

# Distance

$$D: X \times X \rightarrow \mathbb{R}^+$$

$$D(a,b) \rightarrow [0, \infty)$$

(M1)  $D(a,b) \geq 0$  <sup>Metric</sup>

(non-negativity)

↳ choice

(M2)  $D(a,b) = 0$  iff  $a=b$

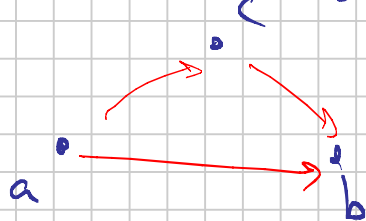
• Unexpected Properties

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

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

M1, M3, M4  $\rightarrow$  pseudo metrics

M1, M2, M4  $\rightarrow$  quasimetrics



## $L_p$ Distance

$$p \in (0, \infty) + \{0\}, \{\infty\}$$

•  $L_2 =$  Euclidean

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

↳ metric

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

$a, b \in \mathbb{R}^d$   
 $a = (a_1, a_2, \dots, a_d)$

(M4)  $\Leftrightarrow p \geq 1$

•  $L_1 =$  Manhattan  
 = SLC Distance

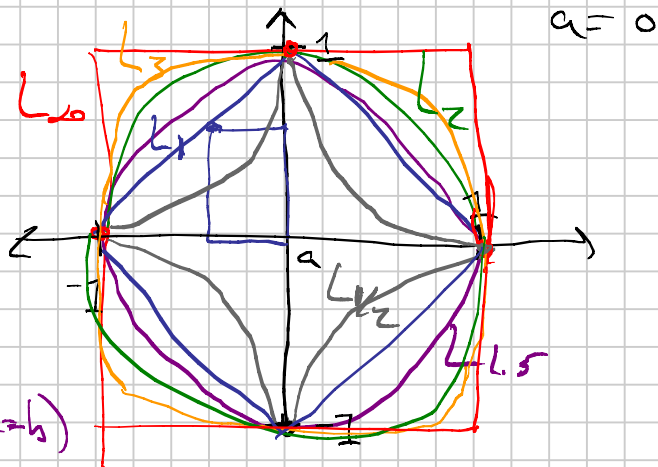
$$D_1(a,b) = \|a-b\|_1 = \sum_{i=1}^d |a_i - b_i|$$

$$(2^{1000} + 1^{1000})^{1/1000}$$

$$\|v\|_1$$

•  $L_{\infty} =$

$$D_{\infty}(a,b) = \max_{i=1..d} |a_i - b_i|$$



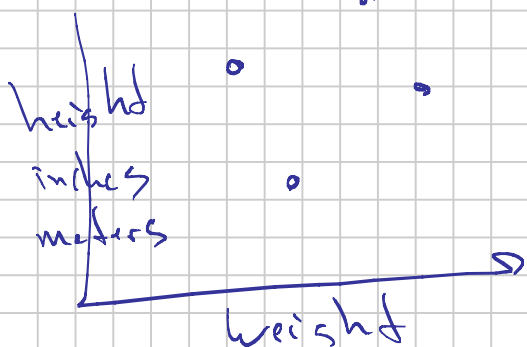
•  $L_0 \approx$  Hamming

$$D_0(a,b) = \|a-b\|_0 = d - \sum_{i=1}^d \mathbb{1}(a_i = b_i)$$

**WARNING**

ONLY use for vectors in  $\mathbb{R}^d$

where all coordinates have same unit.



"Normalize"  $\left\{ \begin{array}{l} \text{all coords in } [0,1] \\ \text{mean} = 0, \text{ var} = 1 \end{array} \right.$

hacks!

Distance Metric Learning (MDS)

Jaccard Distance  $D_J(A,B) = 1 - \frac{|A \cap B|}{|A \cup B|} \in [0,1]$

metric

Triangle Inequality

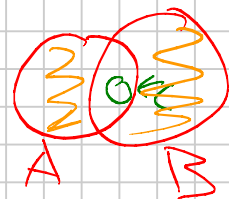
Another set C

$$D_J(A,B) \leq D_J(A,C) + D_J(C,B)$$

$$= \frac{|A \setminus C|}{|A|} + \frac{|B \setminus C|}{|B|}$$

$$\geq \frac{|A \setminus C| + |B \setminus C|}{|A \cup B|}$$

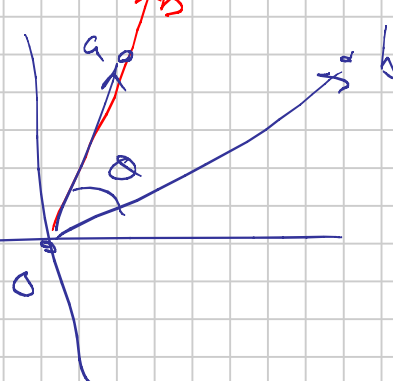
$$\geq \frac{|A \Delta B|}{|A \cup B|} = D_J(A,B)$$



Assume  $C \subset A, C \subset B$   
 $C \subset A \cap B$

# Cosine Distance

$$D_{\cos}(a, b) = 1 - \frac{|\langle a, b \rangle|}{\|a\|_2 \|b\|_2} = 1 - \frac{|\sum_{i=1}^d a_i b_i|}{\|a\|_2 \|b\|_2}$$



↳ "Bag-of-words" data Normalizing

the and up it  $(0, 1, 2, 0)$  ← "up and up"

metric

(M1)  $\rightarrow \mathbb{R}^+$

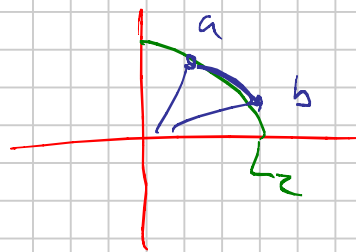
$a = (1, 0, 0, \dots) \Rightarrow a' = \frac{a}{\|a\|}$

$b = (2, 0, 0, \dots) \Rightarrow b' = \frac{b}{\|b\|}$

NOT (M2)  $D_{\cos}(a, b) = 0$  iff  $a = b$

(M3) symmetric

(M4)  $\Delta$  ineq?



# KL Divergence

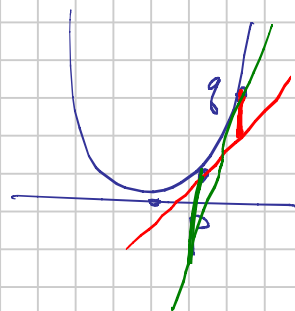
Kullback-Leibler

$P, q$  discrete prob distributions

$P = (P_1, P_2, \dots, P_d)$

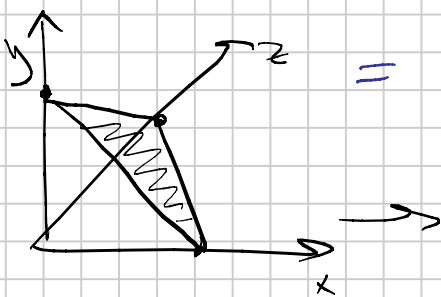
s.t.  $P_i \geq 0$

$\sum_{i=1}^d P_i = 1$



$D_{KL}(P \parallel q)$

$= \sum_{i=1}^d P_i \cdot \ln\left(\frac{P_i}{q_i}\right) = H(P, q) - H(q)$



$\mathbb{R} = \sum_{i=1}^d P_i e_i$

$e_i = (0, \dots, 1, \dots, 0)$