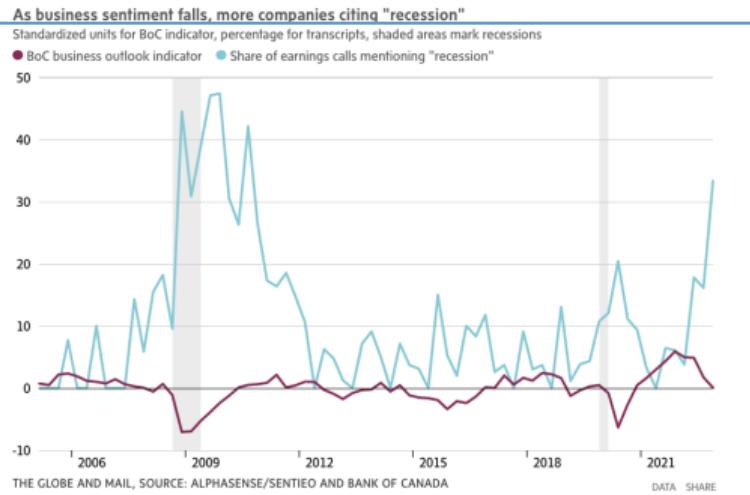


# Mathematical Statistics II

STA2212H S LEC0101

Week 4

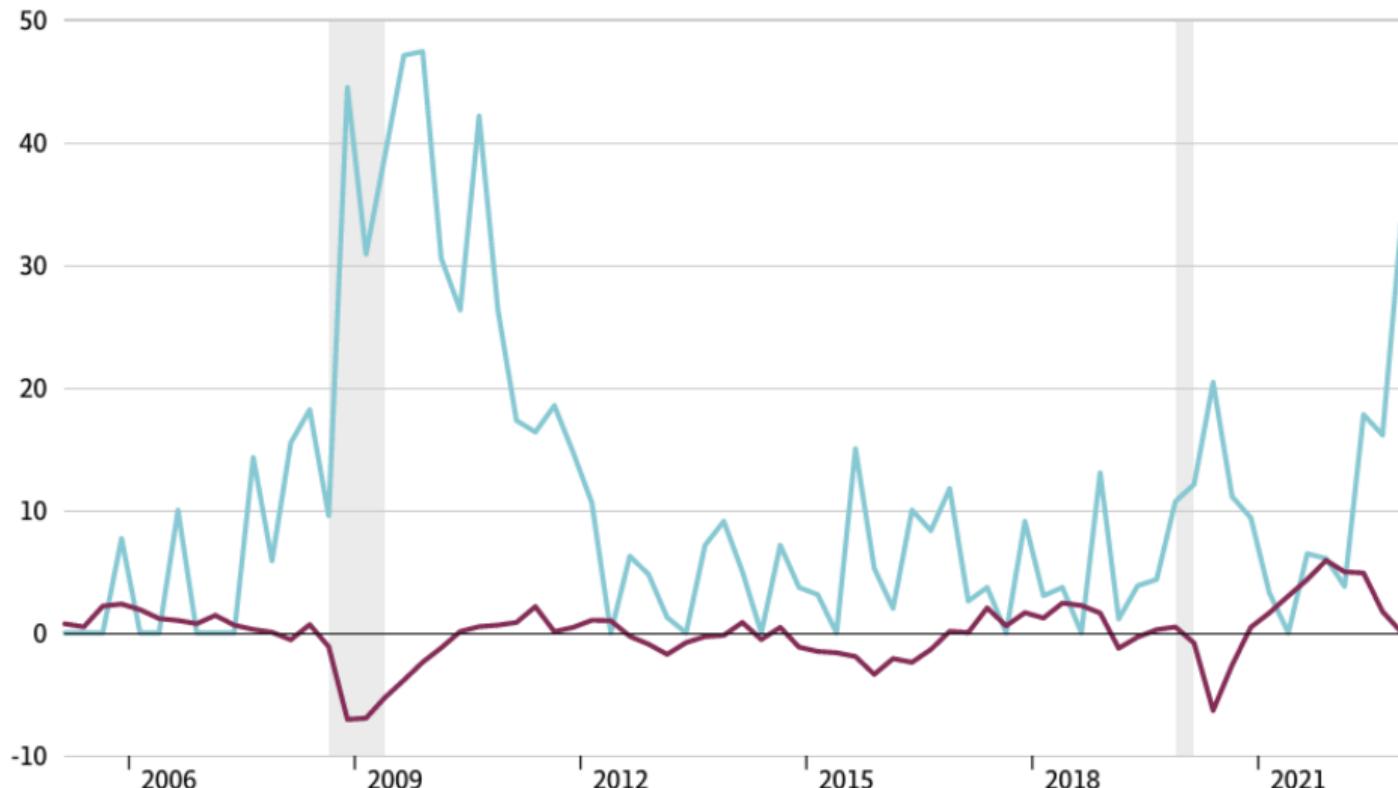
January 31 2023



## As business sentiment falls, more companies citing "recession"

Standardized units for BoC indicator, percentage for transcripts, shaded areas mark recessions

● BoC business outlook indicator ● Share of earnings calls mentioning "recession"



## As business sentiment falls, more companies citing "recession"

Standardized units for BoC indicator, percentage for transcripts, shaded areas mark recessions

● BoC business outlook indicator ● Share of earnings calls mentioning "recession"

50

the portion of companies using the word “recession” in their earnings calls rising at the fastest pace in 15 years

40

30

20

10

0

link

-10

2006

2009

2012

2015

2018

2021

THE GLOBE AND MAIL, SOURCE: ALPHASENSE/SENTIEO AND BANK OF CANADA

DATA SHARE

1. Recap
2. Bayesian hierarchical modelling SM 11.4, AoS 24.5
3. Multi-parameter posteriors AoS 11.7, SM, 11.1-3
4. Interval estimation MS 7.1,2
5. H3: project

## Recap

- priors for Bayesian inference: conjugate, Jeffreys', flat, convenience, weakly informative, [hierarchical](#), matching
- optimality in estimation: Cramer-Rao lower bound:  $\text{Var}\{S(\mathbf{X})\} \geq g'(\theta)^2 / \{nI(\theta)\}$
- matrix version, equality

## Recap

- priors for Bayesian inference: conjugate, Jeffreys', flat, convenience, weakly informative, hierarchical, matching
- optimality in estimation: Cramer-Rao lower bound:  $\text{Var}\{S(\mathbf{X})\} \geq g'(\theta)^2 / \{nI(\theta)\}$
