# Data Mining: Data

Lecture Notes for Chapter 2

Introduction to Data Mining , 2$^{nd}$ Edition
by
Tan, Steinbach, Karpatne, Kumar

# Outline

- Attributes and Objects

- Types of Data

- Data Quality

- Similarity

- Data Preprocessing

# 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**

**Objects**

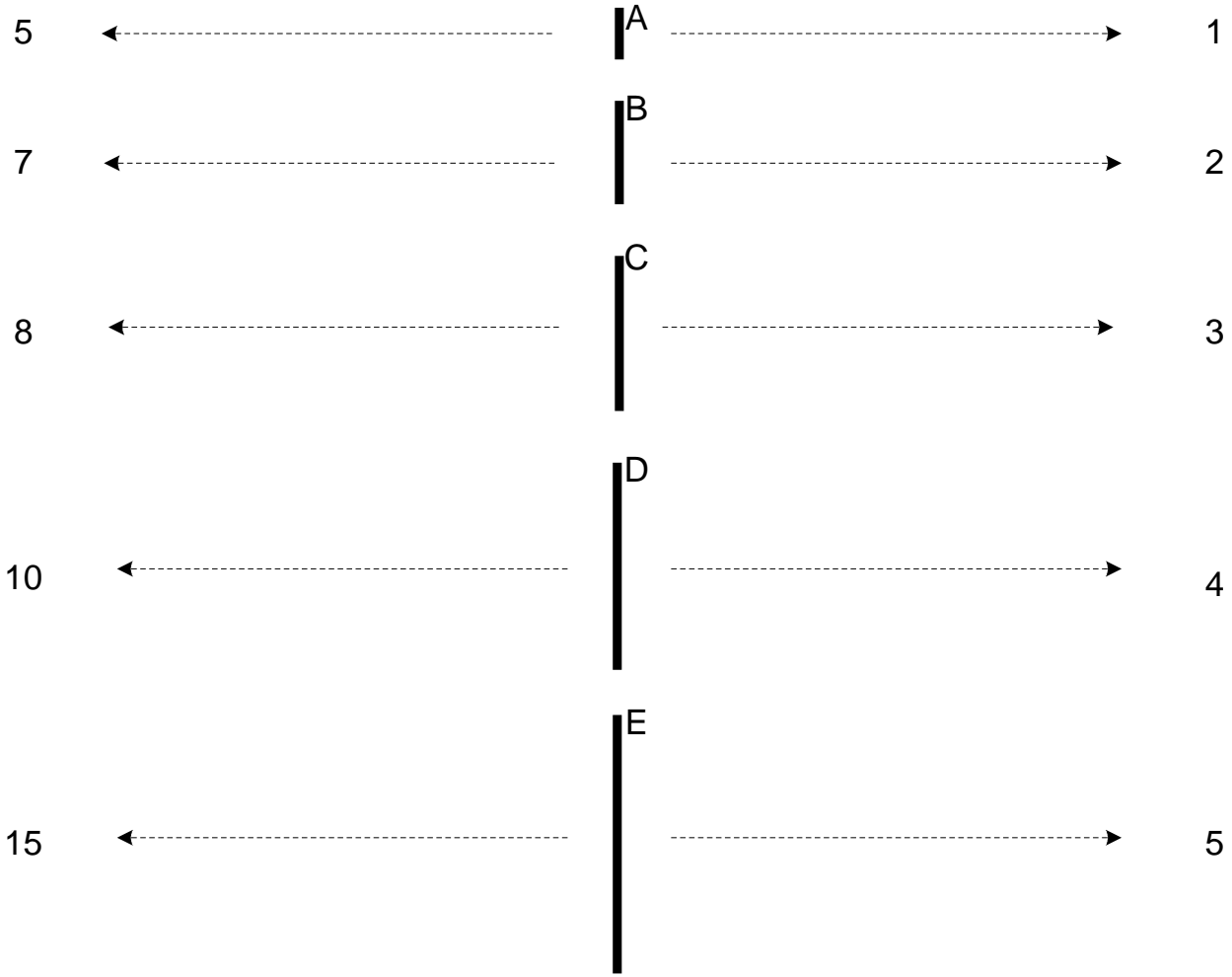| Tid | Refund | Marital Status | Taxable Income | Cheat |
|-----|--------|----------------|----------------|-------|
| 1 | Yes | Single | 125K | **No** |
| 2 | No | Married | 100K | **No** |
| 3 | No | Single | 70K | **No** |
| 4 | Yes | Married | 120K | **No** |
| 5 | No | Divorced | 95K | **Yes** |
| 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** |

# A More Complete View of Data

- Data may have parts

- The different parts of the data may have relationships

- More generally, data may have structure

- Data can be incomplete

- We will discuss this in more detail later

# 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 values can be different

# Measurement of Length

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

| 5 | ← | A | → | 1 |
| 7 | ← | B | → | 2 |

**This scale preserves only the ordering property of length.**

| 8 | ← | C | → | 3 |
| 10 | ← | D | → | 4 |

**This scale preserves the ordering and additvity properties of length.**

| 15 | ← | E | → | 5 |

# 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, time, counts

# Properties of Attribute Values

- The type of an attribute depends on which of the following properties/operations it possesses:
    - Distinctness:                $=  \neq$
    - Order:                          $<  >$
    - Differences are          $+  -$
      meaningful :
    - Ratios are                     $*  /$
      meaningful

    - Nominal attribute: distinctness
    - Ordinal attribute: distinctness & order
    - Interval attribute: distinctness, order & meaningful differences
    - Ratio attribute: all 4 properties/operations

# Difference Between Ratio and Interval

- Is it physically meaningful to say that a temperature of 10 ° is twice that of 5° on
  - the Celsius scale?
  - the Fahrenheit scale?
  - the Kelvin scale?

