# Data Mining: Data

## Introduction to Data Mining

### CS412, 2021 Spring
### Liu Zuozhu

# 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

# Attribute Types

- **Nominal:** categories, states, or "names of things"
  - *Hair_color = {auburn, black, blond, brown, grey, red, white}*
  - marital status, occupation, ID numbers, zip codes
- **Binary**
  - Nominal attribute with only 2 states (0 and 1)
  - Symmetric binary: both outcomes equally important
    - ◆ e.g., gender
  - Asymmetric binary: outcomes not equally important.
    - ◆ e.g., medical test (positive vs. negative)
    - ◆ Convention: assign 1 to most important outcome (e.g., HIV positive)
- **Ordinal**
  - Values have a meaningful order (ranking) but magnitude between successive values is not known.
  - *Size = {small, medium, large},* grades, army rankings

# Numeric Attribute Types

- Quantity (integer or real-valued)
- **Interval**
    - ◆ Measured on a scale of **equal-sized units (meaningful differences)**
    - ◆ Values have order
        - – E.g., *temperature in C˚ or F˚, calendar dates*
    - ◆ No true zero-point
- **Ratio**
    - ◆ Inherent **zero-point**
    - ◆ We can speak of values as being an order of magnitude larger than the unit of measurement (10 K˚ is twice as high as 5 K˚).
        - – e.g., *temperature in Kelvin, length, counts, monetary quantities*

# Types of Attributes

- There are different types of attributes
  - Nominal
    - Examples: ID numbers, eye color, zip codes
  - Ordinal
    - Examples: rankings (e.g., taste of potato chips on a scale from 1-10), grades, height {tall, medium, short}
  - Interval
    - Examples: calendar dates, temperatures in Celsius or Fahrenheit.
  - Ratio
    - Examples: temperature in Kelvin, length, counts, elapsed time (e.g., time to run a race)

# Properties of Attribute Values

- The type of an attribute depends on which of the following properties/operations it possesses:

  – Distinctness:             $= \neq$

  – Order:                    $< >$

  – Differences are           $+ -$
    meaningful :

  – Ratios are                 $* /$
    meaningful

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

# Difference Between Ratio and Interval

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

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

| | 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."*

# Critiques of the attribute categorization

- Incomplete
  - Asymmetric binary
  - Cyclical
  - Multivariate
  - Partially ordered
  - Partial membership
  - Relationships between the data

- Real data is approximate and noisy
  - This can complicate recognition of the proper attribute type
  - Treating one attribute type as another may be approximately correct

# 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 (among many possible) 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 present

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

  – In the end, what is meaningful can be specific to domain

# 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

# Outline

- Attributes and Objects

- **Types of Data**

- 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
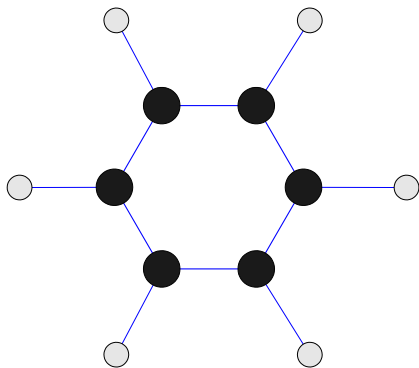
**Book References in Data Mining and Knowledge Discovery**

