

# MS = BAYESIAN STATISTICS

Bayesian reasoning provides rules of logic for updating prior beliefs in the light of new evidence. It is named after Thomas Bayes, an XVIII<sup>e</sup>-century presbyterian minister who did maths on the side. His work "An Essay towards solving a problem in the doctrine of chances", published posthumously in 1763, contains theorems of conditional probability, that set the grounds for what is known now as Bayes's theorem. The Bayesian framework is now widely used across many fields, ranging from machine learning to epidemics, and medical studies.

For two events A and B, Bayes's theorem states that

$$P(A|B) = \frac{P(B|A)P(A)}{P(B)} \propto P(B|A)P(A)$$

"posterior" belief  
= updated belief  
after new evidence  
is collected

new evidence

"prior" probability  
or  
base rate

Update prior belief after  
having observed new evidence

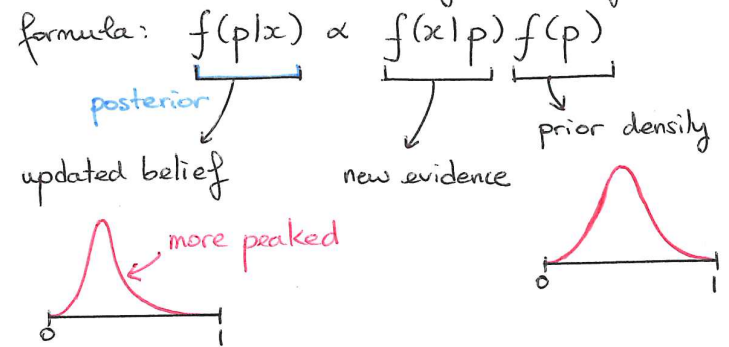
- Ex = A = { presence of antibodies }
- B = { positive test result }

$P(B|A) \neq 1$  since medical tests are imperfect.

Bayes's theorem is a simple but powerful statement about proportions. It provides a way to integrate what

you previously thought with what you have learned and reach a conclusion that incorporates them both, with appropriate weights.

In the formula  $P(A|B) \propto P(B|A)P(A)$ , the prior puts a mass on some probability, and assigns a zero probability to other values. It may be more realistic to spread this belief around this probability, and work instead with the "conditional density" version of Bayes's formula:



The last written formula describes a situation where

- $X|p \sim B(p)$  i.e.  $f(x|p) = p^x(1-p)^{1-x}$   $x \in \{0, 1\}$ .
- $p \sim f(p)$  some prior probability distribution. When nothing is known about  $p$  beforehand, it is common to take a uniform density over  $[0, 1]$ ,

$$p \sim U(0, 1) = B(1, 1)$$

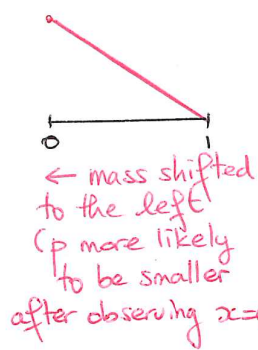
↖ Beta distribution

The posterior is  $f(p|x) = p^x(1-p)^{1-x} = B(1+x, 2-x)$

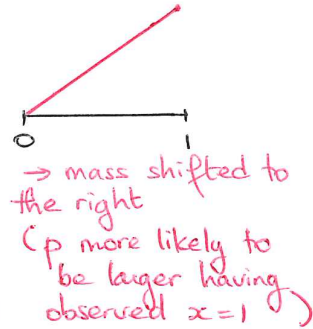
Again, a Beta distribution. We say that Beta is the conjugate prior of the binomial distribution.

In particular,

$x=0$   
 $\sim B(1, 2)$



$x=1$   
 $\sim B(2, 1)$

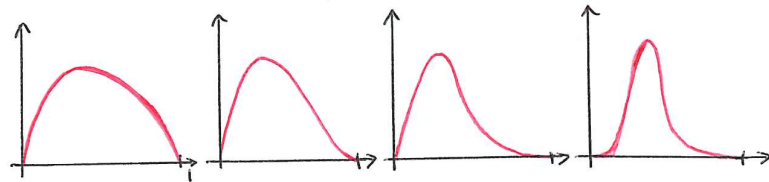


(3)

The  $B(\alpha, \beta)$  distribution has density  $\frac{\Gamma(\alpha+\beta)}{\Gamma(\alpha)\Gamma(\beta)} p^{\alpha-1} (1-p)^{\beta-1}$ ,  
with mean  $\frac{\alpha}{\alpha+\beta}$  and variance  $\frac{\alpha\beta}{(\alpha+\beta)^2(\alpha+\beta+1)}$ .

As  $\alpha$  increases, the distribution is shifted to the left  
 $\beta$  — " — , — " — to the right.  
& more peaked as  $(\alpha+\beta)$  grows.

Ex:



$(\alpha, \beta) = (2, 3) \quad (2, 10) \quad (20, 100) \quad (200, 1000)$

(and uniform with  $(\alpha, \beta) = (1, 1)$ )

Suppose now that we collect  $n$  observations  $\mathcal{L}_n = \{X_1, \dots, X_n\}$ ,  
 $X_1, \dots, X_n$  iid  $B(p)$  :  $P(X_i = 1) = p$ .

→ Consider the conjugate beta prior on  $p \sim B(\alpha, \beta)$ . (4)

The likelihood is  
 $L(X_1, \dots, X_n | p) = \prod_{i=1}^n p^{X_i} (1-p)^{1-X_i}$   
 $= p^{S_n} (1-p)^{n-S_n}$ ,

where  $S_n = \sum_{i=1}^n X_i$ .  
The MLE is  $\hat{p} = \frac{S_n}{n} =$   
proportion of successes.

→ The posterior is proportional to the product of the likelihood with the prior

$p | \mathcal{L}_n \sim p^{S_n} (1-p)^{n-S_n} p^{\alpha-1} (1-p)^{\beta-1}$   
 $= p^{(S_n+\alpha)-1} (1-p)^{(n+\beta-S_n)-1}$   
 $\sim B(\alpha + S_n, n + \beta - S_n)$ .

We retrieve again the conjugate property of the beta distribution.

→ The posterior mean is

$E(p | \mathcal{L}_n) = \frac{\alpha + S_n}{\alpha + \beta + n} = \left( \frac{n}{\alpha + \beta + n} \right) \left[ \frac{S_n}{n} \right] + \left( \frac{\alpha + \beta}{\alpha + \beta + n} \right) \left[ \frac{\alpha}{\alpha + \beta} \right]$

Annotations:   
-  $\frac{S_n}{n}$  is labeled "MLE" with an arrow.  
-  $\frac{\alpha}{\alpha + \beta}$  is labeled "prior mean" with an arrow.  
-  $\frac{n}{\alpha + \beta + n} \rightarrow 1$  as  $n \rightarrow \infty$   
-  $\frac{\alpha + \beta}{\alpha + \beta + n} \rightarrow 0$  as  $n \rightarrow \infty$

It is a general feature of the posterior distribution: centered at a point that represents a compromise between prior info & the data. Also, the data get more weight as  $n$  gets large  $\rightarrow E(p | \mathcal{L}_n) \approx \frac{S_n}{n} = \text{MLE}$ .

5

A uniform prior yields a posterior mean equal to  $\frac{S_n+1}{n+2}$ . Comparing with the MLE  $\frac{S_n}{n}$ , this accounts to augmenting the set of observations  $\mathcal{X}_n = \{X_1, \dots, X_n\}$  with two additional artificial values: one extra success and one extra failure.

↖  $\equiv$  regularization in small sample cases, preventing over-fitting.

Ex:  $\mathcal{X}_n = \{X_1 = X_2 = X_3 = X_4 = X_5 = 1\}$ ,  $n=5$ , the MLE is equal to 1, while the posterior mean is equal to  $6/7$  (prior mean takes into consideration that  $p < 1$  has strictly positive probability)

→ The POSTERIOR PREDICTIVE DISTRIBUTION of a new observation can be calculated by integration of the posterior distribution:

$$P(\tilde{X}=1 | \mathcal{X}_n) = \int_0^1 \underbrace{P(\tilde{X}=1 | p, \mathcal{X}_n)}_{\text{prediction for a new observation}} \underbrace{f(p | \mathcal{X}_n)}_{\text{posterior}} dp$$

$= p \cdot B(\alpha + S_n, n + \beta - S_n)$

$= \text{mean of the } B(\alpha + S_n, n + \beta - S_n) \text{ distribution.}$

$$= \frac{\alpha + S_n}{\alpha + \beta + n}$$

In a Machine Learning context,  $\mathcal{X}_n = \{(X_1, Y_1), \dots, (X_n, Y_n)\}$ . Having observed a new input point  $x$ , the posterior predictive

6

distribution of  $Y$  given  $X=x$  and  $\mathcal{X}_n$  is

$$f(y | x, \mathcal{X}_n) = \int_{\Theta} f(y, \theta | \mathcal{X}_n, x) d\theta$$

$\theta = \text{model parameter(s)}$

$$= \int_{\Theta} \underbrace{f(y | \theta, x)}_{\text{model distribution}} \underbrace{f(\theta | \mathcal{X}_n, x)}_{\text{posterior distribution}} d\theta$$

→ The posterior variance is

$$\text{var}(p | \mathcal{X}_n) = \frac{(\alpha + S_n)(n + \beta - S_n)}{(\alpha + \beta + n)^2 (\alpha + \beta + n + 1)}$$

$$= \frac{E(p | \mathcal{X}_n) (1 - E(p | \mathcal{X}_n))}{\alpha + \beta + n + 1}$$

As  $n \rightarrow \infty$ ,  $E(p | \mathcal{X}_n) \approx \frac{S_n}{n} = \text{MLE}$  (see p. 4), and  $\text{var}(p | \mathcal{X}_n) \approx \frac{1}{n} \frac{S_n}{n} (1 - \frac{S_n}{n}) = \frac{\hat{p}(1-\hat{p})}{n}$ .

Summarizing, as  $n$  gets large,  $E(p | \mathcal{X}_n) \approx \hat{p}$  &  $\text{var}(p | \mathcal{X}_n) \approx n^{-1} \hat{p}(1-\hat{p})$

Familiar expressions in non-Bayesian statistics.

In fact, under some regularity conditions, a CLT in Bayesian context holds:

$$\frac{p - E(p | \mathcal{X}_n)}{\sqrt{\text{var}(p | \mathcal{X}_n)}} | \mathcal{X}_n \xrightarrow{d} \mathcal{N}(0, 1), \text{ as } n \rightarrow \infty.$$

and we recover the usual "frequentist" CLT (7)

$$\frac{n^{1/2}(\hat{p} - p)}{\sqrt{\hat{p}(1-\hat{p})}} \xrightarrow{d} \mathcal{N}(0, 1).$$

Combined with Slutsky theorem, replacing  $p$  by its consistent estimator  $\hat{p}$ .

⇒ Under some regularity conditions, the posterior is usually approximated using a normal distribution.

Remark =  $\text{var } p = \mathbb{E} \text{var}(p | \mathcal{L}_n) + \text{var } \mathbb{E}(p | \mathcal{L}_n)$   
 $\uparrow \qquad \qquad \geq \mathbb{E} \text{var}(p | \mathcal{L}_n)$   
 prior variance posterior variance

The posterior variance is on average smaller than the prior variance ⇒ observations reduce the uncertainty about  $p$ . This property holds more generally, and is not specific about binomial proportions. ■

• Application: Comparing two proportions.

Suppose that we collect two independent learning samples:

$$X_{n_0} := \{X_1, \dots, X_{n_0}\}, \quad X_1, \dots, X_{n_0} \sim B(p_0)$$

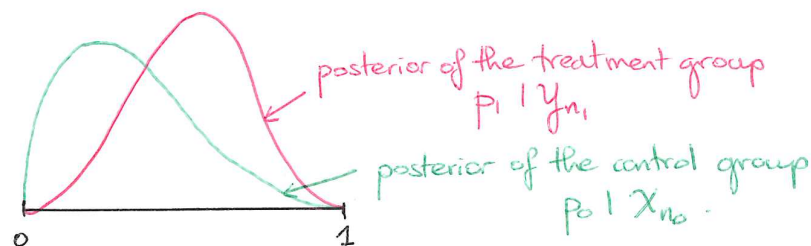
$$Y_{n_1} := \{Y_1, \dots, Y_{n_1}\}, \quad Y_1, \dots, Y_{n_1} \sim B(p_1),$$

and that we wish to infer which one of  $p_0$  or  $p_1$  is the largest. In clinical trials,  $X_{n_0}$  might represent a control group, who was administered a placebo, while individuals in  $Y_{n_1}$  receive a new treatment. The  $X_i/Y_i$  indicate

whether an individual recovered after some time, and the goal is to assess the new treatment's effectiveness. (8)

Prior  $p_0 \sim B(\alpha_0, \beta_0)$ ,  $p_1 \sim B(\alpha_1, \beta_1)$

Posterior  $p_0 | X_{n_0} \sim B(\alpha_0 + S_{n_0}, n_0 + \beta_0 - S_{n_0})$ ,  $S_{n_0} = \sum_{i=1}^{n_0} X_i$   
 $p_1 | Y_{n_1} \sim B(\alpha_1 + S'_{n_1}, n_1 + \beta_1 - S'_{n_1})$ ,  $S'_{n_1} = \sum_{i=1}^{n_1} Y_i$



In particular, we may be interested in the posterior distribution of  $p_0 - p_1 | X_{n_0} \cup Y_{n_1}$ . This can easily be done numerically using the following procedure:

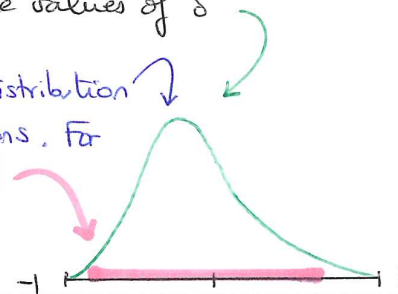
(i) Draw  $p_0 | X_{n_0} \sim B(\alpha_0 + S_{n_0}, n_0 + \beta_0 - S_{n_0})$   
 $p_1 | Y_{n_1} \sim B(\alpha_1 + S'_{n_1}, n_1 + \beta_1 - S'_{n_1})$

(ii) Compute  $\delta := p_0 - p_1$

(iii) Repeat (i)+(ii) as many times as desired

(iv) Draw a histogram of the values of  $\delta$

A summary of this posterior distribution may be used to draw conclusions. For example, is the 95% highest posterior density region strictly contained in some small interval?



Remarks (i) To compute the exact analytical expression 9

of the posterior distribution of  $\delta$ , we proceed as follows.

$$\text{Let } p_0 | X_{n_0} \sim B(\alpha_1, \beta_1)$$

$$p_1 | Y_{n_1} \sim B(\alpha_2, \beta_2),$$

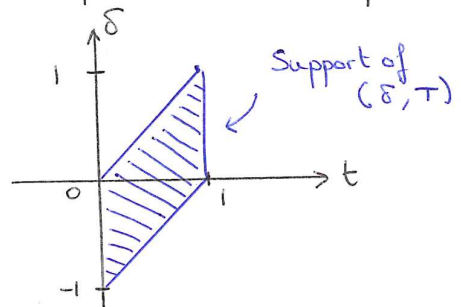
$$\text{where } \alpha_1 = \alpha_0 + S_{n_0}, \quad \beta_1 = n_0 + \beta_0 - S_{n_0}$$
$$\alpha_2 = \alpha_1 + S'_{n_1}, \quad \beta_2 = n_1 + \beta_1 - S'_{n_1}.$$

$$\text{Consider the change of variable } \begin{cases} \delta = p_0 - p_1 \\ T = p_0 \end{cases}$$

The Jacobian of this transformation is equal to 1, so that the joint density of  $(\delta, T)$  is

$$f_{\delta, T}(d, t) = f_{p_0, p_1}(t, t-d)$$
$$= B(t | \alpha_1, \beta_1) B(t-d | \alpha_2, \beta_2),$$

$$\text{where } 0 < t < 1$$
$$t-1 < d < t$$



It "remain" to integrate out the variable  $T$  to get the posterior density of  $\delta_0$  given  $X_{n_0} \cup Y_{n_1}$ .

(ii) Compare the approach with the frequentist paradigm, where we usually test for  $H_0: p_0 = p_1$  (versus some alternative that can be  $H_1: p_0 < p_1$ ),

or construct a confidence interval for  $\hat{\delta} = \hat{p}_0 - \hat{p}_1$ , 10  
based on large sample approximations, see p. 24 in  
MS: HYPOTHESIS TESTING.

(iii) Large-sample correspondence.

A multivariate normal distribution holds for the asymptotic density of the posterior distribution of the bivariate vector  $\begin{pmatrix} p_0 \\ p_1 \end{pmatrix}$ , given  $X_{n_0} \cup Y_{n_1}$ , as  $n_0, n_1 \rightarrow \infty$ , meaning that this holds true as well for the difference  $\delta = p_0 - p_1$ . In view of the results page 6 and 7, we see that the posterior distribution for  $\delta$  is asymptotically the same as the repeated sampling distribution derived from frequentist arguments  $\Rightarrow$  A 95% central posterior interval for  $\delta$  will cover the true value 95% of the time under repeated sampling, for any value of  $-1 < \delta < 1$ .

$\uparrow$  As we shall see later, this holds true as well in a wide range of interesting cases. ■

$\downarrow$  So far, we have derived a full Bayesian inference for a binomial proportion. The general procedure can be summarized as follows:

- (a) Specify a joint probability on the parameter  $\Theta$  & observations (prior + model given  $\Theta$ )
- (b) Compute the posterior distribution of  $\Theta$  given  $X_n$ .
- (c) Derive the posterior predictive distribution of a new observation.

The Chapter is organised as follows.

- Section I = Inference for the Normal Distribution, (I.1, I.2)  
Proper and Improper Priors (I.3) (I.4)
- Section II = Jeffreys Prior
- Section III = Bayesian vs Frequentists =
  - ↳ Bernstein-von-Mises Thm (III.1)
  - ↳ Hypothesis Testing (III.2)
  - ↳ Sequential Testing (III.3)
- Section IV = Empirical Bayes
  - ↳ The James-Stein Estimator (IV.1)
  - ↳ Bayesian FDR (IV.3)
- Section V = Hierarchical Models  
& Application to Multiple Testing (V.1)
- References.

I. BAYESIAN INFERENCE FOR THE NORMAL DISTRIBUTION

I.1. unknown mean, known variance.

Let  $\mathcal{X}_n = \{X_1, \dots, X_n\}$ , iid,  $X_i \sim \mathcal{N}(\mu, \sigma^2)$   
 unknown, with  $p \sim \mathcal{N}(\mu | \mu_0, \tau_0^2)$   
 known  $\tau_0^2$   
 prior distribution

The conjugate prior for  $\mu$ , when  $\sigma^2$  is known.

General definition of conjugacy:

$p(\theta | \mathcal{X}_n) \in \mathcal{P}$      $\forall$      $p(\cdot | \theta) \in \mathcal{F}$  = model distrib  
 $p(\cdot) \in \mathcal{P}$  = prior

parameter of interest    a family of distribution

- ⊕ analytical results usually available
- simplifies computations
- building blocks for more complicated models.

