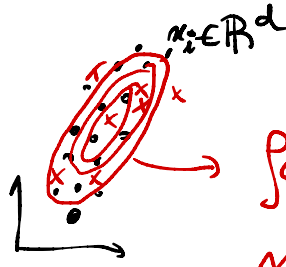
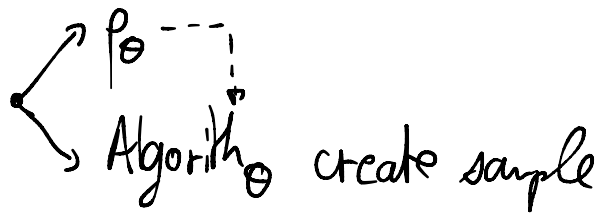


# Density fitting



$p_\theta(x) \in \mathbb{R}^d$  density.  
 $\alpha_\theta$  as a measure

$$\frac{d\alpha_\theta}{dx} = p_\theta$$



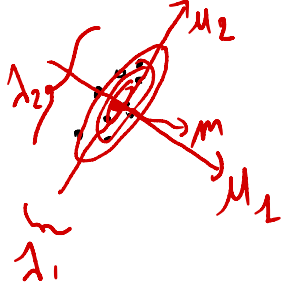
# Max likelihood

$\{x_i\}_i$  are indep. identically dist  
i.i.d  $\alpha_\theta // p_\theta$

$$\max_{\theta} P(\{x_i\}) = \prod_{i=1}^m \frac{P(x_i)}{p_\theta(x_i)}$$

$$\begin{aligned} \min_{\theta} \zeta(\theta) &= -\log(\pi \dots) \\ &= -\frac{1}{m} \sum_i \log(p_\theta(x_i)) \end{aligned}$$

Gaussian case:  $\theta = (m \in \mathbb{R}^d, \Sigma \in \mathbb{R}^{d \times d})$   
 $\Sigma \succeq 0$



$$V = (u_1 | u_2) \quad \mu_i \in \text{eigenvalues}(\Sigma)$$

$$p_{\Theta}(x) = \frac{1}{\sqrt{2\pi \det(\Sigma)}} \exp\left(-\frac{1}{2} \langle \Sigma^{-1}(x-\mu), x-\mu \rangle\right)$$

"Precision"

$$\min_{\mu, \Sigma} \sum_i \frac{1}{2} \log \det(\Sigma) + \frac{1}{2} \langle \Sigma^{-1}(\mu - x_i), \mu - x_i \rangle$$

Convex                      Convex                      Convex

$$\nabla \mathcal{L}\left(\frac{\mu}{\Sigma}\right) = 0$$

$$\mu^* = \frac{1}{n} \sum_i x_i^o$$

$$\Sigma^* = \frac{1}{n} \sum_i \underbrace{(x_i^o - \mu)(x_i^o - \mu)^T}_{\text{rank } 1}$$

Convex<sup>o</sup> - MLE / IT :  $n \rightarrow +\infty$

$$(x_i^o)_i \sim \text{i.i.d. } \hat{\mu}$$

$$\min_{\Theta} \frac{1}{n} \sum \log(p_{\Theta}(x_i)) = \sum \log\left(\frac{1/n}{p_{\Theta}(x_i)}\right) \frac{1}{n}$$

(+  $\log(n)$ )

$$\hat{\tau}(\theta) \xrightarrow{n \rightarrow \infty} \int \log \left( \frac{d\hat{\mu}}{d\mu_{\theta}}(x) \right) d\hat{\mu}(x)$$

$$\triangleq \text{KL}(\hat{\mu} | \mu_{\theta})$$

$\begin{matrix} \uparrow & \uparrow \\ \text{test} & \text{ref.} \\ \text{DATA} & \text{MODEL} \end{matrix}$

Prop:  $\text{KL}(\mu | \nu) = \int \log \left( \frac{d\mu}{d\nu} \right) d\mu \geq 0$

$$= 0 \iff \mu = \nu$$

KL "distance like"

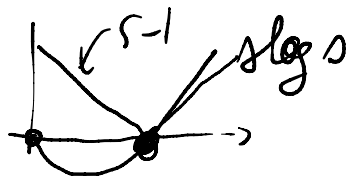
Def:  $\varphi$ -divergence // Csiszar div

$$D_{\varphi}(\mu | \nu) = \int \varphi \left( \frac{d\mu}{d\nu} \right) d\nu$$

KL using  $\varphi(s) = s \log(s)$

Hyp:  $\varphi$  convex +  $\varphi(1) = 0$

$$\varphi(s) = |s - 1|$$



$$D\varphi(\mu | \nu) = \int |\mu - \nu|$$

Jensen-Shannon

Proof:  $D\varphi(\mu | \nu) = \int \varphi\left(\frac{d\mu}{d\nu}\right) d\nu(x) \geq 0$

