

1. Introduction to Bayesian Statistics

→ In Bayesian statistics probabilities represent „states of knowledge“ or „degrees of belief“

→ Different than frequentist interpretation:
Prob.'s express expected frequencies of events

→ Bayesian statistics tells us how to formally update our state of knowledge when new information (a measurement) becomes available

→ Bayesian statistics works with Conditional Probabilities

a) Working with conditional probabilities

X, Y : continuous random variables
or discrete events

$$\text{Conditional: } P(X|Y) = \frac{P(X, Y)}{P(Y)} \quad \leftarrow \text{joint} \quad (1)$$

$$\text{Marginal: } P(Y) = \int dx P(X, Y) \quad (2)$$

OR

$$\sum_{x \in \{x_i\}} P(X_i|Y)$$

b) Bayes Theorem

→ Tells us how to update our state of knowledge as new data becomes available

- prior state of belief $P(x)$
 x could be parameters in a physics model
- additional knowledge gained through measurement d

Likelihood

~~probability~~

$P(d|x)$

- updated knowledge

posterior

$P(x|d)$

Bayes Theorem

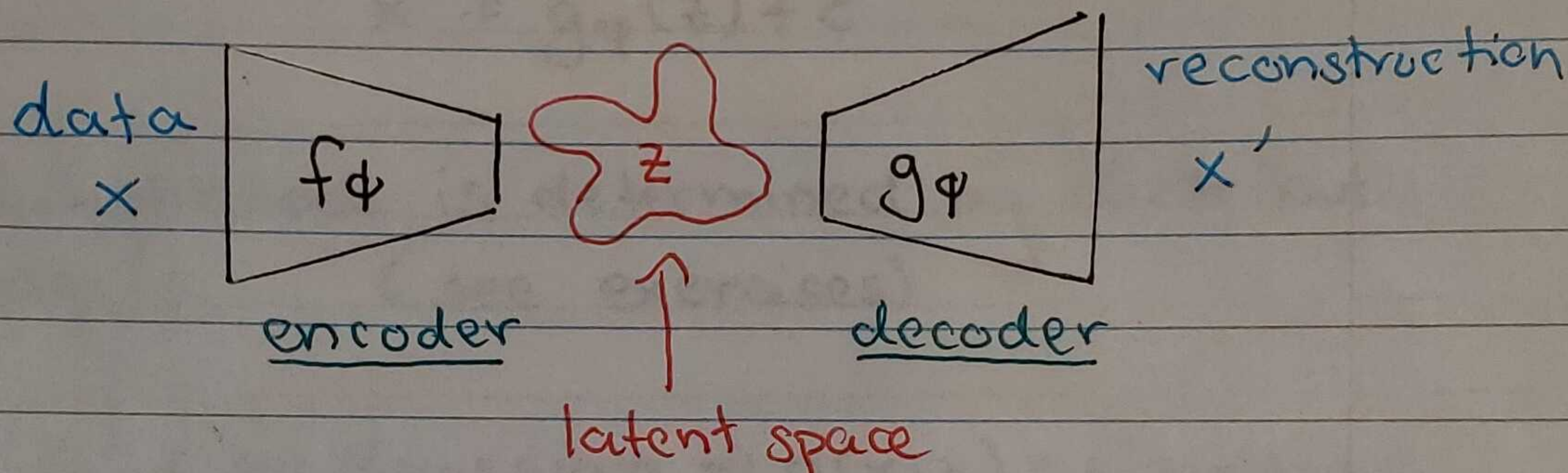
$$P(x|d) = \frac{P(d|x) P(x)}{P(d)}$$

$P(d)$: normalization constant, evidence

2. Why do I care?
OR

Introduction to Variational AEs

a) recap: Autoencoder



$$x' = g_\psi(f_\phi(x))$$

$$\mathcal{L}_{AE} = \underbrace{\|x' - g_\psi(f_\phi(x))\|_2^2}_{\text{loss function}}$$

b) Probabilistic interpretation

-> Information loss in the compression

=> can't recover x precisely

=> implies probabilistic structure

-> data x follows probability distr. $p(x)$

-> all available data points are drawn from this distribution

$$p(x) = \int dz P(x, z) = \int dz P(x|z) P(z)$$

Likelihood $p(x|z)$

→ Information loss in compression implies probabilistic structure

$$x' = g_{\psi}(z) + \epsilon$$

→ Likelihood is determined by distribution of ϵ (see exercises)