• The posterior distribution of  $\mu | \mathcal{X}_n$  is

$$f(\mu | \mathcal{X}_n) \propto f(\mathcal{X}_n | \mu) f(\mu)$$

$$\propto \prod_{i=1}^n \frac{1}{\sigma} \exp\left\{-\frac{1}{2\sigma^2}(x_i - \mu)^2\right\} \frac{1}{\tau_0} \exp\left\{-\frac{1}{2\tau_0^2}(\mu - \mu_0)^2\right\}$$

$$= \frac{1}{\sigma^n \tau_0} \exp\left\{-\frac{1}{2} \left( \frac{1}{\sigma^2} \sum (x_i - \mu)^2 + \frac{1}{\tau_0^2} (\mu - \mu_0)^2 \right)\right\}$$

write this term  $\frac{1}{\tau_n^2} (\mu - \mu_n)^2 + C$   
 & identify  $\mu_n$  and  $\tau_n^2$ .  
 cst indpt of  $\mu$ .

expand & complete the squares:

(13)

$$\begin{aligned} & \left(\frac{1}{\tau_0^2} + \frac{n}{\sigma^2}\right) \mu^2 - 2\left(\frac{n\bar{x}}{\sigma^2} + \frac{\mu_0}{\tau_0^2}\right) \mu + \dots \\ & = \underbrace{\left(\frac{1}{\tau_0^2} + \frac{n}{\sigma^2}\right)}_{=\frac{1}{\tau_n^2}} \left[ \mu^2 - 2\mu \underbrace{\left(\frac{n\bar{x}}{\sigma^2} + \frac{\mu_0}{\tau_0^2}\right)}_{=\mu_n} / \left(\frac{1}{\tau_0^2} + \frac{n}{\sigma^2}\right) + \dots \right] \end{aligned}$$

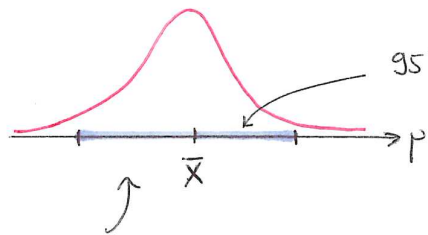
$\Rightarrow f(\mu | \mathcal{L}_n) \sim \mathcal{N}(\mu | \mu_n, \tau_n^2)$ , with

$$\mu_n = \frac{\mu_0/\tau_0^2 + n\bar{x}/\sigma^2}{1/\tau_0^2 + n/\sigma^2}, \quad \frac{1}{\tau_n^2} = \frac{1}{\tau_0^2} + \frac{n}{\sigma^2}$$

= weighted average of prior mean & sample mean

= prior precision + data precision.

x Remark: When  $n$  is large,  $\mu_n \approx \bar{x}$  and  $\tau_n^2 \approx \frac{\sigma^2}{n}$ , and we get the approximation  $\mu | \mathcal{L}_n \sim \mathcal{N}(\mu | \bar{x}, \frac{\sigma^2}{n})$



95% highest posterior density region is given by  $[\bar{x} - z_{1/2} \sigma, \bar{x} + z_{1/2} \sigma]$ ,

where  $z$  = corresponding quantile of the  $\mathcal{N}$  distrib.

Compare this interval with the 95% confidence interval for  $\mu \rightarrow$  the two coincide when  $n$  is large.

The posterior predictive distribution is

(14)

$$f(x | \mathcal{L}_n) = \int \omega(x | \mu, \sigma^2) \omega(\mu | \mu_n, \tau_n^2) d\mu$$

new sample point

model distribution

posterior distribution

(see top of p. 6)

convolution of two normal distribution is still normal.

$$\sim \mathcal{N}(x | m_n, \sigma_n^2), \quad \begin{cases} (a) \mu | \mathcal{L}_n \sim \mathcal{N}(\mu_n, \tau_n^2) \\ (b) X | \mu, \mathcal{L}_n \sim \mathcal{N}(x | \mu, \sigma^2) \end{cases}$$

with

$$m_n := \mathbb{E}(X | \mathcal{L}_n) = \mathbb{E}_\mu \mathbb{E}(X | \mu, \mathcal{L}_n) = \mathbb{E}(\mu | \mathcal{L}_n) = \mu_n$$

$$\sigma_n^2 := \text{var}(X | \mathcal{L}_n) = \mathbb{E} \underbrace{\text{var}(X | \mathcal{L}_n)}_{=\sigma^2} + \underbrace{\text{var} \mathbb{E}(X | \mathcal{L}_n)}_{=\text{var} \mu_n = \tau_n^2}$$

$$\Rightarrow f(x | \mathcal{L}_n) \sim \mathcal{N}(x | \mu_n, \sigma^2 + \tau_n^2)$$

(posterior predictive distribution of a new point)

## I.2. known mean, unknown variance.

This case is rarely useful in practice, but important as a building block for more complicated models.

$$\mathcal{L}_n = \{X_1, \dots, X_n\}, \text{ iid, } X_i \sim \mathcal{N}(\mu, \sigma^2)$$

known

unknown, with  $\sigma^2 \sim \text{Inv-}\chi^2(\nu_0, \sigma_0^2)$   
(scaled inverse  $\chi^2$  distrib with  $\nu_0$  d.o.f.)

Toolbox

15

(a)  $X \sim \chi^2(\nu_0)$  = Chi-square distribution with  $\nu_0$  degrees of freedom,  $\nu_0 > 0$

X has density

$$f(x) = \frac{1}{\Gamma(\nu_0/2)} 2^{-\nu_0/2} x^{\nu_0/2-1} e^{-x/2}, \quad x > 0$$

&  $EX = \nu_0$  ;  $\text{var } X = 2\nu_0$ .

(b) If  $X \sim \chi^2(\nu_0)$ , then  $\frac{1}{X} \sim \text{Inv-}\chi^2(\nu_0)$ ,  $\nu_0 > 0$

A variable  $\Theta \sim \text{Inv-}\chi^2(\nu_0)$  has density

$$f(\theta) = \frac{2^{-\nu_0/2}}{\Gamma(\nu_0/2)} \theta^{-(\nu_0/2+1)} e^{-1/2\theta}, \quad \theta > 0$$

&  $EQ = \frac{1}{\nu_0-2}$  &  $\text{var } \Theta = \frac{2}{(\nu_0-2)^2(\nu_0-4)}$   
 (for  $\nu_0 > 2$ ) (for  $\nu_0 > 4$ )

(c)  $\Theta \sim \text{Inv-}\chi^2(\nu_0, \sigma_0^2)$  has density

$$f(\theta) = \frac{(\sigma_0^2 \nu_0/2)^{\nu_0/2}}{\Gamma(\nu_0/2)} \theta^{-(\nu_0/2+1)} e^{-\nu_0 \sigma_0^2 / 2\theta}, \quad \theta > 0$$

$EQ = \frac{\nu_0}{\nu_0-2} \sigma_0^2$  &  $\text{var } \Theta = \frac{2\nu_0^2 \sigma_0^4}{(\nu_0-2)^2(\nu_0-4)}$

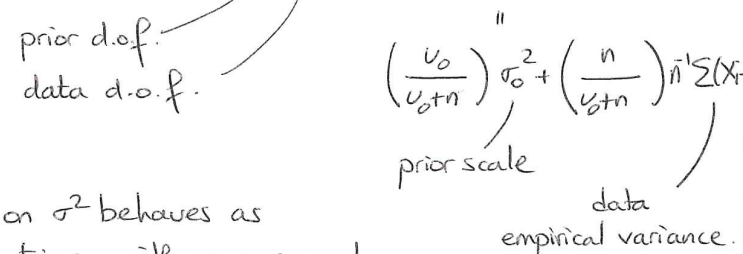
If  $\Theta \sim \text{Inv-}\chi^2(\nu_0, \sigma_0^2)$ , then

$$\frac{\Theta}{\nu_0 \sigma_0^2} \sim \text{Inv-}\chi^2(\nu_0)$$

• The posterior distribution of  $\sigma^2$  given  $\mathcal{X}_n$  is

16

$$\begin{aligned} f(\sigma^2 | \mathcal{X}_n) &\propto f(\mathcal{X}_n | \sigma^2) f(\sigma^2) \\ &\propto \sigma^{-n} \exp\left\{-\frac{1}{2\sigma^2} \sum_{i=1}^n (x_i - \mu)^2\right\} \\ &\quad \times \left(\frac{1}{\sigma^2}\right)^{\nu_0/2+1} \exp\left\{-\frac{\nu_0 \sigma_0^2}{2\sigma^2}\right\} \\ &= \left(\frac{1}{\sigma^2}\right)^{\frac{\nu_0+n}{2}+1} \exp\left\{-\frac{\nu_0 \sigma_0^2 + \sum (x_i - \mu)^2}{2\sigma^2}\right\} \\ &\sim \text{Inv-}\chi^2\left(\nu_0+n, \frac{\nu_0 \sigma_0^2 + \sum (x_i - \mu)^2}{\nu_0+n}\right) \end{aligned}$$



↳ The prior on  $\sigma^2$  behaves as  $\nu_0$  observations, with an averaged squared deviation equal to  $\sigma_0^2$ .

I.3. Proper and Improper prior distributions.

Recall the two examples from sections I.1 and I.2:

$X_1, \dots, X_n$  iid  $\sim \mathcal{N}(\mu, \sigma^2)$

•  $\mu$  unknown,  $\sigma^2$  known

$\mu \sim \mathcal{N}(\mu_0, \tau_0^2) \Rightarrow \mu | \mathcal{X}_n \sim \mathcal{N}(\mu_n, \tau_n^2)$

$$\mu_n = \frac{\mu_0/\tau_0^2 + n\bar{x}/\sigma^2}{1/\tau_0^2 + n/\sigma^2}$$

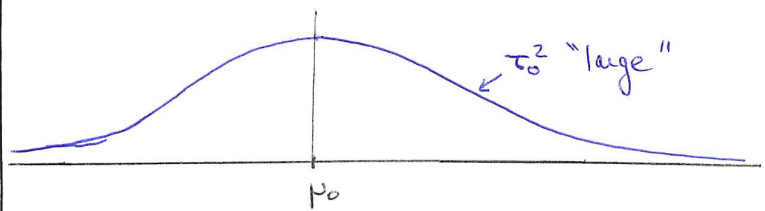
$$\frac{1}{\tau_n^2} = \frac{1}{\tau_0^2} + \frac{n}{\sigma^2}$$



When  $\frac{n}{\tau_0^2} \gg \frac{1}{\tau_0^2}$ , then  $\mu_n \approx \bar{x}$ ,  $\tau_n^2 \approx \frac{n}{\tau_0^2}$  (17)

data precision (green) / prior precision (green) as if  $\tau_0^2 = +\infty$

This situation arises when the prior is "flat"  $\Rightarrow$  little prior knowledge about the most likely values of  $\mu$  are incorporated into the model.



The limiting case  $\tau_0^2 = +\infty$  is not a proper density: it does not integrate to 1. However, provided one data point is observed, the posterior  $\mathcal{N}(\bar{x}, n/\sigma^2)$  is a proper distribution.

"flat" constant prior on  $\mu$ :  $f(\mu) \sim 1$   
 $X_1, \dots, X_n \sim \mathcal{N}(\mu, \sigma^2)$ ,  $\sigma^2$  known

"the likelihood dominates the prior" (green)

The joint distribution is not properly defined, but, provided one data point is observed, the posterior  $\mu | \mathcal{L}_n$  is properly defined.

•  $\mu$  known,  $\sigma^2$  unknown

$$\sigma^2 \sim \text{Inv-}\chi^2(\nu_0, \tau_0^2)$$

$$\Rightarrow \sigma^2 | \mathcal{L}_n \sim \text{Inv-}\chi^2(\nu_0 + n, \frac{\nu_0 \tau_0^2 + \sum (X_i - \mu)^2}{\nu_0 + n})$$

$\nu_0 = \text{prior d.o.f} \equiv$  carries as much information as  $\nu_0$  observations.

If  $n \gg \nu_0$ , then  $\sigma^2 | \mathcal{L}_n \approx \text{Inv-}\chi^2(n, n^{-1} \sum (X_i - \mu)^2)$  (18)

data d.o.f (green) / prior d.o.f. (green)

The limiting case  $\nu_0 = 0$  corresponds to the case where the prior contains as much information as 0 observation  $\Rightarrow$  uninformative. This limiting case is also obtained by defining an improper prior on  $\sigma^2$ , given by  $f(\sigma^2) \sim \frac{1}{\sigma^2}$ .

uninformative prior on  $\sigma^2$  (improper)  $\sim 1/\sigma^2$  } joint distribution is not defined, but,  $\Rightarrow$  posterior  $\sigma^2 | \mathcal{L}_n$  is proper, provided at least one data point is observed. (does not integrate to 1)

I. 4. unknown mean and variance

Consider the case where  $X_1, \dots, X_n$  iid  $\mathcal{N}(\mu, \sigma^2)$ , where both  $\mu$  and  $\sigma^2$  are unknown. We proceed as before:

- (a) Compute the posterior distribution  $f(\mu, \sigma^2 | \mathcal{L}_n)$
- (b) Derive the posterior predictive distribution  $f(x | \mathcal{L}_n)$  of a new observation  $X$ .

$\rightarrow$  Consider first a non-informative prior on  $(\mu, \sigma^2)$ :  
 $f(\mu, \sigma^2) \sim \frac{1}{\sigma^2}$ . (we assume a more general conjugate prior later)

$\rightarrow$  Posterior is

$$f(\mu, \sigma^2 | \mathcal{L}_n) \sim \frac{1}{\sigma^{n+2}} \exp \left\{ -\frac{1}{2\sigma^2} \sum_{i=1}^n (X_i - \mu)^2 \right\}$$

$$f(\mu, \sigma^2 | \mathcal{L}_n) \sim \frac{1}{\sigma^{n+2}} \exp \left\{ -\frac{1}{2\sigma^2} \left[ \underbrace{\sum_{i=1}^n (X_i - \bar{X})^2}_{=: (n-1)s^2} + n(\bar{X} - \mu)^2 \right] \right\} \quad (19)$$

$$= \frac{1}{\sigma^{n+2}} \exp \left\{ -\frac{1}{2\sigma^2} \left[ (n-1)s^2 + n(\bar{X} - \mu)^2 \right] \right\}$$

We decompose this expression further into the product

$$f(\mu, \sigma^2 | \mathcal{L}_n) = f(\mu | \sigma^2, \mathcal{L}_n) f(\sigma^2 | \mathcal{L}_n)$$

This term ( $\sigma^2$  fixed) was considered in section I.1

$$f(\mu | \sigma^2, \mathcal{L}_n) = \mathcal{N}(\bar{X}, \frac{\sigma^2}{n})$$

Make use of

$$f(\sigma^2 | \mathcal{L}_n) = \int f(\mu, \sigma^2 | \mathcal{L}_n) d\mu$$

$$\sim \frac{1}{\sigma^{n+2}} e^{-\frac{(n-1)s^2}{2\sigma^2}} \int e^{-\frac{n}{2\sigma^2} (\bar{X} - \mu)^2} d\mu$$

$$= (2\pi \frac{\sigma^2}{n})^{1/2}$$

$$f(\sigma^2 | \mathcal{L}_n) \sim \frac{1}{\sigma^{n+1}} \exp \left\{ -\frac{(n-1)s^2}{2\sigma^2} \right\}$$

$$= \text{Inv-}\chi^2(n-1, s^2)$$

We get

$$f(\mu, \sigma^2 | \mathcal{L}_n) = \mathcal{N}(\mu | \bar{X}, \frac{\sigma^2}{n}) \text{Inv-}\chi^2(\sigma^2 | n-1, s^2)$$

$$= f(\mu | \sigma^2, \mathcal{L}_n) f(\sigma^2 | \mathcal{L}_n)$$

Next, we derive the posterior predictive distribution of a new observation

$$f(x | \mathcal{L}_n) = \int f(x, \mu, \sigma^2 | \mathcal{L}_n) d\mu d\sigma^2$$

$$f(x | \mathcal{L}_n) = \int \underbrace{f(x | \mu, \sigma^2, \mathcal{L}_n)}_{\text{model distribution } \mathcal{N}(x | \mu, \sigma^2)} \underbrace{f(\mu | \sigma^2, \mathcal{L}_n)}_{\text{derived on page 19 (posterior distribution)}} f(\sigma^2 | \mathcal{L}_n) d\mu d\sigma^2 \quad (20)$$

$$= \int \mathcal{N}(x | \mu, \sigma^2) \mathcal{N}(\mu | \bar{X}, \frac{\sigma^2}{n}) \text{Inv-}\chi^2(\sigma^2 | n-1, s^2) d\mu d\sigma^2$$

$$= \int [\mathcal{N}(x | \mu, \sigma^2) \mathcal{N}(\mu | \bar{X}, \frac{\sigma^2}{n}) d\mu] \text{Inv-}\chi^2 d\sigma^2$$

see top of p.14

$$= \int \mathcal{N}(x | \bar{X}, (1 + \frac{1}{n})\sigma^2) \text{Inv-}\chi^2(\sigma^2 | n-1, s^2) d\sigma^2$$

Note that  $\sigma^2 \sim \text{Inv-}\chi^2(\sigma^2 | n-1, s^2)$

$$\Leftrightarrow \frac{\sigma^2}{(n-1)s^2} \sim \text{Inv-}\chi^2(\sigma^2 | n-1)$$

$$\Leftrightarrow \frac{(n-1)s^2}{\sigma^2} \sim \chi^2\left(\frac{1}{\sigma^2} | n-1\right)$$

$$\Leftrightarrow (n-1)s^2 \lambda \sim \gamma\left(\lambda | \frac{n-1}{2}, \frac{1}{2}\right)$$

$$\Leftrightarrow \lambda \sim \gamma\left(\lambda | \frac{n-1}{2}, \frac{(n-1)s^2}{2}\right)$$

scaling (use MGF)

bottom of p.15

p.15, (b)

with  $\lambda := \frac{1}{\sigma^2}$

& equivalence

$\chi^2$  &  $\gamma$  distributions

$$\Rightarrow f(x | \mathcal{L}_n) = \int \mathcal{N}(x | \bar{X}, (1 + \frac{1}{n})\lambda^{-1}) \gamma(\lambda | \frac{n-1}{2}, \frac{(n-1)s^2}{2}) d\lambda$$

$$= t\left(n-1, \bar{X}, s^2\left(1 + \frac{1}{n}\right)\right)$$

d.o.f. location

scale is  $s\left(1 + \frac{1}{n}\right)^{1/2}$

see p.11/12 in SL: BAYESIAN LINEAR MODELS for details.

$$\left[ X \sim t\left(\nu, \mu, \sigma^2\right) \text{ has density } \frac{\Gamma\left(\frac{\nu+1}{2}\right)}{\Gamma\left(\frac{\nu}{2}\right)} \frac{1}{\sigma\sqrt{\nu\pi}} \left(1 + \frac{1}{\nu} \left(\frac{x-\mu}{\sigma}\right)^2\right)^{-\frac{\nu+1}{2}} \right]$$

• Summary : To draw samples  $X \sim$  posterior predictive distribution, having observed  $\mathcal{L}_n$ , we can proceed in two equivalent ways: (21)

(1) Draw  $\sigma^2 \mid \mathcal{L}_n \sim \text{Inv-}\chi^2(\sigma^2 \mid n-1, s^2)$

(2) Draw  $\mu \mid \sigma^2, \mathcal{L}_n \sim \mathcal{N}(\mu \mid \bar{x}, \frac{\sigma^2}{n})$   $s^2 = \frac{1}{n-1} \sum_{i=1}^n (x_i - \bar{x})^2$

(3) Draw  $X \mid \mu, \sigma^2, \mathcal{L}_n \sim \mathcal{N}(x \mid \mu, \sigma^2)$

or directly using a scaled & shifted  $t$  distribution

(i) Draw  $X \mid \mathcal{L}_n \sim t(n-1, \bar{x}, s^2(1 + \frac{1}{n}))$

$\uparrow$   $X_1, \dots, X_n \sim \mathcal{N}(\mu, \sigma^2)$  + non-informative prior on  $(\mu, \sigma^2)$

Remarks (i) The marginal posterior distribution of  $\mu$  given  $\mathcal{L}_n$  can be analytically computed:

$$f(\mu \mid \mathcal{L}_n) \sim A^{-n/2} \int_0^{+\infty} u^{\frac{n-2}{2}} e^{-u} du$$

A defined on top of page 19

&  $u := A/2\sigma^2$

unnormalized Gamma integral

$$\sim [(n-1)s^2 + n(\mu - \bar{x})^2]^{-n/2}$$

$$\sim \left[ 1 + \frac{(\mu - \bar{x})^2}{(n-1)(s^2/n)} \right]^{-\frac{n}{2}}$$

$$\Rightarrow f(\mu \mid \mathcal{L}_n) = t(n-1, \bar{x}, \frac{s^2}{n})$$

$$\Leftrightarrow \frac{\mu - \bar{x}}{s/\sqrt{n}} \mid \mathcal{L}_n \sim t_{n-1} \leftarrow \begin{array}{l} \text{frequentist \&} \\ \text{Bayesian} \\ \text{intervals are} \\ \text{equal} \end{array}$$

& compare with the frequentist result:  $\frac{\bar{x} - \mu}{s/\sqrt{n}} \mid \mu, \sigma^2 \sim t_{n-1}$   
 compare also with the result on p. 13, with  $\mu$  unknown &  $\sigma^2$  known

(ii) We may assume more generally a conjugate prior for  $(\mu, \sigma^2)$  (the non-informative case can be recovered here as well as a limit). (22)

Motivated by the factorization of the posterior on p. 19,  $f(\mu, \sigma^2 \mid \mathcal{L}_n) = f(\mu \mid \sigma^2, \mathcal{L}_n) f(\sigma^2 \mid \mathcal{L}_n)$

$\mathcal{N} \uparrow$   $\text{Inv-}\chi^2 \uparrow$

a similar factorization must hold for the prior  $f(\mu, \sigma^2) = f(\mu \mid \sigma^2) f(\sigma^2)$

$\mathcal{N} \uparrow$   $\text{Inv-}\chi^2 \uparrow$

the conjugate prior.

