

# Mathematical Statistics II

STA2212H S LEC9101

Week 2

January 17 2023



Sign in



## Two-thirds of world's glaciers expected to disappear by end of the century, study in Science journal says

SETH BORENSTEIN

But if the world can limit future warming to just a few more tenths of a degree and fulfill international goals – technically possible but unlikely according to many scientists – then slightly less than half the globe's glaciers will disappear, said the same study. Mostly small but wellknown glaciers are marching to extinction, study authors said.

In an also unlikely worst-case scenario of several degrees of warming, 83 per cent of the world's glaciers would likely disappear by the year 2100, study authors said.

The study, published Thursday in the journal *Science*, examined all of the globe's 215,000 landbased glaciers – not counting those on ice sheets in Greenland and Antarctica – in a more comprehensive way than past studies. Scientists



Tourists hike to visit the Nigardsbreen glacier in Jostedal, Norway, last August. Scientists project the planet will lose between 38.7 trillion and 64.4 trillion tonnes of glacial ice by the end of the century.

1. Recap
2. Nonparametric Likelihood [MS 5.6](#)
3. Profile Likelihood
4. Bayesian Estimation [MS 5.8](#)

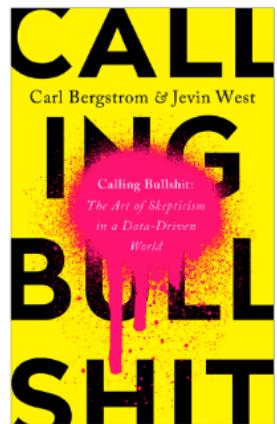
## Upcoming seminars of interest

- **January 23 11.00 –12.00 Jevin West Details**
- “The Art of Skepticism in a Data-Driven World”
- 140 St. George St., 4th floor

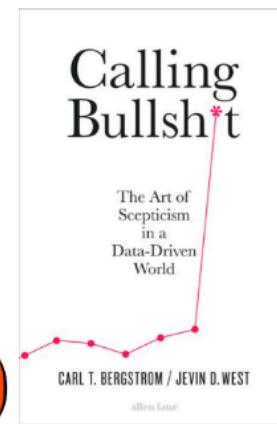


# Calling Bullshit

Data Reasoning in a Digital World



Penguin  
Random House



**Now available!** *Calling Bullshit: The Art of Skepticism in a Data-Driven World*, by Carl Bergstrom and Jevin West. [Available here.](#)

## Recap

iid

- data  $x_1, \dots, x_n$  independent observations; model  $f(\mathbf{x}; \theta) = \prod_i f(x_i; \theta)$ ,  $\theta \in \mathbb{R}$
- limit theorem  $\sqrt{n}(\hat{\theta} - \theta) \xrightarrow{d} N(0, I_1^{-1}(\theta)) = E\{\ell'(\theta; \mathbf{x}_i)^2\}$
- approximation  $\hat{\theta} \sim N\{\theta, I^{-1}(\hat{\theta})\}$ , or  $\hat{\theta} \sim N\{\theta, J^{-1}(\hat{\theta})\}$   
 $I(\theta) = nI_1(\theta)$ ,  $J(\theta) = -\ell''(\theta; \mathbf{x})$   
observed value

## Recap

- data  $x_1, \dots, x_n$  independent observations; model  $f(\mathbf{x}; \theta) = \prod f(x_i; \theta)$ ,  $\theta \in \mathbb{R}$
- limit theorem  $\sqrt{n}(\hat{\theta} - \theta) \xrightarrow{d} N(0, I_1^{-1}(\theta))$
- approximation  $\hat{\theta} \sim N\{\theta, I^{-1}(\hat{\theta})\}$ , or  $\hat{\theta} \sim N\{\theta, J^{-1}(\hat{\theta})\}$   $I(\theta) = nI_1(\theta)$ ,  $J(\theta) = -\ell''(\theta; \mathbf{x})$

- data  $x_1, \dots, x_n$  independent observations; model  $f(\mathbf{x}; \theta) = \prod f(x_i; \theta)$ ,  $\underline{\theta} \in \underline{\mathbb{R}}^p$
- limit theorem  $\sqrt{n}\{I(\theta)\}^{1/2}(\hat{\theta} - \theta) \xrightarrow{d} N_p(\mathbf{0}, I_d)$   $\sqrt{n}(\hat{\theta} - \theta) \xrightarrow{d} N_p(\mathbf{0}, I_1(\theta))$
- approximation  $\hat{\theta} \sim N_p\{\theta, I^{-1}(\hat{\theta})\}$ , or  $\hat{\theta} \sim N_p\{\theta, J^{-1}(\hat{\theta})\}$   $\uparrow p$   
 $\uparrow$  matrix  $\frac{\partial^2 \ell}{\partial \theta \partial \theta^\top}(\hat{\theta}; \bar{x})$

# Recap

- data  $x_1, \dots, x_n$  independent observations; model  $f(\mathbf{x}; \theta) = \prod f(x_i; \theta)$ ,  $\theta \in \mathbb{R}$
- limit theorem  $\sqrt{n}(\hat{\theta} - \theta) \xrightarrow{d} N(0, I_1^{-1}(\theta))$
- approximation  $\hat{\theta} \sim N\{\theta, I^{-1}(\hat{\theta})\}$ , or  $\hat{\theta} \sim N\{\theta, J^{-1}(\hat{\theta})\}$   $I(\theta) = nI_1(\theta)$ ,  $J(\theta) = -\ell''(\theta; \mathbf{x})$

- data  $x_1, \dots, x_n$  independent observations; model  $f(\mathbf{x}; \theta) = \prod f(x_i; \theta)$ ,  $\theta \in \mathbb{R}^p$

- limit theorem  $\sqrt{n}\{I(\theta)\}^{1/2}(\hat{\theta} - \theta) \xrightarrow{d} N_p(\mathbf{0}, I_d)$

- approximation  $\hat{\theta} \sim N_p\{\theta, I^{-1}(\hat{\theta})\}$ , or  $\hat{\theta} \sim N_p\{\theta, J^{-1}(\hat{\theta})\}$

- id. dist'd obs =*
- data  $x_1, \dots, x_n$  independent observations with cdf  $F(\cdot)$
- true model  $F(\mathbf{x}) = \prod F(x_i)$ ,  $\theta \in \mathbb{R}^p$*
- assumed model  $\ell(\theta; \mathbf{x})$ ,*

