

Variational Inference (Wrap-up), Inference via Sampling

Piyush Rai

Topics in Probabilistic Modeling and Inference (CS698X)

Feb 27, 2019



Recap: VI using Monte-Carlo based Gradients of ELBO

- VI = ELBO optimization. Requires ELBO gradients: $\nabla_{\phi} \mathcal{L}(\phi) = \nabla_{\phi} \mathbb{E}_q[\log p(\mathbf{X}, \mathbf{Z}) - \log q(\mathbf{Z}|\phi)]$
- Looked at two approaches that optimize ELBO using its **Monte-Carlo based gradients**
 - Black-box VI (a.k.a. score-function gradients): No model-specific gradient calculations required

$$\mathbf{Z}_s \sim q(\mathbf{Z}|\phi) \quad s = 1, \dots, S$$
$$\nabla_{\phi} \mathcal{L}(q) \approx \frac{1}{S} \sum_{s=1}^S \nabla_{\phi} \log q(\mathbf{Z}_s|\phi) [\log p(\mathbf{X}, \mathbf{Z}_s) - \log q(\mathbf{Z}_s|\phi)]$$

- Reparametrization trick (a.k.a. pathwise gradients)

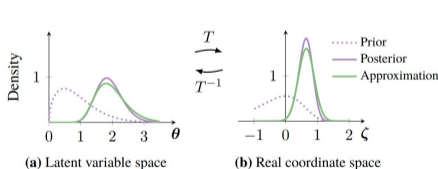
$$\mathbf{Z} = g(\epsilon, \phi)$$
$$\epsilon_s \sim p(\epsilon) \quad s = 1, \dots, S$$
$$\nabla_{\phi} \mathcal{L}(q) \approx \frac{1}{S} \sum_{s=1}^S [\nabla_{\phi} \log p(\mathbf{X}, g(\epsilon_s, \phi)) - \nabla_{\phi} \log q_{\phi}(g(\epsilon_s, \phi))]$$

- Note: We can use minibatches of data (instead of all \mathbf{X}) to compute the above gradients



Automatic Differentiation Variational Inference (ADVI)

- Auto. Diff. (AD): A way to automate differentiation of functions with **unconstrained variables**
- VI is also optimization. However, often the variables are **constrained**, e.g.,
 - Gamma's shape and scale can only be non-negative
 - Beta's parameters can only be non-negative
 - Dirichlet's probability parameter sums to one
- If we can somehow transform our distributions to unconstrained ones, we can use AD for VI



$$T : \text{supp}(p(\theta)) \rightarrow \mathbb{R}^K$$
$$\zeta = T(\theta)$$
$$p(\mathbf{x}, \zeta) = p(\mathbf{x}, T^{-1}(\zeta)) \left| \det J_{T^{-1}}(\zeta) \right|$$

Transformed density Original density Jacobian of inverse of T

- ADVI transforms the variables to real-valued and then does VI with Gaussian variational approx.

* Automatic Differentiation Variational Inference (Kucukelbir et al, 2017)



Amortized Variational Inference



Amortized Variational Inference

- Many latent variable models have one latent variable \mathbf{z}_n for each data point \mathbf{x}_n
- VI finds the optimal ϕ_n for each $q(\mathbf{z}_n|\phi_n)$
- This can be expensive for large datasets (a similar issue which motivated SVI)
- Also slow at test time: Given a new \mathbf{x}_* , finding ϕ_* requires iterative updates
 - Update local ϕ_* , update global λ , and repeat until convergence
- Amortized VI : Learn an “inference network” or “recognition model” to directly get ϕ_n , e.g.,
 - A neural network to directly map \mathbf{x}_n to ϕ_n

$$q(\mathbf{z}_n|\phi_n) \approx q(\mathbf{z}_n|\hat{\phi}_n) \quad \text{where} \quad \hat{\phi}_n = \text{NN}_\phi(\mathbf{x}_n)$$

- The inference network params ϕ can be learned along with the other global vars
- Popular in deep probabilistic models such as variational autoencoders, deep Gaussian Processes, etc



Structured Variational Inference



Structured Variational Inference

- Here “structured” may refer to anything that makes the VI approximation more expressive, e.g.,
 - Removing the independence assumption of mean-field VI
 - Learning more complex forms variational distributions
- To remove the mean-field assumption, various approaches exist
 - Structured mean-field (Saul et al, 1996)
 - Hierarchical VI (Ranganath et al, 2016): Variational params ϕ_1, \dots, ϕ_M “tied” via a [shared prior](#)

$$q(\mathbf{z}_1, \dots, \mathbf{z}_M | \theta) = \int \left[\prod_{m=1}^M q(\mathbf{z}_m | \phi_m) \right] p(\phi | \theta) d\phi$$

