# FDSA

# About Data

Material Adaptation from Introduction to Data Mining : Tan, Steinbach, Karpatne, Kumar

# Outline

- Attributes and Objects

- Types of Data

- Data Quality

- Similarity and Distance

- Data Preprocessing

# Data   : Distinction

- Data
- Information
- Knowledge

# Data   : Distinction

- Data
- Information
- Knowledge

12 , 34 , 16 , 32, 18,  35

What are these ?

*Data   ? Information  ?*

# Data   : Distinction

- Data
- Information
- Knowledge

Last Three days
(June 3,4,5)
Min – Max Temp of
Kathmandu


12 , 34 , 16 , 32, 18,  35


What are these ?

*Data  ? Information ?
Knowledge ?*

# Data : Distinction

- Data
- Information
- Knowledge

Last Three days
Min – Max Temp of
Kathmandu


12 , 34 , 16 , 32, 18,  35


Average Temp ?
23,  24,  26.5

# Data   : Distinction

- Data
- Information
- Knowledge

Last Three days (June 3,4,5)
Min – Max Temp of Kathmandu

12 , 34 , 16 , 32, 18,  35

**A traveler seeks Knowledge regarding traveling to KTM in the month of June!**

*Whether Jacket is needed or not if  the travel is during the month of June ?*

# 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**

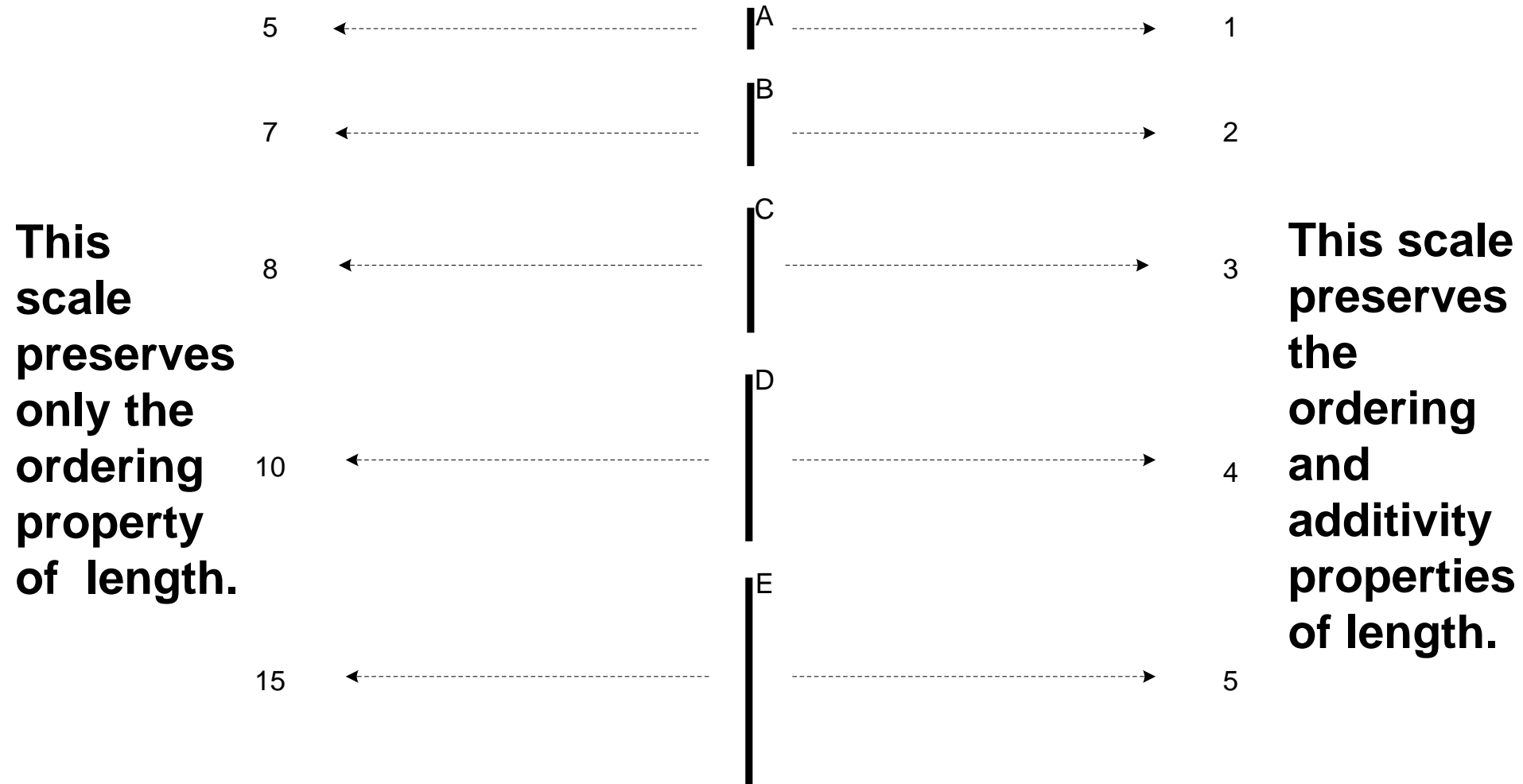| Tid | Refund | Marital Status | Taxable Income | Cheat |
|-----|--------|----------------|----------------|-------|
| 1 | Yes | Single | 125K | No |
| 2 | No | Married | 100K | No |
| 3 | No | Single | 70K | No |
| 4 | Yes | Married | 120K | No |
| 5 | No | Divorced | 95K | Yes |
| 6 | No | Married | 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
    - 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.**