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

- Genomic sequence data

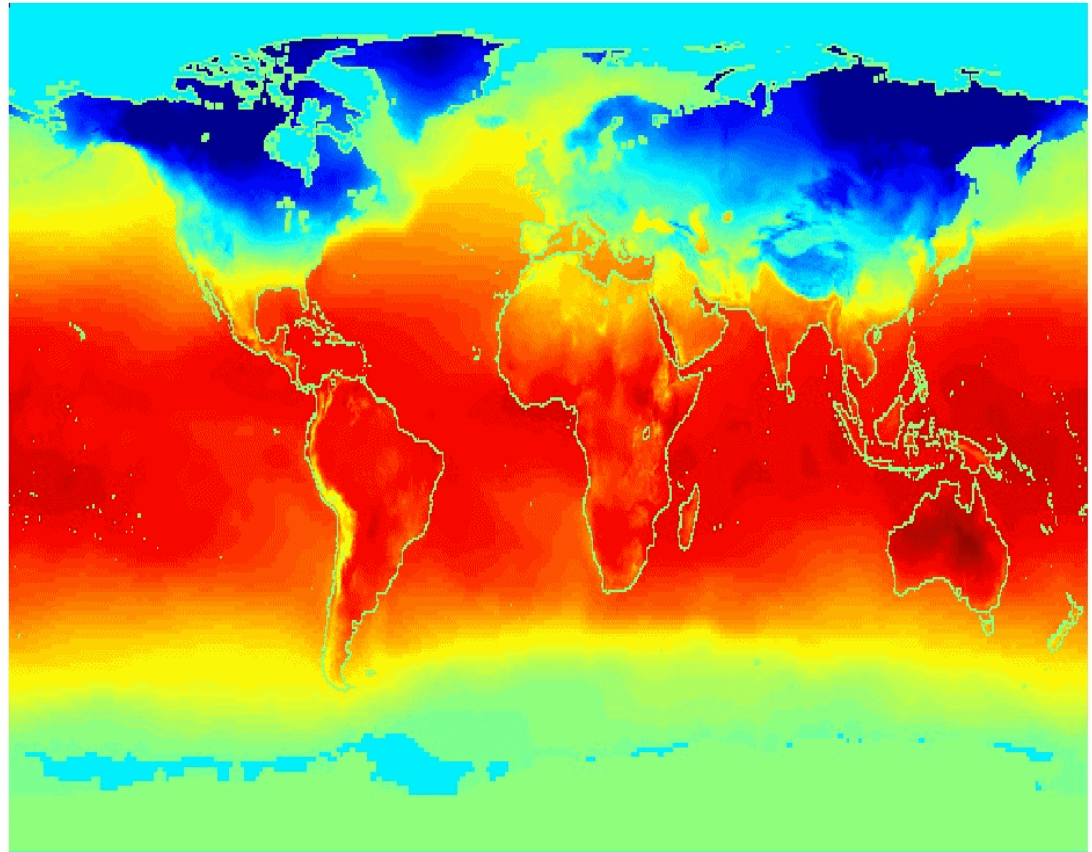GGTTCCGCCTTCAGCCCCGCGCC
CGCAGGGCCCGCCCCGCGCCGTC
GAGAAGGGCCCGCCTGGCGGGCG
GGGGGAGGCGGGGCCGCCCGAGC
CCAACCGAGTCCGACCAGGTGCC
CCCTCTGCTCGGCCTAGACCTGA
GCTCATTAGGCGGCAGCGGACAG
GCCAAGTAGAACACGCGAAGCGC
TGGGCTGCCTGCTGCGACCAGGG

