

COMPSCI 688: Probabilistic Graphical Models

Lecture 21: Learning with Stochastic Variational Inference

Dan Sheldon

Manning College of Information and Computer Sciences
University of Massachusetts Amherst

Partially based on materials by Benjamin M. Marlin (marlin@cs.umass.edu) and Justin Domke (domke@cs.umass.edu)

Black-Box Stochastic Variational Inference

Review: Variational Auto-Encoder

Factor analysis model with non-linear mapping

$$p(\mathbf{z}) = \mathcal{N}(\mathbf{z}; 0, I)$$
$$p(x_j | \mathbf{z}) = \text{Bernoulli}(x_j; (f_\theta(\mathbf{z}))_j), \quad j = 1, \dots, d$$

Example non-linear mapping:

$$f_\theta(\mathbf{z}) = h_2(\mathbf{b}_2 + \mathbf{W}_2 \cdot h_1(\mathbf{b}_1 + \mathbf{W}_1 \mathbf{z}))$$

Exact inference and learning are *intractable*.

Stochastic Variational Inference

Choose variational family, e.g., diagonal Gaussian, to approximate posterior:

$$q_\phi(\mathbf{z}) = \mathcal{N}(\boldsymbol{\mu}, \text{diag}(\boldsymbol{\sigma}^2)), \quad \phi = (\boldsymbol{\mu}, \boldsymbol{\sigma})$$

Stochastic optimization: repeatedly get **unbiased gradient estimate** $\hat{\nabla}_\phi$, update ϕ :

$$\hat{\nabla}_\phi \approx \nabla_\phi \text{ELBO}(\phi)$$

$$\phi \leftarrow \phi + \alpha \hat{\nabla}_\phi$$

How to get $\hat{\nabla}_\phi$?

Warmup: Estimating the ELBO

It's easy to estimate the ELBO via Monte Carlo samples:

$$\text{ELBO}(\phi) \approx \frac{1}{K} \sum_{i=1}^K \log \frac{p(z^{(i)}, x)}{q_\phi(z^{(i)})} \quad z^{(i)} \sim q_\phi$$

This is unbiased: expected value of RHS = ELBO(ϕ) for any value of K .

Estimating the Gradient? (False Start)

What happens if we try to estimate the gradient the same way? We want:

$$\nabla_\phi \text{ELBO}(\phi) = \nabla_\phi \mathbb{E}_{q_\phi} \left[\log \frac{p(Z, x)}{q_\phi(Z)} \right]$$

Consider the estimate:

$$\nabla_\phi \log \frac{p(z, x)}{q_\phi(z)}, \quad z \sim q_\phi.$$

This is not unbiased! It neglects the fact that the *distribution* of z depends on ϕ :

$$\nabla_\phi \text{ELBO}(\phi) = \nabla_\phi \int q_\phi(z) \log \frac{p(z, x)}{q_\phi(z)} dz$$

The false start incorrectly interchanges the gradient and the expectation:

$$\nabla_\phi \mathbb{E}_{q_\phi} \left[\log \frac{p(Z, x)}{q_\phi(Z)} \right] \neq \mathbb{E}_{q_\phi} \left[\nabla_\phi \log \frac{p(Z, x)}{q_\phi(Z)} \right]$$

Reparameterization Trick

The reparameterization trick is a way to convert the ELBO into an expectation with respect to a *fixed* distribution (independent of ϕ) so we can interchange the gradient and expectation. The idea is to draw samples of z by transforming a random variable from a fixed base distribution.

Example: $z = \mu + \sigma\epsilon$, $\epsilon \sim \mathcal{N}(0, 1) \implies z \sim \mathcal{N}(\mu, \sigma^2)$

General case: $z = \mathcal{T}_\phi(\epsilon)$, $\epsilon \sim q(\epsilon) \implies z \sim q_\phi(z)$

We call \mathcal{T}_ϕ and $q(\epsilon)$ a *reparameterization* of q_ϕ .

