

Workshop: Intro to Bayesian Learning

Lecture 4 - Bayesian Classification and Posterior Approximation

Mattias Villani

**Department of Statistics
Stockholm University**



mattiasvillani.com



@matvil



@matvil



mattiasvillani

Overview

- Bayesian logistic regression
- Posterior approximation

Binary regression

■ Logistic regression

$$\Pr(y_i = 1 \mid \mathbf{x}_i) = \frac{\exp(\mathbf{x}_i^\top \boldsymbol{\beta})}{1 + \exp(\mathbf{x}_i^\top \boldsymbol{\beta})}.$$

■ Probit regression

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

■ Multi-class ($c = 1, 2, \dots, C$) logistic regression

$$\Pr(y_i = c \mid \mathbf{x}_i) = \frac{\exp(\mathbf{x}_i^\top \boldsymbol{\beta}_c)}{\sum_{k=1}^C \exp(\mathbf{x}_i^\top \boldsymbol{\beta}_k)}$$

■ Likelihood

$$p(\mathbf{y} \mid \mathbf{X}, \boldsymbol{\beta}) = \prod_{i=1}^n \frac{(\exp(\mathbf{x}_i^\top \boldsymbol{\beta}))^{y_i}}{1 + \exp(\mathbf{x}_i^\top \boldsymbol{\beta})}.$$

- Problem: no conjugate prior. **Posterior is intractable**. Now what?

Likelihood information

- **Observed information** in likelihood $\ln p(\mathbf{x}|\theta)$ for **given** data $\mathbf{x} = (x_1, \dots, x_n)^\top$

$$J_{\mathbf{x}}(\hat{\theta}) = -\frac{\partial^2 \ln p(\mathbf{x}|\theta)}{\partial \theta^2} \Big|_{\theta=\hat{\theta}}$$

where $\hat{\theta}$ is the maximum likelihood estimate.

- Multiparameter **observed information matrix**

$$J_{\mathbf{x}}(\hat{\boldsymbol{\theta}}) = -\frac{\partial^2 \ln p(\mathbf{x}|\boldsymbol{\theta})}{\partial \boldsymbol{\theta} \partial \boldsymbol{\theta}^\top} \Big|_{\boldsymbol{\theta}=\hat{\boldsymbol{\theta}}}$$

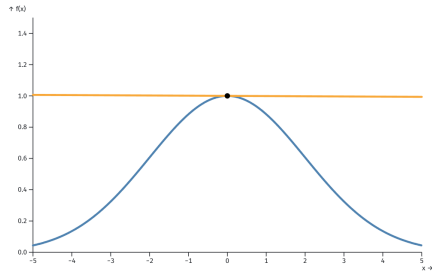
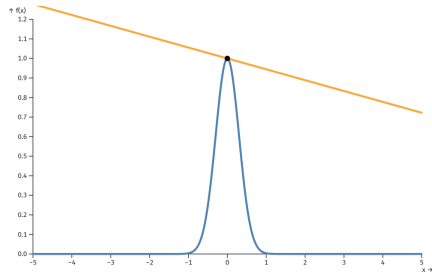
- Example: $\boldsymbol{\theta} = (\theta_1, \theta_2)^\top$

$$\frac{\partial^2 \ln p(\mathbf{x}|\boldsymbol{\theta})}{\partial \boldsymbol{\theta} \partial \boldsymbol{\theta}^\top} = \begin{pmatrix} \frac{\partial^2 \ln p(\mathbf{x}|\boldsymbol{\theta})}{\partial \theta_1^2} & \frac{\partial^2 \ln p(\mathbf{x}|\boldsymbol{\theta})}{\partial \theta_1 \partial \theta_2} \\ \frac{\partial^2 \ln p(\mathbf{x}|\boldsymbol{\theta})}{\partial \theta_1 \partial \theta_2} & \frac{\partial^2 \ln p(\mathbf{x}|\boldsymbol{\theta})}{\partial \theta_2^2} \end{pmatrix}$$

Second derivative measures curvature

Point where
derivative is
evaluated, x

0



Posterior asymptotics

Normal posterior approximation

The posterior can in large samples be approximated by

$$\boldsymbol{\theta}|\mathbf{x} \stackrel{a}{\sim} \mathcal{N}\left(\tilde{\boldsymbol{\theta}}, J_{\mathbf{x}}^{-1}(\tilde{\boldsymbol{\theta}})\right)$$

where $\tilde{\boldsymbol{\theta}}$ is the posterior mode and

$$J_{\mathbf{x}}(\tilde{\boldsymbol{\theta}}) = -\frac{\partial^2 \ln p(\mathbf{x}|\boldsymbol{\theta})p(\boldsymbol{\theta})}{\partial \boldsymbol{\theta} \partial \boldsymbol{\theta}^\top} \Big|_{\boldsymbol{\theta}=\tilde{\boldsymbol{\theta}}}$$

is the $d \times d$ observed posterior information matrix at $\tilde{\boldsymbol{\theta}}$.

