

Example Class 4

LIU Chen

Department of Statistics and Actuarial Science,
The University of Hong Kong

October 11, 2022

Bayesian approach

- X : observable random variate with probability function $f(x | \theta)$
- θ : unobservable random variate with a specified prior probability function $\pi(\theta)$:

$$\int_{\Theta} \pi(\theta) d\theta = 1 \quad (\text{continuous } \theta), \quad \text{or}$$

$$\sum_{\theta \in \Theta} \pi(\theta) = 1 \quad (\text{discrete } \theta)$$

Posterior Probability Function

The posterior probability function of θ given the observed data \mathbf{x} is defined to be the conditional probability function of θ given $\mathbf{X} = \mathbf{x}$, that is

$$\begin{aligned}\pi(\theta | \mathbf{x}) &= \frac{f(\mathbf{x} | \theta)\pi(\theta)}{\int_{\Theta} f(\mathbf{x} | \theta')\pi(\theta') d\theta'} \\ &\propto f(\mathbf{x} | \theta)\pi(\theta)\end{aligned}$$

Bayesian Approach: **Expected Posterior Loss**

Let the prior $\pi(\theta)$ be given for $\theta \in \Theta$. Consider a decision problem with loss function $L(\theta, a)$ for $\theta \in \Theta$ and action $a \in \mathcal{A}$ (action space). Definition. The expected posterior loss given data \mathbf{x} , incurred by taking action a , is

$$\mathbb{E}[L(\theta, a) \mid \mathbf{x}] = \int_{\Theta} L(\theta, a) \pi(\theta \mid \mathbf{x}) d\theta$$

Definition.

A Bayesian decision is to take an action $a \in \mathcal{A}$ which minimises the expected posterior loss $\mathbb{E}[L(\theta, a) \mid \mathbf{x}]$

Writing $f(\mathbf{x}) = \int_{\Theta} \pi(\theta') f(\mathbf{x} \mid \theta') d\theta'$, we have

$$\begin{aligned}\mathbb{E}[L(\theta, a) \mid \mathbf{x}] &= \int_{\Theta} L(\theta, a) \frac{f(\mathbf{x} \mid \theta) \pi(\theta)}{f(\mathbf{x})} d\theta \\ &= \frac{1}{f(\mathbf{x})} \int_{\Theta} L(\theta, a) f(\mathbf{x} \mid \theta) \pi(\theta) d\theta\end{aligned}$$

Thus, minimising $\mathbb{E}[L(\theta, a) \mid \mathbf{x}]$ w.r.t. $a \in \mathcal{A}$ is equivalent to minimising $\int_{\Theta} L(\theta, a) f(\mathbf{x} \mid \theta) \pi(\theta) d\theta$ w.r.t. $a \in \mathcal{A}$

Point estimation of $\boldsymbol{\theta} \in \mathbb{R}^k$

Consider two examples of loss function L :

- ① $L(\boldsymbol{\theta}, \mathbf{a}) = \|\boldsymbol{\theta} - \mathbf{a}\|_2^2$ ($\mathbf{a} \in \mathbb{R}^k$) is minimized by

$$\mathbf{a} = \mathbb{E}[\boldsymbol{\theta} \mid \mathbf{x}] = \begin{bmatrix} \mathbb{E}[\theta_1 \mid \mathbf{x}] \\ \vdots \\ \mathbb{E}[\theta_k \mid \mathbf{x}] \end{bmatrix} = \text{posterior mean of } \boldsymbol{\theta}$$

- ② $L(\boldsymbol{\theta}, \mathbf{a}) = \|\boldsymbol{\theta} - \mathbf{a}\|_1$ ($\mathbf{a} \in \mathbb{R}^k$) is minimized by \mathbf{a} satisfying $\mathbb{P}(\theta_i \leq a_i \mid \mathbf{x}) = 1/2$, $i = 1, \dots, k$. Thus, the Bayesian decision is to set a_i to be the posterior median of θ_i

Hypothesis testing about θ

Consider testing:

$$H_0 : \theta \in \Theta_0 \quad \text{vs} \quad H_1 : \theta \in \Theta_1 \equiv \Theta \setminus \Theta_0$$

