

VE414 Lecture 5

Jing Liu

UM-SJTU Joint Institute

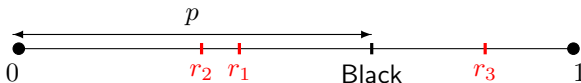
September 19, 2019

- Recall the structure of our Bayesian model so far is very simple, i.e.

$$\begin{aligned} Y &\sim f_Y \\ X | Y &\sim f_{X|Y} \\ Y | X = x &\sim f_{Y|X=x} \end{aligned}$$

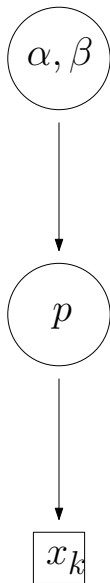
where $f_{Y|X=x} \propto f_{X=x|Y} \cdot f_Y$ according to Bayes theorem.

- For example, the uncertainty in Bayes' original problem



that is, our lack of knowledge on where the black ball, is modelled by,

$$\begin{aligned} P &\sim \text{Beta}(\alpha, \beta) \\ X_k | P &\sim \text{Binomial}(k, p) \\ P | X_k = x_k &\sim \text{Beta}(\alpha + x_k, \beta + (k - x_k)) \end{aligned}$$



- Clearly not all uncertainty can be properly modelled by a simple structure.
- Imagine a study on the effectiveness of cardiac treatments in a hospital,

$$y_{71}/k_{71} = 10/14$$

patients survived. Historically, similar data exist for other hospitals in the city

$$y_j/k_j \quad \text{where } j = 1, 2, \dots, 70.$$

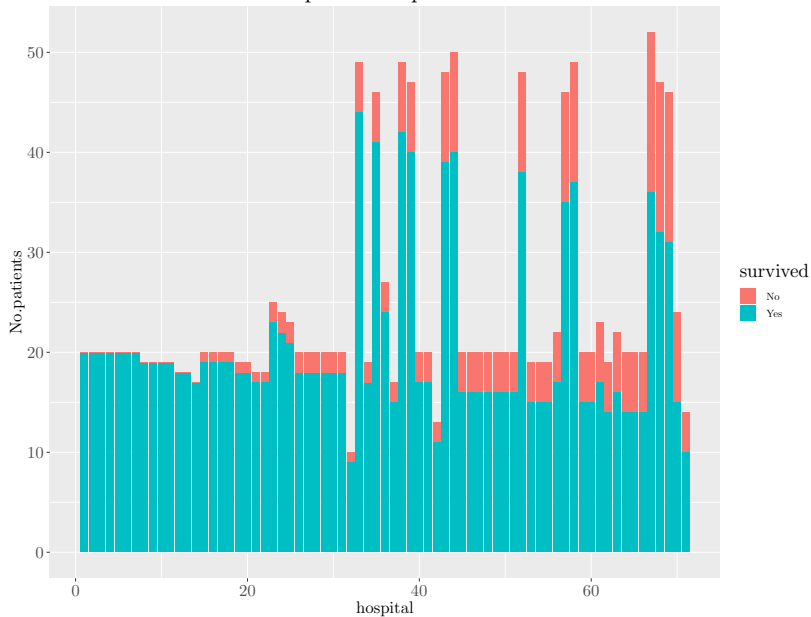
- It is clearly **not** ideal to collapse all data into a single value of y and k

$$y = \sum_{i=1}^{71} y_i; \quad k = \sum_{i=1}^{71} k_i$$

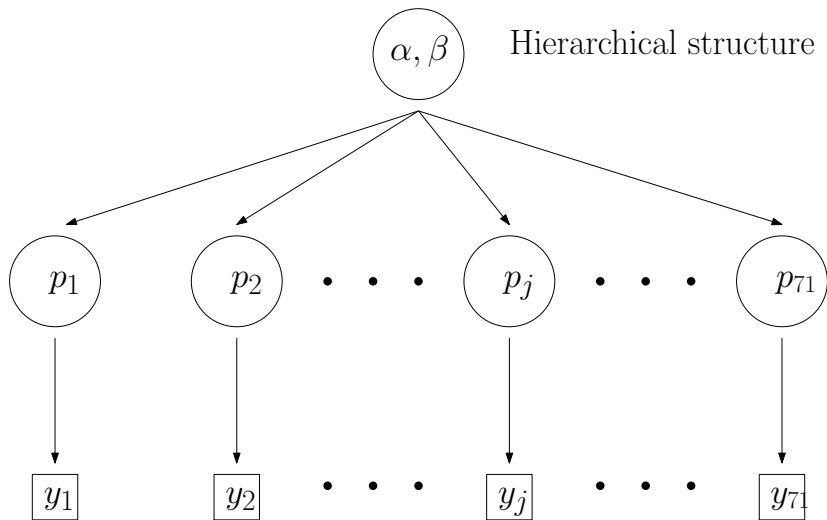
since it is unlikely hospitals have the same equipments, let alone the doctors.

