

Detailed Hand Calculations for Bayesian Models

Part 1: Beta-Binomial Model

1.1 Model Setup

Given:

- Prior: $\theta \sim \text{Beta}(\alpha_0, \beta_0)$
- Data: $y = 2$ successes in $N = 12$ trials
- Likelihood: $y|\theta \sim \text{Binomial}(N, \theta)$

Let's use $\alpha_0 = 1, \beta_0 = 1$ (uniform prior)

1.2 Prior Distribution

The Beta prior density is:

$$\begin{aligned} p(\theta) &= [1/B(\alpha_0, \beta_0)] \times \theta^{\alpha_0-1} \times (1-\theta)^{\beta_0-1} \\ &= [1/B(1, 1)] \times \theta^0 \times (1-\theta)^0 \\ &= 1 \quad (\text{uniform on } [0,1]) \end{aligned}$$

where $B(\alpha_0, \beta_0) = \Gamma(\alpha_0)\Gamma(\beta_0)/\Gamma(\alpha_0+\beta_0) = 1 \times 1/1 = 1$

1.3 Likelihood

The binomial likelihood is:

$$\begin{aligned} p(y|\theta) &= C(N,y) \times \theta^y \times (1-\theta)^{N-y} \\ &= C(12,2) \times \theta^2 \times (1-\theta)^{10} \\ &= 66 \times \theta^2 \times (1-\theta)^{10} \end{aligned}$$

where $C(12,2) = 12!/(2! \times 10!) = (12 \times 11)/(2 \times 1) = 66$

1.4 Joint Distribution

$$\begin{aligned} p(y, \theta) &= p(y|\theta) \times p(\theta) \\ &= 66 \times \theta^2 \times (1-\theta)^{10} \times 1 \\ &= 66 \times \theta^2 \times (1-\theta)^{10} \end{aligned}$$

1.5 Log Joint Distribution

$$\begin{aligned}\log p(y, \theta) &= \log(66) + 2 \times \log(\theta) + 10 \times \log(1-\theta) \\ \dots &= \text{constant} + 2 \times \log(\theta) + 10 \times \log(1-\theta)\end{aligned}$$

More generally:

$$\begin{aligned}\log p(y, \theta) &= (\alpha_0 + y - 1) \times \log(\theta) + (\beta_0 + N - y - 1) \times \log(1-\theta) + \text{constant} \\ \dots &= (1 + 2 - 1) \times \log(\theta) + (1 + 12 - 2 - 1) \times \log(1-\theta) + \text{constant} \\ &= 2 \times \log(\theta) + 10 \times \log(1-\theta) + \text{constant}\end{aligned}$$

1.6 Gradient of Log Joint

$$\begin{aligned}\frac{d}{d\theta} \log p(y, \theta) &= (\alpha_0 + y - 1)/\theta - (\beta_0 + N - y - 1)/(1-\theta) \\ &= 2/\theta - 10/(1-\theta)\end{aligned}$$

1.7 Setting Gradient to Zero (Finding MAP)

$$\begin{aligned}2/\theta - 10/(1-\theta) &= 0 \\ 2/\theta &= 10/(1-\theta) \\ 2(1-\theta) &= 10\theta \\ 2 - 2\theta &= 10\theta \\ 2 &= 12\theta \\ \theta_{\text{MAP}} &= 2/12 = 1/6 \approx 0.167\end{aligned}$$

1.8 Posterior Distribution

By Bayes' theorem:

$$\begin{aligned}p(\theta|y) &\propto p(y|\theta) \times p(\theta) \\ \dots &\propto \theta^2 \times (1-\theta)^{10} \times 1 \\ \dots &\propto \theta^{(2+1-1)} \times (1-\theta)^{(10+1-1)} \\ &\propto \theta^2 \times (1-\theta)^{10}\end{aligned}$$

This is $\text{Beta}(\alpha_0 + y, \beta_0 + N - y) = \text{Beta}(3, 11)$

1.9 Posterior Mean and Variance

$$E[\theta|y] = \alpha/(\alpha + \beta) = 3/(3 + 11) = 3/14 \approx 0.214$$

$$\begin{aligned} \text{Var}[\theta|y] &= \alpha\beta/[(\alpha+\beta)^2(\alpha+\beta+1)] \\ &= (3 \times 11)/[(14)^2(15)] \\ &= 33/(196 \times 15) \\ &= 33/2940 \\ &\approx 0.0112 \end{aligned}$$

Part 2: Bayesian Logistic Regression

2.1 Model Setup

Given:

- Prior: $w \sim N(0, \alpha^{-1}I)$ with $\alpha = 1$
- Likelihood: $y_n|w, x_n \sim \text{Bernoulli}(\sigma(w^T x_n))$
- Sigmoid: $\sigma(z) = 1/(1 + e^{(-z)})$

Example data:

- $x_1 = [1.0, 0.5], y_1 = 1$
- $x_2 = [-0.5, 1.0], y_2 = 0$
- $x_3 = [0.3, -0.8], y_3 = 1$

2.2 Log Prior

For $w = [w_1, w_2]^T$:

$$\begin{aligned} \log p(w) &= -\frac{\alpha}{2} \|w\|^2 + \text{constant} \\ &= -\frac{1}{2}(w_1^2 + w_2^2) + \text{constant} \end{aligned}$$

2.3 Log Likelihood

$$\log p(y|w) = \sum_n [y_n \log \sigma(f_n) + (1-y_n) \log(1-\sigma(f_n))]$$

where $f_n = w^T x_n$

For our data:

- $f_1 = w_1 \times 1.0 + w_2 \times 0.5 = w_1 + 0.5w_2$
- $f_2 = w_1 \times (-0.5) + w_2 \times 1.0 = -0.5w_1 + w_2$

