

VARIATIONAL AUTOENCODERS

P_ϕ & q_θ	$\xrightarrow{\quad} \downarrow$	VARIATIONAL STOCHASTIC MAPPING
TRADITIONAL		
$q_\theta(x) = y$		$q_\theta(y x)$
$p_\phi(y) = \tilde{x}$		$p_\phi(\tilde{x} y)$

Max. Likelihood for the loss

prior to $p(y) \sim N(\phi, 1)$

$$KL(p_\phi(x|y) \parallel \underbrace{N(\phi, 1)})$$

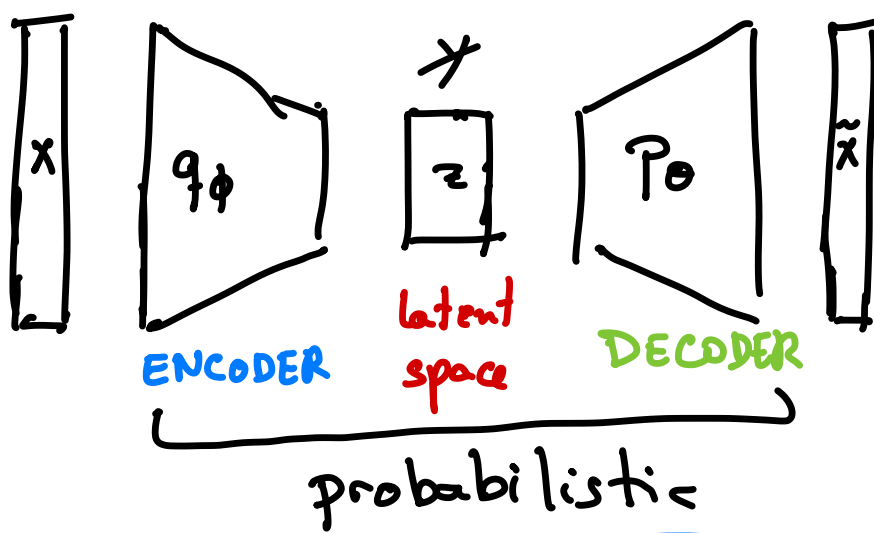
LECTURE 5

28/10/2024

INTRODUCTION TO VAE

INTRINSIC DIMENSION ESTIMATION

DENSITY ESTIMATION



$$q_\phi(x) = z \rightarrow q_\phi(z|x)$$

$$p(z)$$

$$p_\theta(x|z)$$

Generative process

- ① Sample z
- ② Sample x from $p(x|z)$

$$p(x) = \int p(x|z) \underbrace{p(z) dz}_{\text{unknown}}$$

$$MC \rightarrow E_{z \sim p(z)} [p(x|z)] \sim p(x) \approx \frac{1}{K} \sum_{i=1}^K p(x|z_i)$$

Variational inference

$$\{q_\phi(z)\}$$

$$\text{Gaussian } \phi = \{\mu, \sigma^2\}$$

↓ Amortized

$$q_\phi(z|x)$$

$$\log p(x) = \log \int p(x|z) p(z) dz$$

$$\log \int \frac{q_\phi(z|x)}{q_\phi(z|x)} p(x|z) p(z) dz =$$

$$\log E_{z \sim q_\phi(z|x)} \left[\frac{p(x|z) p(z)}{q_\phi(z|x)} \right] \geq$$

$$E_{z \sim q_\phi(z|x)} \left[\log \left(\frac{p(x|z) p(z)}{q_\phi(z|x)} \right) \right] =$$

$$E_{z \sim q_\phi(z|x)} (\log(p(x|z)) -$$
$$E_{z \sim q_\phi(z|x)} [\log(q_\phi(z|x) - \log p(z))])$$

$$\log p(x) \geq E_{z \sim q_\phi(z|x)} [\log(p(x|z))]$$
$$- D_{KL}(q_\phi(z|x) || p(z))$$

$$\log p(x) \geq \mathbb{E}_{z \sim q_\phi(z|x)} [\log(p(x|z))] - D_{KL}(q_\phi(z|x) \parallel p(z))$$

