

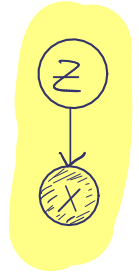
Variational Autoencoders & Lossy Neural Compression

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

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

Recap from last lecture: Variational Inference (VI)

- latent variable model: $P(Z, X) = P(Z)P(X|Z)$
- goal: approximate the posterior: $P(Z|X=x) = \frac{P(Z)P(X=x|Z)}{\int P(Z=z)P(X=x|Z=z)dz}$
- VI turns the inference problem into an optimization problem



- variational distribution: $Q_\phi(Z|X=x)$
← variational parameters

Evidence lower bound (ELBO):

$$\text{ELBO}(\phi) = -\mathbb{E}_{Q_\phi(Z|X=x)}[\tilde{R}_{\text{net}}^{(z)}(\underline{x})] \quad \leftarrow \text{how we motivated it}$$

Problem Set 8

$$\begin{cases} = \mathbb{E}_{Q_\phi(Z|X=x)} [\log P(Z, X=x) - \log Q_\phi(Z|X=x)] & \leftarrow \text{most explicit formulation ("regularized MAP")} \\ & \rightarrow H[Q_\phi(Z|X=x)] \\ = \mathbb{E}_{Q_\phi(Z|X=x)} [\log P(X=x|Z)] - D_{\text{KL}}(Q_\phi(Z|X=x) \| P(Z)) & \leftarrow \text{"regularized maximum likelihood"} \\ = \log P(X=x) - D_{\text{KL}}(Q_\phi(Z|X=x) \| P(Z|X=x)) & \leftarrow \text{connects VI to Bayesian inference} \end{cases}$$

How to maximize the ELBO with stochastic gradient optimization:

- reparameterization gradients: $\nabla_\phi \mathbb{E}_{Q_\phi(Z|X=x)} [\ell(z, \phi)] = \nabla_\phi \mathbb{E}_{\epsilon \sim Q_0} [\ell(g(\epsilon, \phi), \phi)]$
 $z = g(\epsilon, \phi)$ where $\epsilon \sim Q_0$ ← fixed distribution
- score function gradients: $\nabla_\phi \mathbb{E}_{Q_\phi(Z|X=x)} [\ell(z, \phi)] = \mathbb{E}_{Q_\phi(Z|X=x)} [\nabla_\phi \log Q_\phi(Z|X=x) \times \ell(z, \phi)]$
(= REINFORCE method)

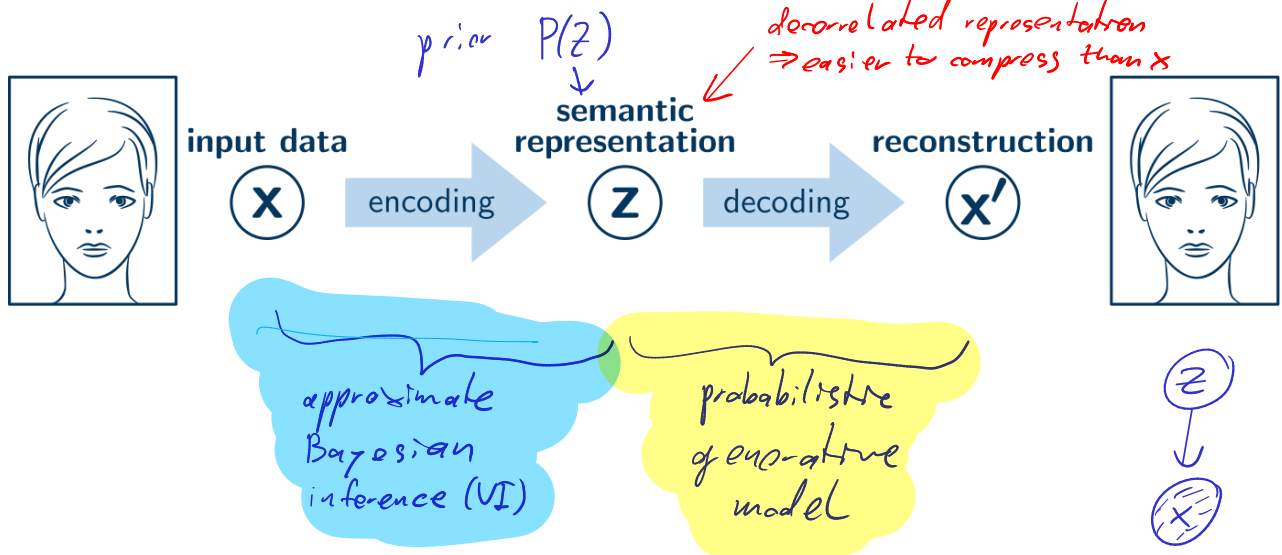
Limitations so far:

- the generative model $P(Z, X)$ is fixed -- and therefore limited to simple models that we can come up with manually; and
- for every concrete message x that we want to compress, we have to run an expensive optimization procedure to find the optimal variational parameters ϕ^* .

