# 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 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**

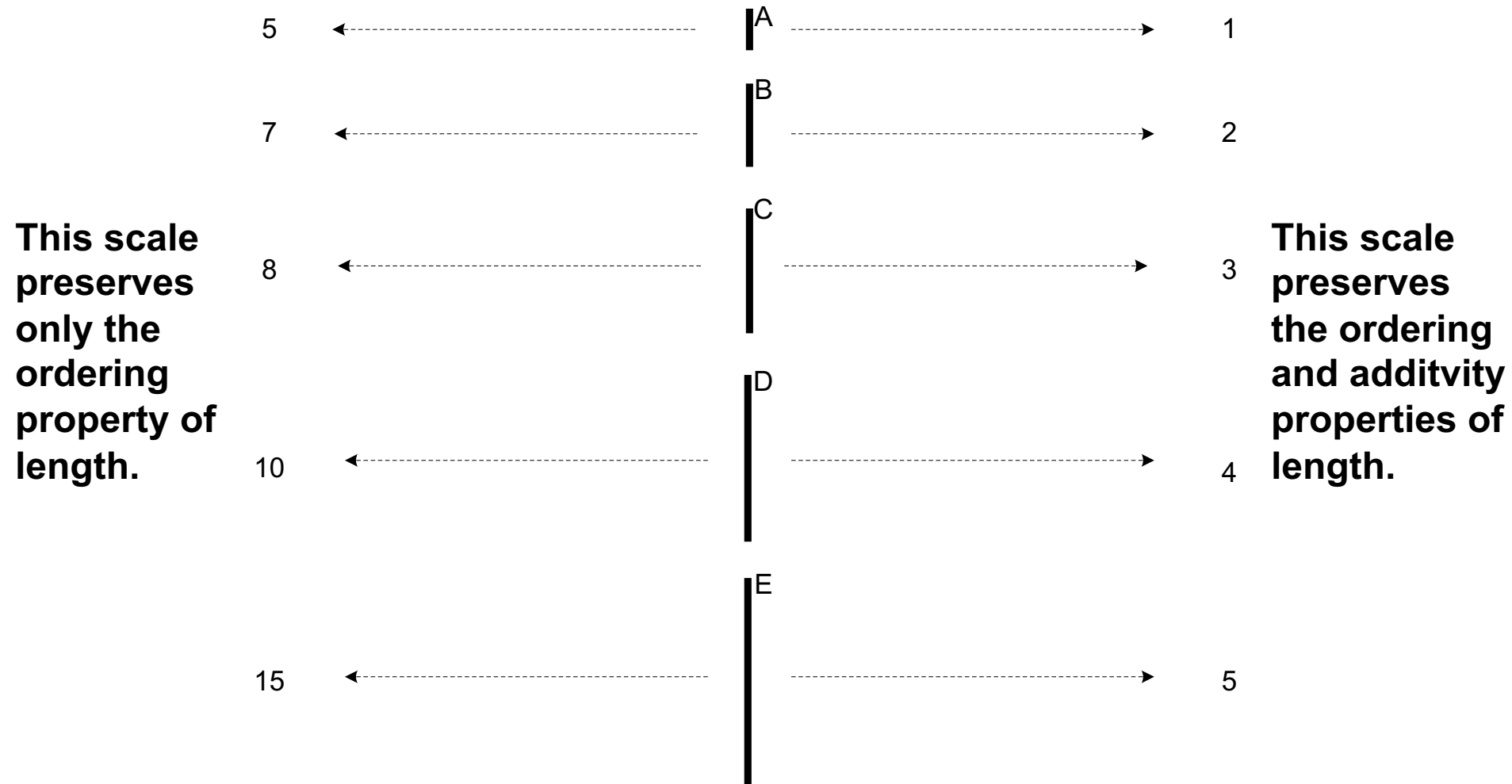| Tid | Refund | Marital Status | Taxable Income | Cheat |
|-----|--------|---------------|----------------|-------|
| 1 | Yes | Single | 125K | No |
| 2 | No | Married | 100K | No |
| 3 | No | Single | 70K | No |
| 4 | Yes | Married | 120K | No |
| 5 | No | Divorced | 95K | Yes |
| 6 | No | Married | 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 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 |
| 8 | C | 3 |
| 10 | D | 4 |
| 15 | E | 5 |

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

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

# 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 Alice'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 Alice?

  - Is this situation analogous to that of temperature?

| | 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? |
| | 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.

# 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

  - In the end, what is meaningful is determined by the domain

# Important Characteristics of Data

– Dimensionality (number of attributes)
  ◆ High dimensional data brings a number of challenges

– Distribution
  ◆ Skewness and sparsity require special handling

– 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

| 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 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 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: $C_6H_6$