It is convenient to parametrize the conjugate prior as follows =

$$\mu \mid \sigma^2 \sim \mathcal{N}(\mu \mid \mu_0, \frac{\sigma^2}{k_0})$$

$$\sigma^2 \sim \text{Inv-}\chi^2(\sigma^2 \mid \nu_0, \sigma_0^2)$$

Then it can be shown after calculations that

$$\begin{aligned} f(\mu, \sigma^2 \mid \mathcal{L}_n) &= f(\mu \mid \sigma^2, \mathcal{L}_n) f(\sigma^2 \mid \mathcal{L}_n) \\ &= \mathcal{N}(\mu \mid \mu_n, \frac{\sigma^2}{k_n}) \text{Inv-}\chi^2(\sigma^2 \mid \nu_n, \sigma_n^2) \end{aligned}$$

where

$$\mu_n = \frac{k_0}{k_0 + n} \mu_0 + \frac{n}{k_0 + n} \bar{x}$$

$$k_n = k_0 + n$$

$$\nu_n = \nu_0 + n$$

$$\nu_n \sigma_n^2 = \nu_0 \sigma_0^2 + (n-1) s^2 + \frac{k_0 n}{k_0 + n} (\bar{x} - \mu_0)^2$$

$\leftarrow$  interpretation?

The marginal posterior distributions can be computed as well,

$$f(\mu | \mathcal{L}_n) = t(\nu_n, \mu_n, \frac{\sigma_n^2}{K_n}) \quad \text{--- (i)}$$

$$f(\sigma^2 | \mathcal{L}_n) = \text{Inv-}\chi^2(\nu_n, \sigma_n^2) \quad \text{--- (ii)}$$

We get from (i) that  $\frac{\mu - \mu_n}{\sigma_n / \sqrt{K_n}} \sim t_{\nu_n}$   
 (compare with the expression p.21)

Summary = inference for  $\mu$ .

	$\mu$ unknown $\sigma^2$ known		$\mu$ unknown $\sigma^2$ unknown	
BAYESIAN	non-informative prior $f(\mu) \sim 1$	conjugate prior $\sim \mathcal{U}(\mu   \mu_0, \tau_0^2)$ $\uparrow$ p.12	non-informative prior $f(\mu, \sigma^2) \sim \frac{1}{\sigma^2}$ $\uparrow$ p.18	conjugate prior $\sim \mathcal{U}(\mu   \mu_0, \frac{\sigma^2}{K_0})$ $\times \text{Inv}\chi^2(\sigma^2   \nu_0, \frac{\sigma_0^2}{K_0})$ $\uparrow$ p.22
	marginal posterior $\mu   \mathcal{L}_n \sim \mathcal{U}(\mu   \bar{x}, \frac{\sigma^2}{n})$	$\mu   \mathcal{L}_n \sim \mathcal{U}(\mu   \mu_n, \frac{\sigma_n^2}{K_n})$ $\uparrow$ p.13  n large: $\mu_n \approx \bar{x}$ $\frac{\sigma_n^2}{K_n} \approx \frac{\sigma^2}{n}$	$\mu   \mathcal{L}_n \sim t(n-1, \bar{x}, \frac{s^2}{n})$ $\uparrow$ p.20	$\mu   \mathcal{L}_n \sim t(\nu_n, \mu_n, \frac{\sigma_n^2}{K_n})$ $\uparrow$ p.23  n large: $\mu_n \approx \bar{x}$ $\frac{\sigma_n^2}{K_n} \approx \frac{s^2}{n}$
FREQUENTIST	$\bar{x} \sim \mathcal{U}(\mu, \frac{\sigma^2}{n})$  $\frac{\bar{x} - \mu}{\sigma/\sqrt{n}} \sim \mathcal{U}(0, 1)$	$\bar{x} \sim t(n-1, \mu, \frac{s^2}{n})$  $\frac{\bar{x} - \mu}{s/\sqrt{n}} \sim t(n-1)$		

*Annotations:*  
 - Credible bounds  $\equiv$  Confidence bounds  
 - Same when n is large

II. JEFFREYS PRIOR

We often want to consider a prior that carries as little information as possible. When the parameter space is discrete and contains exactly  $m$  distinct values, then we just put a mass  $1/m$  at each parameter value. When the parameter space is continuous however, there are several ways of achieving this. We discussed in Section I.3 uniform priors for the normal distribution, taken as limits of conjugate priors:

$$\mathcal{U}(\mu, \sigma^2), \mu \ \& \ \sigma^2 \text{ unknown} \rightarrow (\mu, \sigma^2) \sim \frac{1}{\sigma^2}$$

Similarly, a uniform prior for a binomial proportion  $p$  is  $\mathcal{U}(0, 1)$ .

$\leftarrow$  flat over the parameter space  $[0, 1]$ , which is compact  $\Rightarrow$  proper prior

Although a uniform distribution carries no information about  $p$  itself,  $\sqrt{p}$  has a higher density near 1 than 0.

$\hookrightarrow$  ignorance about  $p$  leads to knowledge about  $\sqrt{p}$ .

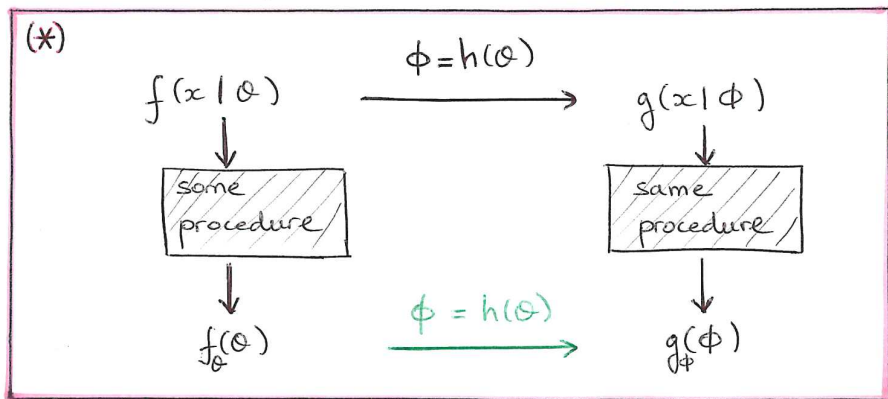
Likewise, under the assumption that  $p \sim \mathcal{U}[0, 1]$ , the transformation  $\phi := \log\left(\frac{p}{1-p}\right)$  (log odds) has

$$\text{density } f_\phi(\phi) = \frac{e^\phi}{(1+e^\phi)^2} \leftarrow \text{not a flat prior over } \mathbb{R} : \text{some values of } \phi \text{ are more likely than others.}$$

This was considered a major drawback in the early days

of Bayesian statistics. The uniform prior cannot be accepted as a "universal" non-informative prior, since it carries information under monotonic transformations of the "canonical" parameter.

Let  $\theta$  denote the original parameter, and consider  $\phi := h(\theta)$ , for some 1-1 monotone mapping  $h$ . The goal is to find a procedure/principle that produces a prior  $f_{\theta}(\theta)$  for any input model density  $f(x|\theta)$  parametrized by  $\theta$ .



Flat priors (proper & improper) fail to satisfy (\*).

Indeed, taking ~~some procedure~~  $\equiv$  ~~flat prior~~, and

applying it to the binomial distribution parametrized by  $\theta = p$  and  $\phi = \log\left(\frac{p}{1-p}\right)$ , we see that we should get  $f_{\theta}(\theta) = \mathcal{U}[0, 1]$  and  $g_{\phi}(\phi) \sim 1$ ,

while we saw on page 24 that transforming  $f_{\theta}(\theta)$  yields the prior  $\frac{e^{\phi}}{(1+e^{\phi})^2}$ . A valid universal

procedure must therefore produce the same prior  $g_{\phi}(\phi)$  whether (a) we first transform  $\phi = h(\theta)$  & apply the procedure or (b) apply the procedure to  $f(x|\theta)$  and then consider the change of variable  $\phi = h(\theta)$ .

Note that for  $\theta \sim f_{\theta}(\theta)$  and the change of variable  $\phi = h(\theta)$  [ $\theta = h^{-1}(\phi)$ ], the density of  $\phi$  is given by

$$g_{\phi}(\phi) = f_{\theta}(\theta) \left| \frac{d\theta}{d\phi} \right|$$

$$= f_{\theta}(h^{-1}(\phi)) \left| \frac{dh^{-1}(\phi)}{d\phi} \right|$$

Likewise, for  $\phi \sim g_{\phi}(\phi)$  and the change of variable  $\theta = h^{-1}(\phi)$  [ $\phi = h(\theta)$ ], the density of  $\theta$  is given by

$$f_{\theta}(\theta) = g_{\phi}(\phi) \left| \frac{d\phi}{d\theta} \right|$$

$$= g_{\phi}(h(\theta)) \left| \frac{dh(\theta)}{d\theta} \right|$$

Fisher information

Harold Jeffreys (1946) proposed  $f_{\theta}(\theta) \sim \sqrt{I_{\theta}(\theta)}$  as the default "objective" choice of prior distribution.

(27)

Recall that Fisher information is given by

$$I_{\theta}(\theta) = \mathbb{E} \left[ \left( \frac{\partial \log f(X|\theta)}{\partial \theta} \right)^2 \right]$$

$$= - \mathbb{E} \left[ \frac{\partial^2 \log f(X|\theta)}{\partial \theta^2} \right]$$

under regularity conditions

(cf MS: MAXIMUM LIKELIHOOD ESTIMATION)

We show that Jeffreys prior is invariant under smooth monotone transformations of the input parameter.

We have  $I_{\theta}(\theta) = |h'(\theta)|^2 I_{\phi}(h(\theta))$  — (\*\*)

Indeed,  $I_{\theta}(\theta) = \mathbb{E} \left( \frac{\partial}{\partial \theta} \log f(X|\theta) \right)^2$

$$= \mathbb{E} \left( \frac{\partial}{\partial \theta} \log g(X|h(\theta)) \right)^2$$

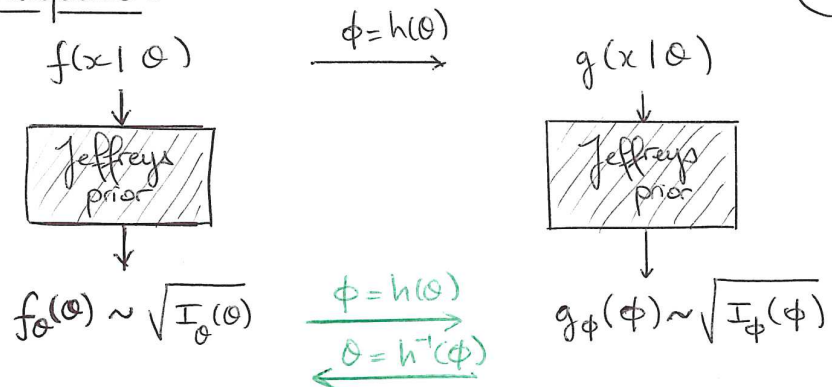
$$= \mathbb{E} \left( \frac{\partial}{\partial \phi} \log g(X|\phi) \frac{\partial \phi}{\partial \theta} \right)^2$$

$$= \left( \frac{\partial \phi}{\partial \theta} \right)^2 \mathbb{E} \left( \frac{\partial}{\partial \phi} \log g(X|\phi) \right)^2$$

$$= |h'(\theta)|^2 I_{\phi}(h(\theta)) \quad \blacksquare$$

(28)

x Consequence =



Indeed, starting with  $g_{\phi}(\phi) \sim \sqrt{I_{\phi}(\phi)}$  and considering the change of variable  $\theta = h^{-1}(\phi)$  yields the density

$$g_{\phi}(\phi) \left| \frac{d\phi}{d\theta} \right| \sim \sqrt{I_{\phi}(\phi)} \left| \frac{d\phi}{d\theta} \right| = \sqrt{I_{\theta}(\theta)} \sim f_{\theta}(\theta),$$

as required.

(\*\*) p. 27

x Remarks (i) If  $\int \sqrt{I_{\theta}(\theta)} d\theta < +\infty$ , Jeffreys prior is proper and given by  $f_{\theta}(\theta) = \frac{1}{C_{\theta}} \sqrt{I_{\theta}(\theta)}$ .

If the integral is infinite, the prior is improper. This is not an issue as long as the posterior distribution is a proper density.

(ii) Binomial distribution.

$$\bullet \theta = p \rightarrow f(x|\theta) = \theta^x (1-\theta)^{1-x} \quad x \in \{0, 1\}$$

$$\phi = h(\theta) = \log\left(\frac{\theta}{1-\theta}\right)$$

$$\begin{aligned} \hookrightarrow g(x|\phi) &= f(x|h^{-1}(\phi)) \\ &= \left(\frac{e^\phi}{1+e^\phi}\right)^x \left(\frac{1}{1+e^\phi}\right)^{1-x} \\ &= \frac{e^{\phi x}}{(1+e^\phi)^2} \end{aligned}$$

$$h'(u) = \frac{1}{u(1-u)}$$

$$h^{-1}(x) = \frac{e^x}{1+e^x}$$

$$\text{Then } \log f(x|\theta) = x \log \theta + (1-x) \log(1-\theta)$$

$$\& \frac{\partial^2}{\partial \theta^2} \log f(x|\theta) = -\frac{x}{\theta^2} - \frac{1-x}{(1-\theta)^2}$$

$$\Rightarrow I_\theta(\theta) = \frac{1}{\theta(1-\theta)} \quad \& \quad \text{Jeffreys prior is } \sim \frac{1}{\sqrt{\theta(1-\theta)}} \\ \text{B}(\frac{1}{2}, \frac{1}{2}).$$

$$\text{Also, } \log g(x|\phi) = \phi x - \log(1+e^\phi)$$

$$\& \frac{\partial^2}{\partial \phi^2} \log g(x|\phi) = -\frac{e^\phi}{(1+e^\phi)^2}$$

$$\Rightarrow I_\phi(\phi) = \frac{e^\phi}{(1+e^\phi)^2} \Rightarrow \text{Jeffreys prior is } \sim \frac{e^{\phi/2}}{1+e^\phi}$$

$$\text{Summarizing, } f_\theta(\theta) \sim \frac{1}{\sqrt{\theta(1-\theta)}} \text{ and } g_\phi(\phi) \sim \frac{e^{\phi/2}}{1+e^\phi}$$

$$\theta = h^{-1}(\phi) = \frac{e^\phi}{1+e^\phi}$$

The change of variable  $\theta = h^{-1}(\phi)$  yields the density  $g_\phi(h(\theta)) \left| \frac{dh(\theta)}{d\theta} \right| \sim \frac{e^{h(\theta)/2}}{1+e^{h(\theta)}} \times \frac{1}{\theta(1-\theta)}$   
 $= \frac{1}{\sqrt{\theta(1-\theta)}}$ , as required  $\square$

(iii)  $\mathcal{N}(\mu, \sigma^2)$ ,  $\mu$  unknown,  $\sigma^2$  known.

Then  $I_\mu(\mu) = \frac{1}{\sigma^2}$  so Jeffreys prior is  $\sim 1$  (improper)

(iv) Multivariate case:  $f_\theta(\theta) \sim \sqrt{\det(I_\theta(\theta))}$

Ex:  $\mathcal{N}(\mu, \sigma^2)$ ,  $\mu$  and  $\sigma^2$  unknown

Then one can show that  $I_\theta(\theta) = \begin{pmatrix} 1/\sigma^2 & 0 \\ 0 & 1/2\sigma^4 \end{pmatrix}$   
 $\theta = (\mu, \sigma^2)$

So that  $\det(I_\theta(\theta)) = \frac{1}{2\sigma^6}$

Jeffreys prior is  $f(\mu, \sigma^2) \sim \frac{1}{\sigma^3}$

differs from the non-informative prior  $\frac{1}{\sigma^2}$  encountered before.

(v) Bernardo (1979) showed that Jeffreys prior maximizes the (average) KL divergence between the prior and the posterior: the one that maximizes the contribution of the data and minimizes the contribution of the prior.

### III - BAYESIAN VS FREQUENTIST

(31)

#### III.1. Bernstein - Von Mises theorem

- We establish in this section some connections between posterior-based inference and the frequentist approach.
- Assume that  $\mathcal{L}_n = \{X_1, \dots, X_n\}$  are independent observations sampled from a common fixed distribution  $P_0$ , with density  $\gamma_0$ .

↖ A common assumption in frequentist statistics. The posterior distribution will be interpreted under this framework.

- The data are modeled using the parametric family  $f(x|\theta)$ ,  $\theta \in \Theta$  and prior  $f(\theta)$ . The posterior distribution is denoted  $f(\theta|\mathcal{L}_n)$ . Put

$$\pi(B|\mathcal{L}_n) := \int_B f(\theta|\mathcal{L}_n) d\theta$$

↖ Posterior Probability of B

↖ Borel set

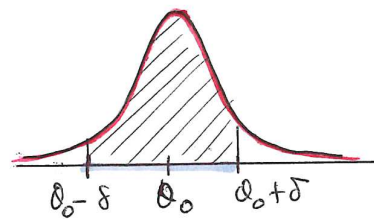
(\*) Assume that the parametric family  $f(\cdot|\theta)$  includes the "true" density  $\gamma_0$ ; i.e.  
 $\exists \theta_0 \in \Theta$  s.t.  $\gamma_0(x) = f(x|\theta_0)$

↖ We relax this assumption later, and consider the misspecified case where there is no such  $\theta_0$ .

Under (\*) and some regularity conditions (given more explicitly below for a more general result), CONSISTENCY follows:

As the number of observations tends to  $+\infty$ , the posterior distribution collapses to a point mass at  $\theta_0$ .  
 $\forall \delta > 0 \quad \pi(\{\theta : \|\theta - \theta_0\| > \delta\} | \mathcal{L}_n) \rightarrow 0$   
as  $n \rightarrow +\infty$ , in probability (under  $P_0$ ).

↖ No that consistency holds in probability under the true  $P_0$ , hence the need to introduce it. In fact, it is otherwise possible for most interesting problems to construct a sequence  $x_1, x_2, \dots$  that violates the convergence to 0.



← We can put as much mass as possible in any neighborhood of  $\theta_0$ , by taking  $n$  large enough. Under some regularity conditions, we can establish

that the posterior distribution approaches normality as  $n \rightarrow \infty$ . This result is known as the Bernstein - Von Mises theorem.

[Reg Conditions]

(1)  $X_1, \dots, X_n$  iid  $\sim P_0$ . (otherwise the data may not bring enough information about the parameter of interest)

(32)



- (2) The prior puts a positive mass around  $\theta_0$ . (otherwise the likelihood has no chance to dominate the prior)
- (3) Smoothness (twice continuously differentiable) & identifiability of the log-likelihood. (otherwise there is not a single point  $\theta_0$  the posterior can converge to)
- (4)  $\theta_0$  is an interior point of  $\Theta$ . (otherwise the normal distribution may not be appropriate)

[BERNSTEIN - VON-MISES THM]

Under the regularity conditions stated above, TV distance

$$\sup_B \left| \Pi(B | \mathcal{L}_n) - \mathcal{N}(\hat{\theta}_n, [n I_1(\theta_0)]^{-1}) \right| \xrightarrow{\text{TV distance}} 0$$

in probability (under  $P_{\theta_0}$ ), as  $n \rightarrow \infty$ .

MLE      Fisher Information of 1 observation

A very powerful result. It shows that for any measurable set  $B$ , the full posterior distribution is close to the normal distribution. It has huge implications when interpreting credible intervals constructed from the posterior density. Indeed, under regularity conditions, the MLE  $\hat{\theta}_n$  is consistent ( $\hat{\theta}_n \rightarrow \theta_0$  in probability) and asymptotically normal  $[n I_1(\theta_0)]^{1/2} (\hat{\theta}_n - \theta_0) \xrightarrow{d} \mathcal{N}(0, 1)$

↖ Same covariance matrix as in the B-VM theorem.

