# Particle Mirror Descent

**0. Reference:**

Dai et al. "Scalable Bayesian Inference via Particle Mirror Descent" arXiv:1506.03101

**1. Solve VFE minimization with MD**

( assume we can do that exactly ...)

We have some prior $p_0(\theta)$, model $p(x|\theta)$, and observation $D = \{x_1, \cdots, x_N\}$.

Goal: find some tractable

$$q(\theta) \approx p(\theta|D)$$

Variational Inference: solve $q$ by minimizing

$$\mathcal{L}(q) = KL[q(\theta) \| p(\theta|D)] - \log p(D)$$

$$= KL[q(\theta) \| p_0(\theta)] - \sum_{n=1}^{N} \mathbb{E}_q[\log p(x_n|\theta)]$$

Variational Free Energy

( subject to $\int_\theta q(\theta) d\mu = 1$ )

Stochastic approximation: $\mathcal{L}(q) = \mathbb{E}_{x \sim D}[l(q;x)]$,

$$l(q;x) = KL[q(\theta) \| p_0(\theta)] - N\mathbb{E}_q[\log p(x|\theta)]$$

Now solve it with MD!

$$q_{t+1}(\theta) = \underset{\hat{q}(\theta) \in \mathcal{P}}{\arg\min} \left\{ \langle \hat{q}(\theta), \nabla l(q_t;x_t) \rangle_{L_2} + \frac{1}{\gamma_t} KL[\hat{q} \| q_t] \right\}$$

$$\hat{\mathcal{L}}_t(\hat{q})$$

where $x_t \sim D$ is the sample at time $t$,

$\gamma_t$ the learning rate at time $t$,

$$\mathcal{P} = \left\{ p(\theta): \int_\theta p(\theta) d\mu = 1 \right\}$$

details:

$$g_t(\theta) \hat{=} \nabla l(q_t;x_t) = \log q_t(\theta) - \log p_0(\theta) - N\log p(x_t|\theta)$$

$$0 = \nabla \hat{\mathcal{L}}_t(\hat{q}) = \nabla l(q_t;x_t) + \frac{1}{\gamma_t}[\log \hat{q}(\theta) - \log q_t(\theta)]$$

$$\Rightarrow \quad \log q_{t+1}(\theta) \leftarrow \log q_t(\theta) - \gamma_t g_t(\theta)$$

$$\Rightarrow \quad q_{t+1}(\theta) = \frac{q_t(\theta) \exp[-\gamma_t g_t(\theta)]}{Z}$$

$$= \frac{q_t(\theta)^{1-\gamma_t} p_0(\theta)^{\gamma_t} p(x_t|\theta)^{N\gamma_t}}{Z}$$

also called "normalised exponential gradient" if want to read more!

MD update rule:

$$q_{t+1}(\theta) \leftarrow \frac{1}{Z} q_t(\theta)^{1-\gamma_t} p_0(\theta)^{\gamma_t} p(x_t|\theta)^{N\gamma_t}, \quad x_t \sim D$$

<u>problem 1</u>: such $q_{t+1}$ is generally intractable!

Solution: restrict $q \in Q$ (tractable),

Compute MD update, then project it
to the $Q$ family:

$$q_{t+1}(\theta) \leftarrow \text{proj}_Q \left[ \frac{1}{Z} q_t(\theta)^{1-\gamma_t} p_0(\theta)^{\gamma_t} p(x_t|\theta)^{N\gamma_t} \right]$$

<u>problem 2</u>: theory of MD stochastic approximations
works only for running average answer
over time!

Solution: empirical verification

## 2. PMD with weighted particle

assume the starting proposal $\pi(\theta) \approx p(\theta|D)$:

define $\quad q_t(\theta) = \sum_{i=1}^{m} \alpha_i \delta(\theta_i)$

with $\theta_i \sim \pi(\theta)$ and <u>fixed</u> over time

$\Rightarrow$ at time $t$,

$$\alpha_i \leftarrow \frac{1}{Z} \alpha_i^{1-\gamma_t} p_0(\theta_i)^{\gamma_t} p(x_t|\theta_i)^{N\gamma_t},$$

$$Z = \sum_{i=1}^{m} \alpha_i^{1-\gamma_t} p_0(\theta_i)^{\gamma_t} p(x_t|\theta_i)^{N\gamma_t}$$

## 3. PMD with weighted KDE

works when we don't have a good guess $\pi(\theta)$.

define $\quad q_t(\theta) = \sum_{i=1}^{m} \alpha_i K_h(\theta-\theta_i), \quad \theta_i \sim q_{t-1}(\theta)$

$\Rightarrow$ at time $t$,

$$\alpha_i \leftarrow \frac{\exp[-\gamma_t g_t(\theta_i)]}{\sum_{i=1}^{m} \exp[-\gamma_t g_t(\theta_i)]}, \quad \theta_i \sim q_t(\theta)$$

$$= \frac{1}{Z} q_t(\theta_i)^{-\gamma_t} p_0(\theta_i)^{\gamma_t} p(x_t|\theta_i)^{N\gamma_t}$$