5 ← - - - - - - - - - - - - - - - - A - - - - - - - - - - - - - - → 1

7 ← - - - - - - - - - - - - - - - - B - - - - - - - - - - - - - - → 2

C

8 ← - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - → 3

D

10 ← - - - - - - - - - - - - - - - - - - - - - - - - - - - - - → 4

E

15 ← - - - - - - - - - - - - - - - - - - - - - - - - - - - - → 5

**This scale preserves the ordering and additivity 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?

- The meaning of 0 (Zero) difference

    - Last night temp reached up to 0 ° C

    - As a full-time student, he is earning 0 Rs per month!

| | 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."*

- Asymmetric attributes typically arise from objects that are sets

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

# Types of data sets  (Organization Aspect)

- 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: C6H6

# Ordered Data

- Sequences of transactions

**Items/Events**

$$( A\ B)\quad (D)\quad (C\ E)$$
$$( B\ D)\quad (C)\quad (E)$$
$$( C\ D)\quad (B)\quad (A\ E)$$
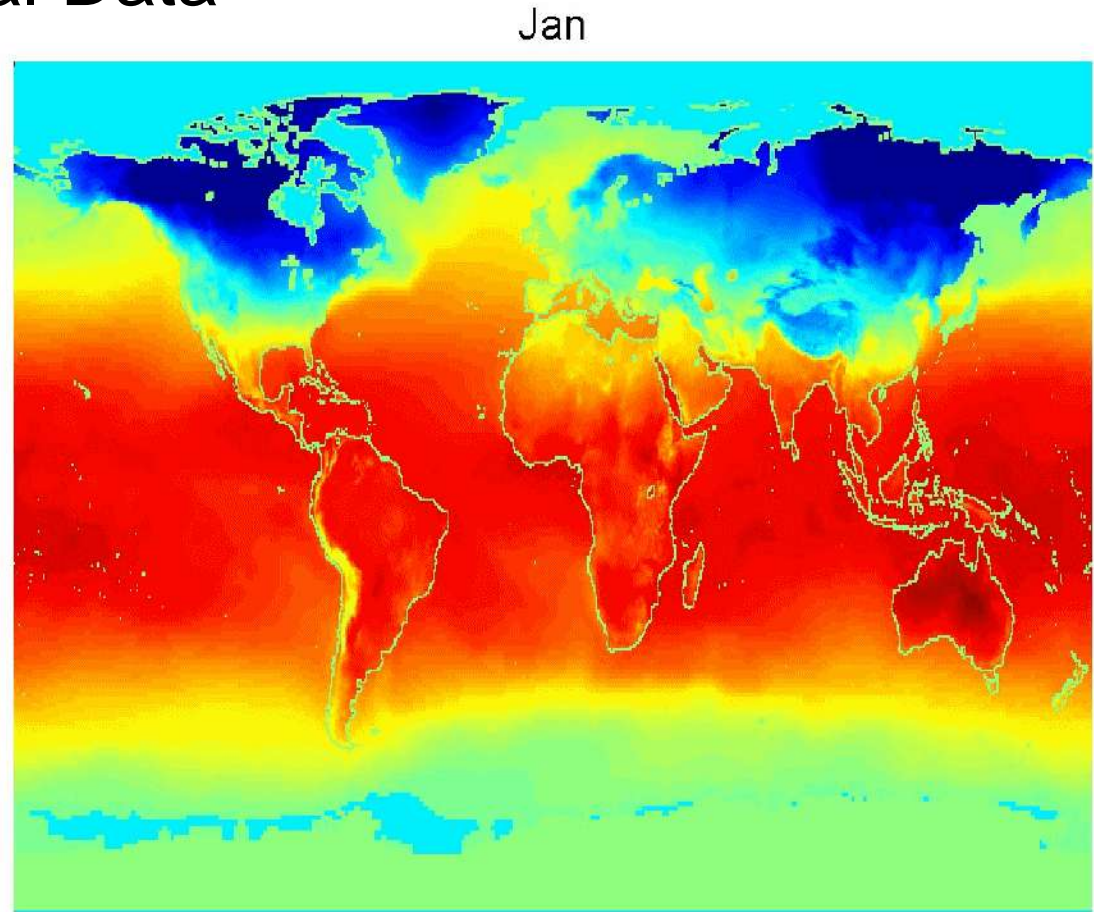
**An element of
the sequence**