- limit theorem  $\sqrt{n}\{\hat{\theta} - \theta(F)\} \xrightarrow{d} N_p\{\mathbf{0}, J^{-1}(F)I(F)J^{-1}(F)\}$

pdf  $\hat{F}(\cdot)$

$$\boxed{\ell'(\hat{\theta}; \mathbf{x}) = 0}$$

$\theta(F), I(F), J(F)$

$$J(F) = \sum_{i=1}^n \ell''(\theta_F; x_i)$$

$$I(F) = \sum_{i=1}^n \ell'(x_i; \theta_F)^2$$

## ... Recap

- proof requires many smoothness conditions on underlying model
- i.i.d. can often be weakened to independent (not i.d.) observations, or even dependent AR(1), AR(p), GLMs (regression) need WLLN and CLT
- MS Theorem 5.3, p.253 has a careful proof for  $\theta \in \mathbb{R}$

- key step is

$$\ell'(\hat{\theta}; \tilde{x}) = 0 \\ = \dots \text{TS}$$

$$\sqrt{n}(\hat{\theta} - \theta) = \frac{n^{-1/2} \sum_{i=1}^n \ell'(X_i; \theta)}{-n^{-1} \sum_{i=1}^n \ell''(X_i; \theta) + (\hat{\theta} - \theta)(2n)^{-1} \sum_{i=1}^n \ell'''(X_i; \theta^*)}$$

CLT      WLLN

$E_\theta \ell''(x_i; \theta)$       0       $R_n$

$(\theta^* - \theta) < |\hat{\theta} - \theta|$

see also MSI, Nov 29, likelihood handout

$$\hat{\theta} \xrightarrow{P} \theta$$

$$\Rightarrow \theta^* \xrightarrow{P} \theta$$

## ... Recap

- proof requires many smoothness conditions on underlying model
- i.i.d. can often be weakened to independent (not i.d.) observations, or even dependent need WLLN and CLT
- MS Theorem 5.3, p.253 has a careful proof for  $\theta \in \mathbb{R}$  see also MSI, Nov 29, likelihood handout
- key step is

$$\sqrt{n}(\hat{\theta} - \theta) = \frac{-n^{-1/2} \sum_{i=1}^n \ell'(X_i; \theta)}{n^{-1} \sum_{i=1}^n \ell''(X_i; \theta) + (\hat{\theta} - \theta)(2n)^{-1} \sum_{i=1}^n \ell'''(X_i; \theta^*)}$$

- vector version is

$$\sqrt{n} \sum_{k=1}^p (\hat{\theta}_k - \theta_k) \left\{ n^{-1} \ell''_{jk}(\hat{\theta}) + (2n)^{-1} \sum_{l=1}^p (\hat{\theta}_l - \theta_l) \ell'''_{jkl}(\theta^*) \right\} = -n^{-1/2} \ell'_j(\theta),$$

$\uparrow \quad \underbrace{\quad}_{j = 1, \dots, p}$

$$\ell'(\hat{\theta}) = \underline{\Omega} \simeq \ell'(\theta) + \dots$$

$\frac{\partial^2 \ell}{\partial \theta_j \partial \theta_k}$

# Loose ends

see also MS I Nov 29

AoS chg

~~A1 = A6~~

B1 - B6

$\theta^*$

} AoS

- proof of consistency (Thms 5.1,2) uses WLLN applied to

$\theta$  fixed (true)  
 $t$  varying

$$\phi(t) = E_\theta \log \frac{f(X_i; t)}{f(X_i; \theta)} \equiv -K(f_t : f_\theta)$$

$$\phi_n(t) \xrightarrow{\text{P}} \phi(t) \quad \text{for any } t \in \mathbb{R}^p$$

$$-K(f_t : f_\theta) = 0 \quad \text{if } t = \theta$$

$$\text{o.w. } -K < 0$$

$$K(f_{\theta^*})$$

$$M_n(\theta) = \frac{1}{n} \sum_{i=1}^n \log \frac{f(X_i; \theta)}{f(X_i; \theta_{\text{true}})}$$

$$E_{\theta_{\text{true}}} \log \frac{f(X_i; \theta)}{f(X_i; \theta_{\text{true}})} \equiv -D(\theta_{\text{true}}, \theta)$$

$$\sup_t \phi(t) = \phi(\theta)$$

$$\sup_{\theta \in \Theta} |M_n(\theta)| < C$$

$$\sup_t \phi_n(t) \not\equiv \text{determines } \hat{\theta}$$

$$\hat{\theta} \subseteq \mathbb{R}$$

$$\sup_{\{\mu, \sigma^2 \in \mathbb{R}^n \times \mathbb{R}^+\}} \left[ \ln(p_{\mu, \sigma^2}(x)) \right]$$

## Loose ends

$$\hat{\mu}_{i,n} = \frac{1}{2}(X_i + Y_i) \quad i=1, \dots, n \quad \text{see also MS I Nov 29}$$

$\hat{\mu}_{i,n} \xrightarrow{\text{b.c. or } n \rightarrow \infty} \mu_i$  doesn't change

- proof of consistency (Thms 5.1,2) uses WLLN applied to

$$\phi_n(t) = \frac{1}{n} \sum_{i=1}^n \log \frac{f(X_i; t)}{f(X_i; \theta)}$$

$$\phi(t) = E_\theta \log \frac{f(X_i; t)}{f(X_i; \theta)} \equiv -K(f_t : f_\theta)$$

$$M_n(\theta) = \frac{1}{n} \sum_{i=1}^n \log \frac{f(X_i; \theta)}{f(X_i; \theta_{true})}$$

$$E_{\theta_{true}} \log \frac{f(X_i; \theta)}{f(X_i; \theta_{true})} \equiv -D(\theta_{true}, \theta)$$

$\trianglelefteq$

$\downarrow$   
true

- by Jensen's  $\phi(t)$  maximized at  $\theta$ , which suggests  $\hat{\theta} \rightarrow \theta$  but functions are tricky
- need sup condition, see Thm 5.1 (a)

- proof of consistency (Thms 5.1,2) uses WLLN applied to

$$\phi_n(t) = \frac{1}{n} \sum_{i=1}^n \log \frac{f(X_i; t)}{f(X_i; \theta)}$$

$$\phi(t) = E_\theta \log \frac{f(X_i; t)}{f(X_i; \theta)} \equiv -K(f_t : f_\theta)$$

$$M_n(\theta) = \frac{1}{n} \sum_{i=1}^n \log \frac{f(X_i; \theta)}{f(X_i; \theta_{true})}$$

$$E_{\theta_{true}} \log \frac{f(X_i; \theta)}{f(X_i; \theta_{true})} \equiv -D(\theta_{true}, \theta)$$

$p$  fixed  
 $n \rightarrow \infty$

$$\int \log \left\{ \frac{f_0(x)}{f(x)} \right\} f_0(x) dx$$

but functions are tricky

- by Jensen's  $\phi(t)$  maximized at  $\theta$ , which suggests  $\hat{\theta} \rightarrow \theta$

