# **Mathematical Statistics II**

STA2212H S LECO101

Week 3

January 24 2023



#### **Drinking less is better**

#### We now know that even a small amount of alcohol can be damaging to health.

Science is evolving, and the recommendations about alcohol use need to change.

Research shows that no amount or kind of alcohol is good for your health. It doesn't matter what kind of alcohol it is—wine, beer, cider or spirits.

Drinking alcohol, even a small amount, is damaging to everyone, regardless of age, sex, gender, ethnicity, tolerance for alcohol or lifestyle. That's why if you drink, it's better to drink less.

#### Alcohol consumption per week

Drinking alcohol has negative consequences. The more alcohol you drink per week, the more the consequences add up.



#### Aim to drink less

Drinking less benefits you and others. It reduces your risk of injury and violence, and many health problems that can shorten life.

#### Here is a good way to do it

Count how many drinks you have in a week.



Set a weekly drinking target. If you're going to drink, make sure you don't exceed 2 drinks on any day.

#### Good to know

You can reduce your drinking in steps! Every drink counts: any reduction in alcohol use has benefits.

#### It's time to pick a new target

What will your weekly drinking target be?



#### Tips to help you stay on target

Stick to the limits you've set for yourself.
 Drink loss of water.
 For every drink of alcohol, have one non-alcoholic drink.
 Choose alcohol-free or low-alcohol beverages.
 Eat before and while you're drinking.
 Have alcohol-free weeks or do alcohol-free activities.



#### ink

#### Today

- 1. Recap
- 2. Bayesian Inference
- 3. Optimality in Estimation MS 6
- 4. H3: comments on HW1, 2; Examples ...

Upcoming seminars of interest

- January 30 3.30 4.30 Chiara Sabatti Details "Human populations and gene mapping"
- January 30 6.00 7.00 pm Vera Liao Details "Introduction to Explainable AI"

Mathematical Statistics II

Jan. 30 OH 5-6 Monday tical Statistics II January 24 2023 (2000 / office)





#### Recap



#### **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
mechanics	7	44	49	59	34	46	0	32	49	52	44
vectors	51	69	41	70	42	40	40	45	57	64	61
	12	13	14	15	16	17	18	19	20	21	22
mechanics	36	42	5	22	18	41	48	31	42	46	63
vectors	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  $\hat{\theta} = 0.498$ 

$$\hat{\theta} = \sum_{i=1}^{22} (m_i - \bar{m})(v_i - \bar{v}) \left/ \left[ \sum_{i=1}^{22} (m_i - \bar{m})^2 \sum_{i=1}^{22} (v_i - \bar{v})^2 \right]^{1/2} \right.$$
(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 according to the true correlation coefficient  $\theta$ , "true" meaning the dorrelation for the hypothetea population of the students, of which we observed only 2

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



Â

579

11.2 · Inference

**Table 11.2** Mortalityrates r/m from cardiacsurgery in 12 hospitals(Spiegelhalter *et al.*,1996b, p. 15). Shown arethe numbers of deaths rout of m operations.

0/47B 18/148 8/119 D 46/810 8/211 13/196 С Ε Η 31/215 14/207 8/97 Κ 29/256 24/360 9/148 J L

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

Ô<sub>A</sub> = O

prior for hospital A Beta(1, 1)

Mathematical Statistics II January 24 2023

#### **Example: Binomial**



## **Choosing priors**

- conjugate priors 🧹
- non-informative priors  $\checkmark$
- convenience priors
- minimally/weakly informative priors
- hierarchical priors 🖌

flat, "ignorance"