ELBO Gradient with Reparameterization

With reparameterization, we can write the ELBO as an expectation over $q(\epsilon)$:

$$\text{ELBO}(\phi) = \mathbb{E}_{q(\epsilon)} \left[\log \frac{p(\mathcal{T}_\phi(\epsilon), x)}{q_\phi(\mathcal{T}_\phi(\epsilon))} \right]$$

Now we can interchange the gradient and expectation

$$\nabla_\phi \text{ELBO}(\phi) = \nabla_\phi \mathbb{E}_{q(\epsilon)} \left[\log \frac{p(\mathcal{T}_\phi(\epsilon), x)}{q_\phi(\mathcal{T}_\phi(\epsilon))} \right] = \mathbb{E}_{q(\epsilon)} \left[\nabla_\phi \log \frac{p(\mathcal{T}_\phi(\epsilon), x)}{q_\phi(\mathcal{T}_\phi(\epsilon))} \right]$$

Reparameterization Gradient Estimate

This gives a simple unbiased Monte Carlo estimate of the gradient:

$$g = \nabla_\phi \left(\frac{1}{K} \sum_{i=1}^K \log \frac{p(\mathcal{T}_\phi(\epsilon^{(i)}), x)}{q_\phi(\mathcal{T}_\phi(\epsilon^{(i)}))} \right), \quad \epsilon^{(1)}, \dots, \epsilon^{(K)} \sim q(\epsilon)$$

We can compute it as follows:

1. Draw $\epsilon^{(1)}, \dots, \epsilon^{(K)} \sim q(\epsilon)$
2. Compute $\widehat{\text{ELBO}}(\phi, \epsilon^{(1)}, \dots, \epsilon^{(K)}) =$ term in parentheses above
3. Use autodiff to get $g = \nabla_\phi \widehat{\text{ELBO}}(\phi, \epsilon^{(1:K)})$

Gradient Estimation: Reparameterization

	Without reparameterization	With reparameterization
Variational distribution	$q_\phi(z)$	$q_\phi(z)$
Sampling	$z \sim q_\phi(z)$	$\epsilon \sim q(\epsilon), z = \mathcal{T}_\phi(\epsilon)$
ELBO	$\mathbb{E}_{q_\phi(Z)} \left[\log \frac{p(Z, x)}{q_\phi(Z)} \right]$	$\mathbb{E}_{q(\epsilon)} \left[\log \frac{p(\mathcal{T}_\phi(\epsilon), x)}{q_\phi(\mathcal{T}_\phi(\epsilon))} \right]$
ELBO estimate	$\log \frac{p(z, x)}{q_\phi(z)}, z \sim q_\phi(z)$	$\log \frac{p(\mathcal{T}_\phi(\epsilon), x)}{q_\phi(\mathcal{T}_\phi(\epsilon))}, \epsilon \sim q(\epsilon)$
Gradient estimate	$\nabla_\phi \log \frac{p(z, x)}{q_\phi(z)}, z \sim q_\phi(z)$ (wrong/biased)	$\nabla_\phi \log \frac{p(\mathcal{T}_\phi(\epsilon), x)}{q_\phi(\mathcal{T}_\phi(\epsilon))}, \epsilon \sim q(\epsilon)$ (unbiased)

Reparameterization with Diagonal Gaussians

Suppose the variational family is a diagonal Gaussian

$$q_\phi(\mathbf{z}) = \mathcal{N}(\boldsymbol{\mu}, \text{diag}(\boldsymbol{\sigma}^2))$$

This can be reparameterized as:

$$\mathbf{z} = \boldsymbol{\mu} + \boldsymbol{\sigma} \odot \epsilon, \quad \epsilon \sim \mathcal{N}(0, I)$$

(\odot = elementwise multiplication)

Aside: Reparameterization with Arbitrary Gaussians

Another choice would be to use a general Gaussian distribution:

$$\epsilon \sim \mathcal{N}(0, I) \implies \mu + L\epsilon \sim \mathcal{N}(\mu, LL^\top).$$

This is a reparameterization with

$$q(\epsilon) = \mathcal{N}(\epsilon|0, I), \quad \mathcal{T}_\phi(\epsilon) = \mu + L\epsilon \quad \phi = (L, \mu)$$

It covers any multivariate Gaussian, since an arbitrary covariance matrix Σ can be written as $\Sigma = LL^\top$ for some L (e.g., a Cholesky factor)

Example: Bernoulli VAE

Let's return to our Bernoulli VAE factor analysis model and use a diagonal Gaussian approximation:

$$\begin{aligned} p(\mathbf{z}) &= \mathcal{N}(\mathbf{z}; 0, I) \\ p(x_j|\mathbf{z}) &= \text{Bernoulli}(x_j; (f_\theta(\mathbf{z}))_j), \quad j = 1, \dots, d \\ q_\phi(\mathbf{z}) &= \mathcal{N}(\mathbf{z}; \boldsymbol{\mu}, \text{diag}(\boldsymbol{\sigma}^2)) \end{aligned}$$