Action space: $\mathcal{A} = \{a_0 (\leftrightarrow \text{accept } H_0), a_1 (\leftrightarrow \text{reject } H_0)\}$

Define, for some $L_0, L_1 > 0$, the **loss function**

$$L(\theta, a_0) = L_1 \mathbf{1}\{\theta \in \Theta_1\}, \quad L(\theta, a_1) = L_0 \mathbf{1}\{\theta \in \Theta_0\}$$

Given \mathbf{x} , the **expected posterior loss**

$$\mathbb{E}[L(\theta, a_j) \mid \mathbf{x}] = L_{1-j} \mathbb{E}[\mathbf{1}\{\theta \in \Theta_{1-j}\} \mid \mathbf{x}] = L_{1-j} \mathbb{P}(\theta \in \Theta_{1-j} \mid \mathbf{x})$$

The Bayesian decision is to choose a_j such that $L_{1-j} \mathbb{P}(\theta \in \Theta_{1-j} \mid \mathbf{x})$ is minimised, or equivalently, that

$$\frac{L_{1-j}}{L_j} < \frac{\mathbb{P}(\theta \in \Theta_j \mid \mathbf{x})}{1 - \mathbb{P}(\theta \in \Theta_j \mid \mathbf{x})}$$

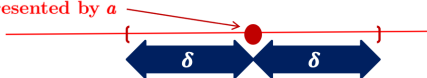
i.e.

$$\text{reject } H_0 \text{ if } \mathbb{P}(H_1 \mid \mathbf{x}) = \mathbb{P}(\theta \in \Theta_1 \mid \mathbf{x}) = \int_{\Theta_1} \pi(\theta \mid \mathbf{x}) d\theta > \frac{L_0}{L_0 + L_1}$$

Interval estimation of θ I

(a) Fix length = 2δ \rightarrow maximise posterior coverage probability

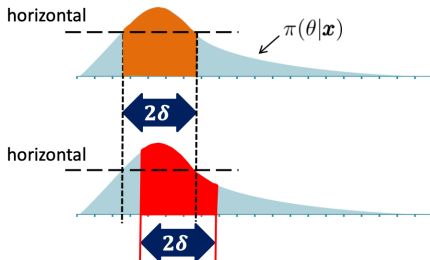
Any candidate can be represented by a



Maximise: $\underbrace{\int_{a-\delta}^{a+\delta} \pi(\theta|\mathbf{x}) d\theta}_{\text{posterior coverage probability of } [a-\delta, a+\delta]} \propto \int_{a-\delta}^{a+\delta} \pi(\theta)f(\mathbf{x}|\theta) d\theta \quad \text{w.r.t. } a$

posterior coverage probability of $[a-\delta, a+\delta]$

Special case: unimodal $\pi(\theta|\mathbf{x})$



optimal solution



$[a-\delta, a+\delta]$ satisfying:

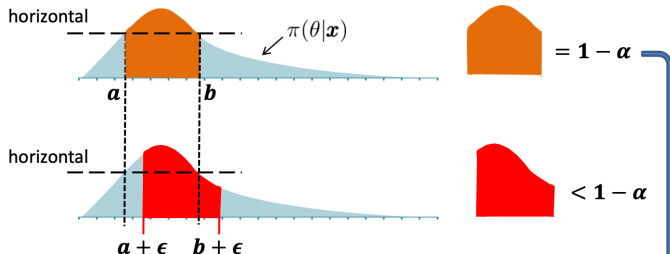
$$\pi(a-\delta|\mathbf{x}) = \pi(a+\delta|\mathbf{x})$$

(assuming $\pi(\theta|\mathbf{x})$ continuous at $\theta = a-\delta, a+\delta$)

Interval estimation of θ II

(b) Desire posterior coverage probability $\geq 1 - \alpha \rightarrow$ minimise length

Special case: unimodal $\pi(\theta|\mathbf{x})$



Optimal solution:

$[a, b]$ satisfying:

