Mathematical Statistics II

STA2212H S LEC9101

Week 5

February 4 2025



Today

- 1. Recap Jan 28 marginal posterior, hierarchical Bayes
- 2. Optimality in estimation: efficiency, CRLB MS §6.4
- 3. Optimality in estimation: decision theory MS §6.2
- 4. HW4
- 5. Office Hour today 3.30 4.30

Department Seminar Thursday February 6 11.00 – 12.00 Hydro Building, Room 9014 " Numerical integration in statistical problems " Alex Stringer, U Waterloo



- conjugate priors
- non-informative priors
- convenience priors
- minimally/weakly informative priors
- hierarchical priors

flat, "ignorance"

Marginal posterior distributions

• Bayes posterior carries all the information about θ , given **x**

by definition

• probabilities for any set A computed using the posterior distribution

•
$$\operatorname{pr}(\boldsymbol{\Theta} \in \mathsf{A} \mid \boldsymbol{x}) = \int_{\mathsf{A}} \pi(\boldsymbol{\theta} \mid \boldsymbol{x}) d\boldsymbol{\theta}$$

• if
$$\boldsymbol{\theta} = (\psi, \boldsymbol{\lambda})$$
, $\pi_{\mathsf{m}}(\psi \mid \boldsymbol{x}) = \int \pi(\psi, \boldsymbol{\lambda} \mid \boldsymbol{x}) d\boldsymbol{\lambda} = \frac{\int L(\psi, \boldsymbol{\lambda}; \boldsymbol{x}) \pi(\psi, \boldsymbol{\lambda}) d\boldsymbol{\lambda}}{\int L(\psi, \boldsymbol{\lambda}; \boldsymbol{x}) \pi(\psi, \boldsymbol{\lambda}) d\psi d\boldsymbol{\lambda}}$

• or, if
$$\psi = \psi(\theta)$$
, $\pi_{\mathsf{m}}(\psi \mid \mathbf{x}) = \int_{A} \pi(\theta \mid \mathbf{x}) d\theta$, $A = \{\theta \in \Theta : \psi(\theta) = \psi\}$

 with marginalization, 'flat' priors can have a large influence on the marginal posterior

Example: many normal means



Normal Circle, k=2, 5, 10

Mathematical Statistics II

- $x_i \mid \theta_i \sim N(\theta_i, v_i)$
- $\theta_i \mid \mu \sim N(\mu, \sigma^2)$
- $\mu \sim N(\mu_0, \tau^2)$
- $f(\mathbf{x} \mid \theta, \mu)$

v_i known

 $\sigma^{\rm 2}~{\rm known}$

hyperparameters

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

• $\pi(\theta, \mu \mid \mathbf{X})$

hyperparameters

 $E(\mu \mid \mathbf{x}) =$ $var(\mu \mid \mathbf{x}) =$ $E(\theta_i \mid \mathbf{x}) =$



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.

11 · Bayesian Models

622

F

.

.

$$E(\theta_i \mid \mathbf{x}) = \mathbf{x}_i \frac{\sigma^2}{\sigma^2 + \mathbf{v}_i} + E(\mu \mid \mathbf{x})(1 - \frac{\sigma^2}{\sigma^2 + \mathbf{v}_i})$$
$$E(\mu \mid \mathbf{x}) = \frac{\mu_0/\tau^2 + \sum \mathbf{x}_i/(\sigma^2 + \mathbf{v}_i)}{1/\tau^2 + \sum 1/(\sigma^2 + \mathbf{v}_i)}$$

- If σ^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

• recall, in regular models,

 $I(\theta)$ definition

$$\sqrt{n}(\hat{\theta} - \theta) \stackrel{d}{\rightarrow} N\{0, I^{-1}(\theta)\}$$

- smaller variance means more precise estimation
- Is $I^{-1}(\theta)$ small?