- Consider measuring the height above average
  - If Bill's height is three inches above average and Bob's height is six inches above average, then would we say that Bob is twice as tall as Bill?
  - Is this situation analogous to that of temperature?

**Introduction to Data Mining, 2nd Edition**

| | Attribute Type | Description | Examples | Operations |
|---|---|---|---|---|
| **Categorical Qualitative** | Nominal | Nominal attribute values only distinguish. (=, ≠) | zip codes, employee ID numbers, eye color, sex: {*male, female*} | mode, entropy, contingency correlation, $\chi 2$ test |
| | Ordinal | Ordinal attribute values also order objects. (<, >) | hardness of minerals, {*good, better, best*}, 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, *t* and *F* 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? |
| Categorical Qualitative | 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). |
| Numeric Quantitative | 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.

# Asymmetric Attributes

- Only presence (a non-zero attribute value) is regarded as important
  - Words present in documents
  - Items present in customer transactions

- If we met a friend in the grocery store would we ever say the following?

  *"I see our purchases are very similar since we didn't buy most of the same things."*

- We need two asymmetric binary attributes to represent one ordinary binary attribute
  - Association analysis uses asymmetric attributes

- Asymmetric attributes typically arise from objects that are sets

# Key Messages for Attribute Types

- The types of operations you choose should be "meaningful" for the type of data you have
  - Distinctness, order, meaningful intervals, and meaningful ratios are only four properties of data

  - The data type you see – often numbers or strings – may not capture all the properties or may suggest properties that are not there

  - Analysis may depend on these other properties of the data
    - Many statistical analyses depend only on the distribution

  - Many times what is meaningful is measured by statistical significance

  - But in the end, what is meaningful is measured by the domain

# 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

# 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

# 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 |

**Introduction to Data Mining, 2nd Edition**

# 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 |

# 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 |

**Introduction to Data Mining, 2nd Edition**

# 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 |

# Graph Data

- Examples: Generic graph, a molecule, and webpages

Benzene Molecule: C6H6

**Introduction to Data Mining, 2nd Edition**

# Ordered Data

● Sequences of transactions

**Items/Events**

$$( A\ B) \quad (D) \quad (C\ E)$$
$$( B\ D) \quad (C) \quad (E)$$
$$( C\ D) \quad (B) \quad (A\ E)$$

**An element of
the sequence**

**Introduction to Data Mining, 2nd Edition**

# Ordered Data

- Genomic sequence data

GGTTCCGCCTTCAGCCCCGCGCC
CGCAGGGCCCGCCCCGCGCCGTC
GAGAAGGGCCCGCCTGGCGGGCG
GGGGGAGGCGGGGCCGCCCGAGC
CCAACCGAGTCCGACCAGGTGCC
CCCTCTGCTCGGCCTAGACCTGA
GCTCATTAGGCGGCAGCGGACAG
GCCAAGTAGAACACGCGAAGCGC
TGGGCTGCCTGCTGCGACCAGGG

# Ordered Data

- Spatio-Temporal Data

Jan

**Average Monthly Temperature of land and ocean**

**Introduction to Data Mining, 2nd Edition**

# Data Quality

- Poor data quality negatively affects many data processing efforts

"The most important point is that poor data quality is an unfolding disaster.

- Poor data quality costs the typical company at least ten percent (10%) of revenue; twenty percent (20%) is probably a better estimate."

Thomas C. Redman, DM Review, August 2004

- 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 ...

- 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
  - Wrong data

# Noise

- For objects, noise is an extraneous object
- For attributes, noise refers to modification of original values
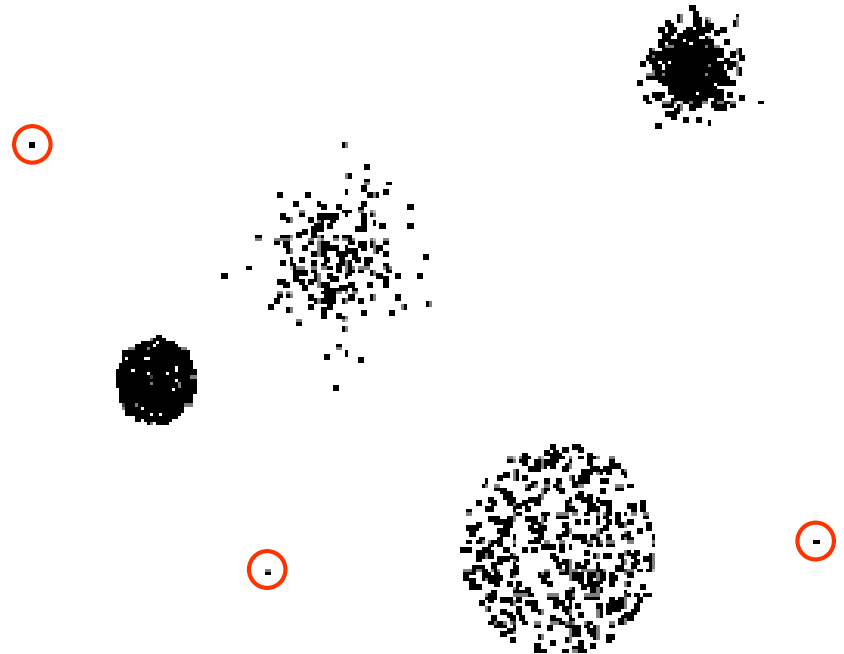  - Examples: distortion of a person's voice when talking on a poor phone and "snow" on television screen

**Two Sine Waves**

**Two Sine Waves + Noise**

**Introduction to Data Mining, 2nd Edition**

# Outliers

- *Outliers* are data objects with characteristics that are considerably different than most of the other data objects in the data set
  - **Case 1:** Outliers are noise that interferes with data analysis

  - **Case 2:** Outliers are the goal of our analysis
    - Credit card fraud
    - Intrusion detection

- Causes?

