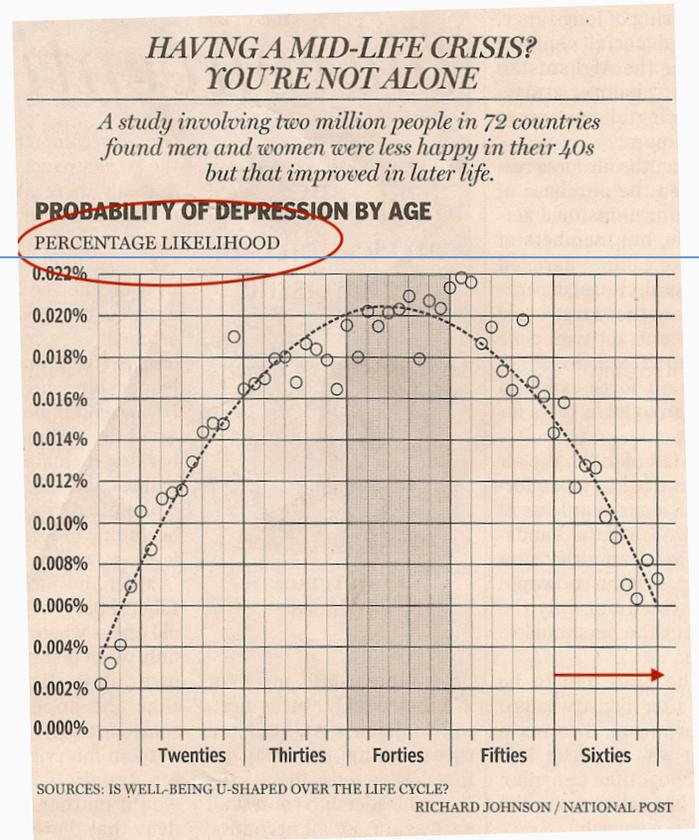


# Mathematical Statistics II

STA2212H S LEC9101

Week 1

January 10 2023



1. Course Overview
2. Review of Likelihood [MS Ch 4](#)
3. Upcoming seminars of interest

January 16 3.30 – 4.30 [Nicholas Horton](#) [Details](#)

“Teaching reproducibility and responsible workflows”



[Link](#)

## **STA 2212S: Mathematical Statistics II**

**Tuesday, 10.00-13.00 Eastern**

**OI 5150**

January 10 – April 4 2023

### **From the calendar:**

This course is a continuation of STA2112H. It is designed for graduate students in statistics and biostatistics. Topics include: Likelihood inference, Bayesian methods, Significance testing, Linear and generalized linear models, Goodness-of-fit, Computational methods

Prerequisite: STA2112H

I will definitely cover the first 3 topics, and the 5th, and we'll see how time goes for the others. “Computational methods” was probably meant to be shorthand for “bootstrap” and “MCMC”, and will be touched on in the other topics.

## STA 2212S: Mathematical Statistics II Syllabus

[Link](#)

Spring 2023

Week	Date	Methods	References
1	Jan 10	Likelihood inference: review of ML estimation; mis-specified models; computation; nonparametric mle	MS §§5.1–7, SM Ch 4
2	Jan 17	Bayesian estimation; Bayesian inference	MS §5.8; AoS §§ 11.1–4; SM §§11.1,2
3	Jan 24	Optimality in estimation	MS Ch 6; AoS Ch 12; SM §7.1, 11.5.2
4	Jan 31	Interval estimation; Confidence bands	MS §§7.1,2; AoS Ch 7; SM §7.1.4
5	Feb 7	Upper-tail test; likelihood ratio test	MS §§7.1,4; AoS Ch 10.6; SM

[Link](#)

**HW Question Week 1**

STA 2212S 2023

**Due January 16**

*MS, Exercise 5.2*

~~i.i.d.~~

Suppose  $(X_1, Y_1), \dots, (X_n, Y_n)$  are independent pairs of random variables where  $X_i$  and  $Y_i$  are i.i.d.  $N(\mu_i, \sigma^2)$  random variables.

- Find the maximum likelihood estimators of  $\mu_1, \dots, \mu_n$  and  $\sigma^2$ .
- Show that the maximum likelihood estimator of  $\sigma^2$  is not consistent. Does this contradict the theory we have established regarding the consistency of maximum likelihood estimators? Why or why not?
- Suppose we observe only  $Z_1, \dots, Z_n$ , where  $Z_i = X_i - Y_i$ . Find the maximum likelihood estimator of  $\sigma^2$  based on  $Z_1, \dots, Z_n$  and show that it is consistent.

[Link](#)

## STA2212: Inference and Likelihood

### A. Notation

**One random variable:** Given a model for  $X$  which assumes  $X$  has a density  $f(x; \theta)$ ,  $\theta \in \Theta \subset \mathbb{R}^k$ , we have the following definitions:

likelihood function

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

log-likelihood function

$$\ell(\theta; x) = \log L(\theta; x) = \log f(x; \theta) + a(x)$$

score function

$$u(\theta) = \partial \ell(\theta; x) / \partial \theta \quad \ell'(x; \theta)$$

observed information function

$$j(\theta) = -\partial^2 \ell(\theta; x) / \partial \theta \partial \theta^T \quad J(\theta) = E_{\theta}\{j(\theta)\}$$

expected information (in one observation)

$$i(\theta) = E_{\theta}\{U(\theta)U(\theta)^T\}^{-1} \quad I(\theta) \text{ (p.245)}$$

**Independent observations:** When we have  $X_i$  independent, identically distributed from  $f(x_i; \theta)$ , then, denoting the observed sample  $\mathbf{x} = (x_1, \dots, x_n)$  we have:

