

# Kullback-Leibler Divergence (KL)

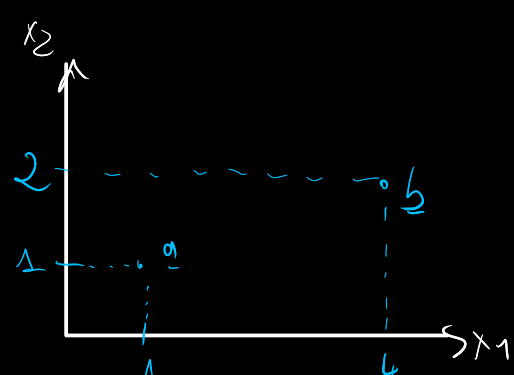
## Introduction

distance between  $a$  &  $b$

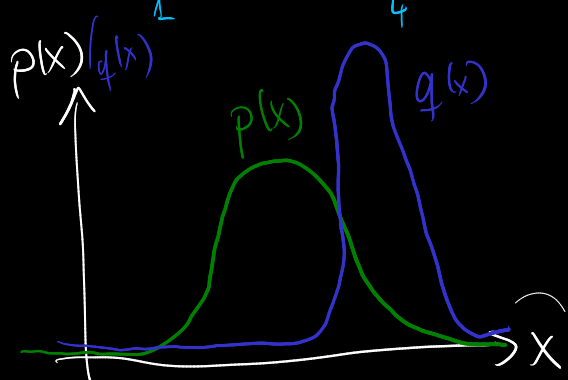


$$d(a,b) = |a-b|$$

distance between  $a$  &  $b$



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



distance between  $p$  &  $q$   
(between two distributions)

↳ KL-Divergence

$$D_{KL}(p \parallel q) = \mathbb{E}_{X \sim p(X)} \left[ \log \left( \frac{p(X)}{q(X)} \right) \right] = \begin{cases} \sum_X p(X) \cdot \log \left( \frac{p(X)}{q(X)} \right) & \text{if } X \text{ is discrete} \\ \int p(X) \cdot \log \left( \frac{p(X)}{q(X)} \right) & \text{if } X \text{ is continuous} \end{cases}$$

KL-Divergence does not have all properties of distance

## Example

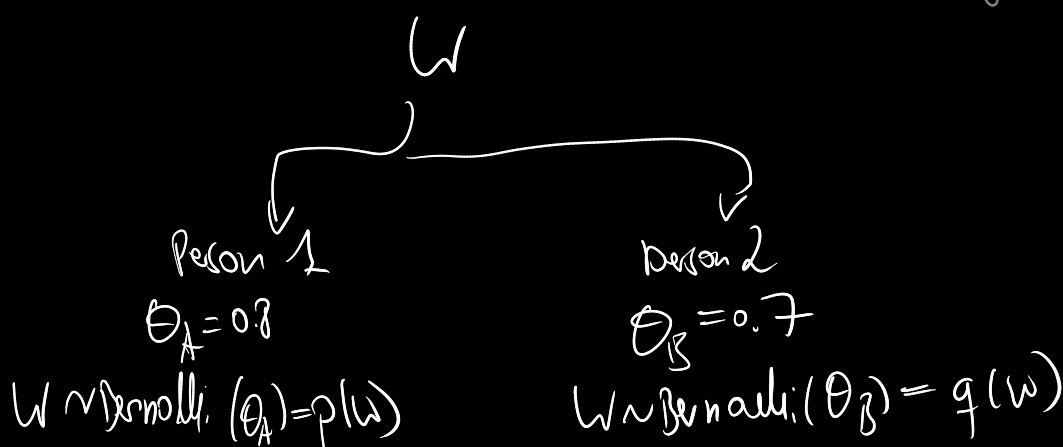
Weather (W)   
 ↗ bad   
 ↘ good

$W \in \{ \text{Bad}, \text{Good} \}$

→ Bernoulli distributions

$$W \sim \text{Bernoulli}(\theta)$$

↑  
prob of good weather



How far apart are the distributions?

$$D_{KL}(p \parallel q) = \sum_W p(W) \log \left( \frac{p(W)}{q(W)} \right) \\ = \sum_{w=0}^1 p(W=w) \log \left( \frac{p(W=w)}{q(W=w)} \right)$$

$$\left[ \text{Bernoulli}(W; \theta) = \theta^W (1-\theta)^{1-W} \right]$$

$$= \underbrace{\theta_A^0 \cdot (1-\theta_A)^{1-0}}_{\theta_A^0} \cdot \log \left( \frac{\theta_A^0 \cdot (1-\theta_A)^{1-0}}{\theta_B^0 \cdot (1-\theta_B)^{1-0}} \right) + \underbrace{\theta_A^1 \cdot (1-\theta_A)^{1-1}}_{\theta_A^1} \cdot \log \left( \frac{\theta_A^1 \cdot (1-\theta_A)^{1-1}}{\theta_B^1 \cdot (1-\theta_B)^{1-1}} \right)$$

$$= (1-\theta_A) \cdot \log \left( \frac{1-\theta_A}{1-\theta_B} \right) + \theta_A \cdot \log \left( \frac{\theta_A}{\theta_B} \right)$$

$$= \underline{\underline{0.0257}}$$