# 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

# Missing Values …

- Missing completely at random (MCAR)
  - Missingness of a value is independent of attributes
  - Fill in values based on the attribute
  - Analysis may be unbiased overall
- Missing at Random (MAR)
  - Missingness is related to other variables
  - Fill in values based other values
  - Almost always produces a bias in the analysis
- Missing Not at Random (MNAR)
  - Missingness is related to unobserved measurements
  - Informative or non-ignorable missingness
- Not possible to know the situation from the data

# 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

- When should duplicate data not be removed?

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

# Correlation measures the linear relationship between objects

$$\text{corr}(\mathbf{x}, \mathbf{y}) = \frac{\text{covariance}(\mathbf{x}, \mathbf{y})}{\text{standard\_deviation}(\mathbf{x}) * \text{standard\_deviation}(\mathbf{y})} = \frac{s_{xy}}{s_x \ s_y}, \quad (2.11)$$

where we are using the following standard statistical notation and definitions

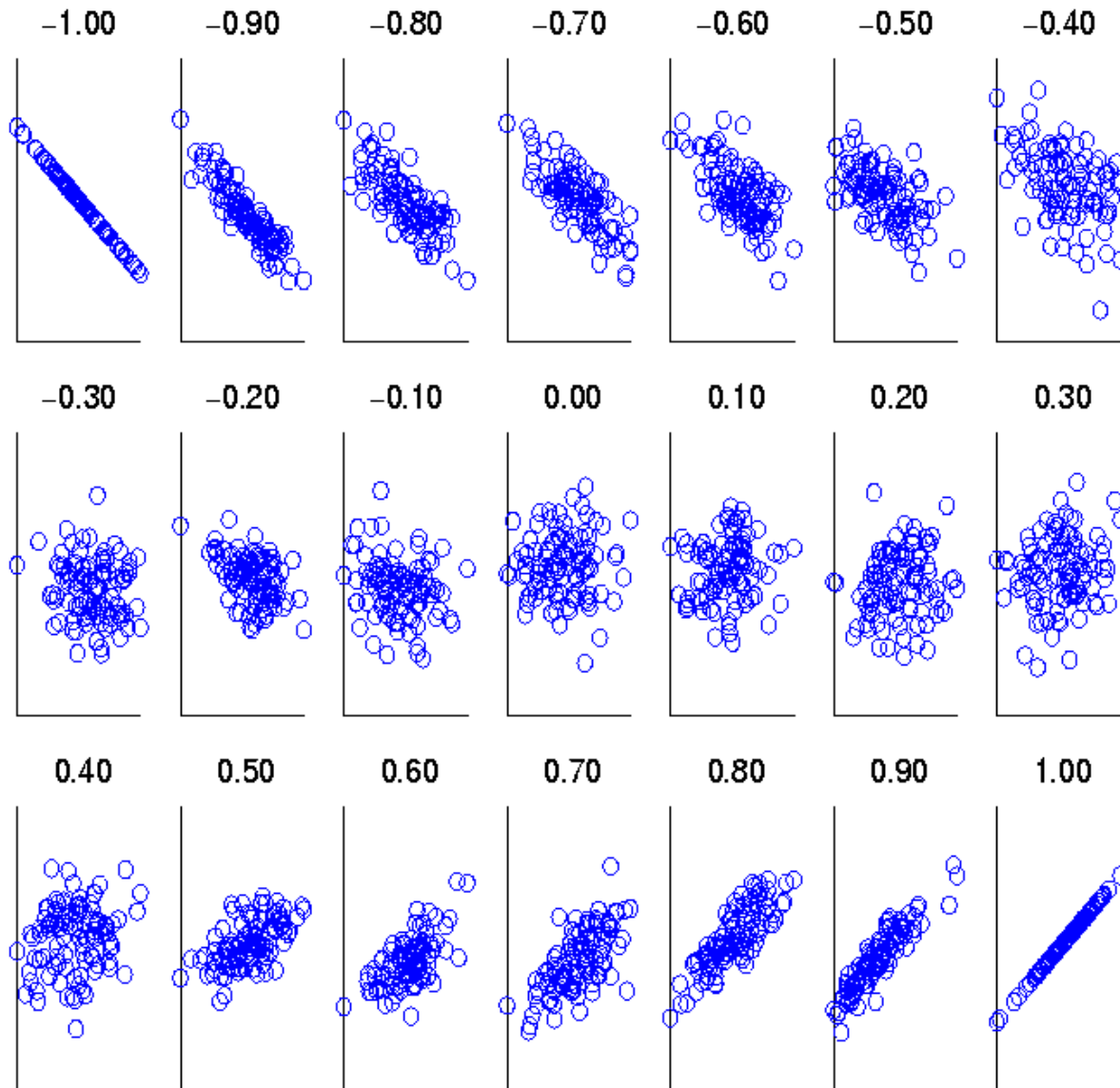$$\text{covariance}(\mathbf{x}, \mathbf{y}) = s_{xy} = \frac{1}{n-1} \sum_{k=1}^{n} (x_k - \overline{x})(y_k - \overline{y}) \qquad (2.12)$$

$$\text{standard\_deviation}(\mathbf{x}) = s_x = \sqrt{\frac{1}{n-1} \sum_{k=1}^{n} (x_k - \overline{x})^2}$$

$$\text{standard\_deviation}(\mathbf{y}) = s_y = \sqrt{\frac{1}{n-1} \sum_{k=1}^{n} (y_k - \overline{y})^2}$$

$$\overline{x} = \frac{1}{n} \sum_{k=1}^{n} x_k \text{ is the mean of } \mathbf{x}$$

$$\overline{y} = \frac{1}{n} \sum_{k=1}^{n} y_k \text{ is the mean of } \mathbf{y}$$

# Visually Evaluating Correlation



Scatter plots showing the similarity from –1 to 1.