# **Exponential families and conjugate priors** MS p.288,9 $\pi(\theta; \alpha, \beta) = K(\alpha, \beta) \exp\{\alpha c(\theta) - \beta d(\theta)\} \quad \int \pi(\theta) d\theta = 1$ $f(x;\theta) = \exp\{c(\theta)T(x) - d(\theta) + S(x)\};$ xeR ac (0) - Ball) $c(\Theta)T(x) - d(\Theta) \neq S(x)$ . OER $\pi(\varphi|z; \alpha, \beta) =$ ·Rab dO $c(0) \{ T(x) \neq x \} - (\beta \neq 1) d(0) + S(x) \}$ $T(\Theta(x; \alpha, \beta) = \mathcal{K}(\chi + T(x), \beta + 1) e^{(\alpha + T(x)/c(\Theta) - (\beta + 1)d(\Theta)}$ Statistics II January 24 2023 d O

Mathematical Statistics II

## Exponential families and conjugate priors

MS p.288,9

$$f(x;\theta) = \exp\{c(\theta)T(x) - d(\theta) + S(x)\}; \qquad \pi(\theta;\alpha,\beta) = K(\alpha,\beta) \exp\{\alpha c(\theta) - \beta d(\theta)\} \triangleq$$
  
Example:  $f(x;\theta) = \theta(1-\theta)^{x}, x = 0, 1, ...; 0 < \theta < 1$   

$$T(\theta_{1}^{*}\alpha,\beta) = k_{\alpha}e^{\alpha \log ((-\theta) + \beta \log \theta)}$$
  

$$T(x) = 2$$
  

$$C(\theta)_{z} = \log(1-\theta)$$
  

$$d(\theta) = -\log(\theta)$$
  

$$T(x) = 2$$
  

$$C(\theta)_{z} = \log(1-\theta)$$
  

$$d(\theta) = -\log(\theta)$$
  

$$T(x) = 2$$
  

$$C(\theta)_{z} = \log(1-\theta)$$
  

$$d(\theta) = -\log(\theta)$$
  

$$T(x) = 2$$
  

$$C(\theta)_{z} = \log(1-\theta)$$
  

$$d(\theta) = -\log(\theta)$$
  

$$T(x) = 2$$
  

$$T(x) = 2$$
  

$$C(\theta)_{z} = \log(1-\theta)$$
  

$$T(x) = 2$$
  

$$T(x) = 2$$

### Exponential families and conjugate priors

5

Conjugate proor

7

**Mathematical Statistics** 

$$f(x;\theta) = \exp\{c(\theta)T(x) - d(\theta) + S(x)\}; \quad \pi(\theta;\alpha,\beta) = K(\alpha,\beta) \exp\{\alpha c(\theta) - \beta d(\theta)\}$$
Example:  $f(x;\theta) = \theta(1-\theta)^{x}, x = 0, 1, ...; 0 < \theta < 1$ 
Example:  $f(x;\mu) = \frac{1}{\sqrt{2\pi}} \exp\{-\frac{1}{2}(x-\mu)^{2}\}$   $N(\mu, 1)$ 

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

$$e^{-\frac{x^{2}}{2\sigma^{2}}} - \log \sigma + \log \sqrt{m^{2}}$$

$$e^{-\frac{x^{2}}{2\sigma^{2}}} - \log \sigma + \log \sqrt{m^{2}}$$

$$e^{-\frac{x^{2}}{2\sigma^{2}}} - \log \sigma + \log \sqrt{m^{2}}$$

$$f(\theta) = -\frac{1}{2\sigma^{2}} A(\theta) = t \log \sigma$$

MS p.288,9

• if parameter space is closed (interval), e.g.  $\Theta = [a, b]$ , then  $\pi(\theta) \sim U(a, b)$  represents 'indifference'  $A = \pi(\phi) = 1$  ,  $0 \le 0 \le 1$  last week • example: Beta (1,1) prior for Bernoulli probability  $T(\mu|z) = \int_{\nabla \pi} e^{-\frac{1}{2}(z-\mu)^2} - \infty < \mu < t\infty$ • example 5.34: X  $\sim$  N( $\mu$ , 1),  $\pi(\mu) \propto$  1 Xin, Xn ind N (Mi)  $= \frac{1}{\sqrt{2\pi}} e^{-\frac{1}{2}(\mu - \pi)^2} \left(\int_{t_0} t_0\right)$  $\pi(\mu(\varkappa) \sim \mathcal{N}(\overline{x}, \frac{1}{n})$  $E(\mu|x) = x var(\mu|x) = 1$ (ntk) Mathematical Statistics II January 24 2023

