# EE 4146 Data Engineering and Learning Systems

## Lecture 2: Data Exploration

Semester A, 2021-2022

# Schedules

| Week | Date | Topics |
|------|------|--------|
| 1 | Sep. 1 | Introduction |
| 2 | Sep. 8 | Data exploration |
| 3 | Sep. 15 | Feature reduction and selection (HW1 out) |
| 4 | Sep. 22 | Mid-Autumn Festival |
| 5 | Sep. 29 | Clustering I: Kmeans based models (HW1 due in this weekend) |
| 6 | Oct. 6 | Clustering II: Hierarchical/density based/fuzzing clustering |
| 7 | Oct. 13 | Midterm (no tutorials this week) |
| 8 | Oct. 20 | Linear classifiers |
| 9 | Oct. 27 | Classification based on decision tree (Tutorial on project) (HW2 out) |
| 10 | Nov. 3 | Bayes based classifier (Tutorial on codes) (HW2 due in this weekend) |
| 11 | Nov. 10 | KNN and classifier ensemble |
| 12 | Nov. 17 | Deep learning based models (Quiz) |
| 13 | Nov. 24 | Summary |

# Outline

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

# 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

# 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)

# Outline

- Attributes and Objects

- <span style="color:purple">Types of Data</span>

- Data Quality

- Similarity and Distance

- Data Preprocessing

# 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

| 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 |

# 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.

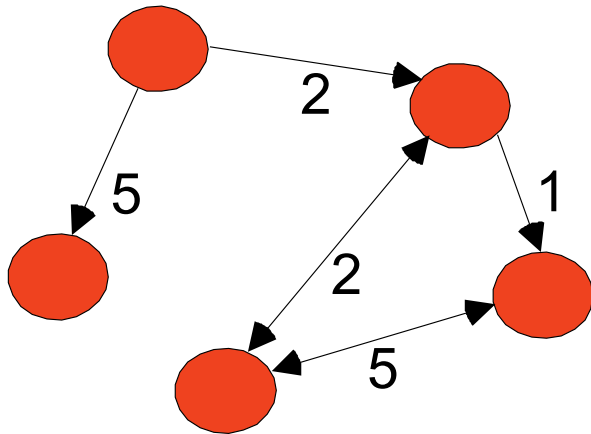|  | team | coach | play | ball | score | game | win | lost | timeout | season |
|---|---|---|---|---|---|---|---|---|---|---|
| Document 1 | 3 | 0 | 5 | 0 | 2 | 6 | 0 | 2 | 0 | 2 |
| Document 2 | 0 | 7 | 0 | 2 | 1 | 0 | 0 | 3 | 0 | 0 |
| Document 3 | 0 | 1 | 0 | 0 | 1 | 2 | 2 | 0 | 3 | 0 |

# 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

| 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



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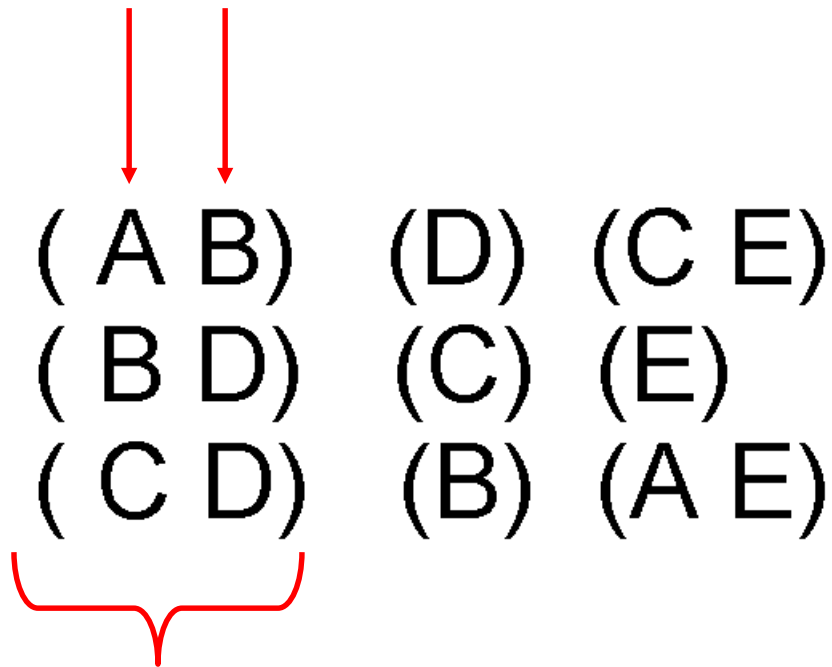
Christopher Matheus, Philip Chan, and Gregory Piatetsky-Shapiro, "Systems for knowledge Discovery in databases", IEEE Transactions on Knowledge and Data Engineering, 5(6):903-913, December 1993.