**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
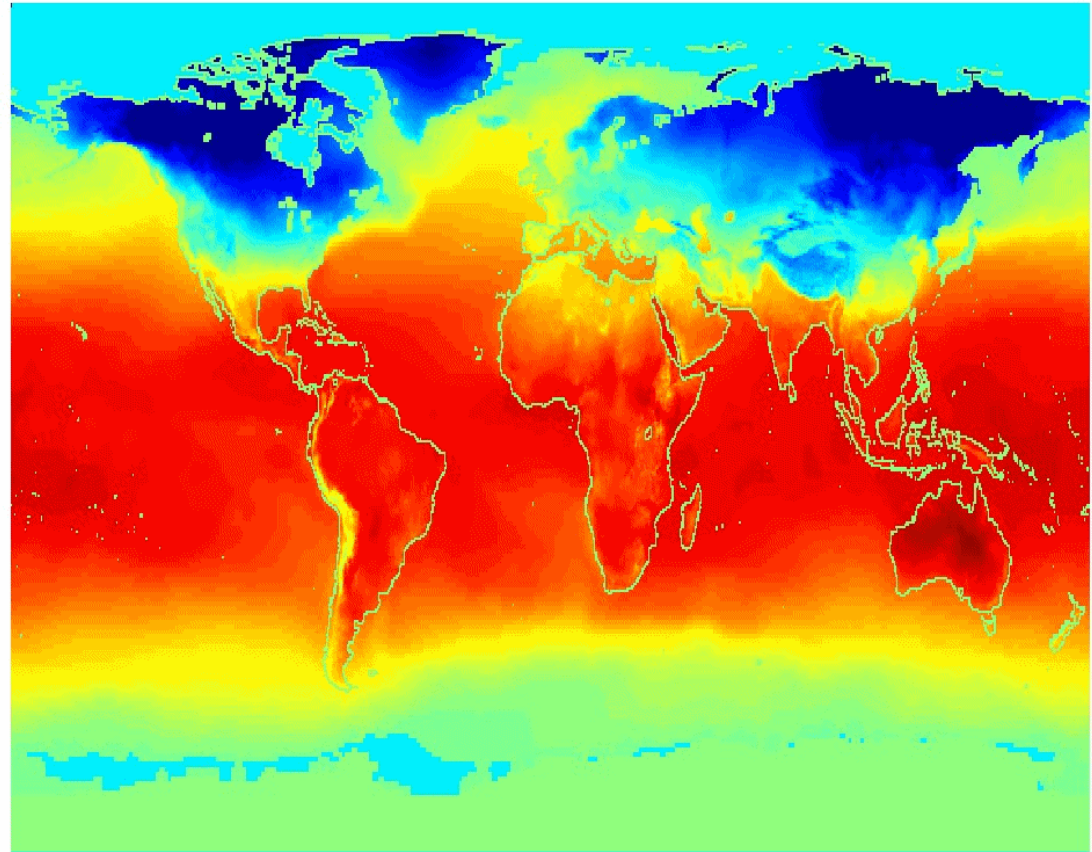the sequence**

# 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

# 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 attribute values
  - Examples: distortion of a person's voice when talking on a poor quality phone and "snow" on television screen



**Two Sine Waves**                **Two Sine Waves + 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 unwanted and interfere 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
  - 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

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

# 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 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, $L_\infty$ 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**

# Visual Interpretation of Distances



- L1-norm(x, y) = a + b

- L2-norm(x, y) = $\sqrt{a^2 + b^2}$

- L∞-norm(x, y) = max(a, b)

● L1-norm is robust to outliers in a few attributes

● L∞-norm is robust to noise in irrelevant attributes

# Mahalanobis Distance

$$\textbf{mahalanobis}(\mathbf{x}, \mathbf{y}) = \sqrt{(x - y)^T \Sigma^{-1} (x - y)}$$



$\Sigma$ **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**

# 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 *$d(\mathbf{x}, \mathbf{y}) = 0$* only if **x** = **y**. (Positive definiteness)
  2. *$d(\mathbf{x}, \mathbf{y}) = d(\mathbf{y}, \mathbf{x})$* for all **x** and **y**. (Symmetry)
  3. *$d(\mathbf{x}, \mathbf{z}) \leq d(\mathbf{x}, \mathbf{y}) + d(\mathbf{y}, \mathbf{z})$* for all points **x**, **y**, and **z**. (Triangle Inequality)

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

- A distance that satisfies these properties is a metric

# Common Properties of a Similarity

- Similarities, also have some well known properties.

  1. $s(\mathbf{x}, \mathbf{y}) = 1$ (or maximum similarity) only if $\mathbf{x} = \mathbf{y}$.

  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, *p* and *q*, have only binary attributes

- Compute similarities using the following quantities

  $f_{01}$ = the number of attributes where *p* was 0 and *q* was 1
  $f_{10}$ = the number of attributes where *p* was 1 and *q* was 0
  $f_{00}$ = the number of attributes where *p* was 0 and *q* was 0
  $f_{11}$ = the number of attributes where *p* was 1 and *q* was 1

- Simple Matching and Jaccard Coefficients

  SMC = number of matches / number of attributes

  $$= (f_{11} + f_{00}) / (f_{01} + f_{10} + f_{11} + f_{00})$$

  J = number of 1-1 matches / number of non-zero attributes

  $$= (f_{11}) / (f_{01} + f_{10} + f_{11})$$

# Cosine Similarity

● If $\mathbf{d}_1$ and $\mathbf{d}_2$ are two document vectors, then

$$\cos(\mathbf{d}_1, \mathbf{d}_2) = <\mathbf{d}_1,\mathbf{d}_2> / \|\mathbf{d}_1\| \|\mathbf{d}_2\|,$$

where $<\mathbf{d}_1,\mathbf{d}_2>$ indicates inner product or vector dot product of vectors, $\mathbf{d}_1$ and $\mathbf{d}_2$, and $\| \mathbf{d} \|$ is the length of vector $\mathbf{d}$.