ELBO

PARAMETERIZING

MNIST $28 \times 28 \rightarrow D = 784$

0 ... 255

$p(x) = \text{Categorical}$

$$P \begin{pmatrix} p(x_1=0) \dots p(x_1=255) \\ \vdots \\ p(x_{784}=0) \dots \end{pmatrix}$$

$p_\theta(x|z) = \text{Categorical}(x | \theta(z))$

$\theta(z) = \text{Softmax}(NN(z))$

$q_\phi(z|x), p(z)$

$q_\phi(z|x) = \mathcal{N}(z; \mu_\phi(x), \text{diag}(\sigma_\phi^2(x)))$

$p(z) = \mathcal{N}(z; 0, 1)$

$(\mu_\phi(x); \sigma_\phi^2(x)) = NN(x)$

Optimizing the ELBO

$$\nabla_{\{\theta, \phi\}}$$

$$\nabla_{\theta} \Rightarrow \nabla_{\theta} \mathcal{L}(\theta, \phi, x) =$$

$$\nabla_{\theta} E_{q_{\phi}(z|x)} \left[\log p_{\theta}(x|z) - \log q_{\phi}(z|x) \right] =$$
$$E_{q_{\phi}(z|x)} \left[\nabla_{\theta} \log p_{\theta}(x|z) - \underbrace{\nabla_{\theta} \log q_{\phi}(z|x)}_{\neq} \right]$$

$$\nabla_{\theta} \mathcal{L}(\theta, \phi, x) = E_{q_{\phi}(z|x)} \left[\nabla_{\theta} \log p_{\theta}(x|z) \right]$$

$$\sim \nabla \log p(x|z)$$

$$\nabla_{\phi}$$

REPARAMETERIZATION TRICK

$$z = g_{\nu}(\epsilon) \quad \epsilon = \hat{q}(\epsilon) \quad \hat{q} \text{ it's independent on } \phi$$

$$\hat{\phi} = (\phi, \nu)$$

$$\mathcal{L}(\hat{\phi}, \theta, x) = E_{\hat{q}(\epsilon)} \left[\log(p(x, g_{\nu}(\epsilon))) - \log q_{\phi}(g_{\nu}(\epsilon)) \right]$$

So $\nabla_{\hat{\phi}}$ -

$$E_{\hat{q}(\epsilon)} \left[\underbrace{\nabla_{\hat{\phi}} \log(p(x, g_{\nu}(\epsilon))) - \nabla_{\hat{\phi}} \log q_{\phi}(g_{\nu}(\epsilon))}_{G(\epsilon)} \right]$$

$$\nabla_{\hat{\phi}} \mathcal{L}(\hat{\phi}, \theta, x) = \frac{1}{S} \sum_i G(\epsilon_i)$$