"Frequentist" confidence intervals constructed from the asymptotic normal distribution of the MLE take the form

$$C_n = \hat{\theta}_n \pm z_{1-\alpha/2} n^{-1/2} I_1(\theta_0)^{-1/2} \leftarrow \text{"Wald" interval}$$

Has approximate coverage  $(1-\alpha)$  i.e. usually estimated by replacing  $\theta_0$  with  $\hat{\theta}_n$ .

$$P_{\theta_0}(\theta_0 \in C_n) \approx 1-\alpha$$

Let  $B_n$  be such that  $\Pi(B_n | \mathcal{L}_n) = 1-\alpha$ . The Bernstein-Von-Mises theorem ensures that  $B_n = C_n + o_{P_{\theta_0}}(1)$ .  
 $\Rightarrow$  Under regularity assumptions, Bayesian credible sets are valid confidence sets in the frequentist sense: central sets of posterior probability  $(1-\alpha)$  cover the parameter at confidence level  $(1-\alpha)$ .

\* Ex:  $X_1, \dots, X_n \sim B(\theta_0)$ .

Then  $I_1(\theta_0) = \frac{1}{\theta_0(1-\theta_0)}$  (see page 29) &  $\hat{\theta}_n = \bar{X}$

Frequentist interval is  $\bar{X} \pm z_{1-\alpha/2} n^{-1/2} [\theta_0(1-\theta_0)]^{-1/2}$

replace with  $\bar{X}$

With a conjugate Beta prior  $B(\alpha, \beta)$ , the posterior dist'n is  $B(\alpha + S_n, n + \beta - S_n)$  (see page 4), with mean  $\approx S_n/n$  and variance  $\approx n^{-1} \bar{X}(1-\bar{X})$  when  $n$  is large. A central interval of posterior probability  $(1-\alpha)$  therefore is of the same form as the frequentist interval. ■

• Sketch of proof = We provide here some heuristics explaining why the posterior distribution converges to the normal distribution as the sample size increases. A Taylor series expansion of the log posterior  $\log f(\theta | \mathcal{L}_n)$ , centered at the posterior mode yields:

$$\log f(\theta | \mathcal{L}_n) = \log f(\hat{\theta} | \mathcal{L}_n) + \frac{1}{2} (\theta - \hat{\theta})^t \left. \frac{\partial^2 \log f(\theta | \mathcal{L}_n)}{\partial \theta^2} \right|_{\theta = \hat{\theta}} + \dots$$

$$f(\theta | \mathcal{L}_n) \propto f(\theta) f(\mathcal{L}_n | \theta)$$

$$= \left. \frac{\partial^2 \log f(\theta)}{\partial \theta^2} \right|_{\theta = \hat{\theta}} + \sum_{i=1}^n \left. \frac{\partial^2 \log f(\mathcal{L}_n^i | \theta)}{\partial \theta^2} \right|_{\theta = \hat{\theta}}$$

constant as a function of  $\theta$

dominating term since the sum of  $n$  quantities

$$\frac{1}{n} \sum_{i=1}^n \left. \frac{\partial^2 \log f(\mathcal{L}_n^i | \theta)}{\partial \theta^2} \right|_{\theta = \hat{\theta}} = \text{empirical Fisher information}$$

for some  $\theta \in \Theta$

$$\xrightarrow{SLLN} \mathbb{E}_{\mathbb{P}_\theta} \left. \frac{\partial^2 \log f(X | \theta)}{\partial \theta^2} \right|_{\theta = \hat{\theta}} = -\mathbb{I}_1(\theta)$$

$$\Rightarrow \frac{1}{n} \log f(\theta | \mathcal{L}_n) \approx \text{constant} - \frac{1}{2} (\theta - \hat{\theta})^t [\mathbb{I}_1(\hat{\theta})]^{-1} (\theta - \hat{\theta}) + \text{something that goes to 0.}$$

We are being sloppy. Not a rigorous proof.

quadratic in  $(\theta - \hat{\theta})$

$$\Rightarrow f(\theta | \mathcal{L}_n) \approx \mathcal{N}(\theta | \hat{\theta}, [n \mathbb{I}_1(\hat{\theta})]^{-1}), \text{ as required.}$$

Under regularity conditions, the posterior mode  $\hat{\theta}$  is consistent for  $\theta_0$ , and so is the MLE  $\hat{\theta}_n$ . This ensures that  $f(\theta | \mathcal{L}_n) \approx \mathcal{N}(\theta | \hat{\theta}_n, [n \mathbb{I}_1(\theta_0)]^{-1})$ , as stated in the Bernstein-Von-Mises theorem.

x Counterexamples =

(a) Let  $X_1, \dots, X_n \sim \mathcal{N}(\theta, 1)$  (true distribution)

Data are modeled with  $\mathcal{N}(\theta, 1)$ , with  $\theta \geq 0$  constrained to be positive & non-informative prior on  $\theta$ . The true parameter thus lies on the edge of the parameter space  $\Theta = [0, +\infty)$ .

The posterior distribution for  $\theta$  is normal, centered at  $\bar{X}$ , and truncated to be positive. As  $n \rightarrow \infty$ , it converges to half of a normal distribution centered around  $\theta$  and truncated to be positive.

(b) No convergence to a single point. (= underidentified model)  
Consider the mixture distribution

$$f(x | \mu_1, \mu_2, \sigma_1^2, \sigma_2^2, \lambda) \leftarrow \theta = (\mu_1, \mu_2, \sigma_1^2, \sigma_2^2, \lambda) \in \mathbb{R}^5$$

$$= \lambda \frac{1}{\sigma_1 \sqrt{2\pi}} \exp\left(-\frac{1}{2\sigma_1^2} (x - \mu_1)^2\right) + (1 - \lambda) \frac{1}{\sigma_2 \sqrt{2\pi}} \exp\left(-\frac{1}{2\sigma_2^2} (x - \mu_2)^2\right)$$

The likelihood remains the same if we interchange  $(\mu_1, \mu_2)$ ,  $(\sigma_1^2, \sigma_2^2)$  and  $(\lambda, 1 - \lambda)$  so the posterior has at least two modes and cannot converge to a single point. A way around this is to impose constraints on the parameters.

x Bernstein-Van-Mises theorem under misspecification. (37)

Kleijn & van der Vaart (2012)

The case where the true generating distribution  $P_0$  is not a member of the parametric family  $f(\cdot | \theta)$  was considered by Kleijn and van der Vaart, extending the original Bernstein-Van-Mises theorem. In the authors' words, "this misspecified version of the Bernstein-Van-Mises theorem has important consequence for the interpretation of Bayesian credible sets. In the misspecified situation the posterior distribution of a parameter shrinks to the point within the model at minimum Kullback-Leibler divergence to the true distribution, a consistency property that it shares with the maximum likelihood estimator. Consequently one can consider both the Bayesian procedure and the MLE as estimates of this minimum KL point. A confidence region for this minimum KL point can be built around the MLE based on its asymptotic normal distribution, involving the sandwich covariance. One might also hope that a Bayesian credible set automatically yields a valid confidence set for the minimum KL point. However, the misspecified Bernstein-Van-Mises theorem shows the latter to be false.  $\rightarrow$

$$\text{Let } \theta^* = \underset{\theta}{\operatorname{argmin}} \operatorname{KL}(P_0 \| f(\cdot | \theta)) = \underset{\theta}{\operatorname{argmin}} \int P_0(x) \log \frac{P_0(x)}{f(x | \theta)} dx$$

Then, under some regularity conditions, (38)

$$\sup_B \left| \Pi(B | \mathcal{L}_n) - \mathcal{N}(\hat{\theta}_n, [n I_1(\theta^*)]^{-1})(B) \right| \rightarrow 0$$

in probability (w.r.t.  $P_0$ ), as  $n \rightarrow \infty$ , where

$$I_1(\theta) = - \mathbb{E}_{P_0} \left\{ \frac{\partial^2 \log f(X | \theta)}{\partial \theta^2} \right\}.$$

However, the MLE  $\hat{\theta}_n$  satisfies (van der Vaart (1998)):

$$n^{1/2}(\hat{\theta}_n - \theta^*) \xrightarrow{d} \mathcal{N}(0, \Sigma(\theta^*))$$

with  $\Sigma(\theta) = I_1^{-1}(\theta) W_1(\theta) I_1^{-1}(\theta)$  [aka the sandwich cov.]

$$W_1(\theta) = \mathbb{E}_{P_0} \left\{ \left( \frac{\partial \log f(X | \theta)}{\partial \theta} \right)^2 \right\}.$$

$\Rightarrow$  The asymptotic covariance matrix of the posterior distribution and of the MLE differ in the misspecified case.

Bayesian credible sets are not valid confidence sets when the model is misspecified —

• Example 1 (example 2.1 in Kleijn & van der Vaart (2012))

Let  $X_1, \dots, X_n$  iid  $\sim P_0 = \mathcal{N}(0, \sigma^2)$  = true distribution

Data are modeled using  $f(\cdot | \theta) = \mathcal{N}(\cdot | \theta, 1)$ , where  $\theta$  is the parameter of interest.

Note that  $KL(\varphi_0 \parallel f(\cdot|\theta))$

$$\text{constant} - \frac{1}{2} \int \varphi_0(x) (x-\theta)^2 dx,$$

which is minimized at  $\theta^* = \mathbb{E}_{P_0} X = 0$ .

In addition,  $\log f(x|\theta) = -\frac{1}{2} \log 2\pi - \frac{1}{2}(x-\theta)^2$

$$\frac{\partial \log f(x|\theta)}{\partial \theta} = x - \theta$$

$$\frac{\partial^2 \log f(x|\theta)}{\partial \theta^2} = -1$$

We get that  $I_1(\theta^*) = 1$

$$W_1(\theta^*) = \mathbb{E}_{P_0} X^2 = \sigma^2,$$

so that  $\Sigma(\theta^*) = I_1^{-1}(\theta^*) W_1(\theta^*) I_1^{-1}(\theta^*) = \sigma^2$ .

and the MLE satisfies  $n^{-1/2} \hat{\theta}_n \xrightarrow{d} \mathcal{N}(0, \sigma^2)$ .

$\Rightarrow$  The  $(1-\alpha)$  confidence set around the mean has radius  $z_{1-\frac{\alpha}{2}} \sigma n^{-1/2}$

On the other hand, according to the misspecified Bernstein-Van-Mises theorem, the posterior distribution yields an  $(1-\alpha)$  Bayesian credible set with radius  $z_{1-\frac{\alpha}{2}} n^{-1/2}$ , since  $f(\cdot|\hat{\alpha}_n) \approx \mathcal{N}(\hat{\theta}_n, 1)$ . The Bayesian credible set thus has a frequentist coverage that can be arbitrarily close to 0 or 1 (depending on  $\sigma^2 < 1$  or  $\sigma^2 > 1$ ).

• Example 2: taken from White (1982), Maximum Likelihood Estimation of Misspecified Models, Econometrica, vol 50, n° 1, pp 1-25.

Let  $X_1, \dots, X_n$  iid  $\sim P_0$  with density  $\varphi_0$  (true distribution)

$$\begin{aligned} \text{Put } \mu_0 &= \mathbb{E}_{P_0} X & \gamma &= \frac{\mathbb{E}_{P_0} (X - \mu_0)^3}{\sigma_0^3} \quad (\text{skewness}) \\ \sigma_0^2 &= \text{Var } X & \gamma &= \frac{1}{\sigma_0^4} \mathbb{E}_{P_0} (X - \mu_0)^4 \quad (\text{kurtosis}) \end{aligned}$$

Data are modeled using the normal distribution  $f(\cdot|\theta) = \mathcal{N}(\cdot | \mu, \sigma^2)$  and the goal is to estimate  $\mu$  and  $\sigma^2$ .  $\theta = (\mu, \sigma^2)$

It is straightforward to see that

$$\begin{aligned} \theta^* &= (\mu^*, (\sigma^*)^2) \\ &= \underset{(\mu, \sigma^2)}{\text{argmin}} KL(\varphi_0 \parallel f(\cdot|\theta)) = (\mu_0, \sigma_0^2) \end{aligned}$$

Indeed,

$$\begin{aligned} KL(\varphi_0 \parallel f(\cdot|\theta)) &= \int \varphi_0(x) \log \frac{\varphi_0(x)}{f(x|\theta)} dx \\ &= \text{constant} - \log \sigma - \frac{1}{2\sigma^2} \int \varphi_0(x) (x-\mu)^2 dx \\ &= \text{constant} - \log \sigma - \frac{1}{2\sigma^2} \left( \mathbb{E}_{P_0} X^2 - 2\mu \mathbb{E}_{P_0} X + \mu^2 \right) \\ &\text{minimized at } \mu^* = \mu_0 (= \mathbb{E}_{P_0} X). \end{aligned}$$

At  $\mu^* = \mu_0$ ,  $KL(\varphi_0 \| f(\cdot | \theta)) = \text{cst} - \log \sigma - \frac{\sigma_0^2}{2\sigma^2}$ , (41)

which is minimized at  $\sigma^* = \sigma_0$ .

With  $\log f(X|\theta) = -\frac{1}{2} \log 2\pi - \frac{1}{2} \log \sigma^2 - \frac{1}{2\sigma^2} (X-\mu)^2$

•  $\frac{\partial}{\partial \mu} \log f(X|\theta) = \frac{1}{\sigma^2} (X-\mu)$

•  $\frac{\partial}{\partial \sigma^2} \log f(X|\theta) = -\frac{1}{2\sigma^2} + \frac{1}{2\sigma^4} (X-\mu)^2$

•  $\frac{\partial^2}{\partial \mu \partial \sigma^2} \log f(X|\theta) = -\frac{X-\mu}{\sigma^4}$

•  $\frac{\partial^2}{\partial \mu^2} \log f(X|\theta) = -\frac{1}{\sigma^2}$

•  $\frac{\partial^2}{(\sigma^2)^2} \log f(X|\theta) = \frac{1}{2\sigma^4} - \frac{1}{\sigma^6} (X-\mu)^2$ ,

we get

$\rightarrow I_1(\theta^*) = \begin{bmatrix} \frac{1}{\sigma_0^2} & 0 \\ 0 & \frac{1}{2\sigma_0^4} \end{bmatrix}$

$\rightarrow W_1(\theta^*) = \begin{bmatrix} 1/\sigma_0^2 & \sqrt{\beta}/2\sigma_0^3 \\ \sqrt{\beta}/2\sigma_0^3 & (\gamma-1)/4\sigma_0^4 \end{bmatrix}$

$\rightarrow \Sigma(\theta^*) := I_1^{-1}(\theta^*) W_1(\theta^*) I_1^{-1}(\theta^*) = \begin{bmatrix} \sigma_0^2 & \sqrt{\beta} \sigma_0^3 \\ \sqrt{\beta} \sigma_0^3 & (\gamma-1) \sigma_0^4 \end{bmatrix}$

Unless  $\beta=0$  &  $\gamma=3$  (ie.  $\varphi_0 = \mathcal{N}$ ),  $\Sigma(\theta^*) \neq I_1^{-1}(\theta^*)$  indeed.

### III.2. Hypothesis testing:

#### Fisher vs Neyman vs Jeffreys

Suppose that  $X \sim f(x|\theta)$ , and we wish to test for  $\theta = \theta_0$ . Fisher, Neyman and Jeffreys would all disagree on how to proceed

#### Fisher's significance testing =

$H_0: \theta = \theta_0$

Select a statistic  $T(X)$  so that large values of  $T$  reflect evidence against  $H_0$

Compute the p-value for the observed data  $x$

$p = \mathbb{P}_{H_0}(T(X) \geq T(x))$   
and reject  $H_0$  if  $p$  is small, e.g.  $p \leq 0.05$

#### Neyman-Pearson hypothesis testing

Test  $H_0: \theta = \theta_0$  vs an alternative  $H_1: \theta = \theta_1$ .

Reject  $H_0$  if  $T > c$ , for a pre-selected critical value  $c$

Compute  $\alpha = \mathbb{P}_{H_0}(\text{rejecting } H_0)$

$\beta = \mathbb{P}_{H_1}(\text{"accepting" } H_0)$

#### Jeffrey's approach

Define Bayes Factor

$B(x) = \frac{f(x|H_0)}{f(x|H_1)}$

Reject  $H_0$  if  $B(x) < 1$ , otherwise "accept" it.

Repeat the posterior proba

$P(H_0|x) = \frac{B(x)}{1+B(x)}$

Fisher's p-values are often poorly interpreted. A p-value less than 0.05 indicates that the chances are less than 5% of having obtained the observed response or any more extreme response if  $H_0$  is true. (43)

p-values \* are not \* an error probability for rejecting  $H_0$   
 p-values \* are not \* the posterior probability that  $H_0$  is true

— p-values are usually misinterpreted —

We present next the derivation of Bayes factor in some simple cases, and highlight severe discrepancies between the Bayesian and frequentist approaches to testing that arise from it (the Jeffrey-Lindley & Bartlett paradoxes). We then show that transformations of Fisher's p-value lead to its interpretation as odds or posterior probability of the hypothesis (Sellke, Bayarri & Berger (2001)). We conclude this section discussing multiple and sequential testing.

### • Bayes Factor

Let  $X|\theta \sim f(x|\theta) \rightarrow$  sample  $\mathcal{L}_n = \{X_1, \dots, X_n\}$

We test  $H_0: \theta \in \Theta_0$

vs  $H_1: \theta \in \Theta_1$

with prior distributions

→ prior probabilities  $P(H_0)$ ,  $P(H_1)$  of the hypothesis

→ prior densities  $f_0(\theta)$  and  $f_1(\theta)$  on  $\Theta_0$  and  $\Theta_1$ .

(can reduce to a point mass)

The marginal likelihoods are  $f(\mathcal{L}_n | H_i) = \int f(\mathcal{L}_n | \theta, H_i) f_i(\theta) d\theta$

The posterior probabilities of the hypotheses are = (44)

$$P(H_0 | \mathcal{L}_n) = \frac{f(\mathcal{L}_n | H_0) P(H_0)}{f(\mathcal{L}_n | H_0) P(H_0) + f(\mathcal{L}_n | H_1) P(H_1)}$$

$$= 1 - P(H_1 | \mathcal{L}_n)$$

Then

$$\frac{P(H_0 | \mathcal{L}_n)}{P(H_1 | \mathcal{L}_n)} = \frac{P(H_0)}{P(H_1)} \times \frac{f(\mathcal{L}_n | H_0)}{f(\mathcal{L}_n | H_1)}$$

posterior odds      prior odds      Bayes Factor  $B_{01}$

$$\Rightarrow P(H_0 | \mathcal{L}_n) = \left( 1 + \frac{P(H_1)}{P(H_0)} \frac{1}{B_{01}} \right)^{-1}$$

Conclusions are drawn based (1) on the posterior odds ( $H_0$  is rejected if  $P(H_0 | \mathcal{L}_n) < P(H_1 | \mathcal{L}_n)$ ) or (2) on Bayes Factor  $B_{01}$ , with the default  $P(H_0) = P(H_1)$ .

Jeffreys (1961) suggested the scale:

$B_{01}$	Strength of evidence
1 to 3	Barely worth mentioning
3 to 10	Substantial
10 to 30	Strong
30 to 100	Very strong
> 100	Decisive

We derive next the expression of Bayes factor in the simplest case of testing for the mean of a normal population.

• Testing for the mean of a normal population.

↳ Case 1:  $X_1, \dots, X_n \sim \mathcal{N}(\mu, \sigma^2)$ ,  $\sigma^2$  known.

Test for  $\begin{cases} H_0: \mu = 0 \text{ (w.p. 1)} \\ H_1: \mu \sim \mathcal{N}(0, \frac{\sigma^2}{m}) = f_1(\mu) \end{cases}$

$m$  controls the width of the prior  $\equiv$  sample size

We compute  $B_{01} = \frac{f(\mathcal{L}_n | H_0)}{f(\mathcal{L}_n | H_1)}$

•  $f(\mathcal{L}_n | H_0) = \frac{1}{(2\pi\sigma^2)^{n/2}} \exp\left\{-\frac{1}{2\sigma^2} \sum_{i=1}^n X_i^2\right\}$

•  $f(\mathcal{L}_n | H_1) = \int f(\mathcal{L}_n | \mu, H_1) f_1(\mu) d\mu$   
 $= \int \frac{1}{(2\pi\sigma^2)^{n/2}} \exp\left\{-\frac{1}{2\sigma^2} \sum_{i=1}^n (X_i - \mu)^2\right\}$   
 $\times \frac{1}{(2\pi \frac{\sigma^2}{m})^{1/2}} \exp\left\{-\frac{m}{2\sigma^2} \mu^2\right\} d\mu$   
 $= \frac{1}{(2\pi\sigma^2)^{n/2}} \frac{1}{(2\pi \frac{\sigma^2}{m})^{1/2}} \int \exp\left\{-\frac{1}{2\sigma^2} \left[\sum (X_i - \mu)^2 + m\mu^2\right]\right\} d\mu$

$(n+m)\mu^2 - 2\mu n\bar{X} + n\bar{X}^2$

$(n+m)(\mu - \mu_n)^2 - \frac{(n\bar{X})^2}{n+m} + n\bar{X}^2$

where  $\mu_n = \frac{n}{n+m} \bar{X}$ .