- need sup condition, see Thm 5.1(a) *also As Chapter 9*

$\downarrow$   $\downarrow$  base density

- $K(f : f_0)$  Kullback-Leibler divergence measures 'closeness' of densities  $f$  and  $f_0$

$$E_0 \{ \log \{ f_0(x) / f(x) \} \}$$

- maximum likelihood estimator minimizes K-L divergence between empirical cdf and model

$$E_{F_n} \log \{ dF_n(\mathbf{x}) / f_\theta(\mathbf{x}) \}$$

# Nonparametric MLE

$x_i$  assume  $\sim f(x_i; \theta)$   $\theta \in \mathbb{R}^P$   $\oplus \subseteq \mathbb{R}^P$  MS 5.6

- sample  $x_1, \dots, x_n$  independent, identically distributed, with cdf  $F$

*if we use this*

- likelihood function  $L(f) = \prod_{i=1}^n f(x_i)$

- assume solution puts mass only at  $x_1, \dots, x_n$

- log-likelihood function  $\ell(p) = \sum_{i=1}^n \log(p_i)$

$$\sum p_i = 1 \quad 0 < p_i \leq 1 \\ - \quad \text{(ntbc)}$$

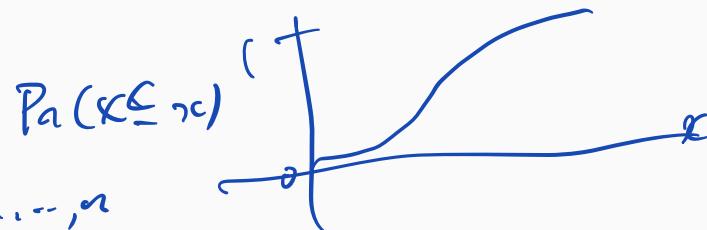
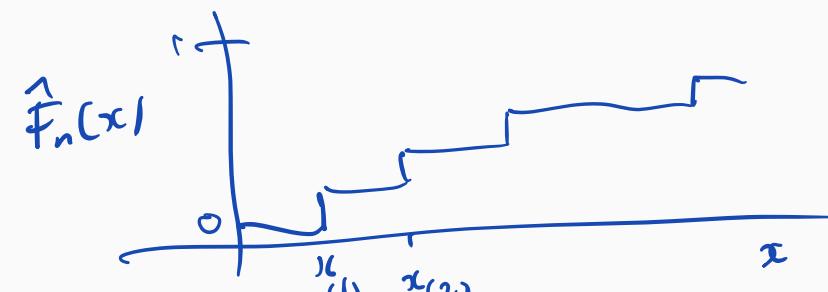
$$\max_p \sum_{i=1}^n \log(p_i) ; \text{ s.t. } \sum p_i = 1 \\ 0 < p_i < 1$$

$$p_i = \frac{1}{n} \quad i = 1, \dots, n$$

*assume  $F$  has density*  
no parametric model assumed

$$\hat{F}_n(x) = \frac{1}{n} \sum 1\{x_i \leq x\}$$

e c d f



- sample  $x_1, \dots, x_n$  independent, identically distributed, with cdf  $F$

no parametric model assumed

- likelihood function  $L(F) = \prod f(x_i)$

$$f_x(\underline{x}) = \ell(\underline{p}) + \lambda(\sum p_i - 1)$$

- assume solution puts mass only at  $x_1, \dots, x_n$

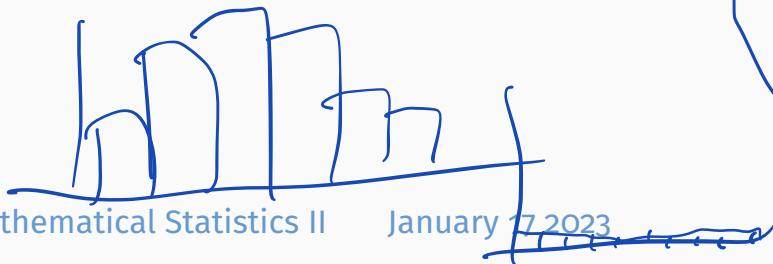
$$\frac{\partial}{\partial \lambda} \quad \frac{\partial}{\partial p_1} \quad \dots \quad \frac{\partial}{\partial p_n}$$

- log-likelihood function  $\ell(p) = \sum_{i=1}^n \log(p_i)$

- maximized at  $p_i = 1/n, i = 1, \dots, n$

Lagrange

- gives empirical cdf



$$\hat{F}_n(x) = \frac{1}{n} \sum_{i=1}^n \mathbf{1}(X_i \leq x)$$

$\hat{F}_n(\cdot)$  is "MLE" of  $F(\cdot)$

$$\hat{f}_n(x) = \frac{1}{nh} \sum_{i=1}^n K\left(\frac{x_i - x}{h}\right) \quad ?? \text{ntbc}$$

# Multi-parameter example: logistic regression

```
Boston$crim2 <- Boston$crim > median(Boston$crim) # define binary response
Boston.glm <- glm(crim2 ~ . - crim, family = binomial,
                    data = Boston) # fit logistic regression
```

summary(Boston.glm)

Coefficients:

	Estimate	Std. Error	z value	Pr(> z )	
(Intercept)	-34.103704	6.530014	-5.223	1.76e-07 ***	
zn	-0.079918	0.033731	-2.369	0.01782 *	
indus	-0.059389	0.043722	-1.358	0.17436	
chas	0.785327	0.728930	1.077	0.28132	
nox	48.523782	7.396497	6.560	5.37e-11 ***	
rm	-0.425596	0.701104	-0.607	0.54383	
age	0.022172	0.012221	1.814	0.06963 .	

$$\sqrt{[i^{-1}(\hat{\beta})]_{jj}} \quad j=1, \dots, p$$

$N(0, 1)$  if  $\beta_j = 0$

Wald st.  
 $\hat{\theta}_j \sim N(\theta_j, \dots)$

$$\partial l(\hat{\theta}) = 0$$

$$\frac{\partial l(\hat{\theta})}{\partial \theta_1} = 0$$

⋮

$$\frac{\partial l(\hat{\theta})}{\partial \theta_p} = 0$$

## ... Example: logistic regression

$$\hat{\beta}_j \pm t.96(\hat{s}_{\epsilon_j})$$

```
Boston.glm <- glm(crim2 ~ . - crim, family = binomial,
                     data = Boston) #fit logistic regression
```

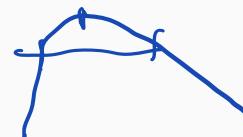
confint(Boston.glm)

Waiting for profiling to be done...