- To learn more expressive variational approximations, various approaches exist, e.g.,
 - Boosting or mixture of simpler distributions, e.g., $q(\mathbf{z}) = \sum_{c=1}^C \rho_c q_c(\mathbf{z} | \phi_c)$
 - [Normalizing flows](#): Turn a simple $q(\mathbf{z})$ into a complex one via series of invertible transformations



Other Divergence Measures



Other Divergence Measures

- VI minimizes $KL(q||p)$ but other divergences can be minimized as well
- A general form of divergence is Renyi's α -divergence defined as

$$D_{\alpha}^R(p(\mathbf{x})||q(\mathbf{x})) = \frac{1}{\alpha - 1} \log \int p(\mathbf{x})^{\alpha} q(\mathbf{x})^{1-\alpha} d\mathbf{x}$$

- $KL(p||q)$ is a special case with $\alpha \rightarrow 1$ (can verify using L'Hopital rule of taking limits)
- An even more general form of divergence is f -Divergence

$$D_f(p(\mathbf{x})||q(\mathbf{x})) = \int q(\mathbf{x}) f\left(\frac{p(\mathbf{x})}{q(\mathbf{x})}\right) d\mathbf{x}$$

- Many recent inference algorithms are based on minimizing such divergences



Variational Inference: Some Comments

- Many probabilistic models (deep/non-deep) nowadays rely on VI to do tractable inference
- Even mean-field for locally-conjugate models has many applications in lots of probabilistic models
 - This + SVI gives excellent scalability
- Stoch. opt., auto. diff., Monte-Carlo gradient of ELBO, contributed immensely to the success
- Note: Most of these ideas apply also to Variational EM
- Many VI and advanced VI algorithms are implemented in probabilistic programming packages (e.g., Stan, Tensorflow Probability, etc), making VI a painless exercise even for complex models
- Still a very active area of research, especially for doing VI in complex models
 - Models with discrete latent variables
 - Reducing the variance in Monte-Carlo estimate of ELBO gradients



Inference via Sampling

(Note that we have already seen Gibbs sampling)



Sampling for Approximate Inference

- Some typical inference tasks

- Compute a (possibly intractable) **posterior distribution**: $p(\theta|\mathcal{D}) = \frac{p(\mathcal{D}|\theta)p(\theta)}{p(\mathcal{D})} = \frac{p(\mathcal{D}|\theta)p(\theta)}{\int p(\mathcal{D}|\theta)p(\theta)d\theta}$
- Compute a difficult **expectation** of a random quantity w.r.t. a distribution (an integral), e.g.,
 - The **posterior predictive** (an expectation w.r.t the posterior over θ)

$$p(\mathcal{D}^{new}|\mathcal{D}) = \int p(\mathcal{D}^{new}|\theta)p(\theta|\mathcal{D})d\theta = \mathbb{E}_{p(\theta|\mathcal{D})}[p(\mathcal{D}^{new}|\theta)]$$

- The **marginal likelihood** or “evidence” (an expectation over the prior)

$$p(\mathcal{D}|m) = \int p(\mathcal{D}|\theta)p(\theta|m)d\theta = \mathbb{E}_{p(\theta|m)}[p(\mathcal{D}|\theta)]$$

- The **expected complete data log-likelihood** needed for doing MLE/MAP in LVMs (recall EM)

$$\text{Exp-CLL} = \int p(\mathbf{z}|\theta, \mathbf{x})p(\mathbf{x}, \mathbf{z}|\theta)d\mathbf{z} = \mathbb{E}_{p(\mathbf{z}|\theta, \mathbf{x})}[p(\mathbf{x}, \mathbf{z}|\theta)]$$

- The **ELBO** in variational inference

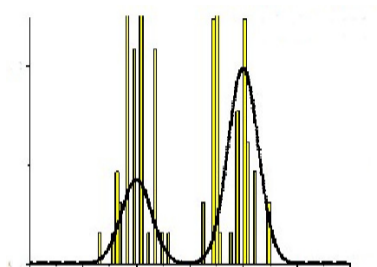
$$\mathcal{L}(q) = \mathbb{E}_q[\log p(\mathbf{x}, \mathbf{z})] - \mathbb{E}_q[\log p(\mathbf{z})]$$

- **Sampling methods** provide a general way to (approximately) solve these problems



The Basic Idea

- Can approximate any distribution using a set of **randomly drawn samples** from it



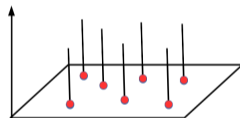
- The samples can also be used for computing expectations (Monte-Carlo averaging)
- Usually straightforward to generate samples if it is a simple/standard distribution
- **The interesting bit:** Even if the distribution is “difficult” (e.g., an intractable posterior), it is often possible to generate random samples from such a distribution, as we will see..