# 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

# 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

# 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.

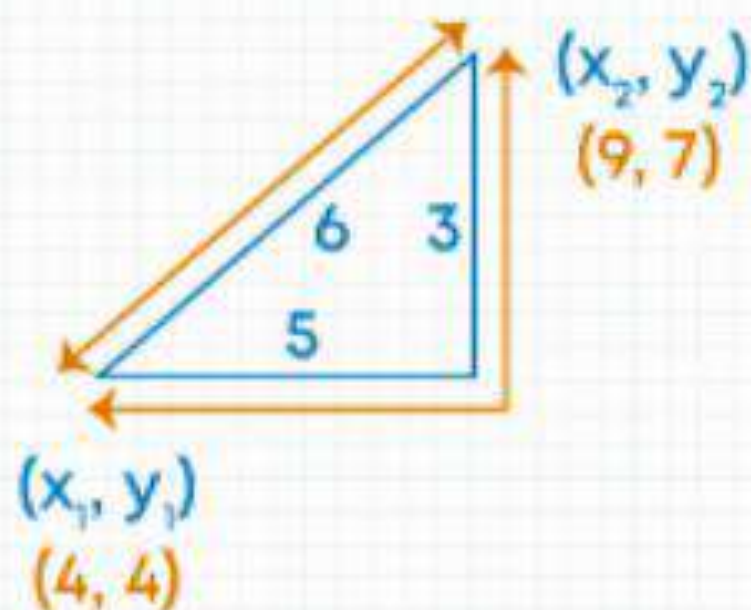| Attribute Type | Dissimilarity | Similarity |
|---|---|---|
| Nominal | $d = \begin{cases} 0 & \text{if } x = y \\ 1 & \text{if } x \neq y \end{cases}$ | $s = \begin{cases} 1 & \text{if } x = y \\ 0 & \text{if } x \neq y \end{cases}$ |
| Ordinal | $d = |x - y|/(n - 1)$ (values mapped to integers 0 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.