## Important Definitions

For  $X \sim f(x; \theta)$ , we define the following:

$$\theta \in \mathbb{R}$$

$$J(\theta) = E_{f(x; \theta)} [J(x; \theta)]$$

$n=1$

★ Score function:  $U(\theta) = \frac{d\ell(x; \theta)}{d\theta}$

★ Observed information function:  $J(\theta) = -\frac{d^2\ell(x; \theta)}{d\theta^2} = -\frac{dU(\theta)}{d\theta}$

★ Expected information function:

$$I(\theta) = E_{\theta}[U^2(X; \theta)]$$

**Note:** Sometimes we add a subscript 1 to these functions to indicate they are based on a single observation.

$$I_1(\theta) \quad J_1(\theta)$$

## Important Definitions

For  $X_1, \dots, X_n \sim^{iid} f(x; \theta)$ , we define the following:

★ Score function:  $\underline{U}_n(\theta) = \frac{d\ell_n(x_i; \theta)}{d\theta}$

★ Observed information function:  $\underline{J}_n(\theta) = -\frac{d^2\ell_n(x_i; \theta)}{d\theta^2} = -\frac{dU_n(\theta)}{d\theta}$

★ Observed (Fisher) information:  $\underline{J}_n(\hat{\theta}_n)$

★ Expected (Fisher) information function:

$$I_n(\theta) = E_{\theta}[U_n^2(X_i; \theta)] = \underline{nl_1(\theta)} \quad n i_1(\theta)$$

**Note:** Typically the subscript  $n$  is removed, but is used to emphasize that we are considering a random sample of size  $n$ .

# Vector parameters

- model  $X \sim f(x; \theta), \theta \in \mathbb{R}^p$        $\theta$  is a **column vector**

$L = f \in (0, 1)$   
 $\ln L < 0$



- $L(\theta; \underline{x}) : \mathbb{R}^p \rightarrow \mathbb{R} \propto f(x; \theta)$

map from  $\mathbb{R}^p \rightarrow \mathbb{R}$

- $l'(\theta; x) = \begin{bmatrix} \frac{\partial \ln L(\theta; x)}{\partial \theta_1} \\ \vdots \\ \frac{\partial \ln L(\theta; x)}{\partial \theta_p} \end{bmatrix}$

$p \times 1$  vector

- $-l''(\theta; x)$

$$= \begin{bmatrix} -\frac{\partial^2 \ln L(\theta; x)}{\partial \theta_1^2} & -\frac{\partial^2 \ln L(\theta; x)}{\partial \theta_1 \partial \theta_2} & \dots & -\frac{\partial^2 \ln L(\theta; x)}{\partial \theta_1 \partial \theta_p} \\ \vdots & \vdots & \ddots & \vdots \\ -\frac{\partial^2 \ln L(\theta; x)}{\partial \theta_2 \partial \theta_1} & -\frac{\partial^2 \ln L(\theta; x)}{\partial \theta_2^2} & \dots & -\frac{\partial^2 \ln L(\theta; x)}{\partial \theta_2 \partial \theta_p} \\ \vdots & \vdots & \ddots & \vdots \\ -\frac{\partial^2 \ln L(\theta; x)}{\partial \theta_p \partial \theta_1} & -\frac{\partial^2 \ln L(\theta; x)}{\partial \theta_p \partial \theta_2} & \dots & -\frac{\partial^2 \ln L(\theta; x)}{\partial \theta_p^2} \end{bmatrix}$$

$p \times p$  matrix

# Vector parameters

- model  $X \sim f(\mathbf{x}; \theta)$ ,  $\theta \in \mathbb{R}^p$        $\theta$  is a **column vector**
- $L(\theta; \mathbf{x})$       map from  $\mathbb{R}^p \rightarrow \mathbb{R}$
- $\ell'(\theta; \mathbf{x})$        $p \times 1$  vector
- $-\ell''(\theta; \mathbf{x})$        $p \times p$  matrix

# Properties of maximum likelihood estimators

$$\underline{x} \sim f(\underline{x}; \theta)$$

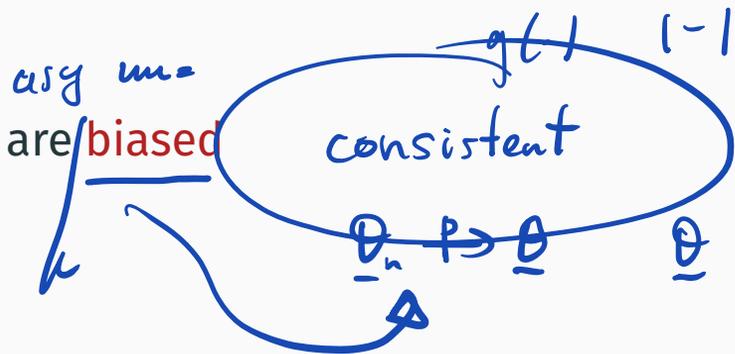
- maximum likelihood estimators are **equivariant**

$$l'(\hat{\theta}_n) = \underline{0} \quad \begin{matrix} \uparrow \\ p \times 1 \end{matrix}$$

$$\hat{\theta}_n \text{ mle of } \theta \Rightarrow g(\hat{\theta}_n) \text{ mle of } g(\theta)$$