• recall, in regular models,

 $I(\theta)$ definition

$$\sqrt{n}(\hat{\theta} - \theta) \stackrel{d}{\rightarrow} N\{0, I^{-1}(\theta)\}$$

- smaller variance means more precise estimation
- Is $I^{-1}(\theta)$ small?
- Yes, there's a sense in which it is "as small as possible"

• recall, in regular models,

 $I(\theta)$ definition

$$\sqrt{n}(\hat{\theta} - \theta) \stackrel{d}{\rightarrow} N\{0, I^{-1}(\theta)\}$$

- smaller variance means more precise estimation
- Is $I^{-1}(\theta)$ small?
- Yes, there's a sense in which it is "as small as possible"
- Step 1: suppose $\mathbf{X} = X_1, \dots, X_n$ is an i.i.d. sample from a density $f(\mathbf{x}; \theta)$
- Let $U = U(\mathbf{X}) = \ell'(\theta; \mathbf{X})$

score function $E_{\theta}{S(\mathbf{X})} = q(\theta)$

proof: Cauchy-Schwarz

- Let S = S(X) be an unbiased estimator of $g(\theta)$
- then $\operatorname{var}_{\theta}(S) \geq {\operatorname{Cov}_{\theta}(S, U)}^2/\operatorname{Var}_{\theta}(U)$

Cramer-Rao lower bound

• Cauchy-Schwartz inequality: Z₁, Z₂, with $E(Z_1^2) < \infty$, $E(Z_2^2) < \infty$, MS Ex 2.7; HW2 STA2112F

 ${Cov(Z_1, Z_2)}^2 \leq var(Z_1)var(Z_2)$

- take $Z_1 = S(\mathbf{X})$, an unbiased estimator of $g(\theta)$
- take $Z_2 = U(\mathbf{X}) = \Sigma \ell'(\theta; X_i)$

score function

then

.

 ${Cov_{\theta}(S, U)}^2 \leq var_{\theta}(S)var_{\theta}(U)$

$$\mathsf{var}_{ heta}(\mathsf{S}) \geq rac{\mathsf{Cov}_{ heta}^2(\mathsf{S}, U)}{I_n(heta)}$$

$$\mathsf{var}_{ heta}(\mathsf{S}) \geq rac{\mathsf{Cov}_{ heta}^2(\mathsf{S}, U)}{I_n(heta)}$$

• $Cov_{\theta}(S, U)$

.

• when would we get equality?

$$\mathsf{var}_{ heta}(\mathsf{S}) \geq rac{\mathsf{Cov}_{ heta}^2(\mathsf{S}, U)}{I_n(heta)}$$

• $Cov_{\theta}(S, U)$

.

- when would we get equality?
- special case, $g(\theta) = \theta$

... Cramer-Rao lower bound

MS: $U_{\theta}(x)$ (p. 323)

- CRLB attained $\iff U(\theta; X) = A(\theta)S(X) + B(\theta)$
- MS Example 6.12: X_1, \ldots, X_n i.i.d. Poisson(λ)
- $\bar{X} = \hat{\lambda}$ has variance = $1/I(\theta)$
- + For estimating $\lambda^{\rm 2}$ three estimators are proposed

an unbiased estimator
$$T_1 = \frac{1}{n} \sum X_i(X_i - 1)$$
 $\frac{4\lambda^3}{n} + \frac{2\lambda^2}{n}$ the best unbiased estimator $T_2 = E(T_1 | \bar{X})$ $\frac{4\lambda^3}{n} + \frac{2\lambda^2}{n^2}$ the MLE $T_3 = \hat{\lambda}^2$ $\frac{4\lambda^3}{n} + \frac{5\lambda^2}{n^2} + \frac{\lambda}{n^3}$

Mathematical Statistics II February 4 2025

in finite samples

... Cramer-Rao lower bound

• A more interesting example, logistic density

$$f(x;\theta) = \frac{\exp(x-\theta)}{\{1+\exp(x-\theta)\}^2}$$

- CRLB of an unbiased estimator of θ is 3/n
- by previous argument, not attained

February 4 2025

• e.g. \bar{X} is unbiased for θ ,

Mathematical Statistics II

$$\operatorname{var}(\bar{X}) = \frac{\pi^2}{3n} > \frac{3}{n}$$

What about maximum likelihood estimator?

• Suppose $\tilde{\theta}_n$ is a sequence of estimators with

$$\sqrt{n}(\tilde{\theta}_n - \theta) \stackrel{d}{\rightarrow} N\{\mathsf{O}, \sigma^2(\theta)\}$$

- Is $\sigma^2(\theta) \ge 1/I(\theta)$?
- Yes, if $\tilde{\theta}_n$ is "regular", and $\sigma^2(\theta)$ continuous in θ

see MS §6.4, and Thm. 6.6

What about maximum likelihood estimator?

