

Variational inference: Re-parametrization trick

Setup: Bayesian inference on a ML model with parameters θ given some data \mathcal{D} : $\theta \in \mathbb{R}^d$

$$p(\theta|\mathcal{D}) = \frac{p(\mathcal{D}|\theta)p(\theta)}{p(\mathcal{D})}$$

→ variational family parametrized by ϕ

Idea: Approximate $p(\theta|\mathcal{D}) \approx q_\phi(\theta|\mathcal{D})$

Mean-field family: $q_\phi(\theta|\mathcal{D}) = \prod_{i=1}^d \mathcal{N}(\theta_i | \mu_i, \sigma_i^2)$

$$\phi := \{\mu_1, \sigma_1^2, \mu_2, \sigma_2^2, \dots, \mu_d, \sigma_d^2\}$$

Train ϕ : $\phi^* = \arg\min_{\phi} \text{KL}[q_\phi(\theta|\mathcal{D}) || p(\theta|\mathcal{D})]$

see bot

⇒

lecture

$$\phi^* = \arg\min_{\phi} \mathcal{L}(\phi) \xrightarrow{\text{G.D.}} \phi^{n+1} = \phi^n - \eta \nabla_{\phi} \mathcal{L}(\phi)$$

$$\mathcal{L}(\phi) := \underbrace{-H[q_\phi(\theta|\mathcal{D})]}_{\text{computed analytically}} - \underbrace{\mathbb{E}_{\theta \sim q_\phi(\theta|\mathcal{D})} [\log p(\mathcal{D}|\theta) + \log p(\theta)]}_{\text{computed via sampling}}$$

Remarks on computing $\nabla_{\phi} \mathcal{L}(\phi)$.

Gradient of the second term wrt ϕ :

$$\begin{aligned} \nabla_{\phi} \mathbb{E}_{\theta \sim q_\phi(\theta|\mathcal{D})} [\log p(\mathcal{D}|\theta) + \log p(\theta)] &= \\ &= \nabla_{\phi} \int \log p(\mathcal{D}|\theta) q_\phi(\theta|\mathcal{D}) d\theta + \nabla_{\phi} \int \log p(\theta) q_\phi(\theta|\mathcal{D}) d\theta \end{aligned}$$

$$= \int \log p(\mathcal{D}|\theta) \nabla_{\theta} q_{\theta}(\theta|\mathcal{D}) d\theta + \int \log p(\theta) \nabla_{\theta} q_{\theta}(\theta|\mathcal{D}) d\theta \quad (1)$$

Recall : $(\log f(x))' = \frac{f'(x)}{f(x)}$

$$\left\{ \nabla_{\theta} \log q_{\theta}(\theta|\mathcal{D}) = \frac{\nabla_{\theta} q_{\theta}(\theta|\mathcal{D})}{q_{\theta}(\theta|\mathcal{D})} \Rightarrow \right.$$

$$\Rightarrow \nabla_{\theta} q_{\theta}(\theta|\mathcal{D}) = \nabla_{\theta} \log q_{\theta}(\theta|\mathcal{D}) \cdot q_{\theta}(\theta|\mathcal{D}) \quad (2)$$

$$\stackrel{(1)}{=} \stackrel{(2)}{=} \int \log p(\mathcal{D}|\theta) \nabla_{\theta} \log q_{\theta}(\theta|\mathcal{D}) \cdot q_{\theta}(\theta|\mathcal{D}) d\theta + \int \log p(\theta) \nabla_{\theta} \log q_{\theta}(\theta|\mathcal{D}) q_{\theta}(\theta|\mathcal{D}) d\theta$$

$$\Rightarrow \nabla_{\theta} \mathbb{E}_{\theta \sim q_{\theta}(\theta|\mathcal{D})} [\log p(\mathcal{D}|\theta) + \log p(\theta)] =$$

$$= \mathbb{E}_{\theta \sim q_{\theta}(\theta|\mathcal{D})} \left[\nabla_{\theta} \log q_{\theta}(\theta|\mathcal{D}) \cdot (\log p(\mathcal{D}|\theta) + \log p(\theta)) \right]$$

$$\approx \frac{1}{n} \sum_{i=1}^n \nabla_{\theta} \log q_{\theta}(\theta_i|\mathcal{D}) [\log p(\mathcal{D}|\theta_i) + \log p(\theta_i)],$$

where $\theta_i \stackrel{i.i.d.}{\sim}_{\theta \in \mathbb{R}^d} q_{\theta}(\theta|\mathcal{D}) \stackrel{m.p.}{=} \mathcal{N}(\theta | \mu_{\theta}, \Sigma_{\theta})$

where $\mu_{\theta} = \begin{bmatrix} \mu_1 \\ \mu_2 \\ \vdots \\ \mu_d \end{bmatrix}$, $\Sigma_{\theta} = \begin{bmatrix} \sigma_1^2 & & & \\ & \sigma_2^2 & & \\ & & \ddots & \\ 0 & & & \sigma_d^2 \end{bmatrix}$

This Monte Carlo estimator depends on θ and in practice exhibits very high variance (i.e. it is very inaccurate unless a very large number

of samples n is considered).

