

Gradient-based Variational Inference

STATS 305C: Applied Statistics

Scott Linderman

May 9, 2023

Where are we?

Model	Algorithm	Application
Multivariate Normal Models	Conjugate Inference	Bayesian Linear Regression
Hierarchical Models	MCMC (MH & Gibbs)	Modeling Polling Data
Probabilistic PCA & Factor Analysis	MCMC (HMC)	Images Reconstruction
Mixture Models	EM & Variational Inference	Image Segmentation
Mixed Membership Models	Coordinate Ascent VI	Topic Modeling
Variational Autoencoders	Gradient-based VI	Image Generation
State Space Models	Message Passing	Segmenting Video Data
Bayesian Nonparametrics	Fancy MCMC	Modeling Neural Spike Trains

Taking stock of posterior inference algorithms thus far

We've covered a number of posterior inference algorithms thus far:

- ▶ **Exact inference:** for simple models (e.g. conjugate exponential family models) where the posterior is available in closed form.
- ▶ **Gibbs sampling:** an MCMC algorithm that iteratively samples conditional distributions for one variable at a time. This works well for conditionally conjugate models with weak correlations.
- ▶ **Metropolis-Hastings:** a very general MCMC algorithm to sample the posterior, and the building block for many other MCMC techniques.
- ▶ **Hamiltonian Monte Carlo:** an MCMC algorithm to draw samples from the posterior by leveraging gradients of the log joint probability. This works well for more general posteriors over continuous variables.
- ▶ **Coordinate Ascent Variational Inference:** Minimizes the KL divergence between an approximate distribution (e.g., the a mean-field distribution) and the true posterior. Yields biased estimates, but potentially lower variance for fixed computational budget.

Coordinate Ascent Variational Inference (CAVI)

- ▶ *What parametric family should we use?*
 - ▶ The **mean-field family**.
- ▶ *How should we measure closeness?*
 - ▶ The **Kullback-Leibler (KL)** divergence.
- ▶ *How do we find the closest distribution in that family?*
 - ▶ **Coordinate ascent**, assuming we have a conditionally conjugate model.

Gradient-based Variational Inference

- ▶ *What parametric family should we use?*
 - ▶ Pretty much any q , as long as we can sample from it and evaluate the log density.
- ▶ *How should we measure closeness?*
 - ▶ The **Kullback-Leibler (KL)** divergence.
- ▶ *How do we find the closest distribution in that family?*
 - ▶ **Stochastic gradient ascent** using Monte Carlo estimates of the ELBO and its gradient.

Gradient-based VI methods go under a few different names: black-box VI (BBVI), automatic differentiation VI (ADVI), fixed-form VI...

Setup

Let θ denote parameters to be inferred and $p(\theta | \mathbf{X})$ the true posterior to be approximated.

Let $\mathcal{Q} = \{q(\theta; \lambda) : \lambda \in \Lambda\}$ denote the variational family – a parametric family of distributions indexed by λ that we will optimize over in variational inference.

The optimal approximation is,

$$q^*(\theta; \lambda) = \arg \min_{q \in \mathcal{Q}} D_{\text{KL}}(q(\theta; \lambda) \parallel p(\theta | \mathbf{X})) \quad (1)$$

or equivalently

$$\lambda^* = \arg \min_{\lambda \in \Lambda} D_{\text{KL}}(q(\theta; \lambda) \parallel p(\theta | \mathbf{X})) \quad (2)$$

$$= \arg \max_{\lambda \in \Lambda} \mathcal{L}(\lambda) \quad (3)$$

where $\mathcal{L}(\lambda)$ denotes the **evidence lower bound (ELBO)**,

$$\mathcal{L}(\lambda) = \mathbb{E}_{q(\theta; \lambda)} [\log p(\mathbf{X}, \theta) - \log q(\theta; \lambda)] \quad (4)$$

