Lecture 9:
# Variational Autoencoders

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: Variational Inference

▶ **Idea:**
  ▶ approximate the (inaccessible) true posterior $P(\mathbf{Z} \mid \mathbf{X}=\mathbf{x})$ with a *variational distribution* $Q_\phi(\mathbf{Z})$.
  ▶ Find the best approximation $\phi^* := \arg\max_\phi \mathrm{ELBO}(\phi, \mathbf{x})$.

▶ **Evidence Lower Bound:** $\boxed{\mathrm{ELBO}(\phi, \mathbf{x}) = \mathbb{E}_{Q_\phi(\mathbf{Z})}\big[\log P(\mathbf{Z}, \mathbf{X}=\mathbf{x}) - \log Q_\phi(\mathbf{Z})\big]}$
  ▶ negative expected net bit rate of bits-back coding: $\mathrm{ELBO}(\phi, \mathbf{x}) = -\mathbb{E}_{\mathbf{s}}\big[R_\phi^{\mathrm{net}}(\mathbf{x} \mid \mathbf{s})\big]$
  ▶ bound on the evidence: $\mathrm{ELBO}(\phi, \mathbf{x}) = \log P(\mathbf{X}=\mathbf{x}) - D_{\mathrm{KL}}\big(Q_\phi(\mathbf{Z}) \,\|\, P(\mathbf{Z} \mid \mathbf{X}=\mathbf{x})\big) \leq \log P(\mathbf{X}=\mathbf{x})$
  ▶ regularized maximum likelihood: $\mathrm{ELBO}(\phi, \mathbf{x}) = \mathbb{E}_{Q_\phi(\mathbf{Z})}\big[\log P(\mathbf{X}=\mathbf{x} \mid \mathbf{Z})\big] - D_{\mathrm{KL}}\big(Q_\phi(\mathbf{Z}) \,\|\, P(\mathbf{Z})\big)$
  ▶ ~~today:~~ rate/distortion-tradeoff: $\boxed{\mathrm{ELBO}_\beta(\phi, \mathbf{x}) = \mathbb{E}_{Q_\phi(\mathbf{Z})}\big[\log P(\mathbf{X}=\mathbf{x} \mid \mathbf{Z})\big] - \beta\, D_{\mathrm{KL}}\big(Q_\phi(\mathbf{Z}) \,\|\, P(\mathbf{Z})\big)}$
    *(actually, next week ☺)*

▶ **Problems:**
  ▶ What's the generative model $P(\mathbf{Z}, \mathbf{X})$? $\longrightarrow$ variational expectation maximization
  ▶ Expensive "$\arg\max_\phi$" for *each message* $\mathbf{x}$ in both encoder & decoder. $\longrightarrow$ amortized inference

---

# Part 1: Learning the Generative Model

▶ **Goal:** learn optimal parameters $\theta^*$ of the *generative* model $P_\theta(\mathbf{Z}, \mathbf{X}) = P_\theta(\mathbf{Z})\, P_\theta(\mathbf{X} \mid \mathbf{Z})$.
  ▶ Thus, the ELBO now depends on $\theta$, i.e., $\mathrm{ELBO}(\theta, \phi, \mathbf{x}) = Q_\phi(\mathbf{Z})\big[\log P_\theta(\mathbf{Z}, \mathbf{X}=\mathbf{x}) - \log Q_\phi(\mathbf{Z})\big]$

  ▶ **Example:** data $\mathbf{X} = (X_i)_i$ are binarized images, i.e., each $X_i$ is a pixel value $\in \{0, 1\}$.
    $\rightarrow$ Prior is fixed: $P(\mathbf{Z}=\mathbf{z}) = \mathcal{N}(\mathbf{z}; 0, I)$ (standard normal distribution)
    $\rightarrow$ Likelihood is parameterized by a (deconvolutional) neural network $g_\theta$: *⎱ see class "Decoder Model" in jupyter notebook.*
       $P_\theta(\mathbf{X} \mid \mathbf{Z}) = \prod_i P_\theta(X_i \mid \mathbf{Z})$ with $P_\theta(X_i=1 \mid \mathbf{Z}=\mathbf{z}) = \sigma(g_{\theta,i}(\mathbf{z}))$

