

# 1 Stochastic Variational Inference in Linear Regression

Model  $y = \mathbf{X}\beta + \mathbf{e}$ ,  $\mathbf{e} \sim \mathcal{N}(\mathbf{0}, \sigma^2 \mathbf{I}_n)$ .

- $\beta \sim \mathcal{N}(\mu_0, \Sigma_0)$
- $\sigma^2 \sim \text{InvGam}(A, B)$

Variational distributions

- $q(\beta) \sim \mathcal{N}(\mu_q, \Sigma_q)$
- $q(\sigma^2) \sim \text{InvGam}(A_q, B_q)$

Let  $h(\theta) = p(y | \theta) p(\theta)$ .

$$\log h(\theta) = -\frac{n}{2} \log(2\pi\sigma^2) - \frac{1}{2\sigma^2} (y - X\beta)' (y - X\beta) - \frac{1}{2} \log |\Sigma_0| - \frac{p}{2} \log(2\pi) \quad (1)$$

$$- \frac{1}{2} (\beta - \mu_0)' \Sigma_0^{-1} (\beta - \mu_0) - (A + 1) \log \sigma^2 - \frac{B}{\sigma^2} + A \log B - \log \Gamma(A) \quad (2)$$

$$\log q_\lambda(\theta) = -\frac{p}{2} \log(2\pi) - \frac{1}{2} \log |\Sigma_q| - \frac{1}{2} (\beta - \mu_q)' \Sigma_q^{-1} (\beta - \mu_q) - (A_q + 1) \log \sigma^2 \quad (3)$$

$$- \frac{B_q}{\sigma^2} + A_q \log B_q - \log \Gamma(A_q) \quad (4)$$

## 1.1 Reparameterization

Let  $\Sigma_q^{-1} = \Omega = CC'$  where  $C$  is the Cholesky factor lower triangular matrix of  $\Omega$ . Then with  $s \sim \mathcal{N}(\mathbf{0}, \mathbf{I}_p)$ , we can recast  $\beta$  as

$$\beta = C'^{-1}s + \mu_q \quad (5)$$

Now, the log-joint density becomes

$$\log h(\theta) = -\frac{n}{2} \log(2\pi\sigma^2) - \frac{1}{2\sigma^2} (y - XC'^{-1}s - X\mu_q)' (y - XC'^{-1}s - X\mu_q) - \frac{1}{2} \log |\Sigma_0| - \frac{p}{2} \log(2\pi) \quad (6)$$

$$- \frac{1}{2} (C'^{-1}s + \mu_q - \mu_0)' \Sigma_0^{-1} (C'^{-1}s + \mu_q - \mu_0) - (A + 1) \log \sigma^2 - \frac{B}{\sigma^2} \quad (7)$$

$$+ A \log B - \log \Gamma(A) \quad (8)$$

$$\log q_\lambda(\theta) = -\frac{p}{2} \log(2\pi) + \log |C| - \frac{1}{2} s's - (A_q + 1) \log \sigma^2 \quad (9)$$

$$- \frac{B_q}{\sigma^2} + A_q \log B_q - \log \Gamma(A_q) \quad (10)$$

## 2 Multivariate Gaussian Distribution

We will assume that the variational posterior distribution  $q_\lambda(\theta)$  is a multivariate Gaussian distribution where  $\lambda$  denotes the natural parameters in the exponential family representation of the density given by

$$\lambda = \begin{bmatrix} \lambda_1 \\ \lambda_2 \end{bmatrix} = \begin{bmatrix} \Sigma^{-1} \mu \\ -1/2 D_d' \text{vec}(\Sigma^{-1}) \end{bmatrix} \quad (11)$$

$D_d$  is the duplication matrix of order  $d$ , which is the unique  $d^2 \times d(d+1)/2$  matrix of zeros and ones such that  $D_d \text{vech}(A) = \text{vec}(A)$  for symmetric  $d \times d$  matrices  $A$  and its Moore-Penrose inverse is

$$D_d^\dagger = (D_d' D_d)^{-1} D_d' \quad (12)$$

On the contrary,  $L_d$  is the elimination matrix such that

$$L_d \text{vec}(A) = \text{vech}(A) \quad (13)$$

With these notations and natural parameters, the original parameters  $\mu, \Sigma$  can be expressed as

$$\mu = -\frac{1}{2} \left\{ \text{vec}^{-1} \left( D_d^\dagger \lambda_2 \right) \right\}^{-1} \lambda_1, \quad \Sigma = -\frac{1}{2} \left\{ \text{vec}^{-1} \left( D_d^\dagger \lambda_2 \right) \right\}^{-1} \quad (14)$$

Then, the exponential family representation becomes

$$q_\lambda(\theta) = \exp(T(\theta)' \lambda - Z(\lambda)) \quad (15)$$

where  $T(\theta)$  is the sufficient statistic

$$T(\theta) = \begin{bmatrix} \theta \\ \text{vech}(\theta\theta') \end{bmatrix} \quad (16)$$