Stochastic gradient ascent on the local ELBO

- ▶ **Idea:** Assume the variational parameters Λ are unconstrained (i.e., $\Lambda = \mathbb{R}^D$), then perform (stochastic) gradient ascent.
- ▶ If the parameters are unconstrained and the ELBO is differentiable, we can use **gradient ascent**. Repeat:

$$\lambda \leftarrow \lambda + \alpha \nabla_{\lambda} \mathcal{L}(\lambda) \quad (5)$$

with **step size** α . Typically, you decrease the step size over iterations so that $\alpha_1 \geq \alpha_2 \geq \dots$

- ▶ More generally, we can use **stochastic gradient ascent** with an estimate of the gradient, $\widehat{\nabla}_{\lambda} \mathcal{L}(\lambda)$, as long as it is unbiased,

$$\mathbb{E}[\widehat{\nabla}_{\lambda} \mathcal{L}(\lambda)] = \nabla_{\lambda} \mathcal{L}(\lambda). \quad (6)$$

Stochastic gradient ascent on the local ELBO

- ▶ **Idea:** Assume the variational parameters Λ are unconstrained (i.e., $\Lambda = \mathbb{R}^D$), then perform (stochastic) gradient ascent.
- ▶ If the parameters are unconstrained and the ELBO is differentiable, we can use **gradient ascent**. Repeat:

$$\lambda \leftarrow \lambda + \alpha \nabla_{\lambda} \mathcal{L}(\lambda) \quad (5)$$

with **step size** α . Typically, you decrease the step size over iterations so that $\alpha_1 \geq \alpha_2 \geq \dots$

- ▶ More generally, we can use *stochastic gradient ascent* with an estimate of the gradient, $\widehat{\nabla}_{\lambda} \mathcal{L}(\lambda)$, as long as it is unbiased,

$$\mathbb{E}[\widehat{\nabla}_{\lambda} \mathcal{L}(\lambda)] = \nabla_{\lambda} \mathcal{L}(\lambda). \quad (6)$$

Stochastic gradient ascent on the local ELBO

- ▶ **Idea:** Assume the variational parameters Λ are unconstrained (i.e., $\Lambda = \mathbb{R}^D$), then perform (stochastic) gradient ascent.
- ▶ If the parameters are unconstrained and the ELBO is differentiable, we can use **gradient ascent**. Repeat:

$$\lambda \leftarrow \lambda + \alpha \nabla_{\lambda} \mathcal{L}(\lambda) \quad (5)$$

with **step size** α . Typically, you decrease the step size over iterations so that $\alpha_1 \geq \alpha_2 \geq \dots$

- ▶ More generally, we can use **stochastic gradient ascent** with an estimate of the gradient, $\widehat{\nabla}_{\lambda} \mathcal{L}(\lambda)$, as long as it is unbiased,

$$\mathbb{E}[\widehat{\nabla}_{\lambda} \mathcal{L}(\lambda)] = \nabla_{\lambda} \mathcal{L}(\lambda). \quad (6)$$

Monte Carlo gradient estimation

No problem! We'll just use ordinary Monte Carlo to estimate the gradient. But we run into a problem...

$$\nabla_{\lambda} \mathcal{L}(\lambda) = \nabla_{\lambda} \mathbb{E}_{q(\theta; \lambda)} [\log p(\mathbf{x}, \theta) - \log q(\theta; \lambda)] \quad (7)$$

$$\neq \mathbb{E}_{q(\theta; \lambda)} [\nabla_{\lambda} (\log p(\mathbf{x}, \theta) - \log q(\theta; \lambda))]. \quad (8)$$

Problem: Why can't we simply bring the gradient inside the expectation?

The “score function” gradient estimator I

The basic problem is that the variational parameters λ determine the distribution we are taking an expectation under. However, there are a few ways to obtain unbiased estimates of the gradient.

One approach is called the **score function gradient estimator** or the **REINFORCE estimator** [Williams, 1992]. It is based on the following identity,

$$\nabla_{\lambda} \log q(\theta; \lambda) = \frac{\nabla_{\lambda} q(\theta; \lambda)}{q(\theta; \lambda)} \quad (9)$$

where the l.h.s. is called the score function of distribution q .