● Example:

$$\mathbf{d}_1 = 3\ 2\ 0\ 5\ 0\ 0\ 0\ 2\ 0\ 0$$

$$\mathbf{d}_2 = 1\ 0\ 0\ 0\ 0\ 0\ 0\ 1\ 0\ 2$$

$<\mathbf{d}_1, \mathbf{d}2> = 3*1 + 2*0 + 0*0 + 5*0 + 0*0 + 0*0 + 0*0 + 2*1 + 0*0 + 0*2 = 5$

$| \mathbf{d}_1 \| = (3*3+2*2+0*0+5*5+0*0+0*0+0*0+2*2+0*0+0*0)^{0.5} = (42)^{0.5} = 6.481$

$\| \mathbf{d}_2 \| = (1*1+0*0+0*0+0*0+0*0+0*0+0*0+1*1+0*0+2*2)^{0.5} = (6)^{0.5} = 2.449$

$\cos(\mathbf{d}_1, \mathbf{d}_2) = 0.3150$

# Extended Jaccard Coefficient (Tanimoto)

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

$$EJ(\mathbf{x}, \mathbf{y}) = \frac{\mathbf{x} \cdot \mathbf{y}}{\|\mathbf{x}\|^2 + \|\mathbf{y}\|^2 - \mathbf{x} \cdot \mathbf{y}}$$

# 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

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

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

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

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

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

# Visually Evaluating Correlation



Scatter plots showing the similarity from –1 to 1.

**Tan, Steinbach, Karpatne, Kumar**

# 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

# Relation b/w Correlation and Cosine

$$\text{corr}(\mathbf{x}, \mathbf{y}) = \frac{\frac{1}{n-1} \sum_{k=1}^{n} (x_k - \overline{x})(y_k - \overline{y})}{\sqrt{\frac{1}{n-1} \sum_{k=1}^{n} (x_k - \overline{x})^2} * \sqrt{\frac{1}{n-1} \sum_{k=1}^{n} (y_k - \overline{y})^2}}$$

● If we transform x and y by subtracting off their means,

   – $x_m$ = x – mean(x)

   – $y_m$ = y – mean(y)

● Then, corr(x, y) = cos($x_m$, $y_m$)

# Differences Among Proximity Measures

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

Proximity Measures
- Cosine
- Correlation
- Euclidean Distance

- Scaling Operator:

$$\mathbf{y_s} = 2 \times \mathbf{y} = (2, 4, 6, 8, 0, 0, 0)$$

- Translation Operator:

$$\mathbf{y_t} = \mathbf{y} + 5 = (6, 7, 8, 9, 5, 5, 5)$$

- Which proximity measure is invariant to scaling?
  - i.e., Proximity (x, y) = Proximity (x, $y_s$)

- Which proximity measure is invariant to translation?
  - i.e., Proximity (x, y) = Proximity (x, $y_t$)

# Differences Among Proximity Measures

$$\mathbf{x} = (1, 2, 4, 3, 0, 0, 0)$$
$$\mathbf{y} = (1, 2, 3, 4, 0, 0, 0)$$
$$\mathbf{y_s} = 2 \times \mathbf{y} = (2, 4, 6, 8, 0, 0, 0)$$
$$\mathbf{y_t} = \mathbf{y} + 5 = (6, 7, 8, 9, 5, 5, 5)$$

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

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

Choice of suitable measure depends on the needs of the application domain

# Mutual Information

- Measures similarity among two objects as the amount of information shared among them

  – How much information does an object $X$ provide about another object $Y$, and vice-versa?

- General and can handle non-linear relationships
- Complicated (especially for objects with continuous attributes) and time-intensive to compute

# Entropy: Measure of Information

- Information often measured using Entropy, H
- Assume objects $X$ and $Y$ contain discrete values
  - Values in $X$ can range in $u_1, u_2, u_3, \ldots u_m$
  - Values in $Y$ can range in $v_1, v_2, v_3, \ldots v_n$

$$H(X) = -\sum_{j=1}^{m} P(X = u_j) \log_2 P(X = u_j)$$

Individual Entropy

$$H(Y) = -\sum_{k=1}^{n} P(Y = v_k) \log_2 P(Y = v_k)$$

$$H(X, Y) = -\sum_{j=1}^{m} \sum_{k=1}^{n} P(X = u_j, Y = v_k) \log_2 P(X = u_j, Y = v_k)$$

Joint Entropy

# Computing Mutual Information

- Mutual Information, $I(X,Y)$, is defined as:

$$I(X,Y) = H(X) + H(Y) - H(X,Y)$$

- Minimum value: $0$ (no similarity)

- Maximum value: $log_2(min(m,n))$
  - Where $m$ and $n$ are the number of possible values of $X$ and $Y$, respectively

- Normalized Mutual Information =

$$I(X,Y)/ log_2(min(m,n))$$

# Mutual Information Example

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

Correlation = 0
Mutual Information = 1.9502
Normalized Mutual Information = 1.9502/log$_2$(4) = 0.9751

**Table 2.14.** Entropy for $\mathbf{x}$

| $x_j$ | $P(\mathbf{x} = x_j)$ | $-P(\mathbf{x} = x_j) \log_2 P(\mathbf{x} = x_j)$ |
|---|---|---|
| -3 | 1/7 | 0.4011 |
| -2 | 1/7 | 0.4011 |
| -1 | 1/7 | 0.4011 |
| 0 | 1/7 | 0.4011 |
| 1 | 1/7 | 0.4011 |
| 2 | 1/7 | 0.4011 |
| 3 | 1/7 | 0.4011 |
| $H(\mathbf{x})$ | | 2.8074 |

**Table 2.15.** Entropy for $\mathbf{y}$

| $y_k$ | $P(\mathbf{y} = y_k)$ | $-P(\mathbf{y} = y_k) \log_2 (P(\mathbf{y} = y_k)$ |
|---|---|---|
| 9 | 2/7 | 0.5164 |
| 4 | 2/7 | 0.5164 |
| 1 | 2/7 | 0.5164 |
| 0 | 1/7 | 0.4011 |
| $H(\mathbf{y})$ | | 1.9502 |