# Ordered Data

- Sequences of transactions

**Items/Events**

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

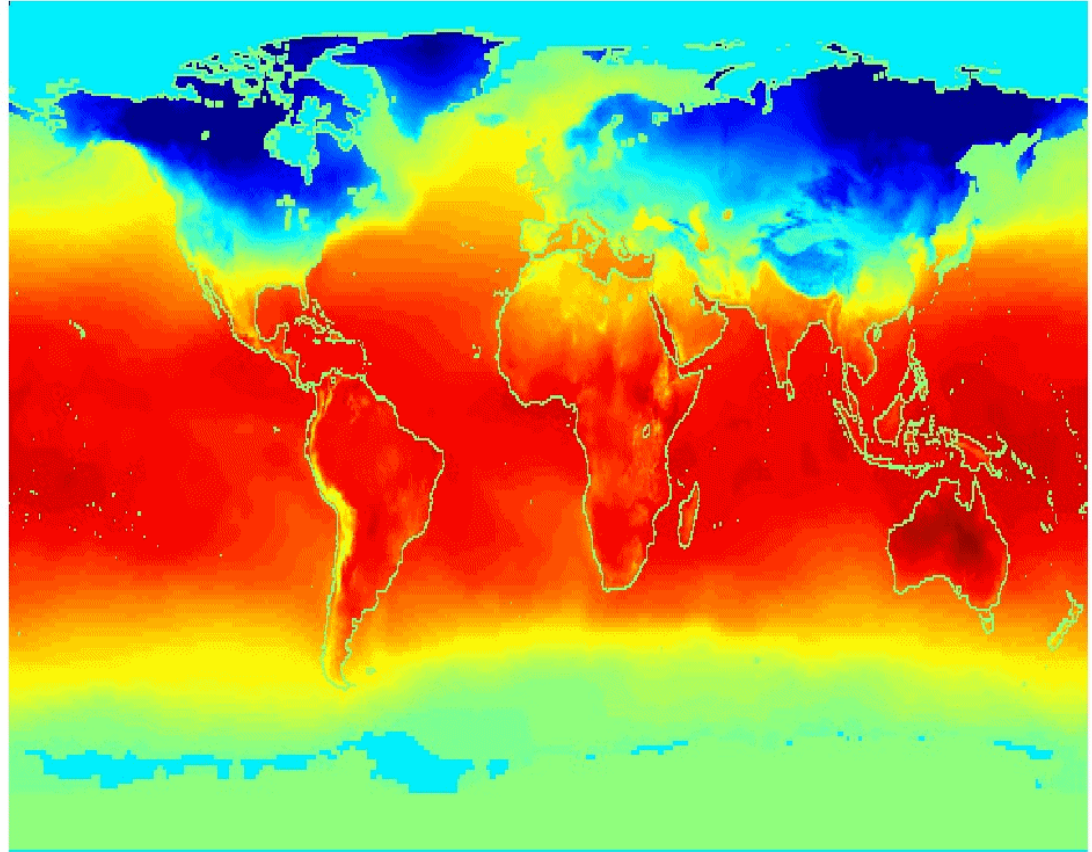**An element of the sequence**

# Ordered Data

- Genomic sequence data

GGTTCCGCCTTCAGCCCCGCGCC
CGCAGGGCCCGCCCCGCGCCGTC
GAGAAGGGCCCGCCTGGCGGGCG
GGGGGAGGCGGGGCCGCCCGAGC
CCAACCGAGTCCGACCAGGTGCC
CCCTCTGCTCGGCCTAGACCTGA
GCTCATTAGGCGGCAGCGGACAG
GCCAAGTAGAACACGCGAAGCGC
TGGGCTGCCTGCTGCGACCAGGG

# Ordered Data

- Spatio-Temporal Data

**Average Monthly Temperature of land and ocean**

Jan

# Converting to Numerical Features

- Often want a real-valued example representation:
- This is called a "1 of k" encoding.
- We can now interpret examples as points in space:
  - E.g., first example is at (23,1,0,0,22000).

| Age | City | Income |
|-----|------|--------|
| 23 | Van | 22,000.00 |
| 23 | Bur | 21,000.00 |
| 22 | Van | 0.00 |
| 25 | Sur | 57,000.00 |
| 19 | Bur | 13,500.00 |
| 22 | Van | 20,000.00 |

| Age | Van | Bur | Sur | Income |
|-----|-----|-----|-----|--------|
| 23 | 1 | 0 | 0 | 22,000.00 |
| 23 | 0 | 1 | 0 | 21,000.00 |
| 22 | 1 | 0 | 0 | 0.00 |
| 25 | 0 | 0 | 1 | 57,000.00 |
| 19 | 0 | 1 | 0 | 13,500.00 |
| 22 | 1 | 0 | 0 | 20,000.00 |

# Approximating Text with Numerical Features

- Bag of words replaces document by word counts:
- Ignores order, but often captures general theme.
- You can compute a "distance" between documents.

The **International Conference on Machine Learning** (ICML) is the leading international <u>academic conference</u> in <u>machine learning</u>

| ICML | International | Conference | Machine | Learning | Leading | Academic |
|------|--------------|------------|---------|----------|---------|----------|
| 1 | 2 | 2 | 2 | 2 | 1 | 1 |

# Approximating Images and Graphs

■ We can think of other data types in this way:

– Images:



graycale intensity →

| (1,1) | (2,1) | (3,1) | ... | (m,1) | ... | (m,n) |
|-------|-------|-------|-----|-------|-----|-------|
| 45    | 44    | 43    | ... | 12    | ... | 35    |

– Graphs:



adjacency matrix →

| N1 | N2 | N3 | N4 | N5 | N6 | N7 |
|----|----|----|----|----|----|----|
| 0  | 1  | 1  | 1  | 1  | 1  | 1  |
| 0  | 0  | 0  | 1  | 0  | 1  | 0  |
| 0  | 0  | 0  | 0  | 0  | 1  | 0  |
| 0  | 0  | 0  | 0  | 0  | 0  | 0  |

# Outline

- Attributes and Objects

- Types of Data

- Data Quality

- Similarity and Distance
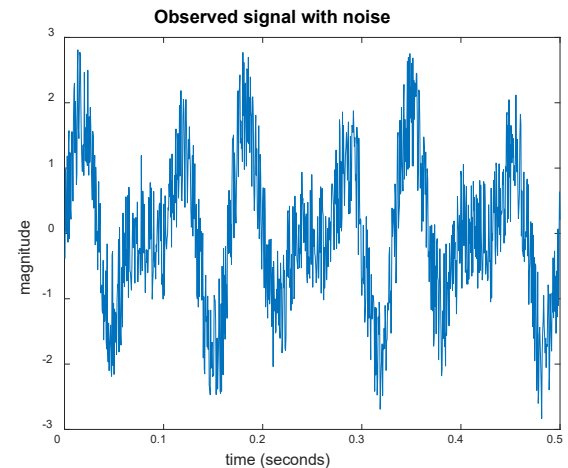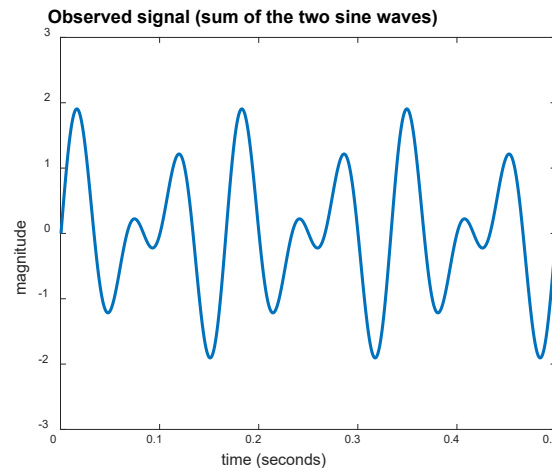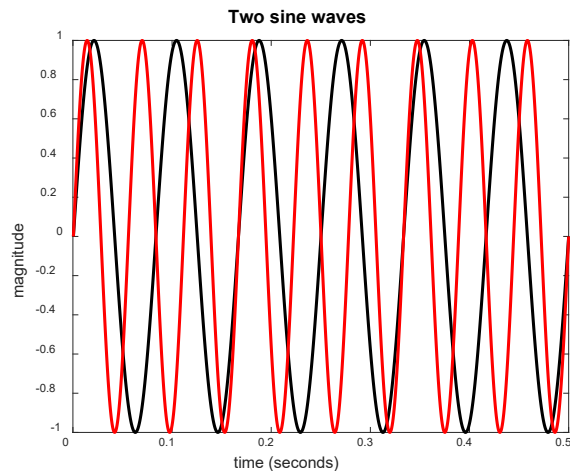
- Data Preprocessing