- matrix version, equality
- maximum likelihood estimators are “BAN” asymptotic relative efficiency
- so are other regular estimators with continuous (in  $\theta$ ) variance functions
- Minimum Variance Unbiased Estimators (MVUE) discussed in MS 6.3
- the theory is elegant, but the application is limited

## Recap 2

- finite sample optimality: loss function, risk function, admissible estimator
- Bayes estimators are admissible proper prior; loss function
- Bayes estimators minimize Bayes risk

- the **Bayes risk** of an estimator is the average of the risk function, over a prior distribution

- 

$$R_B(\hat{\theta}) = \int R_\theta(\hat{\theta})\pi(\theta)d\theta$$

- Optimal **Bayes estimators** minimize the expected posterior loss:

$$\int L\{\hat{\theta}(x), \theta\}\pi(\theta | x)d\theta$$

- Example: squared-error loss  $L(\hat{\theta}, \theta) = (\hat{\theta} - \theta)^2$  need to minimize over  $\hat{\theta}$

$$\int (\hat{\theta} - \theta)^2 \pi(\theta | \mathbf{x})d\theta$$

- solution  $\hat{\theta}(\mathbf{x}) = E(\theta | \mathbf{x})$

posterior mean

- Suppose  $\hat{\theta}$  is a Bayes estimator and is unique
- Suppose we have another estimator  $\tilde{\theta}$  with a smaller frequentist risk function:

$$R_\theta(\tilde{\theta}, \theta) \leq R_\theta(\hat{\theta}, \theta)$$

- The Bayes risk of  $\tilde{\theta}$  is

$$R_B(\tilde{\theta}) = \int$$

- Suppose  $\hat{\theta}$  is a Bayes estimator and is unique
- Suppose we have another estimator  $\tilde{\theta}$  with a smaller frequentist risk function:

