

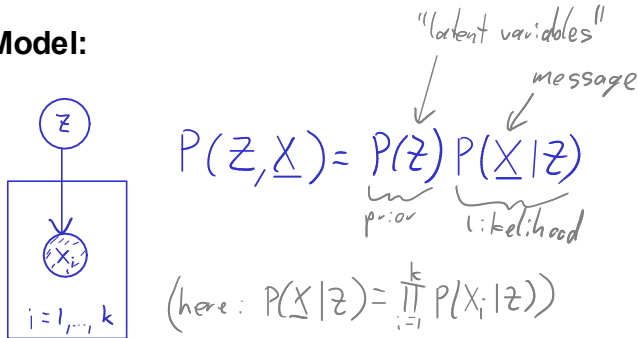
# Bits-Back Coding And Variational Inference

Lecture 8 (23 June 2022); lecturer: Robert Bamler

more course materials online at <https://robamler.github.io/teaching/compress22/>

## Recap from last lecture: Bits-Back Coding With Latent Variable Model

Model:



Note: we want to compress  $\underline{x}$  with marginal model  $P(\underline{x}) = \sum_z P(\underline{x}, z)$

Bits-Back Algorithm (encoder):

1)  $\hat{z} \leftarrow$  decode with ANS from some existing bit string using model  $P(\underline{z} | \underline{x} = \underline{x})$  (posterior)

2) encode message  $\underline{x}$  using ANS and model  $P(\underline{x} | \underline{z} = \hat{z})$

3) encode  $\hat{z}$  using ANS & prior model  $P(\underline{z})$

Net Bit Rate:  $R_{\text{net}}(\underline{x}) = -\log_2 P(\underline{x} = \underline{x} | \underline{z} = \hat{z}) - \log_2 P(\underline{z} = \hat{z}) - (-\log_2 P(\underline{z} = \hat{z} | \underline{x} = \underline{x}))$

$$= -\log_2 \frac{P(\underline{z} = \hat{z}, \underline{x} = \underline{x}) P(\underline{z} = \hat{z})}{P(\underline{z} = \hat{z}) P(\underline{z} = \hat{z}, \underline{x} = \underline{x})} = -\log_2 P(\underline{x} = \underline{x})$$

$\Rightarrow$  optimal

$\Rightarrow$  justifies our use of the posterior in step 1

## Today: Variational Inference

$\rightarrow$  comes from a completely different field of research, unrelated to data compression;

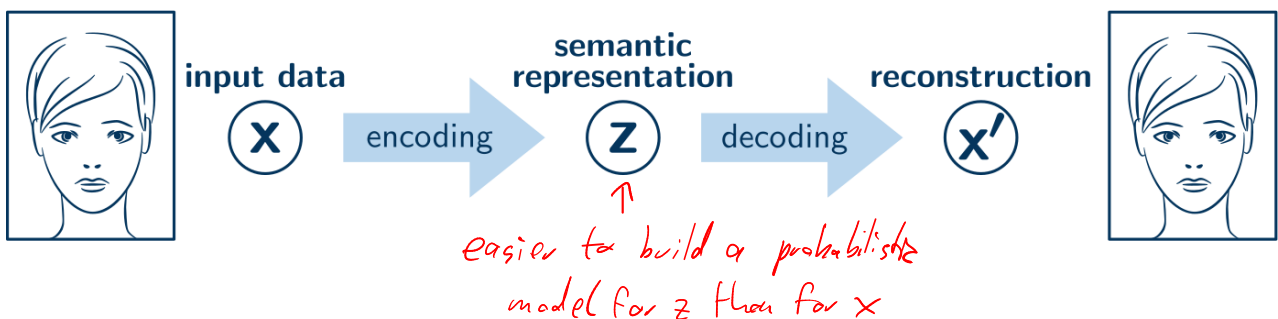
but:

$\rightarrow$  crucial method in modern machine-learning based data compression;

$\rightarrow$  the precise formalism of VI can be motivated most naturally by minimizing the net bit rate of bits-back coding.

### Spoiler: Variational Autoencoders

- $\rightarrow$  popular class of so-called "deep generative models" (use deep neural networks to generate data)
- $\rightarrow$  idea: rather than building a probabilistic model over a complicated message space (e.g., the space of all HD images), design a mapping between the message space and a more abstract semantic representation space and build a probabilistic model over the semantic representation space.