BBSVI would repeat the following steps:

$$\epsilon \sim \mathcal{N}(0, I) \qquad (\mathbf{z} = \boldsymbol{\mu} + \boldsymbol{\sigma} \odot \epsilon)$$

$$\hat{\nabla}_{\boldsymbol{\mu}, \boldsymbol{\sigma}} = \nabla_{\boldsymbol{\mu}, \boldsymbol{\sigma}} \left\{ \begin{aligned} &\log \mathcal{N}(\boldsymbol{\mu} + \boldsymbol{\sigma} \odot \epsilon; 0, I) \\ &+ \sum_{j=1}^d \log \text{Bernoulli}(x_j; (f_\theta(\boldsymbol{\mu} + \boldsymbol{\sigma} \odot \epsilon))_j) \\ &- \log \mathcal{N}(\boldsymbol{\mu} + \boldsymbol{\sigma} \odot \epsilon; \boldsymbol{\mu}, \text{diag}(\boldsymbol{\sigma}^2)) \end{aligned} \right\}$$

$$(\boldsymbol{\mu}, \boldsymbol{\sigma}) \leftarrow (\boldsymbol{\mu}, \boldsymbol{\sigma}) + \alpha \cdot \hat{\nabla}_{\boldsymbol{\mu}, \boldsymbol{\sigma}}$$

With the optimized parameters we could approximate $p(\mathbf{z}|\mathbf{x}) \approx q_\phi(\mathbf{z})$ and lower bound the log-marginal likelihood $\log p(\mathbf{x}) \geq \text{ELBO}(\phi)$. **What about learning?**

Learning with Stochastic Variational Inference

Learning with Stochastic Variational Inference

The basic idea is to jointly maximize the ELBO with respect to model parameters θ and variational parameters ϕ by getting unbiased gradient estimates for both:

$$\log p_{\theta}(\mathbf{x}) \geq \text{ELBO}(\theta, \phi) = \mathbb{E}_{q_{\phi}} \left[\log \frac{p_{\theta}(Z, \mathbf{x})}{q_{\phi}(Z)} \right]$$

$$\hat{\nabla}_{\theta} \approx \nabla_{\theta} \text{ELBO}(\theta, \phi)$$

$$\hat{\nabla}_{\phi} \approx \nabla_{\phi} \text{ELBO}(\theta, \phi)$$

$$(\theta, \phi) \leftarrow (\theta, \phi) + \alpha \cdot (\hat{\nabla}_{\theta}, \hat{\nabla}_{\phi})$$

Learning with IID Data

How do we learn a latent variable model $p_{\theta}(\mathbf{z}, \mathbf{x})$ when we have iid data $\mathbf{x}^{(1)}, \dots, \mathbf{x}^{(N)}$?

Each datum $\mathbf{x}^{(n)}$ has its own:

- ▶ marginal likelihood $p_{\theta}(\mathbf{x}^{(n)})$
- ▶ posterior $p_{\theta}(\mathbf{z}^{(n)} | \mathbf{x}^{(n)})$
- ▶ **variational distribution** $q_{\phi^{(n)}}(\mathbf{z}^{(n)})$

Learning with IID Data

Basic approach: introduce variational parameters $\phi^{(n)}$ for each datum and construct an overall lower bound:

$$\mathcal{L}(\theta) = \frac{1}{N} \sum_{n=1}^N \log p_{\theta}(\mathbf{x}^{(n)}) \geq \frac{1}{N} \sum_{n=1}^N \text{ELBO}(\theta, \phi^{(n)}, \mathbf{x}^{(n)})$$

$$\text{ELBO}(\theta, \phi^{(n)}, \mathbf{x}^{(n)}) = \mathbb{E}_{q_{\phi^{(n)}}} \left[\log p_{\theta}(\mathbf{Z}^{(n)}, \mathbf{x}^{(n)}) - \log q_{\phi^{(n)}}(\mathbf{Z}^{(n)}) \right]$$

Then optimize the lower bound with respect to all parameters. Compute:

$$\hat{\nabla}_{\phi^{(n)}} \approx \nabla_{\phi^{(n)}} \text{ELBO}(\theta, \phi^{(n)}, \mathbf{x}^{(n)}), \quad n = 1, \dots, N,$$

$$\hat{\nabla}_{\theta} \approx \nabla_{\theta} \frac{1}{N} \sum_{n=1}^N \text{ELBO}(\theta, \phi^{(n)}, \mathbf{x}^{(n)})$$

Then update $\theta, \phi^{(1)}, \dots, \phi^{(N)}$ using stochastic gradients.