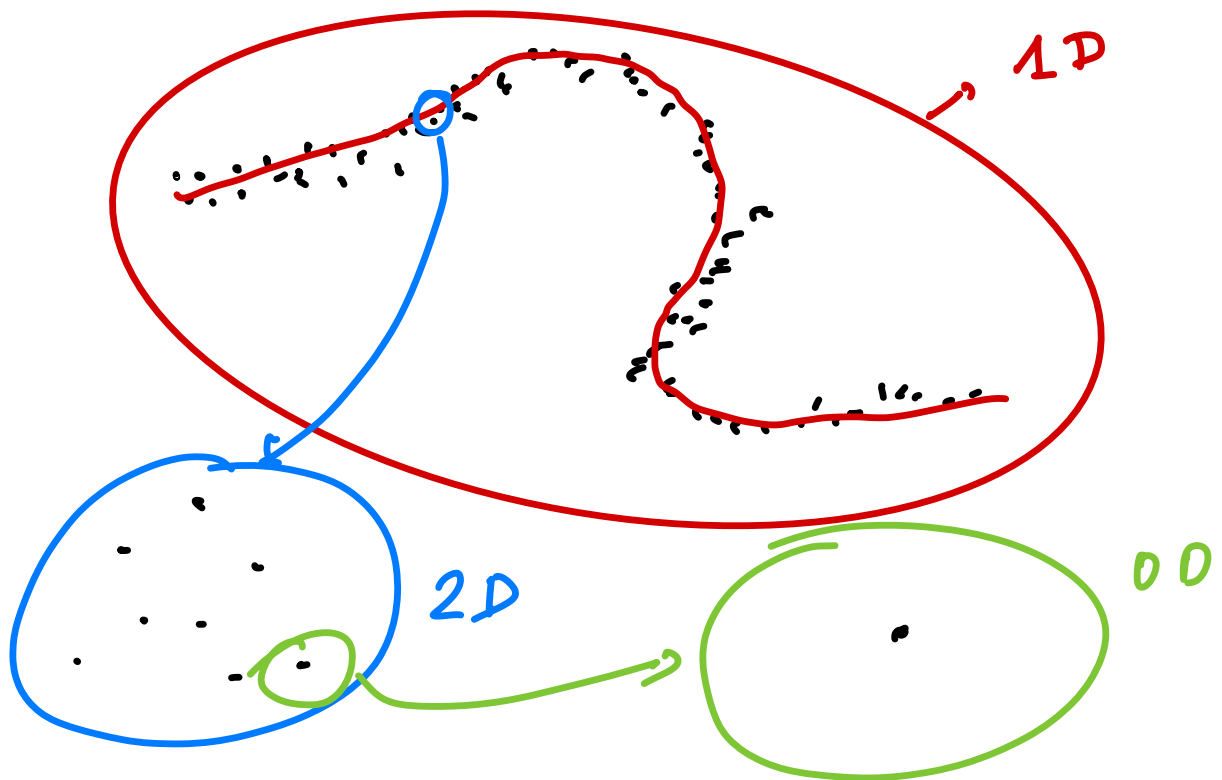
$$z = \mu + \sigma \epsilon \quad \epsilon \sim \mathcal{N}(0, \mathbf{I})$$

BUILDING A VAE

- ① Given x_n apply the encoder $(\mu_{\phi}(x_n), \sigma_{\phi}^2(x_n))$
- ② Use reparam... trick to obtain $z_{\phi,n} = \mu_{\phi}(x_n) + \sigma_{\phi}^2(x_n) \cdot \epsilon$
- ③ Apply decoder to obtain $\theta(z_{\phi,n})$
- ④ Compute the ELBO
- ⑤ Compute gradients and update ϕ, θ

INTRINSIC DIMENSION

- ① Minimum Number of variables needed for represent the data with minimum information loss
- ② The dimension of the manifold in which our data lies

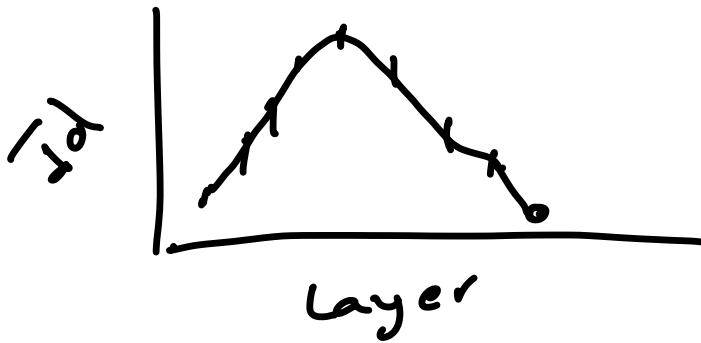
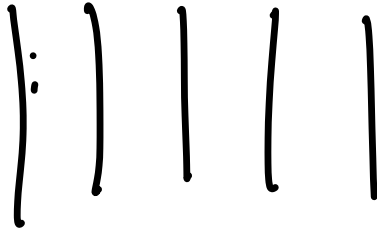