## Back to bits-back coding:

$$R_{\text{net}}(\underline{x}) = -\log_2 P(\underline{x}=\underline{x} | Z=z) - \log_2 P(Z=z) - (-\log_2 P(Z=z | \underline{x}=\underline{x}))$$

$$= -\log_2 \frac{P(Z=z, \underline{x}=\underline{x})}{P(Z=z)} = -\log_2 P(\underline{x}=\underline{x})$$

Problem: obtaining the true posterior is computationally impossible in all but very special models:

$$P(Z | \underline{x}=\underline{x}) = \frac{P(Z, \underline{x}=\underline{x})}{P(\underline{x}=\underline{x})} = \begin{cases} \frac{P(Z, \underline{x}=\underline{x})}{\sum_z P(Z=z, \underline{x}=\underline{x})} \\ \frac{P(Z, \underline{x}=\underline{x})}{\int P(Z=z, \underline{x}=\underline{x}) dz} \end{cases}$$

high-dimensional integral  
(exponentially expensive without tricks)

Idea 1: what if we simply don't use the posterior  $P(Z | \underline{x}=\underline{x})$ , but instead some other distribution  $Q(Z | \underline{x}=\underline{x})$ ?

$$\tilde{R}_{\text{net}}^{(z)}(\underline{x}) = -\log_2 P(\underline{x}=\underline{x} | Z=z) - \log_2 P(Z=z) - (-\log_2 Q(Z | \underline{x}=\underline{x}))$$

$$= -\log_2 \frac{P(Z=z, \underline{x}=\underline{x})}{P(Z=z)} + \log_2 Q(Z=z | \underline{x}=\underline{x}) \stackrel{\text{in general}}{\neq} -\log_2 P(\underline{x}=\underline{x})$$

→ depends on  $z$ , which is out of our control

Recall: if  $Q(Z | \underline{x}=\underline{x}) = P(Z | \underline{x}=\underline{x})$ , then the net bit rate is independent of  $z$  and optimal.

⇒ for any other  $Q(Z | \underline{x}=\underline{x}) \neq P(Z | \underline{x}=\underline{x})$ : <sup>expected</sup> net bit rate is larger

$$\mathbb{E}_{Q(Z | \underline{x}=\underline{x})}[\tilde{R}_{\text{net}}^{(z)}(\underline{x})] \geq R_{\text{net}}(\underline{x}) = -\log_2 P(\underline{x}=\underline{x})$$

→ Problem 8.2 (b): proof of

$$\mathbb{E}_{Q(Z | \underline{x}=\underline{x})}[\tilde{R}_{\text{net}}^{(z)}(\underline{x})] = \underbrace{-\log_2 P(\underline{x}=\underline{x})}_{\text{optimal}} + \underbrace{D_{\text{KL}}(Q(Z | \underline{x}=\underline{x}) \| P(Z | \underline{x}=\underline{x}))}_{\text{overhead} \geq 0}$$

(Expectations here only over  $z$ )

Idea 2: optimize the expected net bit rate over various  $Q(Z | X=x)$

→ parameterize  $Q_\phi(Z | X=\underline{x})$  by some parameters  $\phi$

$$\Rightarrow \boxed{\text{minimize } \mathbb{E}_{Q_\phi(Z | \underline{x}=\underline{x})}[\tilde{R}_{\text{net}}^{(z)}(\underline{x})] \text{ over } \phi}$$

Question: what is the distribution of  $z$  in our modified bits-back algorithm?

1)  $z \leftarrow$  decode with ANS from some existing bit string using model  $Q(z|X=x)$   $\Rightarrow$  claim: thus  $z$  is distributed  $Q(z|X=x)$

$\rightarrow$  Problem Set 9: Prove that decoding a uniformly distributed random bit string with same model  $Q(z|X=x)$   
 $\Leftrightarrow$  drawing  $z \sim Q(z|X=x)$

For historic reasons, one typically talks about maximizing the negative expected net bit rate instead. This is called the Evidence Lower BOUND (ELBO):

$$\begin{aligned} \text{ELBO}(\phi) &= -\mathbb{E}_{Q_\phi(z|X=x)} [\tilde{R}_{\text{net}}^{(z)}(x)] \\ &= \mathbb{E}_{Q_\phi(z|X=x)} [\log P(z, X=x) - \log Q_\phi(z|X=x)] \end{aligned}$$