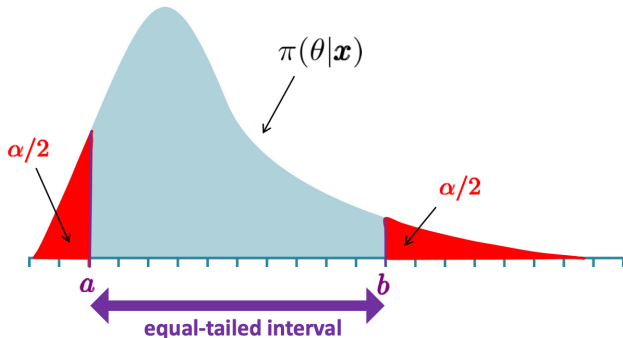
$$\pi(a|\mathbf{x}) = \pi(b|\mathbf{x}) \quad \& \quad \int_a^b \pi(\theta|\mathbf{x}) d\theta = 1 - \alpha$$

(assuming $\pi(\theta|\mathbf{x})$ continuous at $\theta = a, b$)

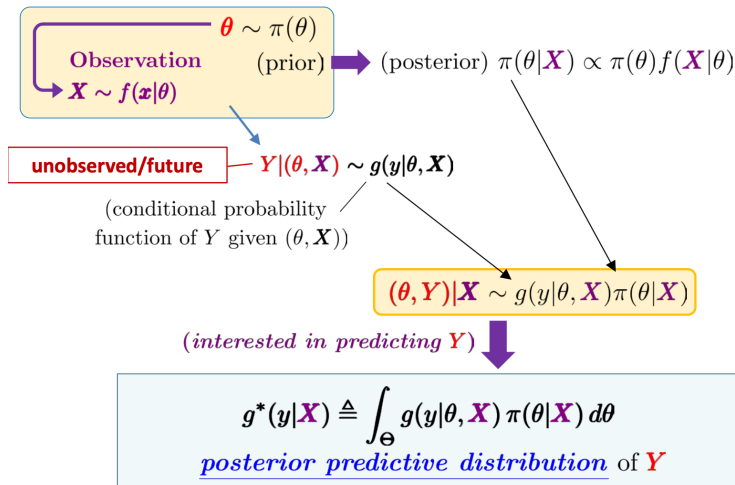
Interval estimation of θ III

(c) Fix posterior coverage probability $= 1 - \alpha$ & require “equal-tailed”

$$\text{i.e. } \left\{ \begin{array}{l} \mathbb{P}(\theta > \mathbf{b} \mid \mathbf{x}) = \int_{\mathbf{b}}^{\infty} \pi(\theta \mid \mathbf{x}) d\theta = \alpha/2 \\ \mathbb{P}(\theta < \mathbf{a} \mid \mathbf{x}) = \int_{-\infty}^{\mathbf{a}} \pi(\theta \mid \mathbf{x}) d\theta = \alpha/2 \end{array} \right.$$



Predictive distribution



Linear Model Basics

- A linear model includes an $n \times 1$ response vector $\mathbf{y} = (y_1, \dots, y_n)$ and an $n \times p$ design matrix (regressors) $\mathbf{X} = [\mathbf{x}_1, \dots, \mathbf{x}_p]$. The relationship between \mathbf{y} and \mathbf{X} has the form,

$$\mathbf{y} = \mathbf{X}\boldsymbol{\beta} + \boldsymbol{\epsilon}, \quad \boldsymbol{\epsilon} \sim N(\mathbf{0}, \sigma^2 \mathbf{I}).$$

- From standard statistical analysis, the classical unbiased and least-square estimates of the regression parameter $\boldsymbol{\beta}$ and σ^2 are

$$\hat{\boldsymbol{\beta}} = (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top \mathbf{y}; \quad \hat{\sigma}^2 = \frac{1}{n - p} (\mathbf{y} - \mathbf{X}\hat{\boldsymbol{\beta}})^\top (\mathbf{y} - \mathbf{X}\hat{\boldsymbol{\beta}})$$

- The predicted value of y_0 given \mathbf{X}_0 is

$$\hat{y}_0 = \mathbf{X}_0 \hat{\boldsymbol{\beta}} = \mathbf{X}_0 (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top \mathbf{y}$$