• Suppose $\tilde{\theta}_n$ is a sequence of estimators with

$$\sqrt{n}(\tilde{ heta}_n - heta) \stackrel{d}{
ightarrow} N\{\mathbf{0}, \sigma^2(heta)\}$$

- Is $\sigma^2(\theta) \ge 1/I(\theta)$?
- Yes, if $\tilde{\theta}_n$ is "regular", and $\sigma^2(\theta)$ continuous in θ

see MS §6.4, and Thm. 6.6

- Is the MLE 'regular'?
- Yes, under the 'usual regularity conditions'
- And, its a.var = lower bound

"BAN"

What about maximum likelihood estimator?

• Suppose $\tilde{\theta}_n$ is a sequence of estimators with

$$\sqrt{n}(\tilde{\theta}_n - \theta) \stackrel{d}{\rightarrow} N\{\mathsf{O}, \sigma^2(\theta)\}$$

- Is $\sigma^2(\theta) \ge 1/I(\theta)$?
- Yes, if $\tilde{\theta}_n$ is "regular", and $\sigma^2(\theta)$ continuous in θ

see MS §6.4, and Thm. 6.6

- Is the MLE 'regular'?
- Yes, under the 'usual regularity conditions'
- And, its a.var = lower bound
- there are other regular estimators that are also asymptotically fully efficient
- and might be better in finite samples

"BAN"

Asymptotic efficiency

· comparison of two consistent estimators

via limiting distributions

- $\sqrt{n}(T_{1n}-\theta) \xrightarrow{d} N\{\mathbf{0}, \sigma_1^2(\theta)\}, \quad \sqrt{n}(T_{2n}-\theta) \xrightarrow{d} N\{\mathbf{0}, \sigma_2^2(\theta)\}$
- asymptotic relative efficiency of T_1 , relative to T_2 is $\frac{\sigma_2^2(\theta)}{\sigma_1^2(\theta)}$

Asymptotic efficiency

· comparison of two consistent estimators

via limiting distributions

- $\sqrt{n}(T_{1n}-\theta) \xrightarrow{d} N\{O, \sigma_1^2(\theta)\}, \quad \sqrt{n}(T_{2n}-\theta) \xrightarrow{d} N\{O, \sigma_2^2(\theta)\}$
- asymptotic relative efficiency of T_1 , relative to T_2 is $\frac{\sigma_2^2(\theta)}{\sigma_1^2(\theta)}$
- if T_{2n} is the MLE $\hat{\theta}_n$, then $\sigma_2^2(\theta) = I^{-1}(\theta)$

as small as possible

- the MLE is fully efficient
- the asymptotic efficiency of T_1 is $1/\sigma_1^2(\theta)I(\theta)$

relative to the MLE implicit

Decision theory and Bayes estimators

- finite-sample approach to optimality in estimation
- start with a loss function $L(\hat{\theta}, \theta)$
- examples: squared error, absolute error, O-1 loss, Kullback-Liebler

Decision theory and Bayes estimators

- finite-sample approach to optimality in estimation
- start with a loss function $L(\hat{\theta}, \theta)$
- examples: squared error, absolute error, O-1 loss, Kullback-Liebler
- Risk function of $\hat{\theta}$ is expected loss:

$$\mathsf{R}_{ heta}(\hat{ heta}) = \mathrm{E}_{ heta}\{\mathsf{L}(\hat{ heta}, heta)\}$$

MSE, MAE, bias/variance trade-off

Decision theory and Bayes estimators

- finite-sample approach to optimality in estimation
- start with a loss function $L(\hat{\theta}, \theta)$
- examples: squared error, absolute error, O-1 loss, Kullback-Liebler
- Risk function of $\hat{\theta}$ is expected loss:

$$R_{\theta}(\hat{\theta}) = E_{\theta}\{L(\hat{\theta}, \theta)\}$$

MSE, MAE, bias/variance trade-off

- Risk function depends on θ , and on the form of the estimator

Examples: squared error loss



FIGURE 12.1. Comparing two risk functions. Neither risk function dominates the other at all values of $\theta.$

 $X \sim N(\theta, 1)$

Examples: squared error loss



$X \sim Binom(n, \theta)$

 $\alpha = \beta = \sqrt{n/4}$

- an estimator is admissible if no other estimator has a smaller risk function
- For a given loss function L, an estimator $\hat{\theta}$ is inadmissible if there is another estimator $\tilde{\theta}$ with