MS p.290

X X

## Flat priors

- if parameter space is closed (interval), e.g.  $\Theta = [a, b]$ , then  $\pi(\theta) \sim U(a, b)$  represents 'indifference'
- example: Beta (1,1) prior for Bernoulli probability
- example 5.34:  $X \sim N(\mu, 1), \pi(\mu) \propto 1$
- improper priors can lead to proper posteriors
- L.g.  $N(0,\sigma^2) \pi(\sigma^2)$  ntbc  $\ll 1/\sigma^2$ • priors flat in one parameterization are not flat in another



last week

... Flat priors

• Example:  $X \sim Bin(n, \theta), 0 < \theta < 1; \theta \sim U(0, 1)^{\text{Re}}$ 

A(0) & SO(1-013-12

(ι-θ)

20

o.

0.10

0.00

prior for psi

 $\mathbf{C}$ 

exp.d. Finfo.

 $T_{4}^{1}(4)$ 

 $T(\varphi) =$ 

TT (4(0)

3.0

2.0

1.0

0.0

0.0

0.4

0.8

prior for theta

 $\int \mathcal{L}(\Theta) \rightarrow \mathcal{L}(\Psi)$   $= \left(-\frac{2^{2}}{2} \mathcal{L}(\Theta(\Psi))\right) = \cdots = \left(-\left(\mathcal{L}''(\Theta(\Psi))\right)\right) = \frac{1}{2} \left(-\frac{2^{2}}{2} \mathcal{L}(\Theta(\Psi))\right)$ 

• log-odds ratio  $\psi \neq \psi(\theta) = \log\{\theta/(1-\theta)\}$ 

•  $\pi(\psi) = \frac{e^{\psi}}{(1+e^{\psi})^2}, -\infty < \psi < \infty$ 

- prior probability  $-3 < \psi < 3 pprox$  0.9

• an invariant prior:  $\pi(\theta) \propto I^{1/2}(\theta)$ 

transf. 4= ψ(θ) → Mathematical Statistics II January 24 2023

5.31

I''-(0)

- $\pi(\theta) \propto l^{1/2}(\theta)$   $\mathcal{D}^{\prime 2}(\Theta)$   $\mathcal{D}^{\prime 2}(\Theta)$
- Example:  $X \sim Bin(n, \theta)$   $I(\theta) = n/\{\theta(1-\theta)\}, \quad 0 < \theta < 1$
- Example 5.35:  $X \sim Poisson(\lambda)$ ,  $I(\lambda) = 1/\lambda$ ,  $\lambda > 0$  (



• Jeffreys' prior for multiparameter  $\theta$ :  $\pi(\theta) \propto |I(\theta)|^{1/2}$  not recommended even by Jeffreys invariant

04

• Example:  $X_1, \ldots, X_n$  i.i.d.  $N(\mu, \sigma^2)$   $I(\mu, \sigma^2) = \frac{\chi e^{-\chi}}{\chi}, \frac{\pi e^{-\chi}}{\chi}, \frac{\pi$ 

## Marginalization

• Bayes posterior carries all the information about  $\theta$ , given **x** 

by definition

52 knom

- probabilities for any set A computed using the posterior distribution
- $\operatorname{pr}(\boldsymbol{\Theta} \in \boldsymbol{A} \mid \boldsymbol{x}) =$
- if  ${oldsymbol{ heta}}=(\psi,{oldsymbol{\lambda}})$ , ...
- or, if  $\psi = \psi(\theta)$
- in this context, 'flat' priors can have a large influence on the marginal posterior

## Not all likelihood functions are regular

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

MS Exercise 5.1

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

## **Optimality of estimators**

 $I(\theta)$ ?