- $f_3 = w_1 \times 0.3 + w_2 \times (-0.8) = 0.3w_1 - 0.8w_2$

2.4 Log Joint

$$\begin{aligned} \log p(y, w) &= \log p(w) + \log p(y|w) \\ &= -\frac{1}{2}(w_1^2 + w_2^2) + \\ &\quad [\log \sigma(w_1 + 0.5w_2) + \\ &\quad \log(1 - \sigma(-0.5w_1 + w_2)) + \\ &\quad \log \sigma(0.3w_1 - 0.8w_2)] \end{aligned}$$

2.5 Gradient Calculation

Using the identity: $d/dz \log \sigma(z) = 1 - \sigma(z)$

$$\nabla_w \log p(y, w) = -\alpha w - \sum_n (\sigma(f_n) - y_n) x_n$$

Component-wise:

$$\begin{aligned} \partial/\partial w_1 \log p &= -w_1 - [(\sigma(f_1) - 1) \times 1.0 + (\sigma(f_2) - 0) \times (-0.5) + (\sigma(f_3) - 1) \times 0.3] \\ &= -w_1 - [\sigma(f_1) - 1 - 0.5\sigma(f_2) + 0.3(\sigma(f_3) - 1)] \end{aligned}$$

$$\begin{aligned} \partial/\partial w_2 \log p &= -w_2 - [(\sigma(f_1) - 1) \times 0.5 + (\sigma(f_2) - 0) \times 1.0 + (\sigma(f_3) - 1) \times (-0.8)] \\ &= -w_2 - [0.5(\sigma(f_1) - 1) + \sigma(f_2) - 0.8(\sigma(f_3) - 1)] \end{aligned}$$

2.6 Hessian Calculation

The Hessian is:

$$H = -X^T S X - \alpha I$$

where S is diagonal with $S_{nn} = \sigma(f_n)(1 - \sigma(f_n))$

For our example:

$$X = \begin{bmatrix} 1.0 & 0.5 \\ -0.5 & 1.0 \\ 0.3 & -0.8 \end{bmatrix}$$

$$S = \text{diag}([\sigma(f_1)(1 - \sigma(f_1)), \sigma(f_2)(1 - \sigma(f_2)), \sigma(f_3)(1 - \sigma(f_3))])$$

2.7 MAP Optimization

To find w_{MAP} , solve $\nabla \log p(y, w) = 0$

This requires iterative optimization. Suppose we find: $w_{\text{MAP}} \approx [0.8, -0.3]^T$

2.8 Laplace Approximation

At w_{MAP} :

1. Evaluate Hessian H
2. Posterior covariance: $S = -H^{-1}$
3. Posterior: $q(w) = N(w|w_{\text{MAP}}, S)$

Example calculation (simplified):

$$H \approx -\begin{bmatrix} 2.5 & -0.8 \\ -0.8 & 2.1 \end{bmatrix} - \begin{bmatrix} 1 & 0 \\ 0 & 1 \end{bmatrix} = \begin{bmatrix} -3.5 & 0.8 \\ 0.8 & -3.1 \end{bmatrix}$$

$$S = -H^{-1} \approx \begin{bmatrix} 0.32 & 0.08 \\ 0.08 & 0.31 \end{bmatrix}$$

2.9 Posterior Predictive for New Point $x^* = [0.7, -0.2]$

Step 1: Distribution of $f^* = w^T x^*$

$$\begin{aligned} E[f^*] &= w_{\text{MAP}}^T x^* = 0.8 \times 0.7 + (-0.3) \times (-0.2) = 0.56 + 0.06 = 0.62 \\ \text{Var}[f^*] &= x^{*T} S x^* = \begin{bmatrix} 0.7 & -0.2 \end{bmatrix} \times \begin{bmatrix} 0.32 & 0.08 \\ 0.08 & 0.31 \end{bmatrix} \times \begin{bmatrix} 0.7 \\ -0.2 \end{bmatrix} \\ &= \begin{bmatrix} 0.7 & -0.2 \end{bmatrix} \times \begin{bmatrix} 0.208 \\ 0.006 \end{bmatrix} \\ &= 0.1456 + 0.0012 = 0.147 \end{aligned}$$

Step 2: Plugin Approximation

$$p(y^*=1|x^*, y) \approx \sigma(E[f^*]) = \sigma(0.62) = 1/(1 + e^{(-0.62)}) \approx 0.650$$

Step 3: Probit Approximation

$$\begin{aligned}
p(y^*=1|x^*, y) &\approx \Phi(E[f^*]/\sqrt{\text{Var}[f^*] + \pi/8}) \\
&= \Phi(0.62/\sqrt{0.147 + 0.393}) \\
&= \Phi(0.62/\sqrt{0.540}) \\
&= \Phi(0.844) \\
&\approx 0.801
\end{aligned}$$

Step 4: Monte Carlo (conceptual)

1. Sample $f^*_1, f^*_2, \dots, f^*_{1000} \sim N(0.62, 0.147)$
2. Compute $\sigma(f^*_i)$ for each sample
3. $p(y^*=1|x^*, y) \approx (1/1000) \sum_i \sigma(f^*_i)$

Example samples:

- $f^*_1 = 0.62 + 0.383 \times \varepsilon_1$ ($\varepsilon_1 \sim N(0,1)$)
- $\sigma(f^*_1) = \sigma(0.62 + 0.383 \times \varepsilon_1)$
- Average over all samples ≈ 0.745

Summary of Predictions

For $x^* = [0.7, -0.2]$:

- Plugin: 0.650 (ignores uncertainty)
- Probit: 0.801 (accounts for uncertainty)
- Monte Carlo: 0.745 (most accurate)

The differences show how accounting for parameter uncertainty affects predictions!