$$R_\theta(\tilde{\theta}, \theta) \leq R_\theta(\hat{\theta}, \theta)$$

- The Bayes risk of  $\tilde{\theta}$  is

$$R_B(\tilde{\theta}) = \int$$

- instead of minimizing the average (over  $\pi(\theta)$ ) of the risk function we could

$$\min \max R_\theta(\hat{\theta})$$

Definition §6.2

- such estimators are called **minimax**

## Marginalization

- Bayes posterior carries all the information about  $\theta$ , given  $\mathbf{x}$  by definition
- probabilities for any set  $A$  computed using the posterior distribution
- $\text{pr}(\Theta \in A | \mathbf{x}) =$
- if  $\theta = (\psi, \lambda)$ , ...
- or, if  $\psi = \psi(\theta)$
- in this context, ‘flat’ priors can have a large influence on the marginal posterior

- parameter  $\theta = (\theta_1, \dots, \theta_p)$
- model  $f(x^n | \theta)$ ,  $x^n = (x_1, \dots, x_n)$
- joint posterior

$$\pi(\theta | x^n) \propto f(x^n | \theta) \pi(\theta), \quad \theta \in \mathbb{R}^p$$

- parameter  $\theta = (\theta_1, \dots, \theta_p)$
- model  $f(x^n | \theta)$ ,  $x^n = (x_1, \dots, x_n)$
- joint posterior

$$\pi(\theta | x^n) \propto f(x^n | \theta) \pi(\theta), \quad \theta \in \mathbb{R}^p$$

- marginal posterior for  $\theta_1$

$$\pi_m(\theta_1 | x^n) = \int \pi(\theta | x^n) d\theta_2 \dots d\theta_p$$

- marginal posterior for  $\psi(\theta)$

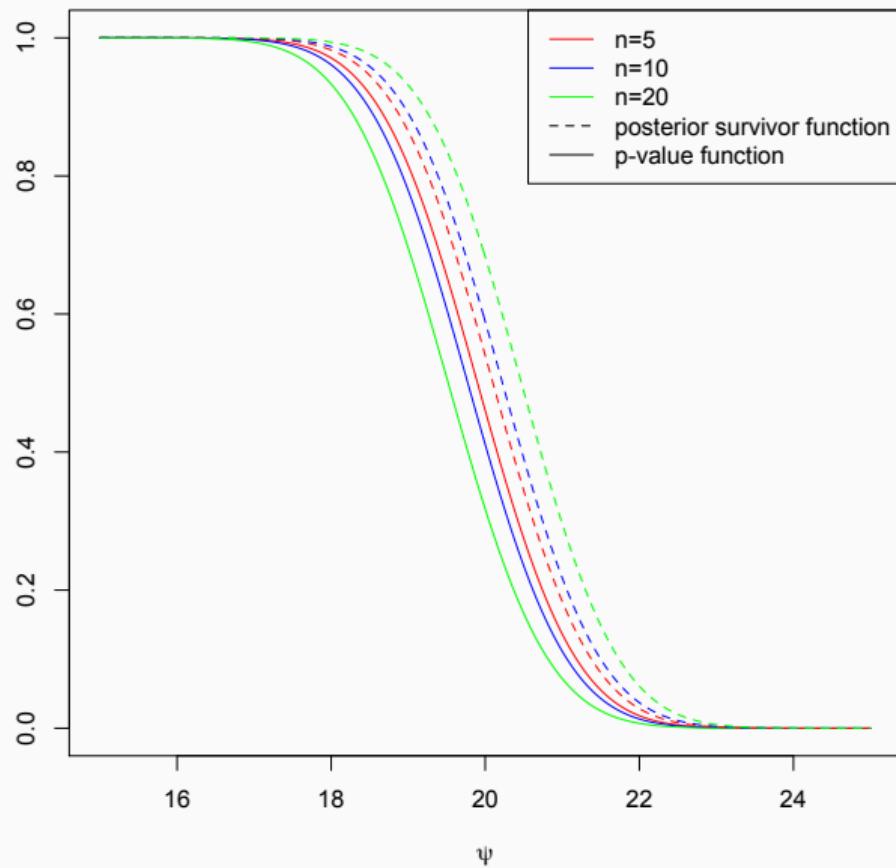
$$\pi_m(\psi | x^n) = \int_{\{\theta : \psi(\theta) = \psi\}} \pi(\theta | x^n) d\theta$$