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

бр :

 $l(\hat{\Theta}, X) = 0$ 

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

Xum, Xn id  $f(x, \theta)$ I E'Finfo.in A obs 2  $I(0) = E_0(3hpf(X_i, 0))$ 

MS Ch 6; AoS Ch 12

 $I(\theta)$ ?

• recall, in regular models,

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

- 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,

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

u(x)

- 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  $X = X_1, \ldots, X_n$  is an i.i.d. sample from a density  $f(x; \theta)$
- and suppose that  $E_{\theta}{S(X)} = g(\theta)$
- then  $var(S) \geq {Cov_{\theta}(S, U)}^2/Var_{\theta}(U)$

S(X) : subvared for g(O) in procepte, proof: Cauchy-Schwarz

 $I(\theta)$ ?

• recall, in regular models,

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

- 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  $\mathbf{X} = X_1, \ldots, X_n$  is an i.i.d. sample from a density  $f(\mathbf{x}; \theta)$
- and suppose that  $E_{\theta}\{S(X)\} = g(\theta)$
- then  $var(S) \ge {Cov_{\theta}(S, U)}^2/Var_{\theta}(U)$

proof: Cauchy-Schwarz

 $I(\theta)$ ?

## ... Optimality of estimators

• Cauchy-Schwartz inequality: for random variables X, Y, with  $E(X^2) < \infty$ ,  $E(Y^2) < \infty$ ,

 ${Cov(X, Y)}^2 \le var(X)var(Y)$ 

- now suppose  $X_1, \ldots, X_n$  i.i.d. with density  $f(x; \theta)$
- and suppose  $S(\mathbf{X})$  is unbiased for  $g(\theta)$
- and recall  $U(\mathbf{X}) = \Sigma \ell'(\theta; X_i)$  score function
- then

$$\{\operatorname{Cov}_{\theta}(S,U)\}^{2} \leq \operatorname{var}_{\theta}(S)\operatorname{var}_{\theta}(U) \qquad \text{(red } \operatorname{var}_{\theta}S(\underline{X}) < \infty \\ \operatorname{var}_{\theta}(S) \gg \begin{array}{c} 2 \\ \int S(\underline{x}) \end{array} = \begin{array}{c} \int S(\underline{x}) \frac{\partial L(\theta, \underline{x})}{\partial \theta} f(\underline{n}; \theta) d\underline{x} \end{array} \right)^{2} \quad \text{Form} \\ F_{\theta}S(\underline{X}) = g(\theta) \\ n \operatorname{I}(\theta) \\ = \int \int S(\underline{x}) \frac{\partial}{\partial \theta} f(\underline{n}; \theta) d\underline{x} \end{array} \right)^{2} / n \operatorname{I}(\Theta) \qquad 16$$



MS Ch 6: AoS Ch 12

... Optimality of estimators

•

MS Ch 6; AoS Ch 12

smooth model any S(X)

17

(0) 3

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

• special case 
$$g(\theta) = \theta$$
  
 $Var_{\phi} S(X) \gg \{g'(\theta)\} [nI(\theta)]^{-1}$ 

80

$$g(\phi) = \phi$$
  
 $E_{\phi}S(\chi) = \phi$  var\_{\phi}S(\chi) >  $nI(\phi)$   
 $\Lambda$ 

GRAMER-RAO L.B.