Hypothesis testing under linear model

$$H_0 : \beta_j = c \quad \text{versus} \quad H_1 : \beta_j \neq c$$

- **Wald t-test**

$$T = \frac{\hat{\beta}_j - c}{\sqrt{\text{var}(\hat{\beta}_j)}} \underset{\text{under } H_0}{\sim} t(n - p)$$

- **Likelihood ratio test**

$$LR = -2 \ln \left(\frac{L(\hat{\beta}_{H_0})}{L(\hat{\beta})} \right) \underset{\text{under } H_0}{\sim} \chi^2(1)$$

Bayesian linear model

$$\mathbf{y} = \mathbf{X}\boldsymbol{\beta} + \boldsymbol{\epsilon}, \quad \boldsymbol{\epsilon} \sim N(\mathbf{0}, \sigma^2 \mathbf{I}).$$

$$\mathbf{y} | \mathbf{X}, \boldsymbol{\beta}, \sigma^2 \sim N(\mathbf{X}\boldsymbol{\beta}, \sigma^2 \mathbf{I})$$

$$\beta_j \sim N(0, \sigma_0^2), \quad j = 1, \dots, p$$

$$\sigma^2 \sim \text{InvGamma}(\xi, \xi)$$

Bayesian linear model (cont'd)

$$\begin{aligned}f(\boldsymbol{\beta}, \sigma^2 | \mathbf{X}, \mathbf{y}) &\propto f(\mathbf{y} | \mathbf{X}, \boldsymbol{\beta}, \sigma^2) \pi(\boldsymbol{\beta}) \pi(\sigma^2) \\&\propto (\sigma^2)^{-\frac{n}{2}} \exp\left(-\frac{(\mathbf{y} - \mathbf{X}\boldsymbol{\beta})^\top (\mathbf{y} - \mathbf{X}\boldsymbol{\beta})}{2\sigma^2}\right) \exp\left(-\frac{\boldsymbol{\beta}^\top \boldsymbol{\beta}}{2\sigma_0^2}\right) \\&\quad (\sigma^2)^{-\xi-1} \exp\left(-\frac{\xi}{\sigma^2}\right) \\f(\sigma^2 | \mathbf{X}, \mathbf{y}, \boldsymbol{\beta}) &\propto (\sigma^2)^{-\xi-n/2-1} \exp\left(-\frac{\xi + \frac{1}{2}(\mathbf{y} - \mathbf{X}\boldsymbol{\beta})^\top (\mathbf{y} - \mathbf{X}\boldsymbol{\beta})}{\sigma^2}\right) \\f(\boldsymbol{\beta} | \mathbf{X}, \mathbf{y}, \sigma^2) &\propto \exp\left(-\frac{\boldsymbol{\beta}^\top \boldsymbol{\beta}}{2\sigma_0^2} - \frac{\boldsymbol{\beta}^\top \mathbf{X}^\top \mathbf{X} \boldsymbol{\beta} - 2\mathbf{y}^\top \mathbf{X} \boldsymbol{\beta}}{2\sigma^2}\right)\end{aligned}$$

Bayesian linear model (cont'd)

Sometimes prior distributions might not be valid distributions. For example,

$$\pi(\beta) \propto 1 \longleftrightarrow N(0, \sigma^2), \sigma^2 \rightarrow \infty$$

$$\pi(\sigma^2) \propto 1/\sigma^2 \longleftrightarrow \text{InvGamma}(\xi, \xi), \xi \rightarrow 0$$

Even if the priors are improper, as long as the resulting posterior distribution are valid we can still conduct legitimate statistical inference on them

Predicting from linear models

- New covariates $\tilde{\mathbf{X}}$, and wish to predict the corresponding outcome \tilde{y} .
- If β and σ^2 are known, then $\tilde{y} \sim N(\tilde{\mathbf{X}}\beta, \sigma^2)$
- When parameters are unknown,

$$f(\tilde{y}|\mathbf{y}) = \int f(\tilde{y}|\beta, \sigma^2)f(\beta, \sigma^2|\mathbf{y})d\beta d\sigma^2$$

- For each posterior draw of $(\beta_{(i)}, \sigma_{(i)}^2)_{i=1}^M$, draw \tilde{y}_i from $N(\tilde{\mathbf{X}}\beta_{(i)}, \sigma_{(i)}^2)$. The resulting samples $(\tilde{y}_{(i)})_{i=1}^M$ represent the predictive distribution.