# Data Quality

- Poor data quality negatively affects many data processing efforts

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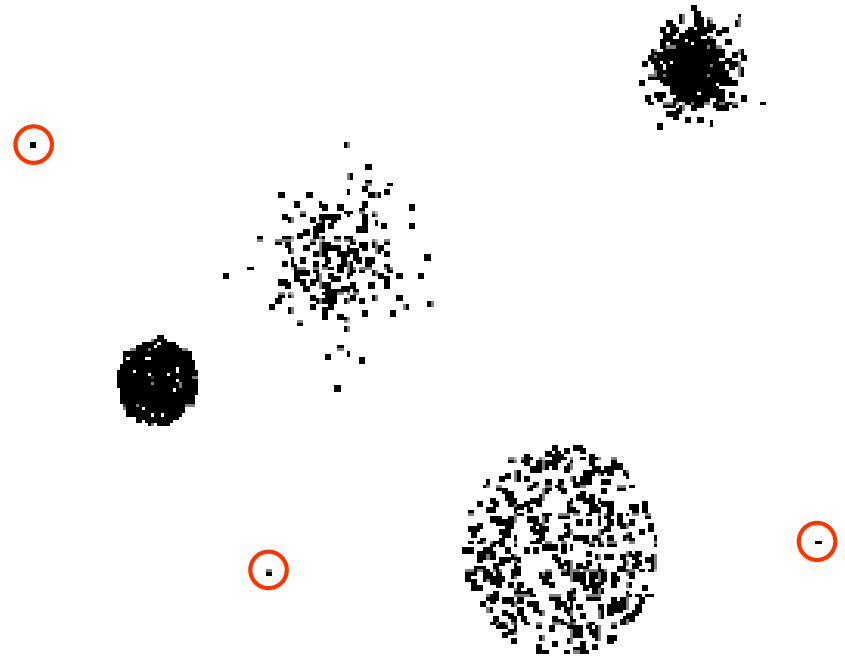  - Missing values
  - Duplicate data

# 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.



Two sine waves

Observed signal (sum of the two sine waves)

Observed signal with noise

# 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

# 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 / census results
  - Ignore the missing value during analysis

# 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?

# Outline

- Attributes and Objects

- Types of Data

- Data Quality

- Similarity and Distance

- Data Preprocessing

# 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

- 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.

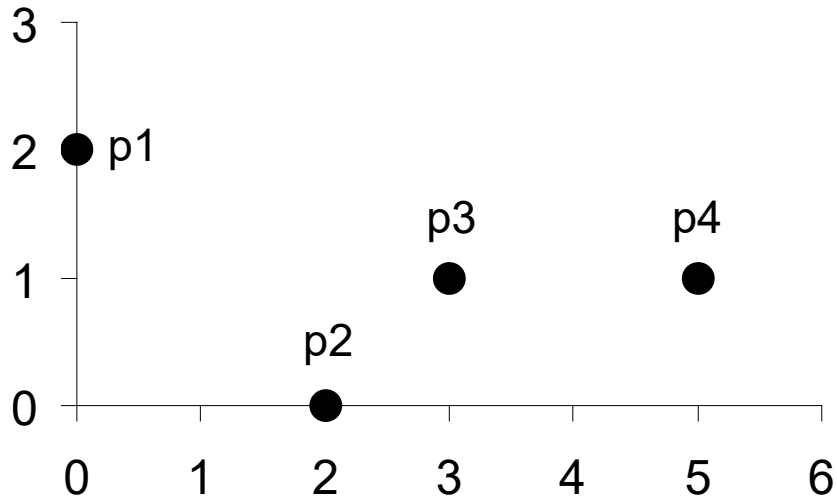| Attribute Type | Dissimilarity | Similarity |
|---|---|---|
| Nominal | $d = \begin{cases} 0 & \text{if } x = y \\ 1 & \text{if } x \neq y \end{cases}$ | $s = \begin{cases} 1 & \text{if } x = y \\ 0 & \text{if } x \neq y \end{cases}$ |
| Ordinal | $d = \lvert x - y \rvert/(n-1)$ (values mapped to integers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

# 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) of 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) of data objects $\boldsymbol{x}$ and $\boldsymbol{y}$.

# 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 <span style="color:red">number of bits that are different between two binary vectors</span>

- *r* = 2.  Euclidean distance

- $r \rightarrow \infty$.  "supremum" ($L_{max}$ norm, $L_\infty$ norm) distance.
  - This is the <span style="color:red">maximum difference between any component of the vectors</span>

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

# Minkowski Distance: examples

| 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**

# Common Properties of a Distance

- Distances, such as the Euclidean distance, have some well known properties.

  *1.* $d(\mathbf{x}, \mathbf{y}) \geq 0$ for all $x$ and $y$ and $d(\mathbf{x}, \mathbf{y}) = 0$ if and only if $\mathbf{x} = \mathbf{y}$.

  *2.* $d(\mathbf{x}, \mathbf{y}) = d(\mathbf{y}, \mathbf{x})$ for all $\mathbf{x}$ and $\mathbf{y}$. (Symmetry)

  *3.* $d(\mathbf{x}, \mathbf{z}) \leq d(\mathbf{x}, \mathbf{y}) + d(\mathbf{y}, \mathbf{z})$ for all points $\mathbf{x}$, $\mathbf{y}$, and $\mathbf{z}$. (Triangle Inequality)

  where $d(\mathbf{x}, \mathbf{y})$ is the distance (dissimilarity) between points (data objects), $\mathbf{x}$ and $\mathbf{y}$.

- A distance that satisfies these properties is a <span style="color:red">metric</span>

# Common Properties of a Similarity

■ Similarities, also have some well known properties.

1.   $s(\mathbf{x}, \mathbf{y}) = 1$ (or maximum similarity) only if $\mathbf{x} = \mathbf{y}$.
     (does not always hold, e.g., cosine)

2.   $s(\mathbf{x}, \mathbf{y}) = s(\mathbf{y}, \mathbf{x})$   for all $\mathbf{x}$ and $\mathbf{y}$. (Symmetry)

where $s(\mathbf{x}, \mathbf{y})$ is the similarity between points (data objects), $\mathbf{x}$ and $\mathbf{y}$.

