# Variational inference: Reparametrization tricks

**Setup**: Bayesian inference for some model with parameters $\theta \in \mathbb{R}^d$ given some data $D$:

$$\underbrace{p(\theta \mid D)}_{\text{intractable}} = \frac{p(D \mid \theta)\, p(\theta)}{p(D)}$$

**Idea**:

Approximate $p(\theta \mid D) \approx q_\varphi(\theta \mid D)$

e.g. mean-field: $q_\theta(\theta \mid D) = \prod_{i=1}^{d} \mathcal{N}(\theta_i \mid \mu_i, \sigma_i^2)$

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

**Training**: $\varphi^* = \underset{\varphi}{\arg\min}\ \mathbb{KL}\left[q_\varphi(\theta \mid D) \,\|\, p(\theta \mid D)\right] := \mathcal{L}(\varphi)$

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

Optimize via SGD: $\varphi^{n+1} = \varphi^n - \eta \underline{\nabla_\varphi \mathcal{L}(\varphi)}$

Compute $\nabla_\varphi \mathcal{L}(\varphi)$:

$$\nabla_\varphi \mathbb{E}_{\theta \sim q_\varphi(\theta \mid D)}\left[\log p(D \mid \theta) + \log p(\theta)\right]$$

$$= \nabla_\varphi \int \log p(D \mid \theta)\, q_\varphi(\theta \mid D)\, d\theta + \nabla_\varphi \int \log p(\theta)\, q_\varphi(\theta \mid D)\, d\theta$$

$$\qquad \qquad \nabla_\varphi \theta \quad (\theta \mid D)\, d\varphi$$

$$= \int \log p(D|\theta) \boxed{\nabla_\varphi q_\varphi(\theta|D)} d\theta + \int \log p(\theta) \nabla_\varphi q_\varphi(\theta|D) \,d\theta \quad \text{①}$$

$$\uparrow$$

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

$$\bullet \quad \nabla_\varphi \log q_\varphi(\theta|D) = \frac{\nabla_\varphi q_\varphi(\theta|D)}{q_\varphi(\theta|D)}$$

$$\implies \nabla_\varphi q_\varphi(\theta|D) = \nabla_\varphi \log q_\varphi(\theta|D) \cdot \underline{q_\varphi(\theta|D)} \quad \text{②}$$

$$\text{①} \xrightarrow{\text{②}} \nabla_\varphi \mathcal{L}(\varphi) = \int \log p(D|\theta) \, \nabla_\varphi \log q_\varphi(\theta|D) \, q_\varphi(\theta|D) \, d\theta$$

$$+ \int \log p(\theta) \, \nabla_\varphi \log q_\varphi(\theta|D) \, q_\varphi(\theta|D) \, d\theta$$

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

$$\approx \frac{1}{S} \sum_{i=1}^{S} \nabla_\varphi \log q_\varphi(\theta_i|D) \left( \log p(D|\theta_i) + \log p(\theta_i) \right),$$

$$\theta_i \sim q_\varphi(\theta|D) \overset{\text{m.F.}}{=} \mathcal{N}(\theta | \mu_\varphi, \Sigma_\varphi)$$

$$\mu_\varphi := \begin{bmatrix} \mu_1 \\ \vdots \\ \mu_d \end{bmatrix}, \quad \Sigma_\varphi = \begin{bmatrix} \sigma_1^2 & & 0 \\ & \ddots & \\ 0 & & \sigma_d^2 \end{bmatrix}$$

The above Monte Carlo estimator of $\nabla_\varphi \mathcal{L}(\varphi)$ depend on $\varphi$ and in practice tends to exhibit very high variance (i.e. it is very innaccrurate unless a very large number of

samples is cons~~idered~~ ...

# Reparametrization trick :

<u>Idea</u> : Introduce a simple "change of variables" such that we compute expectations with respect to distributions that do n depend on $\varphi$.