# Drawback of Correlation

- $\mathbf{x} = (-3, -2, -1, 0, 1, 2, 3)$
- $\mathbf{y} = (9, 4, 1, 0, 1, 4, 9)$

$$y_i = x_i^2$$

- $\text{mean}(\mathbf{x}) = 0, \text{mean}(\mathbf{y}) = 4$
- $\text{std}(\mathbf{x}) = 2.16, \text{std}(\mathbf{y}) = 3.74$

- corr = (-3)(5)+(-2)(0)+(-1)(-3)+(0)(-4)+(1)(-3)+(2)(0)+3(5) / ( 6 * 2.16 * 3.74 )
  = 0

# Comparison of Proximity Measures

- Domain of application
  - Similarity measures tend to be specific to the type of attribute and data
  - Record data, images, graphs, sequences, 3D-protein structure, etc. tend to have different measures

- However, one can talk about various properties that you would like a proximity measure to have
  - Symmetry is a common one
  - Tolerance to noise and outliers is another
  - Ability to find more types of patterns?
  - Many others possible

- The measure must be applicable to the data and produce results that agree with domain knowledge

# Data Preprocessing

- Aggregation

- Sampling

- Dimensionality Reduction

- Feature subset selection

- Feature creation

- Discretization and Binarization

# 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

# 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**

**Introduction to Data Mining, 2nd Edition**

# 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

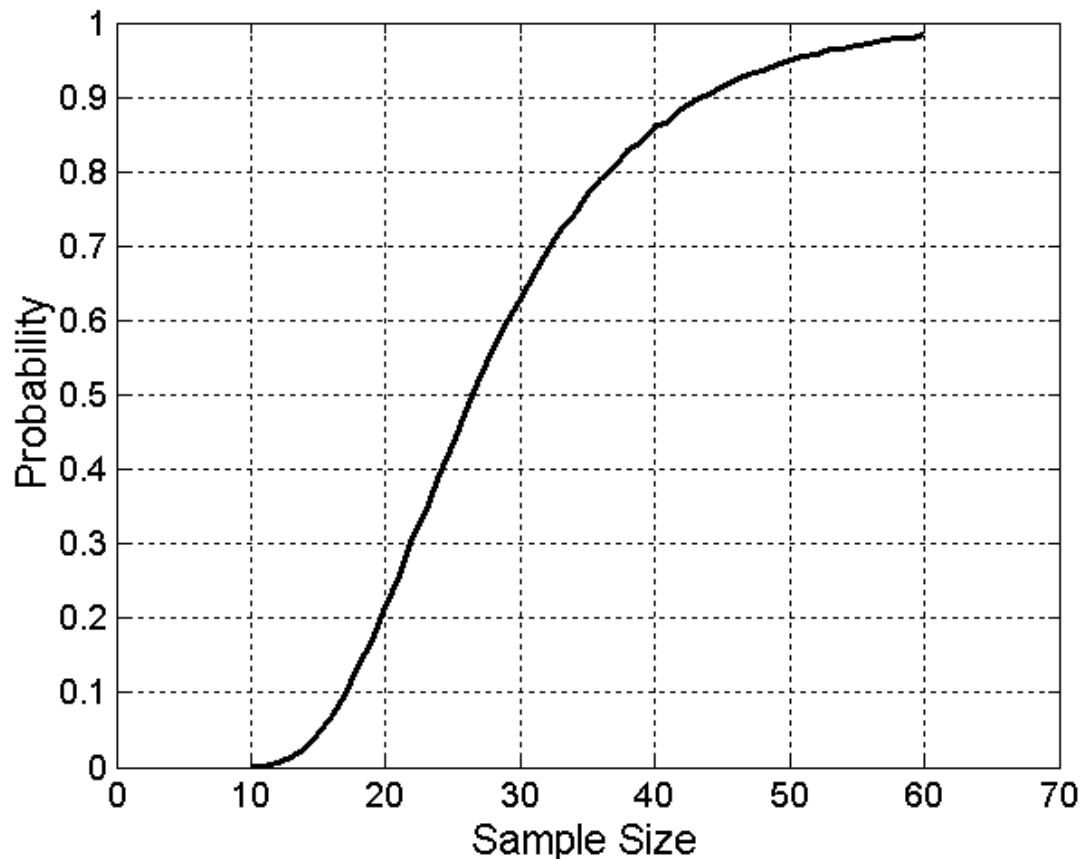# Sample Size



**8000 points**          **2000 Points**          **500 Points**

# 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

# 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
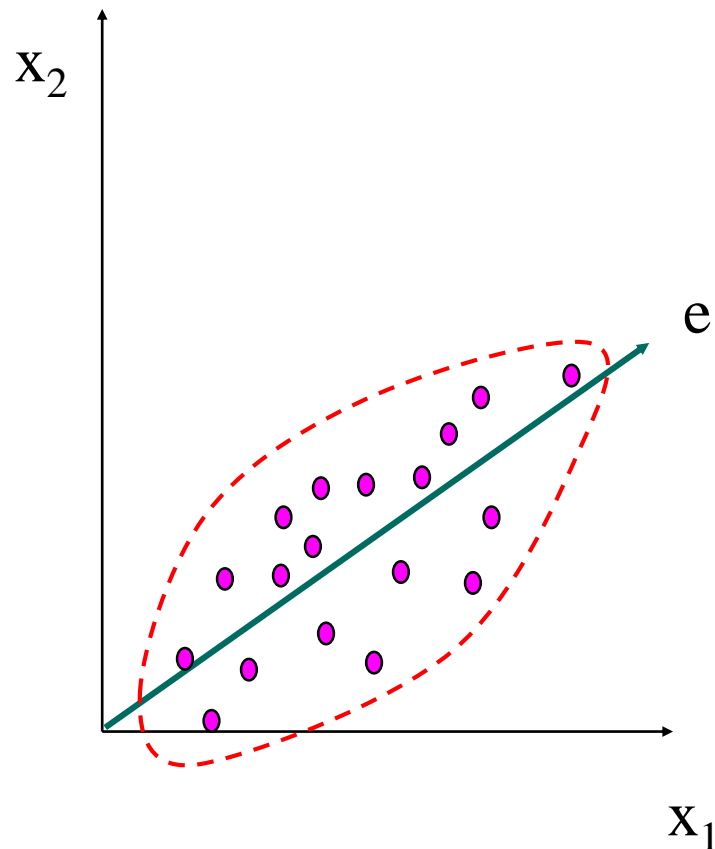


- **Randomly generate 500 points**

- **Compute difference between max and min distance between any pair of points**

# 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

- 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
- 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 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

