

# **BAYESIAN LARGE SAMPLE THEORY AND POSTERIOR APPROXIMATION**

PHD COURSE IN STATISTICAL INFERENCE

MATTIAS VILLANI

**DEPARTMENT OF STATISTICS  
STOCKHOLM UNIVERSITY**

**AND**

**DEPARTMENT OF COMPUTER AND INFORMATION SCIENCE  
LINKÖPING UNIVERSITY**

- **Large sample theory**
- **Variational inference**

# GENERALIZED LINEAR MODELS

- Response  $y$  conditional on a set of covariates  $\mathbf{x}$  belongs to the **exponential family with dispersion parameter  $\kappa$**

$$p(y|\theta, \kappa, \mathbf{x}) = h(y, \kappa) \exp \left( \frac{\phi(\theta)^T \mathbf{t}(y) - A(\theta)}{d(\kappa)} \right)$$

- The conditional mean of  $\mu = E(y|\mathbf{x})$  is a function of a **linear predictor**  $\eta = \mathbf{x}^T \beta$  through a **link function**

$$g(\mu) = \mathbf{x}^T \beta$$

- Example: Poisson regression, where  $y|\mathbf{x} \sim \text{Pois}(\exp(\mathbf{x}^T \beta))$  is a GLM with log-link:  $\log \mu = \mathbf{x}^T \beta$ .
- Logistic regression:  $y_i|\mathbf{x}_i \sim \text{Bern}(\theta_i)$ , where  $\theta_i = \frac{\exp(\mathbf{x}_i^T \beta)}{1 + \exp(\mathbf{x}_i^T \beta)}$ .
- Prior  $\beta \sim N(\mathbf{0}, \tau^2 I)$ .
- Posterior is typically non-standard. What to do?

- **Taylor expansion of log-posterior** around mode  $\theta = \tilde{\theta}$ :

$$\begin{aligned}\ln p(\theta|\mathbf{y}) &= \ln p(\tilde{\theta}|\mathbf{y}) + \frac{\partial \ln p(\theta|\mathbf{y})}{\partial \theta} \Big|_{\theta=\tilde{\theta}} (\theta - \tilde{\theta}) \\ &\quad + \frac{1}{2!} \frac{\partial^2 \ln p(\theta|\mathbf{y})}{\partial \theta^2} \Big|_{\theta=\tilde{\theta}} (\theta - \tilde{\theta})^2 + \dots\end{aligned}$$

- From the definition of the posterior mode:

$$\frac{\partial \ln p(\theta|\mathbf{y})}{\partial \theta} \Big|_{\theta=\tilde{\theta}} = 0$$

- So, in **large samples** (higher order terms negligible):

$$p(\theta|\mathbf{y}) \approx p(\tilde{\theta}|\mathbf{y}) \exp \left( -\frac{1}{2} J_{\mathbf{y}}(\tilde{\theta}) (\theta - \tilde{\theta})^2 \right)$$

where  $J_{\mathbf{y}}(\tilde{\theta}) = -\frac{\partial^2 \ln p(\theta|\mathbf{y})}{\partial \theta^2} \Big|_{\theta=\tilde{\theta}}$  is the **observed information**.

- **Approximate normal posterior** in large samples.

$$\theta|\mathbf{y} \stackrel{approx}{\sim} N[\tilde{\theta}, J_{\mathbf{y}}^{-1}(\tilde{\theta})]$$

## EXAMPLE: GAMMA POSTERIOR

- **Poisson model:**  $\theta|y_1, \dots, y_n \sim \text{Gamma}(\alpha + \sum_{i=1}^n y_i, \beta + n)$

$$\log p(\theta|y_1, \dots, y_n) \propto (\alpha + \sum_{i=1}^n y_i - 1) \log \theta - \theta(\beta + n)$$