Intrinsic dimension is a scale-dependent property

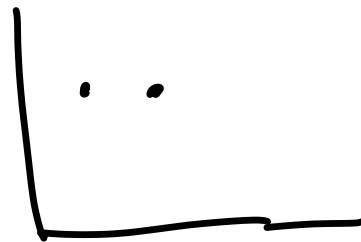
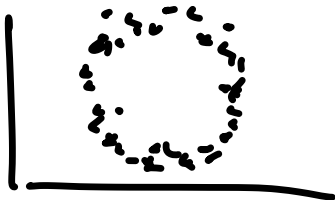
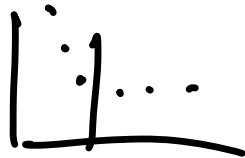
① Decide d

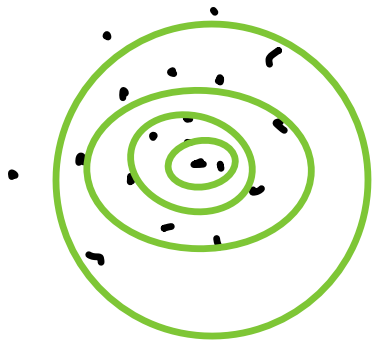
② Information content

I_d

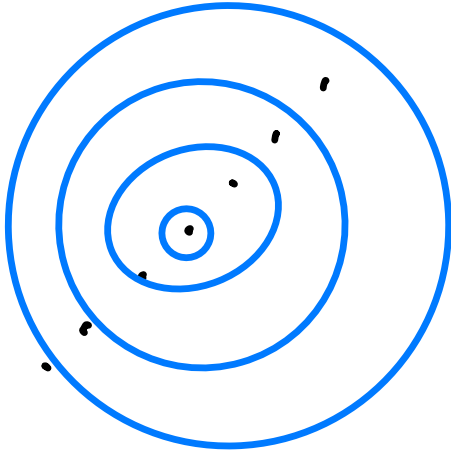
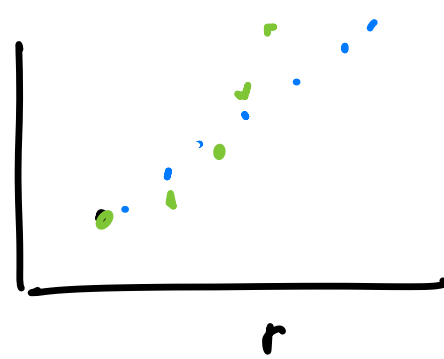