- maximum likelihood estimators are **biased**

$$\left( \begin{array}{c} E\hat{\theta}_n \neq \theta \\ \downarrow \\ \theta \end{array} \right)$$



special exceptions

- maximum likelihood estimators have no explicit formula

in general

$$l'(\hat{\theta}_n) = 0 \quad \checkmark$$

# Asymptotic properties of maximum likelihood estimators

- maximum likelihood estimators are (i) consistent, (ii) **asymptotically normal**
- (ii) TS expansion

p.256

$$l'(\hat{\theta}_n; \underline{x}) = 0 = \underbrace{l'(\theta; \underline{x})}_{p \times 1} + \underbrace{l''(\theta; \underline{x})}_{p \times p} (\hat{\theta}_n - \theta) + \underbrace{R_n}_{p \times 1}$$

$$= \{l''(\theta; \underline{x})\}^{-1} l'(\theta; \underline{x}) + (\hat{\theta}_n - \theta) + R_n'$$

$$\hat{\theta}_n - \theta = \{-l''(\theta; \underline{x})\}^{-1} l'(\theta; \underline{x}) + R_n'$$

$$\hat{\theta}_n(\underline{x}) - \theta = \{-l''(\theta; \underline{x}_n)\}^{-1} l'(\theta; \underline{x}_n) + R_n'$$

↑  
sum of iid's

$$\sqrt{n} (\hat{\theta}_n(x) - \theta) = \left\{ -\frac{1}{n} l''(\theta; \underline{x}) \right\}^{-1} \frac{1}{\sqrt{n}} l'(\theta; X_n) + R_n$$

$\downarrow \omega$

$\downarrow \text{CLT}$

$$E_{\theta} l''(\cdot; X_i)$$

$$(0, \sigma^2)$$

$$I_1(\theta)$$

~~$\frac{d}{d\theta} \log h(\theta)$~~

$$\frac{d}{d\theta} \log h(\theta) = \frac{h'(\theta)}{h(\theta)}$$

$$E_{\theta} l'(\theta; \underline{X}) = 0 = \int l(\theta; \underline{x}) f(\underline{x}; \theta) d\underline{x}$$

$$\int f(\underline{x}; \theta) d\underline{x} = 1 \implies \frac{d}{d\theta} \int f(\underline{x}; \theta) d\underline{x} = 0 =$$

$p \times 1$

$$\begin{aligned} \sigma^2 &= \text{var}_{\theta} \{ \ell'(\theta; x_i) \} = E_{\theta} \{ \ell'(\theta; x_i) \ell'(\theta; x_i)^T \} \\ &= I(\theta), \quad I_i(\theta), \quad i_i(\theta) \\ &\quad \text{exp'd F. inf.} \end{aligned}$$

$$\begin{aligned} \sqrt{n}(\hat{\theta}_n - \theta) &\xrightarrow{d} N(0, J'(\theta) I(\theta) J^{-1}(\theta)) \\ \text{if } I &= J &\equiv N_p(0, I^{-1}(\theta)) \quad (*) \\ &\quad \uparrow \\ &\quad p \times p \end{aligned}$$

left to prove

$$\boxed{R_n \rightarrow 0} \quad (\text{assumptions})$$

$$\begin{aligned} I_i(\theta) &= E_{\theta} \left\{ \frac{\partial \ell}{\partial \theta} \left( \frac{\partial \ell}{\partial \theta} \right)^T \right\} & \mathbb{F}(\theta) &= E_{\theta} \left\{ - \frac{\partial^2 \ell(\theta)}{\partial \theta \partial \theta^T} \right\} \\ &= E_{\theta} \left\{ \left( \frac{\partial \ell(\theta; x_i)}{\partial \theta} \right) \right\}^T \end{aligned}$$

$$\begin{aligned} \frac{\partial}{\partial \theta^T} \int \frac{\partial \ell(\theta; x)}{\partial \theta} f(x; \theta) dx &= 0 \\ \text{p.k.t. vector} & & &= \int \frac{\partial^2 \ell(\theta; x)}{\partial \theta \partial \theta^T} f(x; \theta) + \int \frac{\partial \ell}{\partial \theta} \left( \frac{\partial f}{\partial \theta} \right)^T \end{aligned}$$

# Your friend the delta-method

MS Th.3.4 and p.148

Suppose

$$\frac{1}{\sqrt{n}} \text{ or } \sqrt{n} \dots$$

$$\downarrow$$

$$a_n(\mathbf{x} - \underline{\theta}) \xrightarrow{d} \mathbf{Z},$$

and  $g(\mathbf{x})$  is continuously differentiable at  $\theta$ , then

$$a_n\{g(\mathbf{x}) - g(\theta)\} \xrightarrow{d} D(\theta)\mathbf{Z}$$

where  $D(\theta) =$

$$\left[ \begin{array}{ccc} \frac{\partial g_1(\mathbf{x})}{\partial x_1} & \dots & \frac{\partial g_1(\mathbf{x})}{\partial x_p} \\ \vdots & & \vdots \\ \frac{\partial g_k(\mathbf{x})}{\partial x_1} & \dots & \frac{\partial g_k(\mathbf{x})}{\partial x_p} \end{array} \right]_{k \times p}$$

$$\theta \in \mathbb{R}^p \quad \mathbf{x} = (x_1, \dots, x_p)$$

$$g(\hat{\theta}_n) \quad g(\theta)$$

$$\{g_1(\mathbf{x}), \dots, g_k(\mathbf{x})\}$$

$$\hat{\mu}, \hat{\sigma}^2$$

$$\downarrow$$

$$\begin{pmatrix} \hat{\mu} \\ \hat{\sigma}^2 \end{pmatrix}$$

$$(\hat{\mu}^2 \text{ and } \hat{\sigma}^2)$$

$$\underline{\theta} = \underline{x}$$

$$\sqrt{n}(\hat{\theta}_n - \theta) \xrightarrow{d} N\{0, I_n^{-1}(\theta)\}$$

$$\underbrace{\sqrt{n}\{g(\hat{\theta}_n) - g(\theta)\}} \xrightarrow{d} N\{0, \underbrace{g'(\theta)^T I_n(\theta)^{-1} g'(\theta)}\}$$