# Similarity Between Binary Vectors

- Common situation is that objects, **x** and **y**, have only binary attributes

- Compute similarities using the following quantities
$f_{01}$ = the number of attributes where **x** was 0 and **y** was 1
$f_{10}$ = the number of attributes where **x** was 1 and **y** was 0
$f_{00}$ = the number of attributes where **x** was 0 and **y** was 0
$f_{11}$ = the number of attributes where **x** was 1 and **y** was 1

- Simple <span style="color:red">Matching</span> and <span style="color:red">Jaccard</span> Coefficients
SMC    =  number of matches / number of attributes
$$= (f_{11} + f_{00}) / (f_{01} + f_{10} + f_{11} + f_{00})$$

J    = number of 11 matches / number of non-zero attributes
$$= (f_{11}) / (f_{01} + f_{10} + f_{11})$$

# SMC versus Jaccard: Example1

$\mathbf{x} = $ 1 0 0 0 0 0 0 0 0 0

$\mathbf{y} = $ 0 0 0 0 0 0 1 0 0 1

$f_{01} = 2$   (the number of attributes where $\mathbf{x}$ was 0 and $\mathbf{y}$ was 1)

$f_{10} = 1$   (the number of attributes where $\mathbf{x}$ was 1 and $\mathbf{y}$ was 0)

$f_{00} = 7$   (the number of attributes where $\mathbf{x}$ was 0 and $\mathbf{y}$ was 0)

$f_{11} = 0$   (the number of attributes where $\mathbf{x}$ was 1 and $\mathbf{y}$ was 1)

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

$$\text{J} = (f_{11}) \, / \, (f_{01} + f_{10} + f_{11}) = 0 \, / \, (2 + 1 + 0) = 0$$

# Cosine Similarity

■ If $d_1$ and $d_2$ are two document vectors, then

$$\cos(d_1, d_2) = \langle d_1, d_2 \rangle / ||d_1||\ ||d_2||,$$

where $\langle d_1, d_2 \rangle$ indicates inner product or vector dot product of vectors, $d_1$ and $d_2$, and $||\ d\ ||$ is the length of vector d.

■ Example:

$$d_1 = 3\ 2\ 0\ 5\ 0\ 0\ 0\ 2\ 0\ 0$$
$$d_2 = 1\ 0\ 0\ 0\ 0\ 0\ 0\ 1\ 0\ 2$$

$\langle d_1, d2 \rangle = 3*1 + 2*0 + 0*0 + 5*0 + 0*0 + 0*0 + 0*0 + 2*1 + 0*0 + 0*2 = 5$

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

$||\ d_2\ || = (1*1+0*0+0*0+0*0+0*0+0*0+0*0+1*1+0*0+2*2)^{0.5} = (6)^{0.5} = 2.449$

$\cos(d_1, d_2) = 0.3150$

# 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

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

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

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

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

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

# Drawback of Correlation

- **x** = (-3, -2, -1, 0, 1, 2, 3)
- **y** = (9, 4, 1, 0, 1, 4, 9)

$y_i = x_i^2$

- mean(**x**) = 0, mean(**y**) = 4
- std(**x**) = 2.16, std(**y**) = 3.74



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

# Correlation vs Cosine vs Euclidean Distance

- Consider the example
  - $\mathbf{x}$ = (1, 2, 4, 3, 0, 0, 0), $\mathbf{y}$ = (1, 2, 3, 4, 0, 0, 0)
  - $\mathbf{y_s}$ = $\mathbf{y}$ * 2 (scaled version of y),  $\mathbf{y_t}$ = $\mathbf{y}$ + 5 (translated version)

| Measure | $(x , y)$ | $(x , y_s)$ | $(x , y_t)$ |
|---|---|---|---|
| Cosine | 0.9667 | 0.9667 | 0.7940 |
| Correlation | 0.9429 | 0.9429 | 0.9429 |
| Euclidean Distance | 1.4142 | 5.8310 | 14.2127 |

- Compare the three proximity measures according to their behavior under variable transformation
  - scaling: multiplication by a value
  - translation: adding a constant

| Property | Cosine | Correlation | Euclidean Distance |
|---|---|---|---|
| Invariant to scaling (multiplication) | Yes | Yes | No |
| Invariant to translation (addition) | No | Yes | No |

43

# Correlation vs cosine vs Euclidean distance

- Choice of the right proximity measure depends on the domain
- What is the correct choice of proximity measure for the following situations?
  - Comparing documents using the frequencies of words
    - Documents are considered similar if the word frequencies are similar
    - (Test is good) vs (Test test is is good good)?-> want to find a criteria not sensitive to scale.
  - Comparing the temperature in Celsius of two locations
    - Two locations are considered similar if the temperatures are similar in magnitude
    - ->just want to calculate the direct similarity
  - Comparing two time series of temperature measured in Celsius
    - Two time series are considered similar if their "shape" is similar, i.e., they vary in the same way over time, achieving minimums and maximums at similar times, etc.
    - -> want to find a criteria not sensitive to translation

# 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

# Information Based Measures

- Information theory is a well-developed and fundamental disciple with broad applications

- Some similarity measures are based on information theory
  - Mutual information in various versions
  - Maximal Information Coefficient (MIC) and related measures
  - General and can handle non-linear relationships
  - Can be complicated and time intensive to compute

# Information and Probability

- Information relates to possible outcomes of an event
  - transmission of a message, flip of a coin, or measurement of a piece of data

- <span style="color:red">The more certain an outcome, the less information that it contains and vice-versa</span>
  - For example, if a coin has two heads, then an outcome of heads provides no information
  - More quantitatively, the information is related the probability of an outcome
  - The smaller the probability of an outcome, the more information it provides and vice-versa
  - Entropy is the commonly used measure
  - <span style="color:red">– Low entropy means "very predictable".</span>
  - <span style="color:red">– High entropy means "very random".</span>

# Entropy

- For
  - a variable (event), $X$,
  - with $n$ possible values (outcomes), $x_1$, $x_2$ ..., $x_n$
  - each outcome having probability, $p_1$, $p_2$ ..., $p_n$
  - the entropy of $X$, $H(X)$, is given by

$$H(X) = -\sum_{i=1}^{n} p_i \log_2 p_i$$

- Entropy is between 0 and $\log_2 n$ and is measured in bits
  - Thus, entropy is a measure of how many bits it takes to represent an observation of $X$ on average

# Entropy Examples

- For a coin with probability p of heads and probability q = 1 − p of tails

$$H = -p \log_2 p - q \log_2 q$$

  - For p= 0.5, q = 0.5 (fair coin) H = 1
  - For p = 1 or q = 1, H = 0

- What is the entropy of a fair four-sided die?