**Table 2.16.** Joint entropy for $\mathbf{x}$ and $\mathbf{y}$

| $x_j$ | $y_k$ | $P(\mathbf{x} = x_j, \mathbf{y} = x_k)$ | $-P(\mathbf{x} = x_j, \mathbf{y} = x_k) \log_2 P(\mathbf{x} = x_j, \mathbf{y} = x_k)$ |
|---|---|---|---|
| -3 | 9 | 1/7 | 0.4011 |
| -2 | 4 | 1/7 | 0.4011 |
| -1 | 1 | 1/7 | 0.4011 |
| 0 | 0 | 1/7 | 0.4011 |
| 1 | 1 | 1/7 | 0.4011 |
| 2 | 4 | 1/7 | 0.4011 |
| 3 | 9 | 1/7 | 0.4011 |
| $H(\mathbf{x}, \mathbf{y})$ | | | 2.8074 |

# Data Preprocessing

- Aggregation

- Sampling

- Dimensionality Reduction

- Discretization and Binarization

- Attribute 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

# Example: Precipitation in Australia …

● We want to study the variability in precip for 3,030 0.5∘ by 0.5∘ grid cells in Australia from the period 1982 to 1993.



**Standard Deviation of Average Monthly Precipitation (in cm)**

**Standard Deviation of Average Yearly Precipitation (in cm)**

# 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

- Choosing a sampling scheme
  – Type of sampling technique
  – Sample size

# Types of Sampling

- **Simple Random Sampling**
  - There is an equal probability of selecting any particular object
  - 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



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

# 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
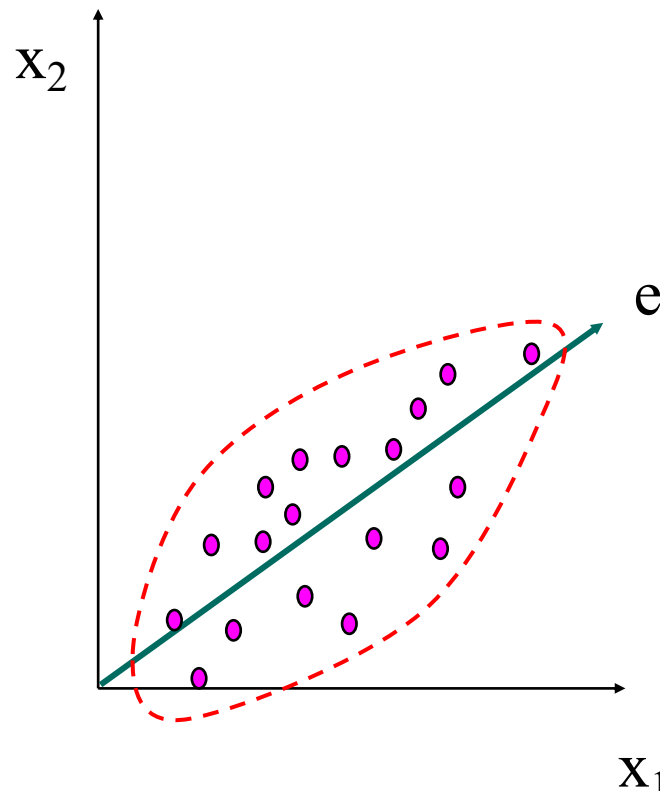


- **Randomly generate 500 points**

- **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



Interactive tool for visualizing PCA:
http://setosa.io/ev/principal-component-analysis/

# Dimensionality Reduction: PCA



256

# 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.
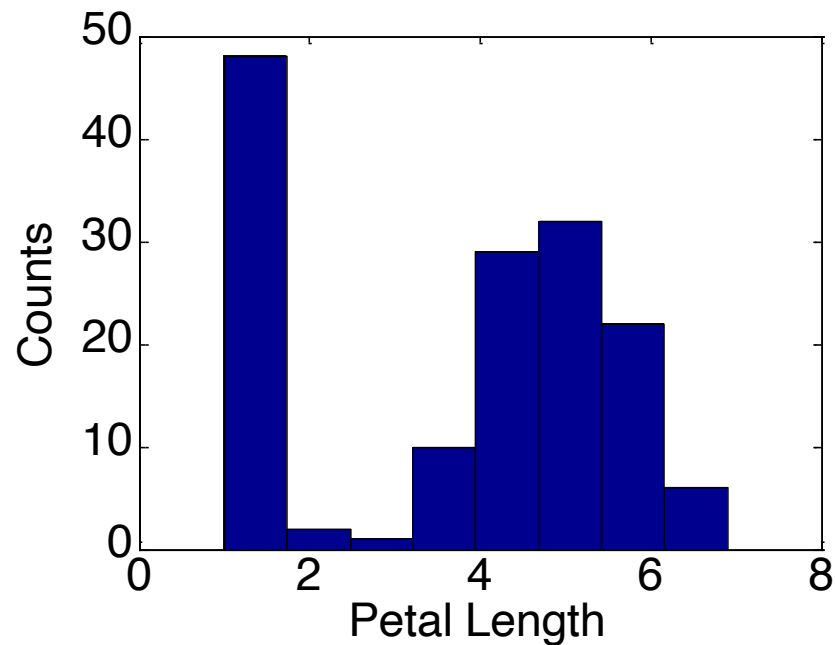Petal width medium or petal length medium implies Versicolour.
Petal width high or petal length high implies Virginica.