Amortized Inference

The basic approach described above introduces a very large number of variational parameters and can be very slow for large data sets.

Amortized inference proposes to use a neural net to *predict* the variational parameters $\phi^{(n)}$ for datum $\mathbf{x}^{(n)}$, e.g.

$$q_\phi(\mathbf{z}^{(n)} | \mathbf{x}^{(n)}) = \mathcal{N}(\mathbf{z}^{(n)}; g_\phi(\mathbf{x}^{(n)}), \tau^2 I)$$

- ▶ The function g_ϕ predicts the mean of the variational posterior approximation for datum $\mathbf{x}^{(n)}$. (We could also model the (co)variance as some function of $\mathbf{x}^{(n)}$.)
- ▶ This is called *amortization* because it shares information across data points for learning the variational approximations.

Amortized Inferences: VAEs

A common choice for g_ϕ is a multi-layer neural network, similar to f_θ , e.g.:

$$f_\theta(\mathbf{z}) = h_2(\mathbf{b}_2 + \mathbf{W}_2 \cdot h_1(\mathbf{b}_1 + \mathbf{W}_1 \mathbf{z}))$$

$$g_\phi(\mathbf{x}) = \mathbf{b}_4 + \mathbf{W}_4 \cdot h_3(\mathbf{b}_3 + \mathbf{W}_3 \mathbf{x})$$

(h_1, h_2, h_3 are elementwise non-linear functions)

Illustration: p and q graphical models

Illustration: "Auto-Encoder"

Example: Inference and Learning in Bernoulli VAE

Putting all the pieces together, stochastic variational inference and learning for a Bernoulli VAE would repeat the following for all n in some order:

$$\epsilon \sim \mathcal{N}(0, I) \qquad (\mathbf{z}^{(n)} = g_\phi(\mathbf{x}^{(n)}) + \tau\epsilon)$$

$$\hat{\nabla}_{\theta, \phi} = \nabla_{\theta, \phi} \left\{ \log \mathcal{N}(g_\phi(\mathbf{x}^{(n)}) + \tau\epsilon; 0, I) + \sum_{j=1}^d \log \text{Bernoulli}(x_j^{(n)}; (f_\theta(g_\phi(\mathbf{x}^{(n)}) + \tau\epsilon)_j) - \log \mathcal{N}(g_\phi(\mathbf{x}^{(n)}) + \tau\epsilon; g_\phi(\mathbf{x}^{(n)}), \tau^2 I) \right\}$$

$$(\theta, \phi) \leftarrow (\theta, \phi) + \alpha \cdot \hat{\nabla}_{\theta, \phi}$$

Bonus: Closed Form Entropy, Etc.

Bonus: Handling Some Terms in Closed Form

The ELBO can be decomposed into several terms with different computational properties:

$$\begin{aligned} \text{ELBO}(\phi) &= \mathbb{E}_{q_\phi} \left[\log \frac{p(Z, x)}{q_\phi(Z)} \right] \\ &= \underbrace{\mathbb{E}_{q_\phi} [\log p(Z)]}_{\text{"cross entropy"}} + \underbrace{\mathbb{E}_{q_\phi} [\log p(x|Z)]}_{\text{"energy"}} - \underbrace{\mathbb{E}_{q_\phi} [\log q_\phi(Z)]}_{\text{"entropy"}} \end{aligned}$$

With simple distributions (esp. Gaussians) the cross entropy and entropy terms can often be computed in closed form.

Example: Closed-Form Cross-Entropy

Example: $p(\mathbf{z})$ is a standard normal and q_ϕ is a diagonal Gaussian:

$$\begin{aligned} p(\mathbf{z}) &= \mathcal{N}(\mathbf{z}; 0, I) \\ q_\phi(\mathbf{z}) &= \mathcal{N}(\mathbf{z}; \boldsymbol{\mu}, \text{diag}(\boldsymbol{\sigma}^2)) \end{aligned} \implies \int q_\phi(\mathbf{z}) \log p(\mathbf{z}) = -\frac{d}{2} \log(2\pi) - \frac{1}{2} \sum_{j=1}^d (\mu_j^2 + \sigma_j^2)$$

When possible, it's usually (but not always) best to compute these terms and their gradients analytically, and only use Monte Carlo estimation for the energy term.

This is because lower variance gradient estimates will make the stochastic optimization converge faster.