- First derivative of log density

$$\frac{\partial \ln p(\theta|\mathbf{y})}{\partial \theta} = \frac{\alpha + \sum_{i=1}^n y_i - 1}{\theta} - (\beta + n)$$

$$\tilde{\theta} = \frac{\alpha + \sum_{i=1}^n y_i - 1}{\beta + n}$$

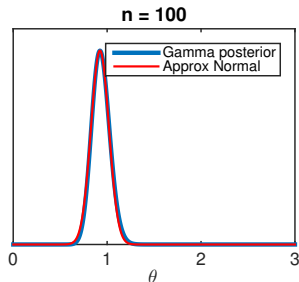
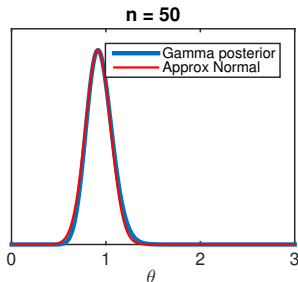
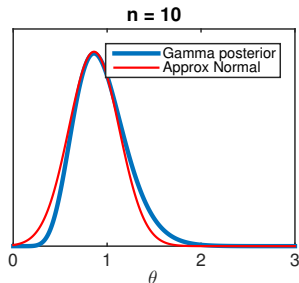
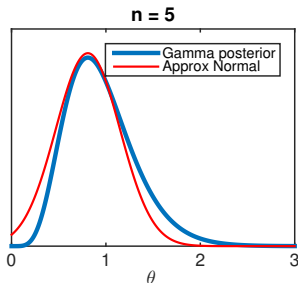
- Second derivative at mode  $\tilde{\theta}$

$$\frac{\partial^2 \ln p(\theta|\mathbf{y})}{\partial \theta^2} \Big|_{\theta=\tilde{\theta}} = -\frac{\alpha + \sum_{i=1}^n y_i - 1}{\left(\frac{\alpha + \sum_{i=1}^n y_i - 1}{\beta + n}\right)^2} = -\frac{(\beta + n)^2}{\alpha + \sum_{i=1}^n y_i - 1}$$

- **Normal approximation**

$$N \left[ \frac{\alpha + \sum_{i=1}^n y_i - 1}{\beta + n}, \frac{\alpha + \sum_{i=1}^n y_i - 1}{(\beta + n)^2} \right]$$

# EXAMPLE: GAMMA POSTERIOR



# NORMAL APPROXIMATION OF POSTERIOR

- $\theta|\mathbf{y} \stackrel{approx}{\sim} N[\tilde{\theta}, J_{\mathbf{y}}^{-1}(\tilde{\theta})]$  works also when  $\theta$  is a vector.
- How to compute  $\tilde{\theta}$  and  $J_{\mathbf{y}}(\tilde{\theta})$ ?
- Standard **optimization routines** may be used. (optim.r).
  - **Input:** expression proportional to  $\log p(\theta|\mathbf{y})$ . Initial values.
  - **Output:**  $\log p(\tilde{\theta}|\mathbf{y})$ ,  $\tilde{\theta}$  and Hessian matrix  $(-J_{\mathbf{y}}(\tilde{\theta}))$ .
- **Re-parametrization** may improve normal approximation. [Don't forget the **Jacobian!**]
  - If  $\theta \geq 0$  use  $\phi = \log(\theta)$ .
  - If  $0 \leq \theta \leq 1$ , use  $\phi = \ln[\theta/(1 - \theta)]$ .
- **Heavy tailed approximation:**  $\theta|\mathbf{y} \stackrel{approx}{\sim} t_v[\tilde{\theta}, J_{\mathbf{y}}^{-1}(\tilde{\theta})]$  for suitable degrees of freedom  $v$ .

# REPARAMETRIZATION - GAMMA POSTERIOR

- Poisson model. Reparameterize to  $\phi = \log(\theta)$ .
- Change-of-variables formula from a basic probability course

$$\log p(\phi|y_1, \dots, y_n) \propto (\alpha + \sum_{i=1}^n y_i - 1)\phi - \exp(\phi)(\beta + n) + \phi$$

- Taking first and second derivatives and evaluating at  $\tilde{\phi}$  gives

$$\tilde{\phi} = \log\left(\frac{\alpha + \sum_{i=1}^n y_i}{\beta + n}\right) \text{ and } \frac{\partial^2 \ln p(\phi|y)}{\partial \phi^2} \Big|_{\phi=\tilde{\phi}} = \alpha + \sum_{i=1}^n y_i$$