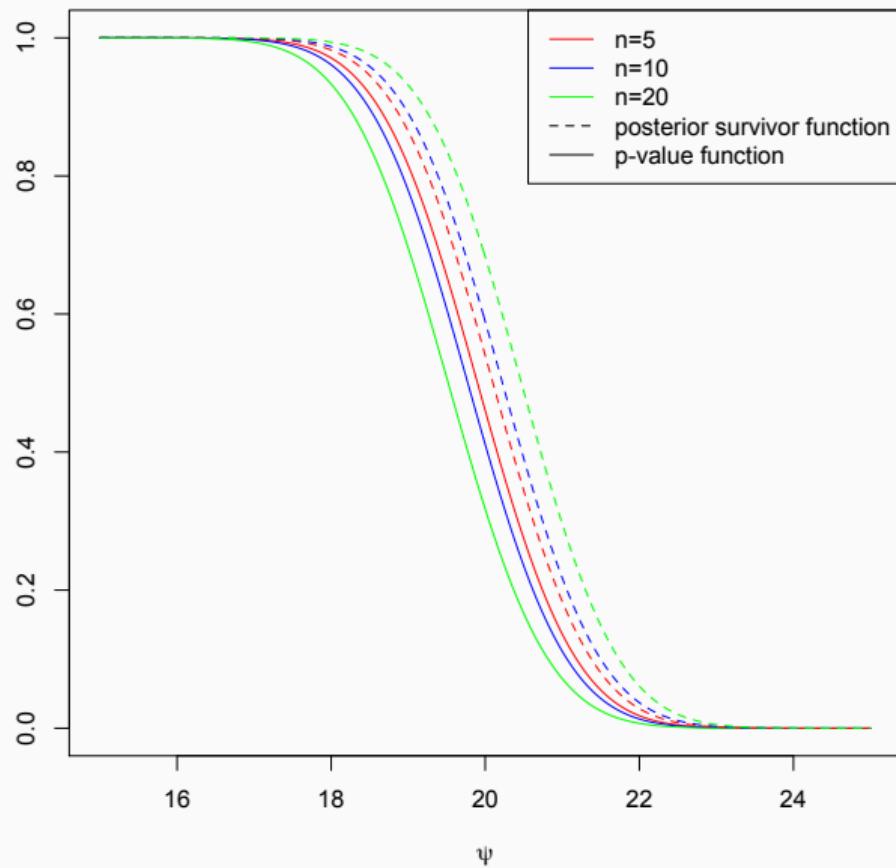
- model:  $x_i \sim N(\mu_i, 1), i = 1, \dots, n$
- prior:  $\pi(\mu) d\mu \propto d\mu$
- posterior  $\pi(\mu | x^n) \propto \prod_{i=1}^n \pi(\mu_i | x_i) = \prod_{i=1}^n \phi(x_i, 1/n)$

- model:  $x_i \sim N(\mu_i, 1), i = 1, \dots, n$
- prior:  $\pi(\mu) d\mu \propto d\mu$
- posterior  $\pi(\mu | x^n) \propto \prod_{i=1}^n \pi(\mu_i | x_i) = \prod_{i=1}^n \phi(x_i, 1/n)$
- $\psi = \sum_{i=1}^n \mu_i^2$  squared length of mean vector

$$\pi(\psi | x^n) = \int_A \pi(\mu | x^n) d\mu$$

- $\mu_i | x_i \sim N(x_i, 1) \implies \sum \mu_i^2 | x^n \sim \chi_n^2(\sum x_i^2)$





- F1 Probability refers to limiting relative frequencies. Probabilities are objective properties of the real world.
- F2 Parameters are fixed, unknown constants. Because they are not fluctuating, no useful probability statements can be made about parameters.
- F3 Statistical procedures should be designed to have well-defined long run frequency properties. For example, a 95 percent confidence interval should trap the true value of the parameter with limiting frequency at least 95 percent.