TODAY: overcoming these limitations

- Variational Expectation Maximization (Variational EM) → "learn the prob. gen. model P from training data"
- Amortized variational inference → "learn how to do inference"

Spoiler: Variational Autoencoders (VAEs) = amortized variational expectation maximization



Variational Expectation Maximization: learning a latent variable model

[Beal & Ghahramani, Bayesian statistics, 2003]

Introduce **free parameters** into the probabilistic model $P(Z, X)$:

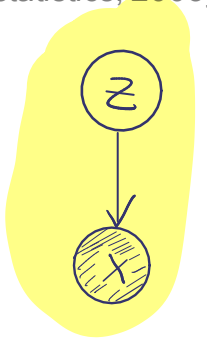
$$P_{\theta}(Z, X) = P_{\theta}(Z) P_{\theta}(X|Z)$$

model parameters θ

$$= \mathbb{E}_{X \sim \text{train set}} [-\log P_{\theta}(X=x)]$$

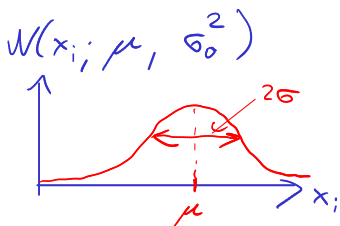
$$= \mathbb{E}_{X \sim \text{train set}} [-\log \left(\sum_z P_{\theta}(Z=z, X=x) \right)]$$

probabilistically expensive



For example, the likelihood could be **parameterized by a neural network with weights θ** :

$$P_{\theta}(X|Z=z) = \mathcal{N}(g_{\theta}(z), \sigma_0^2 I)$$



"a normal distribution with mean $g_{\theta}(z)$ and variance σ_0^2 in each coordinate direction"

where $g_{\theta}: \text{latent space} \rightarrow \text{data space}$ is a neural network with trainable weights θ

Thus, the ELBO now depends both on the variational parameters ϕ and on the model parameters θ :

$$\text{ELBO}(\theta, \phi) = -\mathbb{E} R_{\text{net}}$$

$$= \mathbb{E} [\log P - \log Q]$$

$$= \log P_{\theta}(x=x) - D_{\text{KL}}$$

- optimal bitrate with model P_{θ}

overhead due to VI

• maximize over both of these jointly (at training time)

• at compression time: keep θ fixed & maximize only over ϕ

minimize the expected ^{net} bitrate of modified bits-back

$$\Rightarrow \text{maximize: } ELBO(\phi, \vartheta) = - \mathbb{E}_{Q_{\phi}(z|x)} [\tilde{R}_{\text{net}}^{(\vartheta)}(x)]$$

\nwarrow variational params \nwarrow model params

Alternative:

- store $\phi_x \quad \forall x \in \text{train set}$ on disk

- training loop:

for training-step in $\{1, 2, 3, \dots, n\}$:

sample $(x) \sim \text{train set}$

look up ϕ_x on disk

calculate $S_{\vartheta} = \nabla_{\vartheta} ELBO(\vartheta, \phi_x, x)$

$S_{\phi} = \nabla_{\phi_x} ELBO(\vartheta, \phi_x, x) \leftarrow$

update $\vartheta \leftarrow \vartheta + S_{\vartheta}$

$\phi \leftarrow \phi + S_{\phi} \leftarrow$

store ϕ_x back to disk

Data compression with learnt latent variable models (try 1: without amortization):

1) When designing the compression method:

- collect large (unlabeled) data set of training samples (e.g., a large collection of images)
- come up with a model architecture for the generative model P_θ (that still has free parameters)
- train the model by maximizing the ELBO jointly over both θ and ϕ .

in detail: $\theta^*, \phi^* := \arg \max_{\theta, \phi} \mathbb{E}_{x \sim \text{train set}} [\text{ELBO}(\theta, \phi, x)]$

for training step $p \in \{1, 2, \dots, N\}$
 find $\phi_x^k := \arg \max_{\phi} \text{ELBO}(\theta, \phi, x)$
 set $\hat{\theta} := \nabla_{\theta} \text{ELBO}(\cdot)$
 update $\theta \leftarrow \theta + \eta \hat{\theta}$