# Example:



$(x_2, y_2)$
$(9, 7)$

$(x_1, y_1)$
$(4, 4)$

**Euclidean distance**

$= \sqrt{(x_2 - x_1)^2 + (y_2 - y_1)^2}$

$= \sqrt{(9 - 4)^2 + (7 - 4)^2}$

$= \sqrt{5^2 + 3^2}$

$= \sqrt{25 + 9}$

$= \sqrt{34}$

$= 5.83$

# 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 $x$ and $y$.
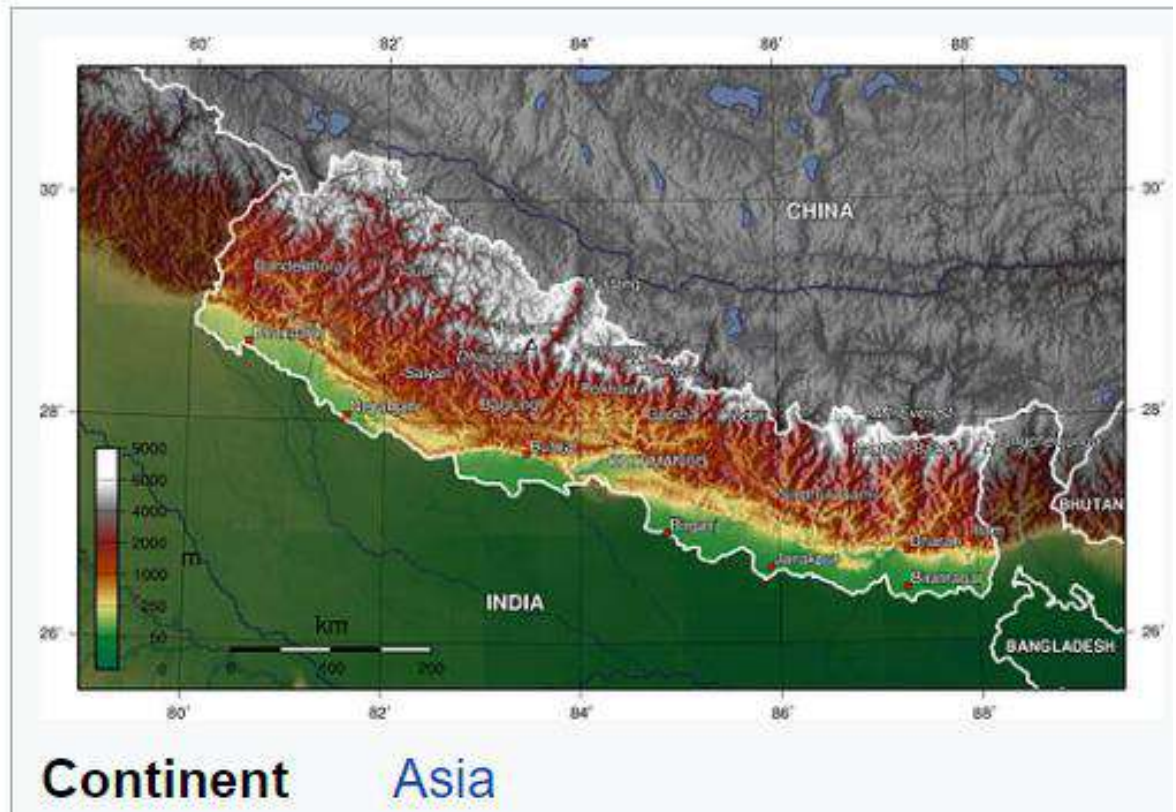
# Minkowski Distance: Examples

- $r = 1$.  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.

Nepal measures about 880 kilometers (547 mi) along its Himalayan axis by 150 to 250 kilometers (93 to 155 mi) across. It has an area of 147,516 km$^2$ (56,956 sq mi).[1]