2.5 % 97.5 %

	(Intercept)	-47.480389822	-21.699753794
zn	-0.152359922	-0.020567540	
indus	-0.149113408	0.024168460	
chas	-0.646429219	2.233443233	
nox	34.967619055	64.088411260	
rm	-1.811639107	0.950196261	
age	-0.001231256	0.046865843	
dis	0.280762523	1.140619391	
rad	0.376833861	0.975898274	
tax	-0.012038221	-0.001324887	

$\ell_p(4)$



$\approx 95\% CI$

for  $\beta_j$

$-.08 \pm .06$

$-0.02 \pm (4)$

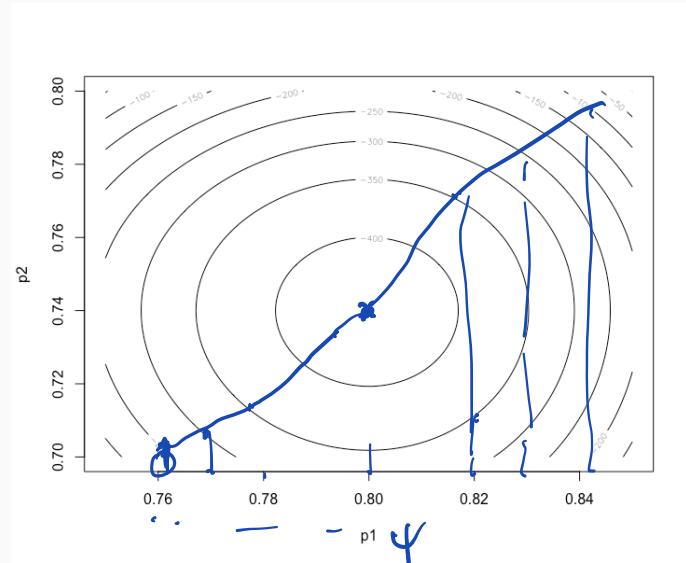
$(-.14, -.02)$

Prof.

$\hat{\beta}_j \pm \dots$

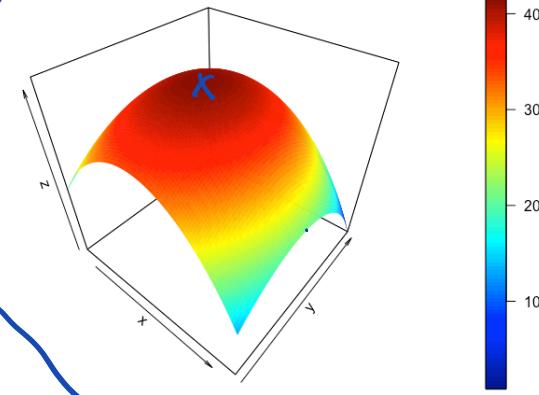
# Multi-parameter setting

AoS §9.10



$$\hat{p}_1 = \frac{160}{200}$$

$$\hat{p}_2 = \frac{180}{200}$$



$$\frac{\partial l(\hat{\theta})}{\partial \theta} = 0$$

$$p_{x1} \quad p_{x1}$$

$Y_1 \sim \text{Binom}(n_1, p_1)$ ,  $Y_2 \sim \text{Binom}(n_2, p_2)$ , independently  
observed values  $y_1 = 160, n_1 = 200, y_2 = 180, n_2 = 200$

$$Y_1 \sim \text{Binom}(200, p_1)$$

$$\frac{\partial l}{\partial p_1} = 0 \quad \dots \quad \hat{p}_1(p_2)$$

$$l(p_1, p_2) = \text{lf}_{y_1}(p_1) p_1^{y_1} (1-p_1)^{n_1-y_1} + \text{lf}_{y_2}(p_2) p_2^{y_2} (1-p_2)^{n_2-y_2}$$

## ... Multi-parameter setting

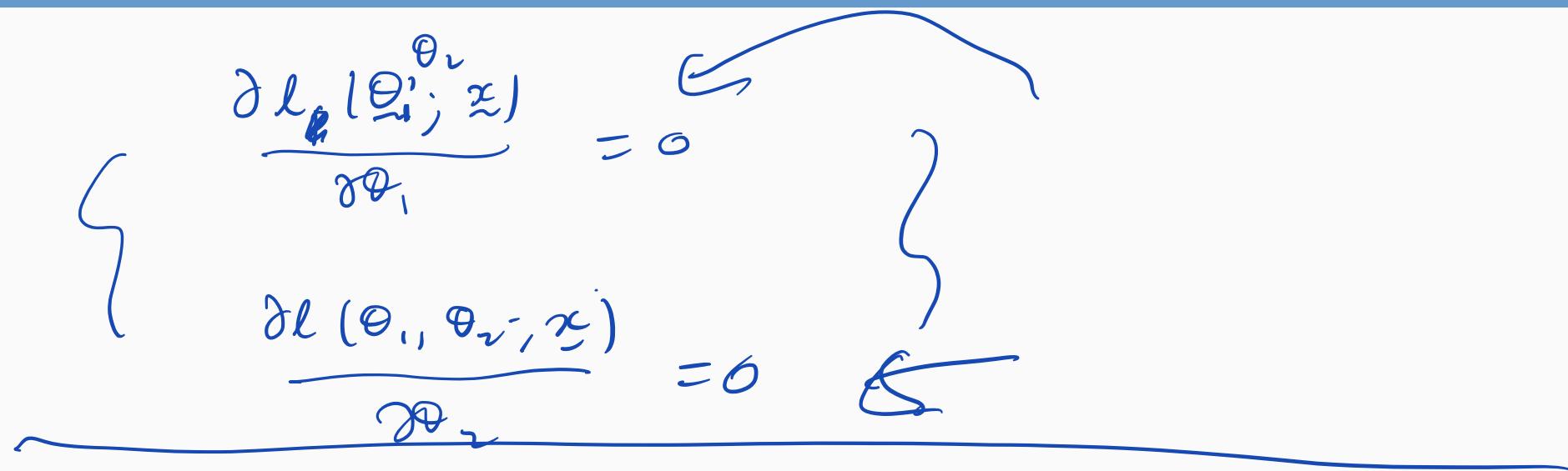
AoS §9.10

$$\frac{\partial l(\theta_1; \underline{x})}{\partial \theta_1}$$

$$= 0$$

$$\frac{\partial l(\theta_1, \theta_2; \underline{x})}{\partial \theta_2}$$

$$= 0$$



## Profile likelihood function

often  $\underline{\theta} = (\psi, \underline{\lambda})$

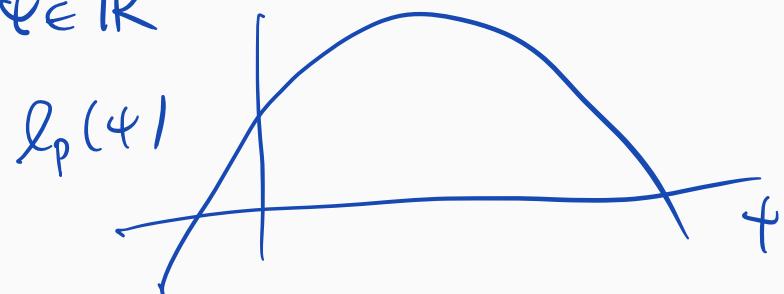
$\psi$  - parameter of interest ( $\sigma^2$ )

$\underline{\lambda}$  - nuisance parameters ( $\mu_1, \dots, \mu_n$ )

define  $l_p(\psi) = l(\psi, \hat{\lambda}_\psi)$

$$\frac{\partial l}{\partial \lambda} (\psi, \hat{\lambda}_\psi) = 0$$

if  $\psi \in \mathbb{R}$



concentrated lik.

