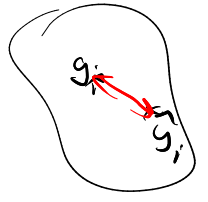


Modeling Tasks & Losses

Choosing physically-meaningful metrics

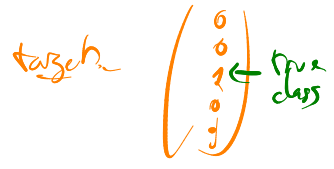


linear regression: $\mathcal{D} = \{(x_i, y_i)\} \rightarrow$

$$\sum_i \|g_i - g_i\|^2$$

MSE

classification: $g_i = \begin{pmatrix} p_1 \\ p_2 \\ \vdots \\ p_c \end{pmatrix}$
sample prob distribution over classes



0.0001
0.0100

$$\Rightarrow -\log p$$

$$-\log \hat{p} \Leftrightarrow -\log p$$

$$\log \hat{p}/p$$

From information theory

0.5
0.5

≈ 0.01

x100

frequency

Kullback-Leibler divergence:

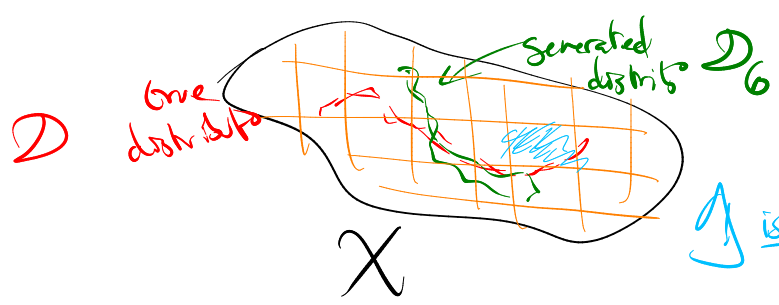
$$KL(p||q) = \int p \log \frac{p}{q} = \sum_{\text{classes}} p_c \log \frac{p_c}{q_c} = -\log \hat{p}_{\text{true class}}$$

not a distance; not symmetric; no triangular inequality

cross-entropy loss

Compare distributions

eg: generative models

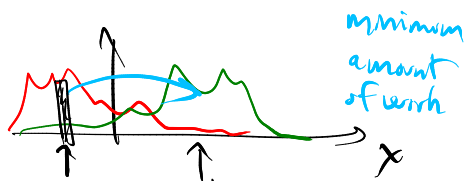
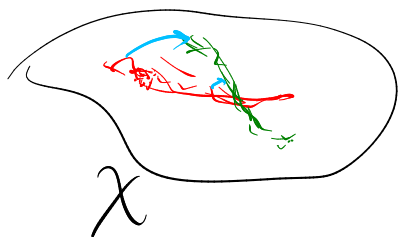


$$KL(D||D-hat)$$

$$KL(D-hat||D)$$

issue: difficult to estimate / empty bins

Optimal transport



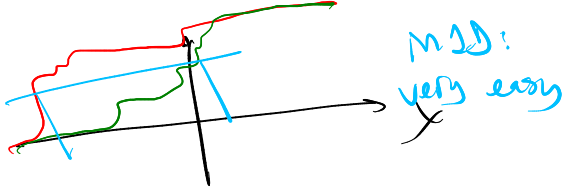
minimum amount of work

$$\text{cost} = \int_{x \rightarrow x'} m_{x \rightarrow x'} d(x, x')$$

$$\text{optimal transport } (D||D-hat) = \inf_M \int \int m_{x \rightarrow x'} d(x, x') dx dx' \left(+ \epsilon H(q) \right)$$

under constraints $\int_x m_{x \rightarrow x'} dx = p(x)$ $\int_{x'} m_{x \rightarrow x'} dx = q(x)$

entropy regularizer



higher dim: very difficult

↳ [Marco Cuturi, 2013] : Sinkhorn algorithm to solve an approximated (very fast)

