

# BAYESIAN LEARNING - LECTURE 9

Mattias Villani

**Division of Statistics and Machine Learning  
Department of Computer and Information Science  
Linköping University**

# LECTURE OVERVIEW

- ▶ Variational Bayes
- ▶ RStan demo

# VARIATIONAL BAYES

- ▶ Let  $\theta = (\theta_1, \dots, \theta_p)$ . Approximate the posterior  $p(\theta|y)$  with a (simpler) distribution  $q(\theta)$ .
- ▶ We have already seen:  $q(\theta) = N[\tilde{\theta}, J_{\mathbf{y}}^{-1}(\tilde{\theta})]$ .
- ▶ **Mean field approximation**

$$q(\theta) = \prod_{i=1}^M q_i(\theta_i)$$

- ▶ **Parametric VB**, where  $q_{\lambda}(\theta)$  is a parametric family with parameters  $\lambda$ .
- ▶ Find the  $q(\theta)$  that **minimizes the Kullback-Leibler distance** between the true posterior  $p$  and the approximation  $q$ :

$$KL(q, p) = \int q(\theta) \ln \frac{q(\theta)}{p(\theta|y)} d\theta = E_q \left[ \ln \frac{q(\theta)}{p(\theta|y)} \right].$$

# MEAN FIELD APPROXIMATION

- ▶ Factorization

$$q(\theta) = \prod_{i=1}^p q_i(\theta_i)$$

- ▶ No specific functional forms are assumed for the  $q_i(\theta)$ .
- ▶ Optimal densities can be shown to satisfy:

$$q_i(\theta) \propto \exp(E_{-\theta_i} \ln p(\mathbf{y}, \theta))$$

where  $E_{-\theta_i}(\cdot)$  is the expectation with respect to  $\prod_{i \neq j} q_j(\theta_j)$ .

- ▶ **Structured mean field approximation.** Group subset of parameters in tractable blocks.

# MEAN FIELD APPROXIMATION - ALGORITHM

- ▶ Initialize:  $q_2^*(\theta_2), \dots, q_M^*(\theta_p)$
- ▶ Repeat until convergence:
  - ▶  $q_1^*(\theta_1) \leftarrow \frac{\exp[E_{-\theta_1} \ln p(\mathbf{y}, \theta)]}{\int \exp[E_{-\theta_1} \ln p(\mathbf{y}, \theta)] d\theta_1}$
  - ▶  $\vdots$
  - ▶  $q_p^*(\theta_p) \leftarrow \frac{\exp[E_{-\theta_p} \ln p(\mathbf{y}, \theta)]}{\int \exp[E_{-\theta_p} \ln p(\mathbf{y}, \theta)] d\theta_p}$
- ▶ Note: we make no assumptions about parametric form of the  $q_i(\theta)$ , but the optimal  $q_i(\theta)$  often turn out to be parametric (normal, gamma etc).
- ▶ The updates above then boil down to just updating of hyperparameters in the optimal densities.

# MEAN FIELD APPROXIMATION - NORMAL MODEL

- ▶ **Model:**  $X_i | \theta, \sigma^2 \stackrel{iid}{\sim} N(\theta, \sigma^2)$ .
- ▶ **Prior:**  $\theta \sim N(\mu_0, \tau_0^2)$  **independent** of  $\sigma^2 \sim \text{Inv} - \chi^2(\nu_0, \sigma_0^2)$ .
- ▶ **Mean-field approximation:**  $q(\theta, \sigma^2) = q_\theta(\theta) \cdot q_{\sigma^2}(\sigma^2)$ .
- ▶ Optimal densities

$$q_\theta^*(\theta) \propto \exp \left[ E_{q(\sigma^2)} \ln p(\theta, \sigma^2, \mathbf{x}) \right]$$
$$q_{\sigma^2}^*(\sigma^2) \propto \exp \left[ E_{q(\theta)} \ln p(\theta, \sigma^2, \mathbf{x}) \right]$$

# NORMAL MODEL - VB ALGORITHM

- Variational density for  $\sigma^2$

$$\sigma^2 \sim \text{Inv} - \chi^2 (\tilde{\nu}_n, \tilde{\sigma}_n^2)$$

where  $\tilde{\nu}_n = \nu_0 + n$  and  $\tilde{\sigma}_n^2 = \frac{\nu_0 \sigma_0^2 + \sum_{i=1}^n (x_i - \tilde{\mu}_n)^2 + n \cdot \tilde{\tau}_n^2}{\nu_0 + n}$

- Variational density for  $\mu$

$$\theta \sim N (\tilde{\mu}_n, \tilde{\tau}_n^2)$$

where

$$\tilde{\tau}_n^2 = \frac{1}{\frac{n}{\tilde{\sigma}_n^2} + \frac{1}{\tau_0^2}}$$