- So, the normal approximation for  $p(\phi|y_1, \dots, y_n)$  is

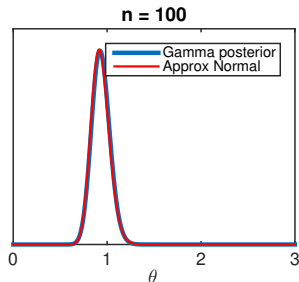
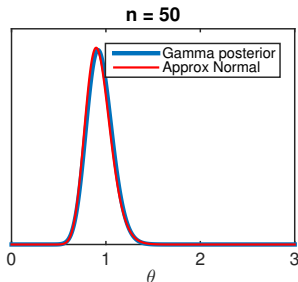
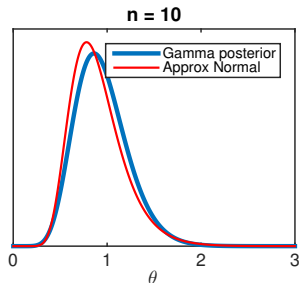
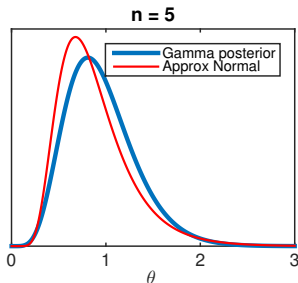
$$\phi = \log(\theta) \sim N\left[\log\left(\frac{\alpha + \sum_{i=1}^n y_i}{\beta + n}\right), \frac{1}{\alpha + \sum_{i=1}^n y_i}\right]$$

which means that  $p(\theta|y_1, \dots, y_n)$  is log-normal:

$$\theta|\mathbf{y} \sim LN\left[\log\left(\frac{\alpha + \sum_{i=1}^n y_i}{\beta + n}\right), \frac{1}{\alpha + \sum_{i=1}^n y_i}\right]$$



# REPARAMETRIZATION - GAMMA POSTERIOR



- Even if the posterior of  $\theta$  is approx normal, **interesting functions** of  $g(\theta)$  may not be (e.g. predictions).
- But approximate posterior of  $g(\theta)$  can be obtained by **simulating** from  $N[\tilde{\theta}, J_{\mathbf{y}}^{-1}(\tilde{\theta})]$ .
- Posterior of **Gini coefficient**
  - Model:  $x_1, \dots, x_n | \mu, \sigma^2 \sim LN(\mu, \sigma^2)$ .
  - Let  $\phi = \log(\sigma^2)$ . And  $\theta = (\mu, \phi)$ .
  - Joint posterior  $p(\mu, \phi)$  may be approximately normal:  
 $\theta | \mathbf{y} \stackrel{approx}{\sim} N[\tilde{\theta}, J_{\mathbf{y}}^{-1}(\tilde{\theta})]$ .
  - Simulate  $\theta^{(1)}, \dots, \theta^{(N)}$  from  $N[\tilde{\theta}, J_{\mathbf{y}}^{-1}(\tilde{\theta})]$ . Compute  $\sigma^{(1)}, \dots, \sigma^{(N)}$ .
  - Compute  $G^{(i)} = 2\Phi(\sigma^{(i)} / \sqrt{2})$  for  $i = 1, \dots, N$ .

- Let  $\theta = (\theta_1, \dots, \theta_p)$ . Approximate the posterior  $p(\theta|y)$  with a (simpler) distribution  $q(\theta)$ .
- Before: **Normal approximation** from optimization:  
 $q(\theta) = N[\tilde{\theta}, J_{\mathbf{y}}^{-1}(\tilde{\theta})]$ .
- **Mean field Variational Bayes (VB)**

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

- 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 VB** is based on factorized approximation:

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

- **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. Similar to Gibbs sampling.

