

FYS5429/9429, April 10 lecture

FYS5429/9429 April 10

Basics of generative models

Latent $q(u) \rightarrow p(y)$

$P(x; \theta)$

$P(y; u)$

$P(x; u)$

Data set

$$P(x; \theta) = \int d\mu p(u) p(x|u)$$

we want $P(x; \theta)$

$$X \in \{ \underbrace{x_0, x_1, \dots, x_{n-1}}_{\mathcal{D}} \}$$

$$P(X_j; \theta) = \prod_{x_i \in \mathcal{D}} P(x_i; \theta) \quad (\text{i.i.d.})$$

$$\underbrace{P(x_i; \theta)}_{\text{Marginal probability}} = \sum_{h_j} P(x_i, h_j; \theta)$$

↑
hidden / latent

$$P(x_i; \theta) = \int dh p(x_i, h; \theta)$$

Boltzmann machine

$$P(x; \theta) = \frac{f(x; \theta)}{Z(\theta)}$$

$$Z(\theta) = \int_{x \in D} dx \int_h dh p(x, h; \theta)$$

normalization constant /
partition function

$$P(X; \theta) = \frac{1}{Z(\theta)} \prod_{x_i \in X} f(x_i; \theta)$$

$$= \frac{1}{Z(\theta)} \prod_{x_i} \left[\sum_{w_j} f(x_i, w_j; \theta) \right]$$