The “score function” gradient estimator II

We can use this identity to obtain an unbiased estimate of the gradient of an expectation,

$$\nabla_{\lambda} \mathbb{E}_{q(\theta; \lambda)} [h(\theta)] = \nabla_{\lambda} \int q(\theta; \lambda) h(\theta) d\theta \quad (10)$$

$$= \int (\nabla_{\lambda} q(\theta; \lambda)) h(\theta) d\theta \quad (11)$$

$$= \int (q(\theta; \lambda) \nabla_{\lambda} \log q(\theta; \lambda)) h(\theta) d\theta \quad (12)$$

$$= \mathbb{E}_{q(\theta; \lambda)} [(\nabla_{\lambda} \log q(\theta; \lambda)) h(\theta)] \quad (13)$$

From this identity, we can obtain an unbiased Monte Carlo estimate,

$$\widehat{\nabla}_{\lambda} \mathbb{E}_{q(\theta; \lambda)} [h(\theta)] = \frac{1}{M} \sum_{m=1}^M [\nabla_{\lambda} \log q(\theta^{(m)}; \lambda) h(\theta^{(m)})]; \quad \theta^{(m)} \stackrel{\text{iid}}{\sim} q(\theta; \lambda) \quad (14)$$

The “score function” gradient estimator III

Notes:

1. The exchange of the gradient and the integral is allowed as long as the dominated convergence theorem holds, and it usually does for ML applications.
2. The score function gradient estimator is broadly applicable; e.g. it works for discrete and continuous latent variables θ . We just need the log density to be continuously differentiable wrt λ and to be able to sample from q .
3. If h is a function of both θ and λ , you need to apply the product rule. This gives another term,

$$\nabla_{\lambda} \mathbb{E}_{q(\theta; \lambda)} [h(\theta, \lambda)] = \mathbb{E}_{q(\theta; \lambda)} [(\nabla_{\lambda} \log q(\theta; \lambda)) h(\theta, \lambda)] + \mathbb{E}_{q(\theta; \lambda)} [\nabla_{\lambda} h(\theta, \lambda)] \quad (15)$$

Control variates

Though broadly applicable, the score function estimator is often too high variance to be useful. This problem can often be mitigated with **control variates**.

Recall that the expectation of the score is zero,

$$\mathbb{E}_{q(\boldsymbol{\theta}; \boldsymbol{\lambda})} [\nabla_{\boldsymbol{\lambda}} \log q(\boldsymbol{\theta}; \boldsymbol{\lambda})] = \int q(\boldsymbol{\theta}; \boldsymbol{\lambda}) \nabla_{\boldsymbol{\lambda}} \log q(\boldsymbol{\theta}; \boldsymbol{\lambda}) d\boldsymbol{\theta} \quad (16)$$

$$= \int \nabla_{\boldsymbol{\lambda}} q(\boldsymbol{\theta}; \boldsymbol{\lambda}) d\boldsymbol{\theta} \quad (17)$$

$$= \nabla_{\boldsymbol{\lambda}} \int q(\boldsymbol{\theta}; \boldsymbol{\lambda}) d\boldsymbol{\theta} \quad (18)$$

$$= \nabla_{\boldsymbol{\lambda}} 1 = 0. \quad (19)$$

Thus, we can subtract off any **baseline** from the function of interest without changing the expectation, but potentially reducing variance substantially,

$$\mathbb{E}_{q(\boldsymbol{\theta}; \boldsymbol{\lambda})} [h(\boldsymbol{\theta}) \nabla_{\boldsymbol{\lambda}} \log q(\boldsymbol{\theta}; \boldsymbol{\lambda})] = \mathbb{E}_{q(\boldsymbol{\theta}; \boldsymbol{\lambda})} [(h(\boldsymbol{\theta}) - b) \nabla_{\boldsymbol{\lambda}} \log q(\boldsymbol{\theta}; \boldsymbol{\lambda})]. \quad (20)$$

The pathwise gradient estimator

- ▶ For example, suppose $q(\theta; \lambda) = \mathcal{N}(\theta; \mu, \text{diag}(\sigma^2))$, where $\lambda = (\mu, \log \sigma^2)$ are the (unconstrained) variational parameters. Then,

$$\theta \sim q(\theta; \lambda) \iff \theta = r(\lambda, \epsilon), \quad \epsilon \sim \mathcal{N}(\mathbf{0}, I) \quad (21)$$

where $r(\lambda, \epsilon) = \mu + \sigma \epsilon$ is a **reparameterization** of θ in terms of parameters λ and “noise” ϵ .