## ... Optimality of estimators

٠



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

Example: Poisson

MS Ex.6.12

In text: see TILE of 
$$\lambda$$
 in a Poisron advance CELE  
Unbiased estimator of  $\lambda^{2}$ :  $S_{1}(\mathbf{X}) = (1/n)\Sigma X_{i}(X_{i} - 1)$   
Maximum likelihood estimator of  $\lambda^{2}$ :  $S_{2}(\mathbf{X}) = \{(1/n)\Sigma X_{i}\}^{2}$   
 $\operatorname{var}(S_{1}) = \begin{pmatrix} 4\lambda^{3} \\ n \end{pmatrix} + \frac{2\lambda^{2}}{n^{2}} ?(n+kc)$   
 $\operatorname{var}(S_{2}) = \begin{pmatrix} 4\lambda^{3} \\ n \end{pmatrix} + \frac{5\lambda^{2}}{n^{2}} + \frac{\lambda}{n^{3}}$   
Cramer-Rao lower bound:  $\{g'(\lambda)\}^{2}/nI(\lambda) = (2\lambda)^{2}/(n/\lambda) = (4\lambda^{3}/n)$   
Note: CRLB cannot be attained even by an unbiased estimator  
 $g'(\lambda) = \lambda^{2}$   
Mathematical Statistics | January 24 2023 | 1000 | 1000 | 1000 | 1000 | 1000 | 1000 | 1000 | 1000 | 1000 | 1000 | 1000 | 1000 | 1000 | 1000 | 1000 | 1000 | 1000 | 1000 | 1000 | 1000 | 1000 | 1000 | 1000 | 1000 | 1000 | 1000 | 1000 | 1000 | 1000 | 1000 | 1000 | 1000 | 1000 | 1000 | 1000 | 1000 | 1000 | 1000 | 1000 | 1000 | 1000 | 1000 | 1000 | 1000 | 1000 | 1000 | 1000 | 1000 | 1000 | 1000 | 1000 | 1000 | 1000 | 1000 | 1000 | 1000 | 1000 | 1000 | 1000 | 1000 | 1000 | 1000 | 1000 | 1000 | 1000 | 1000 | 1000 | 1000 | 1000 | 1000 | 1000 | 1000 | 1000 | 1000 | 1000 | 1000 | 1000 | 1000 | 1000 | 1000 | 1000 | 1000 | 1000 | 1000 | 1000 | 1000 | 1000 | 1000 | 1000 | 1000 | 1000 | 1000 | 1000 | 1000 | 1000 | 1000 | 1000 | 1000 | 1000 | 1000 | 1000 | 1000 | 1000 | 1000 | 1000 | 1000 | 1000 | 1000 | 1000 | 1000 | 1000 | 1000 | 1000 | 1000 | 1000 | 1000 | 1000 | 1000 | 1000 | 1000 | 1000 | 1000 | 1000 | 1000 | 1000 | 1000 | 1000 | 1000 | 1000 | 1000 | 1000 | 1000 | 1000 | 1000 | 1000 | 1000 | 1000 | 1000 | 1000 | 1000 | 1000 | 1000 | 1000 | 1000 | 1000 | 1000 | 1000 | 1000 | 1000 | 1000 | 1000 | 1000 | 1000 | 1000 | 1000 | 1000 | 1000 | 1000 | 1000 | 1000 | 1000 | 1000 | 1000 | 1000 | 1000 | 1000 | 1000 | 1000 | 1000 | 1000 | 1000 | 1000 | 1000 | 1000 | 1000 | 1000 | 1000 | 1000 | 1000 | 1000 | 1000 | 1000 | 1000 | 1000 | 1000 | 1000 | 1000 | 1000 | 1000 | 1000 | 1000 | 1000 | 1000 | 1000 | 1000 | 1000 | 1000 | 1000 | 1000 | 1000 | 1000 | 1000 | 1000 | 1000 | 1000 | 1000 | 1000 | 1000 | 1000 | 1000 | 1000 | 1000 | 1000 | 1000 | 1000 | 1000 | 1000 |

## What about maximum likelihood estimator?

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

$$\sqrt{n}(\tilde{\theta}_n - \theta) \xrightarrow{d} N\{\mathbf{0}, \overset{\mathbf{v}}{\underbrace{\sigma^2(\theta)}}\}$$

asy. Jession of CRLB

√n (ô - 0) → N(0, I-'101)

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

see MS §6.4, and Thm. 6.6

## What about maximum likelihood estimator?

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

$$\sqrt{n}(\widetilde{ heta}_n- heta)\stackrel{d}{
ightarrow} N\{\mathsf{O},\sigma^{\mathsf{2}}( heta)\}$$

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

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"

I(0)