Bayesian Linear Regression (Example)

- Dataset: “data.csv”
- Linear regression model:

$$y = \beta_0 + \beta_1 x_1 + \beta_2 x_2 + \epsilon,$$

where $\epsilon \sim N(0, \sigma^2)$.

- Frequentist MLE:

$$\hat{\beta} = (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{y} \approx (1.085, 0.384, 2.445)^T,$$

$$\hat{\sigma}^2 = \frac{1}{n} \sum_{i=1}^n (y_i - \mathbf{x}_i^T \hat{\beta})^2 \approx 0.326.$$

Bayesian Linear Regression (Example)

- Bayesian paradigm: $\beta \sim N(\mu_0, \sigma_0^2 \mathbf{I}_3)$ and $\sigma^2 \sim IG(\xi_0, \xi_0)$ where \mathbf{I}_n denote the n -dimensional identity matrix, $\mu_0 = \mathbf{0}$, $\sigma_0 = 10$, $\xi_0 = 0.1$
- Gibbs sampler:

$$p(\sigma^2 | \beta, \mathbf{X}, \mathbf{y}) \propto \left(\frac{1}{\sigma^2} \right)^{n/2 + \xi_0 + 1} e^{-\frac{1}{\sigma^2} [(\mathbf{y} - \mathbf{X}\beta)^T (\mathbf{y} - \mathbf{X}\beta) / 2 + \xi_0]}$$
$$\sim IG(n/2 + \xi_0, (\mathbf{y} - \mathbf{X}\beta)^T (\mathbf{y} - \mathbf{X}\beta) / 2 + \xi_0)$$

$$p(\beta | \mu_0, \sigma_0, \sigma, \mathbf{X}, \mathbf{y}) \propto e^{-\frac{1}{2}(\beta - \Sigma\eta)^T \Sigma^{-1}(\beta - \Sigma\eta)}$$
$$\sim N(\Sigma\eta, \Sigma),$$

$$\eta = \mu_0 / \sigma_0^2 + \mathbf{X}^T \mathbf{y} / \sigma^2,$$

$$\Sigma = (\mathbf{I}_3 / \sigma_0^2 + \mathbf{X}^T \mathbf{X} / \sigma^2)^{-1}$$

Bayesian Linear Regression (Example)

	β_0	β_1	β_2	σ^2
mean	1.085	0.382	2.445	0.334
variance	0.00733	0.02331	0.00184	0.00116

Table 1: Posterior means and variances of unknown parameters.

Bayesian Linear Regression (Example, with an interaction)

- MLE results. $y = \beta_0 + \beta_1 x_1 + \beta_2 x_2 + \beta_3 x_1 x_2 + \epsilon$, $\epsilon \sim N(0, \sigma^2)$
- $\hat{\beta}_0 = 1.007$, $\hat{\beta}_1 = 0.518$, $\hat{\beta}_2 = 2.037$, $\hat{\beta}_3 = 0.754$, $\hat{\sigma}^2 = 0.295$
- Posterior means and variances of unknown parameters are listed below.

	β_0	β_1	β_2	β_3	σ^2
mean	1.007	0.518	2.035	0.756	0.305
variance	0.00693	0.02193	0.00966	0.02741	0.00098

Table 2: Posterior means and variances of unknown parameters.

Bayesian Linear Regression (Example, with an interaction)

- Two-sided hypothesis test:

$$H_0 : \beta_3 = 0 \quad \text{versus} \quad H_1 : \beta_3 \neq 0.$$

- Likelihood ratio test (LRT):

$$T_{\text{LRT}} = (n-4)(\Lambda^{-2/n} - 1) = \frac{\mathbf{y}^T (\mathbf{H} - \mathbf{H}_{-4}) \mathbf{y}}{\mathbf{y}^T (\mathbf{I}_n - \mathbf{H}) \mathbf{y} / (n-4)} \sim F(1, n-4), \quad (1)$$

where \mathbf{H} and \mathbf{H}_{-4} denote the hat matrices with and without the interaction term.