- ▶ We can use the **law of the unconscious statistician** to rewrite the expectations as,

$$\mathbb{E}_{q(\theta; \lambda)} [h(\theta, \lambda)] = \mathbb{E}_{\epsilon \sim \mathcal{N}(\mathbf{0}, I)} [h(r(\lambda, \epsilon), \lambda)] \quad (22)$$

The distribution that the expectation is taken under no longer depends on the parameters λ , so we can simply take the gradient inside the expectation,

$$\nabla_{\lambda} \mathbb{E}_{q(\theta; \lambda)} [h(\theta, \lambda)] = \mathbb{E}_{\epsilon \sim \mathcal{N}(\mathbf{0}, I)} [\nabla_{\lambda} h(r(\lambda, \epsilon), \lambda)] \quad (23)$$

- ▶ Now we can use **Monte Carlo** to obtain an unbiased estimate of the final expectation.

$$\widehat{\nabla}_{\lambda} \mathbb{E}_{q(\theta; \lambda)} [h(\theta, \lambda)] = \frac{1}{M} \sum_{m=1}^M \nabla_{\lambda} h(r(\lambda, \epsilon_m), \lambda); \quad \epsilon_m \stackrel{\text{iid}}{\sim} \mathcal{N}(\mathbf{0}, I) \quad (24)$$

The pathwise gradient estimator

- ▶ For example, suppose $q(\theta; \lambda) = \mathcal{N}(\theta; \mu, \text{diag}(\sigma^2))$, where $\lambda = (\mu, \log \sigma^2)$ are the (unconstrained) variational parameters. Then,

$$\theta \sim q(\theta; \lambda) \iff \theta = r(\lambda, \epsilon), \quad \epsilon \sim \mathcal{N}(\mathbf{0}, I) \quad (21)$$

where $r(\lambda, \epsilon) = \mu + \sigma \epsilon$ is a **reparameterization** of θ in terms of parameters λ and “noise” ϵ .

- ▶ We can use the **law of the unconscious statistician** to rewrite the expectations as,

$$\mathbb{E}_{q(\theta; \lambda)} [h(\theta, \lambda)] = \mathbb{E}_{\epsilon \sim \mathcal{N}(\mathbf{0}, I)} [h(r(\lambda, \epsilon), \lambda)] \quad (22)$$

The distribution that the expectation is taken under no longer depends on the parameters λ , so we can simply take the gradient inside the expectation,

$$\nabla_{\lambda} \mathbb{E}_{q(\theta; \lambda)} [h(\theta, \lambda)] = \mathbb{E}_{\epsilon \sim \mathcal{N}(\mathbf{0}, I)} [\nabla_{\lambda} h(r(\lambda, \epsilon), \lambda)] \quad (23)$$

- ▶ Now we can use **Monte Carlo to obtain an unbiased estimate of the final expectation.**

$$\widehat{\nabla}_{\lambda} \mathbb{E}_{q(\theta; \lambda)} [h(\theta, \lambda)] = \frac{1}{M} \sum_{m=1}^M \nabla_{\lambda} h(r(\lambda, \epsilon_m), \lambda); \quad \epsilon_m \stackrel{\text{iid}}{\sim} \mathcal{N}(\mathbf{0}, I) \quad (24)$$