## Geography of Nepal (नेपाल)
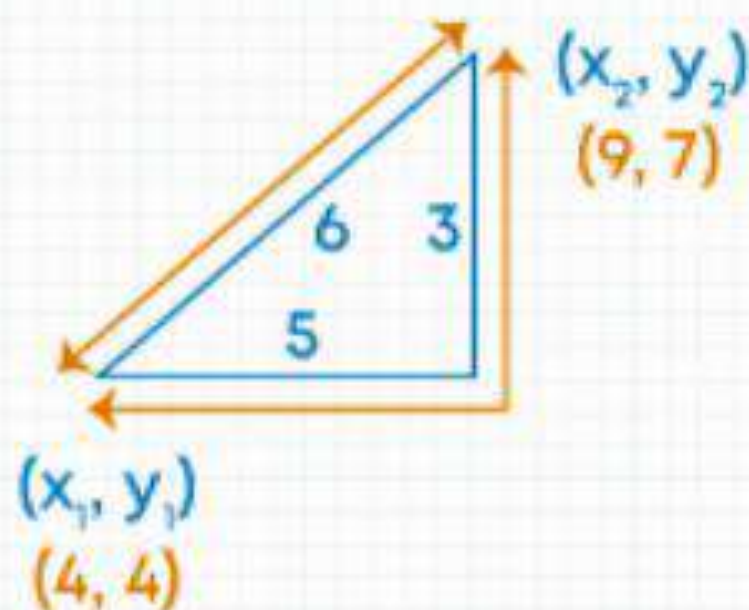


**Continent**     Asia

# 1027.67 km

**East west Highway** of **Nepal** is also known as the Mahendra **highway** is the longest roadway

Length of **east-west highway** is 1027.67 km.   Feb 23, 2016

# Manhattan Street Map

# Example:



**Euclidean distance**

$$= \sqrt{(x_2 - x_1)^2 + (y_2 - y_1)^2}$$

$$= \sqrt{(9 - 4)^2 + (7 - 4)^2}$$

$$= \sqrt{5^2 + 3^2}$$

$$= \sqrt{25 + 9}$$

$$= \sqrt{34}$$

$$= 5.83$$

**Manhattan distance**

$$= |x_2 - x_1| + |y_2 - y_1|$$

$$= |9 - 4| + |7 - 4|$$

$$= 5 + 3$$

$$= 8$$

# Minkowski Distance

| point | x | y |
|-------|---|---|
| p1 | 0 | 2 |
| p2 | 2 | 0 |
| p3 | 3 | 1 |
| p4 | 5 | 1 |

| L1 | p1 | p2 | p3 | p4 |
|----|----|----|----|----|
| p1 | 0 | 4 | 4 | 6 |
| p2 | 4 | 0 | 2 | 4 |
| p3 | 4 | 2 | 0 | 2 |
| p4 | 6 | 4 | 2 | 0 |

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

| $L_\infty$ | p1 | p2 | p3 | p4 |
|------------|----|----|----|----|
| p1 | 0 | 2 | 3 | 5 |
| p2 | 2 | 0 | 1 | 3 |
| p3 | 3 | 1 | 0 | 2 |
| p4 | 5 | 3 | 2 | 0 |