optimize $\log P(X; \theta)$

$$\nabla_{\theta} \log P(X; \theta) = 0 =$$

$$\nabla_{\theta} \left[\sum_{x_i \in X} \log f(x_i; \theta) \right]$$

$$= \boxed{E \left[\log f(x_i; \theta) \right]} \quad \text{(MC)}$$

RBM

$$a^T x + b^T h + x^T w h$$

$$f(x_i | h_j; \theta) = e$$

$$a^T x + b^T h + x^T w h =$$

$$\sum_{i=1}^N a_i x_i + \sum_{j=1}^M b_j h_j + \sum_{i,j}^{MN} x_i w_{ij} h_j$$

Binary-Binary

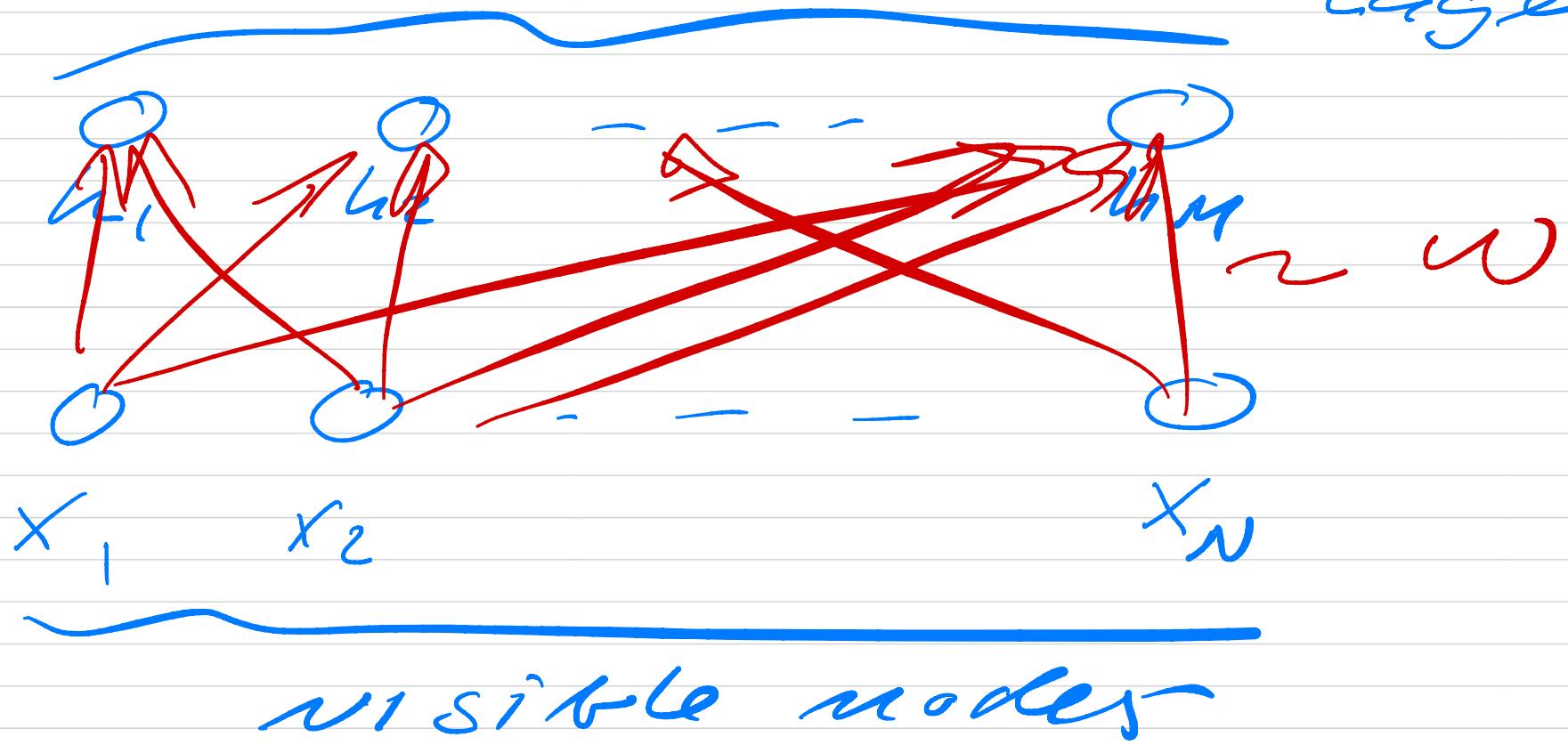
$$\theta = \{a, b, w\}$$

$$x_i = \{0, 1\}$$

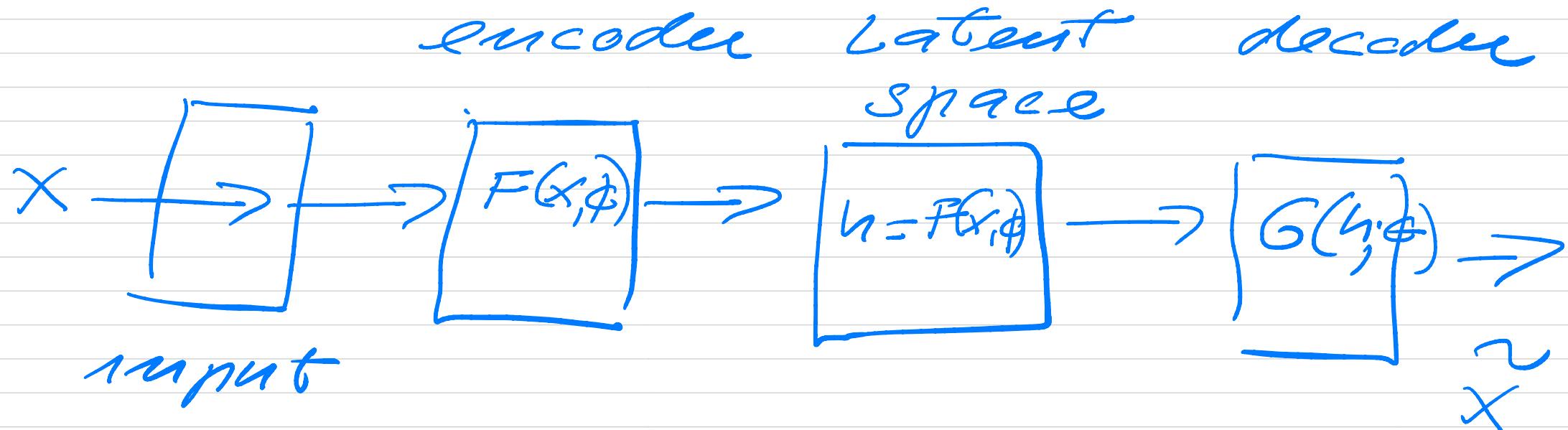
$$h_j = \{0, 1\}$$

NN for RBM

\sim hidden layer