The pathwise gradient estimator

- ▶ For example, suppose $q(\boldsymbol{\theta}; \boldsymbol{\lambda}) = \mathcal{N}(\boldsymbol{\theta}; \boldsymbol{\mu}, \text{diag}(\boldsymbol{\sigma}^2))$, where $\boldsymbol{\lambda} = (\boldsymbol{\mu}, \log \boldsymbol{\sigma}^2)$ are the (unconstrained) variational parameters. Then,

$$\boldsymbol{\theta} \sim q(\boldsymbol{\theta}; \boldsymbol{\lambda}) \iff \boldsymbol{\theta} = r(\boldsymbol{\lambda}, \boldsymbol{\epsilon}), \quad \boldsymbol{\epsilon} \sim \mathcal{N}(\mathbf{0}, \mathbf{I}) \quad (21)$$

where $r(\boldsymbol{\lambda}, \boldsymbol{\epsilon}) = \boldsymbol{\mu} + \boldsymbol{\sigma} \boldsymbol{\epsilon}$ is a **reparameterization** of $\boldsymbol{\theta}$ in terms of parameters $\boldsymbol{\lambda}$ and “noise” $\boldsymbol{\epsilon}$.

- ▶ We can use the **law of the unconscious statistician** to rewrite the expectations as,

$$\mathbb{E}_{q(\boldsymbol{\theta}; \boldsymbol{\lambda})} [h(\boldsymbol{\theta}, \boldsymbol{\lambda})] = \mathbb{E}_{\boldsymbol{\epsilon} \sim \mathcal{N}(\mathbf{0}, \mathbf{I})} [h(r(\boldsymbol{\lambda}, \boldsymbol{\epsilon}), \boldsymbol{\lambda})] \quad (22)$$

The distribution that the expectation is taken under no longer depends on the parameters $\boldsymbol{\lambda}$, so we can simply take the gradient inside the expectation,

$$\nabla_{\boldsymbol{\lambda}} \mathbb{E}_{q(\boldsymbol{\theta}; \boldsymbol{\lambda})} [h(\boldsymbol{\theta}, \boldsymbol{\lambda})] = \mathbb{E}_{\boldsymbol{\epsilon} \sim \mathcal{N}(\mathbf{0}, \mathbf{I})} [\nabla_{\boldsymbol{\lambda}} h(r(\boldsymbol{\lambda}, \boldsymbol{\epsilon}), \boldsymbol{\lambda})] \quad (23)$$

- ▶ Now we can **use Monte Carlo to obtain an unbiased estimate of the final expectation.**

$$\widehat{\nabla}_{\boldsymbol{\lambda}} \mathbb{E}_{q(\boldsymbol{\theta}; \boldsymbol{\lambda})} [h(\boldsymbol{\theta}, \boldsymbol{\lambda})] = \frac{1}{M} \sum_{m=1}^M \nabla_{\boldsymbol{\lambda}} h(r(\boldsymbol{\lambda}, \boldsymbol{\epsilon}_m), \boldsymbol{\lambda}); \quad \boldsymbol{\epsilon}_m \stackrel{\text{iid}}{\sim} \mathcal{N}(\mathbf{0}, \mathbf{I}) \quad (24)$$

Exercises

Exercise: Come up with a reparameterization of an exponential distribution,

$$q(\theta; \lambda) = \text{Exp}(\theta; \lambda)$$

Question: Can you use the pathwise gradient estimator for a Bernoulli posterior,

$$q(\theta; \lambda) = \text{Bern}(\theta; \lambda)?$$

Empirically comparing estimator variances

— Score function
 — Score function + variance reduction
 — Pathwise
 — Measure-valued + variance reduction
— Value of the cost
 - - - Derivative of the cost