z

$$I_n^{-1}(\theta)$$

$$\{I_n(\theta)\}^{-1}$$

See also AoS §9.9

$$\sqrt{n}(\hat{\theta}_n - \theta) \xrightarrow{d} N(0, \Sigma^{-1}(\theta))$$

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

$$f(x_i; \lambda, \alpha) = \frac{1}{\Gamma(\alpha)} \lambda^\alpha x_i^{\alpha-1} \exp(-\lambda x_i)$$

$$L(\lambda, \alpha; \underline{x}) = \prod_{i=1}^n \frac{1}{\Gamma(\alpha)} \quad \downarrow$$

$$= \frac{1}{\Gamma^n(\alpha)} \cdot \lambda^{n\alpha} \cdot \prod x_i^{(\alpha-1)} \cdot e^{-\lambda \sum x_i} \cdot \left( \prod x_i^{-1} \right)$$

$$\ell(\lambda, \alpha; \underline{x}) = \underline{-n \log \Gamma(\alpha)} + n\alpha \log \lambda + (\alpha-1) \sum \log(x_i)$$

$$\frac{\partial}{\partial \alpha}$$

$$-n\psi(\alpha) + n \log x \quad -\lambda \sum x_i$$

$$\hat{\lambda} = \frac{n\alpha}{\sum x_i}$$

$$\text{a. var } \hat{\lambda} = (n\alpha/\lambda^2)??$$

$$\left( \frac{\lambda^2}{n\alpha} \right)$$

duh!

↓

# ... Example

$$\begin{pmatrix} \partial l / \partial \alpha \\ \partial l / \partial \lambda \end{pmatrix} \Big|_{(\hat{\alpha}, \hat{\lambda})} = 0$$

$$\log \Gamma'(\alpha) / \Gamma(\alpha) \equiv \psi(\alpha)$$

find  $\text{a.var}(\hat{\mu})$  via mv delta method

$$\frac{n\hat{\alpha}}{\hat{\lambda}} - \sum x_i = 0 \quad \hat{\lambda} = \left( \frac{\sum x_i}{n\hat{\alpha}} \right)$$

$$-n\psi(\hat{\alpha}) + n \log \hat{\lambda} + \sum \log x_i = 0$$

$$-n\psi(\hat{\alpha}) - n \log \left( \frac{\sum x_i}{n\hat{\alpha}} \right) + \sum \log x_i = 0$$

$$-n\psi(\hat{\alpha}) - n \log(n\hat{\alpha}) + \sum \log x_i + n \log \sum x_i = 0$$

$$\psi(\hat{\alpha}) + \log(\hat{\alpha}) = \frac{1}{n} \sum \log x_i + \log \left( \frac{\sum x_i}{n} \right)$$

$$\underbrace{\psi(\hat{\alpha}) + \log(\hat{\alpha}) = -\frac{1}{n} \sum \log x_i + \log\left(\frac{\sum x_i}{n}\right)}_{\hat{\lambda} = \frac{n\hat{\alpha}}{\sum x_i} \quad \theta}$$

$$\begin{pmatrix} \hat{\lambda} \\ \hat{\alpha} \end{pmatrix} \sim N\left(\begin{pmatrix} \lambda \\ \alpha \end{pmatrix}, \mathbf{I}^{-1}(\theta)\right) \quad \theta = (\lambda, \alpha)$$

$$\frac{\partial^2 \ell}{\partial \lambda^2} = -\frac{n\alpha}{\lambda^2} \quad \frac{\partial^2 \ell}{\partial \lambda \partial \alpha} = -\frac{n}{\lambda}$$

$$\begin{aligned} \frac{\partial^2 \ell}{\partial \alpha^2} &= \frac{\partial}{\partial \alpha} \left[ -n\psi(\alpha) + \log \lambda + \sum \log x_i \right] \\ &= -n\psi'(\alpha) \end{aligned}$$

$$\mathbf{J}(\theta) = \begin{bmatrix} n\alpha/\lambda^2 & n/\lambda \\ n/\lambda & n\psi'(\alpha) \end{bmatrix} = n \begin{bmatrix} \alpha/\lambda^2 & 1/\lambda \\ 1/\lambda & \psi'(\alpha) \end{bmatrix}$$

$\Rightarrow -\ell''(\theta)$

$$\begin{aligned} \mathbf{I}^{-1}(\theta) &= \frac{1}{n} \begin{bmatrix} \psi'(\alpha) - \frac{1}{\lambda} & \\ -\frac{1}{\lambda} & \alpha/\lambda^2 \end{bmatrix} \left( \frac{1}{\frac{\alpha\psi'(\alpha)}{\lambda^2} - \frac{1}{\lambda^2}} \right) \\ &= \frac{\lambda^2}{n} \begin{bmatrix} \psi'(\alpha) - 1/\lambda & \\ -1/\lambda & \alpha/\lambda^2 \end{bmatrix} \left( \frac{1}{\alpha\psi'(\alpha) - 1} \right) \end{aligned}$$

$$a. \text{var}(\hat{\lambda}) = \frac{\lambda^2}{n} \left( \frac{\psi'(\alpha)}{\alpha\psi'(\alpha)-1} \right) \quad a. \text{var}(\hat{\alpha}) = \frac{\lambda^2 \alpha}{n (\alpha\psi'(\alpha)-1)}$$

$$\hat{\lambda} = \frac{n\hat{\alpha}}{\sum x_i}$$

$$a. \text{cov}(\hat{\alpha}, \hat{\lambda}) = -\frac{\lambda}{n} \cdot \frac{1}{\alpha\psi'(\alpha)-1}$$

$$\hat{\lambda} \pm 1.96 \text{ se}(\hat{\lambda}) = \hat{\lambda} \pm 1.96 \sqrt{\frac{\lambda^2}{n} \left( \frac{\psi'(\alpha)}{\alpha\psi'(\alpha)-1} \right)}$$