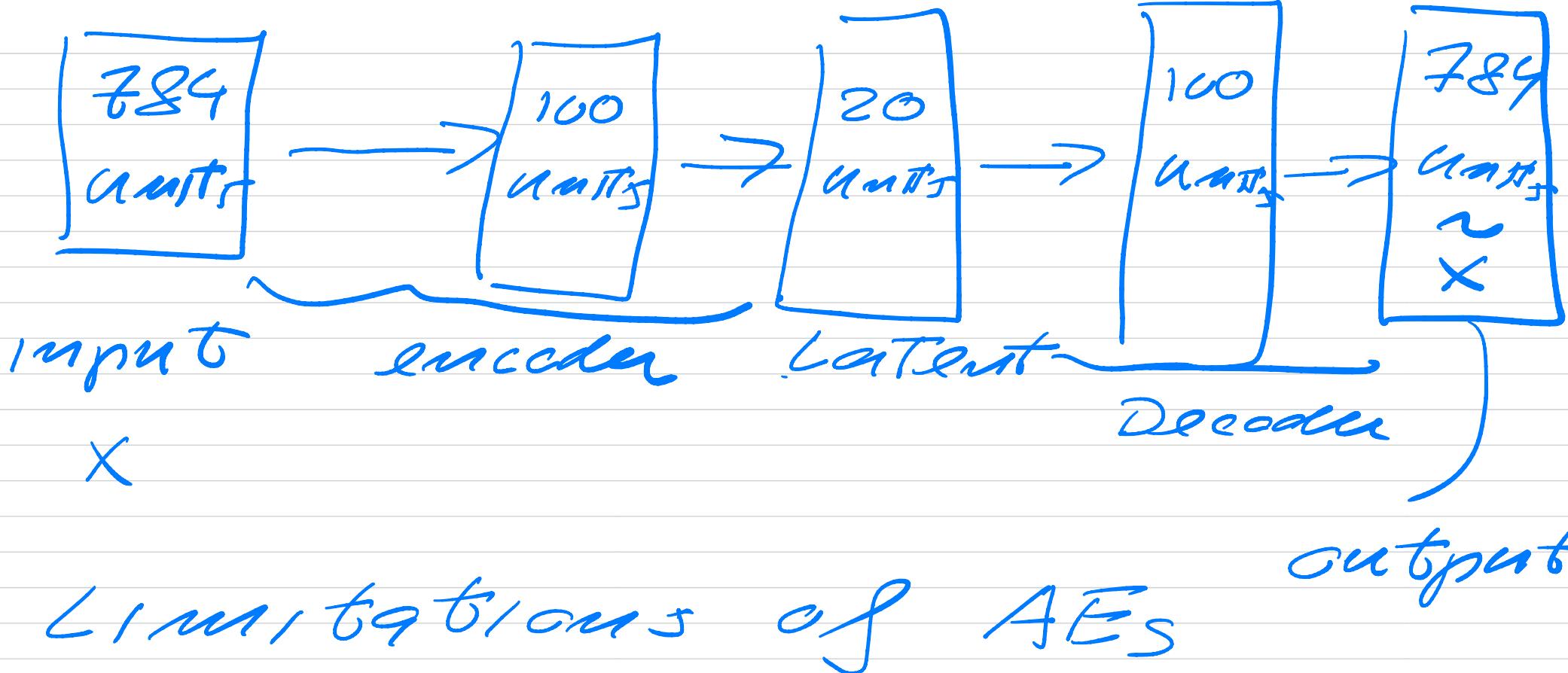
Reminder on Auto encoders



$$\tilde{x} = G(h; \phi)$$

reconstruction error

$$\hat{\epsilon}_{\phi} = \underset{\epsilon_{\phi}}{\text{arg min}} \| (x - \tilde{x}) \|_2^2$$



Limitations of AEs

- Not generative, they don't define a distribution
- How to choose the latent dimension,

why Variational AEs?

- A VAE is a generative model that learns a low-dimensional representation.

$$\begin{aligned} \text{skip } (\epsilon, \phi, \text{etc}) &\quad \xleftarrow{\text{conditional prob}} \\ p(x, h) &= p(x|h) \quad p(h) \\ &= p(h|x) \quad p(x) \end{aligned}$$