# Ordered Data

● Spatio-Temporal Data

**Average Monthly Temperature of land and ocean**

Jan

# 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

- 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
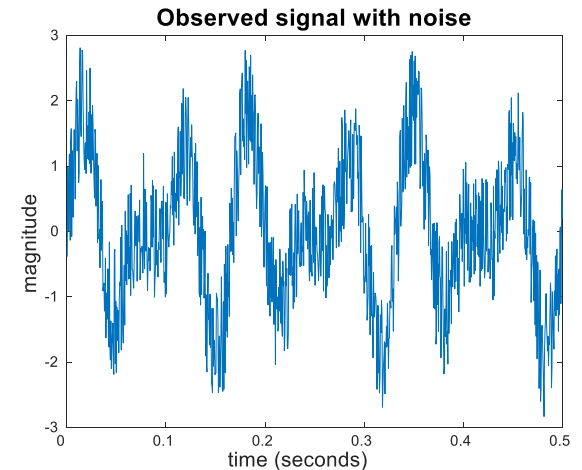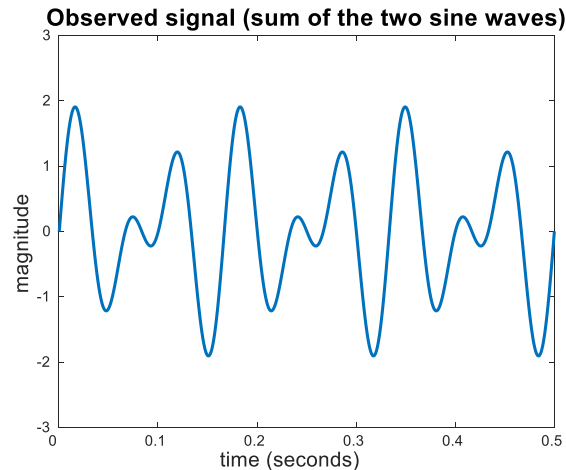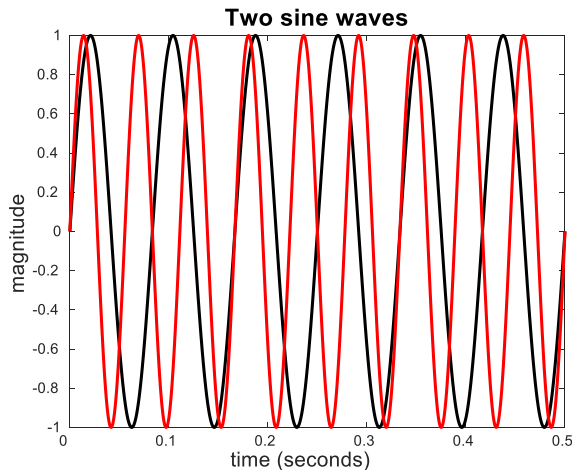    - More loans are given to individuals that default

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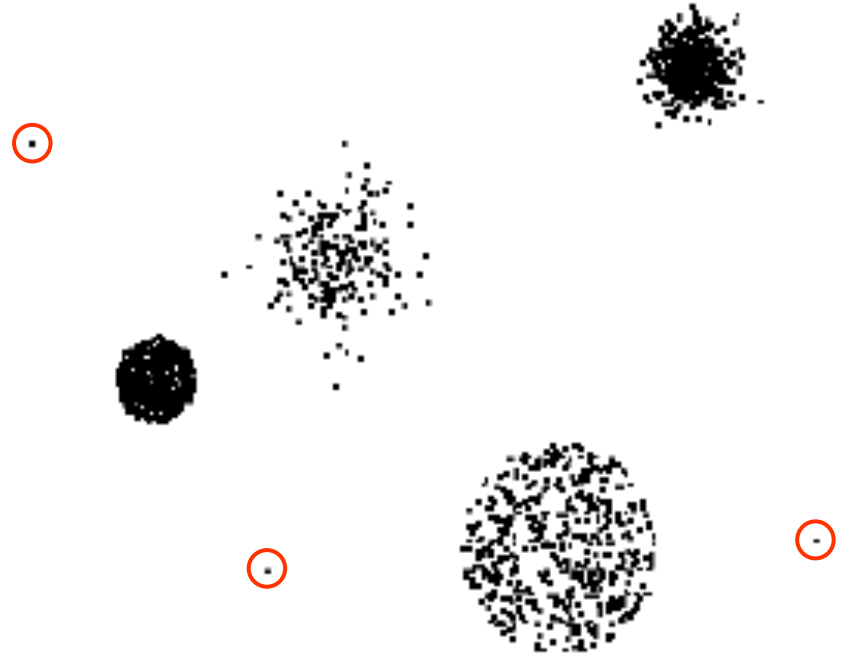  - Missing values
  - Duplicate data

# 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

- Causes?

# Missing Values

- Reasons for missing values
  - Information is not collected
    (e.g., people decline to give their age and weight)
  - Attributes may not be applicable to all cases
    (e.g., annual income is not applicable to children)

- Handling missing values
  - Eliminate data objects or variables
  - Estimate missing values
    - Example: time series of temperature
    - Example: census results
  - Ignore the missing value during analysis

# 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
  - Upper limit varies
- Proximity refers to a similarity or dissimilarity

