

## Lecture 8, Part 1:

# Variational Inference

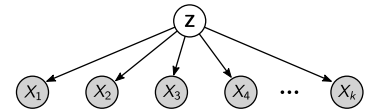
Robert Bamler · Summer Term of 2023

These slides are part of the course “Data Compression With and Without Deep Probabilistic Models” taught at University of Tübingen. More course materials—including video recordings, lecture notes, and problem sets with solutions—are publicly available at <https://robamler.github.io/teaching/compress23/>.

## Recall: Bits-Back Coding & Latent Variable Models



- ▶ **Latent variable model:**  $P(\mathbf{Z}, \mathbf{X}) = P(\mathbf{Z}) \prod_{i=1}^k P(X_i | \mathbf{Z})$



- ▶ **Encoding** a message  $\mathbf{x}$  & side information  $\mathbf{s} \in \{0, 1\}^*$ :

1.  $\mathbf{z} \leftarrow$  decode from  $\mathbf{s}$  with ANS using posterior  $P(\mathbf{Z} | \mathbf{X}=\mathbf{x})$ .
2. Encode  $\mathbf{x}$  using likelihood  $P(\mathbf{X} | \mathbf{Z}=\mathbf{z})$ .
3. Encode  $\mathbf{z}$  using prior  $P(\mathbf{Z})$ .

- ▶ **Net bit rate:**  $R^{\text{net}}(\mathbf{x} | \mathbf{s}) =$

- ▶ optimal & independent of  $\mathbf{s} \implies$  justifies our use of posterior  $P(\mathbf{Z} | \mathbf{X}=\mathbf{x})$  in step 1.

- ▶ **Problem:** calculating the posterior (“Bayesian inference”) is often infeasible:

$$\text{▶ } P(\mathbf{Z} | \mathbf{X}=\mathbf{x}) = \frac{P(\mathbf{Z}, \mathbf{X}=\mathbf{x})}{P(\mathbf{X}=\mathbf{x})} \quad \text{where} \quad P(\mathbf{X}=\mathbf{x}) = \begin{cases} \sum_{\mathbf{z}} P(\mathbf{Z}=\mathbf{z}, \mathbf{X}=\mathbf{x}) & \text{if } \mathbf{Z} \text{ is discrete;} \\ \int p(\mathbf{z}, \mathbf{x}) d\mathbf{z} & \text{if } \mathbf{Z} \text{ is continuous.} \end{cases}$$

## Variational Inference



- ▶ **Problem:** calculating the posterior  $P(\mathbf{Z} | \mathbf{X}=\mathbf{x})$  is often computationally infeasible.

- ▶ **Idea:** use a different distribution  $Q_{\phi}(\mathbf{Z})$  instead of the posterior.

1. Consider the space of *all* probability distributions over  $\mathbf{Z}$ .
2. Choose a subspace  $\mathcal{Q}$  of “simple” distributions.

→ e.g., if  $\mathbf{Z} \in \mathbb{R}^d$ : choose  $\mathcal{Q} := \{Q_{\mu, \sigma}\}_{\mu, \sigma \in \mathbb{R}^d}$

with PDF  $q_{\mu, \sigma}(\mathbf{z}) = \prod_{i=1}^d \mathcal{N}(z_i; \mu_i, \sigma_i^2)$

→ general notation:  $Q_{\phi}(\mathbf{Z})$

3. Find *optimal* variational parameters  $\phi^*$  such that  $Q_{\phi^*}(\mathbf{Z}) \approx P(\mathbf{Z} | \mathbf{X}=\mathbf{x})$  for a given message  $\mathbf{x}$ .



- ▶ **Encoding** a message  $\mathbf{x}$  & side information  $\mathbf{s} \in \{0, 1\}^*$  (copied from first slide):
  1.  $\mathbf{z} \leftarrow$  decode from  $\mathbf{s}$  with ANS using posterior  $P(\mathbf{Z} | \mathbf{X} = \mathbf{x})$ .
  2. Encode  $\mathbf{x}$  using likelihood  $P(\mathbf{X} | \mathbf{Z} = \mathbf{z})$ .
  3. Encode  $\mathbf{z}$  using prior  $P(\mathbf{Z})$ .
- ▶ **Net bit rate:**  $R_{\phi}^{\text{net}}(\mathbf{x} | \mathbf{s}) = -\log_2 P(\mathbf{Z} = \mathbf{z}, \mathbf{X} = \mathbf{x}) + \log_2 Q_{\phi}(\mathbf{Z} = \mathbf{z})$ 
  - ▶ **Naive idea:** find  $\phi^* := \arg \min_{\phi} R_{\phi}^{\text{net}}(\mathbf{x} | \mathbf{s})$ ; then run above encoder using model  $Q_{\phi^*}(\mathbf{Z})$  in step 1.
- ▶ **Decoding:**
  1.  $\mathbf{z} \leftarrow$  decode using prior  $P(\mathbf{Z})$ .
  2.  $\mathbf{x} \leftarrow$  decode using likelihood  $P(\mathbf{X} | \mathbf{Z} = \mathbf{z})$ .
  3. Encode  $\mathbf{z}$  using approximate posterior  $Q_{\phi^*}(\mathbf{Z})$  to reconstruct side information  $\mathbf{s}$ .