↳ The integral is

$\left(\frac{2\pi\sigma^2}{n+m}\right)^{1/2} \exp\left\{-\frac{1}{2\sigma^2} \left[-\frac{(n\bar{X})^2}{n+m} + n\bar{X}^2\right]\right\}$

and

$f(\mathcal{L}_n | H_1) = \frac{1}{(2\pi\sigma^2)^{n/2}} \frac{1}{(2\pi \frac{\sigma^2}{m})^{1/2}} \left(2\pi \frac{\sigma^2}{n+m}\right)^{1/2} \exp\left\{-\frac{1}{2\sigma^2} [\dots]\right\}$

After simplifications,

$B_{01} = \frac{f(\mathcal{L}_n | H_0)}{f(\mathcal{L}_n | H_1)} = \left(\frac{n+m}{m}\right)^{1/2} \exp\left(-\frac{1}{2} \frac{n}{n+m} z^2\right)$ ,  $z = \frac{\bar{X}}{\sigma/\sqrt{n}}$

= function of the prior sample size, the data sample size, and the frequentist z-score  $z$ .

Our intuition is that when the z-score  $z$  is large; i.e. when  $\bar{X}$  is away from the tested value 0,  $H_1$  should be favored over  $H_0$  and  $B_{01}$  should be small.

↳ in frequentist terms, a large z-score would yield a small p-value & we would reject  $H_0$ .

However =

↳ Jeffreys-Lindley paradox = for large  $n$ ,

$B_{01} \sim \sqrt{n} \exp\left(-\frac{1}{2} z^2\right)$ , so that a classical frequentist test can strongly reject the null (which happens with a large  $z$ ), while a Bayesian analysis can strongly support the null (if  $B_{01} \sim \sqrt{n} \exp\left(-\frac{1}{2} z^2\right)$  is large).

↳ Bartlett paradox = Recall that under  $H_1$ ,  $\mu \sim \mathcal{N}(0, \frac{\sigma^2}{m})$ .

As  $\frac{\sigma^2}{m} \rightarrow 0 \iff m \rightarrow \infty$  for  $\sigma^2$  fixed, a non-informative flat prior for  $\mu$  is used. However,  $B_{01} \rightarrow \infty$  as  $m \rightarrow \infty$ , irrespectively how large the z-score is. A Bayesian analysis systematically picks  $H_0$  over  $H_1$ .

Case 2:  $X_1, \dots, X_n \sim \mathcal{N}(\mu, \sigma^2)$ ,  $\sigma^2$  unknown (nuisance) (47)

Test for  $\begin{cases} H_0: \mu = 0 \text{ w.p.1} \\ H_1: \mu | \sigma \sim \mathcal{N}(0, \frac{\sigma^2}{m}) \end{cases}$   $\sigma \sim \frac{1}{\sigma}$  (Jeffrey's suggestion)

Same prior as before, conditionally on  $\sigma$

A common improper distribution is used for the common nuisance parameter  $\sigma$

We compute  $B_{01} = \frac{f(\mathcal{L}_n | H_0)}{f(\mathcal{L}_n | H_1)}$

$f(\mathcal{L}_n | H_0) = \int_0^\infty f(\mathcal{L}_n | \mu=0, \sigma, H_0) f(\sigma | H_0) d\sigma$   
 $\sim (1/\sigma)^n \exp(-\frac{1}{2\sigma^2} \sum X_i^2) \quad 1/\sigma$

we omit the term  $(1/2\pi)^{n/2}$  which will cancel out with the denominator in  $B_{01}$ .

$x = \sigma^2$   $\int_0^\infty (\frac{1}{\sigma})^{n+1} \exp(-\frac{1}{2\sigma^2} T) d\sigma$   $T = n(S^2 + \bar{X}^2)$   
 $S^2 = \frac{1}{n} \sum_{i=1}^n (X_i - \bar{X})^2$

$= \int_0^\infty (\frac{1}{x})^{\frac{n+1}{2}} \exp(-\frac{T}{2x}) \frac{dx}{2\sqrt{x}}$

$= \frac{1}{2} \int_0^\infty (\frac{1}{x})^{\frac{n}{2}+1} \exp(-\frac{T}{2x}) dx$

$= \frac{1}{2} \Gamma(\frac{n}{2}) (\frac{T}{2})^{-n/2}$

inverse gamma distribution  $\frac{\beta^\alpha}{\Gamma(\alpha)} x^{-\alpha-1} e^{-\beta/x}$

$f(\mathcal{L}_n | H_1) = \iint \left( \frac{1}{\sigma} \right)^n \exp(-\frac{1}{2\sigma^2} \underbrace{\sum (X_i - \mu)^2}_{n(S^2 + (\bar{X} - \mu)^2)}) \mathcal{N}(\mu | 0, \frac{\sigma^2}{m}) \frac{1}{\sigma} d\mu d\sigma$

$= \int \left\{ \int \exp(-\frac{n(\bar{X} - \mu)^2}{2\sigma^2}) \mathcal{N}(\mu | 0, \frac{\sigma^2}{m}) d\mu \right\} \left( \frac{1}{\sigma} \right)^{n+1} \times \exp(-\frac{nS^2}{2\sigma^2}) d\sigma$  (48)

Let  $T_2 := \frac{nm}{n+m} \bar{X}^2 + nS^2$

$= \left( \frac{m}{n+m} \right)^{1/2} \int_0^\infty \left( \frac{1}{\sigma} \right)^{n+1} \exp(-\frac{T_2}{2\sigma^2}) d\sigma$

$= \frac{1}{2} \left( \frac{m}{n+m} \right)^{1/2} \Gamma(\frac{n}{2}) \left( \frac{T_2}{2} \right)^{-n/2}$

It follows that

$B_{01} = \frac{f(\mathcal{L}_n | H_0)}{f(\mathcal{L}_n | H_1)} = \left( \frac{n+m}{m} \right)^{1/2} \left( \frac{n + \frac{m\bar{X}^2}{n+m}}{n + T_2} \right)^{n/2}, \quad t = \frac{\bar{X}}{S/\sqrt{n}}$

$t$  plays the same role as  $z$  before, replacing  $\sigma^2$  with its estimate  $s^2$ .

$t$  is the frequentist  $t$ -statistic. A large value of  $t$  yields a small  $p$ -value, and the rejection of the null hypothesis.

• Jeffreys-Lindley paradox shows up again: as  $n \rightarrow \infty$ ,  $B_{01} \sim \sqrt{n} \Rightarrow$  a Bayesian analysis favor  $H_0$  over  $H_1$ , irrespectively how large  $t$  may be.

• Bartlett paradox: A flat prior on  $\mu$  is inadequate here as well, since as  $m \rightarrow 0$ ,  $B_{01} \rightarrow \infty$  and Bayes factor picks  $H_0$  over  $H_1$ , all the time.



In addition to these paradoxes, note that as  $t^2 \rightarrow \infty$ , (49) Bayes Factor  $B_{01}$  tends to a constant value ( $n, m$  fixed), which is also counterintuitive. As  $t^2 \rightarrow \infty$ ,  $\bar{X}$  differs more and more from 0 and one would expect  $B_{01}$  to converge to 0.

One possible solution is to use a Cauchy prior on  $\mu | \sigma$  instead of the conjugate Gaussian prior, as suggested by Jeffreys. While the numerator in  $B_{01}$  is still available in close form, the denominator must be computed numerically. Jeffreys gave a (bad) approximation to  $B_{01}$  in this case,

$B_{01} \approx \sqrt{\frac{\pi n}{2}} \exp\left(-\frac{t^2}{2}\right)$ , which is sufficient to show that  $B_{01} \rightarrow 0$  as  $t^2 \rightarrow \infty$ , while Lindley's paradox remains since  $B_{01} \rightarrow \infty$  as  $n \rightarrow \infty$ .

\* Remark: Jeffreys recommend to use improper priors for common parameters to  $H_0$  and  $H_1$  (such as  $\sigma^2$ ). Kass & Vaidyanathan (1992) showed that Bayes Factors are  $\approx$  insensitive to the choice of a common prior under weaker assumptions than orthogonality (i.e. diagonal Fisher info matrix). Jeffreys advocate the use of proper priors for non-common parameters between  $H_0$  and  $H_1$ . The prior should be centered & symmetric about 0, has scale  $\sigma$ , and no moments (like Cauchy)  $\blacksquare$

\* Remark = To further support Jeffreys recommendation not to use vague or improper priors for common parameters, consider again

the case  $X_1, \dots, X_n \sim \mathcal{N}(\mu, \sigma^2)$ ,  $\sigma^2$  known, (50) where we test for  $H_0: \mu = 0$  (w.p.1).

• vague prior for  $H_1: \mu \sim \mathcal{U}(-c, c)$

$$B_{01} = \frac{f(\mathcal{L}_n | H_0)}{\int_{-c}^c f(\mathcal{L}_n | \mu, H_1) \frac{1}{2c} d\mu} \approx \frac{2c f(\mathcal{L}_n | H_0)}{\int_{-\infty}^{+\infty} f(\mathcal{L}_n | \mu, H_1) d\mu} = 2c f(\mathcal{L}_n | H_0)$$

$\swarrow$   
c large enough

The value of  $B_{01}$  depends heavily on the choice of  $c$ .

• improper prior for  $H_1: \mu \sim 1$  or  $\mu \sim 2$  yields Bayes Factors  $f(\mathcal{L}_n | H_0)$  and  $2f(\mathcal{L}_n | H_0)$ , respectively.  
 $\Rightarrow$  Which normalization to use?  $\blacksquare$

Back to Jeffreys-Lindley's paradox.

$\hookrightarrow$  denotes a situation where a frequentist test would strongly reject the null, while  $B_{01}$  would favor heavily the null over the alternative.

This situation is likely to occur when

- (i) the sample size  $n$  is large
- (ii)  $H_0$  is relatively precise (e.g.  $\mu = 0$  w.p.1), while the alternative is relatively diffuse (vague, improper, but not only).

\* How to make sense of this paradox?

In large sample situations, the Bayesian analysis requires more evidence to detect an effect (reject  $H_0$ );

while the frequentist approach tends to reject too easily the null

(51)

One way around this is to adjust the level  $\alpha$  of the test as  $n$  grows, instead of keeping it constant.

Indeed, as  $n$  grows,  $\bar{x}$  concentrates around its mean. This allows us to move from confidence intervals with fixed coverage  $1-\alpha =$

$$\left[ \bar{x} - z_{\frac{1-\alpha}{2}} \sigma n^{-1/2}, \bar{x} + z_{\frac{1-\alpha}{2}} \sigma n^{-1/2} \right] \rightarrow \text{coverage } 1-\alpha$$

to intervals of the form =

$$\left[ \bar{x} - z_n \sigma n^{-1/2}, \bar{x} + z_n \sigma n^{-1/2} \right].$$

coverage  $1-h_n \rightarrow 1$  as  $n \rightarrow \infty$ , provided  $h_n \rightarrow 0$ , and  $z_n = \Phi^{-1}(1-h_n)$  such that  $z_n n^{-1/2} \rightarrow 0$  (otherwise the confidence interval blows up). Take e.g.  $h_n \sim n^{-p}$ ,  $p > 1$ , see Naaman (2016)

The last written interval has coverage  $1-h_n > 1-\alpha \Rightarrow$  it includes the value 0 more often than the classic interval, and therefore the null is "accepted" more often.

Since  $\bar{x}$  concentrate around its mean as  $n$  grows, the alternative prior under  $H_1$  in the Bayesian setting consider a wide range of alternatives that are meaningless to describe the data.

$B_{01}$  tends to pick up  $H_0$  even if it poorly describes the data. This is even more problematic with diffuse priors.

(52)

The prior is treated very differently in Bayesian posterior inference & Bayesian testing:

$$f(\mu | \mathcal{Z}_n) \propto f(\mathcal{Z}_n | \mu) f(\mu) \quad \text{vs} \quad \int f(\mathcal{Z}_n | \mu) f(\mu) d\mu$$

"update belief about the prior"      "average over all possible values"

\* Summary = We considered so far Bayesian tests where

$\rightarrow$  A point null is used  $H_0: \theta = 0$  w.p.1

$\rightarrow$  An alternative whose prior puts positive mass on  $\theta_0 = \{0\}$   $H_1: \theta \sim \mathcal{N}(0, \frac{\sigma^2}{m})$  or Cauchy.

$\rightarrow$  As noted by Johnson & Rossell (2010), this is a problematic assumption since such alternative priors "do not incorporate any notion of a minimally significant separation between the null and the alternative hypotheses". The authors show that

recall the expression of  $B_{01}$  in p.46

"For a true null hypothesis, the Bayes factor in favor of the alternative hypothesis decreases only at rate  $O_p(n^{-1/2})$ " while

"For a true alternative hypothesis, the Bayes factor in favor of the null hypothesis decreases exponentially fast"

When  $H_0$  is true, evidence for the null accumulates slowly, while when the alternative is true, evidence accumulates exp. fast.

⇒ We can either fix the null, or the alternative & reconcile the Bayesian approach with the frequentist one.

• Johnson & Rossell (2010) consider a class of "non-local alternative priors", that put zero prior mass for all values  $\theta \in \Theta_0$ , and a strictly positive mass on  $\Theta_1$ . "Although this class of prior densities does not provide exponential convergence of Bayes factors in favor of true null hypothesis, it offers a substantial improvement in convergence rates over local alternative priors."

• Casella & Berger (1987) consider non-point nulls  $H_0: \theta \leq 0$  vs  $H_1: \theta > 0$  and show closer agreement to the frequentist approach. They prove that for many classes of prior distributions, the infimum of the Bayesian posterior probability of  $H_0$  is either equal to or bounded above by the p-value.

• Calibration of p-values.

Since frequentist p-values are often misinterpreted, some attempts to transform them so that they are interpretable as odds ratio or posterior probability of the hypothesis have been considered in the literature. Following Sellke, Bayarri & Berger (2001),

under  $H_0: p \sim U(0,1)$   $p = p\text{-value}$   
(see e.g. p.26 in MS = HYPOTHESIS TESTING)

Test  $H_0$  versus  $H_1: p \sim f(p|\alpha) = B(p|\alpha, \beta=1) = \alpha p^{\alpha-1}$ ,  
with some prior  $\pi(\alpha)$  on  $\alpha$ .

Then  $B_{01}(\pi) = \frac{f(p|\alpha=1)}{\int_0^1 f(p|\alpha) \pi(\alpha) d\alpha}$  ↖ since  $U(0,1) = B(1,1)$

One can show that

$$B = \inf_{\pi} B_{01}(\pi) = \begin{cases} -e p \log p & \text{if } p < e^{-1} \\ 1 & \text{o/w} \end{cases}$$

⇒ Calibration  $-e p \log p$ , where  $p$  is the frequentist p-value, can be taken as an objective lower bound for any prior  $\pi$  on  $\alpha$  on the odds of  $H_0$  to  $H_1$ .

• Comparing two means.

Suppose that  $X_1, \dots, X_{n_0} \sim \mathcal{N}(\mu_X, \sigma^2)$  ↖  $\sigma^2$  assumed equal  
 $Y_1, \dots, Y_{n_1} \sim \mathcal{N}(\mu_Y, \sigma^2)$ .

We wish to test for  $H_0: \mu_X = \mu_Y$  vs the alternative  $\mu_X \neq \mu_Y$ .

To proceed, it is convenient to reparametrize  $\mu_X$  and  $\mu_Y$  using the overall mean  $\mu = \frac{1}{2}(\mu_X + \mu_Y)$ , and  $\delta = \mu_X - \mu_Y$ :

$$X \sim \mathcal{N}(\mu + \frac{\delta}{2}, \sigma^2)$$

$$Y \sim \mathcal{N}(\mu - \frac{\delta}{2}, \sigma^2),$$

and testing for  $\delta = 0$  vs  $\delta \neq 0$ .  
 $(H_0) \quad (H_1)$

$\mu, \sigma^2 =$  common parameters for  $H_0$  and  $H_1 \rightarrow$  ok to use a non-informative prior  $f(\mu, \sigma^2 | H_i) \propto \frac{1}{\sigma^2} \quad i=0,1$ .

$\delta =$  diffuse priors are not recommended.

Consider first a conjugate normal prior  $\frac{\delta}{\sigma} \sim \mathcal{N}(0, \sigma_0^2)$ .

Then  $B_{01} = \frac{f(\mathcal{L}_n | H_0)}{f(\mathcal{L}_n | H_1)} = \frac{\iint f(\mathcal{L}_n | \delta=0, \mu, \sigma^2) f(\mu, \sigma^2 | H_0) d\mu d\sigma^2}{\iint f(\mathcal{L}_n | \delta, \mu, \sigma^2) f(\delta | \sigma^2) f(\mu, \sigma^2 | H_1) d\delta d\mu d\sigma^2}$

$B_{01}$  can be calculated exactly in this case, and is expressed in terms of the frequentist pooled-variance two-sample  $t$  statistic,

$$BF_{01} = \left( \frac{1 + t^2/v}{1 + t^2/[v(1+n\sigma_0^2)]} \right)^{-\frac{v+1}{2}} (1 + n\sigma_0^2)^{1/2},$$

where  $t = \frac{\bar{X} - \bar{Y}}{S_p/\sqrt{n}}$ ,  $n = \left( \frac{1}{n_0} + \frac{1}{n_1} \right)^{-1}$

$$S_p = \frac{1}{n_0 + n_1 - 2} \left( (n_0 - 1)S_0^2 + (n_1 - 1)S_1^2 \right)$$

$$v = n_0 + n_1 - 2,$$

see Gonen, Johnson, Lu & Westfall (2005).

↑ This approach shows that  $t$  arises from a Bayesian formulation, but suffers from the same paradoxes as in the single population problem.

- Jeffreys-Lindley: For large  $n$ ,  $B_{01} \sim \sqrt{n} e^{-t^2/2}$ , so a strong frequentist effect (large value of  $t$ ) may lead to strong evidence for  $H_0$ .
- As  $t \rightarrow \infty$ ,  $BF_{01}$  tends to a constant, while we would expect  $B_{01} \rightarrow 0$  to be in agreement with the frequentist approach.

↳ Use a Cauchy prior instead ( $B_{01}$  must be computed numerically)

• Testing based on the posterior distribution.

A more traditional approach can be used for testing, based directly on the posterior distribution directly.

Consider the problem of comparing two binomial proportions,

revisiting the situation on p.7/8:

$$\begin{aligned} X &\sim B(p_0) & p_0 &\sim B(\alpha_0, \beta_0) & p_0 | X &\sim B(\alpha_0 + S_{n_0}, n_0 + \beta_0 - S_{n_0}) \\ Y &\sim B(p_1) & p_1 &\sim B(\alpha_1, \beta_1) & p_1 | Y &\sim B(\alpha_1 + S_{n_1}, n_1 + \beta_1 - S_{n_1}) \end{aligned}$$

From these, the posterior distribution of the difference  $\delta := p_0 - p_1$  can easily be computed (numerically).

To decide which one of  $H_0: \delta \leq 0$  or  $H_1: \delta > 0$  is more appropriate, we may reject  $H_0$  if the 95% posterior interval is fully contained in the positive half space. Other similar criteria may be considered.

This approach is similar in spirit to the frequentist method which rejects  $H_0$  at 95% confidence level if the one-sided confidence interval excludes 0 ( $\Leftrightarrow$  p-value less than 0.05). The two approaches coincide when the posterior credible interval has a frequentist interpretation (or approximately in large sample situations via the Bernstein-von Mises theorem)

A more Bayesian treatment introduces a loss, characterizing the cost associated with making a mistake. Let

$$\begin{aligned} l_0(p_0, p_1) &:= \max(p_1 - p_0, 0) \\ l_1(p_0, p_1) &:= \max(p_0 - p_1, 0) \end{aligned}$$

If  $H_0$  is true,  $p_0 \leq p_1$ , and  $p_1 - p_0$  is the price to pay to draw erroneous conclusions (linear in  $p_1 - p_0$ ). The cost is 0 if the right decision is made.