## ... Profile likelihood function

$$l_p(\psi) = \underline{l}(\psi, \hat{\lambda}_\psi)$$

$$\underline{l'_p(\tilde{\psi})} = 0$$

$$\tilde{\psi} \equiv \hat{\psi}$$

$$0 = \left. l'_p(\tilde{\psi}) \right|_{\psi=\tilde{\psi}} = \frac{\partial}{\partial \psi} \{ l(\psi, \hat{\lambda}_\psi) \} + \frac{\partial}{\partial \lambda} l(\psi, \hat{\lambda}_\psi) \frac{d\hat{\lambda}_\psi}{d\psi}$$

by def. of  $\hat{\lambda}_\psi$

$$\begin{cases} \frac{\partial l}{\partial \lambda}(\tilde{\psi}, \hat{\lambda}_{\tilde{\psi}}) = 0 \\ \frac{\partial l}{\partial \psi}(\tilde{\psi}, \hat{\lambda}_{\tilde{\psi}}) = 0 \end{cases} \quad \equiv \quad \left( \begin{array}{c} \frac{\partial l(\psi, \lambda)}{\partial \psi} \\ \frac{\partial l(\psi, \lambda)}{\partial \lambda} \end{array} \right) \Big|_{(\hat{\psi}, \hat{\lambda})} = 0$$

Profile log-lik. work like log-lik  
(asymptotically)

theorems

p fixed

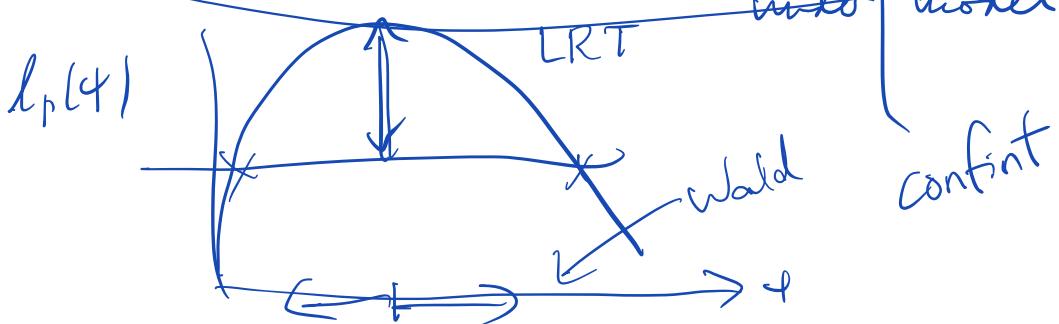
$n \rightarrow \infty$

(approx<sup>=</sup> may  
be poor, p large)

$$(i) \quad l_p'(\hat{\psi}) = 0 \quad \hat{\psi} \text{ m.l.e.}$$

$$(ii) \quad \widehat{\text{a.var.}}(\hat{\psi}) = -l_p''(\hat{\psi}) \quad \begin{array}{l} \text{obs'd} \\ \text{Fisher info} \\ \text{in profile} \end{array}$$

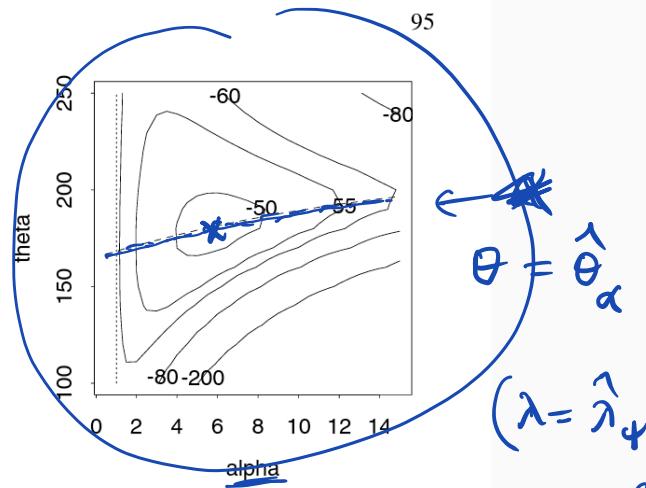
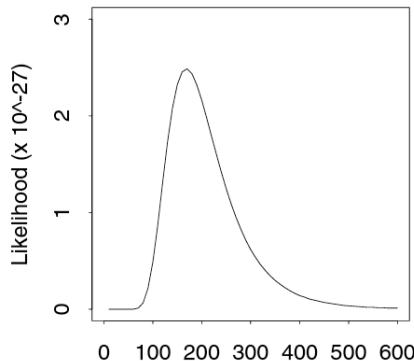
$$(iii) \quad 2\{l_p(\hat{\psi}) - l_p(\psi)\} \sim \chi^2_1$$



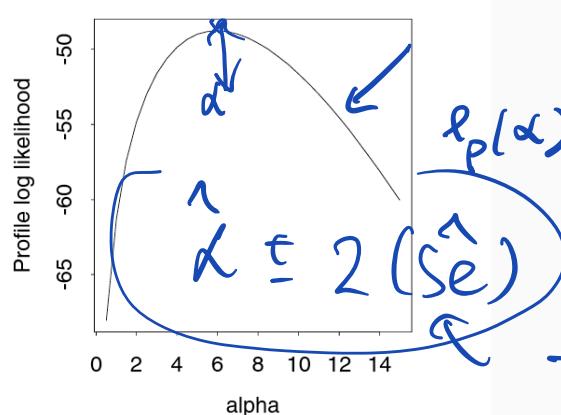
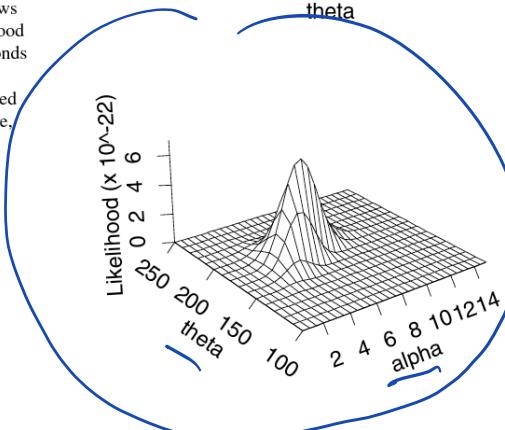
# ... Profile likelihood function

