

# L4: Distances

Aug 27, 2025  
Data Mining

Jeff M. Phillips

Data  $X \subset \mathcal{X} \equiv \mathbb{R}^d$

$X = \{x_1, x_2, \dots, x_n\}$

$x_i = (x_{i1}, x_{i2}, \dots, x_{id})$

input

Distance  $D : \mathcal{X} \times \mathcal{X} \rightarrow \mathbb{R}_{\geq 0}$

$D(a, b)$  small

$\vdash$

$a, b$

close

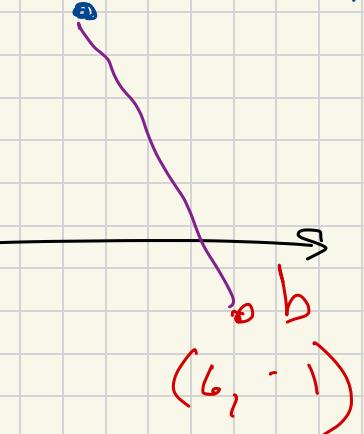
$D(a, b)$  large

$\vdash$

$a, b$

far

$$a = (3, 5)$$



$$D_{Euc}(a, b) = \sqrt{(a_1 - b_1)^2 + (a_2 - b_2)^2}$$

$$(3 - 6)^2 + (5 - (-1))^2 = \sqrt{45}$$

# Metric Distances

(M1) Positivity  $D(a, b) \geq 0$

(M2) Identity  $D(a, b) = 0$  if and only if  $a = b$

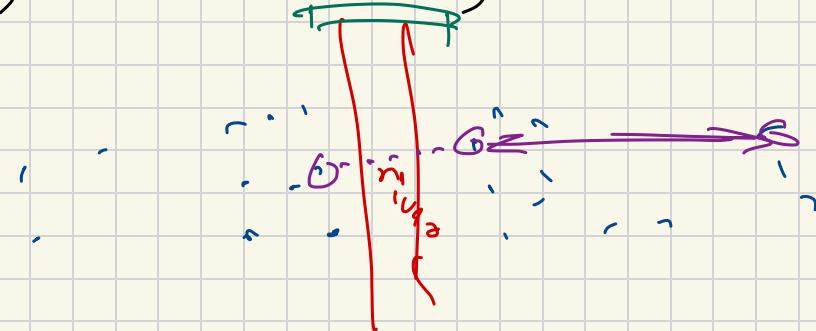
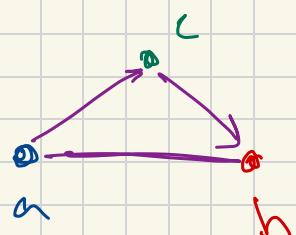
(M3) Symmetry  $D(a, b) = D(b, a)$

(M4) triangle inequality

$$D(a, b) \leq D(a, c) + D(c, b)$$

pseudometric

else  $\Rightarrow$  quasimetric



L<sub>p</sub> - Distances

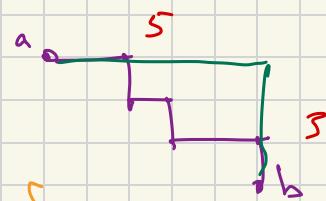
D<sub>p</sub>(a, b)