- Wald test:

$$T_{\text{Wald}} = \frac{\hat{\beta}_3^2}{\hat{\sigma}^2 \mathbf{e}_4^T (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{e}_4}, \quad (2)$$

where $\mathbf{e}_4 = (0, 0, 0, 1)^T$.

Bayesian Linear Regression (Example, with an interaction)

$$T_{\text{Wald}} = \frac{\hat{\beta}_3}{\hat{\sigma} \sqrt{(\mathbf{X}^\top \mathbf{X})_{3,3}^{-1}}} \sim t(200 - 4) = t(196)$$

$$p\text{-value}_{\text{Wald}} = 1.002 \times 10^{-5}$$

$$-2 \ln \left(\frac{L(\hat{\beta}_{H_0})}{L(\hat{\beta})} \right) \underset{\text{under } H_0}{=} -2 \ln \frac{L(\hat{\beta}_{H_0}, \hat{\sigma}_{H_0})}{L(\hat{\beta}, \hat{\sigma})} \sim \chi^2(1)$$

$$L(\boldsymbol{\beta}, \sigma) = (2\pi)^{-n/2} \sigma^{-n} \exp \left(-\frac{(\mathbf{y} - \mathbf{X}\boldsymbol{\beta})^\top (\mathbf{y} - \mathbf{X}\boldsymbol{\beta})}{2\sigma^2} \right)$$

$$p\text{-value}_{\text{LRT}} = 7.926 \times 10^{-6}$$

Bayesian Linear Regression (Example, with an interaction)

- The p -value under the two-sided hypothesis test is

$$p\text{-value}_2 = 2 - 2\Phi(|Z|) = 2[1 - \max\{\Phi(Z), \Phi(-Z)\}].$$

We define the two-sided posterior probability (PoP_2) as

$$\text{PoP}_2 = 2[1 - \max\{\Pr(\beta_3 > 0|y), \Pr(\beta_3 < 0|y)\}].$$

- $\text{PoP}_2 = 0$

Coin Tossing Problem

- We consider an experiment in which a coin was tossed 12 times, with 9 heads and 3 tails observed.
- Let θ be the probability of observing a head for a toss of the coin
- Suppose that we conduct a one-sided hypothesis test,

$$H_0 : \theta \leq 0.5 \quad \text{versus} \quad H_1 : \theta > 0.5.$$

- Based on the binomial or negative binomial likelihood, the frequentist hypothesis test yields conflicting results under the significance level of $\alpha = 0.05$: The null hypothesis is accepted under the binomial distribution, but it is rejected under the negative binomial distribution.

Coin Tossing Problem (Binomial)

- For $Y \sim \text{Bin}(n, \theta)$, we have

$$\begin{aligned}\Pr(Y \geq y \mid \theta) &= \sum_{k=y}^n \binom{n}{k} \theta^k (1 - \theta)^{n-k} \\ &= \frac{\Gamma(n+1)}{\Gamma(n-y+1)\Gamma(y)} \int_0^\theta t^{y-1} (1-t)^{n-y} dt = I_\theta(y, n-y+1)\end{aligned}$$

where $I_x(a, b)$ is the regularized incomplete beta function defined as

$$\begin{aligned}I_x(a, b) &= \frac{B(x; a, b)}{B(a, b)} \\ B(x; a, b) &= \int_0^x t^{a-1} (1-t)^{b-1} dt \\ B(a, b) &= \int_0^1 t^{a-1} (1-t)^{b-1} dt\end{aligned}$$

$$p\text{-value}_{\text{Bin}} = I_{0.5}(y, n-y+1)$$

Coin Tossing Problem (Negative-Binomial)