# Entropy for Sample Data: Example

| Hair Color | Count | $p$ | $-p\log_2 p$ |
|---|---|---|---|
| Black | 75 | 0.75 | 0.3113 |
| Brown | 15 | 0.15 | 0.4105 |
| Blond | 5 | 0.05 | 0.2161 |
| Red | 0 | 0.00 | 0 |
| Other | 5 | 0.05 | 0.2161 |
| Total | 100 | 1.0 | 1.1540 |

Maximum entropy is $\log_2 5 = 2.3219$

# General Approach for Combining Similarities

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

- For the $k^{\text{th}}$ attribute, compute a similarity, $s_k(\mathbf{x}, \mathbf{y})$, in the range [0, 1].

- Define an indicator variable, $\delta_k$, for the $k^{\text{th}}$ attribute as follows:

  $\delta_k$ = 0 if the $k^{\text{th}}$ attribute is an asymmetric attribute and both objects have a value of 0, or if one of the objects has a missing value for the kth attribute

  $\delta_k$ = 1 otherwise

- Compute

$$\text{similarity}(\mathbf{x}, \mathbf{y}) = \frac{\sum_{k=1}^{n} \delta_k s_k(\mathbf{x}, \mathbf{y})}{\sum_{k=1}^{n} \delta_k}$$

# Using Weights to Combine Similarities

- May not want to treat all attributes the same.
  - Use non-negative weights $\omega_k$

  - $similarity(\mathbf{x}, \mathbf{y}) = \dfrac{\sum_{k=1}^{n} \omega_k \delta_k s_k(\mathbf{x}, \mathbf{y})}{\sum_{k=1}^{n} \omega_k \delta_k}$

- Can also define a weighted form of distance

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

# Outline

- Attributes and Objects

- Types of Data

- Data Quality

- Similarity and Distance

- Data Preprocessing

# Data Preprocessing

- Aggregation

- Sampling

- Discretization and Binarization

- Attribute Transformation

- Dimensionality Reduction (Illustrated in the next lecture)

- Feature subset selection

# 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: Feature Aggregation

■ Combine features to form new features:

| Van | Bur | Sur | Edm | Cal |
|---|---|---|---|---|
| 1 | 0 | 0 | 0 | 0 |
| 0 | 1 | 0 | 0 | 0 |
| 1 | 0 | 0 | 0 | 0 |
| 0 | 0 | 0 | 1 | 0 |
| 0 | 0 | 0 | 0 | 1 |
| 0 | 0 | 1 | 0 | 0 |

| BC | AB |
|---|---|
| 1 | 0 |
| 1 | 0 |
| 1 | 0 |
| 0 | 1 |
| 0 | 1 |
| 1 | 0 |

# 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, $0.5^{\circ}$ 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**
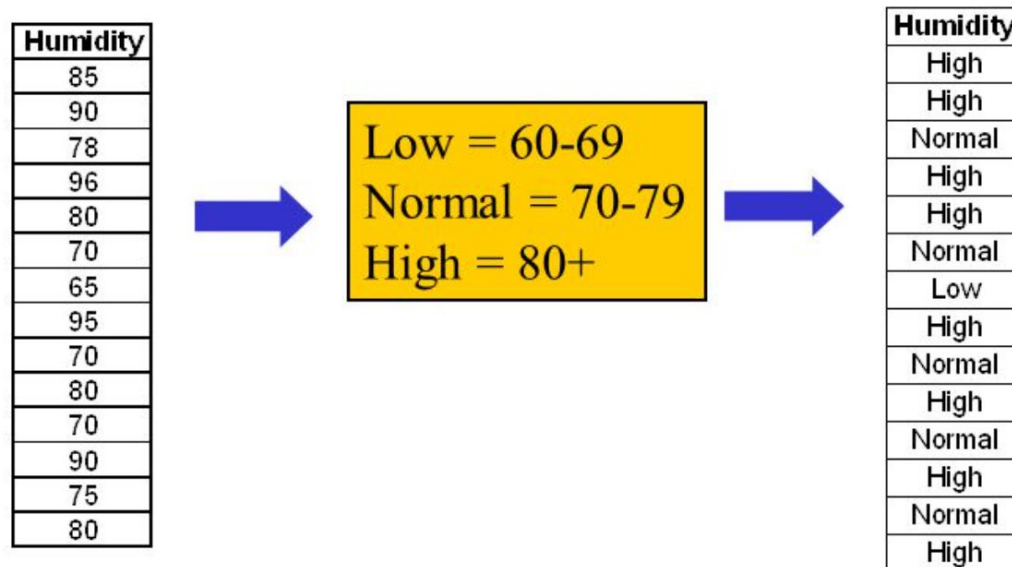
# Types of Sampling