a, b ∈ ℝ<sup>d</sup>

$$D_p(a, b) = \left( \sum_{i=1}^d |a_i - b_i|^p \right)^{1/p} = \|a - b\|_p$$

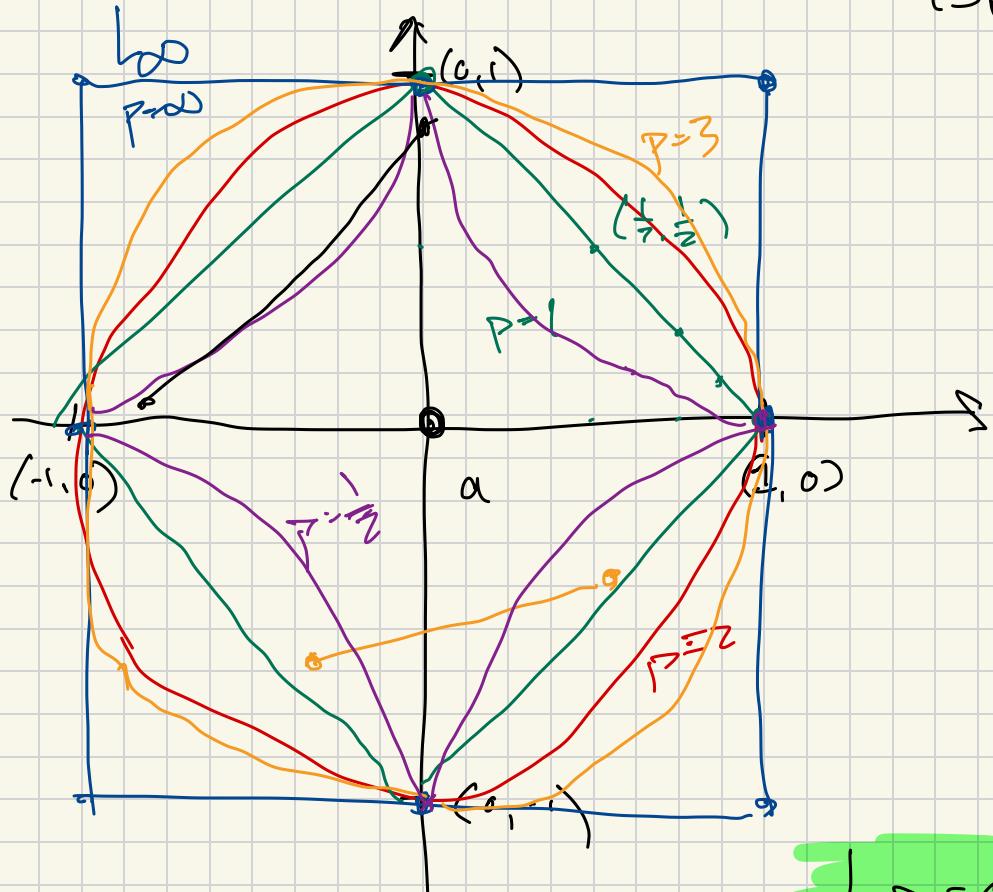
$$D_2(a, b) = \sqrt{\sum_{i=1}^d (a_i - b_i)^2} = \|a - b\|_2 = \|a - b\| = D_{Euc}(a, b)$$

$$D_1(a, b) = \sum_{i=1}^d |a_i - b_i|$$



$$D_\infty(a, b) = \lim_{p \rightarrow \infty} \left( \sum_{i=1}^d |a_i - b_i|^p \right)^{1/p}$$

$$\max_{i \in [d]} |a_i - b_i| = \|a - b\|_\infty$$



Draw all points b.c.r  
 $D_p(a, b) = 1$

$L_p$ -ball  $\subset L_{p'}\text{-ball}$   
subset

if  $p \leq p'$

$p = 1/2 ?$

not convex

$L_p$ -dist  $\Rightarrow$  metric  
 $p \in [1, \infty) + \infty$

# NEW CUYAMA

Population

562

Ft. above sea level

2150

Established

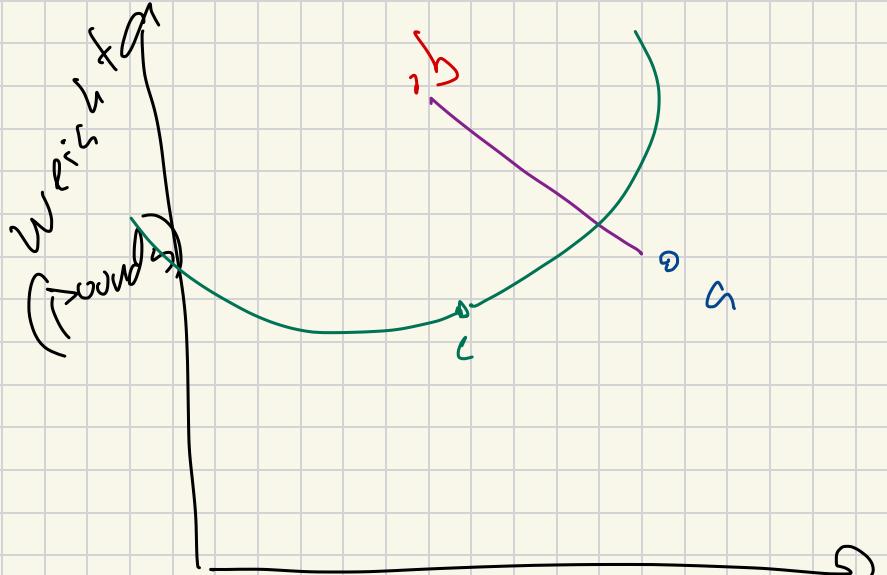
1951

TOTAL

---

4663

Don't do this!