**Distance Matrix**
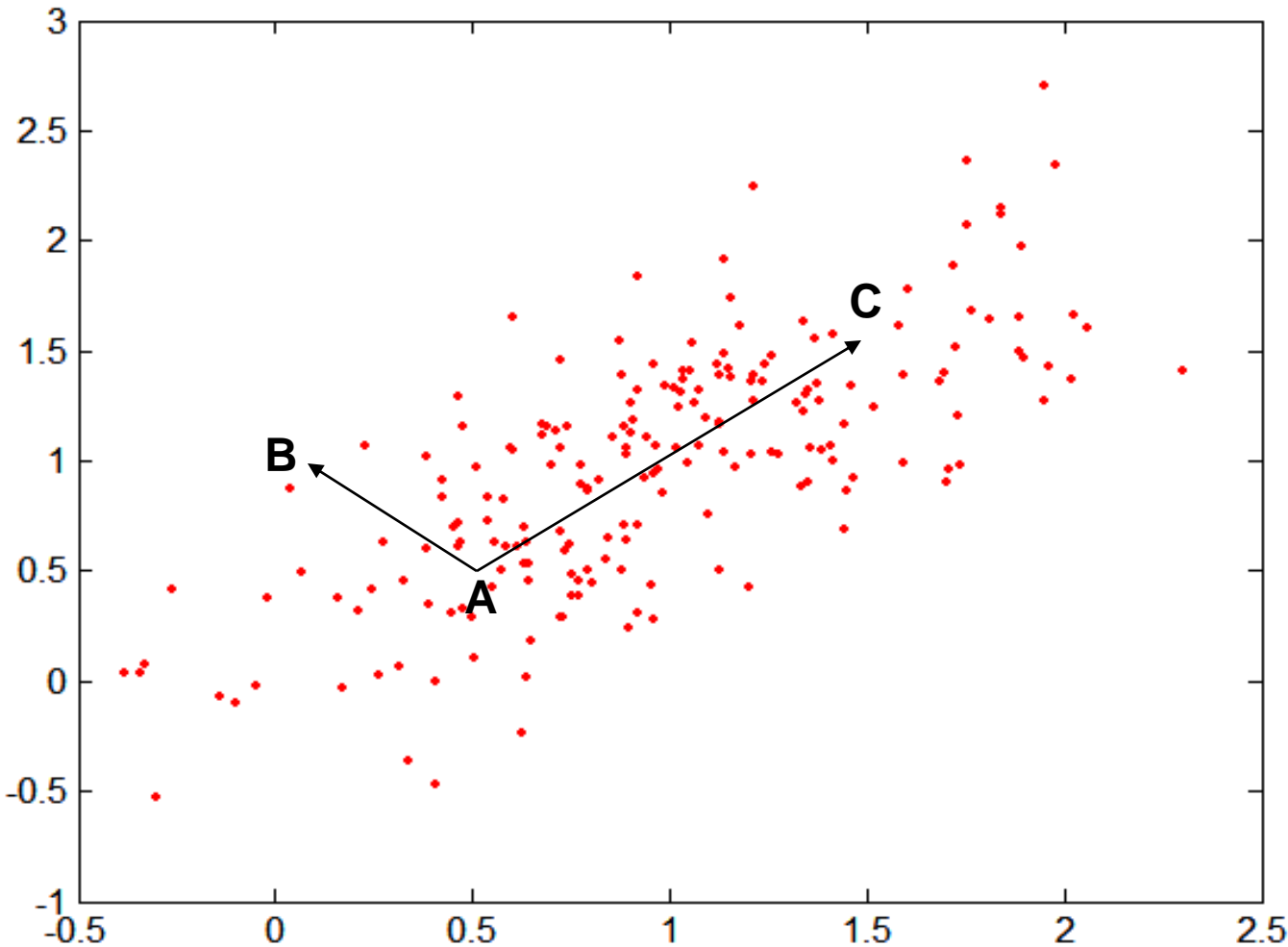
# Mahalanobis Distance

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



**$\Sigma$ 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**

# 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 11 matches / number of non-zero attributes

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

# SMC versus Jaccard: Example

$\mathbf{x} =$ 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 $p$ was 0 and $q$ was 1)

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

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

$f_{11} = 0$ (the number of attributes where $p$ was 1 and $q$ 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 $\mathbf{d}_1$ and $\mathbf{d}_2$ are two document vectors, then

$$\cos(\mathbf{d}_1, \mathbf{d}_2) = \langle\mathbf{d}_1,\mathbf{d}_2\rangle / \|\mathbf{d}_1\| \|\mathbf{d}_2\| ,$$

where $\langle\mathbf{d}_1,\mathbf{d}_2\rangle$ 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:

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

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

$\langle\mathbf{d}_1, \mathbf{d2}\rangle = 3*1 + 2*0 + 0*0 + 5*0 + 0*0 + 0*0 + 0*0 + 2*1 + 0*0 + 0*2 = 5$