- throw away ϕ^* and share θ^* between sender and receiver

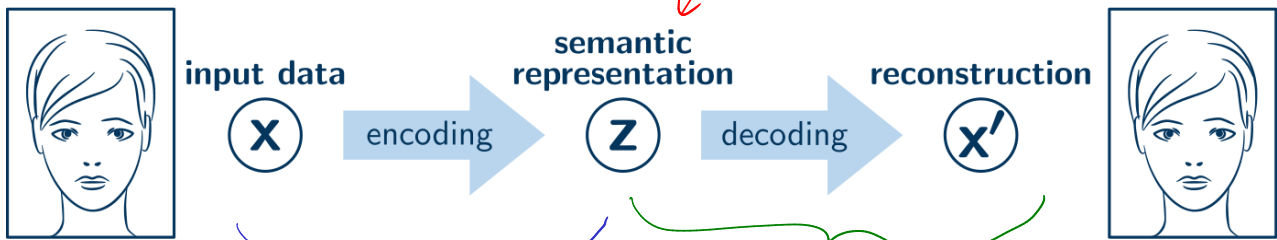
ϕ_x are the variational params for training point x

2) When compressing some given data x (i.e., on the sender side):

- perform variational inference, i.e., maximize $\text{ELBO}(\theta^*, \phi, x)$ over ϕ but keep θ^* fixed at the agreed-upon values.
- use probabilistic generative model $P_\theta(Z, X)$ and the resulting variational distribution $Q_\phi(Z | X=x)$ to compress x .

3) When decompressing data (i.e., on the receiver side):

- needs the exact same probabilistic generative model $P_\theta(Z, X)$ that the sender used for compression.
- if the data was compressed with bits-back coding, then the receiver also needs to perform variational inference once it has reconstructed x (i.e., maximize $\text{ELBO}(\theta, \phi)$ over ϕ but keep θ^*



compress z using prior $P_{\theta^*}(z)$

$Q_{\phi^*}(Z | X=x)$
 where $\phi^* := \arg \max_{\phi} \text{ELBO}(\phi, \theta^*, x)$
 expensive operation

$P_{\theta^*}(X | Z)$
 where θ^* is known at (de-)compression time

Amortized Variational Inference: learn how to do inference

Variational inference maps data x to a variational distribution $Q_\phi(Z | X=x)$:

Idea: learn this mapping from x to $Q_\phi(Z | X=x)$:

- rename the parameters of $Q_\lambda(Z)$ from ϕ to λ
- rather than optimizing over λ , learn a function f that maps x to λ (and that is parameterized by some neural network weights ϕ):

for example: $Q_\lambda(z) = \mathcal{N}(\mu, \text{diag}(\sigma^2))$
 $\lambda = (\mu, \sigma^2)$

$\lambda = f_\phi(x)$ where f_ϕ is a neural network with weights ϕ

$Q_\phi(z|x=x) := Q_\lambda(z)$
 where $\lambda = f_\phi(x)$

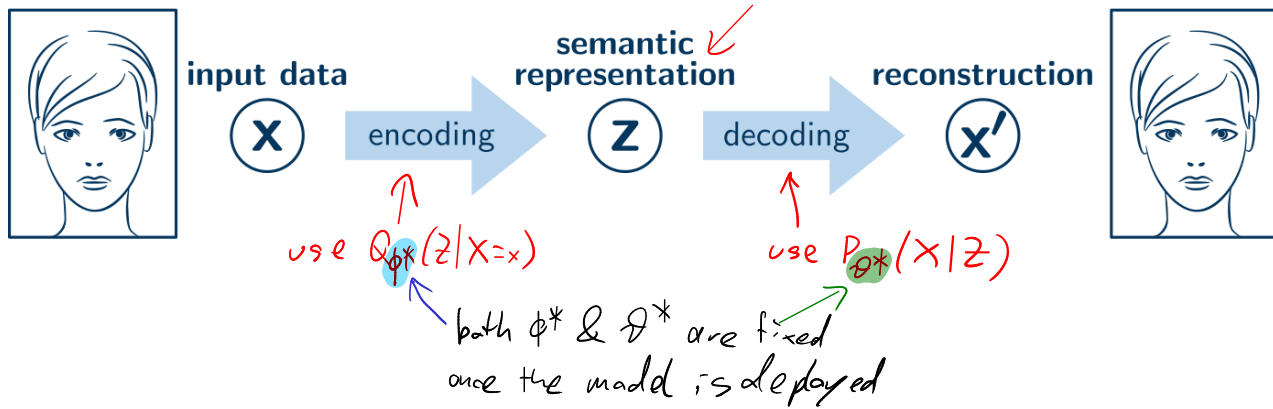
\Rightarrow variational parameters ϕ are now shared between all data points x

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

Combining amortized inference with variational expectation maximization results in:

Variational Autoencoders (VAEs):

[Kingma and Welling, 2015]



training objective: $E[\text{ELBO}] = \dots$

Problem Set: implement a variational autoencoder for simple images (MNIST)

Note: Variational expectation maximization (EM) is not limited to VAEs. Even without amortized inference, variational EM is a very useful algorithm that is very simple and allows you to treat some model parameters (Z) probabilistically while using point estimates for others (θ).