## 4.1 · Likelihood

**Figure 4.1** Likelihoods for the spring failure data at stress  $950 \text{ N/mm}^2$ . The upper left panel is the likelihood for the exponential model, and below it is a perspective plot of the likelihood for the Weibull model. The upper right panel shows contours of the log likelihood for the Weibull model; the exponential likelihood is obtained by setting  $\alpha = 1$ , that is, slicing  $L$  along the vertical dotted line. The lower right panel shows the profile log likelihood for  $\alpha$ , which corresponds to the log likelihood values along the dashed line in the panel above, plotted against  $\alpha$ .



$$\ell_p(\hat{\alpha}) = \ell(\hat{\alpha}, \hat{\lambda}_4)$$



$$\ell_p(\hat{\alpha}) - \ell_p(\alpha)$$

## Bayesian estimation

inference more relevant than point

MS 5.8; AoS 11

estimation

model  $f(x; \theta)$  ← density for the dist<sup>r</sup> of  $X$ , given  $\theta$

$\theta$  viewed as another r.v.  $\begin{cases} \text{④ - random} \\ \theta - \text{fixed} \end{cases}$

prior  $\pi(\theta)$

density = prior sp.

posterior  $\pi(\theta | x) = f(x | \theta) \pi(\theta) / \int \int f(x | \theta) \pi(\theta) d\theta$

sample

$x_1, \dots, x_n$

note  $f(\tilde{x}; \theta) = \prod_{i=1}^n f(x_i; \theta)$

Bayes th.  
cond'l prob.

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

$$\pi(\theta | \tilde{x}) = \frac{\prod_{i=1}^n f(x_i; \theta) \pi(\theta)}{\int \dots d\theta}$$

← marginal for  $\tilde{x}$

# Frequentist and Bayesian contrast

$$L(\theta; \underline{x}) = c(\underline{x}) f(\underline{x}; \theta)$$

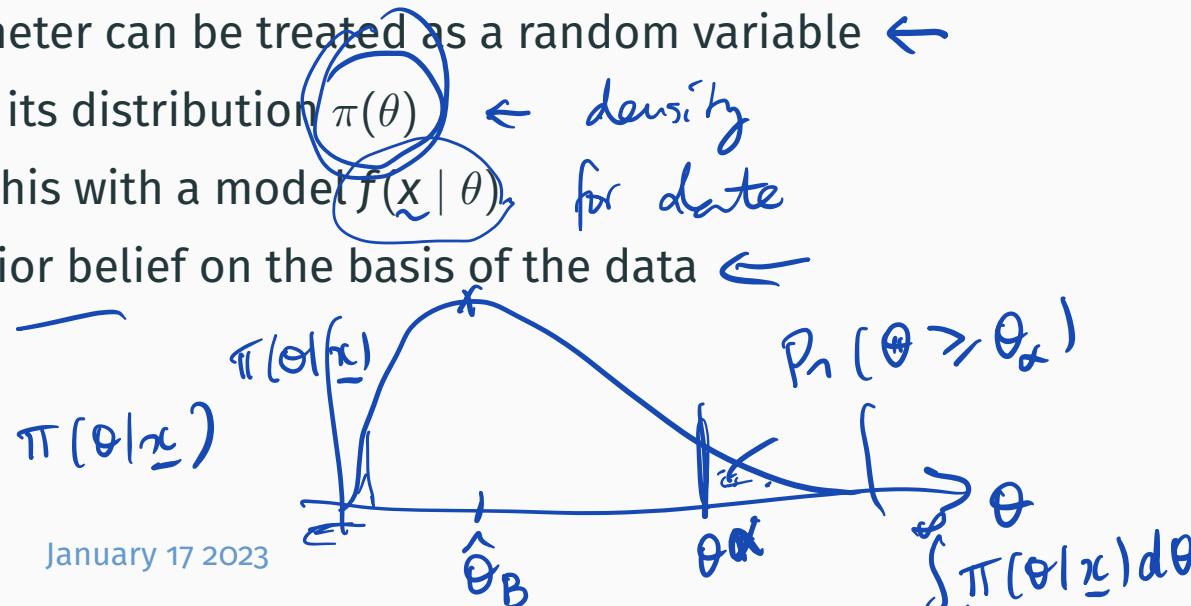
$$L(\theta_1; \underline{x}) / L(\theta_2; \underline{x})$$

Frequentist:

- There is a fixed parameter (unknown) we are trying to learn
- Our methods are evaluated using probabilities based on  $f(x; \theta)$

Bayesian:

- The parameter can be treated as a random variable ←
- We model its distribution  $\pi(\theta)$  ← density
- Combine this with a model  $f(x | \theta)$  for data
- Update prior belief on the basis of the data ←



$$E_{\theta} \{ l'(\theta; \underline{x}) \} = 0$$

$$\int l'(\theta; \underline{x}) f(\underline{x}; \theta) d\underline{x}$$

$$l'(\hat{\theta}; \underline{x}) = 0$$

$$\hat{\theta} \rightarrow \theta$$

↑  
under  $f(\underline{x}; \theta)$

## Example: Binomial

MS 5.26; AoS Ex.11.2

$X_1, \dots, X_n$  i.i.d. Bernoulli ( $\theta$ )

$$\pi(\theta; \underline{\alpha}, \underline{\beta}) = \frac{\Gamma(\alpha + \beta)}{\Gamma(\alpha)\Gamma(\beta)} \theta^{\alpha-1} (1-\theta)^{\beta-1}, 0 < \theta < 1$$

posterior mean, mode

$$f(\underline{x}; \theta) = \prod_{i=1}^n \theta^{x_i} (1-\theta)^{1-x_i} = L(\theta; \underline{x}) = \theta^{\sum x_i} (1-\theta)^{n-\sum x_i}$$

$$0 \leq \theta \leq 1$$

$$\begin{aligned} \pi(\theta | \underline{x}) &= \frac{\theta^{x_+} (1-\theta)^{n-x_+} \cdot \theta^{\alpha-1} (1-\theta)^{\beta-1}}{\int_0^1 \dots d\theta} \quad 0 \leq \theta \leq 1 \\ &= \theta^{x_+ + \alpha - 1} (1-\theta)^{n-x_+ + \beta - 1} / \int_0^1 \dots d\theta \end{aligned}$$

$$\text{Be}(x_+ + \alpha, n - x_+ + \beta)$$

$$E(\text{Beta?}) =$$

$$\frac{\alpha}{\alpha + \beta}$$

# Example: censored exponential

$$\theta_B = E_{\pi(\theta|x)}$$

## MS 5.27

$X_1, \dots, X_n$  i.i.d. Exponential ( $\lambda$ )