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

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

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

# 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

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

where we are using the following standard statistical notation and definitions

$$\mathrm{covariance}(\mathbf{x}, \mathbf{y}) = s_{xy} = \frac{1}{n-1} \sum_{k=1}^{n} (x_k - \overline{x})(y_k - \overline{y}) \quad (2.12)$$
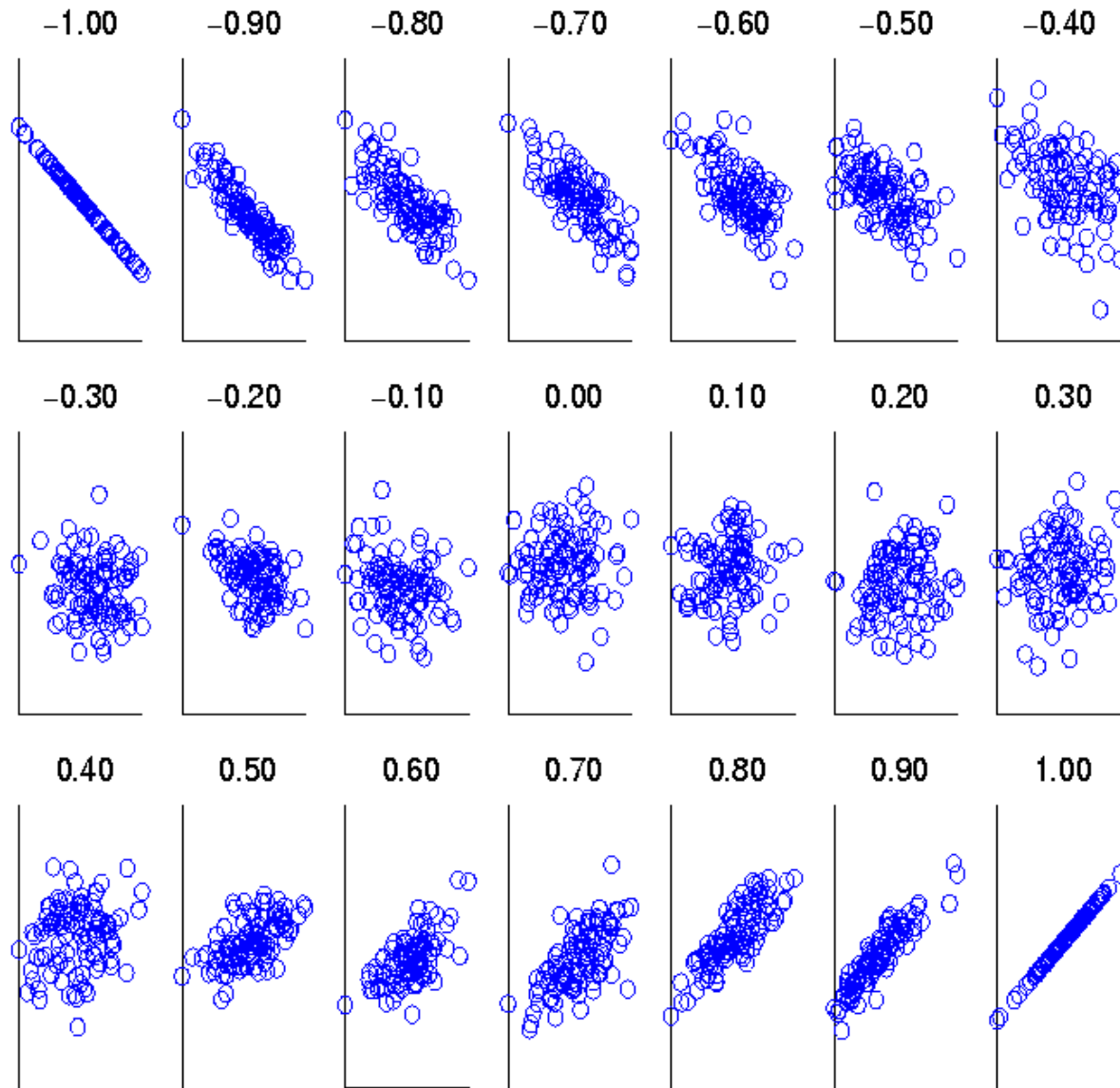
$$\mathrm{standard\_deviation}(\mathbf{x}) \; = \; s_x = \sqrt{\frac{1}{n-1} \sum_{k=1}^{n} (x_k - \overline{x})^2}$$

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

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

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

# 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

# General Approach for Combining Similarities

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

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

2: 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

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

# 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

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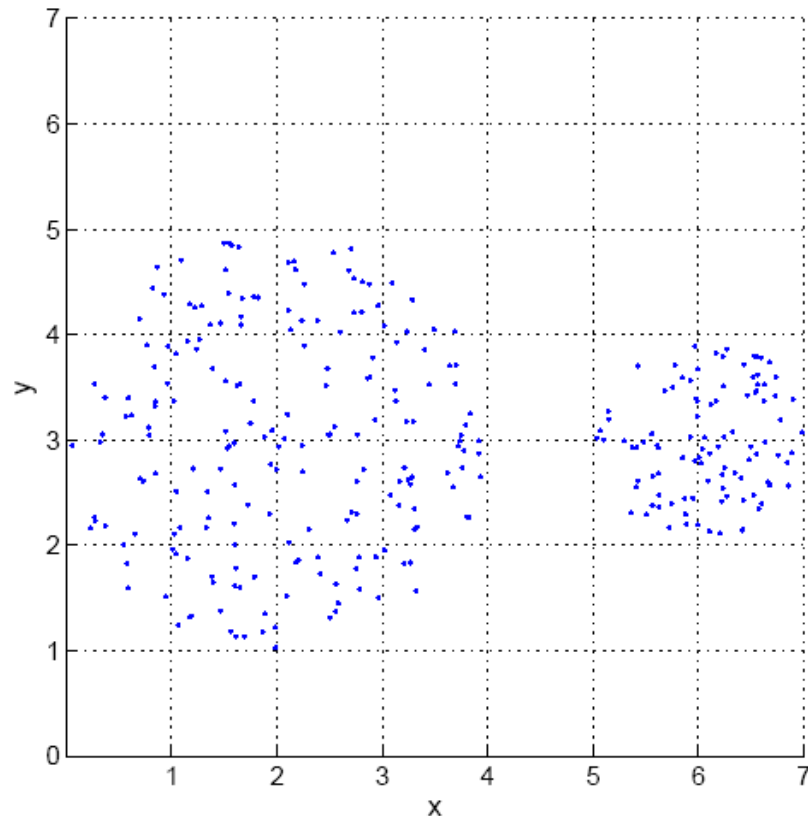
# Density