- 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

# 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 used in both unsupervised and supervised settings
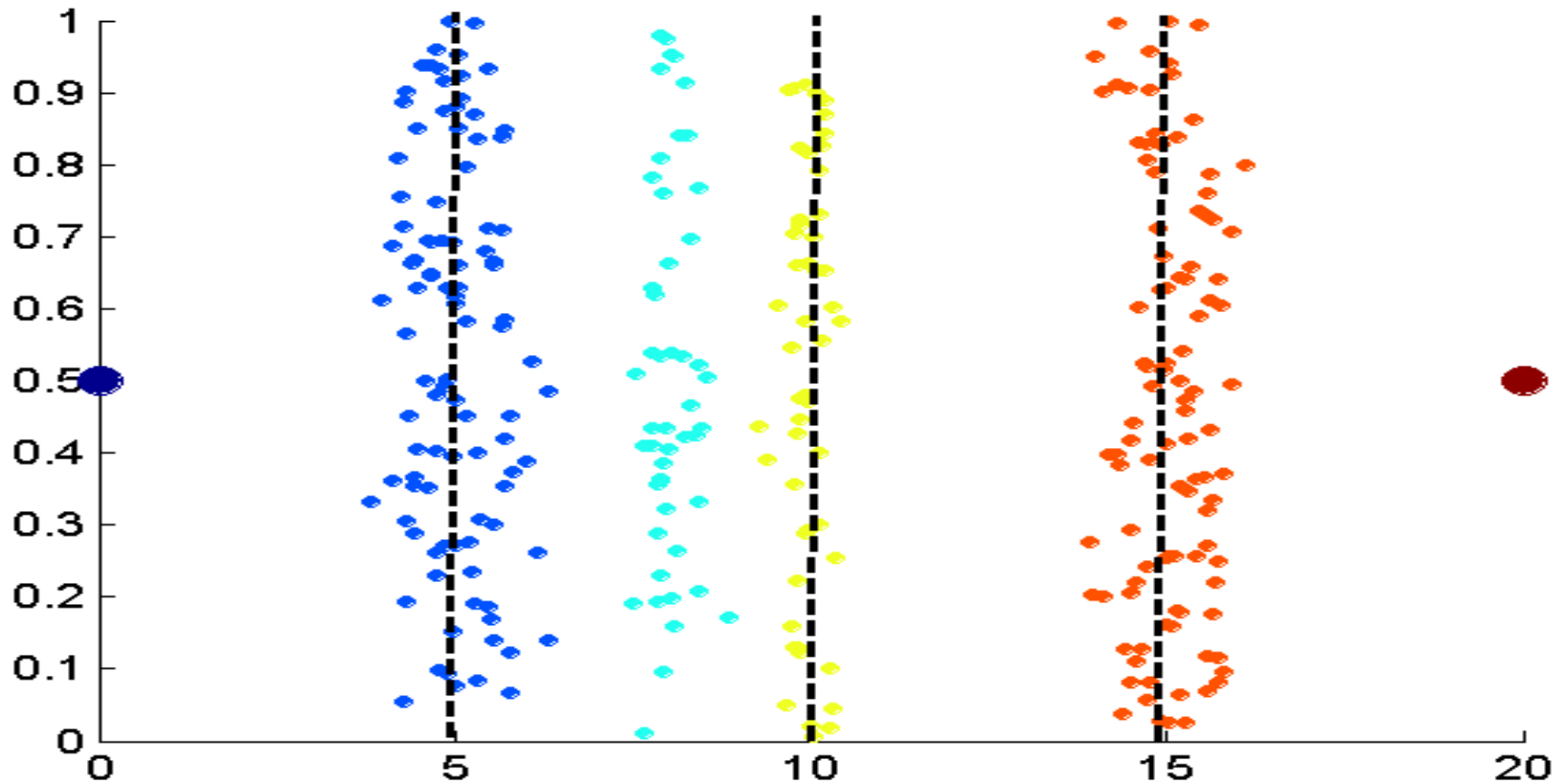
| Humidity |
|----------|
| 85 |
| 90 |
| 78 |
| 96 |
| 80 |
| 70 |
| 65 |
| 95 |
| 70 |
| 80 |
| 70 |
| 90 |
| 75 |
| 80 |

Low = 60-69
Normal = 70-79
High = 80+

| Humidity |
|----------|
| High |
| High |
| Normal |
| High |
| High |
| Normal |
| Low |
| High |
| Normal |
| High |
| Normal |
| High |
| Normal |
| High |

# 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.
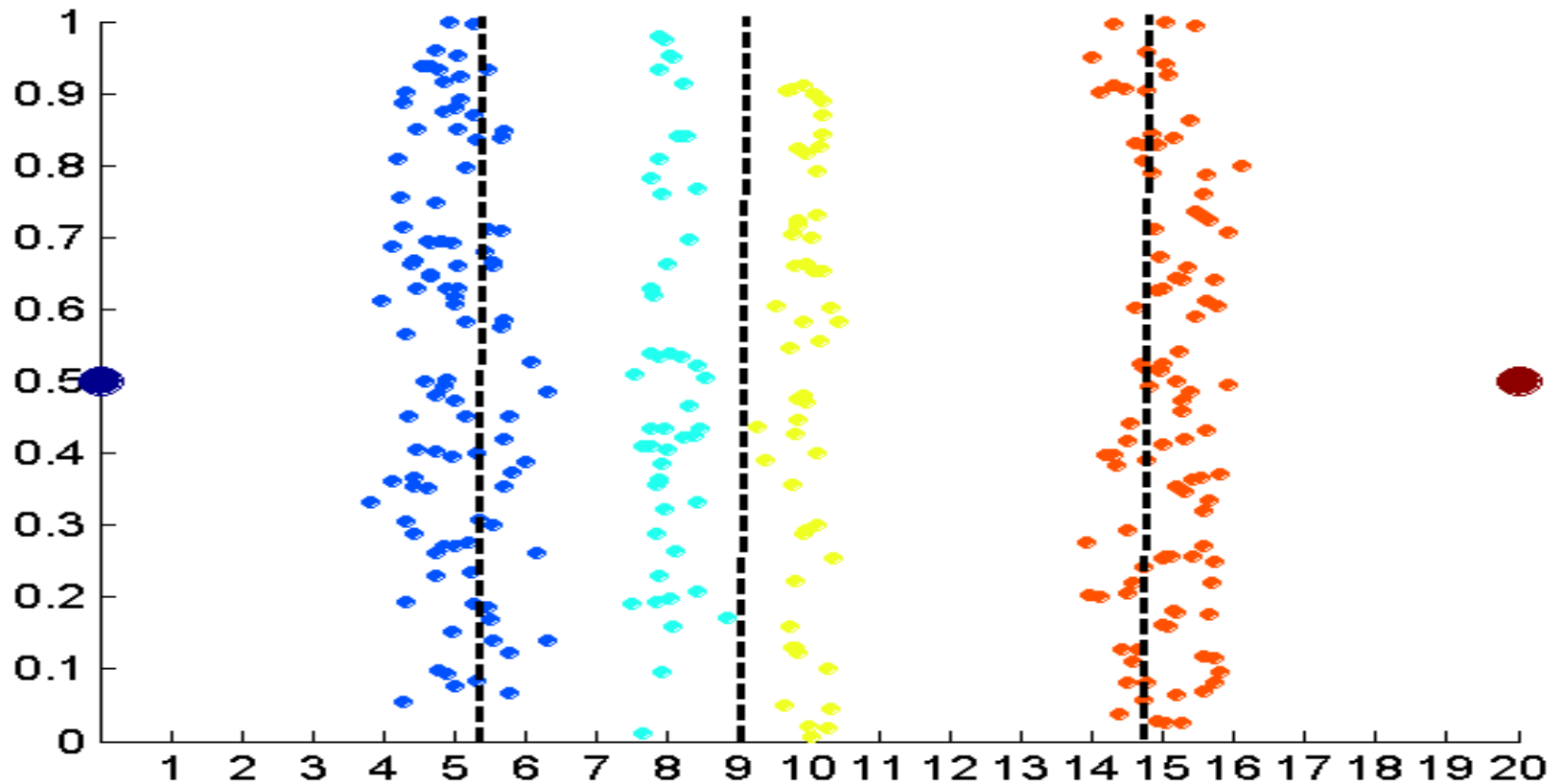
# Unsupervised Discretization