censored at  $r$  smallest  $x$ ; let  $Y_i = X_{(i)}$ ,  $i = 1, \dots, r$

$$\pi(\lambda) \sim \text{Exp}(\alpha)$$

$$\alpha = \beta = 1$$

$$\hat{\theta}_B = \frac{\sum x_i + 1}{n + 2}$$

$$f(\mathbf{y} \mid \lambda) = \prod_{i=1}^r \lambda^r \exp(-\lambda y_i) \prod_{i=r+1}^n \exp(-\lambda y_i) = \lambda^r \exp\{-\lambda \sum_{i=1}^r y_i + (n-r)y_r\}$$

$$f(x_i; \lambda) = \lambda^{x_i} e^{-\lambda x_i}$$

$$1 - F(x; \lambda) = e^{-\lambda x}$$

$$\pi(\lambda | \mathbf{y}) = \frac{\lambda^r e^{-\lambda}}{\lambda^r} \times \alpha e^{-\alpha}$$



$$\pi(\theta|x) = \text{Be}(x_+ + \alpha, n - x_+ + \beta)$$

$$E_{\pi(\theta|x)}(\theta) = \frac{x_+ + \alpha}{n + \alpha + \beta}$$

a possible est. of  $\theta$  is

or  $\hat{\theta}_B = \frac{x_+}{n}$

or mode of  $\pi(\theta|x) \equiv \arg \sup \pi(\theta|x)$   
 $\hat{\theta}_{B_m}$

$$\hat{\theta}_B = \left( \frac{x_+}{n} \right) w + \left( \frac{\alpha}{\alpha + \beta} \right) (1-w)$$

$\uparrow$   
sample average

$\downarrow$   
prior exp'd value  
of  $\theta$

$$\alpha = \beta = 1 \quad \pi(\theta) = \begin{cases} 1 & , 0 \leq \theta \leq 1 \\ 0 & \text{o.w.} \end{cases}$$

$$\frac{x_+ + 1}{n + 2}$$

$$f(x; \theta) = \exp\{c(\theta)T(x) - d(\theta) + S(x)\}; \quad \pi(\theta; \alpha, \beta) = K(\alpha, \beta) \exp\{\alpha c(\theta) - \beta d(\theta)\}$$

$$f(x; \theta) = \exp\{c(\theta)T(x) - d(\theta) + S(x)\}; \quad \pi(\theta; \alpha, \beta) = K(\alpha, \beta) \exp\{\alpha c(\theta) - \beta d(\theta)\}$$

Example:  $f(x; \theta) = \theta(1 - \theta)^x, x = 0, 1, \dots; 0 < \theta < 1$

$$f(x; \theta) = \exp\{c(\theta)T(x) - d(\theta) + S(x)\}; \quad \pi(\theta; \alpha, \beta) = K(\alpha, \beta) \exp\{\alpha c(\theta) - \beta d(\theta)\}$$

Example:  $f(x; \theta) = \theta(1 - \theta)^x, x = 0, 1, \dots; 0 < \theta < 1$

Example:  $f(x; \mu) = \frac{1}{\sqrt{2\pi}} \exp\{-\frac{1}{2}(x - \mu)^2\}$

# Example: Bivariate normal

EH §3.1

**Table 3.1** Scores from two tests taken by 22 students, **mechanics** and **vectors**.

	1	2	3	4	5	6	7	8	9	10	11
<b>mechanics</b>	7	44	49	59	34	46	0	32	49	52	44
<b>vectors</b>	51	69	41	70	42	40	40	45	57	64	61
	12	13	14	15	16	17	18	19	20	21	22
<b>mechanics</b>	36	42	5	22	18	41	48	31	42	46	63
<b>vectors</b>	59	60	30	58	51	63	38	42	69	49	63

Table 3.1 shows the scores on two tests, **mechanics** and **vectors**, achieved by  $n = 22$  students. The sample correlation coefficient between the two scores is  $\hat{\theta} = 0.498$ ,

$$\hat{\theta} = \frac{\sum_{i=1}^{22} (m_i - \bar{m})(v_i - \bar{v})}{\sqrt{\left[ \sum_{i=1}^{22} (m_i - \bar{m})^2 \sum_{i=1}^{22} (v_i - \bar{v})^2 \right]^{1/2}}}, \quad (3.10)$$

with  $m$  and  $v$  short for **mechanics** and **vectors**,  $\bar{m}$  and  $\bar{v}$  their averages. We wish to assign a Bayesian measure of posterior accuracy to the true correlation coefficient  $\theta$ , “true” meaning the correlation for the hypothetical population of all students, of which we observed only 22.

If we assume that the joint  $(m, v)$  distribution is bivariate normal (as

$$\begin{pmatrix} x_i \\ y_i \end{pmatrix} \sim N_2 \left( \begin{pmatrix} \mu_x \\ \mu_y \end{pmatrix}, \Sigma \right)$$

$$\Theta = \text{corr} \text{ bet. } x \text{ & } y$$

$$(\Psi)$$

$$\frac{\sigma_{12}}{\sqrt{\sigma_1^2 \sigma_2^2}}$$

# ... Bivariate normal

EH §3.1

marginal for  $\hat{\theta}$  in  $N_2 \left( \begin{pmatrix} \bar{x} \\ \bar{\mu}_2 \end{pmatrix}, \Sigma \right)$

$$f(\hat{\theta} | \theta) = \frac{1}{\pi} (n-2)(1-\theta^2)^{(n-1)/2} (1-\hat{\theta}^2)^{(n-4)/2} \int_0^\infty \frac{1}{\cosh(w) - \theta\hat{\theta}} dw$$

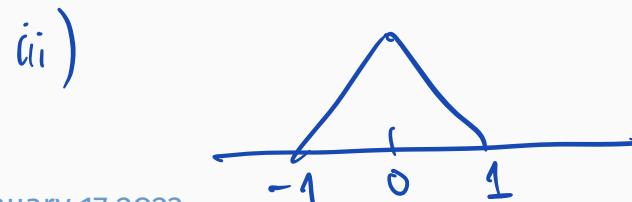
$\uparrow \times \pi(\theta) / \int \dots d\theta$