# Similarity/Dissimilarity for Simple Attributes

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

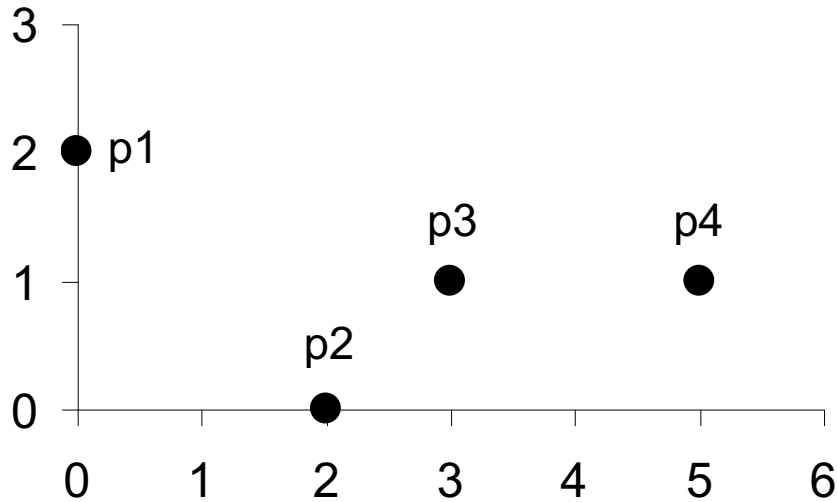| Attribute Type | Dissimilarity | Similarity |
|---|---|---|
| Nominal | $d = \begin{cases} 0 & \text{if } x = y \\ 1 & \text{if } x \neq y \end{cases}$ | $s = \begin{cases} 1 & \text{if } x = y \\ 0 & \text{if } x \neq y \end{cases}$ |
| Ordinal | $d = |x - y|/(n - 1)$ (values mapped to integers $0$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

# Euclidean Distance

- Euclidean Distance

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

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

- Standardization is necessary, if scales differ.

# Euclidean Distance



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

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

**Distance Matrix**

# Minkowski Distance

- Minkowski Distance is a generalization of Euclidean Distance

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

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

# Minkowski Distance: Examples

- *r* = 1. City block (Manhattan, taxicab, $L_1$ norm) distance.
  - A common example of this for binary vectors is the Hamming distance, which is just the number of bits that are different between two binary vectors

- *r* = 2. Euclidean distance

- $r \rightarrow \infty$. **"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 |