$$\left( \frac{\lambda^2}{n\alpha} \right)$$

$$\frac{\lambda^2}{n\alpha} \left[ \frac{\psi'(\alpha)}{\alpha\psi'(\alpha)-1} \right]$$

$$\alpha > 0$$

$$\frac{1}{P(\alpha)} \lambda^\alpha x_i^{\alpha-1} e^{-\lambda x_i}$$

$$\mu = \alpha/\lambda$$

$$\lambda = \alpha/\mu$$

$$f(x_i; \alpha, \mu) = \frac{1}{P(\alpha)} \left( \frac{\alpha}{\mu} \right)^\alpha x_i^{\alpha-1} e^{-x_i(\alpha/\mu)}$$

$$\lambda, \alpha \quad \mu = \frac{\alpha}{\lambda} \Rightarrow g(\lambda, \alpha) = \frac{\alpha}{\lambda}$$

$$I(\alpha, \mu) = \begin{bmatrix} \sim & 0 \\ 0 & \sim \end{bmatrix}$$

asympt  
 $\hat{\mu}, \hat{\alpha} \perp$

$$\hat{\mu} = \frac{\sum x_i}{n}$$

$$\mu = \mathbb{E} X_i \quad \hat{\mu} = \bar{X}$$

$$\hat{\lambda} = \frac{n\hat{\alpha}}{\sum x_i}$$

$$p \times 1 \quad \underline{l}'(\hat{\theta}_n) = \underline{0} = l'(\theta) + l''(\theta)(\hat{\theta}_n - \theta) + \frac{1}{2}(\hat{\theta}_n - \theta)^T l'''(\theta)(\hat{\theta}_n - \theta)$$

$$\begin{aligned} \{ -l''(\theta) \}^{-1} l'(\theta) &= (\hat{\theta}_n - \theta) + \frac{1}{2} (-l''(\theta))^{-1} \left( \begin{matrix} \text{---} \\ \text{---} \\ \text{---} \end{matrix} \right) \\ &= (\hat{\theta}_n - \theta) \left[ \mathbb{I} + \frac{1}{2} \dots \right] \end{aligned}$$

$$\frac{\partial l(\theta)}{\partial \theta_j} = 0 = \frac{\partial l}{\partial \theta_j}(\theta) + \sum_k (\hat{\theta}_k - \theta_k) \frac{\partial^2 l}{\partial \theta_j \partial \theta_k} (\hat{\theta}_k - \theta_k) + \sum_{k,l} (\hat{\theta}_k - \theta_k) \frac{\partial^3 l}{\partial \theta_j \partial \theta_k \partial \theta_l} (\hat{\theta}_k - \theta_k)$$

$$-\frac{\partial l(\theta)}{\partial \theta_j} = (\hat{\theta}_j - \theta_j) \sum_{k=1}^p \left( \frac{\partial^2 l}{\partial \theta_j \partial \theta_k} \right) (\hat{\theta}_k - \theta_k) + (\hat{\theta}_j - \theta_j) \sum_{k,l} \left( \frac{\partial^3 l}{\partial \theta_j \partial \theta_k \partial \theta_l} \right) (\hat{\theta}_k - \theta_k)$$

$$\begin{aligned} (*) &= \sum_{k,l} \frac{\partial^3 l}{\partial \theta_j \partial \theta_k \partial \theta_l} \cdot (\hat{\theta}_k - \theta_k) (\hat{\theta}_l - \theta_l) \\ &\quad (p+p+p) \cdot (p \times p) \end{aligned}$$

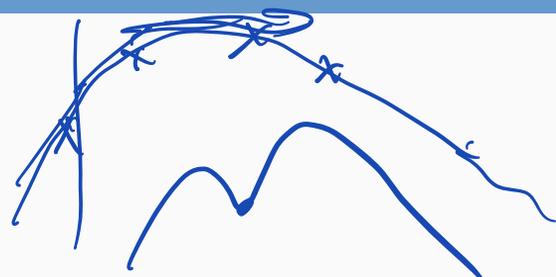
# Calculating maximum likelihood estimators

Newton-Raphson:

$$f(x) = 0 = f(x_0) + \nabla f(x_0) \overset{\text{any start}}{\text{value}}$$

$$0 = l'(\hat{\theta}_n) \approx l'(\theta_0) + l''(\theta_0)(\hat{\theta} - \theta_0)$$

$$\hat{\theta} \approx \theta_0 - \{l''(\theta_0)\}^{-1} l'(\theta_0)$$



$$l'(\hat{\theta}_n) = 0$$

$$= l'(\theta_0) +$$

- suggests iteration

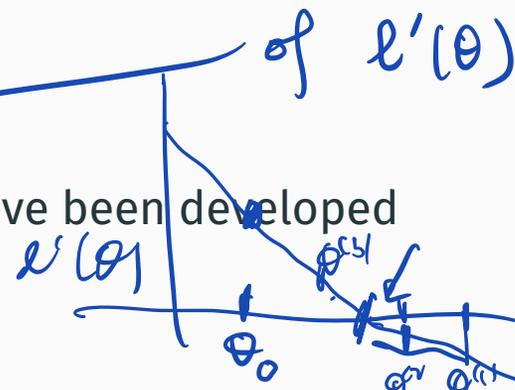
$$\hat{\theta}^{(k+1)} = \hat{\theta}^{(k)} + \{-l''(\hat{\theta}^{(k)})\}^{-1} l'(\hat{\theta}^{(k)}) =$$

Scalar  $\theta$

$$\hat{\theta}^{(k)} + \frac{S(\hat{\theta}^{(k)}) = l'}{H(\hat{\theta}^{(k)}) \leftarrow -l''}$$