The expected loss is

$$\int_0^1 \int_0^1 l_i(p_0, p_1) f(p_0, p_1 | X \cup Y) dp_0 dp_1, \quad i=0,1.$$

↖ can easily be calculated numerically

In practice, we select a threshold  $\varepsilon > 0$  and stop the experiment when the expected loss for  $i=0$  or  $i=1$  falls below  $\varepsilon$ .

↖ For example, suppose that historical data indicate that  $p_0 \approx 0.20$ . If we want to detect a relative effect of about 2% improvement from this base rate, leading to  $(p_1 - p_0) / p_0 = 2\%$ ,  $p_1 = 0.204$ , we set  $\varepsilon = 0.004$ .

x Remark = Superiority of BF vs posterior-based testing for a particular application should be established (at least numerically via simulations) by controlling for the type I and type II error rates.

### III.3. Sequential Testing

In sequential testing, samples are collected until enough evidence for  $H_0$  or  $H_1$  is obtained, and a decision is made.

In frequentist statistics, repeated peeks at the data is known to inflate the type I error. Boundaries for optimal stopping must be updated accordingly, see e.g. p.37-47 in MS = HYPOTHESIS TESTING.

We explore here how

Bayesian techniques for sequential testing compare, and discuss testing based on

- the posterior distribution
- Bayes Factors.

(57)

### (i) Posterior distribution for sequential testing.

(58)

Consider the simple case where  $X_1, \dots, X_n$  are received sequentially &  $N(\mu, \sigma^2)$  distributed, testing for  $\mu$  ( $H_0: \mu = 0$  vs  $H_1: \mu \neq 0$ ),  $\sigma^2$  unknown.

With a non-informative prior  $f(\mu, \sigma^2) \sim \frac{1}{\sigma^2}$ , the posterior distribution of  $\mu | \mathcal{X}_n$  is  $t(n-1, \bar{x}, \frac{1}{n} s^2)$ , see page 21.

In this particular case, credible bounds & frequentist bounds exactly coincide, and

95% credible bound excludes 0

⇔ 95% confidence interval excludes 0

⇔ p-value < 0.05

↖ of  $z = \frac{\bar{x}}{s/\sqrt{n}} \sim t(n-1)$

⇒ Repeated peeks at the data will inflate the type I error with this procedure. With an informative prior, credible & confidence intervals do not match exactly, and the equivalence does not hold. However, in large sample situations, Bernstein von Mises ensures that the posterior distribution has a frequentist interpretation, and inflation of the type I error follows. In fact, even in the non-asymptotic case, the type I error rate needs to be controlled, see Shi & Yin (2019). The authors present a procedure for adjusting the bounds, following the work of Pocock & O'Brien in the frequentist case (see MS = HYPOTHESIS TESTING).

What makes this approach prone to type I error inflation (59) is the vanishing effect on the prior, resulting in the posterior being "too frequentist", therefore inheriting issues associated with sequential testing

$$\text{posterior} \propto \text{likelihood} \times \text{prior}$$

The prior is treated very differently in Bayes Factors. Its effect is persistent as new data are collected, since the likelihood is integrated with respect to the prior distribution:

$$\int \text{likelihood} \times \text{prior} \quad \leftarrow \text{no feedback loop}$$

Although this is a limitation of BF to some extent (the prior has a huge effect on the value of BF & on the decision), it turns out to be a nice feature for sequential testing, being more robust to the inflation of the type I error.

(iii) Bayes Factors for sequential testing.

"The rules of governing when data collection stops are irrelevant to data interpretation. It is entirely appropriate to collect data until a point has been proven or disproven, or until the data collector runs out of time, money, or patience" *Edwards, Lindman & Savage (1963)*

Suppose that  $X_1, \dots, X_n \sim f_0$  under  $H_0$   
 $\sim f_1$  under  $H_1$   
 (with appropriate priors on the parameters)

$$BF_{01} = \frac{f(\mathcal{L}_n | H_0)}{f(\mathcal{L}_n | H_1)}$$

Consider the following procedure, described in (60) *Schönbrodt, Wagenmakers, Zehetleitner & Perugini (2015)*

- (i) Select a threshold indicating the decisiveness of evidence  
 Ex:  $B_{01}$  of 10 for  $H_0$  & the reciprocal value of  $1/10$  for  $H_1$ .
- (ii) Collect a minimum of  $n_{\min}$  observations before stopping the experiment; compute  $BF_{01}$  after each observation is collected.  
 (to avoid misleading evidence leading to early termination)
- (iii) When  $BF_{01}$  crosses one of the boundaries, stop the experiment & report the final  $BF_{01}$  & mean and 95% credible bounds.  
 posterior

x Remarks: (i) As  $n \rightarrow \infty$ ,  $BF_{01}$  is consistent (it converges to  $\infty$  if  $H_0$  is true or to 0 if  $H_1$  is true) and the procedure is guaranteed to converge

(ii) The false positive rate ( $\alpha$ ) is not controlled by the experimenter (a major drawback of this approach, compared to the frequentist procedure). The FPR is a function of the priors, the chosen threshold, and the true effect size in the population. It cannot be computed in advance, but simulations can be performed to get an idea of its magnitude. The FNR magnitude can also be estimated via simulations.

(iii) The procedure suffers from the same drawbacks as in the non-sequential case: evidence accumulates faster when a true effect is present than when there is none (recommended to "separate" the null from  $H_1$ , see page 52).

(iv) Simulations indicate that sample size requirements are much lower here than in the frequentist case (e.g. when calculated from a usual power study).

⇒ less data are needed to achieve better performance (FP and FN remain small in %, well below 5% for  $\alpha$ , and well below 5% for  $\beta$ , even for relatively small true effect sizes)

↳ Simulations performed by the authors in the context of comparison of two means.

(v) Downwards bias reported in the point estimates. The procedure tends to underestimate the true effect.

(vi)  $B_{F_0}$  provides evidence \* for \*  $H_0$ , and not solely against  $H_0$ , as happens in the frequentist approach.

(vii) **Rauder (2014)** shows that in simple situations at least,  $B_{F_0}$  keeps its interpretation whether optimal stopping is performed or not.

IV.1. The James-Stein Estimator

Let  $X_1, \dots, X_n$  denote  $n$  iid univariate random variables, with mean  $\theta$ .

A common justification for taking  $\bar{X} = \frac{1}{n} \sum_{i=1}^n X_i$  is its desirable properties. It is unbiased for  $\theta$  ( $E\bar{X} = \theta$ ), and no other unbiased function of the data, linear or non-linear, can estimate  $\theta$  more accurately than  $\bar{X}$  in terms of variance (= mean square error).

In the 30's, Neyman, Pearson & Wald considered a different approach to this estimation problem & dropped unbiasedness as a criterion.

↳ Statistical Decision Theory

They examine all functions  $\hat{\theta} = \theta(X_1, \dots, X_n)$  of the data & estimators are compared through a risk function, such as the square loss:  $R(\hat{\theta}) = E(\hat{\theta} - \theta)^2$ .

\* Ex = Let  $X_1, \dots, X_n \sim \mathcal{N}(\theta, \sigma^2)$ ,  $\sigma^2$  known, and consider the following estimators of  $\theta$ :

$\hat{\theta}_1 = \bar{X}$  ,  $\hat{\theta}_2 = \frac{1}{2} \bar{X}$  ,  $\hat{\theta}_3 = \text{median}(X_1, \dots, X_n)$

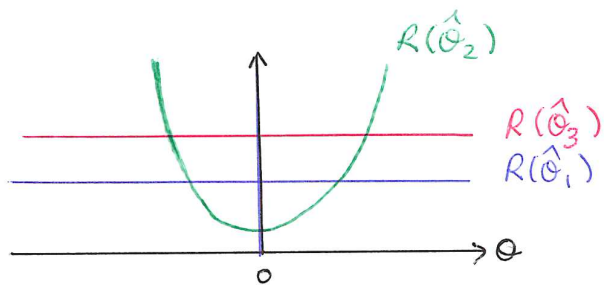
We compare the risk of  $\hat{\theta}_1$ ,  $\hat{\theta}_2$  and  $\hat{\theta}_3$ .

•  $R(\hat{\theta}_1) = E(\bar{X} - \theta)^2 = \text{var } \bar{X} = \frac{\sigma^2}{n}$   
"  $E\hat{\theta}_1$

$$\begin{aligned}
 R(\hat{\theta}_2) &= \mathbb{E}\left(\frac{1}{2}\bar{X} - \theta\right)^2 = \frac{1}{4} \mathbb{E}(\bar{X} - 2\theta)^2 \\
 &= \frac{1}{4} \left( \mathbb{E}(\bar{X}^2) - 4\theta \mathbb{E}X + 4\theta^2 \right) \\
 &= \frac{1}{4n^2} \mathbb{E}\left( \sum_i X_i^2 + \sum_{i \neq j} X_i X_j \right) \\
 &= \frac{1}{4} \theta^2 + \frac{1}{4n} \sigma^2
 \end{aligned}$$

63

• For large  $n$ , it is possible to show that the sample median is  $\approx \mathcal{N}\left(\theta, \frac{\pi}{2} \frac{\sigma^2}{n}\right)$ , so that  $R(\hat{\theta}_3) = \frac{\pi}{2} \frac{\sigma^2}{n}$ .



- Estimator  $\hat{\theta}_1$  has smaller risk than  $\hat{\theta}_3$  for any value of  $\theta$ . We say that  $\hat{\theta}_3$  is INADMISSIBLE.
- Estimator  $\hat{\theta}_2$  has smaller risk than  $\hat{\theta}_1$  for small values of  $\theta$  only. A natural question arises: can we beat the sample mean estimator  $\hat{\theta}_1 = \bar{X}$  for any value of  $\theta$ ? An answer to this question was provided by Blyth, Lehmann and Hodges in 1950: there are no such estimator. In other words,  $\hat{\theta}_1$  is ADMISSIBLE.

Stein considered the estimation of several means at one time, where the risk is taken as the sum of the expected values of

of square errors of estimation for all individual means:  $R(\hat{\theta}) = \mathbb{E} \|\hat{\theta} - \theta\|^2 = \sum_{i=1}^d \mathbb{E}(\hat{\theta}_i - \theta_i)^2$ , where  $\theta = (\theta_1, \dots, \theta_d)^t \in \mathbb{R}^d$  is the mean vector to be estimated.

64

In 1955, Stein showed that the vector of sample means is also admissible if we estimate simultaneously 2 means, but inadmissible if there are 3 means or more.

↖ In fact, Stein showed this in 1955 only if the number of means is very large. He extended the result together with James in 1961 for all  $n \geq 3$ , and did so in a constructive manner.

### • James-Stein estimator

Let  $X_1, \dots, X_d$  be  $d$  independent variables with

$$X_i | \theta_i \sim \mathcal{N}(\theta_i, \sigma_0^2).$$

Put  $X = (X_1, \dots, X_d)^t \in \mathbb{R}^d$

$\theta = (\theta_1, \dots, \theta_d)^t \in \mathbb{R}^d$ , so that  $X | \theta \sim \mathcal{N}(\theta, \sigma_0^2 I)$

We are looking for estimators  $\hat{\theta} = \theta(X)$  of  $\theta$ .

Consider the case of a single observation  $X$ . Each component of  $\theta$  is therefore estimated using a single noisy measurement.

The maximum likelihood approach gives  $\hat{\theta}_{ML} = X$ . It has risk

$$R(\hat{\theta}_{ML}) = \sum_{i=1}^d \mathbb{E}(X_i - \theta_i)^2 = d \sigma_0^2.$$

↖ since the  $i$ -th component of  $\hat{\theta}_{ML}$  is  $X_i$ .



James & Stein (1961) introduced a novel estimator of  $\theta$  that has a strictly smaller risk than  $\hat{\theta}_{ML}$  for any  $d \geq 3$ . It can be derived in a Bayesian context.

(65)

Consider the model

$$\begin{cases} X | \theta \sim \mathcal{N}(x | \theta, \sigma_0^2 \mathbf{I}) \\ \theta \sim \mathcal{N}(\theta | 0, \alpha \mathbf{I}) \end{cases}$$

but analysed using frequentist tools.

general covariance matrix (not necessarily diagonal)

More generally ...  
 ... mean  $m$   
 ... known or unknown ( $\alpha > 0, \sigma_0^2 > 0$ )  
 ... not necessarily normal

We consider the simplest case here, where both  $\alpha$  and  $\sigma_0^2$  are known.

The posterior distribution is  $\theta | X \sim \mathcal{N}(\theta | \beta X, \beta \mathbf{I})$ , where  $\beta = \frac{\alpha}{\alpha + \sigma_0^2}$ .

The posterior mean is then  $\hat{\theta}_{post} = \beta X = \frac{\alpha}{\alpha + \sigma_0^2} X$

$$\hat{\theta}_{post} = \left( 1 - \frac{\sigma_0^2}{\alpha + \sigma_0^2} \right) X$$

All components of  $X$  are shrunk towards 0

When  $\alpha$  is unknown, it can be estimated from the data. Estimating the prior parameters using the data lies at the heart of Empirical Bayes (EB) methods. We explain the approach more generally in the next subsection. In the context of the James-Stein estimator, note that the marginal distribution of

$X$  is  $\mathcal{N}(x | 0, (\alpha + \sigma_0^2) \mathbf{I})$  - indpt components

(66)

$$\Rightarrow \sum_{i=1}^d \frac{X_i^2}{\alpha + \sigma_0^2} \sim \chi_d^2$$

Denoting  $\|X\|^2 = \sum_{i=1}^d X_i^2$ , we see that  $\frac{\|X\|^2}{\alpha + \sigma_0^2} \sim \chi_d^2$

$$\frac{\alpha + \sigma_0^2}{\|X\|^2} \sim \text{Inv-}\chi_d^2$$

$$\Rightarrow \mathbb{E} \left( \frac{\alpha + \sigma_0^2}{\|X\|^2} \right) = \frac{1}{d-2}$$

$$\mathbb{E} \left( \frac{d-2}{\|X\|^2} \right) = \frac{1}{\alpha + \sigma_0^2}$$

We can replace the term  $\frac{1}{\alpha + \sigma_0^2}$  appearing in the expression of the posterior mean with  $\frac{d-2}{\|X\|^2}$ , an unbiased estimator of it. We obtain precisely the definition of the James-Stein estimator of  $\theta$ :

$$\hat{\theta}_{JS} = \left( 1 - \frac{(d-2)\sigma_0^2}{\|X\|^2} \right) X$$

shrinkage factor

- The term  $\|X\|^2$  assesses if each component of  $X$  is far from 0.
- If  $X$  is close to zero,  $\|X\|^2$  is small & the shrinkage is large: the JS estimator shrinks each component of  $X$  a lot since there is confidence that  $X$  is close to 0.
- If  $X$  is far from the origin,  $\|X\|^2$  is large & the shrinkage is small (since not much confidence in shrinking values towards 0)

### Theorem [James & Stein (1961)]

(67)

Suppose  $X | \theta \sim \mathcal{N}(x | \theta, \sigma_0^2 I)$ ,  $\sigma_0^2 > 0$  known,  $X, \theta \in \mathbb{R}^d$ .

For  $d \geq 3$ , the James-Stein estimator  $\hat{\theta}_{JS}$  everywhere dominates the maximum likelihood estimate  $\hat{\theta}_{ML}$  under the square loss,

$$\mathbb{E} \|\hat{\theta}_{JS} - \theta\|^2 < \mathbb{E} \|\hat{\theta}_{ML} - \theta\|^2,$$

for every fixed  $\theta$ .

A frequentist result! ( $\theta$  fixed)  
Not a Bayesian one!

$\Rightarrow$  The JS estimator is biased for  $\theta$ , but has lower MSE than the MLE for any fixed value of  $\theta$  (the increase in bias is compensated by the large decrease in variance).

This makes the MLE inadmissible for  $d \geq 3$ .

Remark: The result holds as well for  $m$  iid observations

$X_1, \dots, X_m$  since if  $X_i | \theta \sim \mathcal{N}(x | \theta, \sigma_0^2 I_d)$ , then  $\bar{X} = \frac{1}{m} \sum_{i=1}^m X_i \sim \mathcal{N}(x | \theta, \frac{\sigma_0^2}{m} I_d)$ , and the previous analysis holds.

$\Rightarrow$  Repeat the derivation of the JS estimator replacing  $\sigma_0^2$  with  $\frac{\sigma_0^2}{m}$ ,

$$\hat{\theta}_{JS} = \left( 1 - \frac{(d-2) \frac{\sigma_0^2}{m}}{\|\bar{X}\|^2} \right) \bar{X}$$

Each component of  $\bar{X} \in \mathbb{R}^d$  is shrunk by the same factor.

$$\text{proof} = R(\hat{\theta}_{JS}) = \mathbb{E} \left\| X - \theta - \frac{(d-2)\sigma_0^2}{\|X\|^2} X \right\|^2$$

(68)

$$= \mathbb{E} \|X - \theta\|^2 - 2(d-2)\sigma_0^2 \sum_{i=1}^d \mathbb{E} \left\{ \frac{X_i(X_i - \theta_i)}{\|X\|^2} \right\} + (d-2)^2 \sigma_0^4 \mathbb{E} \left\{ \frac{1}{\|X\|^2} \right\}.$$

$$\begin{aligned} (*) &= \int \frac{x_i}{\|x\|^2} (x_i - \theta_i) \frac{1}{(2\pi\sigma_0^2)^{d/2}} \exp \left\{ -\frac{1}{2\sigma_0^2} \|x - \theta\|^2 \right\} dx \\ &= \int \frac{\|x\|^2 - 2x_i}{\|x\|^4} \left[ \sigma_0^2 \frac{1}{(2\pi\sigma_0^2)^{d/2}} \exp \left\{ -\frac{\|x - \theta\|^2}{2\sigma_0^2} \right\} \right] dx \\ &= \sigma_0^2 \mathbb{E} \left\{ \frac{\|X\|^2 - 2X_i}{\|X\|^4} \right\} \end{aligned}$$

$$\begin{aligned} \sum_{i=1}^d \mathbb{E} \left\{ \frac{X_i(X_i - \theta_i)}{\|X\|^2} \right\} &= d\sigma_0^2 \mathbb{E} \left\{ \frac{1}{\|X\|^2} \right\} - 2\sigma_0^2 \mathbb{E} \left\{ \frac{\sum_{i=1}^d X_i}{\|X\|^4} \right\} \\ &= (d-2)\sigma_0^2 \mathbb{E} \left\{ \frac{1}{\|X\|^2} \right\}. \end{aligned}$$

We get

$$\begin{aligned} R(\hat{\theta}_{JS}) &= d\sigma_0^2 - 2(d-2)\sigma_0^4 \mathbb{E} \left\{ \frac{1}{\|X\|^2} \right\} + (d-2)^2 \sigma_0^4 \mathbb{E} \left\{ \frac{1}{\|X\|^2} \right\} \\ &= d\sigma_0^2 - \underbrace{(d-2)^2 \sigma_0^4 \mathbb{E} \left\{ \frac{1}{\|X\|^2} \right\}}_{> 0 \text{ for } d \geq 3} \\ &< d\sigma_0^2 = R(\hat{\theta}_{ML}). \end{aligned}$$

Remarks (i) The James-Stein estimator derived on page 69 shrinks the components of  $X$  towards 0. Other shrinking value can be considered. For example, Efron & Morris (1977) shrink the components of  $X \in \mathbb{R}^d$  towards its mean  $\bar{X} = \frac{1}{d} \sum_{i=1}^d X_i$ :

$$\hat{\theta}_{EM} = \bar{X} + \left( 1 - \frac{(d-3)\sigma_0^2}{\|X - \bar{X}\|^2} \right) (X - \bar{X})$$

$\hat{\theta}_{EM}$  has strictly smaller risk than  $\hat{\theta}_{ML}$  when  $d \geq 4$ .

It can be derived in a similar manner as before, under the assumption that  $X_i | \theta_i \sim \mathcal{N}(\theta_i, \sigma_0^2)$ , and  $\theta_i \sim \mathcal{N}(m, \alpha)$ .

(ii) In the expression of the original JS estimator, the shrinkage factor  $\left( 1 - \frac{(d-2)\sigma_0^2}{\|X\|^2} \right)$  may be less than 0, and arbitrary large negatively, which is counterintuitive. This can be avoided by taking the positive-part of the shrinkage factor.