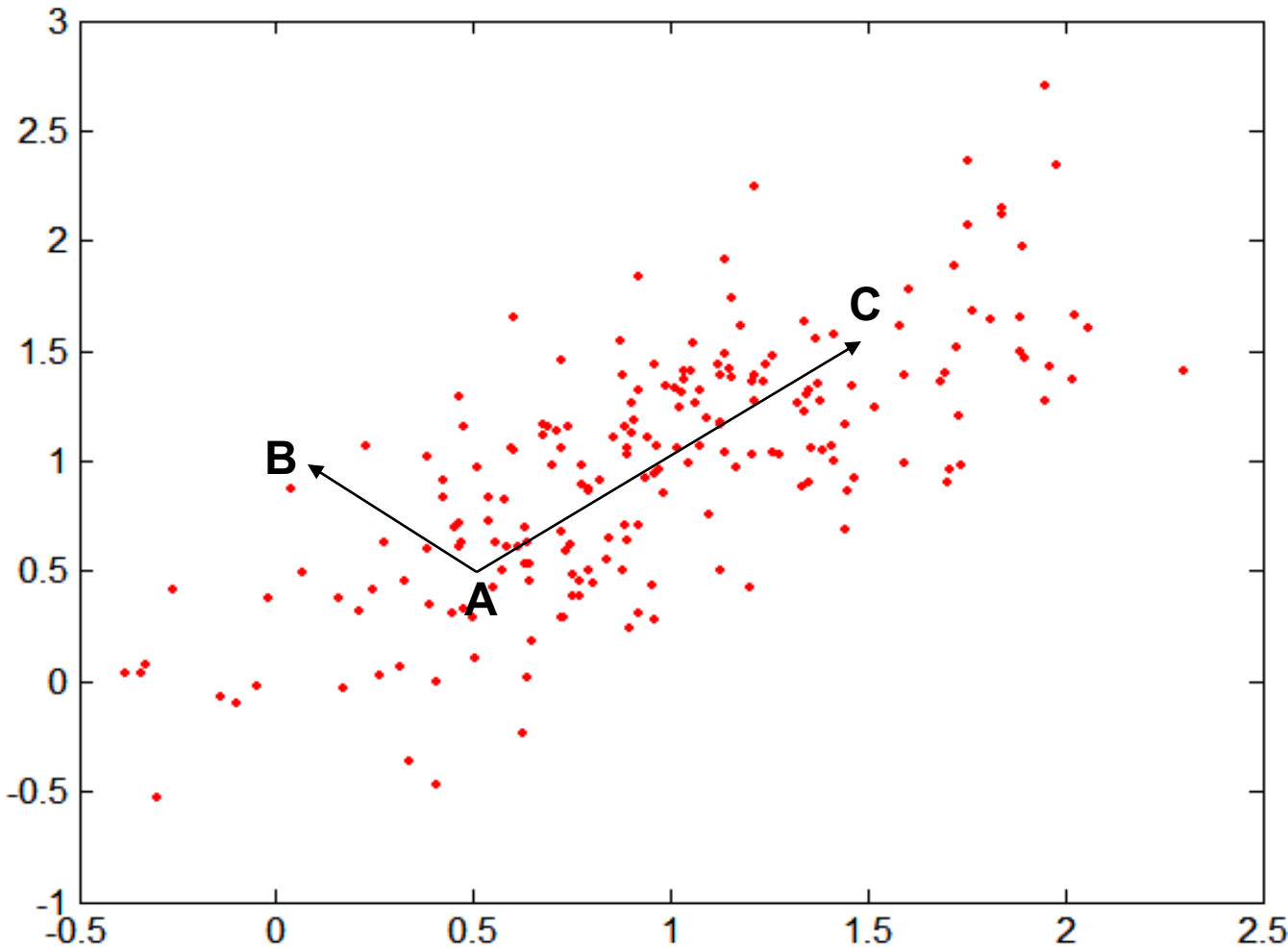**Distance Matrix**

# Mahalanobis Distance

$$\mathbf{mahalanobis}(\mathbf{x}, \mathbf{y}) = ((\mathbf{x} - \mathbf{y})^T \, \Sigma^{-1} (\mathbf{x} - \mathbf{y}))^{-0.5}$$



$\Sigma$ **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 wellknown 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 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 $\mathbf{x} = \mathbf{y}$. (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, $\mathbf{x}$ and $\mathbf{y}$, have only binary attributes

- Compute similarities using the following quantities

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

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

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

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

- Simple Matching and Jaccard Coefficients

  SMC   =  number of matches / number of attributes

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

  J   = number of 11 matches / number of non-zero attributes

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

# 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 $\mathbf{x}$ 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)

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

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

# Cosine Similarity

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

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

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:

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

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

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

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

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

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

# 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

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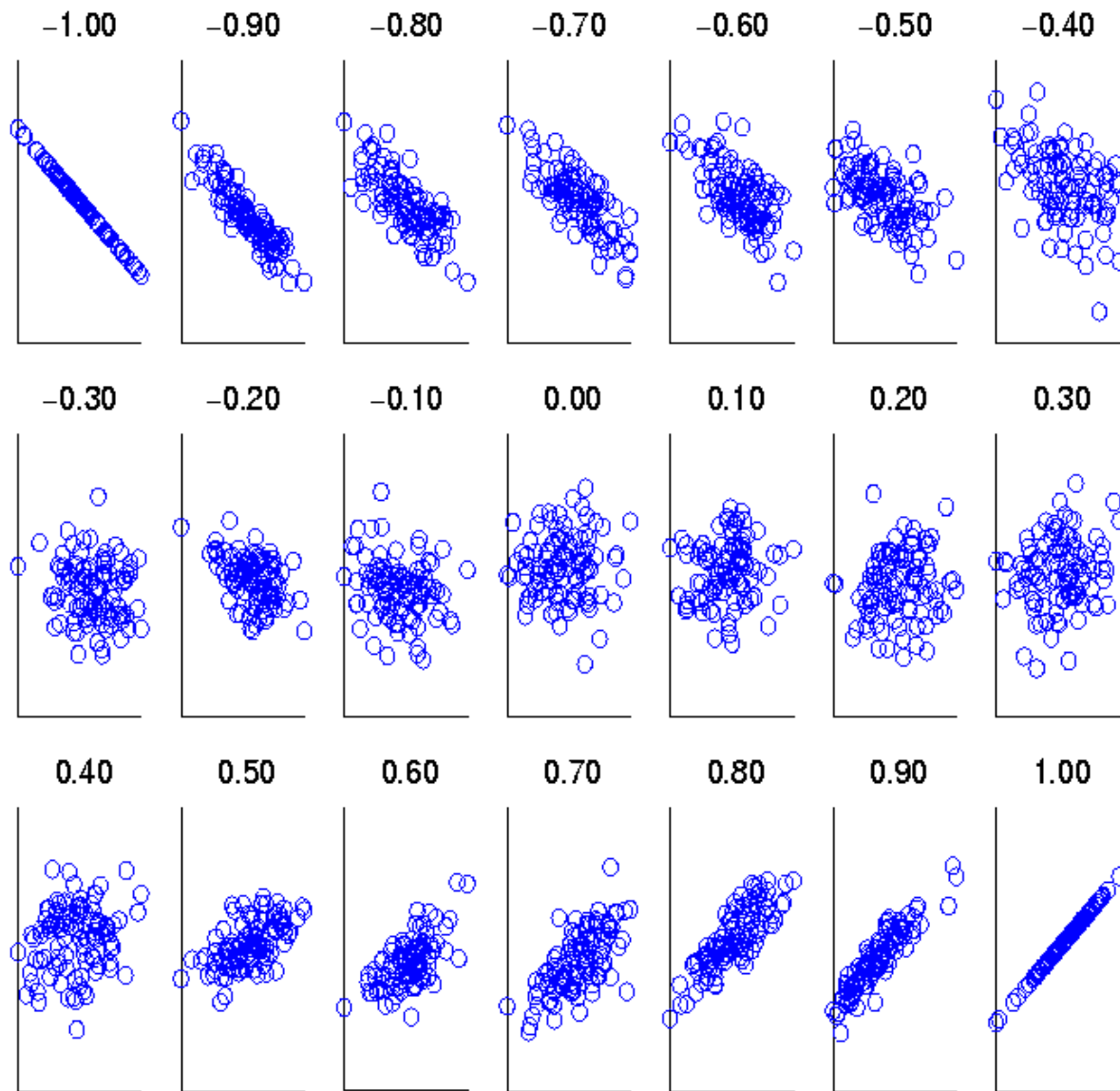
$$\text{standard\_deviation}(\mathbf{x}) = s_x = \sqrt{\frac{1}{n-1} \sum_{k=1}^{n} (x_k - \overline{x})^2}$$

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

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

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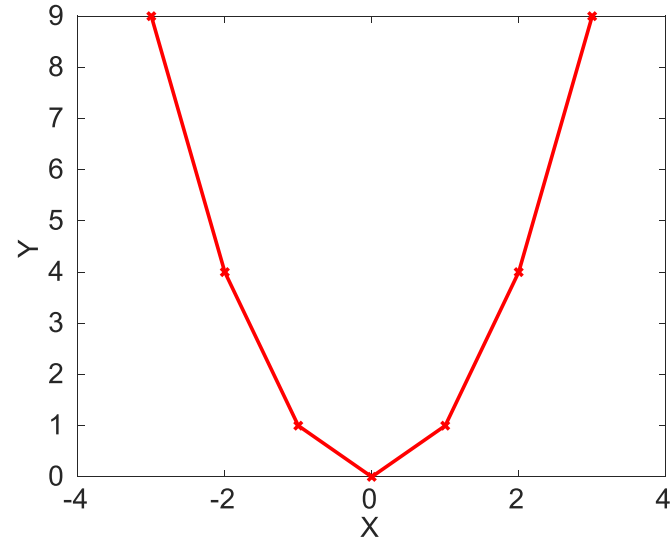
# Visually Evaluating Correlation



Scatter plots showing the similarity from –1 to 1.

# 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

- 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 |

- Consider the example

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

  - $y_s$ = $y$ * 2 (scaled version of y),  $y_t$ = $y$ + 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 |

# 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

  - Comparing the temperature in Celsius of two locations
    - Two locations are considered similar if the temperatures are similar in magnitude

  - 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.