JENSEN:  $\varphi\left(\int \frac{d\mu}{d\nu}(x) d\nu(x)\right) \leq \int \varphi\left(\frac{d\mu}{d\nu}\right) d\nu$

$\varphi(1) = 0$

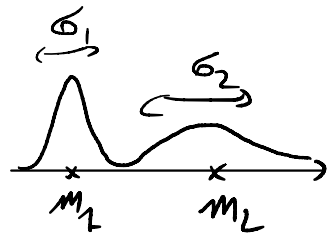
$\varphi\left(\int \frac{d\mu}{d\nu}\right) = \varphi(1) = 0$

Mixture models:

$\pi = (\pi_k)_{k=1}^K \rightarrow$  draw  $k$  according to  $\pi$

$\rightarrow x_i \sim p_{\theta_k}(x)$

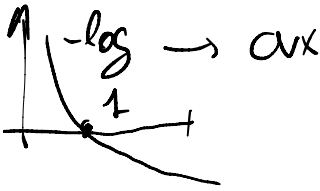
$$p_{\theta}(x) = \sum_k \pi_k p_{\theta_k}(x)$$



$$\ell(\theta) = - \sum_i \log \sum_k \pi_k p_{\theta_k}(x_i) \quad \text{non convex}$$



EM Expectation Maximization



PDF:  $p_{\theta} \rightarrow f(\cdot | \theta)$

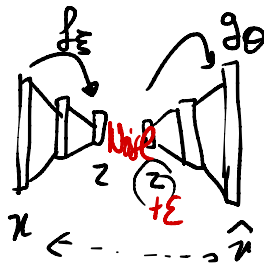
k-means  $\rightarrow$  hard cluster  $P_{ik} \in \{0, 1\}$   
 $\rightarrow \Sigma \alpha Id$

$\Sigma = \sigma Id, \sigma \rightarrow 0$

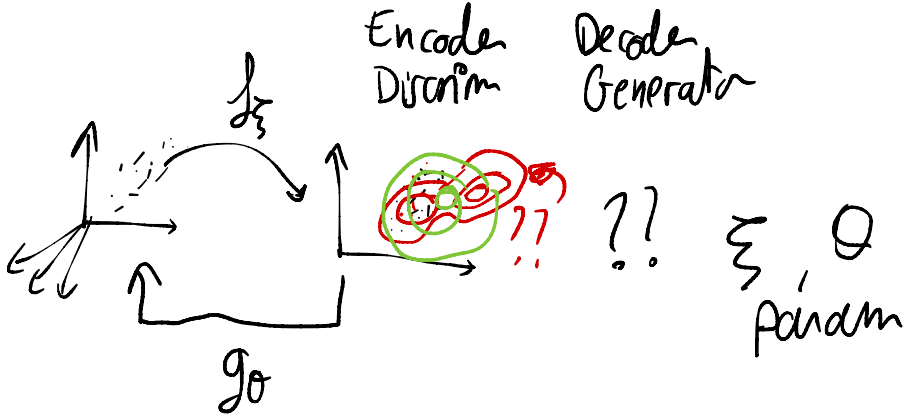
EM  $\rightarrow$  k-means

VAE:

AE

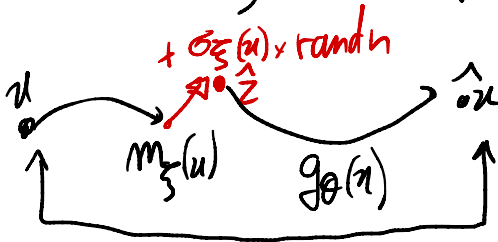


Non-linear  
PCA



$$AE(\theta, \xi) = \frac{1}{n} \sum_i \|x_i - g_\theta(f_\xi(x_i))\|^2$$

$$f_\xi(x) = \begin{pmatrix} m_\xi(x) \\ \sigma_\xi(x) = \exp(S_\xi(x)) > 0 \end{pmatrix} \begin{matrix} \leftarrow \text{mean} \\ \leftarrow \text{noise level} \end{matrix}$$

