

Reminder: Variance reduction in MC methods [See M. Hoffman's RCC slides] ^①

$$E_{p(x)}[f(x)] \approx \frac{1}{N} \sum_{n=1}^N \overbrace{f(x^i)}^F, x^i \sim p(x)$$

Problem: HIGH Variance

Want: F' s.t. $E[F'] = E[F] = E[f(x)]$

$$\text{var}[F'] \leq \text{var}[F]$$

Use $g(x)$ where $\mu_g = E[g(x)]$ is tractable

$$\rightarrow E[f(x)] = E[f(x) - g(x)] + \mu_g$$

\rightarrow SAME Expectations, LOWER Variance if f & g are STRONGLY CORRELATED

(2)

Now consider $\hat{f}(x) = f(x) - a g(x) + a \mu_g$
 a : some constant

$$\text{var}(\hat{f}) = \text{var}(f) + a^2 \text{var}(g) - 2a \text{Cov}(f, g)$$

$$\text{optimal } a = \frac{\text{Cov}(f, g)}{\text{var}(g)} \quad \left. \vphantom{\frac{\text{Cov}(f, g)}{\text{var}(g)}} \right\} \text{Use } \underline{\text{empirical estimate}}$$

Questions : How to choose g ?

For VI, g could be

- + another lower bound [deterministic]
- + Taylor series approximation
- + . . .

Back to VI

③

Consider

$$p(\theta)$$

$$p(y|\theta) = \prod_{n=1}^N p(y_n|\theta)$$

$$\left. \begin{array}{l} p(\theta) \\ p(y|\theta) = \prod_{n=1}^N p(y_n|\theta) \end{array} \right\} \Rightarrow \text{want } p(\theta|y) \propto p(\theta) \cdot p(y|\theta)$$

VI

$$F(q(\theta)) = \left\langle \log \frac{p(y|\theta) \cdot p(\theta)}{q(\theta)} \right\rangle_{q(\theta)}$$

$$= -KL(q(\theta) || p(\theta)) + \sum_{n=1}^N \left\langle \log p(y_n|\theta) \right\rangle_{q(\theta)}$$

However $\left\langle \log p(y_n|\theta) \right\rangle_{q(\theta)}$ ~~is~~ typically not analytically tractable

① approx. \downarrow by another function, e.g. Jaakkola & Jordan 2000
Martin + Murphy 2011

② go stochastic...

Optimising $F(q(\theta))$ requires $\frac{dF}{d\lambda}$ ④

$$\frac{dF}{d\lambda} = - \frac{d}{d\lambda} KL(q(\theta) \| p(\theta)) + \underbrace{\sum_{n=1}^N \int \frac{dq(\theta)}{d\lambda} \log p(y_n | \theta) d\theta}_{= \int q(\theta) \log p(y_n | \theta) \frac{d}{d\lambda} \log q(\theta) d\theta}$$

$$\approx - \underbrace{\frac{d}{d\lambda} KL(q(\theta) \| p(\theta))}_{\frac{dF_1}{d\lambda}} + \sum_{n=1}^N \sum_{m=1}^M \log p(y_n | \theta_m) \frac{d}{d\lambda} \log q(\theta_m)$$

$\theta_m \sim q(\theta)$

SGD: $\lambda \leftarrow \lambda + \eta \frac{dF}{d\lambda}$

But need to lower the variance of $\log p(y_n | \theta_m) \frac{d}{d\lambda} \log q(\theta_m)$

$$\Rightarrow \frac{dF}{d\lambda} = \frac{dF_1}{d\lambda} + \sum_{n=1}^N \sum_{m=1}^M \left(\log p(y_n | \theta_m) - \hat{a}_n(\theta_m) \right) \frac{d}{d\lambda} \log q(\theta_m)$$

i.e. $f \equiv \log p(y_n | \theta_n) \frac{d}{d\lambda} \log q(\theta_n)$

$$g \equiv g_n(\theta_n) \cdot \frac{d}{d\lambda} \log q(\theta_n)$$

$$\hat{a} = \frac{\text{Cov}(f, g)}{\text{Var}(g)}$$

Example: Logistic regression $\log p(y_n | \theta_n) = \log \sigma(y_n x_n^T \theta_n)$
 $q(\theta) = N(\theta; \mu, \Sigma)$

⊕ T & J bound $g_n(\theta) = \ln \sigma(\varepsilon_n) + \frac{1}{2} (y_n x_n^T \theta - \varepsilon_n) + \dots$

Merlin bound $g_n(\theta) = \text{linear} + \text{quadratic pieces}$

⊕ 2nd order Taylor of $\log p(y_n | \theta_n)$ around μ :

$$g_n(\theta) = \ln \sigma(y_n x_n^T \mu) + y_n (1 - \sigma_n) (\theta - \mu)^T x_n + \dots$$

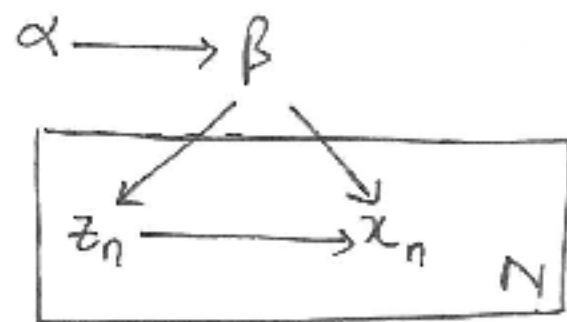
Results See paper

Review of SVI (Hoffman et al)