# 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

# Understand Entropy

**Information function**

$$I_A(\omega) = \begin{cases} -\log_2(\mu(A_1)), & if \;\; \omega \in A_1; \\ -\log_2(\mu(A_2)), & if \;\; \omega \in A_2; \\ \vdots \\ -\log_2(\mu(A_n)), & if \;\; \omega \in A_n. \end{cases}$$

One bit of information is equivalent to acquiring an answer to a **binary question**, i.e., to a question asking for a choice between two possibilities.

If the question is not binary, answer **a series of questions**.

The number of questions $N(\omega)$ (i.e., the amount of information in bits) needed to determine the location of the point $\omega$ within the partition may vary from point to point. The best arrangement satisfies

$$N(\omega) \leq I_A(\omega) + 1$$

# Understand Entropy

The information function equals $-\log_2\left(\frac{1}{8}\right) = 3$ on $A_1$ and $A_3$, $-\log_2\left(\frac{1}{4}\right) = 2$ on $A_2$ and $-\log_2\left(\frac{1}{2}\right) = 1$ on $A_4$. The entropy of $\mathcal{A}$ equals

$$H(\mathcal{A}) = \frac{1}{8} \cdot 3 + \frac{1}{4} \cdot 2 + \frac{1}{8} \cdot 3 + \frac{1}{2} \cdot 1 = \frac{7}{4}.$$

The arrangement of questions that optimizes the expected value of the number of questions asked is the following
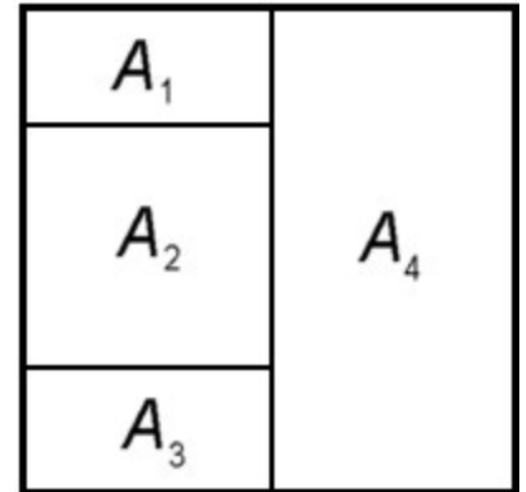
- Question 1. *Are you in the left half?*

The answer *no*, locates $\omega$ in $A_4$ using one bit. Otherwise the next question is:

- Question 2. *Are you in the central square of the left half?*

The *yes* answer locates $\omega$ in $A_2$ using two bits. If not, the last question is:

- Question 3. *Are you in the top half of the whole square?*

Now *yes* and *no* locate $\omega$ in $A_1$ or $A_3$, respectively. This takes three bits.



**Entropy is the expected amount of information needed to locate a point in the partition**

# Understand Entropy 2

Let X be a **random variable** defined on the probability space Ω and assuming values in a finite set {x1,x2,…,xn} .

Suppose an experimenter knows the distribution of X and tries to guess the outcome of X before performing the experiment, i.e., before picking some ω∈Ω and reading the value X(ω)

His **uncertainty** about the outcome is the expected value of the information he is missing to be certain.

Notice that the entropy does not depend on the metric structure of the set of values of X , so entropy cannot be compared to variance. Variance measures a different kind of *uncertainty* of the outcome of a real random variable X , which takes into account the distances between the outcome values.

# Entropy Properties

- The entropy of a partition is nonnegative and equal to zero if and only if one of the elements $A_i$ of the partition has measure 1 (and all other elements have measure zero).

- The entropy of a partition into n sets is highest for the measure which assigns equal values 1n to these sets. The entropy then equals log2n .