Marginal prob

$$p(x) = \int dh \, p(x|h) \, p(h)$$

in case of discrete prob's

$$P(x_i) = \sum_{h_j} P(x_i | h_j) P(h_j)$$

$$P(h_j) = \sum_{x_i} P(h_j | x_i) P(x_i)$$

Bayes' theorem — likelihood

$$\boxed{P(x|h)} = \frac{\boxed{P(h|x)} \boxed{P(x)}}{\sum_{x_i} P(x_i | x_i) P(x_i)}$$

posterior prior

we want to optimize

$$\log P(X) \quad (\log P(X; \theta))$$

$$= \log \left[\int p(h) P(X|h) dh \right]$$

$$= \log \left[\underbrace{\int p(x, h)}_{P(X|h)P(h)} dh \right]$$

$$P(X|h)P(h)$$

insert

$$1 = \underline{q(h)}$$

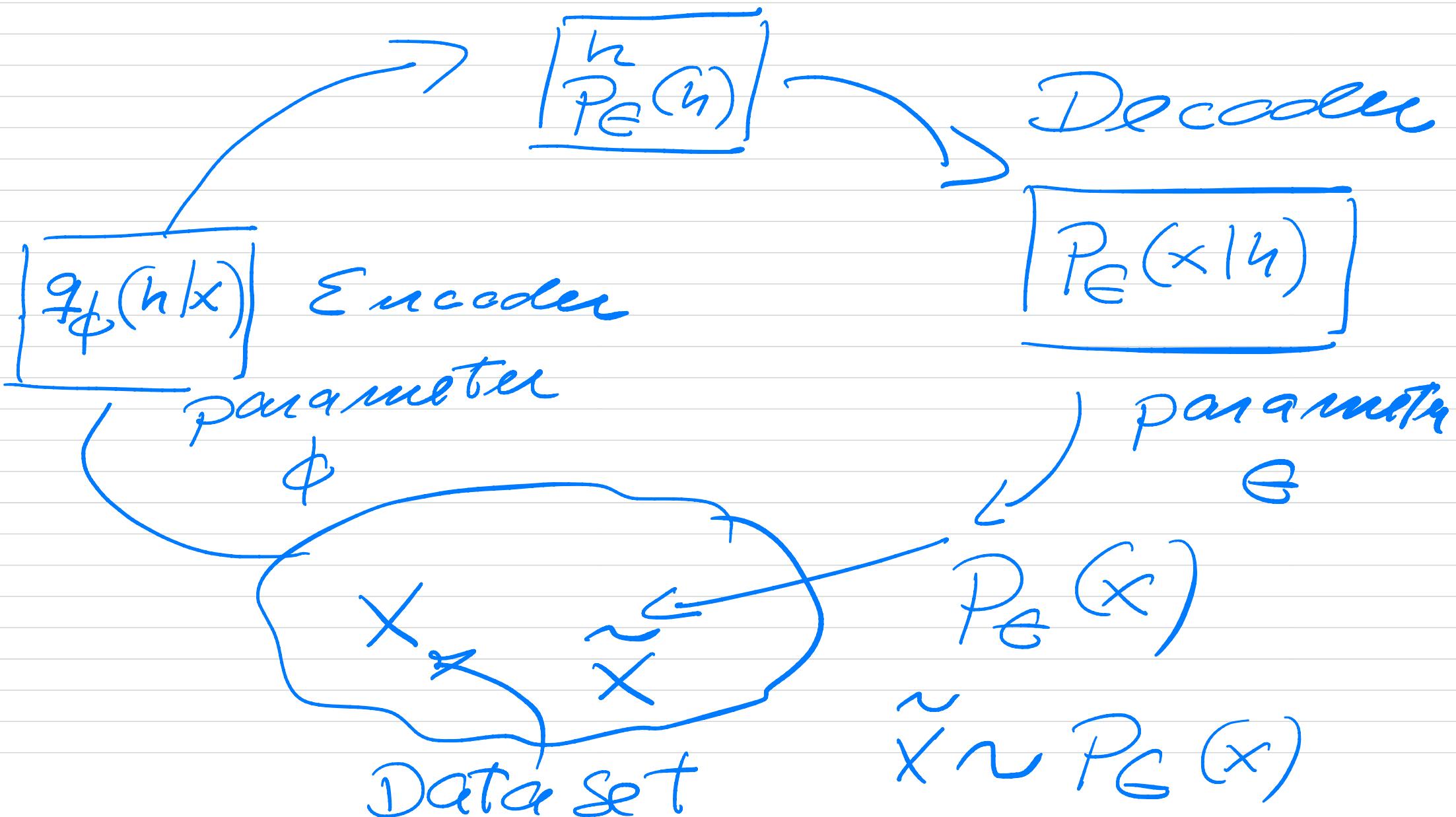
Decoder

encoder

$$\underline{q(h)}$$

Latent Space

VAE latent space



$$\log p(x) = \log \left[\int p(u) p(x|u) du \right]$$

q is probability

$$\text{insert } 1 = \frac{q(u)}{q(u)}$$

Jensen's inequality

$$\geq \int q(u) \log \left[\frac{p(u)}{q(u)} p(x|u) \right] du$$

$$= \int q(u) \log \left[\frac{p(u)}{q(u)} \right] du +$$

$$\int q(h) \log(p(x|h)) dh$$

$$= E_{h \sim q} \left[\log \frac{p(x)}{q(h)} \right] \quad (1)$$

\downarrow
 $(h \sim q(h) \quad x \sim N(\mu_x, \Sigma_x))$

$$+ E_{h \sim q} [\log p(x|h)] \quad (2)$$

let us look at the 2nd term
first

assume

$$\log p(x|u) = \log N(x; G_e(u), \Sigma)$$

$$= \log \left[\frac{1}{(2\pi\Sigma)^{d/2}} \exp \left[-\frac{1}{2\Sigma} \|x - G_e(u)\|_2^2 \right] \right]$$

$$= -\frac{1}{2\Sigma} \|x - G_e(u)\|_2^2 + \text{const}$$

(think back to autoencoder)

$$\tilde{x} = G_e(u)$$

This is the reconstruction error

what about the first term?

- intermezzo : kullback-leibler divergence

$$KL(P||Q) = D_{KL}(P||Q)$$

$$= \int_{x \in D} P(x) \log \frac{P(x)}{Q(x)} dx$$