Equal interval width approach used to obtain 4 values.
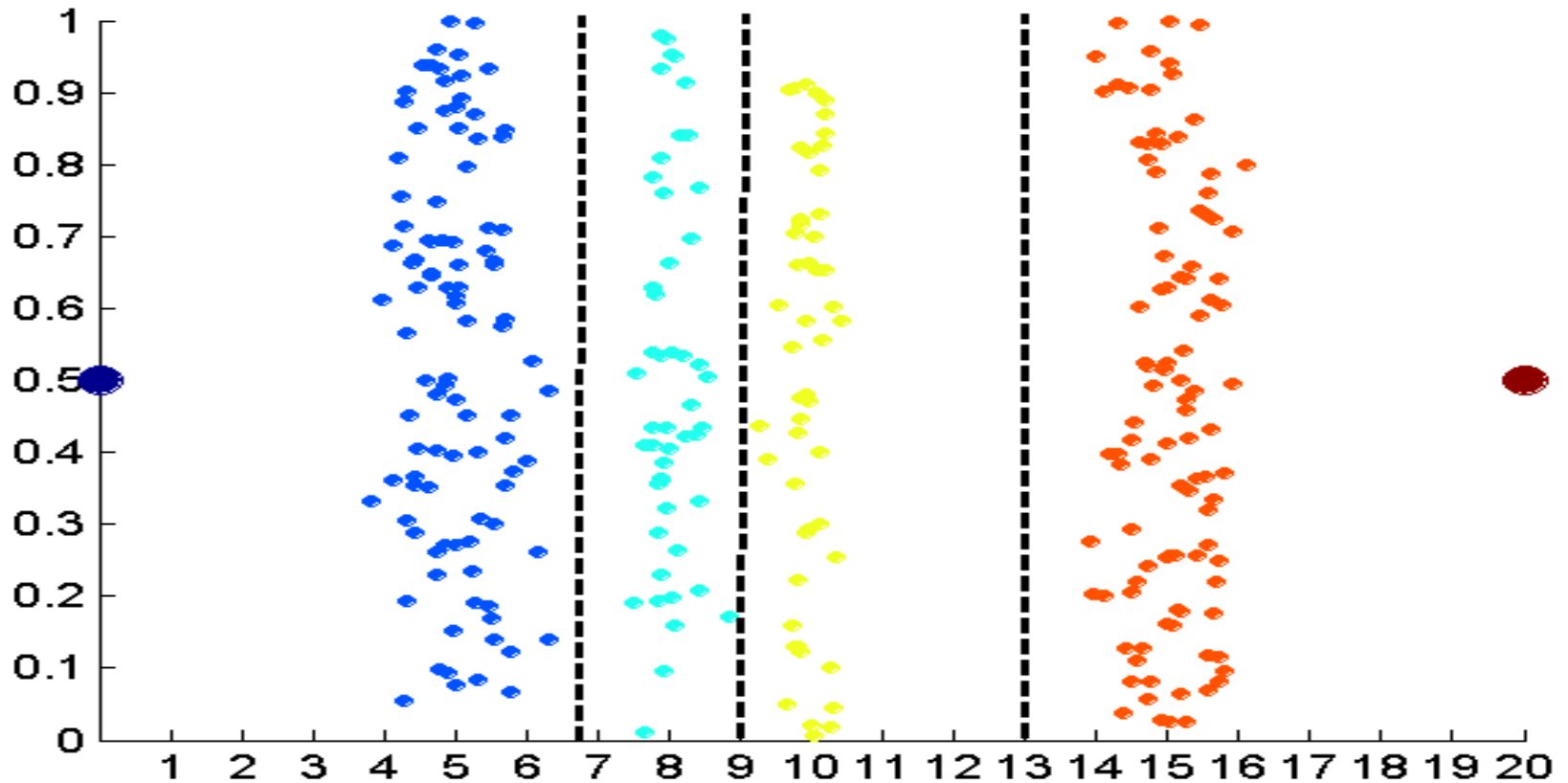
# Unsupervised Discretization



Equal frequency approach used to obtain 4 values.

# Unsupervised Discretization



K-means approach to obtain 4 values.

# Discretization in Supervised Settings

- Many classification algorithms work best if both the independent and dependent variables have only a few values
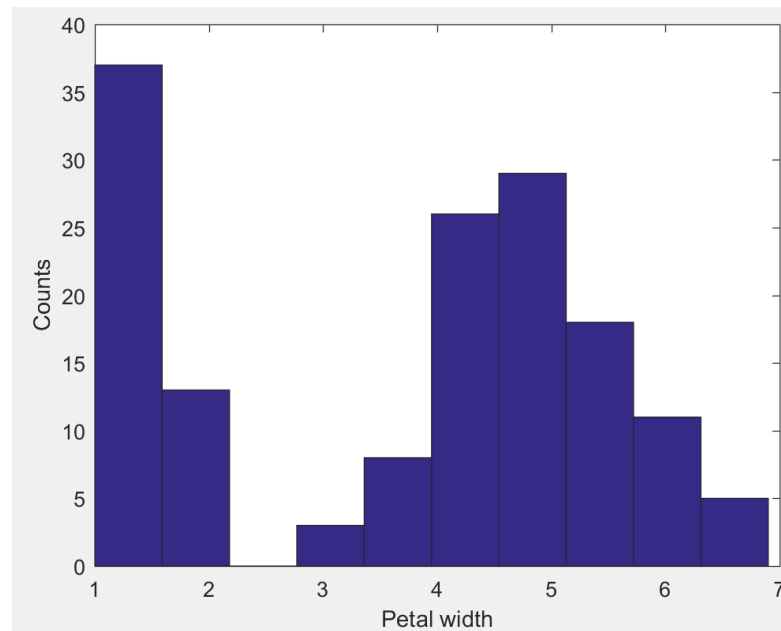  - We give an illustration of the usefulness of discretization using the Iris data set

# 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
        - Sepal length
        - Petal width
        - Petal length

# Discretization: Iris Example

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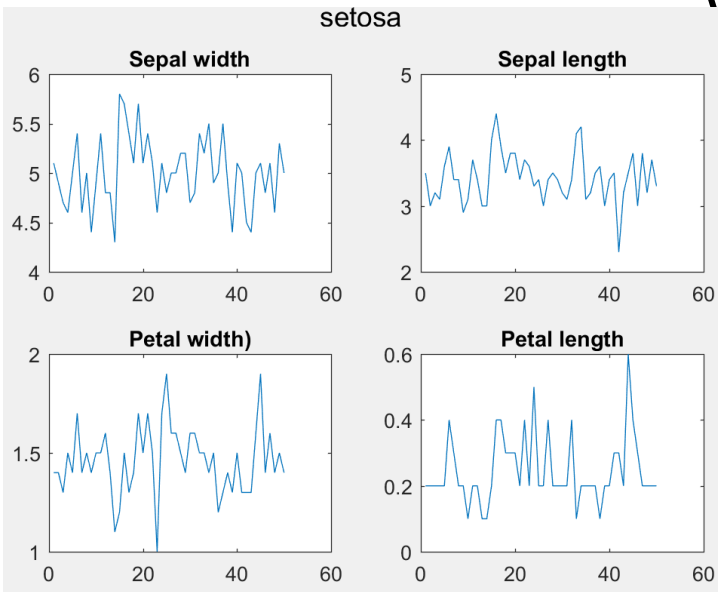