# Discretization

- Discretization is the process of converting a continuous attribute into an ordinal attribute
    - A potentially infinite number of values are mapped into a small number of categories
    - Discretization is commonly used in classification
    - Many classification algorithms work best if both the independent and dependent variables have only a few values
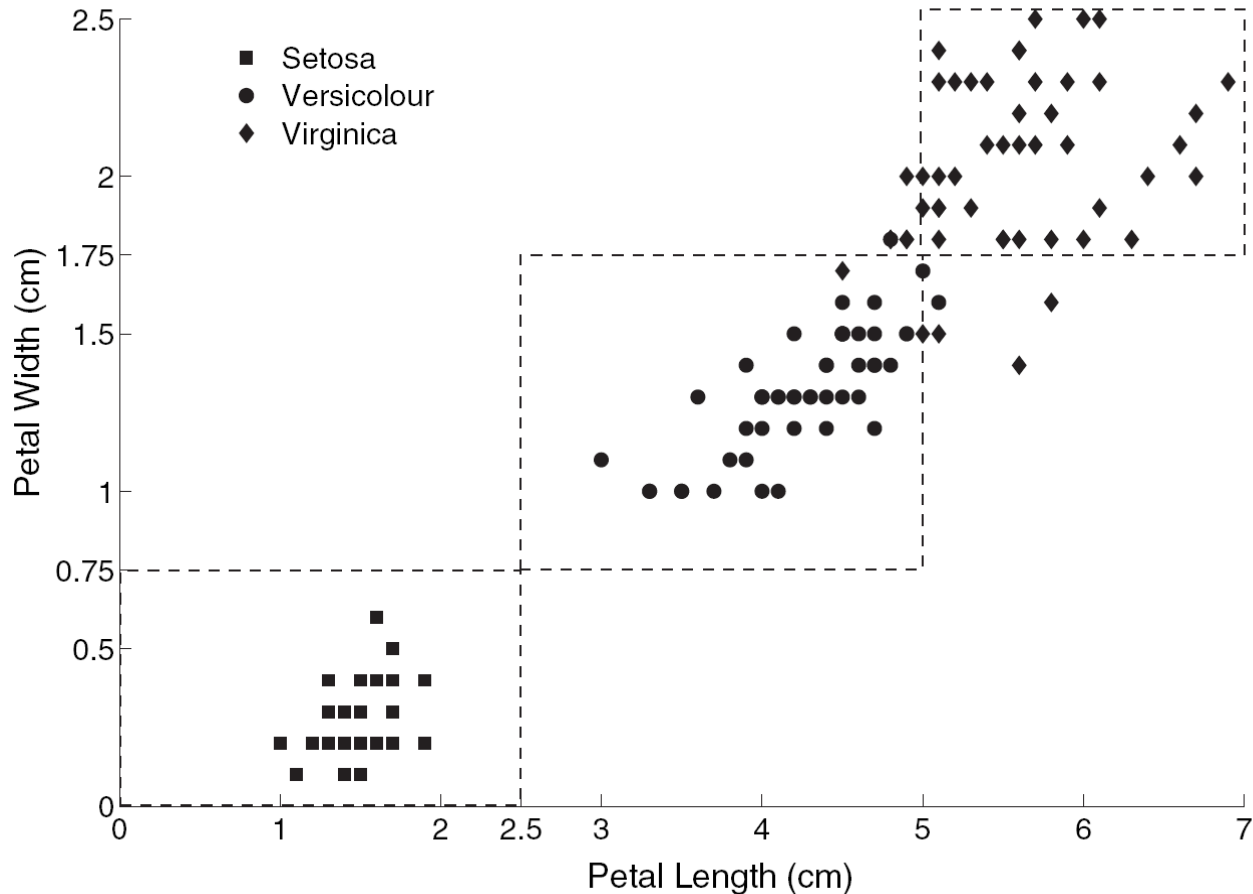    - We give an illustration of the usefulness of discretization using the Iris data set

# Iris Sample Data Set

- Iris Plant data set.
  - Can be obtained from the UCI Machine Learning Repository http://www.ics.uci.edu/~mlearn/MLRepository.html
  - From the statistician Douglas Fisher
  - Three flower types (classes):
    - Setosa
    - Versicolour
    - Virginica
  - Four (non-class) attributes
    - Sepal width and length
    - Petal width and length



Virginica. Robert H. Mohlenbrock. USDA NRCS. 1995. Northeast wetland flora: Field office guide to plant species. Northeast National Technical Center, Chester, PA. Courtesy of USDA NRCS Wetland Science Institute.