Problem 8.1 (b):

$$\underbrace{\text{ELBO}(\phi)}_{\text{lower bound on the evidence}} = \underbrace{\log P(X=x)}_{\text{"evidence"}} - \underbrace{D_{\text{KL}}(Q_\phi(z|X=x) \| P(z|X=x))}_{\leq 0}$$

Thus, the following three are equivalent:

minimizing the expected net bit rate of our modified bits-back algorithm

$\Uparrow$

maximizing the ELBO

$\Downarrow$

minimizing  $D_{\text{KL}}(Q_\phi(z|X=x) \| P(z|X=x))$

$\Rightarrow$  result is an approximate posterior

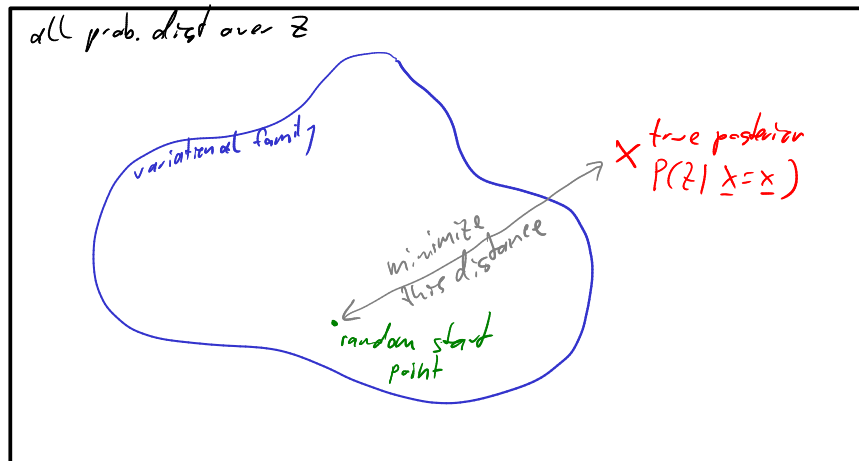
$$Q_{\phi^*}(z|X=x) \approx P(z|X=x)$$

$$\phi^* := \arg \max_{\phi} \text{ELBO}(\phi)$$

use  $Q_\phi(z|X=x)$  instead of posterior

# How Can We Maximize The ELBO?

## 1) Choosing a Variational Family = $\{Q_\phi(z|x=x)\}_\phi$



In practice: often  $z \in \mathbb{R}^d$

$$Q_\phi(z|x=x) = \prod_{i=1}^d Q_{\phi_i}(z_i|x=x)$$

where  $Q_{\phi_i}(z_i|x=x)$  is, for example, a normal distribution with some mean  $\mu_i$  and std. deviation  $\sigma_i$   $\leftarrow \phi_i = (\mu_i, \sigma_i)$

This is called the "mean field approximation" due to an analogy to physics.

## 2) Performing the Maximization

Three methods:

→ "coordinate ascent variational inference (CAVI)": fastest optimization algorithm, but only possible in special models (mostly so-called "conditional conjugate" models; see references)

details: Problem 8.3

- "reparameterization gradients": very simple in practice and relatively widely applicable, but not possible for all variational distributions  $Q$  (in particular, not for discrete  $Q$ )
- "score function gradients" = "REINFORCE method": works also in some cases where reparameterization gradients don't work, but typically slower in practice unless additional tricks are used.

$$ELBO(\phi) = \mathbb{E}_{Q_\phi(z|x=x)} [\log P(z, x=x) - \log Q_\phi(z|x=x)]$$

Goal: find  $\phi^* := \arg \max_{\phi} ELBO(\phi)$

→ stochastic gradient descent

↑  
i.e., estimate  $\mathbb{E}_{z \sim Q} [\dots]$   
by sampling

$g := \nabla_{\phi} ELBO(\phi)$   $S \geq 0$   
update  $\phi \leftarrow \phi + \frac{1}{S} g$   
repeat

Problem: we have to estimate the gradient  $\nabla_{\phi} ELBO(\phi) = \nabla_{\phi} \mathbb{E}_{z \sim Q} [\dots]$

Solutions: Problem 8.2

this dependency on  $\phi$  makes estimation by samples tricky