→ if $\epsilon \leftarrow \text{Gaussian} \Rightarrow p(x|z) = \text{Gaussian}$

$$p_{\psi}(x|z) = \mathcal{G}(x | \mu = g_{\psi}(z), \Sigma = \Sigma)$$

→ Likelihood depends ^{on} / is parameterized by network parameters

Prior $p(z)$

→ prior knowledge (before looking at data sample x) is average distribution of z

$$p(z) = \int dx p(z|x) p(x)$$

c) Derivation of Evidence Lower Bound

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

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

$$\log p(x) = \log \int dz \, q(z|x) p(x|z) \frac{p(z)}{q(z|x)}$$

$$E_{q(z|x)} \left[p(x|z) \frac{p(z)}{q(z|x)} \right]$$

$$\log p(x) \geq \int dz \, q(z|x) \left[\log p(x|z) - \log \left[\frac{q(z|x)}{p(z)} \right] \right]$$

likelihood term = measure of reconstruction quality

$$\geq \int dz \, q(z|x) \log p(x|z)$$

Gaussian likelihood $\log p(x|z) \hat{=} \frac{(x - g(z))^2}{\sigma^2}$

$$- \int dz \, q(z|x) \log \left[\frac{q(z|x)}{p(z)} \right]$$

measures similarity between posterior & prior
Kullback-Leibler Divergence

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

$$= \text{ELBO}$$

d) In practice

→ $q_{\phi}(z|x)$ is usually chosen to be a Gaussian with diagonal covariance *
Mean and covariance are parameterized by the encoder network

$$q_{\phi}(z|x) = \mathcal{G}(z | (\mu, \sigma) = f_{\phi}(x))$$

* note that this decorrelates latent variables
→ $p(z)$ can be chosen. Typically a standard Normal distribution

$$p(z) = \mathcal{G}(z | \mu=0, \sigma=1)$$

We want to be able to sample from $p(z)$

→ The ELBO is evaluated stochastically

$$\text{ELBO} = \mathbb{E}_{z \sim q_{\phi}(z|x)} \log p_{\psi}(x|z)$$

$$+ \mathbb{E}_{z \sim q_{\phi}} \log \left[\frac{q_{\phi}(z|x)}{p(z)} \right]$$

Training requires derivative $\frac{\partial \text{ELBO}}{\partial (\psi, \phi)}$

→ How do we take this derivative?

Reparameterization Trick

sample: $z' \sim \mathcal{Q}(z' | 0, 1)$

compute: $z = \sigma_{\phi} z' + \mu_{\phi}$

e) What is the VAE good for?

① The VAE is a generative model

→ it allows us to generate artificial data

→ it allows us to approximately sample from $p(x)$

How?

$$z \sim p(z)$$

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

② Different to the AE, we obtain a continuous latent space.

This allows us to interpolate between data points.

How?

$$z_1 = \mu_{\phi}(x_1) \quad z_2 = \mu_{\phi}(x_2)$$

$$z = z_1 + t \cdot (z_2 - z_1), \quad t \in [0, 1]$$

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

③ VAE allows us to explore the space of likely reconstructions

$$z \sim q_\phi(z|x) \quad x' = g_\psi(z)$$

"posterior analysis"

④ Take into account physical data properties (mask, noise)

\Rightarrow tomorrow's exercise!