# Discretization: Iris Example



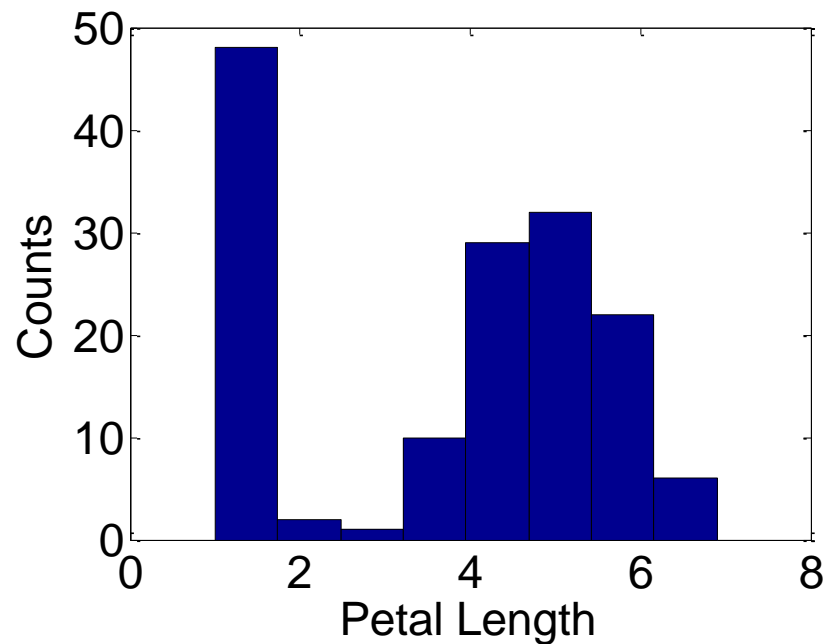Petal width low or petal length low implies Setosa.
Petal width medium or petal length medium implies Versicolour.
Petal width high or petal length high implies Virginica.
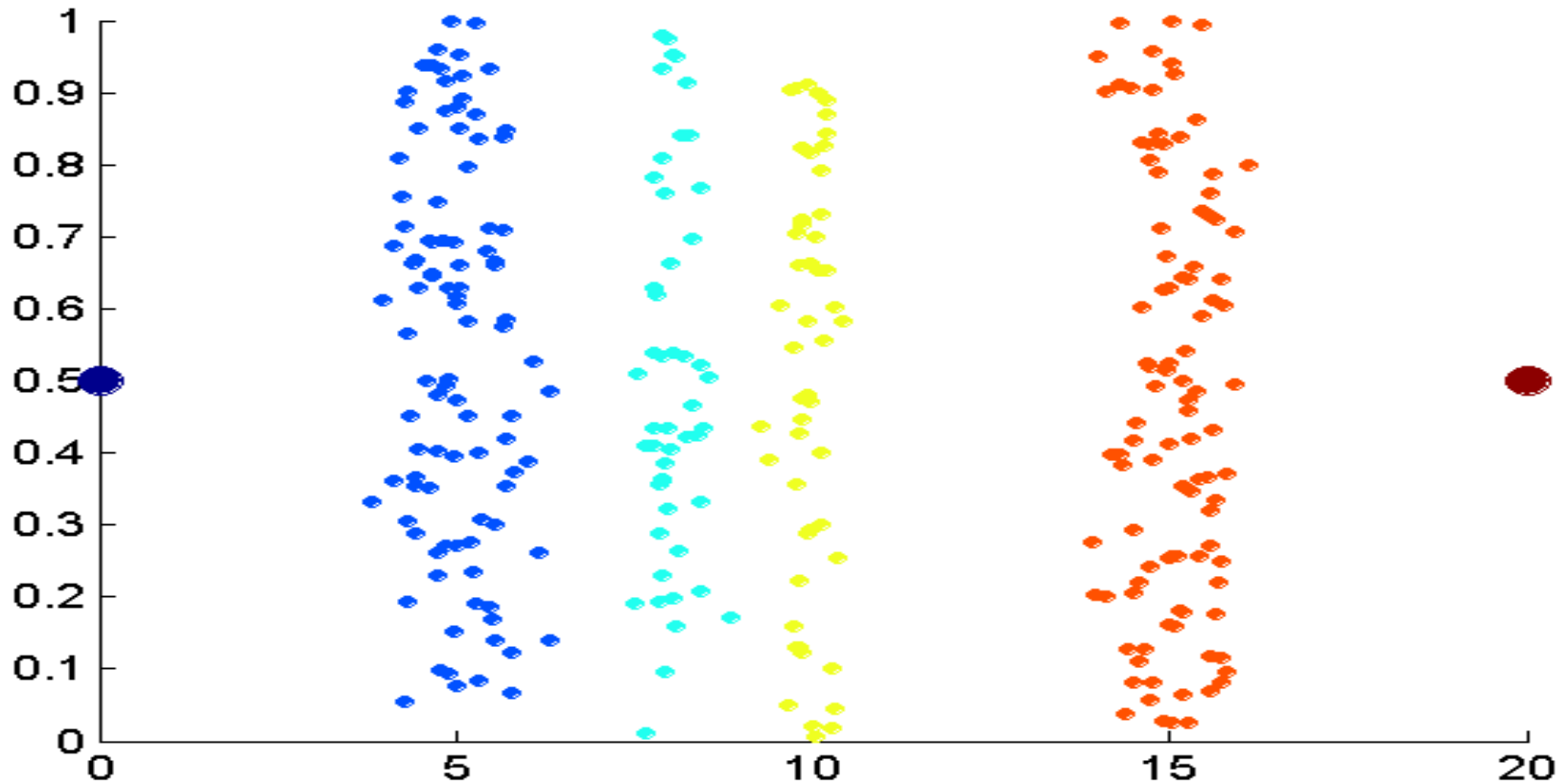
# Discretization: Iris Example ...