 $R_{ heta}(ilde{ heta}) \leq R_{ heta}(\hat{ heta}), \quad ext{for all } heta \in \Theta,$

and

$$R_{ heta_{
m o}}(ilde{ heta}) < R_{ heta_{
m o}}(\hat{ heta}), \quad ext{for some } heta_{
m o} \in \Theta.$$

- an estimator is admissible if no other estimator has a smaller risk function
- For a given loss function L, an estimator $\hat{\theta}$ is inadmissible if there is another estimator $\tilde{\theta}$ with

 $R_{ heta}(ilde{ heta}) \leq R_{ heta}(\hat{ heta}), \quad ext{for all } heta \in \Theta,$

and

$${\it R}_{ heta_{
m o}}(ilde{ heta}) < {\it R}_{ heta_{
m o}}(\hat{ heta}), \quad {
m for \ some \ } heta_{
m o} \in \Theta.$$

• MS Ex 6.1; $X \sim \lambda \exp(-\lambda x)$: under squared-error loss, $\hat{\lambda}$ is inadmissible: Beat by $\tilde{\lambda} = (n-1)\hat{\lambda}/n$ But under a different loss function the MLE has smaller risk than $\tilde{\lambda}$

 $L(\hat{\theta}, \theta) = \log(\frac{\theta}{\hat{\theta}}) - 1 - \frac{\theta}{\hat{\theta}}$

Optimal Bayes estimators

.

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

$${\sf R}_{\sf B}(\hat{ heta}) = \int {\sf R}_{ heta}(\hat{ heta}) \pi(heta) {\sf d} heta$$

• Optimal Bayes estimators minimize the expected posterior loss:

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

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

$$\int (\hat{ heta} - heta)^2 \pi(heta \mid \mathbf{x}) d heta$$

• solution $\hat{\theta}(\mathbf{x}) = \mathrm{E}(\theta \mid \mathbf{x})$

• Suppose $\hat{\theta}$ is a Bayes estimator

Bayes estimators are admissible

• 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

$$R_B(ilde{ heta}) = \int$$

and is unique

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

 ${\sf R}_ heta(ilde{ heta}, heta) \leq {\sf R}_ heta(\hat{ heta}, heta)$

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

$$R_B(ilde{ heta}) = \int$$

• 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

Mathematical Statistics II February 4 2025

Decision theory

- finding the 'best' point estimator $\hat{\theta}$
- best = smallest expected loss
- no asymptotic theory involved
- can find these using a Bayesian argument
- but the justification is not Bayesian
- another non-asymptotic approach to 'best' estimators: UMVU

MS 6.3

- 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 θ by producing a probability distribution for θ . Inferences, such as point estimates and interval estimates, may then be extracted from this distribution.

QUESTION 1: Interpreting probability



P(Heads) = 0.5 means...

- a. If I flip this coin over and over, roughly 50% will be Heads.
- b. Heads and Tails are equally plausible.
- c. Both a and b make sense.

QUESTION 2: Interpreting probability (again)



P(candidate A wins) = 0.8 means...

- a. If we observe this election over & over, candidate A will win roughly 80% of the time.
- b. Candidate A is 4 times more likely to win than to lose.
- c. The pollster's calculation is wrong.
 Candidate A will either win or lose, thus their probability of winning can only be 1 or 0.

QUESTION 3: Bigger picture



I claim that I can predict the outcome of a coin flip.

Mine claims she can distinguish between non-vegan and vegan poutine. We both succeed in 10 of 10 trials! What do you conclude?



- a. My claim is ridiculous. You're still more confident in Mine's claim than in my claim.
- b. 10-out-of-10 is 10-out-of-10 no matter the context. Thus the evidence supporting my claim is just as strong as the evidence supporting Mine's claim.

QUESTION 4: Asking questions



You've tested positive for a very rare genetic trait. If you only get to ask the doctor **one** question, which would it be?

- a. P(rare trait | +)
 Given the positive test result, what's the probability I actually have the trait?
- b. P(+ | rare trait)

If I *don't* have the trait, what's the chance I would have tested positive anyway?

Some history

LII. An Effay towards folving a Problem in the Doctrine of Chances. By the late Rev. Mr. Bayes, F. R. S. communicated by Mr. Price, in a Letter to John Canton, A. M. F. R. S.