MS p.270; note change in notation

- requires reasonably good starting values for convergence
- need  $-l''(\hat{\theta}^{(k)})$  to be non-negative definite
- Fisher scoring replaces  $-l''(\cdot)$  by its expected value  $J(\cdot)$
- N-R and F-S are gradient methods; many improvements have been developed
- solution is a global max only if  $l(\theta)$  is concave



E-M algorithm:

procedure

- complete data  $\mathbf{X} = (X_1, \dots, X_n)$ ,  $X_i$  i.i.d.  $f_X(x; \theta) \rightarrow \underline{\ell(\theta; \mathbf{x})}$
- observed data  $\mathbf{y} = (y_1, \dots, y_m)$ , with  $y_i = \underline{g_i(\mathbf{x})}$  many-to-one
- joint density  $\underline{f_Y(y; \theta)} = \int_{A(y)} f_X(x; \theta) dx$   $A(y) = \{x; y_i = g_i(x), i = 1, \dots, m\}$
- algorithm:
  1. (E step) estimate the complete data log-likelihood function for  $\theta$  using current guess  $\hat{\theta}^{(k)}$
  2. (M step) maximize that function over  $\theta$  and update to  $\hat{\theta}^{(k+1)}$  usually by N-R or Fisher scoring
- likelihood function increases at each step  $\leftarrow$
- can be implemented in complex models  $\leftarrow$
- doesn't automatically provide an estimate of the asymptotic variance  $\leftarrow$   
but methods exist to obtain this as a side-product

# Example

$N(\mu, \sigma_1^2)$        $N(\mu_2, \sigma_2^2)$



- $f_X(x_i; \lambda, \mu, \theta) = \alpha \frac{e^{-\lambda} \lambda^{x_i}}{x_i!} + (1 - \alpha) \frac{e^{-\mu} \mu^{x_i}}{x_i!}, \quad x = 1, 2, \dots; \lambda, \mu > 0, 0 < \theta < 1$
- Observed data:  $x_1, \dots, x_n$
- Complete data:  $(x_1, y_1), \dots, (x_n, y_n); y_i \sim \text{Bernoulli}(\alpha)$
- Complete data log-likelihood function:

$f(x_i | y_i = 1) = \lambda^{x_i} e^{-\lambda} / x_i!$   
 $p(y_i = 1) = \alpha$   
 $\mu^{x_i} e^{-\mu} / x_i!$

$\hat{\theta}^{(k)} = \begin{bmatrix} \hat{\alpha}^{(k)} \\ \hat{\lambda}^{(k)} \\ \hat{\mu}^{(k)} \end{bmatrix}$

$l_c(\alpha, \lambda, \mu; \mathbf{y}, \mathbf{x}) = \sum_{i=1}^n y_i \{ \log(\alpha) + x_i \log(\lambda) - \lambda \} + \sum_{i=1}^n (1 - y_i) \{ \log(1 - \alpha) + x_i \log(\mu) - \mu \}$

$E_{\hat{\theta}^{(k)}} \{ l_c(\alpha, \lambda, \mu; \mathbf{y}, \mathbf{x}) | \mathbf{x} \} = \sum_{i=1}^n \hat{y}_i \{ \log(\alpha) + x_i \log(\lambda) - \lambda \} + \sum_{i=1}^n (1 - \hat{y}_i) \{ \log(1 - \alpha) + x_i \log(\mu) - \mu \}$

$\hat{y}_i = E(Y_i | x_i; \hat{\theta}^{(k)})$

- maximizing values of  $\alpha, \lambda, \mu$  can be obtained in closed form

see p.280 for exact value  
p.281

AoS likes to work with  $\log \mathcal{L}_n(\theta) / \mathcal{L}_n(\hat{\theta}^{(k)})$

# ... Example

$$\theta = \begin{pmatrix} \alpha \\ \mu \\ \lambda \end{pmatrix}$$

$$\hat{\alpha}^{(k+1)} = \frac{1}{n} \sum y_i^{(k)}$$

$$\hat{\lambda}^{(k+1)} = \frac{\sum_{i=1}^n x_i y_i^{(k)}}{\sum y_i^{(k)}}$$

$$\hat{\mu}^{(k+1)} = \frac{\sum x_i (1 - y_i^{(k)})}{\sum \{1 - y_i^{(k)}\}}$$

$$E(Y_i | X_i = x_i; \alpha, \lambda, \mu) = \frac{\alpha e^{-\lambda x_i} \lambda^{x_i}}{\alpha e^{-\lambda x_i} \lambda^{x_i} + (1-\alpha) e^{-\mu x_i} \mu^{x_i}}$$

$\hat{\alpha}^{(k)}$     $\hat{\mu}^{(k)}$     $\hat{\lambda}^{(k)}$

## General-purpose Optimization

?optim

### Description

General-purpose optimization based on Nelder–Mead, quasi-Newton and conjugate-gradient algorithms. It includes an option for box-constrained optimization and simulated annealing.