# 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**

# Understand Mutual Information

$$H(X) = -\sum_x P_X(x) \log P_X(x) = -E_{P_X} \log P_X$$

**Conditional entropy** is the average uncertainty about $X$ after observing a second random variable $Y$

$$H(X|Y) = \sum_y P_Y(y) \left[ -\sum_x P_{X|Y}(x|y) \log\left(P_{X|Y}(x|y)\right) \right] = E_{P_Y}\left[ -E_{P_{X|Y}} \log P_{X|Y} \right]$$

$$I(X;Y) = H(X) - H(X|Y).$$

Mutual information: the *reduction* in uncertainty about variable $X$, or the expected reduction in the number of yes/no questions needed to guess $X$ after observing $Y$

# Understand Mutual Information

- $I(X;Y) = H(Y) - H(Y|X)$

- $H(X,Y) = H(X) + H(Y|X) \rightarrow I(X;Y) = H(X) + H(Y) - H(X,Y)$

- $I(X;X) = H(X) - H(X|X) = H(X)$
  Entropy is "self-information"

# Understand Mutual Information

$$D_{KL}\left(P(z)\|Q(z)\right) \equiv \sum_z P(z) \log\left[\frac{P(z)}{Q(z)}\right]$$

$$I(X;Y) = D_{KL}\left(P_{XY}(x,y)\|P_X(x)P_Y(y)\right)$$

$$I(X;Y) = I(Y;X)$$
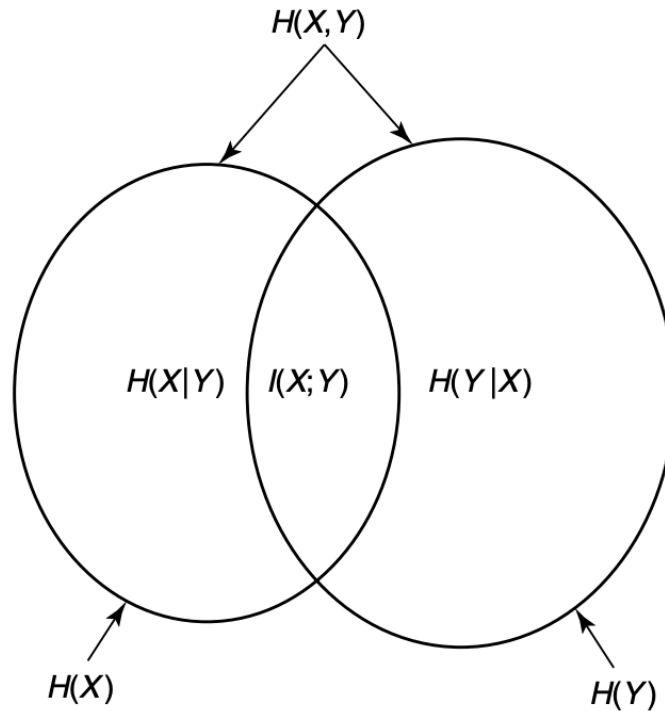
Another way to think about mutual information:

It is a measure of **how close** the true joint distribution of X and Y is to the independent joint distribution.

Because the Kullback-Leibler distance, and thus the mutual information, is zero if and only if P(x,y)=P(x)P(y), it follows that the mutual information captures *all* dependencies between random variables, not just, say, second order ones, as captured by the covariance.

# Understand Mutual Information

**Vien diagram**

$$H(X,Y)$$

$$H(X|Y) \quad I(X;Y) \quad H(Y|X)$$

$$H(X) \qquad H(Y)$$

$I(X;Y)$ is the intersection of information in $X$ with information in $Y$

# Understand Mutual Information

X: blood type

|  | A | B | AB | O |
|---|---|---|---|---|
| **Very Low** | 1/8 | 1/16 | 1/32 | 1/32 |
| **Low** | 1/16 | 1/8 | 1/32 | 1/32 |
| **Medium** | 1/16 | 1/16 | 1/16 | 1/16 |
| **High** | 1/4 | 0 | 0 | 0 |

Y: chance for skin cancer

X: marginal (1/2, 1/4, 1/8, 1/8)

Y: marginal (1/4, 1/4, 1/4, 1/4)

H(X) = 7/4 bits      H(Y) = 2 bits

Conditional entropy: H(X|Y) = 11/8 bits, H(Y|X) = 13/8 bits

$$H(Y|X) \neq H(X|Y)$$

Mutual information: I(X; Y) = H(X) − H(X|Y) = 0.375 bit

# 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

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

# Outline

- Attributes and Objects

- Types of Data

- Data Quality

- Similarity and Distance

- **Data Preprocessing**