176      11. Bayesian Inference

B1 Probability describes degree of belief, not limiting frequency. As such, we can make probability statements about lots of things, not just data which are subject to random variation. For example, I might say that “the probability that Albert Einstein drank a cup of tea on August 1, 1948” is .35. This does not refer to any limiting frequency. It reflects my strength of belief that the proposition is true.

B2 We can make probability statements about parameters, even though they are fixed constants.

B3 We make inferences about a parameter  $\theta$  by producing a probability distribution for  $\theta$ . Inferences, such as point estimates and interval estimates, may then be extracted from this distribution.

- $x_i \mid \theta_i \sim N(\theta_i, v_i)$   $v_i$  known
- $\theta_i \mid \mu \sim N(\mu, \sigma^2)$   $\sigma^2$  known
- $\mu \sim N(\mu_0, \tau^2)$  hyperparameters

- $x_i \mid \theta_i \sim N(\theta_i, v_i)$
- $\theta_i \mid \mu \sim N(\mu, \sigma^2)$
- $\mu \sim N(\mu_0, \tau^2)$  hyperparameters
- $\pi(\theta, \mu \mid x)$

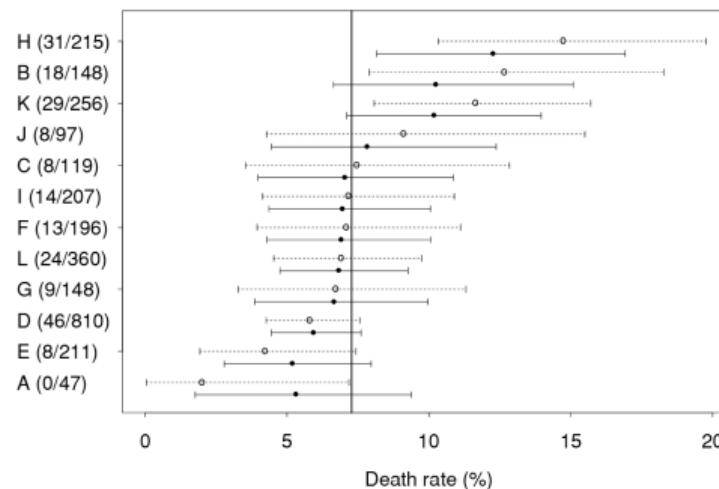
$$E(\mu | x) =$$

$$\text{var}(\mu | x) =$$

$$E(\theta_i | x) =$$

622

II · Bayesian Models



**Figure 11.11** Posterior summaries for mortality rates for cardiac surgery data. Posterior means and 0.95 equitailed credible intervals for separate analyses for each hospital are shown by hollow circles and dotted lines, while blobs and solid lines show the corresponding quantities for a hierarchical model. Note the shrinkage of the estimates for the hierarchical model towards the overall posterior mean rate, shown as the solid vertical line; the hierarchical intervals are slightly shorter than those for the simpler model.

- $$E(\theta_i | x) = x_i \frac{\sigma^2}{\sigma^2 + v_i} + E(\mu | x) \left(1 - \frac{\sigma^2}{\sigma^2 + v_i}\right)$$
- $$E(\mu | x) = \frac{\mu_0 / \tau^2 + \sum x_i / (\sigma^2 + v_i)}{1 / \tau^2 + \sum 1 / (\sigma^2 + v_i)}$$
- If  $\sigma^2$  unknown, then need to sample from the posterior, no closed form available
- Figure 11.11 applies similar ideas, plus sampling from the posterior, in logistic regression

- $X_1, \dots, X_n$  i.i.d.  $f(x; \theta), \theta \in \mathbb{R}$
- a  $100(1 - \alpha)\%$  confidence interval for  $\theta$  is a random interval  $[L(\mathbf{X}), U(\mathbf{X})]$  with

$$\text{pr}_\theta\{L(\mathbf{X}) \leq \theta \leq U(\mathbf{X})\} \geq 1 - \alpha,$$

for all  $\theta$ , with equality for some  $\theta$

- in continuous cases usually can get  $= 1 - \alpha$  for all  $\theta$
- similarly, upper and lower  $(1 - \alpha)$ -confidence bounds:

$$\text{pr}\{\theta \geq L(\mathbf{X})\} = (1 - \alpha); \quad \text{pr}\{\theta \leq U(\mathbf{X})\} = 1 - \alpha$$

