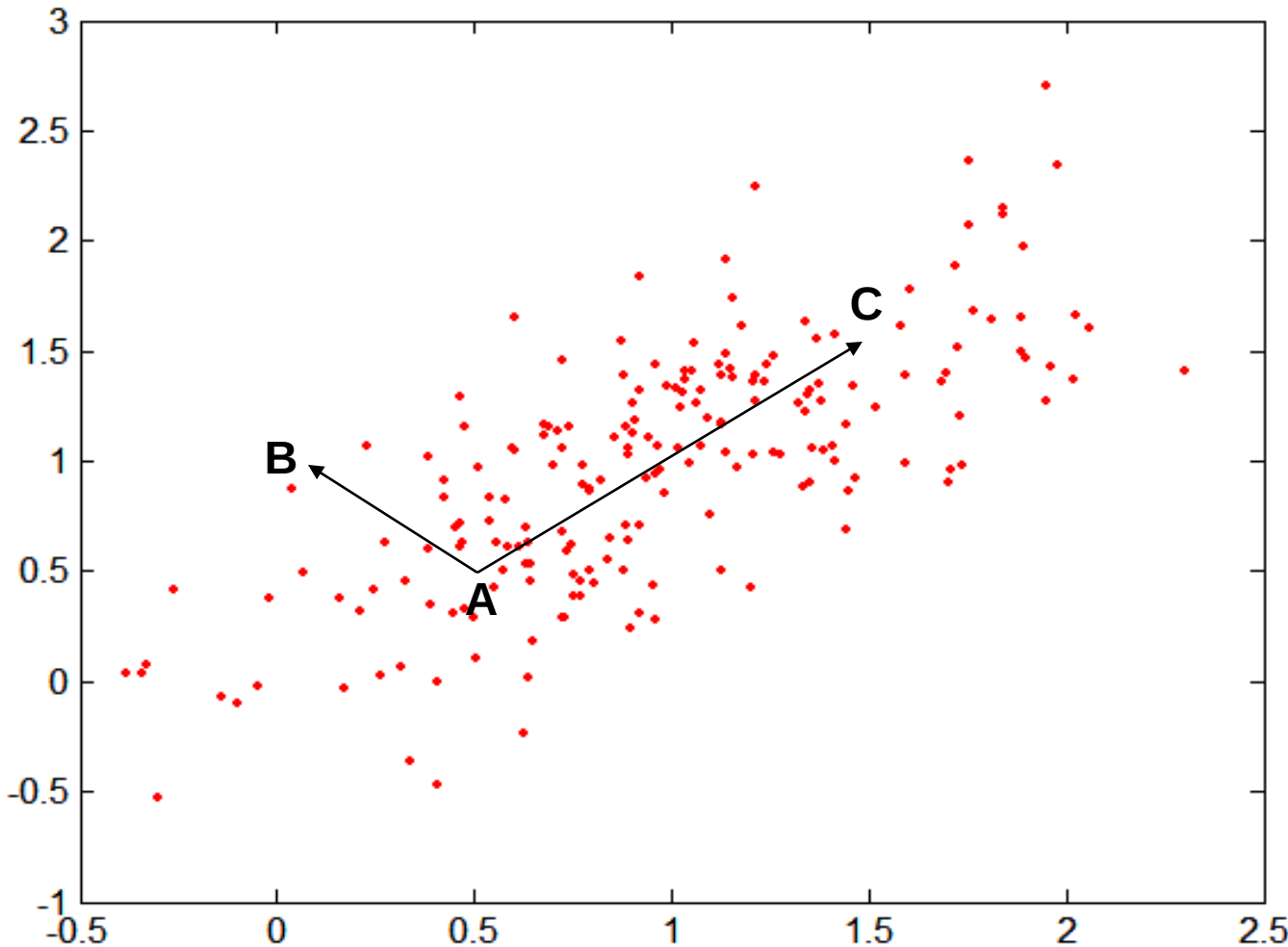


is the covariance matrix

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



# Mahalanobis Distance



**Covariance  
Matrix:**

$$\Sigma = \begin{bmatrix} 0.3 & 0.2 \\ 0.2 & 0.3 \end{bmatrix}$$

**A: (0.5, 0.5)**

**B: (0, 1)**

**C: (1.5, 1.5)**

**Mahal(A,B) = 5**

**Mahal(A,C) = 4**