- How can we tell what the best discretization is?
  - Unsupervised discretization: find breaks in the data values
    - Example: Petal Length
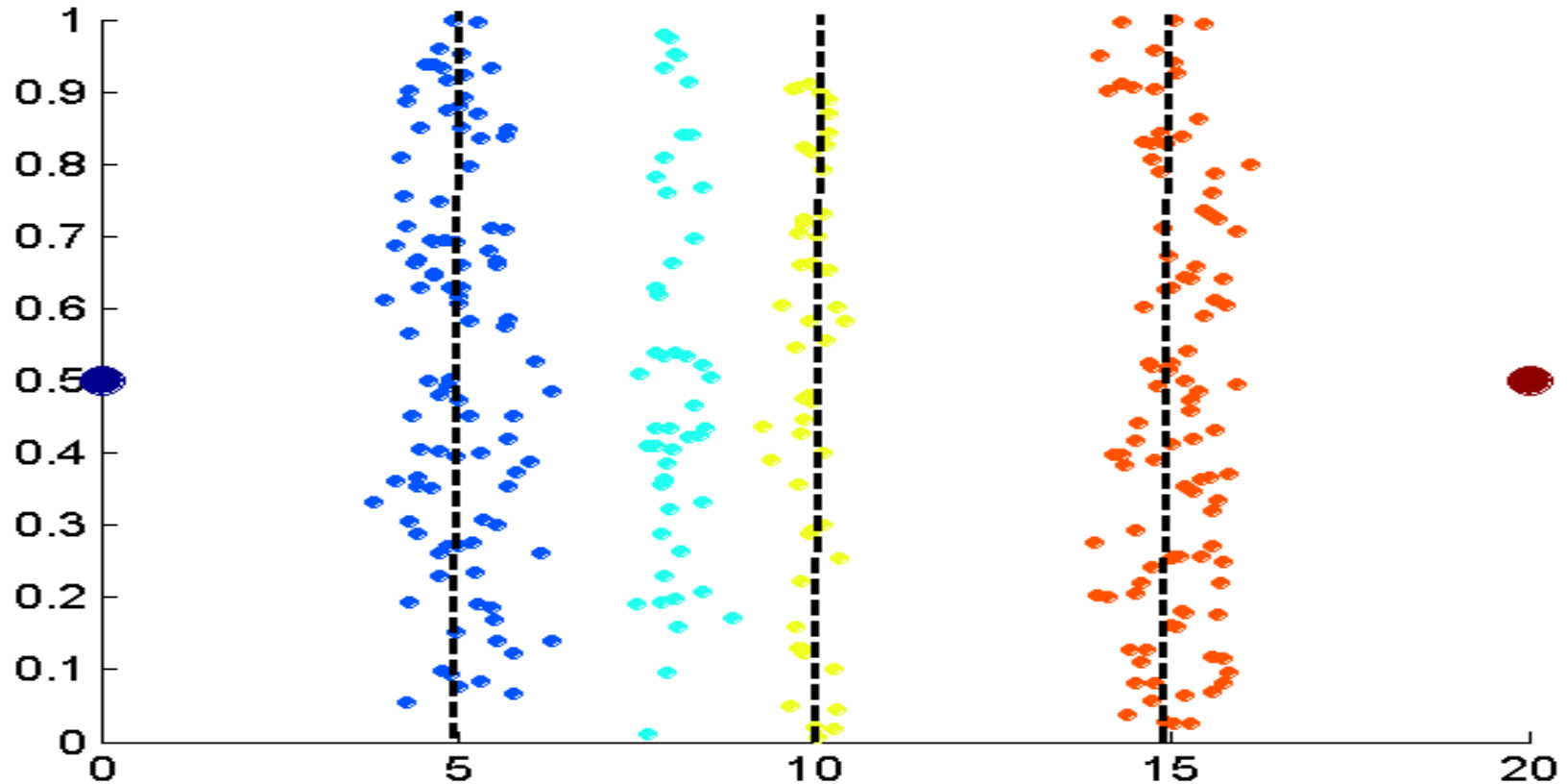


  - Supervised discretization: Use class labels to find breaks
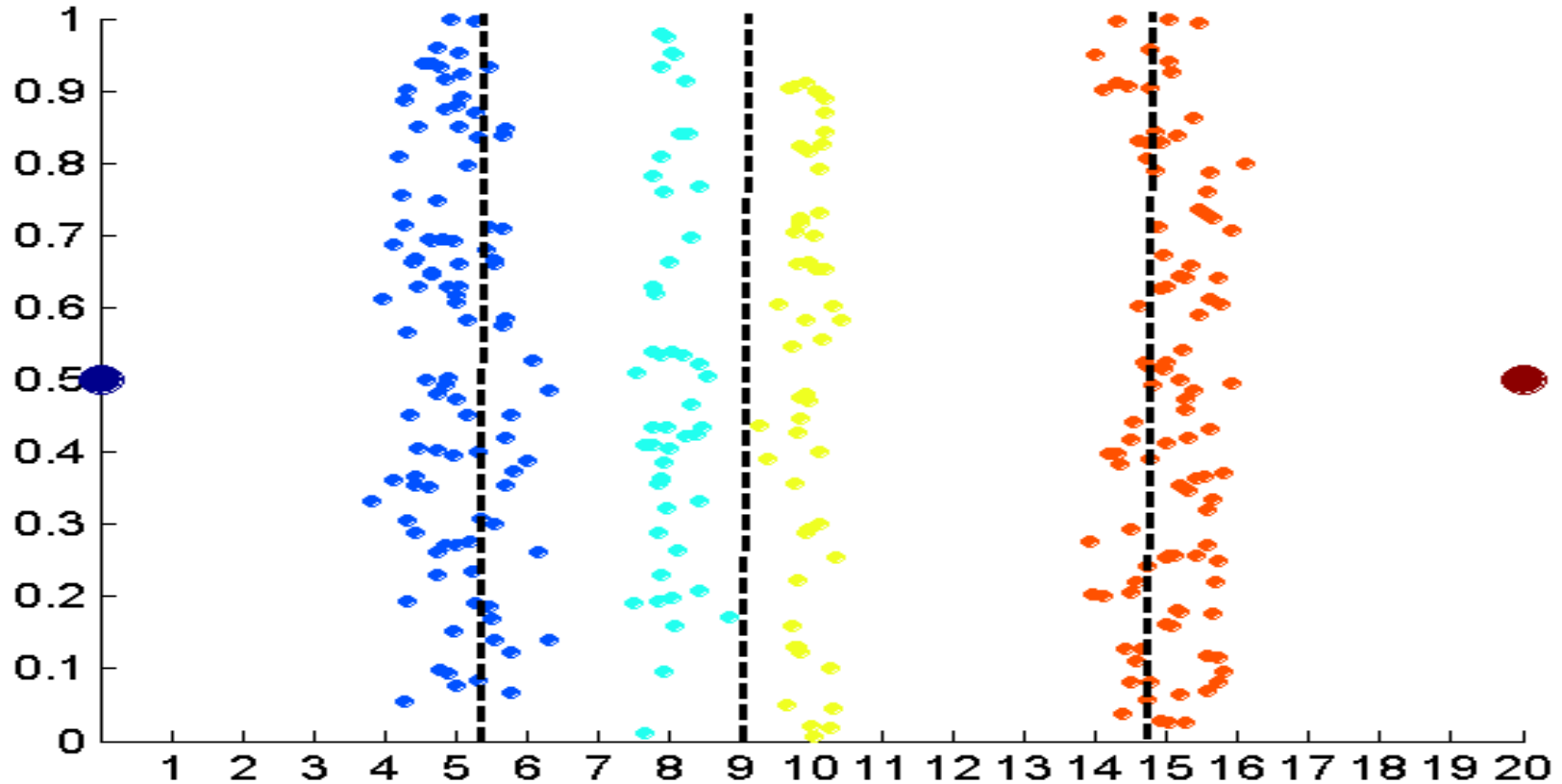
# Discretization Without Using Class Labels



