

Measuring distance/  
similarity of data objects

# Multiple data types

- Records of users
- Graphs
- Images
- Videos
- Text (webpages, books)
- Strings (DNA sequences)
- Timeseries
- **How do we compare them?**

# Feature space representation

- Usually data objects consist of a set of attributes (also known as **dimensions**)
- J. Smith, 20, 200K
- If all **d** dimensions are **real-valued** then we can **visualize** each data point as points in a **d-dimensional space**
- If all **d** dimensions are **binary** then we can think of each data point as a **binary vector**

# Distance functions

- The distance  $d(x, y)$  between two objects  $x$  and  $y$  is a **metric** if
  - $d(i, j) \geq 0$  (**non-negativity**)
  - $d(i, i) = 0$  (**isolation**)
  - $d(i, j) = d(j, i)$  (**symmetry**)
  - $d(i, j) \leq d(i, h) + d(h, j)$  (**triangular inequality**)
- The definitions of distance functions are usually different for **real**, **boolean**, **categorical**, and **ordinal** variables.
- Weights may be associated with different variables based on applications and data semantics.

# Data Structures

- **data** matrix

$$\begin{array}{c} \text{attributes/dimensions} \\ \left\{ \begin{array}{c} x_{11} \quad \dots \quad x_{1l} \quad \dots \quad x_{1d} \\ \dots \quad \dots \quad \dots \quad \dots \quad \dots \\ x_{i1} \quad \dots \quad x_{il} \quad \dots \quad x_{id} \\ \dots \quad \dots \quad \dots \quad \dots \quad \dots \\ x_{n1} \quad \dots \quad x_{nl} \quad \dots \quad x_{nd} \end{array} \right\} \\ \text{tuples/objects} \end{array}$$

- **Distance** matrix

$$\begin{array}{c} \text{objects} \\ \left\{ \begin{array}{cccc} 0 & & & \\ d(2,1) & 0 & & \\ d(3,1) & d(3,2) & 0 & \\ \vdots & \vdots & \vdots & \\ d(n,1) & d(n,2) & \dots & \dots & 0 \end{array} \right\} \\ \text{objects} \end{array}$$

# Distance functions for real-valued vectors

- $L_p$  norms or **Minkowski** distance:

$$L_p(x, y) = \left( \sum_{i=1}^d |x_i - y_i|^p \right)^{\frac{1}{p}}$$

- $p = 1$ ,  $L_1$ , **Manhattan (or city block)** or **Hamming** distance:

$$L_1(x, y) = \left( \sum_{i=1}^d |x_i - y_i| \right)$$

# Distance functions for real-valued vectors

- $L_p$  norms or **Minkowski** distance:

$$L_p(x, y) = \left( \sum_{i=1}^d |x_i - y_i|^p \right)^{\frac{1}{p}}$$

- $p = 2$ ,  $L_2$ , **Euclidean** distance:

$$L_2(x, y) = \left( \sum_{i=1}^d (x_i - y_i)^2 \right)^{1/2}$$

# Distance functions for real-valued vectors

- Dot product or cosine similarity

$$\cos(x, y) = \frac{x \cdot y}{||x|| ||y||}$$

- Can we construct a distance function out of this?
- When use the one and when the other?