Bainchick (1964, 1970) showed that the JS-positive-part estimator

$$\hat{\theta}_{JS}^+ = \left( 1 - \frac{(d-2)\sigma_0^2}{\|X\|^2} \right)_+ X$$

dominates  $\hat{\theta}_{JS}$  for any  $\theta$ .  $(u)_+ = \begin{cases} u & \text{if } u \geq 0 \\ 0 & \text{otherwise} \end{cases}$   
smaller MSE

In other words, the JS estimator is inadmissible.

Brown (1971) showed that  $\hat{\theta}_{JS}^+$  also is inadmissible (but improving on  $\hat{\theta}_{JS}^+$  is known to be hard)

### (iii) James-Stein & Risk Minimization.

Since the JS estimator achieves lower MSE than the MLE, and is of the form  $cX$  for  $c = \left( 1 - \frac{(d-2)\sigma_0^2}{\|X\|^2} \right)$ , we may directly consider the class of estimators  $cX$  for  $c \in \mathbb{R}$ , and look for the optimal value of  $c$ .

$$\begin{aligned} R(cX) &= \mathbb{E} \| \theta - cX \|^2 \\ &= \| \theta \|^2 + c \underbrace{\mathbb{E} \| X \|^2}_{d\sigma_0^2 + \| \theta \|^2} - 2c \underbrace{\sum_{i=1}^d \mathbb{E}(X_i \theta_i)}_{\| \theta \|^2} \\ &= \| \theta \|^2 (1-c)^2 + c^2 d \sigma_0^2. \end{aligned}$$

The minimum is achieved at  $c^* = \frac{\| \theta \|^2}{\| \theta \|^2 + d \sigma_0^2}$ .

This value cannot be computed in practice since it relies on the unknown  $\theta$ . However, note that  $\mathbb{E} \| X \|^2 = \| \theta \|^2 + d \sigma_0^2$ , so that  $\| X \|^2$  is an unbiased estimator for  $(\| \theta \|^2 + d \sigma_0^2)$ . This leads to

$$\hat{\theta} = \left( 1 - \frac{d \sigma_0^2}{\| X \|^2} \right) X, \text{ which has the same form as } \hat{\theta}_{JS}, \text{ except that the optimal value of } (d-2) \text{ is replaced with } d.$$

### (iv) James-Stein & Ridge Regression.

[Ridge Regression] see SL: RIDGE REGRESSION & LASSO

The linear model is  $y = X\beta + \varepsilon$ ,  $\varepsilon \sim \mathcal{N}(0, \sigma^2 I_n)$ , so that  $y \sim \mathcal{N}(X\beta, \sigma^2 I_n)$  + Bayesian prior  $\beta \sim \mathcal{N}(0, \tau I_d)$ .

The posterior is  $\beta | \mathcal{D}_n \sim \mathcal{N}(\beta | \mu_n, \Sigma_n)$ , with

$$\begin{cases} \mu_n = \frac{1}{\sigma^2} \sum_n X^t y \\ \Sigma_n^{-1} = \frac{1}{\sigma^2} I + \frac{1}{\sigma^2} X^t X \end{cases}$$

The posterior mean is thus  $\frac{1}{\sigma^2} \left( \frac{1}{\tau} I + \frac{1}{\sigma^2} X^t X \right)^{-1} X^t y$   
 $= \left( \underbrace{\frac{\sigma^2}{\tau} I + X^t X}_{\lambda} \right)^{-1} X^t y$

⇒ With  $\tau = \frac{\sigma^2}{\lambda}$ , the Bayesian view of RR is equivalent to the penalized optimization problem

min  $RSS_2(\lambda)$ , where  $RSS_2(\lambda) = \|y - X\beta\|^2 + \lambda \|\beta\|^2$ ,

since the value  $\hat{\beta}_\lambda$  minimizing  $RSS_2(\lambda)$  is precisely given by  $\hat{\beta}_\lambda = (\lambda I + X^t X)^{-1} X^t y$ .

Note that when  $\lambda=0$ ,  $\hat{\beta}_0 = \hat{\beta} = \text{least squares solution} = (X^t X)^{-1} X^t y$ .

One can show that  $\hat{\beta} \sim \mathcal{N}(\beta, \sigma^2 (X^t X)^{-1})$ .

[James-Stein] Assume first that the columns of  $X$  are orthogonal, i.e.  $X^t X = I_d$ . Then the Bayesian linear model is equivalent to  $\hat{\beta} \sim \mathcal{N}(\beta, \sigma^2 I_d)$   
 $\beta \sim \mathcal{N}(0, \tau I_d)$

and the James-Stein estimator of  $\beta$  is

$$\hat{\beta}_{JS} = \left( 1 - \frac{(d-2)\sigma^2}{\|\hat{\beta}\|^2} \right) \hat{\beta} \quad \text{vs} \quad \hat{\beta}_\lambda = \frac{1}{1+\lambda} \hat{\beta}$$

↑  
Entries of  $\hat{\beta}$  are shrunk in both cases. However,

the JS procedures enforce a specific value of the shrinkage factor that depends on  $\hat{\beta}$  itself, while the RR approach determines the optimal  $\lambda$  via other techniques such as cross-validation.

In the general case, when  $\hat{\beta} \sim \mathcal{N}(\beta, \sigma^2 (X^t X)^{-1})$   
 $\beta \sim \mathcal{N}(0, \tau I_d)$ ,

the JS-like estimator can be taken as

$$\tilde{\beta}_{JS} = \left( 1 - \frac{(d-2)\sigma^2}{\hat{\beta}^t (X^t X) \hat{\beta}} \right) \hat{\beta},$$

see for example [Bock \(1975\)](#).

(v) Stein's paradox.

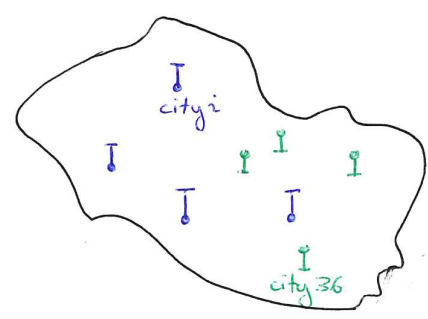
Ex: Toxoplasmosis, [Efron & Morris \(1977\)](#).

↳ disease of the blood endemic in Central America

↳ study in El Salvador on ~5000 people in 36 cities.

↳ interested in the incidence rate (w.r.t. the national average)

↳ incidence rates are ≈ normally distributed.



$$\theta = \begin{pmatrix} \theta_1 \\ \vdots \\ \theta_{36} \end{pmatrix},$$

where  $\theta_i =$  incidence rate of city  $i$  with respect to the national average.

Frequentist estimator  $\bar{X}$ , whose  $i$ -th coordinate corresponds to the sample average for city  $i$

James-Stein estimator, which shrinks the components of  $\bar{X}$  by some factor.

Which one to choose?

If interested in one city in particular, then consider the MLE i.e. the observed rate for that city

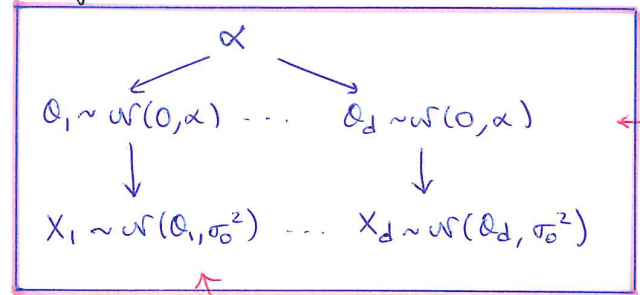
If interested in the precision at the national level, then select the JS estimator. For example, if the Ministry of Health wants to build local hospitals. The hospitals may end up at the wrong place or have the wrong size, but the sum of the mismatches will be smaller than for the observed rates.

There is a flaw in the reasoning. Suppose that you augment the vector of unknowns  $\theta$  with an extra component  $\theta_{37}$  corresponding to the incidence rate of a tumor in a group of rats. Including this extra coordinate may reduce the overall risk compared with the MLE. Combining two unrelated problems into one, and reducing the global risk, is at the origin of what is known as STEIN'S PARADOX.

In particular, note that the shrinkage factor in the expression of the JS estimator depends on  $\|X\|^2$ , i.e. on all the components of  $X = (X_1, \dots, X_d)^t$ , despite the assumption that the  $X_i$  are independent  $\Rightarrow$  as if the incidence rate estimate of a tumor in a group of rats carries information for estimating the incidence rate of toxoplasmosis in El Salvador.

The Empirical Bayes formulation of the JS estimator brings light to the picture. The components of  $\theta$  are assumed to be  $\mathcal{N}(0, \alpha)$  distribution.

$\Rightarrow$  The common variance amongst the  $\theta_i$  indicate that they are all similar to each other, up to some random fluctuations  $\Rightarrow \theta_1$  can learn from  $(\theta_2, \dots, \theta_d)$  and  $X_1$  together.



sample of the population distrib  $\sim \mathcal{N}(0, \alpha)$  = the  $\theta_1, \dots, \theta_d$  are connected via  $\alpha$ .

This picture emphasizes the underlying structure/hierarchy in the approach.

Works well when studying clusters of individuals. In particular, using the data to estimate the prior parameters is inherent to the Empirical Bayes (EB) approach, and is described further in the next section.

IV.2. The Empirical Bayes Approach.

Empirical Bayes method are usually not considered "full Bayesian", since the prior parameters are estimated from the data. We illustrate the procedure in the binomial case with

Beta prior:  $X_i | p_i \sim \text{Bi}(x_i | n_i, p_i)$  iid  
 $p_i | \alpha, \beta \sim \text{B}(p_i | \alpha, \beta)$  iid

EB uses point estimates  $\hat{\alpha}, \hat{\beta}$  of the prior distribution, maximizes the MARGINAL LIKELIHOOD:

$$f(\mathcal{X}_n | \alpha, \beta) = f(x_1, \dots, x_n | \alpha, \beta)$$

↑

Same idea as the frequentist Maximum Likelihood estimation, except that maximization occurs on the prior parameters instead of the parameters themselves.

$$= \int \underbrace{f(x_1, \dots, x_n, p_1, \dots, p_n | \alpha, \beta)}_{\text{full joint distribution}} dp_1 \dots dp_n$$

$$= \int f(x_1, \dots, x_n | p_1, \dots, p_n, \alpha, \beta) \times f(p_1, \dots, p_n | \alpha, \beta) dp_1 \dots dp_n$$

(independence)

$$= \prod_{i=1}^n \int f(x_i | p_i) f(p_i | \alpha, \beta) dp_i$$

Focusing on a single term

$$\int_0^1 \text{Bi}(x | n, p) B(p | \alpha, \beta) dp = \text{The Beta-Binomial distribution}$$

$$= \dots = \binom{n}{x} \frac{B(x+\alpha, n-x+\beta)}{B(\alpha, \beta)}$$

↑ the beta function (not distribution)

↳ Then consider  $(\hat{\alpha}, \hat{\beta}) = \underset{(\alpha, \beta)}{\text{argmax}} f(\mathcal{X}_n | \alpha, \beta)$

Or, alternatively, use the method of moments

EB violates the Bayesian principle that the prior should be chosen independently of the data. A "more Bayesian" approach, described in Section IV, puts a prior distribution on  $(\alpha, \beta)$ .

• Robbin's formula, or Non Parametric Empirical Bayes

Consider the case of a Poisson distribution

$$\begin{cases} X_i | \theta_i \sim P(\theta_i) \\ \theta_i \sim g \end{cases}$$

where  $g$  is left unspecified.

We are interested in the posterior mean  $E(\theta_i | X_i = x_i)$

We drop the subscript  $i$  for convenience:

$$E(\theta | X=x) = \frac{\int_0^\infty \theta p(x | \theta) g(\theta) d\theta}{\int_0^\infty p(x | \theta) g(\theta) d\theta}$$

=  $f(x)$ , the marginal likelihood

$$= \frac{\int_0^\infty e^{-\theta} \frac{\theta^{x+1}}{x!} g(\theta) d\theta}{\int_0^\infty e^{-\theta} \frac{\theta^x}{x!} g(\theta) d\theta} = (x+1) \frac{\int_0^\infty e^{-\theta} \frac{\theta^{x+1}}{(x+1)!} g(\theta) d\theta}{\int_0^\infty e^{-\theta} \frac{\theta^x}{x!} g(\theta) d\theta}$$

$$= (x+1) \frac{f(x+1)}{f(x)}$$

We estimate the marginal likelihood non-parametrically (instead of maximizing the marginal likelihood in a parametric context)

Let  $\hat{f}(x+1) = \frac{1}{n} (\# \text{ observations} = x+1) =: \frac{n_{x+1}}{n}$

Then  $E(\theta | X=x) \approx (x+1) \frac{n_{x+1}}{n_x}$

ROBBIN'S FORMULA

### IV.3. Bayesian False Discovery Rate

(77)

Context: Multiple testing;  $d$  tests, where each test assumes that  $H_0$  has prior  $\pi_0$ ;  $X \sim f_0 \sim \mathcal{N}(0, 1)$  on  $H_0$   
 $H_1$  — "—  $\pi_1 = 1 - \pi_0$ ;  $X \sim f_1$  "away from 0" on  $H_1$ .

The cumulative distribution functions are denoted  $F_0$  and  $F_1$ , respectively.

Then  $X \sim f =$  mixture distribution;  $f(x) = \pi_0 f_0(x) + \pi_1 f_1(x)$ .

Let  $F(X) := \int_X f(x) dx$  for some measurable set  $X$ .

We consider the quantity

$$\phi(X) := \mathbb{P}(H_0 | x \in X) = \frac{\mathbb{P}(x \in X | H_0) \pi_0}{\mathbb{P}(x \in X)} = \frac{F_0(X) \pi_0}{F(X)}$$

aka the Bayesian False Discovery Rate; also usually denoted  $Fdr(X)$ .

[if we report  $x \in X$  as non-null, then  $\phi(X)$  represents the probability of false-discovery].

Remark = Typical application where  $X$  is a  $z$ -score i.e.

$z_i = \frac{\bar{X}}{s/\sqrt{n}}$  for the  $i$ -th test  $\sim \mathcal{N}(0, 1)$  under  $H_0$ .

Empirical Bayes estimate of  $\phi(X)$ .

Assumptions:

$\rightarrow f_0$  is known since it is the distribution of a statistic of interest under the null.

$\rightarrow f_1$  is unlikely to be known.

$\rightarrow \pi_0$  is unknown.

(78)

$\rightarrow F$  can be easily estimated with  $\hat{F}(X) = \frac{1}{d} \#\{X_i \in X\}$ .

$\Rightarrow$  We get the estimator  $\hat{\phi}(X) = \frac{F_0(X) \pi_0}{\hat{F}(X)}$

How good of an estimator of  $\phi$  is  $\hat{\phi}$ ?

Let

$n_0(X) =$  # of null  $X_i$  falling into  $X$  (unobservable)

$n_1(X) =$  # of non-null  $X_i$  falling into  $X$  (unobservable)

$n_+(X) = n_0(X) + n_1(X)$  (observable)

Then

$$e_0(X) := \mathbb{E} n_0(X) = \mathbb{E} \sum_{i=1}^d \mathbb{1}(X_i \in X | X_i \in H_0) \mathbb{1}(X_i \in H_0) = d \pi_0 F_0(X)$$

$$e_1(X) := \mathbb{E} n_1(X) = d \pi_1 F_1(X)$$

$$e_+(X) := \mathbb{E} n_+(X) = d F(X).$$

With this notation,

$$\phi(X) = \frac{F_0(X) \pi_0}{F(X)} = \frac{e_0(X)}{e_+(X)}$$

$$\hat{\phi}(X) = \frac{\mathbb{E} n_0(X)}{n_+(X)} = \frac{e_0(X)}{n_+(X)}$$

In addition, we consider

$$Fdp(X) := \frac{n_0(X)}{n_+(X)} = \text{proportion of } X_i \text{ from the null in } X$$

Assume that  $e_0(x) = \mathbb{E}(n_0(x) | n_1(x))$ ; which holds when  $n_0(x) \perp n_1(x)$ . Then

(79)

$$\mathbb{E}(\text{FdP}(x) | n_1(x)) \leq \mathbb{E}(\hat{\phi}(x) | n_1(x))$$

"For any value of  $n_1(x)$ , the EB estimate  $\hat{\phi}(x)$  is a conservatively biased estimate of the actual  $\text{FdP}(x)$ ".

• proof:  $\text{FdP}(x) = \frac{n_0(x)}{n_0(x) + n_1(x)} = \varphi(n_0(x)) - \varphi = \text{concave}$

Jensen inequality  $\Rightarrow$

$$\mathbb{E}(\text{FdP}(x) | n_1(x)) \leq \frac{e_0(x)}{e_0(x) + n_1(x)}$$

On the other hand,

$$\hat{\phi}(x) = \frac{e_0(x)}{n_0(x) + n_1(x)} = \tilde{\varphi}(n_0(x)) - \tilde{\varphi} = \text{convex}$$

Applying Jensen again:

$$\mathbb{E}(\hat{\phi}(x) | n_1(x)) \geq \frac{e_0(x)}{e_0(x) + n_1(x)}$$

Combining the two inequalities yields the desired result.  $\square$

### • Mean and Variance of $\hat{\phi}(x)$

Suppose that the  $X_1, \dots, X_d$  are independent.

Then  $n_+(x) \sim \text{Bi}(d, F(x))$ .

$$\hat{\phi}(x) = \frac{e_0(x)}{n_+(x)} = \frac{e_0(x)}{e_+(x)} \times \frac{1}{1 + \frac{n_+(x) - e_+(x)}{e_+(x)}}$$

$$\hat{\phi}(x) \approx \frac{e_0(x)}{e_+(x)} \left[ 1 - \frac{n_+(x) - e_+(x)}{e_+(x)} + \left( \frac{n_+(x) - e_+(x)}{e_+(x)} \right)^2 \dots \right] \quad (80)$$

where  $\frac{n_+(x) - e_+(x)}{e_+(x)}$  has mean 0 and variance

$$\gamma(x) := \frac{\text{var } n_+(x)}{e_+^2(x)}$$

It follows that  $\frac{\hat{\phi}(x)}{\phi(x)}$  has approximately mean  $1 + \gamma(x)$  and variance  $\gamma(x)$ .

Under the independence assumption,  $\text{var } n_+(x) = d F(x)(1-F(x))$

$$\Rightarrow \gamma(x) = \frac{d F(x)(1-F(x))}{d^2 F^2(x)} = \frac{1-F(x)}{d F(x)} = \frac{1-F(x)}{e_+(x)}$$

since we are mostly interested in regions  $x$  where  $F(x)$  is small.  $\approx \frac{1}{e_+(x)}$

$$\Rightarrow \frac{\hat{\phi}(x)}{\phi(x)} \text{ has mean } \approx 1 + \frac{1}{e_+(x)}$$

$$\& \text{ variance } \approx \frac{1}{e_+(x)}$$

These expressions allow us to evaluate if  $\hat{\phi}$  is a biased estimator of  $\phi$ , and to construct approx confidence bounds for  $\phi$ .

The accuracy of  $\hat{\phi}(x)$  heavily depends on  $e_+(x)$  = the expected number of observations falling into  $x$



\* Illustration =  $X$  = tail interval of the form  $(-\infty, x)$  or  $(x, +\infty)$ . (81)

Take e.g.  $X = (3, +\infty)$ , and suppose that  $n_+(X) = 50$ .

Then  $\hat{\phi}(X) = \frac{d \pi_0 (1 - \Phi(3))}{n_+(X)} = 0.1$  (say)

take  $\pi_0 = 1$  to get an upper bound on  $\hat{\phi}(X)$  (conservative approach)

meaning that  $\approx 10\%$  of the 50 discoveries are false discoveries.

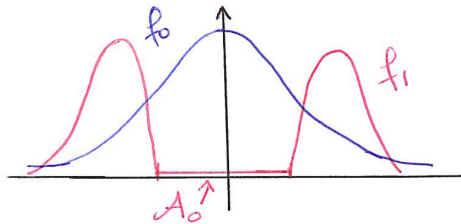
- Remarks (i) This approach is particularly well suited when testing thousands of hypothesis where many are expected to have zero effect ( $\pi_0 \approx 1$ ).  
Ex: effect of a treatment on gene expression.

(ii) Estimation of  $\pi_0$

$\pi_0$  can be estimated under additional assumptions on  $f_1$ .