The Fisher's information matrix  $I_F(\lambda) = \text{Cov}_\lambda(T(\theta))$  has an inverse

$$I_F(\lambda)^{-1} = \begin{bmatrix} \Sigma^{-1} + M' S^{-1} M & -M' S^{-1} \\ -S^{-1} M & S^{-1} \end{bmatrix} \quad (17)$$

where  $M = 2D_d^\dagger(\mu \otimes I_d)$  and  $S = 2D_d^\dagger(\Sigma \otimes \Sigma)D_d^{\dagger'}$ . Now

$$\nabla_\lambda \log q_\lambda(\theta) = \begin{bmatrix} \theta - \mu \\ \text{vech}(\theta\theta' - \Sigma - \mu\mu') \end{bmatrix} \quad (18)$$

### 3 Gradients

$$\nabla_{\mu_q} \log q_\lambda(\theta) = -\Sigma_q^{-1} (\mu_q - \beta) \quad (19)$$

$$\nabla_{\Sigma_q} \log q_\lambda(\theta) = -\frac{1}{2} \left( \Sigma_q^{-1} + (\beta - \mu_q)(\beta - \mu_q)' \right) \quad (20)$$

$$\nabla_{A_q} \log q_\lambda(\theta) = -\log \sigma^2 + \log B_q - \psi(A_q) \quad (21)$$

$$\nabla_{B_q} \log q_\lambda(\theta) = -\frac{1}{\sigma^2} + \frac{A_q}{B_q} \quad (22)$$

Lower bound, log-derivative trick

$$\nabla_\lambda \mathcal{L}(\lambda) = \nabla_\lambda \int (\log h(\theta) - \log q_\lambda(\theta)) q_\lambda(\theta) \theta \quad (23)$$

$$= \int \log h(\theta) \nabla_\lambda \log q_\lambda(\theta) q_\lambda(\theta) d\theta - \int \log q_\lambda(\theta) \nabla_\lambda \log q_\lambda(\theta) q_\lambda(\theta) \theta \quad (24)$$

$$= \int \nabla_\lambda \log q_\lambda(\theta) (\log h(\theta) - \log q_\lambda(\theta)) q_\lambda(\theta) d\theta \quad (25)$$

$$\approx \frac{1}{S} \sum_{s=1}^S \nabla_\lambda \log q_\lambda(\theta^{(s)}) (\log h(\theta^{(s)}) - \log q_\lambda(\theta^{(s)})) \quad (26)$$

where  $\theta^{(s)} \sim q_\lambda(\theta)$ ,  $s = 1, \dots, S$ .

### 4 Step Size

For the step size sequence  $\rho_t$ , the Robbins-Monro approximation states that it should satisfy two conditions

- $\sum_{t=1}^{\infty} \rho_t = \infty$
- $\sum_{t=1}^{\infty} \rho_t^2 < \infty$

which the first indicates that the step size should be large enough that the stochastic search is able to sweep through all possible zones of the parameter space, and the second indicates that despite the first condition, the variation should not be so big that it is beyond control. The recommended Robbins-Monro sequence is

$$\rho = (t + \tau)^{-\kappa} \quad (27)$$

where  $\kappa \in (0.5, 1]$ , the forgetting rate, controls how quickly old information decays and  $\tau \geq 0$ , the delay, downweights early iterations. In many cases,  $\tau$  is set to 1.

## 5 Inverse Gamma Distribution

The log-density of an inverse-gamma distribution is

$$A \log B - \log \Gamma(A) - (A + 1) \log \theta - B/\theta \quad (28)$$

Therefore for  $A, B$ , the gradient is given as

$$\nabla_{\lambda} q_{\lambda}(\theta) = \begin{bmatrix} \log B - \psi(A) - \log \theta \\ A/B - 1/\theta \end{bmatrix} \quad (29)$$

Now to get the Fisher's information matrix, we need the following

$$\text{Var}(-\log \theta) = \psi_1(A) \quad (30)$$

$$\text{Var}(-1/\theta) = A/B^2 \quad (31)$$

$$\text{Cov}(-\log \theta, -1/\theta) = \mathbf{E} \left( \frac{1}{\theta} \log \theta \right) - \mathbf{E}(\log \theta) \mathbf{E} \left( \frac{1}{\theta} \right) \quad (32)$$

$$= -\mathbf{E} \left( \frac{1}{\theta} \log \frac{1}{\theta} \right) + \mathbf{E} \left( \log \frac{1}{\theta} \right) \mathbf{E} \left( \frac{1}{\theta} \right) \quad (33)$$

$$= \frac{A}{B} (\psi(A) - \psi(A + 1)) \quad (34)$$

$$G_{\theta} = \begin{bmatrix} \psi_1(A) & -\frac{1}{B} \\ -\frac{1}{B} & A/B^2 \end{bmatrix} \quad (35)$$

since  $\psi(a + 1) = \psi(a) + 1/a$

## 6 Reparameterization

If  $q(\beta) = \mathcal{N}(\mu_q, \Sigma_q)$  and if we let  $\Sigma_q^{-1} = C C'$  and  $\phi(z) = \mathcal{N}(\mathbf{0}, \mathbf{I}_p)$ ,

- $\beta = C^{-1'} z + \mu_q$
- $dz = |C| d\beta$
- $q(\beta) = |C| \phi(z)$

are all true. Then

$$\int q(\beta) \log \frac{h(\theta)}{q(\beta)} d\beta = \int \phi(z) \log \frac{h(C^{-1'} z + \mu_q)}{\phi(z) |C|} dz \quad (36)$$

$$= \mathbf{E}_{\phi(z)} \left[ h(C^{-1'} z + \mu_q) \right] - \log |C| - \mathcal{H}(z) \quad (37)$$