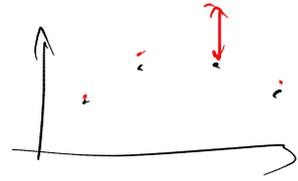


Choosing physically-meaningful metrics

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

ℓ_2 loss: $\sum_{i \in \mathcal{D}} \|\hat{g}_i - y_i\|^2$

↑ prediction
↑ target label



↑ sensitive to outliers

(is ℓ_1 loss: $\sum_{i \in \mathcal{D}} |\hat{g}_i - y_i| \rightarrow$ inducing sparsity \rightarrow on weights rather than on outputs)

- classifier: cross-entropy $\sum_{i \in \mathcal{D}} -\ln p$

Kullback-Leibler divergence

$KL(p||q) = \int p \log \frac{p}{q} \rightarrow \sum_c p_c \log \frac{p_c}{q_c}$

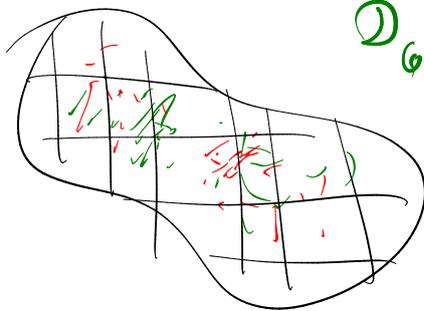
↑ true label $\rightarrow \begin{pmatrix} 0 \\ 1 \\ 0 \end{pmatrix}$
↑ your prediction

↑ not a distance
not symmetric

$\hookrightarrow -\log q$ true class

↑ on a continuous space

$q = \begin{pmatrix} 0.1 \\ 0.2 \\ 0.01 \\ 0.5 \\ \vdots \end{pmatrix}$



\mathcal{D}_G generated distrib.

\mathcal{D} target distrib. (training set)

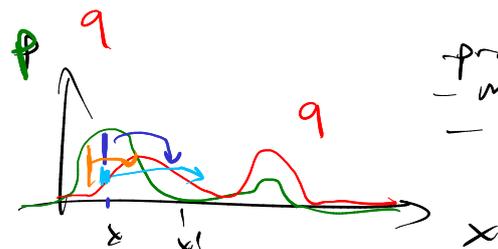
- no notion of geometry
- if $q=0$ for some part where p is not \rightarrow too cost

in practice: continuous & estimate density

↑ high-dim setting

Optimal transport

- Earth Mover distance
- Wasserstein
- Monge-Kantorovic

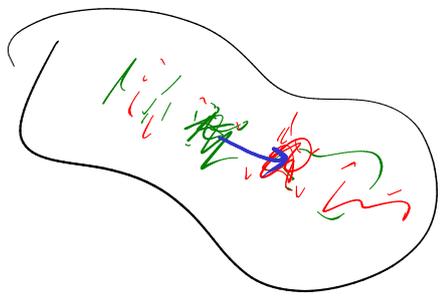


- proba distrib.^o
- mass = 1
- ≥ 0

$\inf_{M \rightarrow} C(M) = \int_{x'} \int_x m_{x \rightarrow x'} \|x-x'\| dx dx'$

↑ mass eff. \rightarrow
↑ amount of work to do

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

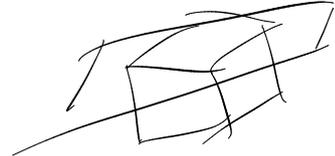


X

- much more robust to displacements than KL
 - maths: distance

1D: closed-form solution

higher-dims: ?



approx: regularize by entropy

$$OT(p, q) + \epsilon H(m)$$

small value

efficient algo. to compute the global opt

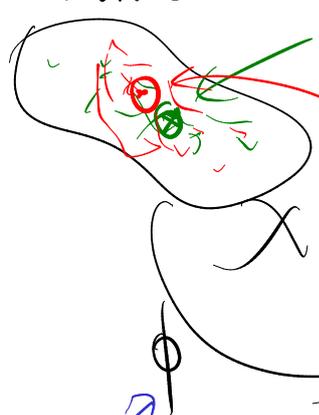
Subhorn algo

→ consists in a sequence of matrix-vector products

Δ sense in high-dim?