- It is more realistic to assume that the patents in hospital j having their own survival probability p_j , which might be related to each other in some sense.

Bar plot of hospital data



- If we assume p_j 's follow a common distribution, then we have the following



- It is reasonable to assume the following

$$P_j \sim \text{Beta}(\alpha, \beta)$$

$$Y_j \mid P_j \sim \text{Binomial}(k_j, p_j)$$

Q: How to specify α and β , and how to obtain the posterior of P_{71} ?

- We have discussed collapsing historical data, $j = 1, 2, \dots, 70$, with current data $j = 71$ is **not** ideal, that is the following estimation is not ideal

$$P \sim \text{Beta}(1, 1)$$

$$Y \mid P \sim \text{Binomial}(k = 1739, p)$$

$$P \mid X_{1739} = 1472 \sim \text{Beta}(1 + 1472, 1 + (1739 - 1472))$$

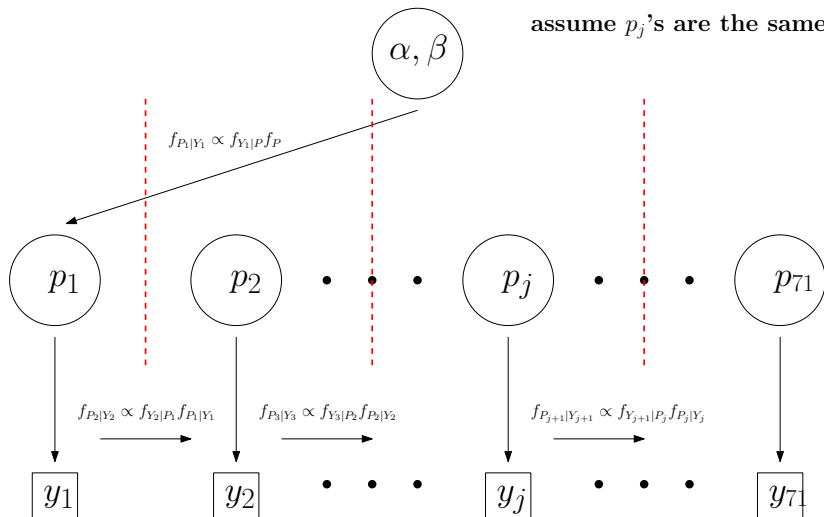
which treat it as Bayes' original problem, and should **not** be taken seriously.

- Using historical data sequentially as before to obtain a prior is just as bad.

uniform prior $\xrightarrow{y_1}$ posterior as prior $\xrightarrow{y_2} \dots \xrightarrow{y_{71}}$ posterior

- It is just as bad since it is equivalent to collapsing the data, in which we

assume p_j 's are the same



Q: How can we incorporate the information in the historical data into our prior?

- Perhaps one rather simple and nature way in this case is to combine
frequentist and Bayesian approaches

- The historical data, $j = 1, 2, \dots, 70$, can be treated as observed proportions

$$P \sim \text{Beta}(\alpha, \beta)$$

- That is, instead of treating p_j 's as unobserved and having the lay of y_j 's, let

$$p_j = y_j/k_j \quad \text{for } j = 1, 2, \dots, 70$$

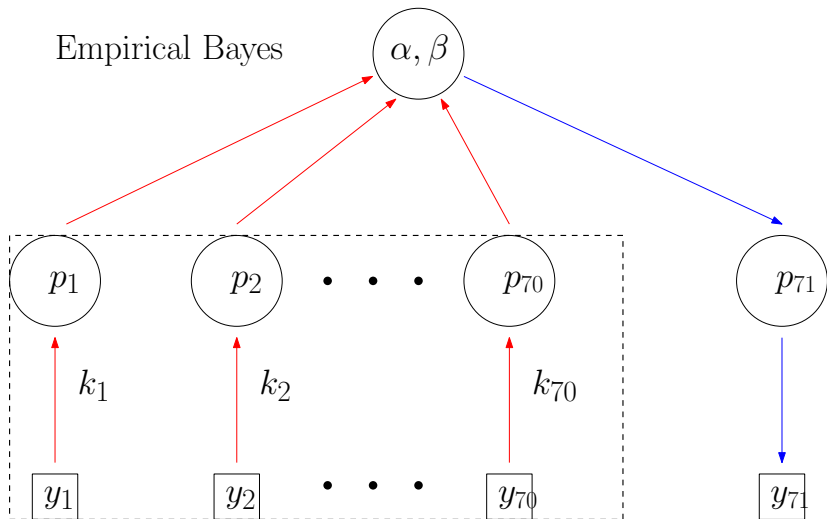
- However, we keep the current data on hospital 71 as it is

$$y_{71} = 10, \quad k_{71} = 14$$

and still in search for the posterior of the random variable P_{71} given the data.