**Problem:**  $\phi^*$  depends on  $\mathbf{s}$ .  
 $\Rightarrow$  cyclic dependency

## Expected Net Bit Rate



- ▶ **Idea:** minimize the *expected* net bit rate for random  $\mathbf{s}$  (but still for a fixed message  $\mathbf{x}$ ):
 
$$\phi^* := \arg \min_{\phi} \left( \mathbb{E}_{\mathbf{s}} [R_{\phi}^{\text{net}}(\mathbf{x} | \mathbf{s})] \right) = \arg \min_{\phi} \left( \mathbb{E}_{\mathbf{s}} [ -\log_2 P(\mathbf{Z} = \mathbf{z}, \mathbf{X} = \mathbf{x}) + \log_2 Q_{\phi}(\mathbf{Z} = \mathbf{z}) ] \right)$$
- ▶ What is the distribution of the bit string  $\mathbf{s}$ ?
  - ▶ Generic argument: we have no idea  $\Rightarrow$  uniform distribution
  - ▶ Example: assume  $\mathbf{s}$  is the compressed representation of some previously encoded data.
- ▶ What distribution does this induce for  $\mathbf{z}$ ?
  - ▶ Assume decoding with model  $Q_{\phi}(\mathbf{Z})$  results in a value  $\mathbf{z}$ .
  - ▶ For an optimal coder, only 1 bit string corresponds to  $\mathbf{z}$ .  
 $\Rightarrow$  probability that each one of the  $-\log_2 Q_{\phi}(\mathbf{Z} = \mathbf{z})$  consumed bits matches:

Decoding from a uniform random bit string with a code that is optimal for some model probabilistic model



sampling from the same probabilistic model

## Evidence Lower Bound (ELBO)



- ▶ Minimize the *expected net bit rate* for encoding a given message  $\mathbf{x}$ :
 
$$\phi^* := \arg \min_{\phi} \mathbb{E}_{\mathbf{s}} [R_{\phi}^{\text{net}}(\mathbf{x} | \mathbf{s})] = \arg \min_{\phi} \mathbb{E}_{Q_{\phi}(\mathbf{Z})} [ -\log_2 P(\mathbf{Z}, \mathbf{X} = \mathbf{x}) + \log_2 Q_{\phi}(\mathbf{Z}) ]$$
- ▶ In the (non-compression) literature, one typically *maximizes* the *negative* net bit rate:

$$\phi^* = \arg \max_{\phi} \text{ELBO}(\phi, \mathbf{x}) \quad \text{where} \quad \text{ELBO}(\phi, \mathbf{x}) := \mathbb{E}_{Q_{\phi}(\mathbf{Z})} [ \log P(\mathbf{Z}, \mathbf{X} = \mathbf{x}) - \log Q_{\phi}(\mathbf{Z}) ]$$

- ▶ **Problem 8.1:** derive and interpret three equivalent expressions for the ELBO:
  - ▶ maximum a-posteriori (MAP) + entropy:  $\text{ELBO}(\phi, \mathbf{x}) = \mathbb{E}_{Q_{\phi}(\mathbf{Z})} [ \log P(\mathbf{Z}, \mathbf{X} = \mathbf{x}) ] + H_{Q_{\phi}}[\mathbf{Z}]$
  - ▶ regularized maximum likelihood:  $\text{ELBO}(\phi, \mathbf{x}) = \mathbb{E}_{Q_{\phi}(\mathbf{Z})} [ \log P(\mathbf{X} = \mathbf{x} | \mathbf{Z}) ] - D_{\text{KL}}(Q_{\phi}(\mathbf{Z}) \| P(\mathbf{Z}))$
  - ▶ evidence lower bound:  $\text{ELBO}(\phi, \mathbf{x}) = \log P(\mathbf{X} = \mathbf{x}) - D_{\text{KL}}(Q_{\phi}(\mathbf{Z}) \| P(\mathbf{Z} | \mathbf{X} = \mathbf{x})) \leq \log P(\mathbf{X} = \mathbf{x})$   
 $\Rightarrow Q_{\phi^*}(\mathbf{Z})$  minimizes KL-divergence from true posterior  $\Rightarrow$  "variational inference"