- For $Y \sim \text{NB}(r, \theta)$, we have

$$\begin{aligned}\Pr(Y \geq y \mid r, \theta) &= \sum_{k=y}^{\infty} \binom{k+r-1}{k} \theta^k (1-\theta)^r \\ &= \frac{\Gamma(y+r)}{\Gamma(r)\Gamma(y)} \int_0^{\theta} t^{y-1} (1-t)^{r-1} dt \\ &= l_{\theta}(y, r)\end{aligned}$$



$$p\text{-value}_{\text{NB}} = l_{0.5}(y, r) = l_{0.5}(y, n - y)$$

Coin Tossing Problem (Example)

- Under Bayesian paradigm, we assume a symmetric beta prior ($\alpha = \beta$) for θ , i.e., $\theta \sim \text{Beta}(\alpha, \beta)$.
- For $\alpha = 2, 1, 0.5, 0.1, 0.01, 0.001, 0.0001, 0.00001$, calculate $P(H_0|n = 12, y = 3)$ and comment on your findings in comparison with the p -values obtained under the binomial and negative binomial likelihood.
- Plot (i) $p\text{-value}_B$ against $P(H_0|n, y)$ (ii) $p\text{-value}_{NB}$ against $P(H_0|n, y)$ as sample size n increases while fixing $y/n = 0.75$. For the Bayesian paradigm, we set $\alpha = \beta = 0.01$.

Coin Tossing Problem (Example)

$\alpha = \beta$	$P(H_0 n, y)$
2	0.0592346
1	0.0461426
0.5	0.0394446
0.1	0.0340598
0.01	0.0328493
0.001	0.0327283
0.0001	0.0327162
0.00001	0.0327150

Coin Tossing Problem (Example)

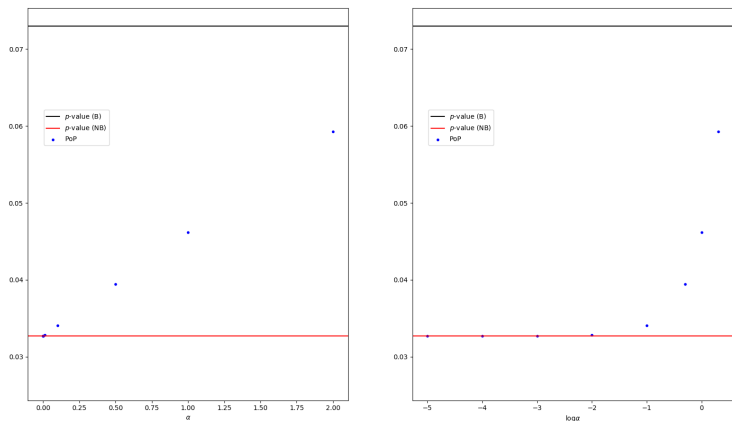


Figure 1: Plots of $P(H_0|n, y)$ against $\alpha=(2, 1, 0.5, 0.1, 0.01, 0.001, 0.0001, 0.00001)$.

Coin Tossing Problem (Example)

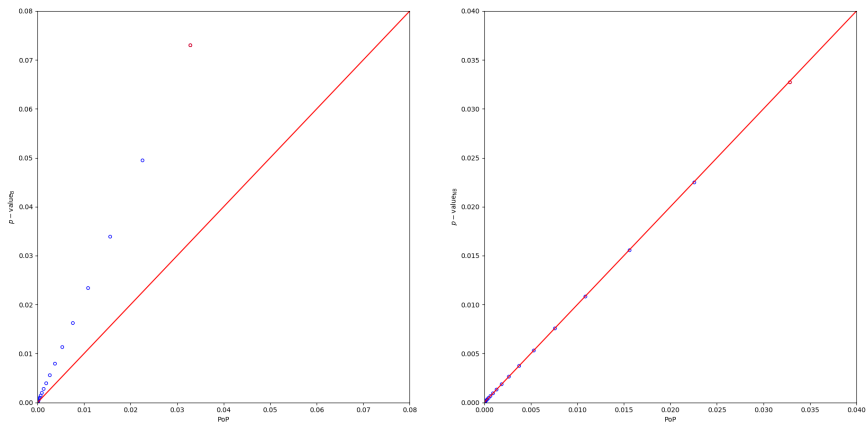


Figure 2: Plots of $P(H_0|n, y)$ against $p\text{-value}_B$ (left) and $p\text{-value}_{NB}$ (right)