Re-parametrization trick:

If we can find a function $h: (\varepsilon, \phi) \rightarrow \theta$, where ε is a random variable $\varepsilon \sim p(\varepsilon)$, then we can write:

$$\theta_i = h_\phi(\varepsilon), \quad \varepsilon \sim p(\varepsilon) \quad \text{such that} \quad \theta_i \sim q_\phi(\theta|\mathcal{D})$$

i.e. we will try to find a function such that samples from $q_\phi(\theta|\mathcal{D})$

can be written as $\theta = h_\phi(\varepsilon)$, $\varepsilon \sim p(\varepsilon)$

e.g. Re-parametrize a Gaussian:

$$q_\phi(\theta|\mathcal{D}) = \mathcal{N}(\theta | \mu_\phi, \Sigma_\phi)$$

In fact, we can generate samples θ using the following re-parametrization

$$\theta \in \mathbb{R}^d = \underbrace{\mu_\phi + \varepsilon \Sigma_\phi^{\frac{1}{2}}}_{h_\phi(\varepsilon)}, \quad \text{where} \quad \underline{\varepsilon \sim p(\varepsilon) = \mathcal{N}(0, \mathbf{I})}$$

\hookrightarrow this does not depend on ϕ !

Now recall the trouble-some gradient term:

$$\nabla_\phi \mathbb{E}_{\theta \sim q_\phi(\theta|\mathcal{D})} [\log p(\mathcal{D}|\theta) + \log p(\theta)] \stackrel{\text{proof?}}{=}$$

$$\underbrace{\mathbb{E}_{\varepsilon \sim p(\varepsilon)} \left[\nabla_\phi \log p(\mathcal{D} | h_\phi(\varepsilon)) + \nabla_\phi \log p(h_\phi(\varepsilon)) \right]}$$

Now the gradient is not related to the variational parameter.

hence it is not related to the distribution with respect to which the

expectation is taken.

~

Summary:

$$\phi^* = \underset{\phi}{\operatorname{argmin}} \left(\operatorname{KL}[q_{\phi}(\theta|\mathcal{D}) \| p(\theta|\mathcal{D})] \right) := \mathcal{L}(\phi)$$

$$\mathcal{L}(\phi) := -H[q_{\phi}(\theta|\mathcal{D})] - \mathbb{E}_{\theta \sim q_{\phi}(\theta|\mathcal{D})} [\log p(\mathcal{D}|\theta) + \log p(\theta)]$$

To solve this problem we typically employ SGD:

$$\phi^{n+1} = \phi^n - \eta \nabla_{\phi} \mathcal{L}(\phi)$$

For a mean-field approximation, i.e. $q_{\phi}(\theta|\mathcal{D}) = \mathcal{N}(\theta | \mu_{\phi}, \Sigma_{\phi})$

$$\text{where } \mu_{\phi} = \begin{bmatrix} \mu_1 \\ \vdots \\ \mu_d \end{bmatrix}, \Sigma_{\phi} = \begin{bmatrix} \sigma_1^2 & & 0 \\ & \ddots & \\ 0 & & \sigma_d^2 \end{bmatrix}, \phi := \{\mu_1, \sigma_1^2, \dots, \mu_d, \sigma_d^2\}$$

1st term:

$$-H[q_{\phi}(\theta|\mathcal{D})] := \mathbb{E}_{\theta \sim q_{\phi}(\theta|\mathcal{D})} [\log q_{\phi}(\theta|\mathcal{D})] = -\sum_{i=1}^d \log \sigma_i + \text{constant}$$

$$-\nabla_{\phi} H[q_{\phi}(\theta|\mathcal{D})] = -\sum_{i=1}^d \frac{\partial}{\partial \sigma_i} \log \sigma_i = -\sum_{i=1}^d \frac{1}{\sigma_i}$$

2nd term:

$$\nabla_{\phi} \mathbb{E}_{\theta \sim q_{\phi}(\theta|\mathcal{D})} [\log p(\mathcal{D}|\theta) + \log p(\theta)]$$

$$= \mathbb{E}_{\varepsilon \sim \mathcal{N}(0, \mathbf{I})} \left[\nabla_{\phi} \log p(\mathcal{D} | \mu_{\phi} + \varepsilon \Sigma_{\phi}^{\frac{1}{2}}) + \nabla_{\phi} \log p(\mu_{\phi} + \varepsilon \Sigma_{\phi}^{\frac{1}{2}}) \right]$$

$$\approx \frac{1}{n} \sum_{i=1}^n \left[\nabla_{\phi} \log p(\mathcal{D} | \mu_{\phi} + \varepsilon_i \Sigma_{\phi}^{\frac{1}{2}}) + \nabla_{\phi} \log p(\mu_{\phi} + \varepsilon_i \Sigma_{\phi}^{\frac{1}{2}}) \right],$$

$$\text{where } \varepsilon_i \sim \mathcal{N}(0, \mathbf{I})$$

⊛ Typically in practice $n = 1, 10, 20, 50$ samples.