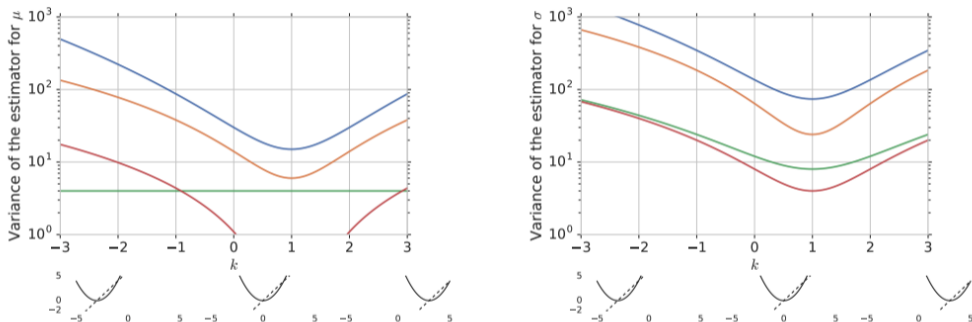


Figure 2: Variance of the stochastic estimates of $\nabla_{\theta} \mathbb{E}_{\mathcal{N}(x|\mu, \sigma^2)} [(x-k)^2]$ for $\mu = \sigma = 1$ as a function of k for three different classes of gradient estimators. Left: $\theta = \mu$; right: $\theta = \sigma$. The graphs in the bottom row show the function (solid) and its gradient (dashed) for $k \in \{-3, 0, 3\}$.

Working with mini-batches of data

Often, the ELBO involves a sum over data points,

$$\mathcal{L}(\lambda) = \mathbb{E}_q[\log p(\mathbf{X}, \boldsymbol{\theta}) - \log q(\boldsymbol{\theta}; \lambda)] \quad (25)$$

$$= \mathbb{E}_{q(\boldsymbol{\theta}; \lambda)} \left[\sum_{n=1}^N \log p(\mathbf{x}_n | \boldsymbol{\theta}) + \log p(\boldsymbol{\theta}) - \log q(\boldsymbol{\theta}; \lambda) \right] \quad (26)$$

$$= \sum_{n=1}^N \mathbb{E}_{q(\boldsymbol{\theta}; \lambda)} [\log p(\mathbf{x}_n | \boldsymbol{\theta})] - D_{\text{KL}}(q(\boldsymbol{\theta}; \lambda) \| p(\boldsymbol{\theta})) \quad (27)$$

We can view the sum as an “expectation” over data indices,

$$\sum_{n=1}^N \mathbb{E}_{q(\boldsymbol{\theta}; \lambda)} [\log p(\mathbf{x}_n | \boldsymbol{\theta})] = N \mathbb{E}_{n \sim \text{Unif}(1, N)} [\mathbb{E}_{q(\boldsymbol{\theta}; \lambda)} [\log p(\mathbf{x}_n | \boldsymbol{\theta})]], \quad (28)$$

and we can use Monte Carlo to approximate both expectations! (The same is true for Monte Carlo estimators of the gradient of the ELBO.)

SGD convergence and extensions

When does SGD work? This is a well studied problem in stochastic optimization [Bottou et al., 1998, Robbins and Siegmund, 1971].

Under relatively mild conditions, SGD converges to a **local minimum** if the step sizes obey the **Robbins-Monro conditions**,

$$\sum_{i=0}^{\infty} \alpha_i = \infty \quad \text{and} \quad \sum_{i=0}^{\infty} \alpha_i^2 < \infty \quad (29)$$

There have been dozens of extensions to basic SGD including,

- ▶ SGD with momentum
- ▶ AdaGrad [Duchi et al., 2011]
- ▶ RMSProp
- ▶ Adam [Kingma and Ba, 2014]

References I

- Ronald J Williams. Simple statistical gradient-following algorithms for connectionist reinforcement learning. *Machine learning*, 8(3):229–256, 1992.
- Shakir Mohamed, Mihaela Rosca, Michael Figurnov, and Andriy Mnih. Monte Carlo gradient estimation in machine learning. *Journal of Machine Learning Research*, 21(132):1–62, 2020.
- Léon Bottou et al. Online learning and stochastic approximations. *On-line learning in neural networks*, 17(9):142, 1998.
- Herbert Robbins and David Siegmund. A convergence theorem for non negative almost supermartingales and some applications. In *Optimizing methods in statistics*, pages 233–257. Elsevier, 1971.
- John Duchi, Elad Hazan, and Yoram Singer. Adaptive subgradient methods for online learning and stochastic optimization. *Journal of machine learning research*, 12(7), 2011.

References II

Diederik P Kingma and Jimmy Ba. Adam: A method for stochastic optimization. *arXiv preprint arXiv:1412.6980*, 2014.