⑥

Prob of everything.

$$p(x, z, \beta) = p(\beta) \cdot \prod_{n=1}^N p(x_n, z_n | \beta)$$



Want $p(z, \beta) \propto p(x, z, \beta)$

VI $F(q(z, \beta)) = -E_q[\log(q)]$
 $+ E_q[\log p(x, z, \beta)]$

α : fixed (hyper)-params

β : global parameters

$\{z_n\}$: local parameters

$\{x_n\}$: observations

Mean field $q(z, \beta) = q(\beta | \lambda) \prod_{n=1}^N \prod_{j=1}^J q(z_{nj} | \phi_{nj})$

$$\Rightarrow \begin{cases} \lambda = E_q[\eta_g(x, z, \alpha)] \\ \phi_{nj} = E_q[\eta_e(x_n, z_{n,-j}, \beta)] \end{cases}$$

VI - Mean field

Iterate : $t = 1:T$

$$\phi_{nj,t} \leftarrow E_{q_{t-1}}[\eta_e(x_n, z_{n,j}, \beta)]$$

$$\lambda_t \leftarrow E_{q_{t-1}}[\eta_g(x, z, \alpha)]$$

SVI - Mean field

Iterate $t = 1:T$

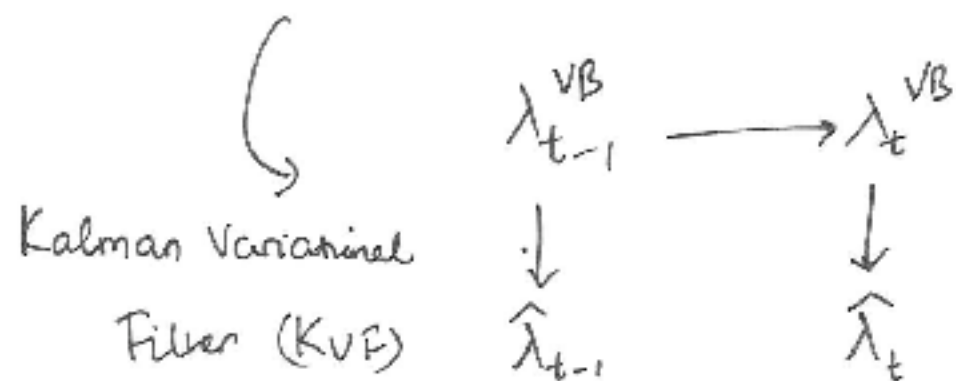
$$\phi_{nj} \leftarrow E_q[\eta_e(x_n, z_{n,j}, \beta)]$$

$$\lambda_t \leftarrow (1 - \rho_t) \lambda_{t-1} + \rho_t \hat{\lambda}_t$$

$$\hat{\lambda}_t = E_{q_{t-1}}[\eta_g(x_m, z_m, \alpha)]$$

mini batch

$$\lambda_t = (1 - \rho_t) \lambda_{t-1} + \rho_t \hat{\lambda}_t$$



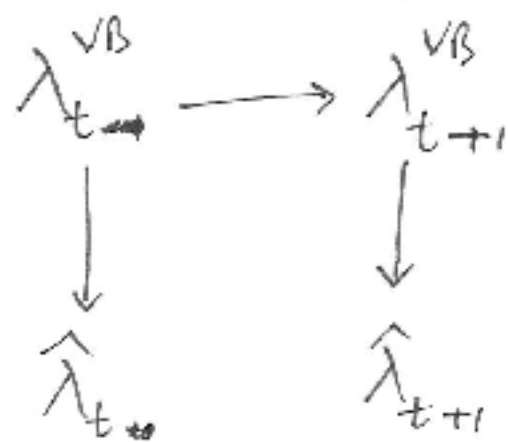
ie observe $\hat{\lambda}_t$: noisy params
infer λ_t^{VB} : true params

More on $\lambda_t = (1 - \rho_t) \cdot \lambda_{t-1} + \rho_t \hat{\lambda}_t$

$$E[\lambda_t - \lambda_t^{VB}] = (1 - \rho_t) \cdot (\lambda_{t-1} - \lambda_t^{VB}) \quad \text{as } E[\hat{\lambda}_t] = \lambda_t^{VB}$$

$$\text{Var}[\lambda_t] = \rho_t^2 \text{Var}[\hat{\lambda}_t]$$

KVF automatically handles bias/variance trade off + step-size



$$P(\lambda_{t+1}^{VB} | \lambda_t^{VB}) = N(\lambda_{t+1}^{VB}; \lambda_t^{VB}, Q)$$

$$P(\hat{\lambda}_t | \lambda_t^{VB}) = N(\hat{\lambda}_t; \lambda_t^{VB}, R)$$

\Rightarrow Every SVI iteration, perform Kalman update

$$\mu_{t|t} = (1 - \rho_t) \mu_{t-1|t-1} + \rho_t \hat{\lambda}_t$$

$$\Sigma_{t|t} = (1 - \rho_t)^{-1} \cdot (\Sigma_{t-1} + Q)$$

$$\text{where } \rho_t = [\Sigma_{t-1} + Q][\Sigma_{t-1} + Q + R]^{-1}$$