Dear Sir,

Read Dec. 25, I Now fend you an effay which I have 1763. I found among the papers of our deceafed friend Mr. Bayes, and which, in my opinion, has great merit, and well deferves to be preferved. Experimental philosophy, you will find, is nearly interefted in the fubject of it; and on this account there feems to be particular reason for thinking that a communication of it to the Royal Society cannot be improper.



... Some history

528 Dr Fisher, Inverse probability

Inverse Probability. By R. A. FISHER, Sc.D., F.R.S., Gonville and Caius College; Statistical Dept., Rothamsted Experimental Station.

[Received 23 July, read 28 July 1930.]

I know only one case in mathematics of a doctrine which has been accepted and developed by the most eminent men of their time, and is now perhaps accepted by men now living, which at the same time hrs appeared to a succession of sound writers to be fundamentally false and devoid of foundation. Yet that is quite exactly the position in respect of inverse probability. Bayes, who seems to have first attempted to apply the notion of probability, not only to effects in relation to their causes but also to causes in relation to their effects, invented a theory, and evidently doubted its soundness, for he did not publish it during his life. It was posthumously published by Price, who seems to have felt no doubt of its soundness. It and its applications must have made great headway during the next 20 years, for Laplace takes for granted in a highly generalised form what Bayes tentatively wished to postulate in a special case.

Before going over the formal mathematical relationships in



This just in

Journal of the Royal Statistical Society Series A: Statistics in Society, 2025, **188**, 181–187 https://doi.org/10.1093/jrsssa/qnae044 Advance access publication 17 May 2024 **Original Article**



Bayesian issues in the 1950s: an episode involving Karl Popper and Jimmie Savage

Stephen M. Stigler

Department of Statistics, University of Chicago, Chicago, IL, USA

Address for correspondence: Stephen M. Stigler, Department of Statistics, University of Chicago, 5747 S. Ellis Ave., Chicago, IL 60637, USA. Email: stigler@uchicago.edu

... This just in

^cThe trifle in my hand that I wanted to mention is that you may remember, that in the dittoed draft of my book [Savage, 1954], one of the earliest ditto drafts, I attributed to R. A. Fisher the expression of the idea that since the a priori distribution washes out in a large sample, that there ought to be some intrinsic way of analyzing the data in itself without ever postulating a prior distribution at all. I don't remember whether I criticized that argument on the spot, but it's not valid, of course, because the prior distribution does wash out, does so only exponentially, and the rate at which it washes out does depend considerably on what prior distribution it is. Thus for example, since I'm firmly convinced that extrasensory perception does not exist, it would take tremendous amounts of data, of relevant opposing data, to bring me to the opposite point of view. Well, the thing was, we couldn't find this passage anywhere in Fisher and, when I wrote him, he said i twas ridiculous, he newer could've said any such thing, but Bob Schlaifer has found the reference for me, and it's in Paper 24 of Fisher's collected papers, it's the passage that stradles pages 286 and 287 and I just thought you might like to look at it for yourself.²

Here is the relevant paragraph from Fisher (1934):

'As an axiom this supposition [a uniform prior distribution] of Bayes fails, since the truth of an axiom should be manifest to all who clearly apprehend its meaning, and to many writers, including, it would seem, Bayes himself, the truth of the supposed axiom has not been apparent. It has, however, been frequently pointed out that, even if our assumed form for f(x)dx be somewhat inaccurate, our conclusions, if based on a considerable sample of observations, will not greatly be affected; and, indeed, subject to certain restrictions as to the true form of f(x)dx, it may be shown that our errors from this cause will tend to zero as the sample of observations is increased indefinitely. The conclusions drawn will depend more and more entirely on the facts observed, and less and less upon the supposed knowledge a priori introduced into the argument. This property of increasingly large samples has been sometimes put forward as a reason for accepting the postulate of knowledge a priori. It appears, however, more natural to infer from it that it should be possible to draw valid conclusions from the data alone, and without a priori assumptions.—If the justification for any particular form of f(x) is merely that it makes no difference whether the form is right or wrong, we may well ask what the expression is doing in our reasoning at all, and whether. if it were altogether omitted, we could not without its aid draw whatever inferences may, with validity, be inferred from the data. In particular we may question whether the whole difficulty

Mathematical Statistics II

Febru