  - Supervised discretization: Use class labels to find breaks
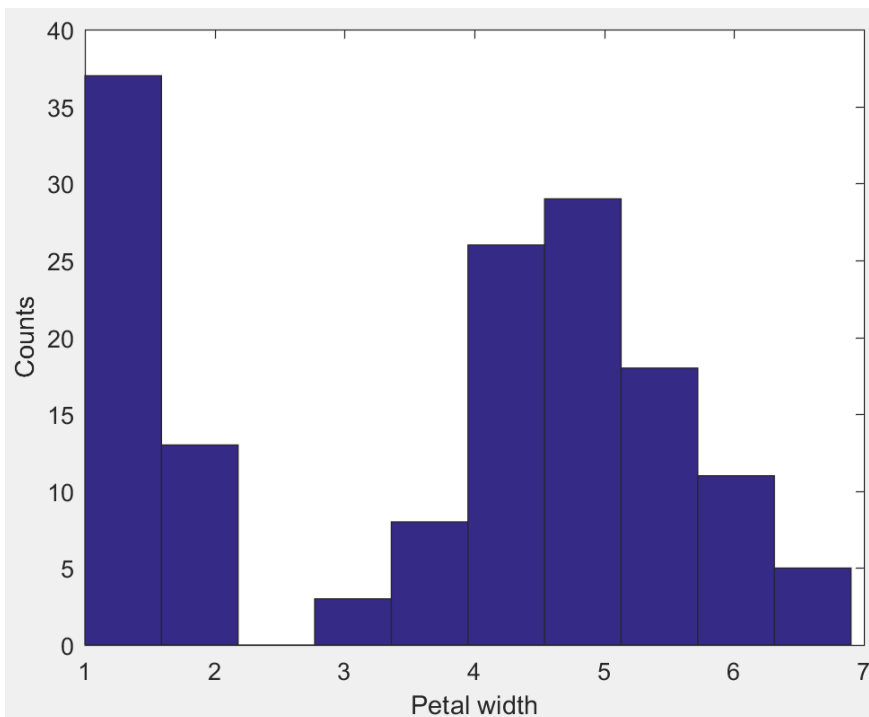
# Discretization: Iris Example

Plot distribution:
Different properties



```matlab
clc;
clear;
close all;
load fisheriris.mat ;
str = unique(species) ;
%% visulization of distribution
i=1;% for type 1
type=str{i};
idx = find(strcmp(species,str{i}));
figure;
subplot(2,2,1);
plot(meas(idx(:),1));title ('Sepal width');
subplot(2,2,2);
plot(meas(idx(:),2));title ('Sepal length');
subplot(2,2,3);
plot(meas(idx(:),3));title 'Petal width');
subplot(2,2,4);
plot(meas(idx(:),4));title ('Petal length');
suptitle(type);
```
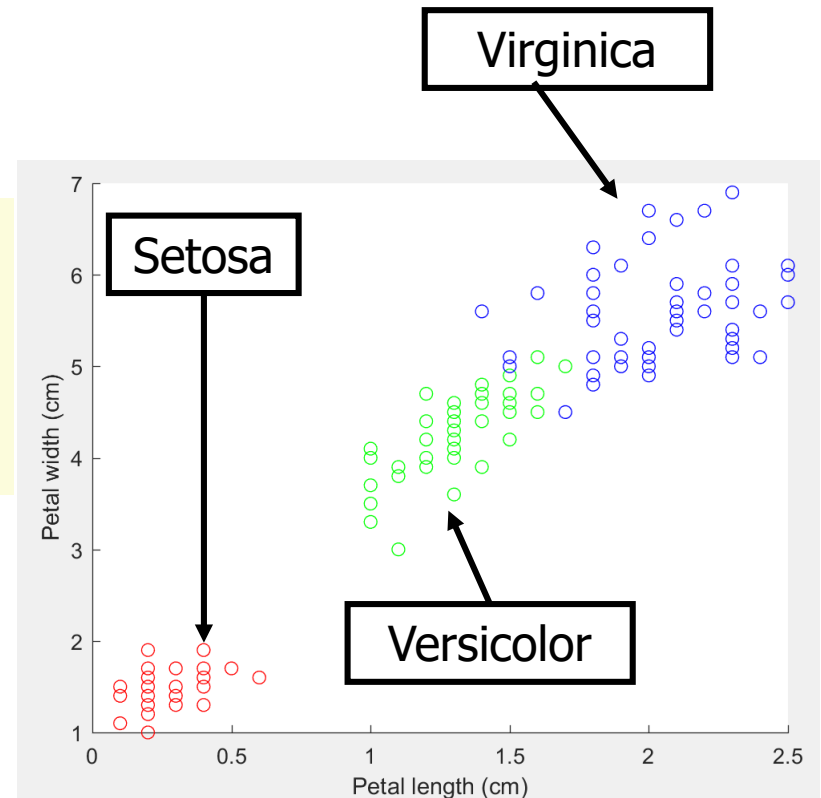
# Discretization: Iris Example

Plot histogram:
'Petal width'
properties

```
%% petal width distribution
figure;
test=meas(:,4);
hist(test);
figure;
test=meas(:,3);
hist(test);
xlabel('Petal width')
ylabel('Counts')
```

# Discretization: Iris Example

```
clc;
clear;
close all;
load fisheriris.mat ;
%% category distribution
figure;
scatter(meas(1:50,4),meas(1:50,3),'r');hold on;
scatter(meas(51:100,4),meas(51:100,3),'g');hold on;
scatter(meas(101:150,4),meas(101:150,3),'b');
xlabel('Petal length (cm)');
ylabel('Petal width (cm)');
```

# 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}

# 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: $x^k$, $\log(x)$, $e^x$, $|x|$

  - Normalization

    - Refers to various techniques to adjust to differences among attributes in terms of frequency of occurrence, mean, variance, range

    - Take out unwanted, common signal, e.g., seasonality

  - In statistics, standardization refers to subtracting off the means and dividing by the standard deviation