# Discretization: Iris Example ...
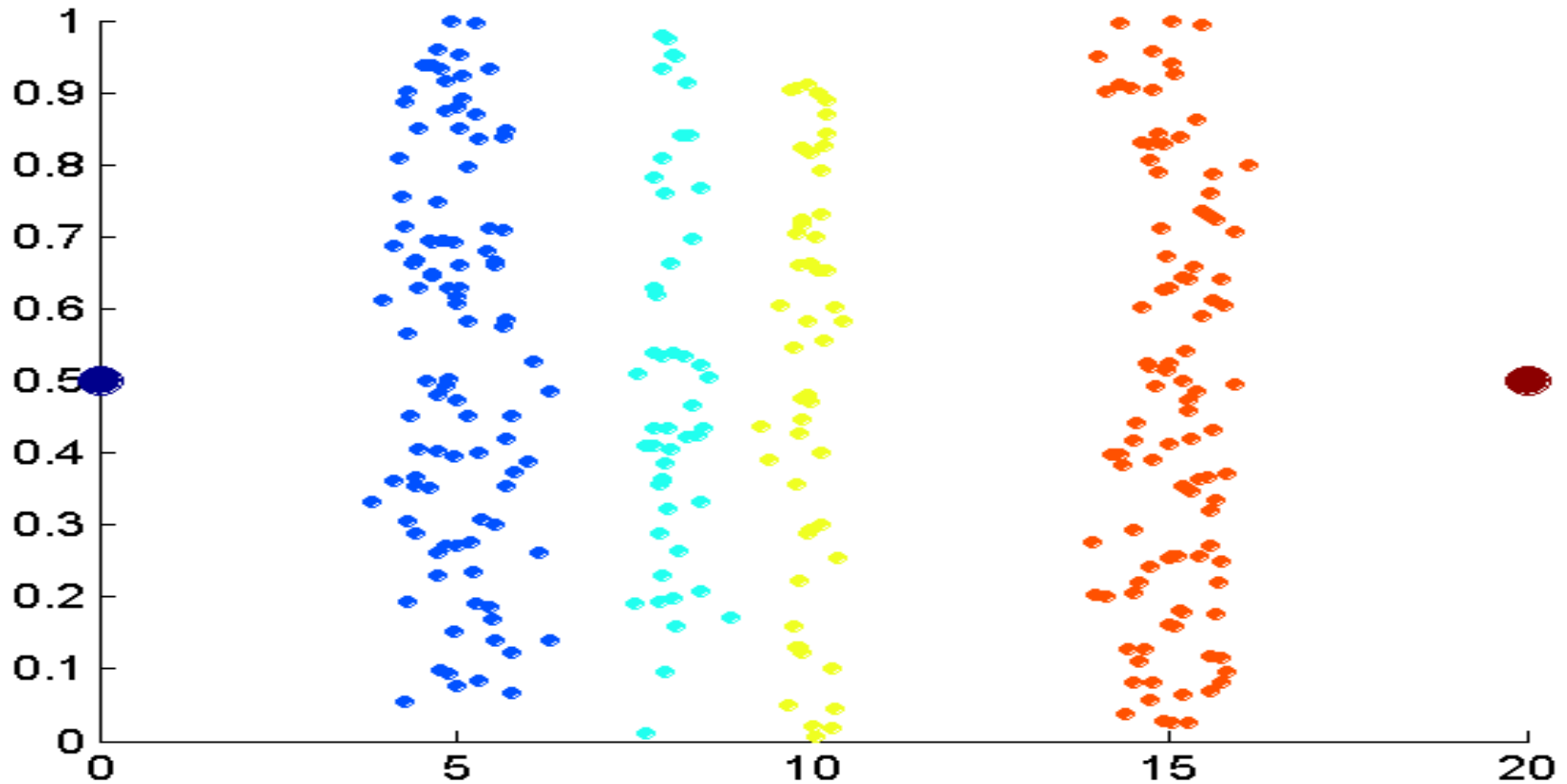
● How can we tell what the best discretization is?

– Unsupervised discretization: find breaks in the data values
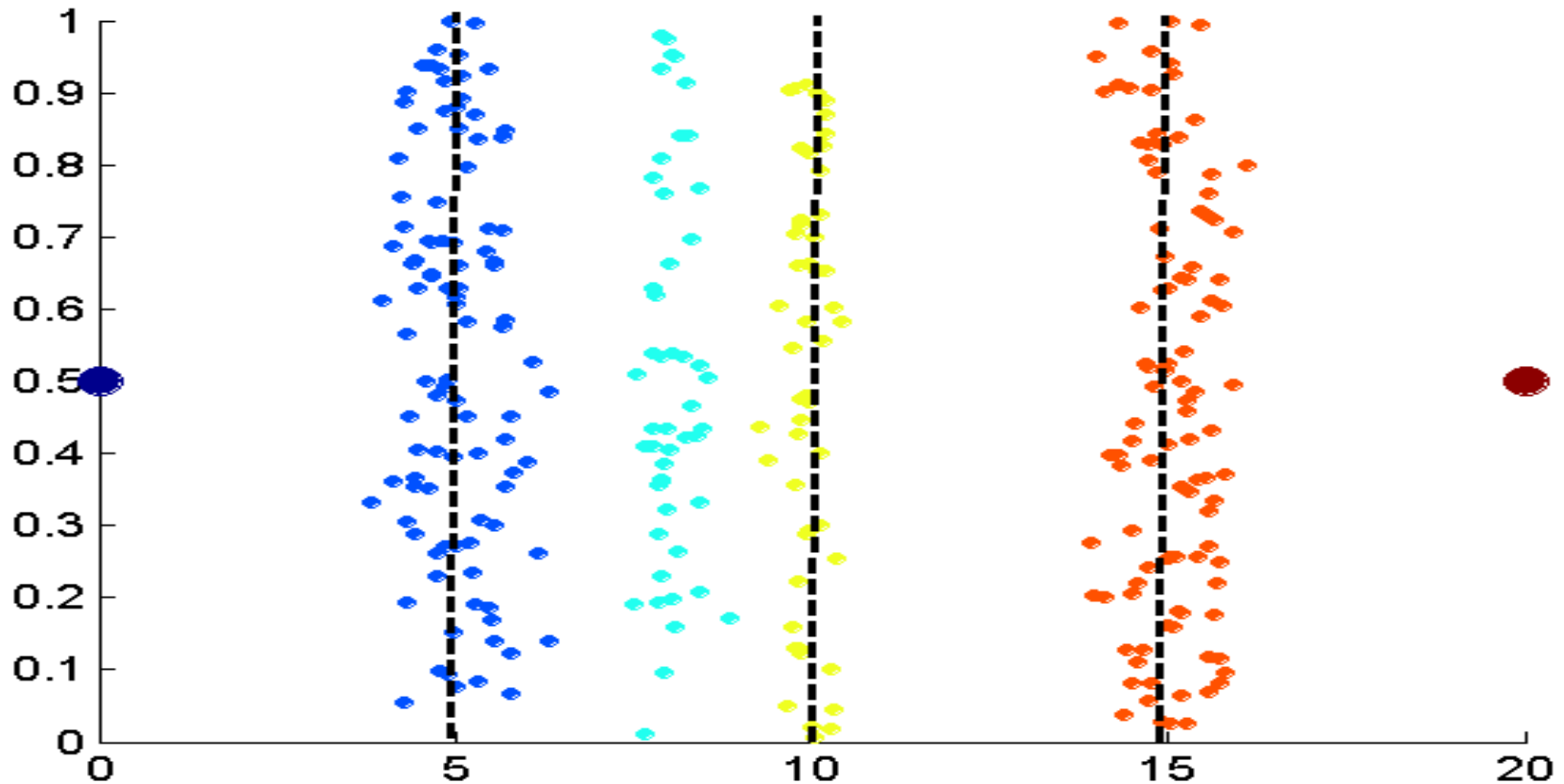
◆ Example:
Petal Length