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

**Introduction to Data Mining, 2nd Edition**
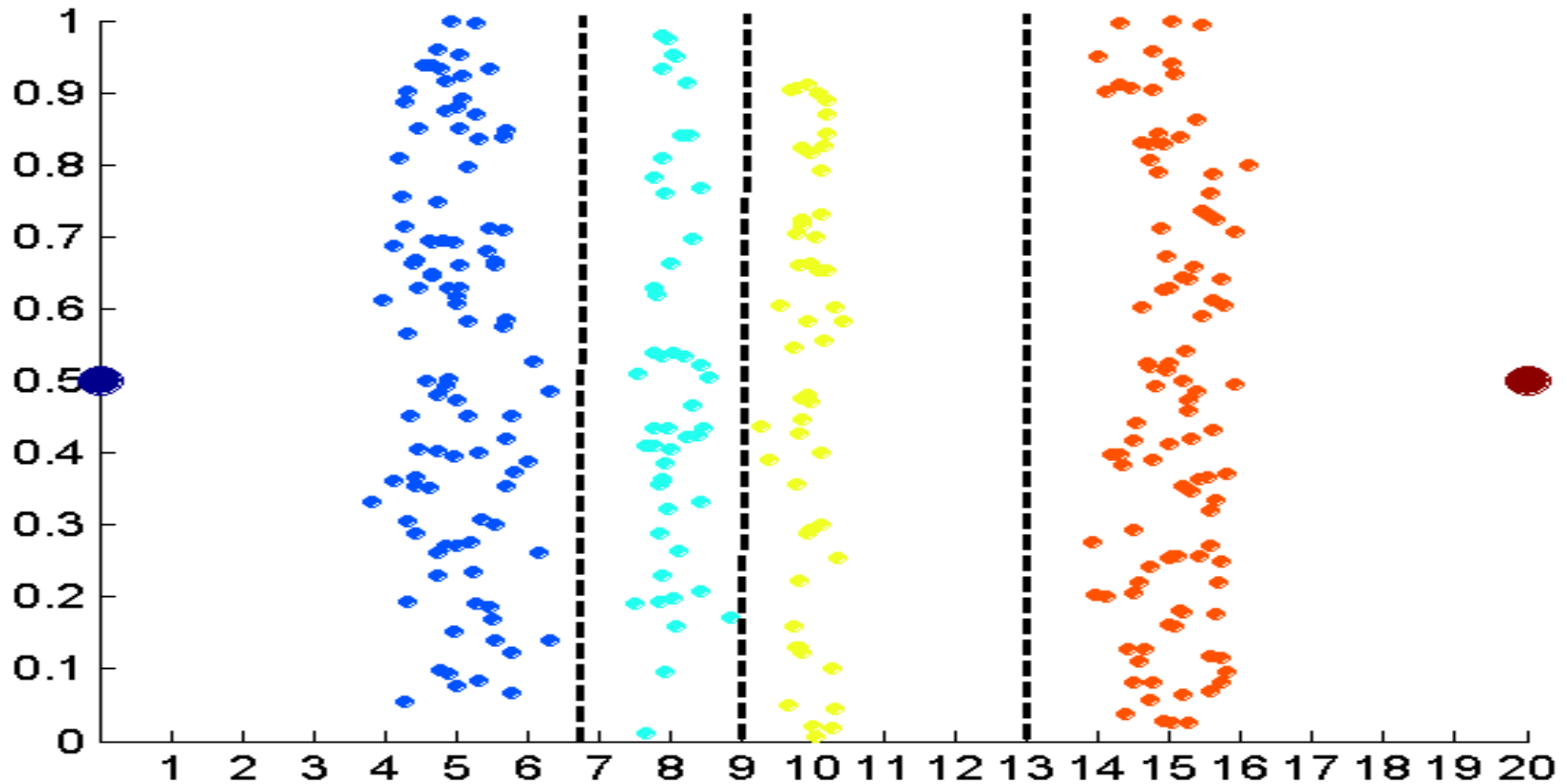
# Discretization Without Using Class Labels



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

**Introduction to Data Mining, 2nd Edition**

# Discretization Without Using Class Labels



**Equal frequency** approach used to obtain 4 values.

# Discretization Without Using Class Labels



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

# Binarization

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

- Typically used for association analysis

- Often convert a continuous attribute to a categorical attribute and then convert a categorical attribute to a set of binary attributes
    - Association analysis needs asymmetric binary attributes
    - Examples: eye color and height measured as {low, medium, high}