Suppose that  $f_1$  is zero for some subset  $A_0$  around 0

("separability" condition)



Then  $n_+(A_0) = \# \{X_i \in A_0\}$

and  $E n_+(A_0) = \pi_0 d F_0(A)$ , which yields

$\hat{\pi}_0 = \frac{n_+(A_0)}{d \times F_0(A_0)}$  ← You may vary the size of  $A_0$  and check visually the stability of  $\hat{\pi}_0$ .

(iii) Relation to the JS estimator (82)

The generative process  $X \sim f = \pi_0 f_0 + \pi_1 f_1$  can be rewritten

(\*) 
$$\begin{cases} X | \mu \sim \mathcal{N}(\mu, 1) \\ \mu \sim \pi_0 \delta_0(\mu) + (1 - \pi_0) g_1(\mu) \end{cases}$$

leading to  $f_1(x) = \int \varphi(x - \mu) g_1(\mu) d\mu$ .

In comparison with the JS estimator, which is constructed from the model  $\begin{cases} X | \mu \sim \mathcal{N}(\mu, 1) \\ \mu \sim \mathcal{N}(0, \alpha) \end{cases}$  (\*\*)

"learning from others" (p. 65)

Expression (\*) and its comparison with the smooth version (\*\*) highlights that each test is learning from one another to control the FDR.

- JS  $\rightarrow$  performs well even when  $d$  is not too large
- Bayesian FDR  $\rightarrow$  the variance of  $\hat{\phi}(X)$  is of order  $1/e_+(X)$  so in comparison  $d$  needs to be quite large to achieve an accurate estimation of  $\phi(X)$ .

(iv) Link with the Benjamini-Hochberg (BH) procedure

The BH procedure was introduced p. 29-33 in the Chapter MS = HYPOTHESIS TESTING.

Consider  $d$  null hypotheses  $H_{0,1}, \dots, H_{0,d}$ , with z-scores  $z_{1,-}, z_d$  and p-values  $p_1, \dots, p_d$ .

Let  $p_{(1)} \leq \dots \leq p_{(d)}$  be the ordered p-values, and  $H_{0,(i)}$  be the null hypothesis associated with  $p_{(i)}$ .

Let  $k$  be the largest  $i$  for which  $p_{(i)} \leq \frac{i\alpha}{d}$ .

- \* BH: reject all  $H_{0,(i)}$  for  $i=1, \dots, k$ .
- \* Guarantee: if the  $p_1, \dots, p_d$  are independent, then

$$FDR = \mathbb{E} \left\{ \frac{FP}{FP+TP} \mathbb{1}(FP+TP \geq 1) \right\} \leq \alpha.$$

Suppose that the  $p_i$  correspond to the left p-values, so that  $p_i = F_0(z_i)$ . Then  $p_{(i)} = F_0(z_{(i)})$ .

On the other hand, ordering  $z_{(1)} \leq \dots \leq z_{(d)}$  yields

$$\hat{F}(z_{(i)}) = \frac{i}{d}.$$

Then

$$p_{(i)} \leq \frac{i\alpha}{d} \Leftrightarrow \pi_0 \frac{F_0(z_{(i)})}{\hat{F}(z_{(i)})} \leq \pi_0 \alpha \leq \alpha$$

$$\Leftrightarrow \Phi(-\infty, z_{(i)}) \leq \alpha$$

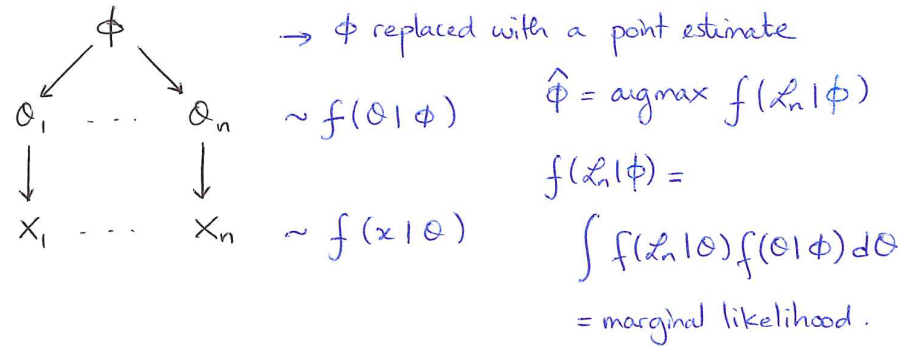
The BH criterion can be equivalently re-expressed in terms of the Bayesian FDR.

(Major) Consequence = For the first time  $\alpha$  has the intuitive meaning that a rejected hypothesis  $H_{0,i}$  is actually correct (i.e. a Bayesian meaning).

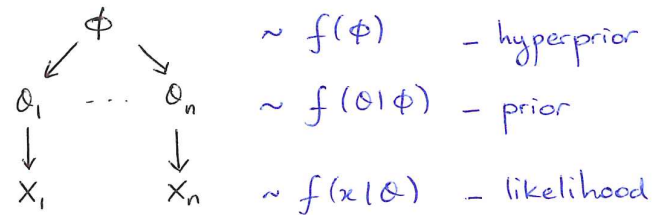
V - HIERARCHICAL MODELS

While Empirical Bayes considers point estimates for the prior's parameters, Hierarchical Models put a prior on them

• Empirical Bayes



• Hierarchical Models



- (i) How to select the hyperprior  $f(\phi)$ ?
- (ii) How to (compute) sample from the joint posterior  $(\theta, \phi) | L_n$  and the posterior predictive distribution?

- A general procedure:
- (1) Sample  $\phi$  from  $\phi | \mathcal{L}_n$
  - (2) Sample  $\theta$  from  $\theta | \phi, \mathcal{L}_n$
  - (3) Sample a new  $X$  from  $X | \theta, \mathcal{L}_n$

↳ Step (3) is easy

Note that step (1) is feasible when  $f(\phi | \mathcal{L}_n)$  is known.

$$f(\phi | \mathcal{L}_n) = \frac{f(\theta, \phi | \mathcal{L}_n)}{f(\theta | \phi, \mathcal{L}_n)}$$

⇒ We need to compute  $f(\theta | \phi, \mathcal{L}_n)$  and  $f(\theta, \phi | \mathcal{L}_n)$ .

this posterior distribution can easily be computed analytically for conjugate models:

Also useful for step 2

$$f(\theta | \phi, \mathcal{L}_n) \propto \underbrace{f(\mathcal{L}_n | \phi, \theta)}_{\text{likelihood}} \underbrace{f(\theta | \phi)}_{\text{prior}}$$

$$f(\theta, \phi | \mathcal{L}_n) \propto f(\mathcal{L}_n | \phi, \theta) f(\theta | \phi) f(\phi)$$

common terms (need to be computed only once)

We illustrate the procedure for the Beta-Binomial model.

- $X_1, \dots, X_n$  ;  $X_i | \theta_i \sim \text{Bi}(n_i, \theta_i)$
- $\theta_i \sim B(\alpha, \beta)$
- $(\alpha, \beta) \sim \text{some distribution } f(\alpha, \beta)$ .

step (2)  $f(\theta | \mathcal{L}_n, \alpha, \beta) \propto f(\mathcal{L}_n | \theta, \alpha, \beta) f(\theta | \alpha, \beta)$

Assuming the  $\theta_i$  are iid generated

$$\begin{aligned} &= \prod_{i=1}^n \text{Bi}(x_i | n_i, \theta_i) \prod_{i=1}^n B(\theta_i | \alpha, \beta) \\ &= \prod_{i=1}^n \frac{\Gamma(\alpha + \beta)}{\Gamma(\alpha)\Gamma(\beta)} \theta_i^{\alpha-1} (1-\theta_i)^{\beta-1} \binom{n_i}{x_i} \theta_i^{x_i} (1-\theta_i)^{n_i-x_i} \\ &= \prod_{i=1}^n B(\theta_i | \alpha + x_i, \beta + n_i - x_i) \end{aligned}$$

The joint posterior is

$$f(\theta, \alpha, \beta | \mathcal{L}_n) \propto f(\alpha, \beta) \prod_{i=1}^n \text{Bi}(x_i | n_i, \theta_i) \prod_{i=1}^n B(\theta_i | \alpha, \beta)$$

The ratio of these two terms allow us to derive step (1):

$$f(\alpha, \beta | \mathcal{L}_n) \propto f(\alpha, \beta) \prod_{i=1}^n \frac{\Gamma(\alpha + \beta)}{\Gamma(\alpha)\Gamma(\beta)} \frac{\Gamma(\alpha + x_i) \Gamma(\beta + n_i - x_i)}{\Gamma(\alpha + \beta + n_i)}$$

After simplifications.

$f(\alpha, \beta | \mathcal{L}_n)$  is usually computed numerically & normalized to get a proper distribution.

Remains the choice of the prior  $f(\alpha, \beta)$ .

Gelman et al argue that the prior  $f(\alpha, \beta) \propto (\alpha + \beta)^{-5/2}$  yield a proper posterior distribution  $f(\alpha, \beta | \mathcal{L}_n)$  as long as  $0 < x_i < n_i$  for at least one  $i = 1, \dots, n$ .

## V.1. Application to multiple testing

87

Hierarchical models prove very useful when performing multiple tests. It complements alternative methods based on the Bayesian FDR which are more adapted to cases where the number of tests is large, and the number of nulls is also believed to be large.

The ideas developed here have been successfully applied in social sciences application, or psychology. See for example Gelman, Hill and Yajima (2012).

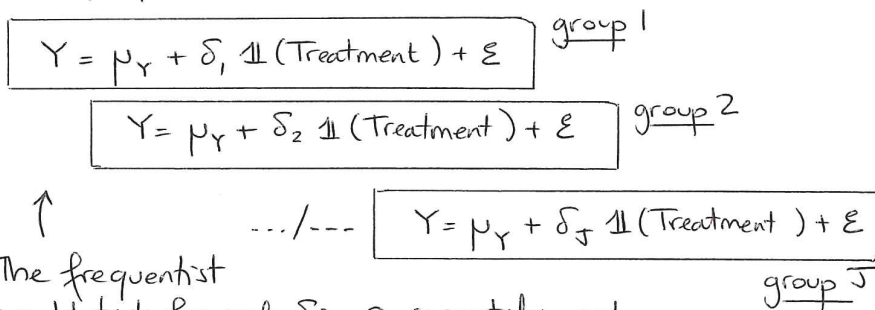
### \* 1 group

Suppose we are interested in estimating the treatment effect in a single group of population, divided in control and treatment sub-groups.

The variable of interest is denoted  $Y$ , and may be modeled using  $Y = \mu + \delta \mathbb{1}(\text{Treatment}) + \varepsilon$ ,  $\varepsilon \sim \mathcal{N}(0, \sigma_\varepsilon^2)$

$\delta$  represents the treatment effect

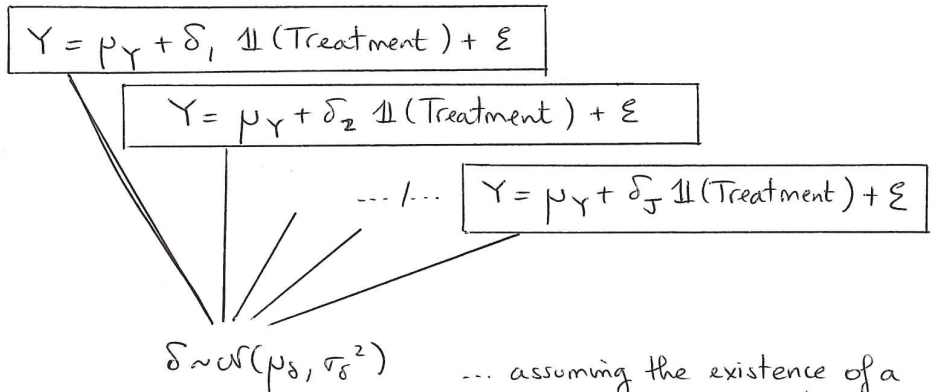
### \* J groups



The frequentist would test for each  $\delta_j = 0$  separately, and probably adjusting the thresholds for rejecting the nulls accordingly (e.g. Bonferroni).

The Bayesian would take a different path, ...

88



This approach allows for some heterogeneity across the groups; while allowing each group to learn from one another. The procedure is commonly referred to as "partial pooling".

### • Effect on the z-scores

We derive some algebra in a simple case, highlighting the main difference between the frequentist and the Bayesian approach.

In a pairwise set-up, each value  $Y$  is itself a difference between paired values ( $Y_1, Y_2$ ) i.e.  $Y = Y_1 - Y_2$ ; and we are testing for the departure of the mean value from 0. We can write  $Y = \delta_i + \varepsilon$  for group  $i$ , and we are testing for  $\delta_i = 0$ .

Then  $\underline{Y} = (Y_1, \dots, Y_d)^t \sim d^p(\underline{\delta}, \sigma_Y^2 \underline{I}_d)$   
 $\underline{\delta} = (\delta_1, \dots, \delta_d)^t \sim w^p(\mu, \sigma_\delta^2 \underline{I}_d)$

The posterior distribution of  $\delta | Y_1, \dots, Y_n$  is normal with mean  $\mu_n$  and covariance matrix  $\Sigma_n$ ,  $\mathcal{L}_n$

$\mu_n = E(\delta | \mathcal{L}_n) = \frac{\frac{\mu}{\sigma_\delta^2} + \frac{n}{\sigma_Y^2} \bar{Y}}{\frac{1}{\sigma_\delta^2} + \frac{n}{\sigma_Y^2}}$   $\leftarrow \bar{Y} = \frac{1}{n} \sum_{i=1}^n Y_i$

$\Sigma_n^{-1} = \left( \frac{1}{\sigma_\delta^2} + \frac{n}{\sigma_Y^2} \right) \underline{I}_d \Rightarrow \text{var}(\delta | \mathcal{L}_n) = \left( \frac{1}{\sigma_\delta^2} + \frac{n}{\sigma_Y^2} \right)^{-1}$

In such a set-up, we are also interesting in comparing groups two by two, evaluating the pairwise differences  $E(\delta_j - \delta_k | \mathcal{L}_n)$ . Straight forward calculations show that

$E(\delta_j - \delta_k | \mathcal{L}_n) = \frac{n/\sigma_Y^2}{1/\sigma_\delta^2 + n/\sigma_Y^2} (\bar{Y}_j - \bar{Y}_k)$

$\text{var}(\delta_j - \delta_k | \mathcal{L}_n) = \frac{2\sigma_\delta^2 \frac{\sigma_Y^2}{n}}{\sigma_\delta^2 + \sigma_Y^2/n}$

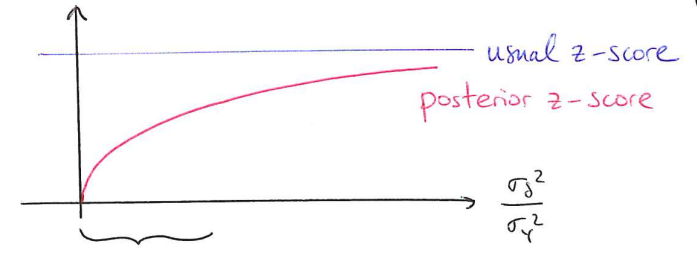
The posterior z-score is

$\frac{E(\delta_j - \delta_k | \mathcal{L}_n)}{\sqrt{\text{var}(\delta_j - \delta_k | \mathcal{L}_n)}} = \frac{\bar{Y}_j - \bar{Y}_k}{\sqrt{2} \sqrt{\sigma_Y^2/n}} \times \frac{1}{\sqrt{\left(1 + \frac{\sigma_Y^2/n}{\sigma_\delta^2}\right)}}$

The usual z-score  $\times$  shrinkage factor  $< 1$  (correction from partial pooling)

The posterior means are pulled together

[Plot]



When  $\sigma_\delta^2 \ll \sigma_Y^2$ , the groups are homogeneous  $\Rightarrow$  large shrinkage.

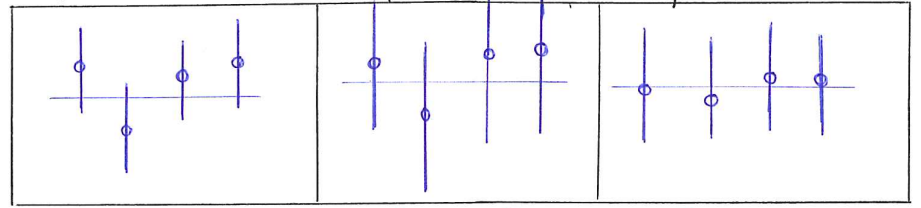
Conclusion: Frequentist vs Bayesian

- z-scores untouched
- rejection region adjusted
- z-scores shrunked towards 0
- no adjustment of the critical region (or coverage of the credible bands)

[Freq]

[Freq]

[Bay]



95% confidence int.

95% + Bonferroni

95% HDI with partial pooling

Credible intervals are larger

Point Estimates are pulled closer together.

## References

- x Baranchick, A. (1964). Multiple regression and estimation of the mean of a multivariate normal distribution. Technical Report No 51, Department of Statistics, Stanford University.
- x Baranchick, A. (1970). A family of minimax estimators of the mean of a multivariate normal distribution. Ann. Math. Stat., 41, p. 642-645.
- x Bernardo, J.M. (1979). Reference Posterior Distributions for Bayesian Inference. JRSS B, vol 41, no 2, p. 113-147.
- x Bock, M.E. (1975). Minimax Estimators of the Mean of a Multivariate Normal Distribution. Ann. of Statistics, vol 3, no 1, pp. 209-218.
- x Brown, L.D. (1977). Admissible Estimators, Recurrent Diffusions, and Insoluble Boundary Value Problems. Ann. Math. Stats, vol 42, pp. 855-903.
- x Casella G. & Berger R.L. (1987). Reconciling Bayesian and Frequentist Evidence in the One-Sided Testing Problem. JASA, vol 82, no 397, pp. 106-111.
- x Edwards W., Lindman H. & Savage L.J. (1963). Bayesian Statistical Inference for Psychological Research. Psychological Review, vol 70, no 3, pp. 193-242.
- x Efron B. (2010). Large Scale Inference: Empirical Bayes Methods for Estimation, Testing and Prediction. Lecture Notes. Stanford University.

- x Efron B. & Morris C. (1977). Stein's Paradox in Statistics. Scientific American, vol 236, no 5, pp. 119-127.
- x Gelman A. et al (2014). Bayesian Data Analysis, 3rd Edition. CRC-Press.
- x Gelman A., Hill J. & Yajima M. (2012). Why we (usually) don't have to worry about multiple comparisons. Journal of Research on Educational Effectiveness. Vol 5, Issue 2, pp 189-211
- x Gönen M., Johnson W.O, Lu Y. & Westfall P. (2005). The Bayesian two-sample t-test. The American Statistician. vol 59, pp 252-257.
- x James W. & Stein C. (1961). Estimation with quadratic loss, Proc. 4th Berkeley Symp. Math. Stat. Prob, vol 1, pp. 361-379.
- x Jeffreys H. (1961). Theory of Probability, 3rd Ed. Oxford Classic Texts in the Physical Sciences.
- x Johnson V. & Rossell D. (2010). On the use of non-local prior densities in Bayesian hypothesis tests. JRSS B. vol 72, Part 2, pp. 143-170.
- x Kass R.E. & Vaidyanathan (1992). Approximate Bayes Factors and Orthogonal Parameters, with Application to Testing Equality of two Binomial Proportions. JRSS B, vol 54, no 1, pp. 129-144.
- x Kleijn B.J.K. and van der Vaart A.W. (2012). The Bernstein-Von-Mises Theorem under misspecification. Electronic-J. Stat, vol 6, pp. 354-381.

- x Rouder J.N (2014). Optional Stopping: No Problem for Bayesians. Psychonomic Bulletin & Review, vol 21, pp. 301-308.
- x Schönbrodt et al (2017). Sequential Hypothesis Testing with Bayes Factors = Efficiently Testing Mean Differences, vol 22, no 2, pp. 322-339. in Psychol. Methods.
- x Sellke T., Bayarri M.J. & Berger J.O. (2001). Calibration of p-values for Testing Precise Null Hypotheses. The American Statistician, vol 55 pp. 62-71.
- x Shi H. & Yin G. (2019). Control of Type I. Error Rates in Bayesian sequential Design. Bayesian Analysis, vol 14, no 2, pp. 391 - 425.
- x Van der Vaart (1998). Asymptotic Statistics. Cambridge University Press.
- x White H. (1982). Maximum likelihood Estimation of Misspecified Models, Econometrica, vol 50, no 1, pp. 1-25.