

L6: Distances

Jeff M. Phillips

January 27, 2020

Distance : bivariate function

$$d: \mathcal{X} \times \mathcal{X} \rightarrow \mathbb{R}^+ \text{ or } \mathbb{R}_{\geq 0}$$
$$d(a, b) =$$

metric

(M1) $d(a, b) \geq 0$ (non-negative)

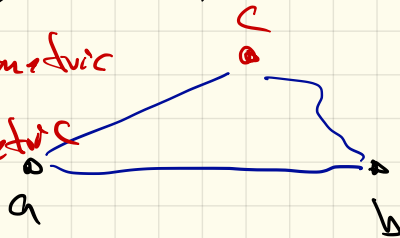
(M2) $d(a, b) = 0$ iff $a = b$ (identity)

(M3) $d(a, b) = d(b, a)$ (symmetry)

(M4) $d(a, b) \leq d(a, c) + d(c, b)$ (triangle inequality)

• M1, M3, M4 pseudometric

• M1, M2, M4 quasimetric



L_p Distances

$$X := \mathbb{R}^d$$

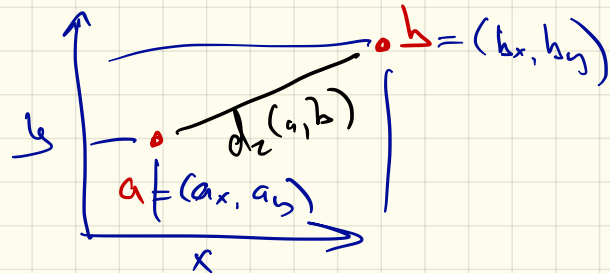
$$a, b \in \mathbb{R}^d$$

$$a = (a_1, a_2, \dots, a_d)$$

$$L_2(a, b) = d_2(a, b) = \|a - b\|_2 = \|a - b\|$$

Euclidean

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

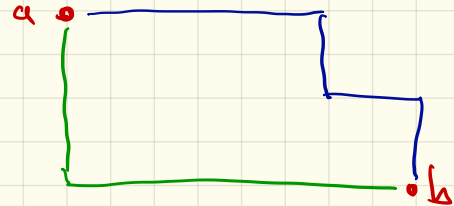


$$L_1(a, b) = d_1(a, b) = \|a - b\|_1$$

$$= \sum_{i=1}^d |a_i - b_i|$$

Manhattan dist.

SLC dist



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

Every L_p dist for $p \in [1, \infty)$
is a metric.

$$L_0 = \|a - b\|_0 = d - \sum_{i=1}^d \mathbb{1}(a_i = b_i)$$

if $a, b \in \{0, 1\}^d$ bit strings

↳ Hamming dist

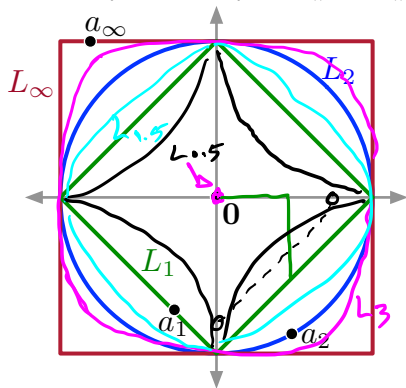
$$\begin{aligned} L_\infty = \|a - b\|_\infty &= \lim_{p \rightarrow \infty} L_p(a, b) \\ &= \max_{i \in \{1, \dots, d\}} |a_i - b_i| \end{aligned}$$

Lp Distances and Unit Balls

For $a = (a_1, a_2, \dots, a_d)$ and $b = (b_1, b_2, \dots, b_d) \in \mathbb{R}^d$

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

Let $b = (0, 0, \dots, 0)$ and $\|a - b\|_p = 1$.



L_{0.5}

Lp Distances and Units

For $a = (a_1, a_2, \dots, a_d)$ and $b = (b_1, b_2, \dots, b_d) \in \mathbb{R}^d$,

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



Rule
All coords
most have
same units.