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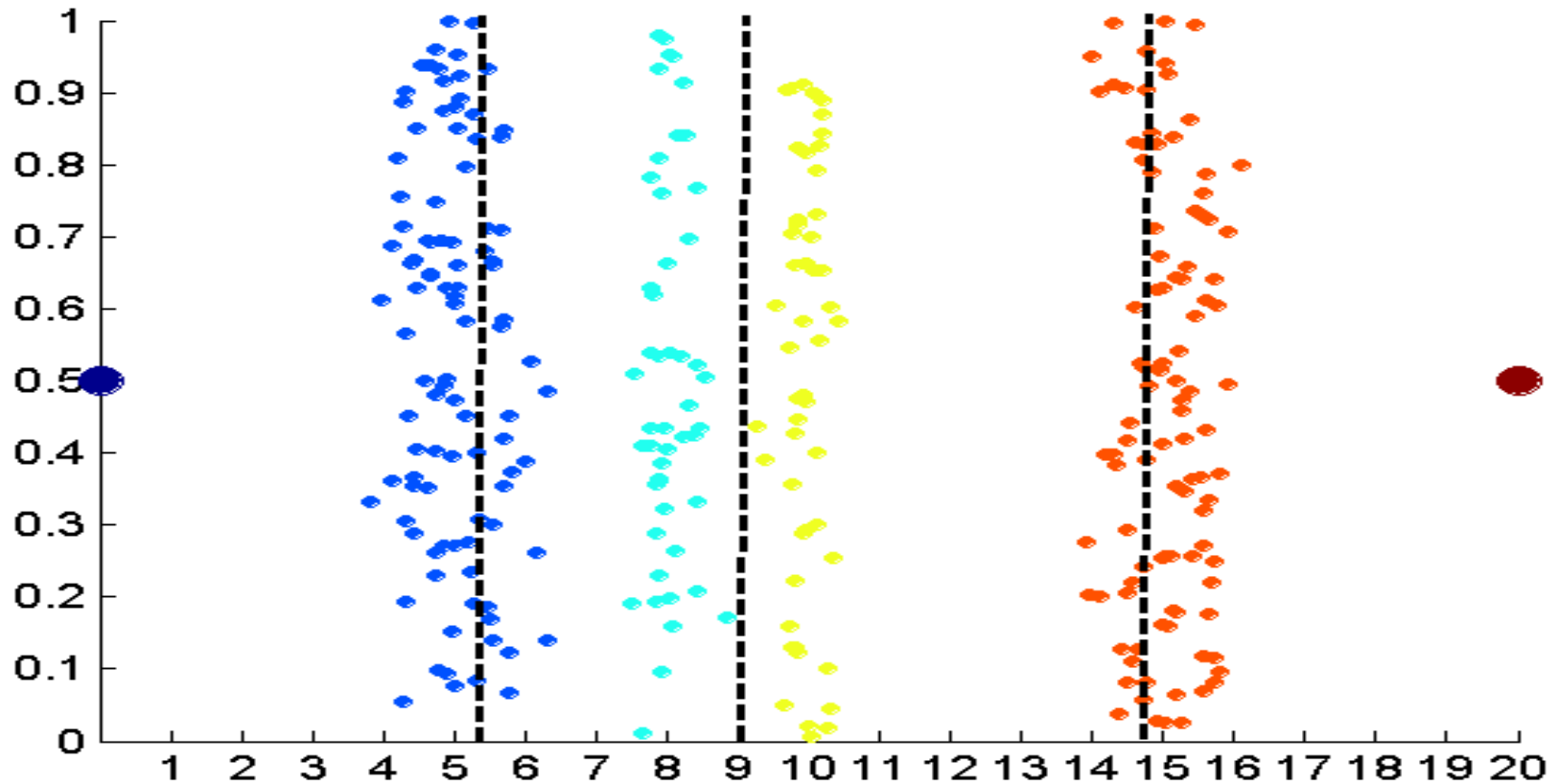
# Discretization Without Using Class Labels