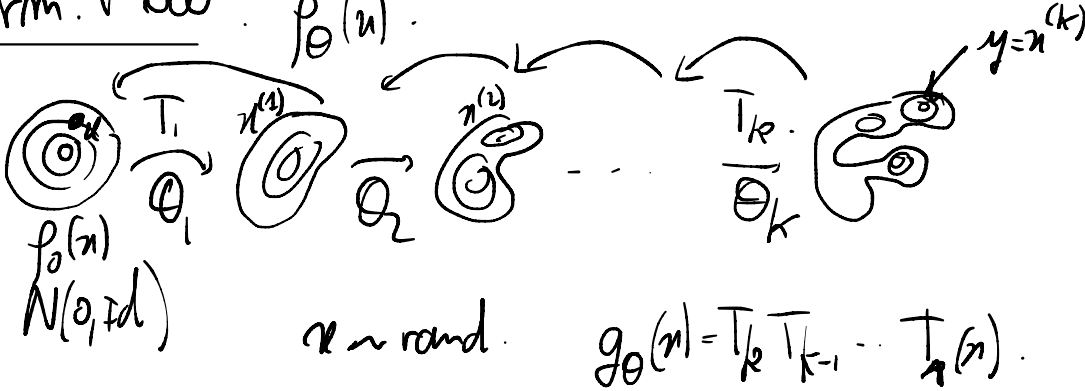


$$VAE(\theta, \xi) = \mathbb{E} \sum_{\text{noise } x_i} \|x_i - \hat{x}_i\|^2 + KL(W_i, \mathcal{N}(0, 1))$$

$$W_i = \mathcal{N}(m_\xi(x_i), \sigma_\xi(x_i))$$

$$KL(N(\mu_i, \Sigma_i) \| N(\sigma, \Sigma)) = \sigma^2(\mu_i) + \|\Sigma^{-1}(\mu_i)\|^2 + \log(\sigma^2)$$

Norm. Flow:  $p_\theta(x)$ .



$$* \min_{\theta} \frac{1}{n} \sum \log(p_\theta(x_i))$$

$$p_0(x) \xrightarrow{y = T_1^{-1}(x)} p_1(y)$$

$$p_1(y) = p_0\left(\frac{T_1^{-1}(y)}{x}\right) \times \frac{1}{|\det \partial T_1^{-1}(y)|}$$

$$\min_{\theta_1, \dots, \theta_k} \frac{1}{n} \sum_i \sum_k \log |\det \partial T(x_i^{(k)})| + \frac{1}{2} \|x_i^{(0)}\|^2$$

$$T_\theta^{-1} \circ x = \begin{bmatrix} x_1 \\ x_2 \end{bmatrix} \rightarrow \begin{bmatrix} y_1 \\ y_2 \end{bmatrix} = \begin{bmatrix} x_1 \\ x_2 \exp(A_\theta(x_1)) + B_\theta(x_1) \end{bmatrix}$$

# Real-Valued Non-Volume Preserving Flows

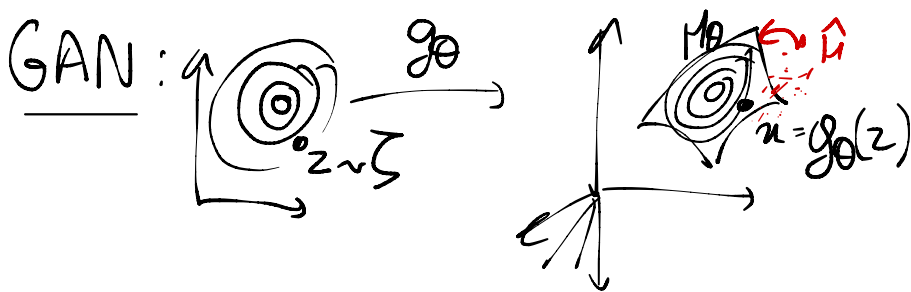
## R-NVP

$$y \rightarrow x = \begin{bmatrix} x_1 = y_1 \\ x_2 = (y_2 - B_\theta(y_1)) \cdot \exp(-A_\theta(y_1)) \end{bmatrix}$$

$$\frac{\partial T}{\partial \theta} = \begin{bmatrix} \text{Id} & | & 0 \\ \dots & | & \text{diag}(\exp A) \end{bmatrix}$$

$$\det \frac{\partial T}{\partial \theta} = \prod_d \exp(A(x_d, d))$$

$$\min_{\theta} \sum_i \sum_k \sum_d A(x_i^{(k)}, d) + \|x_i^{(0)}\|^2 = \text{NF}(\theta)$$



$$\inf_{\theta} D_{\varphi}(\hat{\mu} | M_{\theta})$$

$$\varphi(r) = \log(r)$$

$$\text{JS}(\mu | \nu) = \text{JS}(\nu | \mu) = \text{KL}\left(\mu \middle| \frac{\mu + \nu}{2}\right) + \text{KL}\left(\nu \middle| \frac{\mu + \nu}{2}\right)$$