← non sense

Mahalanobis Dist $a, b \in \mathbb{R}^d$

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

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

M p.d. \rightarrow dm metric

If $M = I$ (identity matrix)

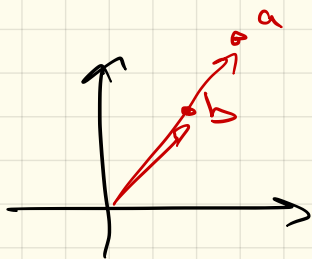
$$d_I(a, b) = \|a - b\|_2$$

$$(a-b)^T (a-b) = \langle a-b, a-b \rangle = \|a-b\|_2^2$$

$$M = \text{diag}(m_1, m_2, \dots, m_d) = \begin{bmatrix} m_1 & & & 0 \\ & m_2 & & \\ & & \dots & \\ 0 & & & m_d \end{bmatrix}$$

Cosine dist $a, b \in \mathbb{R}^d$

$$d_{\cos}(a, b) = 1 - \frac{\langle a, b \rangle}{\|a\| \cdot \|b\|} = 1 - \frac{\sum_{i=1}^d a_i b_i}{\|a\| \cdot \|b\|} \leq 1$$

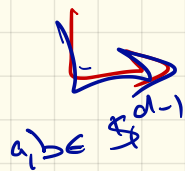
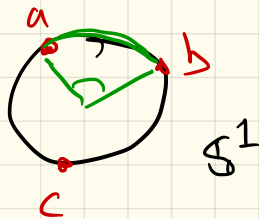


$$d_{\cos}(a, b) = 0 \quad a \neq b$$

only measure direction

$$a \rightarrow \bar{a} = \frac{a}{\|a\|} \quad b \rightarrow \bar{b} = \frac{b}{\|b\|}$$

$$a, b \in \mathbb{S}^{d-1} = \{x \in \mathbb{R}^d \mid \|x\| = 1\}$$



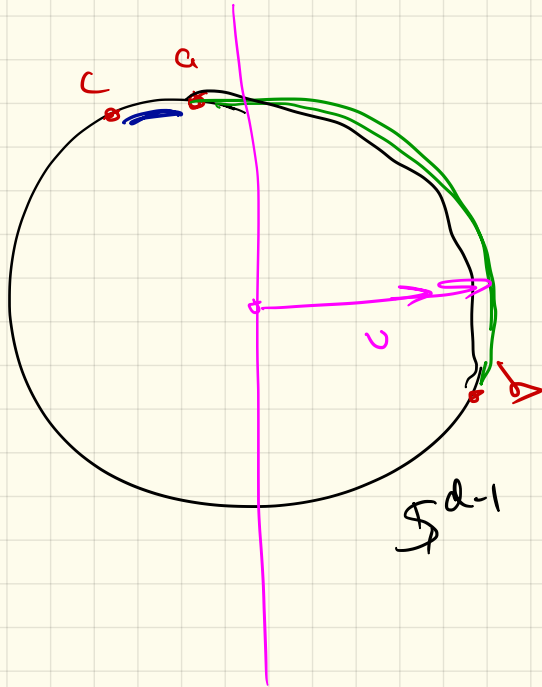
map to angular dist

$$d_{\text{ang}}(a, b) = \cos^{-1}(\langle a, b \rangle) = \arccos(\langle a, b \rangle)$$

metric

Angular dist

LSH



pick random
unit vector
 $u \in \mathbb{S}^{d-1}$

$$h_u(a) = \text{sign}(\langle a, u \rangle) \\ = -1$$

$$h_u(b) = +1$$

LSH

$$\frac{\text{dang}(a, b)}{\pi} = \mathbb{P} \left[\begin{array}{l} h_u(a) \\ \neq \\ h_u(b) \end{array} \right]$$

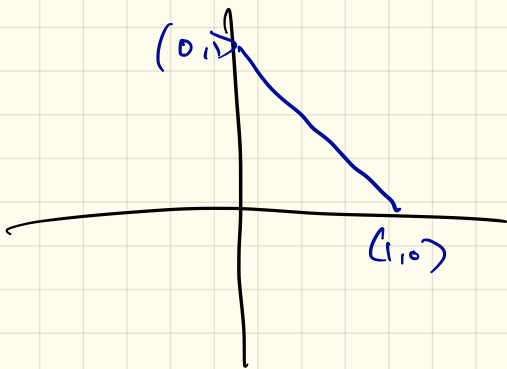
KL Divergence

$$a, b \in \Delta_+^d$$

$$d_{KL}(a \parallel b) = \sum_{i=1}^d a_i \ln(a_i / b_i)$$

$$\Delta^d = \{x \in \mathbb{R}^d \mid \|x\|_1 = 1 \ \& \ \forall_i x_i \geq 0\}$$

$$\Delta_+^d = \{x \in \mathbb{R}^d \mid \|x\|_1 = 1 \ \& \ \forall_i x_i > 0\}$$



Probability distribution