### Usage

```
optim(par, fn, gr = NULL, ...,  
      method = c("Nelder-Mead", "BFGS", "CG", "L-BFGS-B", "SANN",  
                 "Brent"),  
      lower = -Inf, upper = Inf,  
      control = list(), hessian = FALSE)
```

quasi-Newton

```
optimHess(par, fn, gr = NULL, ..., control = list())
```

## Notes on optimization: Tibshirani, Pena, Kolter CO 10-725 CMU

- Goal:  $\max_{\theta} \ell(\theta; \mathbf{x})$  ←
- Solve:  $\ell'(\hat{\theta}; \mathbf{x}) = \mathbf{0}$  ←
- Iterate:  $\hat{\theta}^{(t+1)} = \hat{\theta}^{(t)} + \{j(\hat{\theta}^{(t)})\}^{-1} \ell'(\hat{\theta}^{(t)})$  ←
- Rewrite:  $j(\hat{\theta}^{(t)})(\hat{\theta}^{(t+1)} - \hat{\theta}^{(t)}) = \ell'(\hat{\theta}^{(t)})$   $B\Delta\theta = -\nabla\ell(\theta)$
- Quasi-Newton:
  - approximate  $j(\hat{\theta}^{(t)})$  with something easy to invert
  - use information from  $j(\hat{\theta}^{(t)})$  to compute  $j(\hat{\theta}^{(t+1)})$
- optimization notes add a step size to the iteration  $\hat{\theta}^{(t+1)} = \hat{\theta}^{(t)} + \epsilon_t \{j(\hat{\theta}^{(t)})\}^{-1} \ell'(\hat{\theta}^{(t)})$

```
optim(par, fn, gr = NULL, ...,
      method = c("Nelder-Mead", "BFGS", "CG", "L-BFGS-B", "SANN", "Brent"),
      lower = -Inf, upper = Inf, control = list(), hessian = FALSE)
```

# Regularity conditions

$$\frac{d}{d\theta} \int \dots = \int \frac{d}{d\theta} \dots$$

• (B1) The parameter space  $\Theta$  is an open subset of  $\mathbb{R}^p$

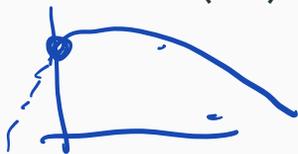
\* • (B2) The set  $A = \{x : f(x; \theta) > 0\}$  does not depend on  $\theta$

• (B3)  $l(\theta)$  is three times continuously differentiable on  $A$

? • (B4)  $E_{\theta}\{l'(\theta; X_i)\} = 0 \forall \theta$  and  $\text{Cov}\{l'(\theta; X_i)\} = I(\theta)$  is positive definite  $\forall \theta$

• (B5)  $E_{\theta}\{-l''(\theta; X_i)\} = J(\theta)$  is positive definite  $\forall \theta$

• (B6) For each  $\theta, \delta > 0, 1 \leq j, k, l, \leq p,$



$$\left| \frac{\partial^3 l(\theta^*; x_i)}{\partial \theta_j \partial \theta_k \partial \theta_l} \right| \leq M_{jkl}(\theta^*),$$

for  $\|\theta - \theta^*\| \leq \delta$ , where  $E_{\theta}\{M_{jkl}(X_i)\} < \infty$

$$\sqrt{n} (\hat{\theta}_n - \theta) = \left\{ \sum_{i=1}^n -l''(\theta; X_i) \right\}^{-1} \sum_{i=1}^n l'(\theta; X_i) / \sqrt{n} + R_n$$

$$l(\theta) = l(\theta^*) + l'(\theta^*)(\hat{\theta} - \theta^*) + \frac{1}{2} l''(\theta^*)(\hat{\theta} - \theta^*)^2$$

$$\frac{1}{2} l'''(\theta_n^*)(\hat{\theta}_n - \theta)^2$$

$$\frac{1}{2} (\hat{\theta} - \theta) l'''(\theta_n^*) (\hat{\theta} - \theta)^T$$

$$\frac{1}{2} (\hat{\theta}_n - \theta)^2 l'''(\theta_n^*)$$

$$\|\theta_n^* - \theta\| < \|\hat{\theta}_n - \theta\|$$

- (B1) The parameter space  $\Theta$  is an open subset of  $\mathbb{R}^p$
- (B2) The set  $A = \{x : f(x; \theta) > 0\}$  does not depend on  $\theta$
- (B3)  $\ell(\theta)$  is three times continuously differentiable on  $A$
- (B4)  $E_\theta\{\ell'(\theta; X_i)\} = 0 \forall \theta$  and  $\text{Cov}\{\ell'(\theta; X_i)\} = I(\theta)$  is positive definite  $\forall \theta$
- (B5)  $E_\theta\{-\ell''(\theta; X_i)\} = J(\theta)$  is positive definite  $\forall \theta$
- (B6) For each  $\theta, \delta > 0, 1 \leq j, k, l, \leq p,$

$$\left| \frac{\partial^3 \ell(\theta_n^*; X_i)}{\partial \theta_j \partial \theta_k \partial \theta_l} \right| \leq \underline{M_{jkl}(\theta^*)},$$

for  $\|\theta - \theta^*\| \leq \delta$ , where  $E_\theta\{\underline{M_{jkl}(X_i)}\} < \infty$

$$\frac{n(\hat{\theta}_n - \theta)^2}{n} \ell'''(\theta_n^*)$$

↓ p

$$E_\theta \ell'''(\theta; X_i)$$

$$M_{jkl}(\theta)$$

$$0 = l'(\hat{\theta}_n) = l'(\theta) + (\hat{\theta}_n - \theta) l''(\theta) + \frac{1}{2} (\hat{\theta}_n - \theta)^2 l'''(\theta^*)$$

- (B1) The parameter space  $\Theta$  is an open subset of  $\mathbb{R}^p$
- (B2) The set  $A = \{x : f(x; \theta) > 0\}$  does not depend on  $\theta$
- (B3)  $l(\theta)$  is three times continuously differentiable on  $A$
- (B4)  $E_\theta\{l'(\theta; X_i)\} = 0 \forall \theta$  and  $\text{Cov}\{l'(\theta; X_i)\} = I(\theta)$  is positive definite  $\forall \theta$
- (B5)  $E_\theta\{-l''(\theta; X_i)\} = J(\theta)$  is positive definite  $\forall \theta$
- (B6) For each  $\theta, \delta > 0, 1 \leq j, k, l, \leq p,$