- So we pretty much remove a layer from the structure, and use it separately

Empirical Bayes



- Using the 70 proportions,

$$p_j = y_j/k_j \quad \text{for } j = 1, 2, \dots, 70$$

we could estimate α and β using a frequentist approach, e.g.

$$\tilde{\alpha} = \bar{p} \left(\frac{\bar{p}(1 - \bar{p})}{s_p^2} - 1 \right); \quad \tilde{\beta} = (1 - \bar{p}) \left(\frac{\bar{p}(1 - \bar{p})}{s_p^2} - 1 \right)$$

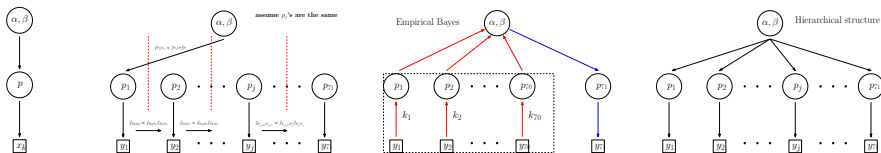
- Here the method of moments (MoM) is used since the MLE is not analytic.
- The prior is set to be $P_{71} \sim \text{Beta}(\tilde{\alpha}, \tilde{\beta})$, from which we obtain

$$P_{71} \mid Y_{71} = 10 \sim \text{Beta}(\tilde{\alpha} + 10, \tilde{\beta} + 4)$$

which is known as the [pseudo-posterior](#), using Bayes theorem as usual.

- This half frequentist half Bayesian method is known as [empirical Bayes](#).

- To derive a full Bayesian analysis for the hierarchical structure, we have to understand the difference in how data provide information



- In a simple structure, the data x_k , provide the information directly on the distribution of p , however, in a hierarchical structure, the historical data,

$$\mathbf{y}_h = [y_1 \quad y_2 \quad \cdots \quad y_{70}]^T$$

provide no information on p_{71} directly; it is only through α and β , that \mathbf{y}_h is useful to improve our knowledge on p_{71} in a way similar to empirical Bayes.

- Since \mathbf{y}_h improves our understanding on α and β , we have to treat α and β as random as well. The distribution $f_{\alpha, \beta}$ is known as **hyperprior**.

- Therefore the joint posterior distribution

$$f_{\{\mathbf{P}, \alpha, \beta\} | \mathbf{Y} = \mathbf{y}}$$

where $\mathbf{y} = [\mathbf{y}_h \ y_{71}]^T$ denotes all the data, and \mathbf{P} the random vector

$$\mathbf{P} = \begin{bmatrix} P_1 \\ \vdots \\ P_{71} \end{bmatrix}$$

would capture our most up to date understanding of the system according to the Bayesian hierarchical model, in which we have to specify a likelihood

$$\mathcal{L}(\mathbf{P}, \mathbf{y}) = f_{\mathbf{y} | \{\mathbf{P}, \alpha, \beta\}}$$

and a joint prior in terms of a conditional prior and a hyperprior

$$f_{\mathbf{P}, \alpha, \beta} = f_{\mathbf{P} | \{\alpha, \beta\}} \cdot f_{\alpha, \beta}$$

- Recall we have been using a single-parameter version of Bayes theorem so far

$$f_{Y|X} \propto f_{X|Y} f_Y$$

and the continuous version is based on the definition conditional probability

$$\begin{aligned} F_{Y|X}(y | x) &= \lim_{\varepsilon \rightarrow 0^+} \Pr(Y \leq y | X \in (x, x + \varepsilon]) \\ &= \lim_{\varepsilon \rightarrow 0^+} \frac{\Pr(Y \leq y, X \in (x, x + \varepsilon])}{\Pr(x < X \leq x + \varepsilon)} \\ &= \lim_{\varepsilon \rightarrow 0^+} \frac{F_{X,Y}(x + \varepsilon, y) - F_{X,Y}(x, y)}{F_X(x + \varepsilon) - F_X(x)} = \frac{\partial F_{X,Y}(x, y) / \partial x}{f_X(x)} \\ \implies f_{Y|X}(y | x) &= \frac{\partial}{\partial y} \frac{\partial F_{X,Y}(x, y) / \partial x}{f_X(x)} \\ &= \frac{1}{f_X(x)} \frac{\partial^2 F_{X,Y}(x, y)}{\partial y \partial x} \\ &= \frac{f_{X,Y}(x, y)}{f_X(x)} \propto f_{X|Y}(x | y) \cdot f_Y(y) \end{aligned}$$

- The same approach can be used to derive a multi-parameter version, e.g.