# Hamming distance for 0-1 vectors

**x**      0   1   0   0   1   0   0   1   0

**y**      1   0   0   0   0   1   0   1   1

$$L_1(x, y) = \left( \sum_{i=1}^d |x_i - y_i| \right)$$

# Hamming distance for 0-1 vectors

$x$	0	1	0	0	1	0	0	1	0
$y$	1	0	0	0	0	1	0	1	1

$$L_1(x, y) = \left( \sum_{i=1}^d |x_i - y_i| \right)$$

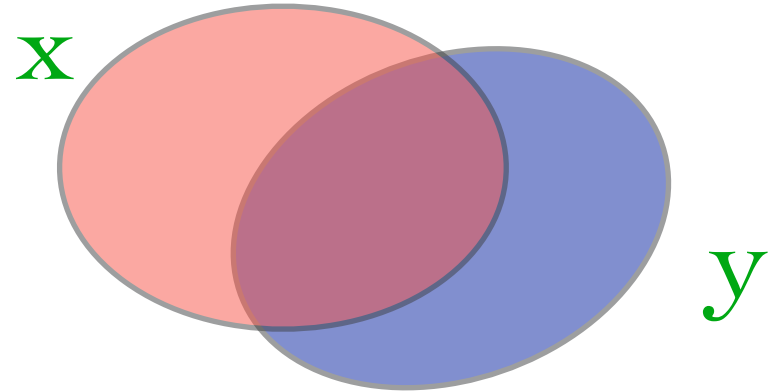
# How good is Hamming distance for 0-1 vectors?

- **Drawback**
- Documents represented as sets (of words)
- Two cases
  - Two **very large** documents -- almost identical -- but for 5 terms
  - Two **very small** documents, with 5 terms each, disjoint

# Distance functions for binary vectors or sets

- **Jaccard** similarity between binary vectors  $x$  and  $y$   
(Range?)

$$\text{JSim}(x, y) = \frac{|x \cap y|}{|x \cup y|}$$

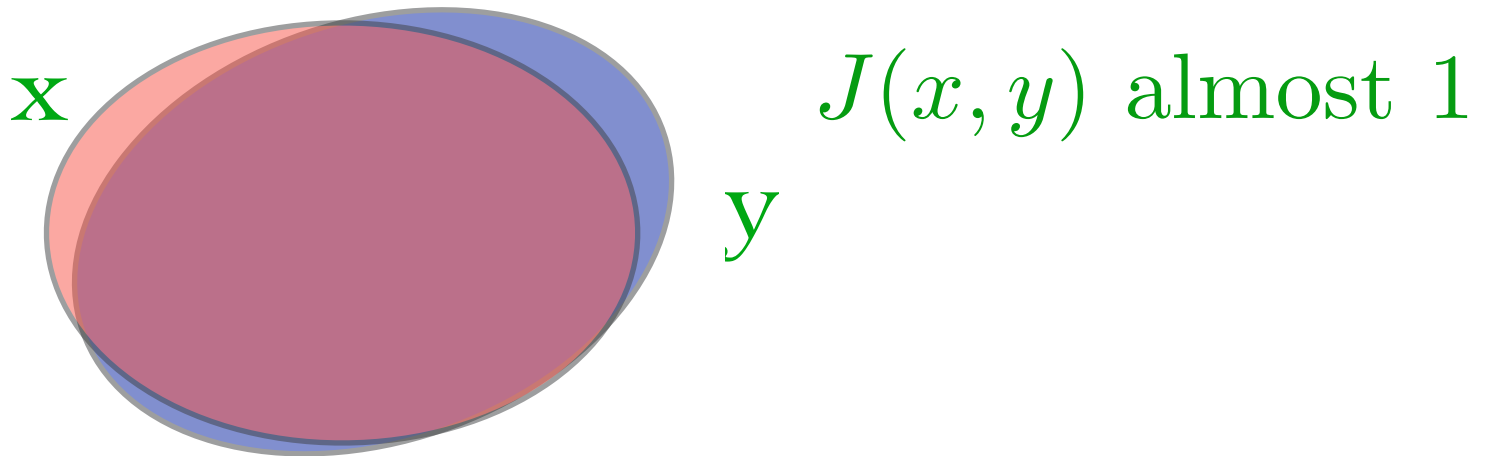


- **Jaccard** distance (Range?):

$$\text{JDist}(x, y) = 1 - \frac{|x \cap y|}{|x \cup y|}$$

# The previous example

- Case 1 (very large almost identical documents)



- Case 2 (small disjoint documents)



# Jaccard similarity/distance

- Example:
  - JSim =  $1/6$
  - Jdist =  $5/6$

	Q1	Q2	Q3	Q4	Q5	Q6
X	1	0	0	1	1	1
Y	0	1	1	0	1	0