**Data consists of four groups of points and two outliers. Data is one-dimensional, but a random y component is added to reduce overlap.**
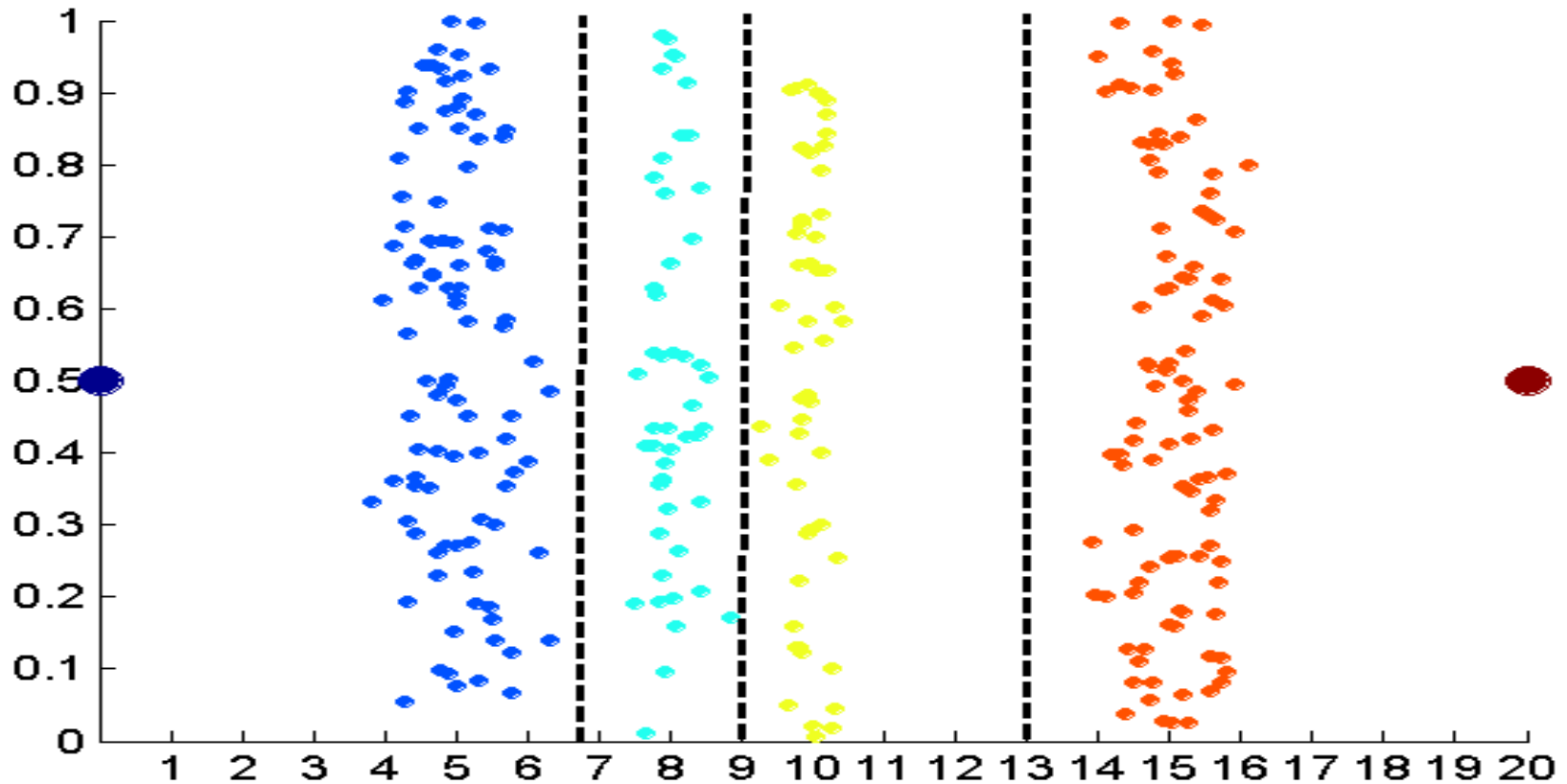
# Discretization Without Using Class Labels