■ Initialize:  $q_2^*(\theta_2), \dots, q_M^*(\theta_p)$

■ Repeat until convergence:

$$\begin{aligned} \bullet \quad q_1^*(\theta_1) &\leftarrow \frac{\exp[E_{-\theta_1} \ln p(\mathbf{y}, \theta)]}{\int \exp[E_{-\theta_1} \ln p(\mathbf{y}, \theta)] d\theta_1} \\ \bullet \quad &\vdots \\ \bullet \quad q_p^*(\theta_p) &\leftarrow \frac{\exp[E_{-\theta_p} \ln p(\mathbf{y}, \theta)]}{\int \exp[E_{-\theta_p} \ln p(\mathbf{y}, \theta)] d\theta_p} \end{aligned}$$

■ Note: no assumptions about parametric form of the  $q_i(\theta)$ .

■ Optimal  $q_i(\theta)$  often **turn out** to be parametric (normal etc).

■ Just update hyperparameters in the optimal densities.

- **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

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

## ■ Variational density for $\sigma^2$

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

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

## ■ Variational density for $\theta$

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

where

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

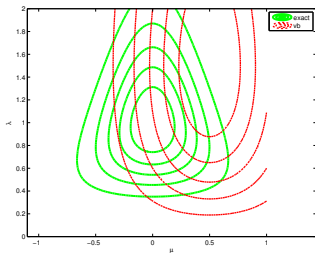
$$\tilde{\mu}_n = \tilde{W} \bar{X} + (1 - \tilde{W}) \mu_0,$$

where

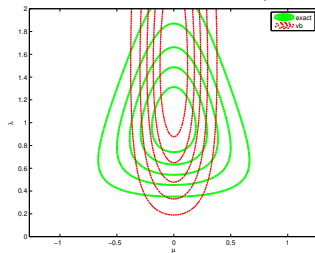
$$\tilde{W} = \frac{\frac{n}{\tilde{\sigma}_n^2}}{\frac{n}{\tilde{\sigma}_n^2} + \frac{1}{\tau_0^2}}$$

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

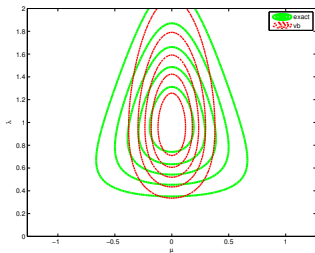
Initial values



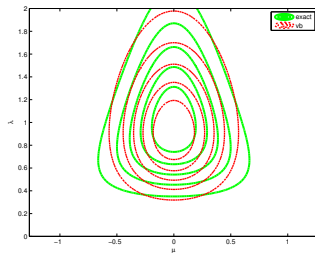
After updating  $q_\mu$



After updating  $q_{\sigma^2}$



At convergence





- **Model:**

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

- **Prior:**  $\beta \sim N(\mathbf{0}, \Sigma_\beta)$ . For example:  $\Sigma_\beta = \tau^2 I$ .

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

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

and

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

- Factorized **variational approximation**

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

## ■ VB posterior

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

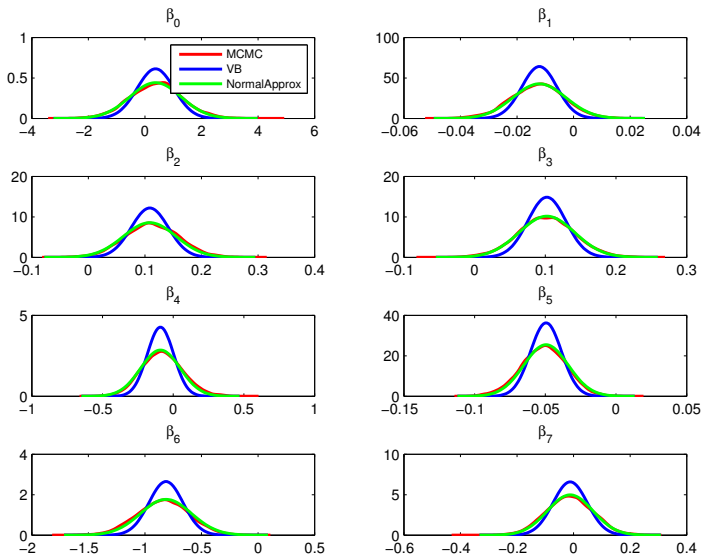
where

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

and

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

# PROBIT EXAMPLE (N=200 OBSERVATIONS)



# PROBIT EXAMPLE