$$\begin{aligned}
 F_{\{Y,Z\}|X} &= \lim_{\varepsilon \rightarrow 0^+} \Pr(Y \leq y, Z \leq z \mid X \in (x, x + \varepsilon]) = \frac{\partial F_{X,Y,Z}/\partial x}{f_X} \\
 \Rightarrow f_{\{Y,Z\}|X} &= \frac{\partial^2}{\partial y \partial z} \frac{\partial F_{X,Y,Z}/\partial x}{f_X} = \frac{f_{X,Y,Z}}{f_X} = \frac{f_{X|\{Y,Z\}} \cdot f_{Y,Z}}{f_X} \\
 &\propto f_{X|\{Y,Z\}} \cdot f_{Y,Z} \propto f_{X|\{Y,Z\}} \cdot f_{Y|Z} \cdot f_Z
 \end{aligned}$$

- In this case, with the assumption of independence, the joint posterior is

$$\begin{aligned}
 f_{\{\mathbf{P}, \alpha, \beta\}|\mathbf{Y}} &\propto f_{\mathbf{Y}|\{\mathbf{P}, \alpha, \beta\}} \cdot \textcolor{red}{f_{\mathbf{P}|\{\alpha, \beta\}}} \cdot \textcolor{blue}{f_{\alpha, \beta}} \\
 &= \prod_{j=1}^{71} \binom{k_j}{y_j} p_j^{y_j} (1 - p_j)^{k_j - y_j} \prod_{j=1}^{71} \frac{\Gamma(\alpha + \beta)}{\Gamma(\alpha)\Gamma(\beta)} \textcolor{red}{p_j^{\alpha-1} (1 - p_j)^{\beta-1}} \textcolor{blue}{f_{\alpha, \beta}} \\
 &= f_{\alpha, \beta} \prod_{j=1}^{71} \binom{k_j}{y_j} \frac{\Gamma(\alpha + \beta)}{\Gamma(\alpha)\Gamma(\beta)} p_j^{y_j + \alpha - 1} (1 - p_j)^{k_j - y_j + \beta - 1}
 \end{aligned}$$

- The conditional posterior of \mathbf{P} given α and β is

$$\begin{aligned}
 f_{\mathbf{P}|\{\alpha,\beta,\mathbf{Y}\}} &= \frac{f_{\{\mathbf{P},\alpha,\beta\}|\mathbf{Y}}}{f_{\{\alpha,\beta\}|\mathbf{Y}}} \propto f_{\{\mathbf{P},\alpha,\beta\}|\mathbf{Y}} \\
 &\propto \prod_{j=1}^{71} p_j^{y_j+\alpha-1} (1-p_j)^{k_j-y_j+\beta-1} \\
 &= \prod_{j=1}^{71} \underbrace{\frac{\Gamma(\alpha+\beta+k_j)}{\Gamma(\alpha+y_j)\Gamma(\beta+k_j-y_j)}}_{\text{normalisation constant } c} p_j^{y_j+\alpha-1} (1-p_j)^{k_j-y_j+\beta-1}
 \end{aligned}$$

where the normalisation constant c is due to the fact that components of \mathbf{P} have independent posterior densities of the form

$$p_j^{\alpha_j^*-1} (1-p_j)^{\beta_j^*-1}$$

thus c is found by finding the normalisation constant of a beta distribution.

- The **marginal posterior** of α and β can be obtained since the joint posterior

$$f_{\{\mathbf{P}, \alpha, \beta\}|\mathbf{Y}} \propto f_{\alpha, \beta} \prod_{j=1}^{71} \binom{k_j}{y_j} \frac{\Gamma(\alpha + \beta)}{\Gamma(\alpha)\Gamma(\beta)} p_j^{y_j + \alpha - 1} (1 - p_j)^{k_j - y_j + \beta - 1}$$

and the conditional posterior of \mathbf{P} given α and β

$$f_{\mathbf{P}|\{\alpha, \beta, \mathbf{Y}\}} = \prod_{j=1}^{71} \frac{\Gamma(\alpha + \beta + k_j)}{\Gamma(\alpha + y_j)\Gamma(\beta + k_j - y_j)} p_j^{y_j + \alpha - 1} (1 - p_j)^{k_j - y_j + \beta - 1}$$

have been worked out, thus the marginal posterior of α and β is given by

$$\begin{aligned} f_{\{\alpha, \beta\}|\mathbf{Y}} &= \frac{f_{\{\mathbf{P}, \alpha, \beta\}|\mathbf{Y}}}{f_{\mathbf{P}|\{\alpha, \beta, \mathbf{Y}\}}} \\ &\propto f_{\alpha, \beta} \prod_{j=1}^{71} \frac{\Gamma(\alpha + \beta)}{\Gamma(\alpha)\Gamma(\beta)} \frac{\Gamma(\alpha + y_j)\Gamma(\beta + k_j - y_j)}{\Gamma(\alpha + \beta + k_j)} \end{aligned}$$