Minimum Mean Discrepancy (MMD)

- compare statistics over distrib^o



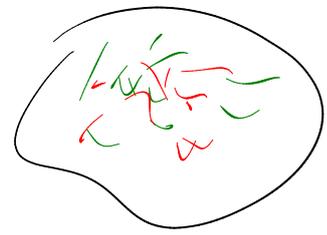
mean: $\mu_p = \mathbb{E}[x_p]$
 avg

$\mu_q = \mathbb{E}[x_q]$
 avg

$$\|\mu_p - \mu_q\|$$

$$\|\mu_p^2 - \mu_q^2\|$$

choose represent^o space (features) that are very rich



Z high-dim

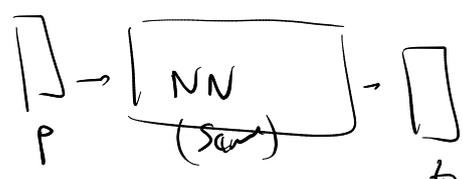
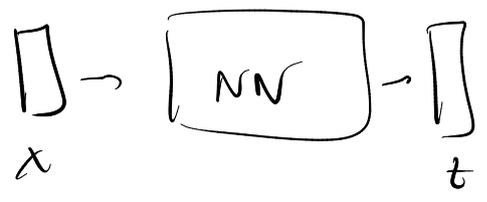
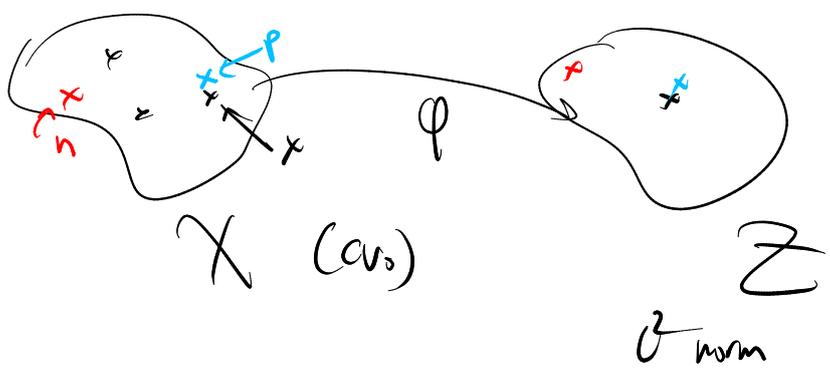
kernelize

$$k(x_1, x_2) = \phi(x_1) \cdot \phi(x_2) \quad \dim = D$$

→ fill discriminative power

Matrix learning

[Siamese networks, GAN]



$d(x_1, x_2) := \| \phi(x_1) - \phi(x_2) \|_2$

Desired property:

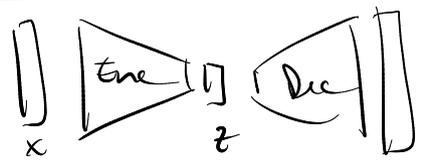
$d(x, p) \leq d(x, n) - m$

contrastive learning margin

loss: $\sum_{(x, p, n) \sim D} \max(d(x, p) - d(x, n) + m, 0)$

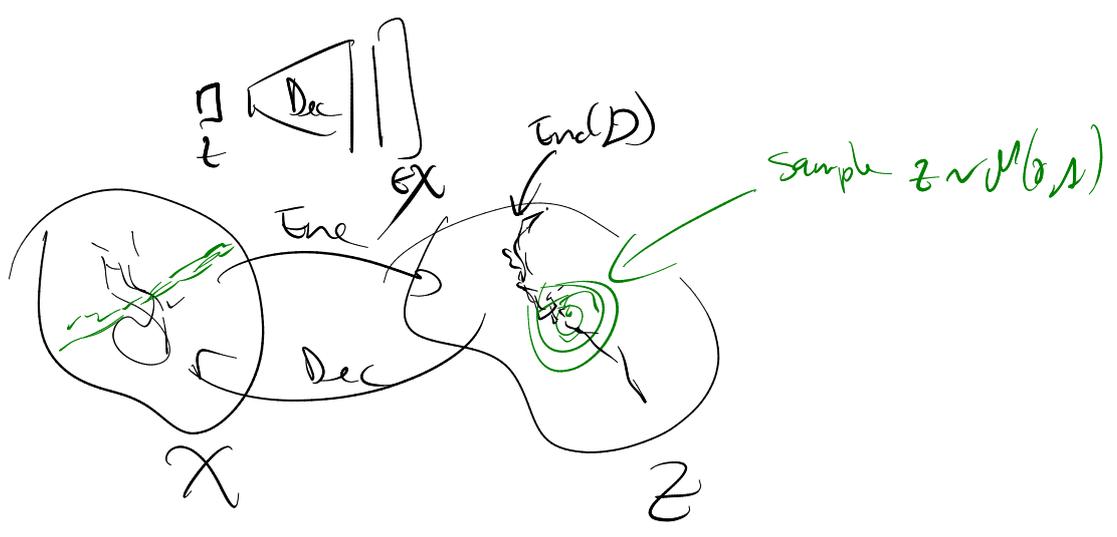
Generative models

Auto-encoder:



$\|x - x'\|$

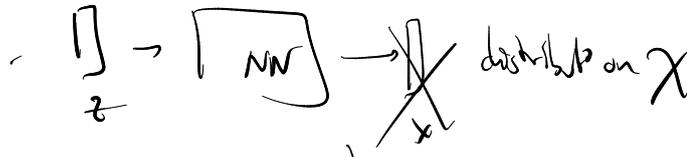
low-dim representation space



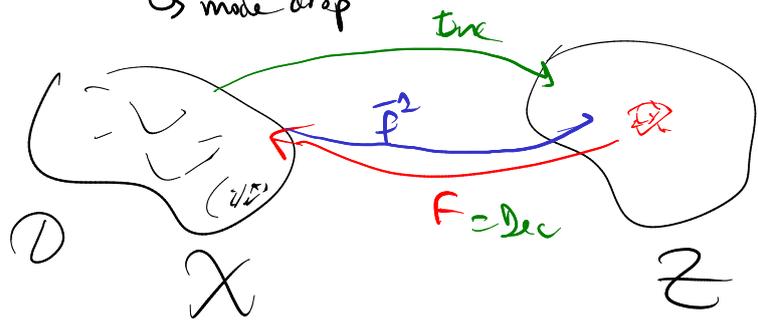
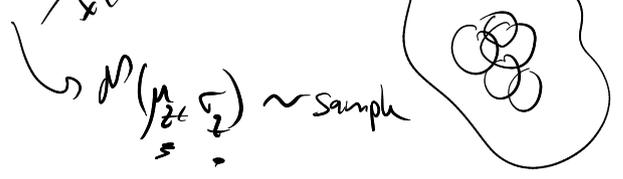
Solve: VAE

Loss = $KL(D || Dec(M(z, \mu)))$ or $KL(Enc(D) || M(z, \mu))$

estimate through sampling? $\int p \sim \sum_{i \in D}$ Monte Carlo
 ELBO lower bound



~~GANs~~: generated samples to be realistic
 \hookrightarrow mode drop



"normalizing" flows

invertible: learn F^2 : send D to a N
 \hookrightarrow just more to opt

sample $u \sim \mathcal{N}(0, 1)$

goal: $F(u, \mu) \sim D$

proposed complex variables $p(x)$
 N closed form $\hookrightarrow KL$

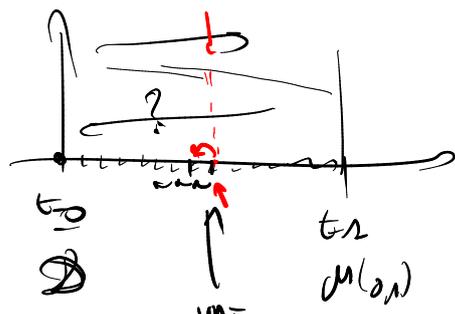
importance sampling

Diffusion models

normalize: explicitly by adding noise (Gaussian)



diffusion $\rightarrow M(t, x)$
 renormalized

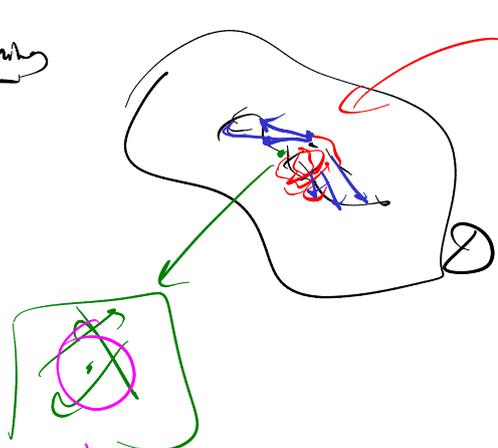


$\frac{\partial I_t}{\partial t} = NN_t(I_t)$
 PDE \downarrow NeuralODE

attention \downarrow 10000 steps ≈ 10

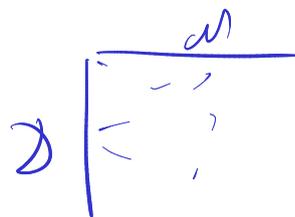
$x_t \rightarrow x_{t_2}$
 $\hat{x}_t = x_t$, $\| \hat{x} - x_t \|$

Flow matching



- N samples of \mathcal{U} (as many of \mathcal{D})
all pairs

draw
all
vectors
from
points of \mathcal{U}
to points of \mathcal{D}



→ vector field

↓
average
of these
vectors
(locally)

& follow that flow: $\mathcal{D} \rightarrow \mathcal{D}$ th

instead: learn a NN to estimate

$x \mapsto$ average flow