Proof: Sym, satisfies triang. ineq. (aTV).

$$0 \leq JS \leq 1$$

$$JS = D_{\varphi} \quad \varphi(s) = s \log \left( \frac{2s}{s+1} \right) + \log \left( \frac{2}{s+1} \right)$$

Legendre transform:  $\varphi^*(t) = \sup_s st - \varphi(s)$

$$(\varphi^*)' = [\varphi']^{-1} \quad t = \varphi'(s)$$

$$s = (\varphi')^{-1}(t)$$

Prop: if  $\varphi$  is convex,  $(\varphi^*)^* = \varphi$

$$\varphi(s) = \sup_t st - \varphi^*(t)$$

Prop/Computation

$$D_{\varphi}(\mu | \nu) = \int \varphi \left( \frac{d\mu(x)}{d\nu(x)} \right) d\nu(x)$$

$$= \int \left[ \sup_{t(x)} t(x) \cdot \frac{d\mu(x)}{d\nu(x)} - \varphi^*(t(x)) \right] d\nu(x)$$

$$= \sup_{t(x)} \int t(x) d\hat{\mu}(x) - \int \varphi^*(t(x)) d\nu(x)$$

$$\min_{\Theta} D_{\varphi}(\hat{\mu} | \mathbb{H}_{\Theta})$$

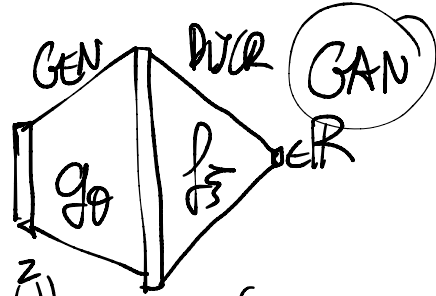
$$x = g_{\Theta}(z)$$

$$\min_{\Theta} \max \int t(x) d\hat{\mu}(x) - \int \varphi^*(t(g_{\Theta}(z))) d\zeta(z)$$

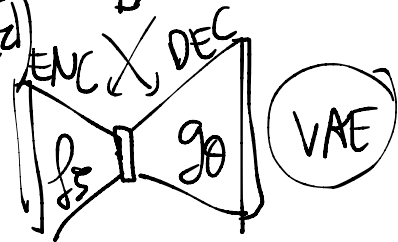
GAN  
move

$$\frac{1}{n} \sum_i t(x_i)$$

$$t(x) = \int_{\mathcal{Z}} \varphi^*(z)$$



$$\min_{\Theta} \max_{\xi} \frac{1}{n} \sum_i \int_{\mathcal{Z}} \varphi^*(z_i) - \mathbb{E}_{z \sim N} [\varphi^*(g_{\Theta}(z))] \text{ ENC } \text{ DEC}$$



in practice  $\rightarrow$  sample  $\frac{1}{N}$

SGD  $\rightarrow$   $k_1$  step desc  $\Theta$

$k_2$  step descent  $\xi$

$$k_1 \gg k_2 = 1$$

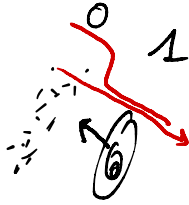
KL  $\varphi(r) = -\log r$   $\varphi^*(t) = e^{t-1}$

JS  $\varphi(r) = \log\left(\frac{2r}{r+1}\right) + \log\left(\frac{2}{r+1}\right)$

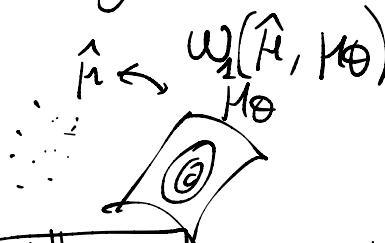
$\varphi^*(t) = \log(t) + \log(1-t)$

$0 \leq t \leq 1$

$\Rightarrow$  discr.  $f_{\mathcal{X}}(x) \in [0, 1]$



WGAN :



$\|\nabla f_{\mathcal{X}}(x)\| \leq 1 \rightarrow |\text{weight}| \leq 1$