Empirical Distribution

- Sampling based approximation of a distribution can be represented using an **empirical distribution**
- Given L “points” $\mathbf{z}^{(1)}, \dots, \mathbf{z}^{(L)}$, the **empirical distribution** of these points is defined as

$$p_L(A) = \sum_{\ell=1}^L w_{\ell} \delta_{\mathbf{z}^{(\ell)}}(A)$$



- Here w_1, \dots, w_L are weights that sum to 1, i.e., $\sum_{\ell=1}^L w_{\ell} = 1$ (for uniform weights, $w_{\ell} = 1/L$)
- Here $\delta_{\mathbf{z}}(A)$ denotes the **Dirac distribution** defined as

$$\delta_{\mathbf{z}}(A) = \begin{cases} 0 & \text{if } \mathbf{z} \notin A \\ 1 & \text{if } \mathbf{z} \in A \end{cases}$$

- $p_L(A)$ is a discrete distribution with finite support $\mathbf{z}^{(1)}, \dots, \mathbf{z}^{(L)}$ (can think of it as a histogram)



Approximate Inference: VI vs Sampling-based

- VI approximates a posterior distribution $p(\mathbf{Z}|\mathbf{X})$ by another distribution $q(\mathbf{Z}|\phi)$
- Sampling uses S (typically large number) samples $\{\mathbf{Z}_s\}_{s=1}^S$ to approximate $p(\mathbf{Z}|\mathbf{X})$
- Sampling can be used within VI (already saw ELBO approximations using Monte-Carlo)
- Also possible (though less common) to use VI in sampling algorithms (will talk about it later)
- In terms of “comparison” between VI and sampling, a few things to be noted
 - **Convergence:** VI only has local convergence, sampling (in theory) can give posterior (more on it later)
 - **Storage requirements:** Sampling-based approximation requires more storage (why?)
 - **Prediction time cost (also related to storage requirement):** Sampling always requires Monte-Carlo averaging for posterior predictive; with VI, sometimes we can get closed form posterior predictive
 - Sampling based posterior predictive: $p(x_*|\mathbf{X}) \approx \frac{1}{S} \sum_{s=1}^S p(x_*|\theta_s)p(\theta_s|\mathbf{X})$
 - VI based posterior predictive: $p(x_*|\mathbf{X}) \approx \int p(x_*|\theta)q(\theta|\phi)d\theta$
 - There is some work on “compressing” sampling-based approximations (e.g., see “Compact approximations to Bayesian predictive distributions” by Snelson and Ghahramani, 2005; and “Bayesian Dark Knowledge” by Korattikara et al, 2015)



Sampling: Some Basic Methods

- Most of these basic methods are based on the idea of transformation
- Given a sample x from an “easy” distribution $p(x)$, transform it into a random sample z from a “less easy” distribution $p(z)$
- Some popular examples of transformation methods

- Inverse CDF method

$$x \sim \text{Unif}(0, 1) \Rightarrow z = \text{Inv-CDF}_{p(z)}(x) \sim p(z)$$

- Reparametrization method

$$x \sim \mathcal{N}(0, 1) \Rightarrow z = \mu + \sigma x \sim \mathcal{N}(\mu, \sigma^2)$$

- Box-Muller method: Given (x_1, x_2) from $\text{Unif}(-1, +1)$, generate (z_1, z_2) from 2D Gaussian $\mathcal{N}(0, \mathbf{I})$
- Transformation Methods are simple but have limitations
 - Mostly limited to standard distributions and/or distributions with very few variables

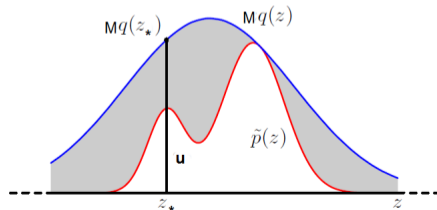


Rejection Sampling

- Want to **sample** from $p(z) = \frac{\tilde{p}(z)}{Z_p}$. Suppose we can only evaluate the numerator $\tilde{p}(z)$ at any z
- Suppose we have a **proposal distribution** $q(z)$ that we can **generate samples from**, and

$$Mq(z) \geq \tilde{p}(z) \quad \forall z \quad (\text{where } M > 0 \text{ is some const.})$$