$$\pi(\theta) \quad \text{prior} \quad -1 \leq \theta \leq 1$$

i)  $\pi(\theta) = \frac{1}{2}$

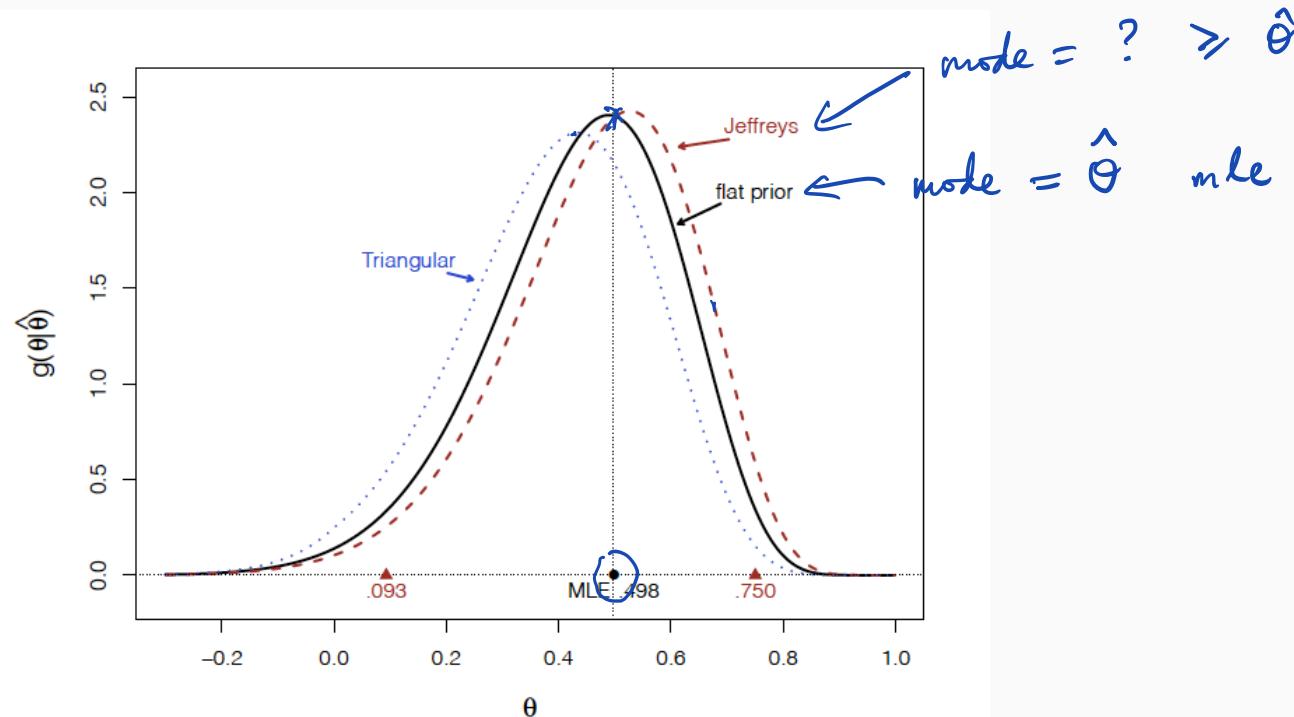
← "ignorance"

ii)  $\pi(\theta) = \frac{1}{1-\theta^2}, \quad -1 \leq \theta \leq 1$  ← Jeffreys'



# Example: Bivariate normal

EH §3.1



**Figure 3.2** Student scores data; posterior density of correlation  $\theta$  for three possible priors.

**Table 11.2** Mortality rates  $r/m$  from cardiac surgery in 12 hospitals (Spiegelhalter *et al.*, 1996b, p. 15). Shown are the numbers of deaths  $r$  out of  $m$  operations.

<i>A</i>	0/47	<i>B</i>	18/148	<i>C</i>	8/119	<i>D</i>	46/810	<i>E</i>	8/211	<i>F</i>	13/196
<i>G</i>	9/148	<i>H</i>	31/215	<i>I</i>	14/207	<i>J</i>	8/97	<i>K</i>	29/256	<i>L</i>	24/360

provided the mode lies inside the parameter space. Here  $\tilde{J}(\theta)$  is the second derivative matrix of  $-\tilde{\ell}(\theta)$ . This expansion corresponds to a posterior multivariate normal

prior for hospital A  $\text{Beta}(1, 1)$

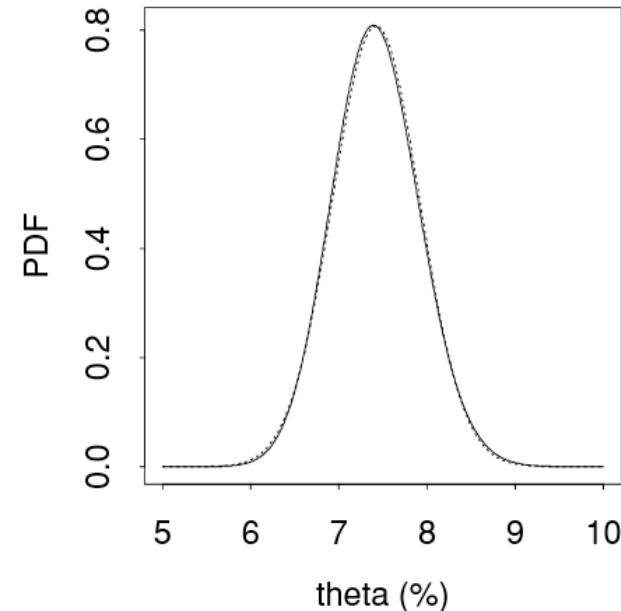
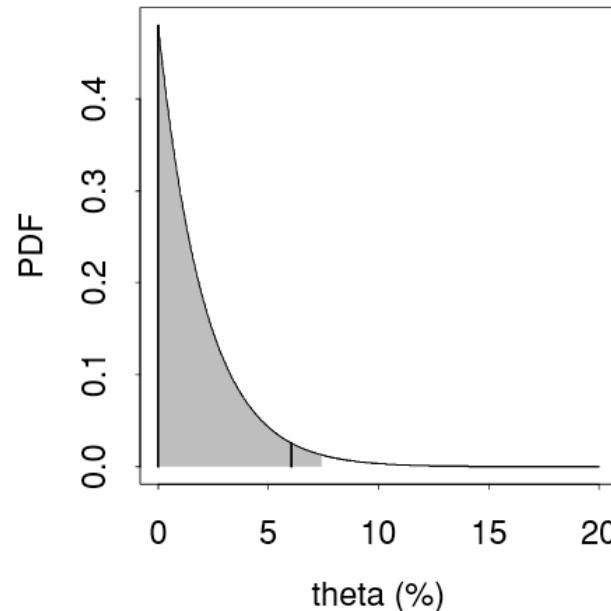
posterior mean

# Example: Binomial

SM Ex.11.11

580

11 · Bayesian Models



**Figure 11.1** Cardiac surgery data. Left panel: posterior density for  $\theta_A$ , showing boundaries of 0.95 highest posterior credible interval (vertical lines) and region between posterior 0.025 and 0.975 quantiles of  $\pi(\theta_A | y)$  (shaded). Right panel: exact posterior beta density for overall mortality rate  $\theta$  (solid) and normal approximation (dots).

put all hospitals together; 208 failures ‘



# Marginalization

# Not all likelihood functions are regular

Example:  $X_1, \dots, X_n$  i.i.d.  $U(0, \theta)$

## ... Not all likelihood functions are regular

MS Exercise 5.1

$$X_1, \dots, X_n \text{ i.i.d. } f(x; \theta) = a(\theta_1, \theta_2)h(x), \quad \theta_1 \leq x \leq \theta_2$$