$$D(a, b) = \sqrt{\sum_{i=1}^n (a_i - b_i)^2}$$

Do not  
do  
this!

height  
inches  
meters

## Mahalanobis Distance

$$D_M(a, b) = \sqrt{(a-b)^T M (a-b)}$$

if  $M = \Sigma = \begin{bmatrix} 1 & c \\ c & 1 \end{bmatrix}$

$$= \sqrt{(a-b)^T (a-b)}$$

$$= \sqrt{\langle a-b, a-b \rangle}$$

$$= \sqrt{\sum_{i=1}^d (a_i - b_i)(a_i - b_i)}$$

$$= \sqrt{\sum_{i=1}^d (a_i - b_i)^2}$$

$$= \sqrt{\sum_{i=1}^d (a_i - b_i)^2}$$

$M \in \mathbb{R}^{d \times d}$

positive definite

$\rightarrow$  metric

## Cosine Distance

$$D_{\cos} = 1 - \frac{\langle a, b \rangle}{\|a\| \cdot \|b\|} \in [0, 2]$$

## Metric

(M1) positivity

(M2) identity  $D(a, b) = 0 \text{ iff } a = b$

but if can if  $\mathcal{X} = \mathbb{S}^{d-1}$

(M3) symmetry

(M4) triangle.

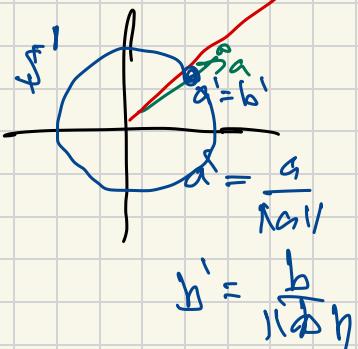
NO,  $\Rightarrow \text{Dang}(a, b) = \cos(\langle a, b \rangle)$   
 $= \arccos(\langle a, b \rangle)$

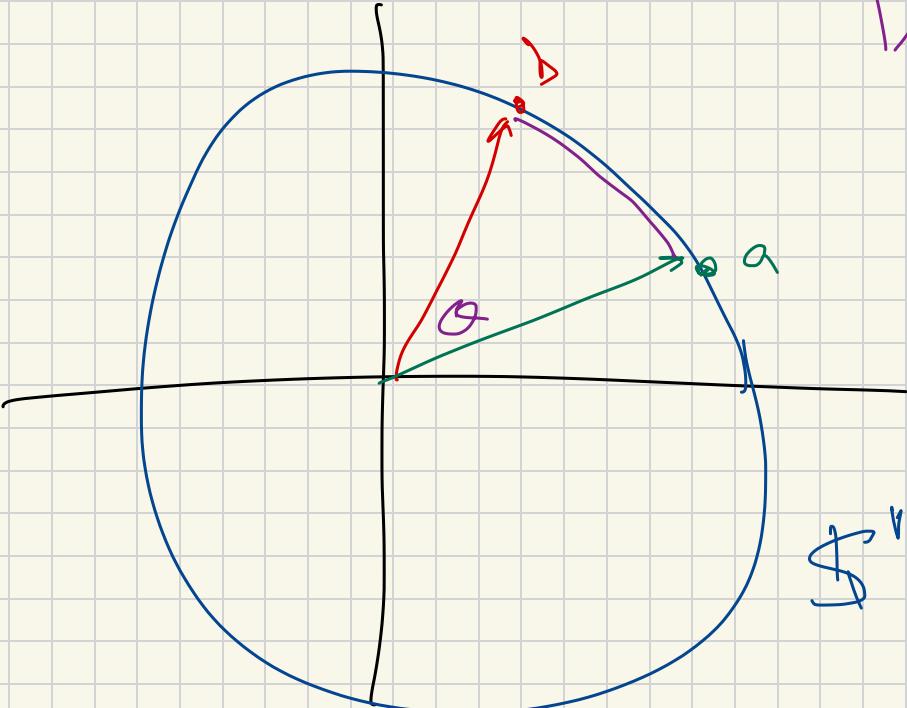
$$a, b \subset \mathbb{R}^d = \mathcal{X}$$

cosine  
similarity

$$S_{\cos}(a, b) = \frac{\langle a, b \rangle}{\|a\| \cdot \|b\|}$$

$$= \left\langle \frac{a}{\|a\|}, \frac{b}{\|b\|} \right\rangle$$





$D_{\text{ang}}(a, b)$

$$= \arccos(\langle a, b \rangle)$$

= angl

between

$\vec{a}, \vec{b}$

in radians

$$\in [0, \pi]$$