- Important: sufficient with proportional form

$$\log p(\boldsymbol{\theta}|\mathbf{x}) = \log p(\mathbf{x}|\boldsymbol{\theta}) + \log p(\boldsymbol{\theta})$$

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

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

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

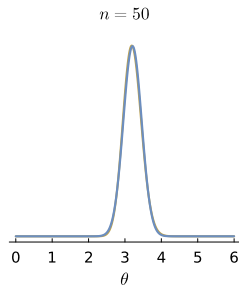
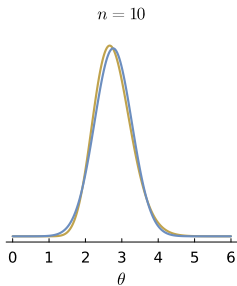
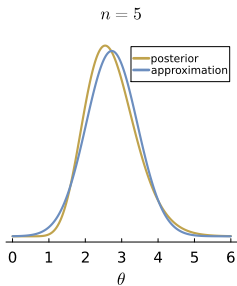
- Second derivative at mode $\tilde{\theta}$

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

- **Normal approximation**

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

Example: gamma posterior for eBay bidders data

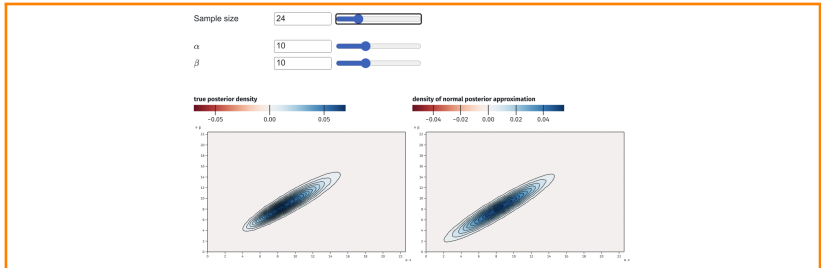


Numerical Normal approximation of posterior

- Standard **optimization routines** may be used to obtain the normal approximation (e.g. `optim` in R).
 - ▶ **Input**: expression proportional to $\log p(\boldsymbol{\theta}|\mathbf{x})$. Initial values.
 - ▶ **Output**: $\log p(\tilde{\boldsymbol{\theta}}|\mathbf{x})$, $\tilde{\boldsymbol{\theta}}$ and Hessian matrix ($-J_{\mathbf{x}}(\tilde{\boldsymbol{\theta}})$).
- **Automatic differentiation** - efficient derivatives on computer.
- **Re-parametrization** may improve normal approximation:
 - ▶ If $\theta \geq 0$ use $\phi = \log(\theta)$.
 - ▶ If $0 \leq \theta \leq 1$, use $\phi = \ln[\theta/(1 - \theta)]$.
 - ▶ Don't forget the **Jacobian**!
- Posterior approximation of functions $g(\boldsymbol{\theta})$ by simulation from

$$\boldsymbol{\theta}|\mathbf{y} \stackrel{\text{a}}{\sim} N\left(\tilde{\boldsymbol{\theta}}, J_{\mathbf{x}}^{-1}(\tilde{\boldsymbol{\theta}})\right)$$

Normal approx of posterior in Beta regression



Logistic regression - who survived the Titanic?

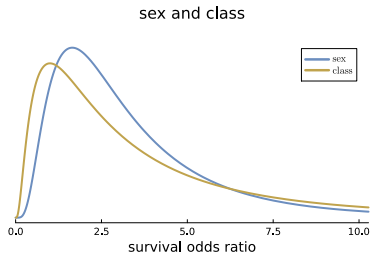
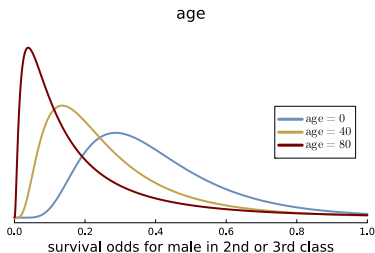
■ Prior

$$\beta \sim N(\mu, \Omega)$$

with

$$\mu = (-1, -1/80, 1, 1)^\top$$

$$\Omega = \begin{pmatrix} 0.25 & 0 & 0 & 0 \\ 0 & 1/(80^2) & 0 & 0 \\ 0 & 0 & 0.5 & 0 \\ 0 & 0 & 0 & 1 \end{pmatrix}$$

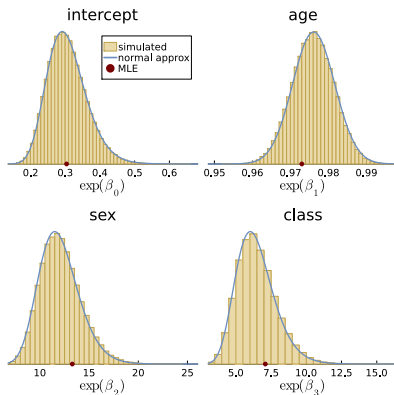


Logistic regression - who survived the Titanic?

Normal posterior approximation

$$\beta | \mathbf{y} \sim N\left(\tilde{\beta}, J_{\mathbf{y}}^{-1}(\tilde{\beta})\right).$$

- Means that the posterior of each β_j is univariate normal.
- Marginal posterior for each $\exp(\beta_j)$ is **LogNormal**.



Logistic regression - who survived the Titanic?

- Comparison with non-informative prior $\beta \sim N(0, 10^2 I_p)$.