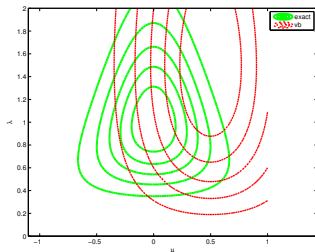
$$\tilde{\mu}_n = \tilde{w} \bar{x} + (1 - \tilde{w}) \mu_0,$$

where

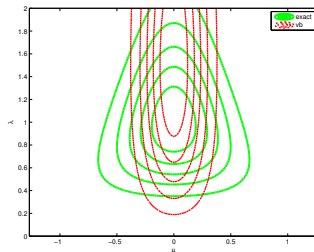
$$\tilde{w} = \frac{\frac{n}{\tilde{\sigma}_n^2}}{\frac{n}{\tilde{\sigma}_n^2} + \frac{1}{\tau_0^2}}$$

# NORMAL EXAMPLE FROM MURPHY ( $\lambda = 1/\sigma^2$ )

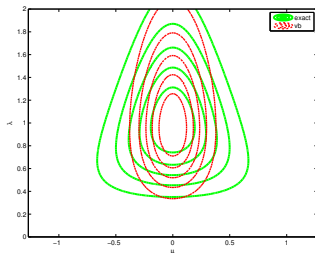
Initial values



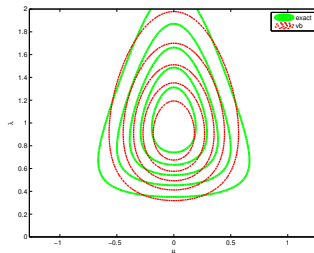
After updating  $q_{\mu}$



After updating  $q_{\sigma^2}$



At convergence





# PROBIT REGRESSION

- **Model:**

$$\Pr(y_i = 1 | \mathbf{x}_i) = \Phi(\mathbf{x}_i^T \boldsymbol{\beta})$$

- **Prior:**  $\boldsymbol{\beta} \sim N(\boldsymbol{\mu}_\beta, \boldsymbol{\Sigma}_\beta)$

- **Latent variable formulation** with  $\mathbf{u} = (u_1, \dots, u_n)'$

$$\mathbf{u} | \boldsymbol{\beta} \sim N(\mathbf{X}\boldsymbol{\beta}, \mathbf{1})$$

and

$$y_i = \begin{cases} 0 & \text{if } u_i \leq 0 \\ 1 & \text{if } u_i > 0 \end{cases}$$

- **Factorized variational approximation**

$$q(\mathbf{u}, \boldsymbol{\beta}) = q_{\mathbf{u}}(\mathbf{u})q_{\boldsymbol{\beta}}(\boldsymbol{\beta})$$

# PROBIT REGRESSION - UPDATING $\beta$

- It can be shown that the VB posterior is

$$\beta \sim N \left( \tilde{\mu}_\beta, \left( \mathbf{X}^T \mathbf{X} + \Sigma_\beta^{-1} \right)^{-1} \right)$$

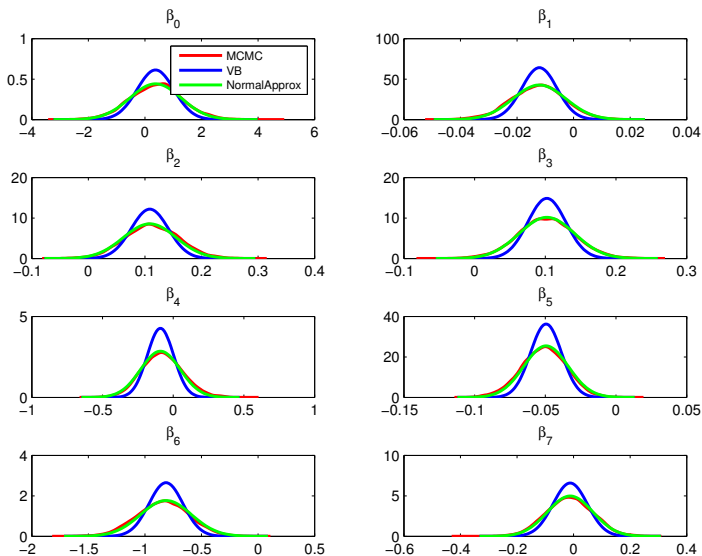
where

$$\tilde{\mu}_\beta = \left( \mathbf{X}^T \mathbf{X} + \Sigma_\beta^{-1} \right)^{-1} \left( \mathbf{X}^T \tilde{\mu}_\mathbf{u} + \Sigma_\beta^{-1} \mu_\beta \right)$$

and

$$\tilde{\mu}_\mathbf{u} = \mathbf{X} \tilde{\mu}_\beta + \frac{\phi(\mathbf{X} \tilde{\mu}_\beta)}{\Phi(\mathbf{X} \tilde{\mu}_\beta)^{\mathbf{y}} [\Phi(\mathbf{X} \tilde{\mu}_\beta) - \mathbf{1}_n]^{\mathbf{1}_n - \mathbf{y}}}.$$

# PROBIT EXAMPLE (N=200 OBSERVATIONS)



# PROBIT EXAMPLE