- exact limits if we have exact distribution of  $\mathbf{X}$
- approximate limits if  $\text{pr}_\theta\{L(\mathbf{X}) \leq \theta \leq U(\mathbf{X})\} \approx 1 - \alpha$

- Example:  $X_1, \dots, X_n$  i.i.d.  $N(\mu, 1)$

- Example:  $X_1, \dots, X_n$  i.i.d.  $N(\mu, 1)$
- Example  $X_1, \dots, X_n$  i.i.d.  $U(0, \theta)$

- Example:  $X_1, \dots, X_n$  i.i.d.  $N(\mu, 1)$
- Example  $X_1, \dots, X_n$  i.i.d.  $U(0, \theta)$
- Example  $X_1, \dots, X_n$  i.i.d.  $N(\mu, \sigma^2)$

- Example:  $X \sim \text{Binom}(n, \theta)$ ,  $\hat{\theta} \sim N(\theta, \theta(1-\theta)/n)$

$$\text{pr}_{\theta} \left[ -1.96 \leq \frac{\sqrt{n}(\hat{\theta} - \theta)}{\{\theta(1-\theta)\}^{1/2}} \leq 1.96 \right] \approx 0.95$$

- Example:  $X \sim \text{Binom}(n, \theta)$ ,  $\hat{\theta} \sim N(\theta, \theta(1-\theta)/n)$

$$\text{pr}_{\theta} \left[ -1.96 \leq \frac{\sqrt{n}(\hat{\theta} - \theta)}{\{\theta(1-\theta)\}^{1/2}} \leq 1.96 \right] \approx 0.95$$

$$\text{pr}_{\theta} \left[ -1.96 \leq \frac{\sqrt{n}(\hat{\theta} - \theta)}{\{\hat{\theta}(1-\hat{\theta})\}^{1/2}} \leq 1.96 \right] \approx 0.95$$

- Example:  $X \sim \text{Binom}(n, \theta)$ ,  $\hat{\theta} \sim N(\theta, \theta(1-\theta)/n)$

$$\text{pr}_{\theta} \left[ -1.96 \leq \frac{\sqrt{n}(\hat{\theta} - \theta)}{\{\theta(1-\theta)\}^{1/2}} \leq 1.96 \right] \approx 0.95$$

$$\text{pr}_{\theta} \left[ -1.96 \leq \frac{\sqrt{n}(\hat{\theta} - \theta)}{\{\hat{\theta}(1-\hat{\theta})\}^{1/2}} \leq 1.96 \right] \approx 0.95$$

- $\hat{\theta}_n$  maximum likelihood estimate  $\hat{\theta} \sim N[\theta, \{nI(\theta)\}^{-1}]$
- approximate 95% confidence interval AoS Thm 6.16

- $\text{pr}_{\theta}\{\boldsymbol{\theta} \in R(\mathbf{X})\} \geq 1 - \alpha,$   
for all  $\theta$ , with equality for some  $\theta$

- pivotal method:

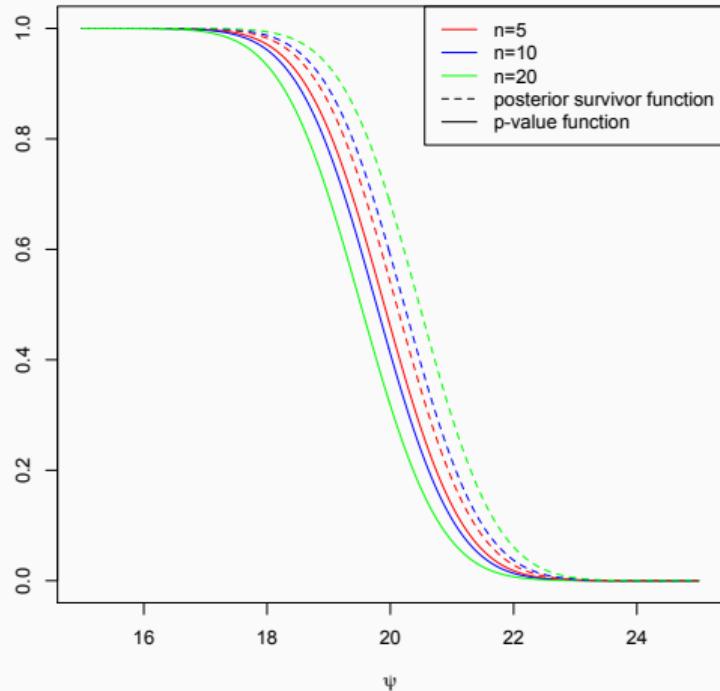
$$1 - \alpha = \text{pr}_{\theta}\{a \leq g(\mathbf{X}; \theta) \leq b\} = \text{pr}_{\theta}\{\boldsymbol{\theta} \in R(\mathbf{X})\}$$

- Example:  $\mathbf{X}_1, \dots, \mathbf{X}_n$  i.i.d.  $N_p(\boldsymbol{\mu}, \boldsymbol{\Sigma})$

MS Ex.7.8

# Bayesian credible intervals

MS 7.2



- upper and lower bounds

- upper and lower bounds
- highest posterior density

- upper and lower bounds
- highest posterior density
- equi-tailed posterior intervals

## Approximate normality of posterior

$$\bullet X_1, \dots, X_n \sim f(x^n \mid \theta), \quad \theta \sim \pi(\theta), \quad \pi(\theta \mid x^n) = \frac{f(x^n \mid \theta)}{f(x^n)} \quad x^n = (x_1, \dots, x_n)$$

## Approximate normality of posterior

- $X_1, \dots, X_n \sim f(x^n | \theta), \quad \theta \sim \pi(\theta), \quad \pi(\theta | x^n) = \frac{f(x^n | \theta)}{f(x^n)}$   $x^n = (x_1, \dots, x_n)$
- $\pi(\theta | x^n) \approx N\{\hat{\theta}, j^{-1}(\hat{\theta})\}; \quad \pi(\theta | x^n) \approx N\{\tilde{\theta}, \tilde{j}(\tilde{\theta})\}$

## Approximate normality of posterior

- $X_1, \dots, X_n \sim f(x^n | \theta), \quad \theta \sim \pi(\theta), \quad \pi(\theta | x^n) = \frac{f(x^n | \theta)}{f(x^n)} \quad x^n = (x_1, \dots, x_n)$
- $\pi(\theta | x^n) \approx N\{\hat{\theta}, j^{-1}(\hat{\theta})\}; \quad \pi(\theta | x^n) \approx N\{\tilde{\theta}, \tilde{j}(\tilde{\theta})\}$
- careful statement Berger, 1985; Ch.4
- For any  $a, b \in \mathbb{R}, a < b$
- let  $a_n = \hat{\theta}_n + aj^{-1/2}(\hat{\theta}_n), b_n = \hat{\theta}_n + bj^{-1/2}(\hat{\theta}_n)$
- $\hat{\theta}_n$  is the solution of  $\ell'(\theta; x^n) = 0$ , assumed unique, and  $j(\theta) = -\ell''(\theta; x^n)$

Then

$$\int_{a_n}^{b_n} \pi(\theta | x^n) d\theta \longrightarrow \Phi(b) - \Phi(a), \quad n \rightarrow \infty.$$

need  $\pi(\theta) > 0, \pi'(\theta)$  continuous

## Approximate normality of posterior

- $X_1, \dots, X_n \sim f(x^n | \theta), \quad \theta \sim \pi(\theta), \quad \pi(\theta | x^n) = \frac{f(x^n | \theta)}{f(x^n)}$   $x^n = (x_1, \dots, x_n)$
- $\pi(\theta | x^n) \approx N\{\hat{\theta}, j^{-1}(\hat{\theta})\}; \quad \pi(\theta | x^n) \approx N\{\tilde{\theta}, \tilde{j}(\tilde{\theta})\}$
- approximate posterior probability intervals
- exact posterior probability intervals  $\tilde{\theta} \approx \hat{\theta}$