#### What about maximum likelihood estimator?

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

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

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

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"

comparison of two consistent estimators

via limiting distributions

MS 4.8

• 
$$\sqrt{n}(T_{1n}-\theta) \xrightarrow{d} N\{O, \sigma_1^2(\theta)\}, \quad \sqrt{n}(T_{2n}-\theta) \xrightarrow{d} N\{O, \sigma_2^2(\theta)\}$$

• asymptotic relative efficiency of  $T_1$ , relative to  $T_2$  is  $\frac{\sigma_2^2(\theta)}{\sigma_1^2(\theta)}$ 

## Asymptotic efficiency

- comparison of two consistent estimators
- $\sqrt{n}(T_{1n} \theta) \xrightarrow{d} N\{O, \sigma_1^2(\theta)\}, \quad \sqrt{n}(T_{2n} \theta) \xrightarrow{d} N\{O, \sigma_2^2(\theta)\}$

• the asymptotic, relative to MLE efficiency of  $T_2$  is  $\sigma_2^2(\theta)I(\theta)$ 

• the MLE is fully efficient & as small a. var. as possible

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

a.vor  $(\hat{\psi}_1) = -g_1(\theta)$ a.vor  $(\hat{\psi}_2) = -g_1(\theta)$ 

ô ± 2 se d-m.

via limiting distributions

as small as possible

Mathematical Statistics II January 24 2023

20

### **Decision theory and Bayes estimators**

#### MS 6.2, AoS Ch 12

- finite-sample approach to optimality in estimation
- start with a loss function  $L(\hat{\theta}, \theta)$
- undel  $X \sim f(x, 0)$  data  $X_{1}, \dots, X_{n}$ • examples: squared error, absolute error, 0-1 loss, K-L divergence

$$L_{1}(\hat{\theta}, \theta) = (\hat{\theta} - \theta)^{2}$$

$$L_{2}(\hat{\theta}, \theta) = |\hat{\theta} - \theta|$$

$$L_{3}(\hat{\theta}, \theta) = \begin{cases} 0 & \hat{\theta} = \theta \\ 1 & \hat{\theta} \neq \theta \end{cases}$$

$$L_{4}(\hat{\theta}, \theta) = (see A_{0}S Ch(r))$$

### **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, 0-1 loss, K-L divergence
- Risk function of  $\hat{\theta}$  is expected loss:

$$R_{\theta}(\hat{\theta}) = E_{\theta}\{L(\hat{\theta},\theta)\} = \int L(\hat{\theta}(\underline{x}), \theta) f(\underline{x}; \theta) d\underline{x}$$

$$MSE \equiv Vor \hat{\theta} + (bias \hat{\theta})^{2} \qquad MSE, MAE, bias/variance trade-off$$

$$e.g. = \int \{\hat{\theta}(\underline{x}) - \theta \hat{f}(\underline{u}; \theta) d\underline{x} \qquad Nisk \quad f \quad \hat{\theta}$$

$$= \int \{\hat{\theta}(\underline{x}) - E_{\theta}\hat{\theta}(\underline{x}) + E_{\theta}\hat{\theta}(\underline{x}) - \theta \hat{f}(\underline{u}; \theta) d\underline{x} \qquad X$$

$$Mathematical Statistics II \qquad January 24,2023$$

$$= \int \{\hat{\theta}(\underline{x}) - E_{\theta}\hat{\theta}(\underline{x}) + E_{\theta}\hat{\theta}(\underline{x}) - \theta \hat{f}(\underline{u}; \theta) d\underline{x} \qquad X$$

 $\hat{\Theta} = \hat{\Theta}(\underline{X})$ 

# 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, K-L divergence
- Risk function of  $\hat{\theta}$  is expected loss:

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

I A BIAC

MSE, MAE, bias/variance trade-off

MS 6.2, AoS Ch 12

• Risk function depends on  $\theta$ , and on the form of the estimator

#### **Examples: squared error loss**

#### AoS 12.2, 12.3; MS Ex.6.1





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

#### **Examples: squared error loss**



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



• an estimator is admissible if no other estimator has a smaller risk function

## Optimality

- an estimator is admissible if no other estimator has a smaller risk function  $\forall \rho_{\mathcal{E}} \oplus$
- For a given loss function *L*, an estimator  $\hat{\theta}$  is inadmissible if there is another estimator  $\tilde{\theta}$  with  $R_{\theta}(\tilde{\theta}) \leq R_{\theta}(\hat{\theta}), \quad \text{for all } \theta \in \Theta,$

and

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

- 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

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

• 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

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

• 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  $L(\hat{\theta}, \theta) = (\hat{\theta} - \theta)^2$  need to minimize over  $\hat{\theta}$ 

$$\int (\hat{\theta} - \theta)^2 \pi(\theta \mid \mathbf{X}) d\theta$$

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

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

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

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

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

Mathematical Statistics II January 24 2023

and is unique

#### Bayes estimators are admissible

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

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

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

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

• 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 January 24 2023

and is unique

 $\mathcal{R}_{o}(\hat{\Theta}) = \mathcal{E}_{o}L(\hat{\Theta}, \sigma)\hat{z}$ 

 $R_{B}(\Theta) = \int R_{O}(\widehat{\Theta}) \pi(\Theta) d\Theta$ 

se e Aos **Decision theory** varg(x)? + - - -(hoose to • finding the 'best' point estimator  $\hat{\theta}$ mining over for best = smallest expected loss  $R_{B}(\Theta) = \int R_{O}(\widehat{\Theta}) \pi(\Theta) d\Theta$  no asymptotic theory involved =  $\int \left( L\left(\hat{\Theta}(\underline{x}), \Theta\right) \cdot f(\underline{u}, \Theta) d\underline{x} \pi \left[ \Theta \right] d\Theta \right)$  can find these using a Bayesian argument =  $\int (\hat{o}(\underline{x}); o) \pi (o|\underline{x}) d\underline{x} f(\underline{u}) d\theta$  but the justification is not Bayesian • another non-asymptotic approach to 'best' estimators: UMVU MS 6.3  $E_{\mu}L(\hat{\theta};\theta) = E\{L(\hat{\theta},\theta) \mid \chi\}$   $\pi(\theta \mid \chi) \qquad (under post.)$ Mathematical Statistics II January 24 2023

## Multi-parameter models

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

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

#### Multi-parameter models

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

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

• marginal posterior

 $\pi_m(\theta_1 \mid \mathbf{x}^n) = \int \pi(\theta \mid \mathbf{x}^n) d\theta_2 \dots d\theta_p$ 

marginal posterior

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

for  $\psi(\theta)$ 

#### **Bayesian inference: Multi-parameter models**

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

#### **Bayesian inference: Multi-parameter models**

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

• 
$$\psi = \sum_{i=1}^{n} \mu_i^2$$

squared length of mean vector

$$\pi(\psi \mid \mathbf{x}^n) = \int_{\mathsf{A}} \pi(\mu \mid \mathbf{x}^n) d\mu$$

•  $\mu_i \mid \mathbf{X}_i \sim N(\mathbf{X}_i, \mathbf{1}) \implies \sum \mu_i^2 \mid \mathbf{X}^n \sim \chi_n^2(\sum \mathbf{X}_i^2)$ 



ψ

- 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.

# QUESTION 1: Interpreting probability



P(Heads) = 0.5 means...

- F a. If I flip this coin over and over, roughly 50% will be Heads.
- B b. Heads and Tails are equally plausible.
- P 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?

- 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?

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

v<sub>i</sub> known

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

hyperparameters

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

hyperparameters

•  $\pi(\theta, \mu \mid \mathbf{X})$ 

#### **Bayesian hierarchical models**

 $egin{aligned} & \mathsf{E}(\mu \mid \mathsf{X}) = \ & \mathsf{var}(\mu \mid \mathsf{X}) = \ & \mathsf{E}( heta_i \mid \mathsf{X}) = \end{aligned}$ 



.

.

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

- If  $\sigma^2$  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