$$\text{ELBO}(\phi, \mathbf{x}) := \mathbb{E}_{\mathbf{z} \sim Q_\phi(\mathbf{z})} [\ell(\phi, \mathbf{x}, \mathbf{z})] \quad \text{where} \quad \ell(\phi, \mathbf{x}, \mathbf{z}) = \log P(\mathbf{Z}=\mathbf{z}, \mathbf{X}=\mathbf{x}) - \log Q_\phi(\mathbf{Z}=\mathbf{z})$$

- ▶ **Goal:** find  $\phi^* := \arg \max_\phi \text{ELBO}(\phi, \mathbf{x})$  (or at least an approximate maximum).
- ▶ **Gradient Descent:** (technically here: *gradient ascent*)
  1. initialize  $\phi \leftarrow \text{random}$
  2. repeat:
    - calculate gradient  $g = \nabla_\phi \text{ELBO}(\phi, \mathbf{x})$
    - update  $\phi \leftarrow \phi + \rho g$  with some “learning rate”  $\rho > 0$
- ▶ **Stochastic Gradient Descent (SGD):** the naive way
  1. initialize  $\phi \leftarrow \text{random}$
  2. repeat:
    - draw a random  $\mathbf{z} \sim Q_\phi(\mathbf{Z})$
    - calculate *gradient estimate*  $\hat{g}(\mathbf{z}) = \nabla_\phi \ell(\phi, \mathbf{x}, \mathbf{z})$
    - update  $\phi \leftarrow \phi + \rho g$  with some (decaying) learning rate  $\rho > 0$
- ▶ **Problem:** SGD only works if  $\hat{g}(\mathbf{z})$  is an *unbiased gradient estimate*: if  $\mathbb{E}_{Q_\phi(\mathbf{z})} [\hat{g}(\mathbf{z})] = g$ .
  - ▶ But we have:  $\mathbb{E}_{Q_\phi(\mathbf{z})} [\hat{g}(\mathbf{z})] =$
  - and:  $g =$

## Black-Box Variational Inference (BBVI)

- ▶ **Problems 8.2 & 8.3:** derive 2 *unbiased* gradient estimators for variational inference:
  - (8.2) **Reparameterization gradients:** [Kingma and Welling, 2014]
 

In SGD, express  $\mathbf{z} \sim Q_\phi(\mathbf{Z})$  as  $\mathbf{z} = f(\epsilon, \phi)$  where  $\epsilon \sim Q_0$ .

$$\implies g = \nabla_\phi \mathbb{E}_{\mathbf{z} \sim Q_\phi(\mathbf{z})} [\ell(\phi, \mathbf{x}, \mathbf{z})] = \nabla_\phi \mathbb{E}_{\epsilon \sim Q_0} [\ell(\phi, \mathbf{x}, f(\epsilon, \phi))] = \mathbb{E}_{\epsilon \sim Q_0} [\nabla_\phi \ell(\phi, \mathbf{x}, f(\epsilon, \phi))]$$

😊 Typically low gradient variance  $\mathbb{E}[(\hat{g} - g)^2] \implies$  fast(-ish) convergence of SGD.

😞 Doesn't work for discrete  $\mathbf{Z}$  (without an additional approximation called “Gumbel-softmax”).
  - (8.3) **Score function gradients** (aka “REINFORCE method”): [Ranganath et al., 2014]
 
$$\hat{g}(\mathbf{z}) = (\nabla_\phi \log Q_\phi(\mathbf{Z}=\mathbf{z})) \ell(\phi, \mathbf{x}, \mathbf{z})$$

😞 Typically higher gradient variance than reparameterization gradients.

😊 Works for all (explicit) variational families.
- ▶ **Remark:** for some models,  $\mathbb{E}_{\mathbf{z} \sim Q_\phi(\mathbf{z})} [\ell(\phi, \mathbf{x}, \mathbf{z})]$  can be calculated analytically.
 

$\longrightarrow$  *much* faster optimization with “coordinate ascent variational inference” (CAVI) [Blei et al. (2017)]

## Outlook

- ▶ **Next Week:**
  - ▶ How can we learn the generative model?  $\rightarrow$  variational expectation maximization
  - ▶ How can we learn to do inference?  $\rightarrow$  amortized variational inference
  - ▶ **Combined:** Variational autoencoders (VAEs)
- ▶ **Afterwards:** lossy data compression (in theory & in practice)