▶ **Distinguish:**
  ▶ *global* parameters $\theta^*$ ("model parameters"):
    $\rightarrow$ specify the *generative model* $P_{\theta^*}(\mathbf{Z}, \mathbf{X})$
    $\rightarrow$ same for all data points $\mathbf{x} \implies$ known to both sender & receiver
  ▶ *local* parameters $\phi^*$ ("variational parameters"):
    $\rightarrow$ specify an approximation $Q_{\phi^*}(\mathbf{Z})$ to the posterior $P_{\theta^*}(\mathbf{Z} \mid \mathbf{X}=\mathbf{x})$ for a *specific data point* $\mathbf{x}$
    $\rightarrow$ different for each data point $\mathbf{x} \implies$ not available to the receiver until it has decoded $\mathbf{x}$

# Variational Expectation Maximization

1. **In order to develop a new compression method:**
   - learn optimal parameters $\theta^*$ of the generative model $P_\theta(\mathbf{Z}, \mathbf{X})$:

   <span style="color:blue">$\phi^*(\vartheta, \underline{x}) = \arg\max_\phi \text{ELBO}(\vartheta, \phi, \underline{x})$</span>

   <span style="color:blue">$\vartheta^* \leftarrow \arg\max_\vartheta \mathbb{E}_{\underline{x} \sim \text{training set}}\left[\text{ELBO}(\vartheta, \overbrace{\phi^*(\vartheta, \underline{x})}, \underline{x})\right]$</span>

   <span style="color:blue">$= \arg\max_\vartheta \mathbb{E}_{\underline{x} \sim \text{training set}}\left[\max_\phi \text{ELBO}(\vartheta, \phi, \underline{x})\right]$</span>

2. **Share the learned generative model $P_{\theta^*}(\mathbf{Z}, \mathbf{X})$ between sender & receiver.**

3. **In deployment:** encode / decode a given data point $\mathbf{x}$
   - Use entropy model $Q_{\phi^*}(\mathbf{Z})$.

   <span style="color:blue">$\phi^*$ depends on $\underline{x}$, so both encoder & decoder need to find it by running (expensive) SGD:</span>

   <span style="color:blue">$\phi^* \leftarrow \arg\max_\phi \text{ELBO}(\vartheta^*, \phi, \underline{x})$</span>

   <span style="color:red">training algorithm:</span>
   <span style="color:red">initialize $\vartheta \leftarrow$ random</span>
   <span style="color:red">repeat until convergence:</span>
   <span style="color:red">draw $\underline{x} \sim$ training set</span>
   <span style="color:red">initialize $\phi \leftarrow$ prior or random</span>
   <span style="color:red">repeat until convergence:</span>
   <span style="color:red">update $\phi \leftarrow \phi + \varsigma_1 \nabla_\phi \text{ELBO}(\vartheta, \phi, \underline{x})$</span>
   <span style="color:red">update $\vartheta \leftarrow \vartheta + \varsigma_2 \nabla_\vartheta \text{ELBO}(\vartheta, \phi, \underline{x})$</span>

   <span style="color:magenta">expensive inner loop for each training step on $\vartheta$</span>

# Part 2: Learning How to Do Inference (Fast)

- **Problems:**
  1. Learning the generative model requires an expensive inner loop *for every training step.*
  2. Expensive optimization over $\phi$ for *each message* $\mathbf{x}$ *we want compress / decompress.*