$$= E_{x \sim P} \left[\log \frac{P(x)}{Q(x)} \right]$$

Suppose $p(x)$ is the real
prob. distribution and
 $q(x)$ is a model

$$D_{KL}(p||q) = 0 \text{ when } p(x) = q(x)$$

when $p(x)$ tends to zero

$$\lim_{\substack{x \rightarrow 0^+ \\ x > 0^+}} x \log x = 0$$

$$D_{KL}(p||q) \geq 0$$

$$\log p(x) = D_{KL}(q||p)$$

- reconstruction error

For aS :

$$\tilde{x} \sim P(x; \epsilon) = P_\epsilon(x)$$

want to optimize

$$\log P_\epsilon(x) = \log \left[\int P_\epsilon(x|h) P_\epsilon(h) \right]$$

a latent space whose distribution $q_\phi(h|x)$ links the input data with the latent space $-h-$

$$\hat{E}_\phi = \underset{G_\phi}{\operatorname{argmax}} \log P_G(x)$$

$$= \underline{\underline{\log \left[\int P_G(x, h) dh \right]}}$$

$$= \log \left[\int dh P_G(x|h) P_\phi(h) \right]$$

insert 1 = $\frac{q_\phi(h|x)}{q_\phi(h|x)}$

$$= \log \left[\int \frac{P_G(x|h)}{q_\phi(h|x)} \frac{P_\phi(h)}{q_\phi(h|x)} dh \right]$$

Jensens inequality

$$\geq \mathbb{E}_{h \sim q_\phi} \left[\log p_{\text{G}}(x|h) \right]$$

$$(h \sim q_\phi(h|x))$$

$$+ \mathbb{E}_{h \sim q_\phi} \left[\log \frac{p_{\text{G}}(u)}{q_\phi(u|x)} \right]$$

$$= \mathbb{E}_{h \sim q_\phi} \left[\log p_{\text{G}}(x|h) \right]$$

reconstruction score

Decoder

$$- D_{KL} \left(q_\phi(u|x) || p_e(u) \right)$$

encoder part

approximation \downarrow multivariate

$$- q_\phi(u|x) = N(u; \mu_\phi(x), \Sigma_\phi(x))$$

$$- p_e(u) = N(u; 0, 1)$$

mean

We can evaluate the DKL
analytically.

$$\arg \max_{\theta, \phi} \left\{ \mathbb{E}_{h \sim q_\phi} [\log p_\theta(x|h)] \right\}$$

$$- D_{KL}[q_\phi(h|x) || p_\theta(h)]$$

analytical evaluation

(MC)² - sampling