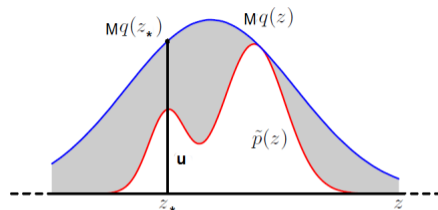
- Basic idea: Generate samples from the proposal $q(z)$ and **accept/reject based on some condition**
 - Sample an r.v. z_* from $q(z)$
 - Sampling a uniform r.v. $u \sim \text{Unif}[0, Mq(z_*)]$



- If $u \leq \tilde{p}(z_*)$ then accept z_* else reject



Rejection Sampling



- Why $\mathbf{z} \sim q(\mathbf{z}) + \text{accept/reject rule}$ is equivalent to $\mathbf{z} \sim p(\mathbf{z})$?
- Let's look at the pdf of \mathbf{z} 's that were accepted, i.e., $p(\mathbf{z}|\text{accept})$

$$p(\text{accept}|\mathbf{z}) = \int_0^{\tilde{p}(\mathbf{z})} \frac{1}{Mq(\mathbf{z})} du = \frac{\tilde{p}(\mathbf{z})}{Mq(\mathbf{z})}$$

$$p(\mathbf{z}, \text{accept}) = q(\mathbf{z})p(\text{accept}|\mathbf{z}) = \frac{\tilde{p}(\mathbf{z})}{M}$$

$$p(\text{accept}) = \int \frac{\tilde{p}(\mathbf{z})}{M} d\mathbf{z} = \frac{Z_p}{M}$$

$$p(\mathbf{z}|\text{accept}) = \frac{p(\mathbf{z}, \text{accept})}{p(\text{accept})} = \frac{\tilde{p}(\mathbf{z})}{Z_p} = p(\mathbf{z})$$



Sampling for Approximating Expectations

- Suppose $f(\mathbf{z})$ is function of a random variable $\mathbf{z} \sim p(\mathbf{z})$
- Wish to compute $\mathbb{E}[f] = \mathbb{E}_{p(\mathbf{z})}[f(\mathbf{z})] = \int f(\mathbf{z})p(\mathbf{z})d\mathbf{z}$
- Given L independent samples $\{\mathbf{z}^{(\ell)}\}_{\ell=1}^L$ from $p(\mathbf{z})$, we can approximate the above as

$$\mathbb{E}[f] \approx \frac{1}{L} \sum_{\ell=1}^L f(\mathbf{z}^{(\ell)}) \quad (\text{Monte Carlo sampling})$$

- What if we can't generate samples from $p(\mathbf{z})$? Answer: Use [Importance Sampling](#)
 - If we can generate L indep. samples $\{\mathbf{z}^{(\ell)}\}_{\ell=1}^L$ from a different "proposal" distribution $q(\mathbf{z})$ then

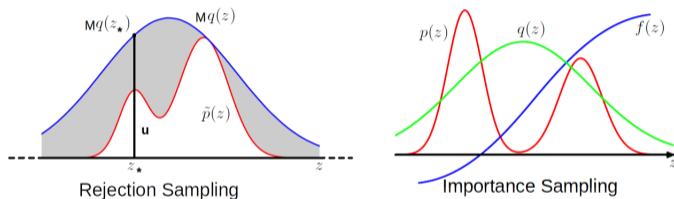
$$\mathbb{E}[f] = \int f(\mathbf{z})p(\mathbf{z})d\mathbf{z} = \int f(\mathbf{z})\frac{p(\mathbf{z})}{q(\mathbf{z})}q(\mathbf{z})d\mathbf{z} \approx \frac{1}{L} \sum_{\ell=1}^L f(\mathbf{z}^{(\ell)})\frac{p(\mathbf{z}^{(\ell)})}{q(\mathbf{z}^{(\ell)})}$$

- IS only requires that we can evaluate $p(\mathbf{z})$ at any \mathbf{z} (in fact, with a small modification to the above, IS works even when we can evaluate $p(\mathbf{z})$ only up to a proportionality constant)
- Note: IS is NOT a sampling method (doesn't generate samples from a desired distribution; just a way to approximate expectations)



Limitations of Basic Sampling Methods

- Transformation based methods: Usually limited to drawing from standard distributions
- Rejection Sampling and Importance Sampling: Require good proposal distributions



- Difficult to find good prop. distr. especially when \mathbf{z} is high-dim. (e.g., models with many params)
 - In high dimensions, most of the mass of $p(\mathbf{z})$ is concentrated in a tiny region of the \mathbf{z} space
 - Difficult to *a priori* know what those regions are, thus difficult to come up with good proposal dist.
- A solution to these: MCMC methods