↑ regularized

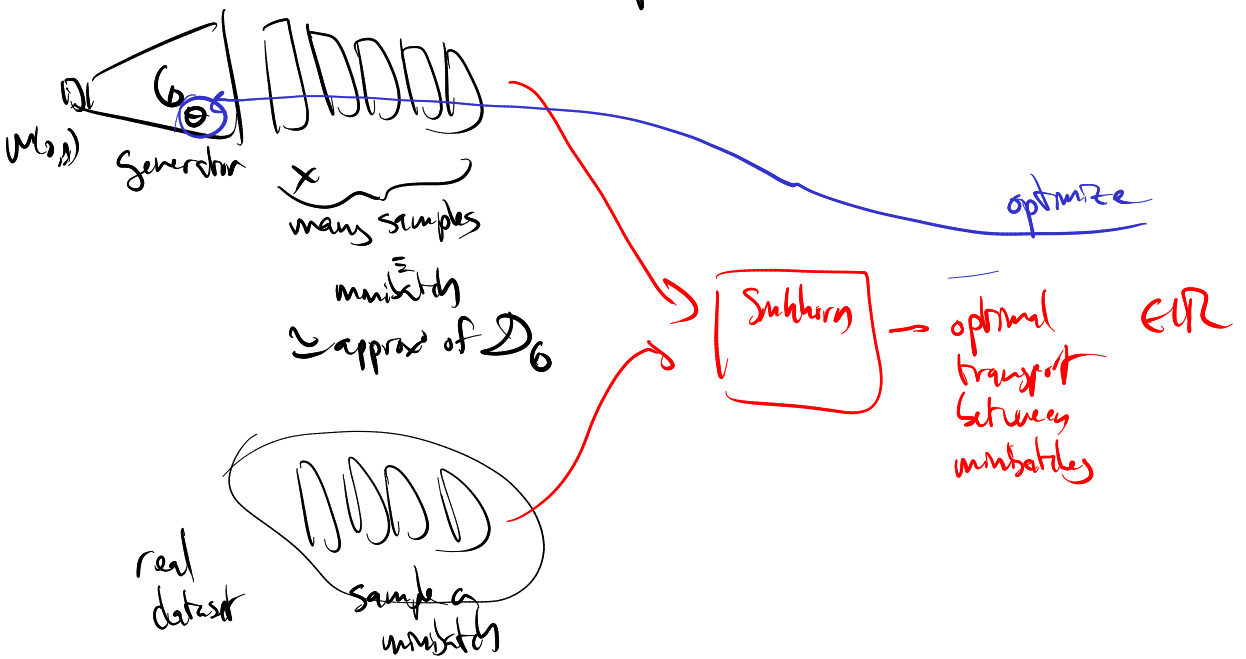
$$K = e^{-d(x_i, x_j) / \epsilon}$$

note: $u, v = (1, \dots, 1)$

$$\text{loop } \left\{ \begin{array}{l} u \leftarrow Kv \\ v \leftarrow Ku \end{array} \right.$$

in Python
for a fixed # of iterations
or until convergence

⇒ optimizable by ∇ descent

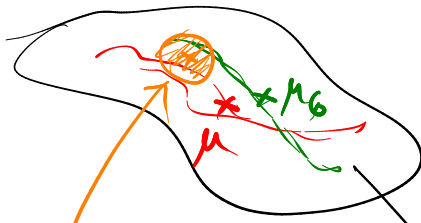


Beware!

of samples $\propto \exp(\text{intrinsic dimension of } \mathcal{D})$

Minimum Mean Discrepancy (MMD)

→ to compare statistics of p & q



$$\|\mu - \mu_G\|^2$$

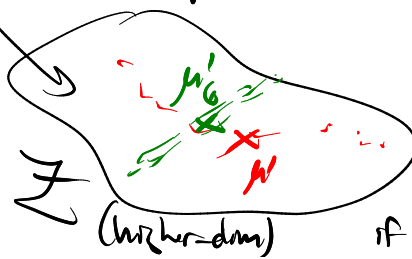
distance between statistics

→ minimize

$$\mu = E[x]_{x \sim \mathcal{D}(p)}$$

$$\mu_G = E[x]_{x \sim \mathcal{D}(p^G)}$$

→ sufficient statistics



$$\|\mu'_G - \mu\|^2$$

↑ very high dim
if \mathcal{Q} well chosen

kernel trick:

$$h(x, x') = \phi(x) \cdot \phi(x')$$

don't design ϕ
but choose k

similarity

typically: $h(x, x') = e^{-\frac{\|x-x'\|^2}{2\sigma^2}}$

rewrite
using only k

$$\left(\underbrace{\mathbb{E}[\phi]}_{\mu} - \underbrace{\mathbb{E}[\phi]}_{\mu_0} \right)^2$$