

Bayesian inference

Setting: We are given a model with unknown parameters $\theta \in \mathbb{R}^d$.
 Some data \mathcal{D} distributed according to the likelihood of the model $p(\mathcal{D}|\theta)$, and a prior $p(\theta)$.

Goal: Infer the posterior:

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

} i.) Variational inference
 ii.) MCMC

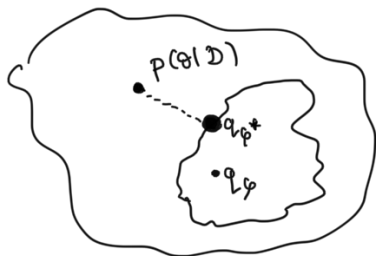
VI: Idea: Approximate the intractable using a family of distributions that is easy to work with.

Most common choice: $p(\theta|\mathcal{D}) \approx q_{\varphi}(\theta|\mathcal{D}) = \prod_{i=1}^d \mathcal{N}(\theta_i | \mu_i, \sigma_i^2)$

(Mean-field approximation)

$$\varphi := \{\mu_1, \sigma_1^2, \dots, \mu_d, \sigma_d^2\}$$

variational parameters



Goal: Find/estimate $\varphi^* := \{\mu_1, \sigma_1^2, \dots, \mu_d, \sigma_d^2\}$ such that

$q_{\varphi^*}(\theta|\mathcal{D})$ is "as close as possible" to $p(\theta|\mathcal{D})$.

$$\phi^* = \underset{\phi}{\operatorname{argmin}} \operatorname{KL}[q_{\phi}(\theta|D) \parallel p(\theta|D)]$$

Recall: $\operatorname{KL}[q_{\phi}(\theta|D) \parallel p(\theta|D)] = \int \log \frac{q_{\phi}(\theta|D)}{p(\theta|D)} q_{\phi}(\theta|D) d$

$$= \mathbb{E}_{\theta \sim q_{\phi}(\theta|D)} \left[\log \frac{q_{\phi}(\theta|D)}{p(\theta|D)} \right]$$

Why KL? • It's easy to work with

$$\bullet \operatorname{KL}[q_{\phi} \parallel p] = 0 \iff q_{\phi} = p$$

... but it's not an actual distance:

$$\operatorname{KL}[q_{\phi} \parallel p] \neq \operatorname{KL}[p \parallel q_{\phi}]$$

Remarks:

i.) Typically mean-field VI tends to favor approximation that capture well the mean of $p(\theta|D)$, but underestimate the variance!

ii.) It's hardly ideal, but in cases where other methods don't scale, it still provide useful inference.

How to estimate $\phi^* = \underset{\phi}{\operatorname{argmin}} \operatorname{KL}[q_{\phi}(\theta|D) \parallel p(\theta|D)]$

→ Old-school approaches had to derive coordinate ascent rules }
for minimizing the KL (see ch. 10 Bishop).

→ New-schoolers use Automatic Differentiation Variational Inference

ADVI: It is a "black-box" approach that is agnostic to any details about $p(\theta|D)$: any model for which we can evaluate (and differentiate) its log-likelihood and log-prior distribution works!

Let's see how it works!

$$KL[q_p(\theta|D) || p(\theta|D)] = \mathbb{E}_{\theta \sim q_p(\theta|D)} \left[\underbrace{\log q_p(\theta|D)}_{\#1} - \underbrace{\log p(\theta)}_{\#2} \right]$$

1st term: $\mathbb{E}_{\theta \sim q_p(\theta|D)} [\log q_p(\theta|D)] = \int \log q_p(\theta|D) q_p(\theta|D) d\theta$
 $= -H[q_p(\theta|D)]$

Notice that we are free to choose any $q_p(\theta|D)$ that suits us, hence we can choose one for which $H[q_p]$ is computable:

e.g. MF-VI: $\mathbb{E}_{\theta \sim q_p(\theta|D)} \left[\underbrace{q_p(\theta|D)}_{\prod_{i=1}^d \mathcal{N}(\theta_i | \mu_i, \sigma_i^2)} \right] = - \sum_{i=1}^d \log \sigma_i + \text{constant}$

Term #2:

$$\mathbb{E}_{\theta \sim q_p(\theta|D)} [\log p(\theta|D)] = \mathbb{E}_{\theta \sim q_p(\theta|D)} [\log p(D|\theta) + \log p(\theta) - \cancel{\log p(\theta)}]$$

$$\mathcal{L}_b(\theta) := -H[q_p(\theta|D)] - \mathbb{E} [\log p(D|\theta) + \log p(\theta)]$$

Now we can use gradient descent to estimate ϕ^* :
 $\phi^* = \underset{\phi}{\operatorname{argmin}} \mathcal{L}(\phi)$
 $\phi^{n+1} = \phi^n - \eta \nabla_{\phi} \mathcal{L}(\phi)$

All terms in $\mathcal{L}(\phi)$ can now be evaluated, however, we still need to compute :

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

exhibits high-variance when approximated via Monte-Carlo sampling

Example tutorial :

Given $p(x)$, then try to fit $q_{\phi}(x)$ such that

$\text{KL}[q_{\phi}(x) || p(x)]$ is minimized.

$$\phi^* = \underset{\phi}{\operatorname{argmin}} \mathbb{E}_{x \sim q_{\phi}(x)} [\log q_{\phi}(x) - \overbrace{\log p(x)}^{\text{known}}]$$

$$\approx \frac{1}{n} \sum_{i=1}^n [\log q_{\phi}(x_i) - \log p(x_i)],$$

$x_i \sim q_{\phi}(x)$