- **Solution:** *amortized* variational inference
  - learn a mapping $f$ from $\mathbf{x}$ to variational parameters such that setting $\phi \leftarrow f(\mathbf{x})$ approximately maximizes $\text{ELBO}(\theta^*, \phi, \mathbf{x})$ for a given $\mathbf{x}$.
  - **Notation:** inference network $f_\phi(\mathbf{x})$; variational distribution $Q_\phi(\mathbf{Z} \mid \mathbf{X} = \mathbf{x})$ <span style="color:blue">← in the notation we've used so far, this would be $Q_{f_\phi(\underline{x})}(\underline{Z})$</span>
  - **Example:** Gaussian mean field variational distribution:
    → inference network $f_\phi(\mathbf{x}) = (\boldsymbol{\mu}_\phi(\mathbf{x}), \log \boldsymbol{\sigma}_\phi^2(\mathbf{x}))$ outputs means and (log) variances
    → these parameterize a variational distribution $Q_\phi(\mathbf{Z} \mid \mathbf{x}) = \mathcal{N}\big(\boldsymbol{\mu}_\phi(\mathbf{x}), \text{diag}(\sigma_{\phi,1}^2(\mathbf{x}), \ldots, \sigma_{\phi,k}^2(\mathbf{x}))\big)$

  <span style="color:blue">training algorithm now:</span>
  <span style="color:red">initialize $\vartheta, \phi \leftarrow$ random</span>
  <span style="color:red">repeat until convergence:</span>
  <span style="color:red">draw $\underline{x} \sim$ training set</span>
  <span style="color:red">update $(\vartheta, \phi) \leftarrow (\vartheta, \phi) + \varsigma \nabla_{\vartheta,\phi} \text{ELBO}(\vartheta, \phi, \underline{x})$</span>

  <span style="color:blue">ELBO for amortized inference:</span>
  <span style="color:blue">$\text{ELBO}(\vartheta, \phi, \underline{x}) = \mathbb{E}_{Q_\phi(\underline{Z}\mid\underline{X}=\underline{x})}\left[\log P(\underline{Z}, \underline{X}=\underline{x}) - \log Q_\phi(\underline{Z}\mid\underline{X}=\underline{x})\right]$</span>

  <span style="color:green">⟹ no expensive inner loop; $\vartheta$ and $\phi$ are learned concurrently</span>

# Variational Autoencoders (VAEs)

Combine variational expectation maximization with amortized variational inference.

- **Lossless compression with variational autoencoders:**
  - use bits-back trick → Problem 9.2

- **Lossy compression with variational autoencoders:**
  - **Example:** data $\mathbf{X} = (X_i)_i$ are color images, i.e., each $X_i$ is a continuous RGB value $\in [0, 1]$.
    → Prior may be learned, e.g.: $P_\theta(\mathbf{Z} = \mathbf{z}) = \mathcal{N}\big(\mathbf{z}; 0, \text{diag}(\sigma_1^2, \ldots, \sigma_{\text{num\_channels}}^2)^{\otimes \text{spatial\_dim}}\big)$
    → Likelihood is parameterized by a (deconvolutional) neural network $g_\theta$:
    $P_\theta(\mathbf{X} \mid \mathbf{Z}) = \prod_i P_\theta(X_i \mid \mathbf{Z})$ with density function $p_\theta(x_i \mid \mathbf{Z} = \mathbf{z}) = \mathcal{N}(x_i; g_{\theta,i}(\mathbf{z}), \frac{\beta}{2} I)$
  - **Idea:** just use $g_{\theta,i}(\mathbf{z})$ as the reconstruction of an image.
    (Don't bother using the likelihood $P_\theta(\mathbf{X} \mid \mathbf{Z} = \mathbf{z})$ to encode the true image.)

    <span style="color:magenta">$\beta$ controls a "rate/distortion trade-off" (see next week)</span>
  - Likelihood no longer has a probabilistic meaning. But $-\log P_\theta(\mathbf{X} \mid \mathbf{Z} = \mathbf{z})$ is a *distortion* metric.
    ⟹ ELBO becomes a *rate-distortion* trade-off
  - *next week*