Handwritten notes and derivations:

$$0 = l'(\theta) + (\hat{\theta}_n - \theta) \left\{ l'' + \frac{1}{2} (\hat{\theta}_n - \theta) l''' \right\}$$

$$\frac{l'(\theta)}{-l''} + (\hat{\theta}_n - \theta) \left\{ 1 + \frac{1}{2} (\hat{\theta}_n - \theta) \frac{l'''}{-l''} \right\}$$

$$\sqrt{n}(\hat{\theta}_n - \theta) = \dots$$

$$\left| \frac{\partial^3 l(\theta^*; x_i)}{\partial \theta_j \partial \theta_k \partial \theta_l} \right| \leq M_{jkl}(\theta^*),$$

for  $\|\theta - \theta^*\| \leq \delta$ , where  $E_\theta\{M_{jkl}(X_i)\} < \infty$

$$+ \frac{1}{2} (\hat{\theta}_n - \theta)^2 l'''(\theta^*; x_i) = (\hat{\theta}_n - \theta)^2 \frac{\sum l'''(\theta^*; x_i)}{n^{3/2}}$$

$\sqrt{n}$

- model assumption  $X_1, \dots, X_n$  i.i.d.  $f(x; \theta), \theta \in \Theta$
- true distribution  $X_1, \dots, X_n$  i.i.d.  $F(x)$
- maximum likelihood estimator based on model:

$$\sum_{i=1}^n \ell'(\hat{\theta}_n; X_i) = 0$$

- what is  $\hat{\theta}_n$  estimating?

$\subseteq \mathbb{R}^p$   
 $\leftarrow \ell(\theta; \underline{x})$

$\hat{\theta}_n \rightarrow \theta_{\neq}$   
 notation  
 }  
 under  
 model  
 $f(x; \theta_{\neq})$

$$E_{\theta_{\neq}} \{ \ell'(\theta_{\neq}; X_i) = 0 \}$$

$$E_{\theta}$$

- model assumption  $X_1, \dots, X_n$  i.i.d.  $f(x; \theta), \theta \in \Theta \subseteq \mathbb{R}$  (scalar)
- true distribution  $X_1, \dots, X_n$  i.i.d.  $F(x)$   $g(x)$  density
- maximum likelihood estimator based on model:

notation

$$\sum_{i=1}^n \ell'(\hat{\theta}_n; X_i) = 0$$

- what is  $\hat{\theta}_n$  estimating?
- define the parameter  $\theta(F)$  by

$$\int_{-\infty}^{\infty} \ell'\{x; \theta(F)\} dF(x) = 0$$

$$\sqrt{n}\{\hat{\theta}_n - \theta(F)\} \xrightarrow{d} N(0, \sigma^2)$$

$$\sigma^2 = \frac{\int [\ell'\{x; \theta(F)\}]^2 dF(x)}{\int [\ell''\{x; \theta(F)\}]^2 dF(x)}$$

$$E_{\uparrow F} \ell'(\theta(F); X_i) = 0$$

wrong model

$$\left(\frac{d}{d\theta} \log f(x; \theta)\right)$$

•

$$\sqrt{n}\{\hat{\theta}_n - \theta(F)\} \xrightarrow{d} N(0, \sigma^2)$$

•

$$\sigma^2 = \frac{\int [l'\{x; \theta(F)\}]^2 dF(x)}{(\int [l''\{x; \theta(F)\}]^2 dF(x))^2}$$

$E_{\mathbb{P}} [l''(\theta; x)] \neq E_{\mathbb{P}} [l'(\theta; x)^2]$

• more generally,

$$\sqrt{n}\{\hat{\theta}_n - \theta(F)\} \xrightarrow{d} N\{0, \underline{\underline{G^{-1}(F)}}\}$$

$\theta \in \mathbb{R}^p$

•

Godambe information

$$\hat{G}(F) = \hat{J}(F) \hat{I}^{-1}(F) \hat{J}(F)$$

sandwich estimator

$I^{-1}(\theta)$

•

$$J(F) = \int \underbrace{-l''\{\theta(F); x_i\}}_{p \times p} dF(x_i)$$

$$I(F) = \int \underbrace{\{l'(\theta(F); x_i)\} \{l'(\theta(F); x_i)\}^T}_{p \times p} dF(x_i)$$

Godambe information

# Multi-parameter example: logistic regression

$$\sqrt{n}(\hat{\theta}_n - \theta) \xrightarrow{d} N_p(0, I^{-1}(\theta))$$

suggests  $\hat{\theta}_n \sim N(\theta, (nI(\theta))^{-1})$

```
Boston$scrim2 <- Boston$scrim > median(Boston$scrim) # define binary response
Boston.glm <- glm(crim2 ~ . - crim, family = binomial,
                 data = Boston) #fit logistic regression
summary(Boston.glm)
```

$$\hat{\theta}_n \sim N_p(\theta, I(\hat{\theta}_n)^{-1})$$

$$\hat{\theta}_{n,j} \sim N(\theta_j, \{I^{-1}(\hat{\theta}_n)\}_{jj})$$

↑  
var.

Coefficients:

	Estimate	Std. Error	z value	Pr(> z )	
(Intercept)	-34.103704	6.530014	-5.223	1.76e-07	***
zn $\theta_j$	-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	

testing  $\theta_{age} = 0$

## ... Example: logistic regression

```
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

Waiting for profiling to be done – what's profiling?

Waiting for profiling to be done – what's profiling?

Waiting for profiling to be done – what's profiling?

Waiting for profiling to be done – what's profiling?

Waiting for profiling to be done – what's profiling?

## 4.1 · Likelihood

95

**Figure 4.1** Likelihoods for the spring failure data at stress 950 N/mm<sup>2</sup>. 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  $\alpha$ .