If we can find a function $h : (\varepsilon, \varphi) \longrightarrow \vartheta$, where $\varepsilon \sim p(\varepsilon)$, then we can write :

$$\vartheta_i = h_\varphi(\varepsilon), \quad \varepsilon \sim p(\varepsilon), \quad \text{such that} \quad \vartheta_i \sim q_\varphi(\vartheta | D)$$

e.g. re-parametrize a Gaussian :

$$\vartheta_i \sim q_\varphi(\vartheta | D) = \mathcal{N}(\vartheta_i | \mu_\varphi, \Sigma_\varphi)$$

In fact, we can generate samples $\vartheta$ by sampling $\varepsilon \sim p($

$$\vartheta = \mu_\varphi + \varepsilon \Sigma_\varphi^{\frac{1}{2}}, \quad \text{where} \quad \varepsilon \sim p(\varepsilon) = \mathcal{N}(0, I)$$

$$\vartheta = h_\varphi(\varepsilon)$$

Now recall the troublesome gradient term :

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

$$= \mathbb{E}_{\vartheta \sim q_\varphi(\vartheta | D)} \left[ \cancel{\nabla_\varphi \log q_\varphi(\vartheta | D)} \cancel{(\log p(D | \vartheta) + \log p(\vartheta))} \right]$$

$$\underset{\hookrightarrow}{} = \nabla_\varphi \ \mathbb{E}_{\varepsilon \sim p(\varepsilon)} \left[ \log p(D \mid h_\varphi(\varepsilon)) + \log p(h_\varphi(\varepsilon)) \right]$$

M.F.
$$= \nabla_\varphi \ \mathbb{E}_{\varepsilon \sim \mathcal{N}(0, I)} \left[ \log p(D \mid \underbrace{\mu_\varphi + \varepsilon \Sigma_\varphi^{\frac{1}{2}}}) + \log p(\mu_\varphi + \varepsilon \Sigma_\varphi^{\frac{1}{2}} \right.$$

Now the gradient **is not** related to the variational parau
(i.e. the distribution with respect to which the expectation is
taken **does not** depend on $\varphi$.)

Summary: $\quad \varphi^* = \underset{\varphi}{\arg\min} \ \mathbb{KL}\left[ q_\varphi(\theta \mid D) \ \| \ p(\theta \mid D) \right].$

$$\downarrow \varphi$$

$$\mathcal{L}(\varphi) := -\underline{H\left[ q_\varphi(\theta \mid D) \right]} - \underline{\mathbb{E}_{\theta \sim q_\varphi(\theta \mid D)}\left[ \log p(D \mid \theta) + \log p(\theta) \right]}$$

Evidence
lower
bound (ELBO)

$$\text{If} \ \boxed{q_\varphi(\theta \mid D) = \mathcal{N}(\theta_i \mid \mu_\varphi, \Sigma_\varphi)}$$

$1^{st}$ term: $\quad -H\left[ q_\varphi(\theta \mid D) \right] = -\sum_{i=1}^{d} \log \sigma_i + \text{constant}$

$$-\nabla_\varphi H\left[ q_\varphi(\theta \mid D) \right] = -\sum_{i=1}^{d} \frac{\partial}{\partial \sigma_i} \log \sigma_i = -\sum_{i=1}^{d} \frac{1}{\sigma_i}$$

$2^{nd}$ term:

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

$$\nabla_\varphi \ \mathbb{E} \quad \left[ \log p(D \mid \mu_\varphi + \varepsilon \Sigma_\varphi^{\frac{1}{2}}) + \log p(\mu_\varphi + \varepsilon \Sigma_\varphi \right.$$

$$= \quad \nabla_\varphi \; \underset{\varepsilon \sim \mathcal{N}(0,I)}{\mathbb{E}} \; L \quad \sigma \; 1$$

$$\approx \frac{1}{S} \sum_{i=1}^{S} \left[ \nabla_\varphi \log p \left( D \mid \underbrace{\mu_\varphi + \varepsilon_i \Sigma_\varphi^{\frac{1}{2}}}_{\theta} \right) + \nabla_\varphi \log p(\mu_\varphi + \varepsilon_i \Sigma \right.$$

where $\quad \varepsilon_i \sim \mathcal{N}(0,I).$