- Measures the degree to which data objects are close to each other in a specified area
- The notion of density is closely related to that of proximity
- Concept of density is typically used for clustering and anomaly detection
- Examples:
  - Euclidean density
    - Euclidean density = number of points per unit volume
  - Probability density
    - Estimate what the distribution of the data looks like
  - Graph-based density
    - Connectivity

# Euclidean Density: Grid-based Approach

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



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

# Euclidean Density: Center-Based

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



**Illustration of center-based density.**

# 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

- The more certain an outcome, the less information that it contains and vice-versa
  - 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

# 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 $\log_2 5 = 2.3219$

# Entropy for Sample Data

- Suppose we have
  - a number of observations ($m$) of some attribute, $X$, e.g., the hair color of students in the class,
  - where there are $n$ different possible values
  - And the number of observation in the $i$th category is $m_i$
  - Then, for this sample

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

- For continuous data, the calculation is harder

# Mutual Information

- Information one variable provides about another

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

$H(X,Y)$ is the joint entropy of $X$ and Y,

$$H(X,Y) = -\sum_i \sum_j p_{ij}\log_2 p_{ij}$$

Where $p_{ij}$ is the probability that the $i^{\text{th}}$ value of $X$ and the $j^{\text{th}}$ value of $Y$ occur together

- For discrete variables, this is easy to compute

- Maximum mutual information for discrete variables is $\log_2(\min(n_X, n_Y))$, where $n_X$ ($n_Y$) is the number of values of $X$ ($Y$)

# Mutual Information Example

| Student Status | Count | $p$ | $-p\log_2 p$ |
|---|---|---|---|
| Undergrad | 45 | 0.45 | 0.5184 |
| Grad | 55 | 0.55 | 0.4744 |
| Total | 100 | 1.00 | 0.9928 |

| Grade | Count | $p$ | $-p\log_2 p$ |
|---|---|---|---|
| A | 35 | 0.35 | 0.5301 |
| B | 50 | 0.50 | 0.5000 |
| C | 15 | 0.15 | 0.4105 |
| Total | 100 | 1.00 | 1.4406 |

| Student Status | Grade | Count | $p$ | $-p\log_2 p$ |
|---|---|---|---|---|
| Undergrad | A | 5 | 0.05 | 0.2161 |
| Undergrad | B | 30 | 0.30 | 0.5211 |
| Undergrad | C | 10 | 0.10 | 0.3322 |
| Grad | A | 30 | 0.30 | 0.5211 |
| Grad | B | 20 | 0.20 | 0.4644 |
| Grad | C | 5 | 0.05 | 0.2161 |
| Total | | 100 | 1.00 | 2.2710 |

**Mutual information of Student Status and Grade = 0.9928 + 1.4406 - 2.2710 = 0.1624**

# Maximal Information Coefficient

- Reshef, David N., Yakir A. Reshef, Hilary K. Finucane, Sharon R. Grossman, Gilean McVean, Peter J. Turnbaugh, Eric S. Lander, Michael Mitzenmacher, and Pardis C. Sabeti. "Detecting novel associations in large data sets." *science* 334, no. 6062 (2011): 1518-1524.

- Applies mutual information to two continuous variables

- Consider the possible binnings of the variables into discrete categories
  - $n_X \times n_Y \leq N^{0.6}$ where
    - $n_X$ is the number of values of $X$
    - $n_Y$ is the number of values of $Y$
    - $N$ is the number of samples (observations, data objects)

- Compute the mutual information
  - Normalized by $\log_2(\min(n_X, n_Y)$

- Take the highest value