

2-data

data = data objects + their attributes

attribute = property / characteristics of an object

object = collection of attributes

attribute values = numbers/symbols assigned to an attribute for a particular object

↳ same attribute can be mapped to different attribute values (ex: height in feet or meters)

	attribute type	description	operations	examples	transformation
categorical (qualitative)	nominal	labeled	$=, \neq$	color of..., types of..., zip codes, to no, list of popular...	any permutation of vals (assign all again)
	ordinal	labeled + order	$= \neq$ $> <$	street numbers, school letter grades, political orientation	new.val = f(old.val) f → monotonic increase (preserving order)
numerical (quantitative)	interval	labeled + order + equal interval	$= \neq$ $> <$ $+ -$	date, temp (in celsius / fahrenheit) test scores, time in clock, iq	new.v = a * old.v + b
	ratio	labeled + order + equal interval + true zero	$= \neq$ $> <$ $+ -$ $\times \div$	weight, height, age, time it takes, income, temp in kelvin	new.val = a * old.val

discrete attributes = finite/continuously infinite set of val → counts, zip codes, binary var

continuous attributes = real number val → temp, height, weight

symmetric attributes = all its vals are equally valuable

asymmetric attributes = only presence (non-zero val) is regarded as important

↳ even the group does not include 0, it is still ratio → ages of the people in the university

important characteristic of data

↳ dimensionality = number of attributes

↳ resolution = pattern depends on the scale

↳ sparsity = only presence counts

↳ size

types of dataset

↳ record data = each object has fixed set of attributes

↳ data matrix = if data objects have same fixed set of numeric attributes

↳ document data = each object is a component-attribute vector → each document becomes term vector
↳ ex: frequency of words in documents

↳ transaction data = each transaction involves set of items → set of purchased products (no quantity only existence)

↳ graph = word wide web, molecular structures

↳ ordered = sequence of transactions → avg monthly temp of world → spatial (related to location) + temporal (rel to time)

data quality problems = noise, outliers, missing values, duplicate data, wrong / fake data

↳ for objects, it is an extraneous object; for attributes, it is modification of original values

outliers case 1 = outliers are noise that interferes with data analysis, ignore/remove

case 2 = outliers are the goal of analysis, ex: cancer cells

similarity - dissimilarity measures = proximity

distance = dissimilarity $s \in [0,1]$ $d \in [0,k]$

↳ nominal if $x=y$ sim=1 dis=0, $x \neq y$ sim=0 dis=1

↳ ordinal dis = $(x-y)/(n-1)$ values mapped to int 0 to $n-1$ sim = $1-d$

↳ interval or ratio $d = |x-y|$ $s = -d, s = 1/(1+d), s = e^{-d}, s = 1 - \frac{d - \min}{\max - \min}$

• minkowski distance = $d(x,y) = \left(\sum_{k=1}^n |x_k - y_k|^r \right)^{1/r}$ r = parameter
 n = number of attributes

↳ $r=1$ → manhattan dist → ex: used to find the distance between binary vectors → hamming distance

↳ $r=2$ → euclidean dist } n can be anything, n and r are independent

↳ $r=\infty$ → supremum (max dist between any component of the vectors)

• mahalanobis distance = difference between a point and a distribution

common properties of distance metrics

↳ $d(x,y) \geq 0$, $d(x,y)=0 \Leftrightarrow x=y$

↳ $d(x,y) = d(y,x)$

↳ $d(x,z) \leq d(x,y) + d(y,z)$

common properties of similarity metrics

↳ $s(x,y) = 1$ only if $x=y$ (cassine does not hold)

↳ $s(x,y) = s(y,x)$

similarity between binary vectors

↳ f_{00} : # att $x=0, y=0$ SMC: simple matching = # matches / # attributes = $(f_{11} + f_{00}) / (f_{01} + f_{10} + f_{11} + f_{00})$

f_{10} : # att $x=1, y=0$ J: Jaccard coefficients = # f_{11} / # non-zero attributes = $f_{11} / (f_{01} + f_{10} + f_{11})$

↳ focusing only existence of attributes

$x = 1000000000$ $f_{00}=7$ $f_{11}=0$ $f_{10}=1$ $f_{01}=2$ SMC = $7/10$ J = 0 → more meaningful mostly

• cosine similarity = $\cos(d_1, d_2) = \langle d_1, d_2 \rangle / \|d_1\| \cdot \|d_2\|$

↳ $d_1 = 3205000200$ $d_2 = 100000102$ $\langle d_1, d_2 \rangle = 3 \cdot 1 + 2 \cdot 0 + \dots + 0 \cdot 2 + 5 = 5$
 $\|d_1\| = \sqrt{3^2 + 2^2 + \dots + 0^2} = 6.481$ $\|d_2\| = 2.459$ $\cos(d_1, d_2) = 0.315$

$\text{corr}(x,y) = \frac{\text{covariance}(x,y)}{\text{standard_deviation}(x) \cdot \text{standard_deviation}(y)} = \frac{s_{xy}}{s_x s_y}$ (2.11)

where we are using the following standard statistical notation and definitions

$\text{covariance}(x,y) = s_{xy} = \frac{1}{n-1} \sum_{i=1}^n (x_i - \bar{x})(y_i - \bar{y})$ (2.12)

$\text{standard_deviation}(x) = s_x = \sqrt{\frac{1}{n-1} \sum_{i=1}^n (x_i - \bar{x})^2}$

$\text{standard_deviation}(y) = s_y = \sqrt{\frac{1}{n-1} \sum_{i=1}^n (y_i - \bar{y})^2}$

$\bar{x} = \frac{1}{n} \sum_{i=1}^n x_i$ is the mean of x

$\bar{y} = \frac{1}{n} \sum_{i=1}^n y_i$ is the mean of y

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

↳ not affected

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

$\rightarrow y = y * 2$ (scaled version of y). $y_1, x, y + 5$ (translated version)

↳ since cos is angle btw vectors, scaling does not affect

Measure	$\langle x, y \rangle$	$\langle x, y_2 \rangle$	$\langle x, y_3 \rangle$
Cosine	0.9667	0.9667	0.7940
Correlation	0.9429	0.9429	0.9429
Euclidean Distance	1.4142	5.6310	14.2127

information based measures → information relates possible outcomes of an event

↳ the more certain an outcome, the less information that it contains

↳ throwing a dice (1/6) contains more information than flipping a coin (1/2)

Entropy = between 0 and $\log_2 n$

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

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

ex: coin → $H = -0.5 \log(0.5) - 0.5 \log(0.5) = 1$

max entropy = $\log_2 n$ → when all outcomes have equal probability

mutual information = $\max \log_2 (\min(n_x, n_y))$

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

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

$$H(X,Y) = - \sum_{i,j} p_{ij} \log_2 p_{ij}$$

where p_{ij} is the probability that the P -value of X and the P -value of Y occur together

maximal information coefficient = to compute the

mutual info of two continuous variables

general approach for combining similarities

• $s_k(x,y) \in [0,1]$

• $s_k = 0$ if k^{th} attribute is asymmetric, or one of the objects has a missing val for k^{th} att } indicator variable

↳ $s_k = 1$ else

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

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

weighted distance

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

↳ weights must be non-negative