PCA \rightarrow Spectrum





nn



$$\# nn \propto r^d$$

$$\# nn = p \cdot r^d$$

$$\log(\#) = \log(p) + d \log(r)$$

for each data point

a) Compute the distances from its K-NN

b) Plot the log of the rank as function of the logarithm of the distances

c) linear fit

d) slope will be d

$$\log(\#) = \log(p) + d \log(r)$$

$$i \quad r_i^1 \quad r_i^2 \quad \dots \quad r_i^k$$

$$\begin{matrix} \log(1) & \log(r_i^1) \\ \log(2) & \log(r_i^2) \end{matrix}$$

Id from the p

two-NN

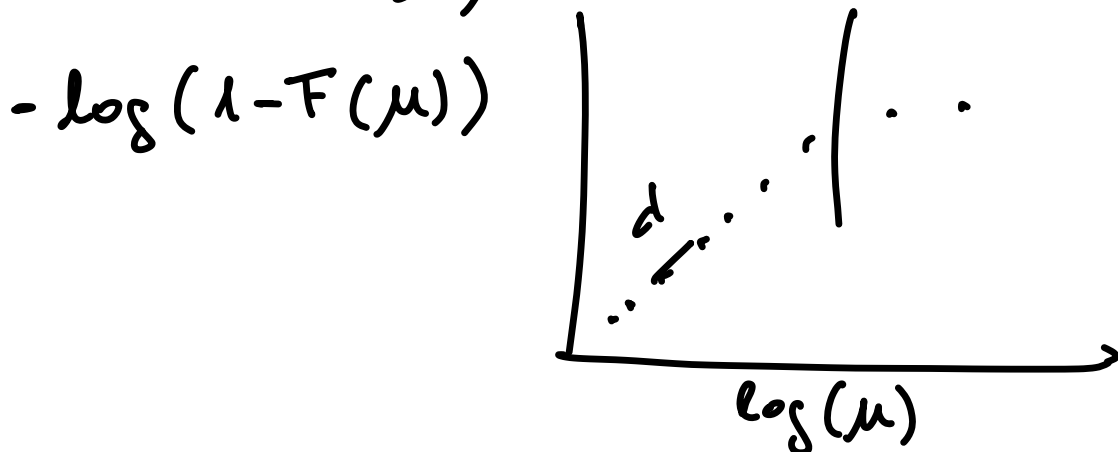
$$\mu_i = \frac{r_2^i}{r_1^i} \in [1, \infty) \quad \text{i.i.d}$$

$$p(\mu) = \underbrace{\mu^{-d-1} \cdot d}_{\text{pareto}}$$

p within the second NN \sim constant

$$\textcircled{a} \quad F(\mu) = \int_1^\mu \mu^{*-d-1} \cdot d \cdot d\mu^* = 1 - \mu^{-d}$$

$$\frac{-\log(1-F(\mu))}{\log(\mu)} = d$$



$$\textcircled{b} \quad \text{ML} \quad p(\mu) = \mu^{-d-1} d$$

$$\log \mathcal{L} = \sum_i \log(\mu_i^{-d-1} \cdot d) =$$

$$N \log(d) - (d+1) \sum_i \log(\mu_i)$$

$$\frac{\partial \log \mathcal{L}}{\partial d} = \frac{N}{d} - \sum_i \log(\mu_i) = 0$$

$$d = \frac{N}{\sum_i \log(\mu_i)}$$

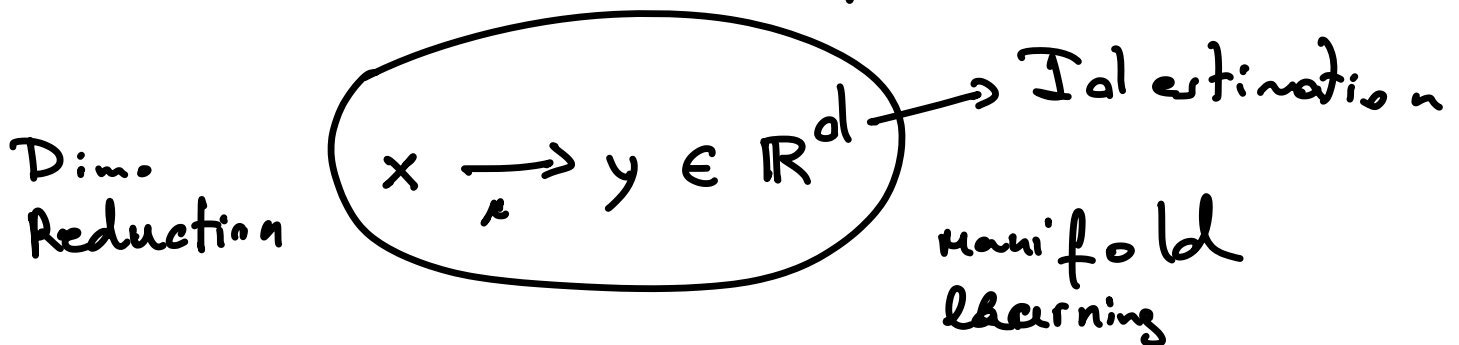
- Max. Likelihood
- DANCO



DENSITY ESTIMATION

$P(y|x)$
Supervised

$P(x)$
Unsupervised



Obtain $p(x)$ task density estimation

Parametric: Assume a functional form
learn the parameters by fitting
the data

Non-parametric: There's no assumption
about the functional
form

Parametric: Powerful but "rigid"

Non-parametric:

Gaussian Mixture Model (Parametric):

$$p(x) \sim p(x)$$

$$p(x) = \sum_i^k \pi_i \psi(x, \theta_i)$$

$$\theta_i = (\mu_i, \Sigma_i)$$

$\mu \rightarrow \mu, \Sigma \rightarrow \Sigma$

Max. Likelihood

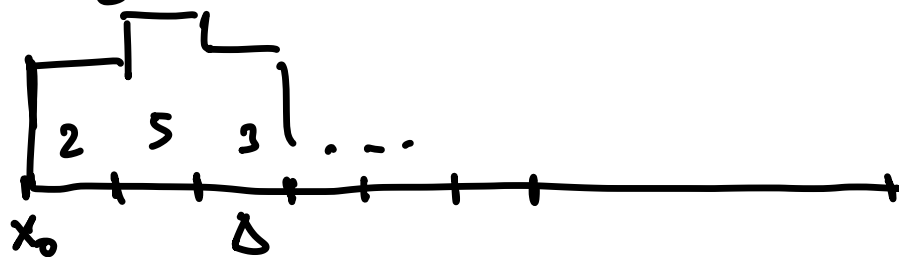
$$\mathcal{L} = \prod_l^N \left(\sum_i^k \pi_i \psi(x^l; \theta_i) \right)$$

$$\log \mathcal{L} = \sum_l^N \log \left(\sum_i^k \pi_i \psi(x^l; \theta_i) \right)$$

Expectation - Maximization

Non-parametric density estimation

• Histograms



$x_{\min}; \Delta; N_{\text{bins}} \text{ or } x_{\max}$

$$p(x) \sim p(x_j) = \frac{n_j}{N \Delta}$$

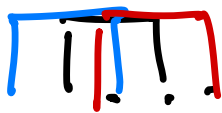
$$\mathcal{E}(p(x_j)) \propto \frac{p(x_j)}{\sqrt{n_j}}$$



$$\Delta \sim \frac{2 \text{ IQR}}{\sqrt[3]{N}}$$

Freedman-diaconis Rule

$$\text{IQR} = X_{(3/4)} - X_{(1/4)}$$



$$p(x) \approx \hat{p}(x) = \sum_i f(x_i, h)$$


$$\frac{1}{N h}$$

kernel density estimation

$$p(x) \approx \hat{p}(x) = \sum_i^N K(x_i, h) \frac{1}{N h}$$

$$K = \text{Uniform} \quad \boxed{\cdot}^h_{x_i} = \text{parzen window}$$

$$K = \text{Gaussian} \quad K(x, x_i, h) = \frac{1}{\sqrt{2\pi}} \exp\left(-\frac{1}{2} \left(\frac{x - x_i}{h}\right)^2\right)$$

K  triangle

Cosine 

Epanechnikov  cubic

$h \rightarrow$ Silverman's Rule

$$h = 0.9 \min \left(\sigma, \frac{IQR}{1.34} \right) \cdot N^{-1/5}$$

Derived assuming gaussianity

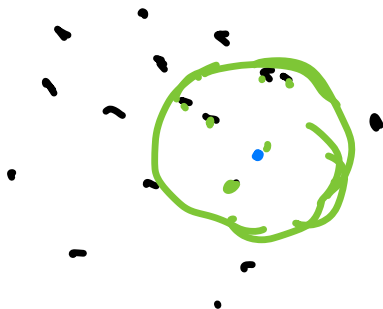
$$h \begin{cases} \rightarrow p_1(x) \\ \rightarrow p_2(x) \end{cases} \quad \Delta = KL(p_1 \| p_2)$$

$$h_{opt} = \operatorname{argmin}_h \Delta(h)$$

K-Nearest Neighbor density estimator

$$P(x) = \frac{k}{N V_k} = \frac{k}{V^d r_k^d N}$$

Volume of the hypersphere of dimension d



$r_k \rightarrow$ distance from the k -th NN

$$E(P) = \frac{P}{\sqrt{k}}$$