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

# Discretization Without Using Class Labels



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

# Discretization Without Using Class Labels



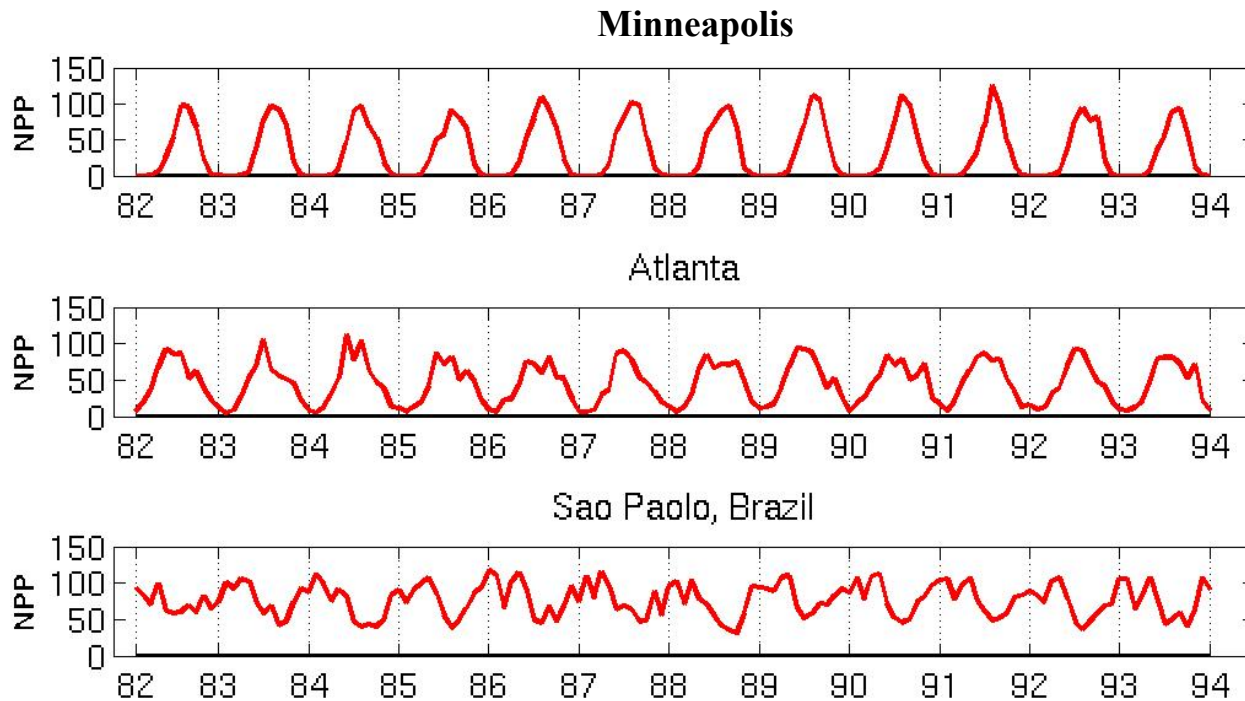**K-means** approach to obtain 4 values.

# 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}

# 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 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
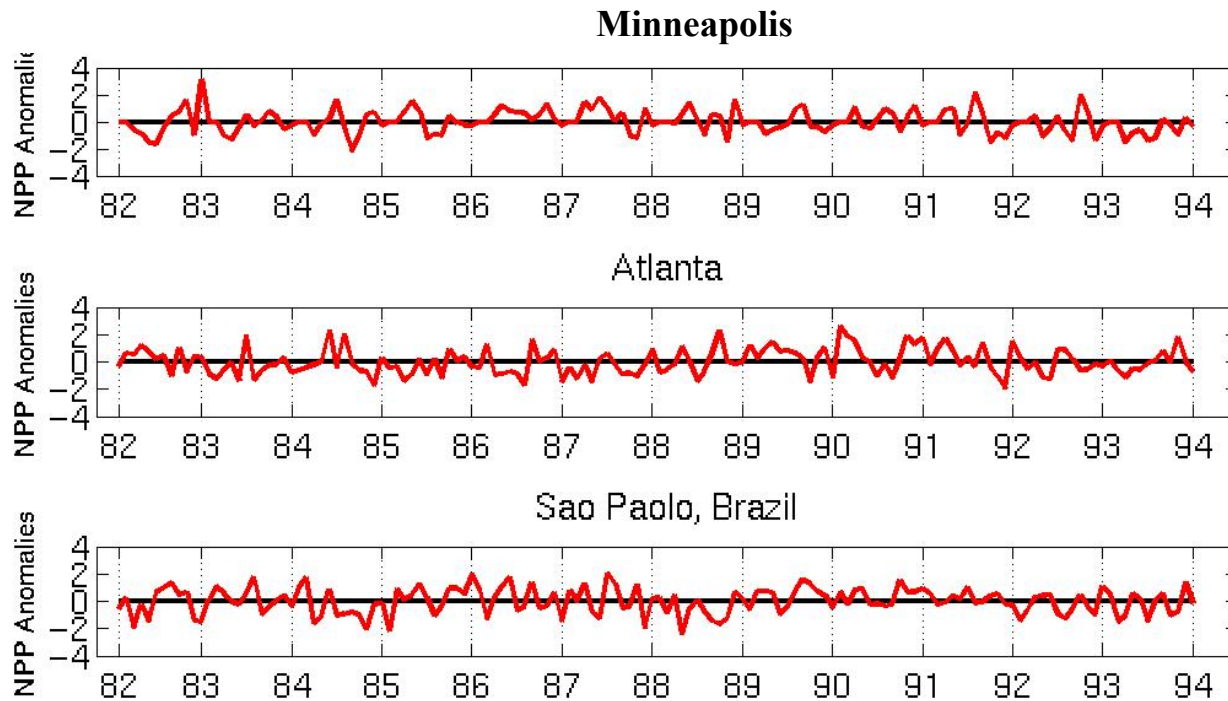
# Example: Sample Time Series of Plant Growth

**Minneapolis**



**Atlanta**

**Sao Paolo, Brazil**

**Net Primary Production (NPP) is a measure of plant growth used by ecosystem scientists.**

## Correlations between time series

|             | Minneapolis | Atlanta | Sao Paolo |
|-------------|-------------|---------|-----------|
| Minneapolis | 1.0000      | 0.7591  | -0.7581   |
| Atlanta     | 0.7591      | 1.0000  | -0.5739   |
| Sao Paolo   | -0.7581     | -0.5739 | 1.0000    |

# Seasonality Accounts for Much Correlation



Minneapolis

Atlanta

Sao Paolo, Brazil

Normalized using monthly Z Score:

Subtract off monthly mean and divide by monthly standard deviation

## Correlations between time series

|  | Minneapolis | Atlanta | Sao Paolo |
|---|---|---|---|
| Minneapolis | 1.0000 | 0.0492 | 0.0906 |
| Atlanta | 0.0492 | 1.0000 | -0.0154 |
| Sao Paolo | 0.0906 | -0.0154 | 1.0000 |