

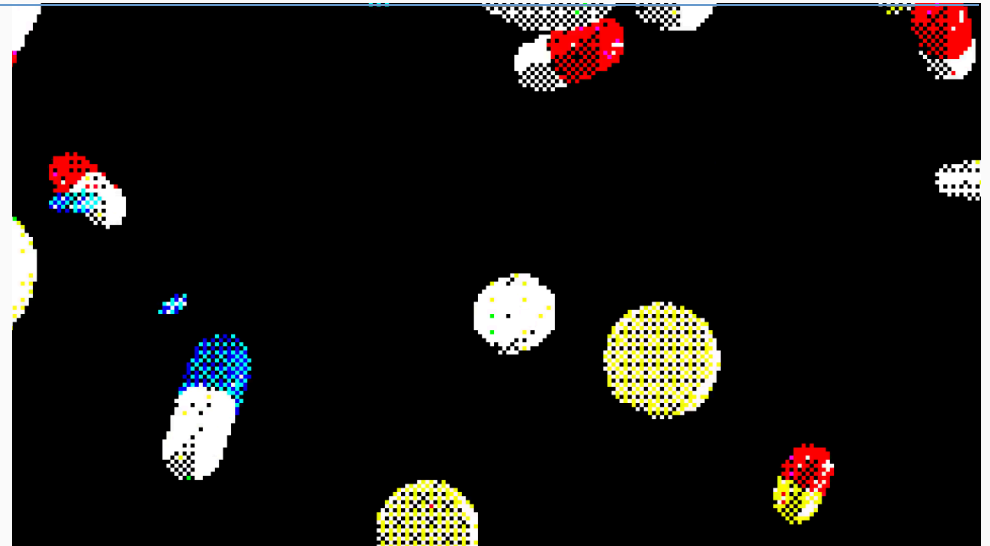
# Statistical Theory for Data Science

STA2212H S LEC9101

The  
Economist

Week 3

January 20 2026



1. Recap: delta method, Bayesian inference
2. Misspecified models; point and interval estimation
3. Nonparametric estimation
4. Statistics in the News
5. HW Questions

## Upcoming

- Toronto Data Workshop, Wednesday 21 January 2026, noon (EST) on [Zoom](#)  
Rachel Porter, University of Notre Dame  
“The CampaignView collection of records about candidates and their policies”
- Thursday 22 January 2026, 11am Hydro 9195/9  
Isaac Gibbs, Job Talk

# Recap: Delta method

$$X \sim N(0, \sigma^2) \quad \sum_{i=1}^n X_i^2 / n = \hat{\sigma}^2$$

- variance-stabilizing transformations

↑

$\text{var}(\bar{X})$  depends on  $E(\bar{X})$   
 $\mu$

then normal approx<sup>n</sup> could be poor

Example  $X_i \sim \text{Po}(\mu)$   $i=1, \dots, n$  iid  $E(\bar{X}) = \mu$   $\text{var}(\bar{X}) = \mu$

$$\text{var}\{g'(\bar{X})\} = c = \frac{1}{n} g'(\mu)^2 \cdot \mu \propto c$$

$$\text{var}(\bar{X}^{1/2}) \approx \frac{1}{4n} E\bar{X} \approx \mu^{1/2} \quad g'(\mu) \propto \frac{1}{\sqrt{\mu}} \quad g(\mu) \propto \sqrt{\mu}$$

$$\text{var}g(\hat{\theta}) \approx g'(E\hat{\theta})^2 \cdot \text{var}\hat{\theta}$$

$g(\cdot)$  1-1 f<sup>n</sup> of  $\theta$

$$\frac{\partial g}{\partial \theta^T}(E\hat{\theta}) \cdot \text{var}\hat{\theta} \frac{\partial g(E\hat{\theta})}{\partial \theta^T}$$

## HIGH DIMENSIONAL GAUSSIAN AND BOOTSTRAP APPROXIMATIONS IN GENERALIZED LINEAR MODELS

MAYUKH CHOUDHURY AND DEBRAJ DAS

$$d = o(n^{1/2})$$
$$\frac{d}{\sqrt{n}} \rightarrow 0$$
$$n \rightarrow \infty$$

Link

ABSTRACT. Generalized Linear Models (GLMs) extend ordinary linear regression by linking the mean of the response variable to covariates through appropriate link functions. This paper investigates the asymptotic behavior of GLM estimators when the parameter dimension  $d$  grows with the sample size  $n$ . In the first part, we establish Gaussian approximation results for the distribution of a properly centered and scaled GLM estimator uniformly over class of convex sets and Euclidean balls. Using high-dimensional results from Fang and Koike (2024) for the leading Bahadur term, bounding remainder terms as in He and Shao (2000), and applying Nazarov's (2003) Gaussian isoperimetric inequality, we show that Gaussian approximation holds when  $d = o(n^{2/5})$  for convex sets and  $d = o(n^{1/2})$  for Euclidean balls—the best possible rates matching those for high-dimensional sample means. We further extend these results to the bootstrap approximation when the covariance matrix is unknown. In the second part, when  $d \gg n$ , a natural question is to answer whether all covariates are equally important. To answer that, we employ sparsity in GLM through the Lasso estimator. While Lasso is widely used for variable selection, it cannot

(C.1)  $y_i \in \mathbf{R}$  for all  $i$ ,  $h$  is the identity function and  $b(u) = u^2/2$  (which is the case for linear regression) or,  $y_i \geq 0$  for all  $i$  and  $-h$  &  $h_1$  are strictly convex (which is the case for logistic, poisson, gamma regression etc).

(C.2)  $h$  is thrice continuously differentiable and  $g^{-1}$  is twice continuously differentiable.

(C.3)  $\mathbf{S}_n$  and  $\mathbf{L}_n$  are positive definite matrices.

(C.4) For some  $\alpha_1, \alpha_2, \alpha_3 > 0$ , satisfying  $\alpha_1 + 2\alpha_2 + 2\alpha_3 < 1/2$  with  $0 < \alpha_2 < 1/14$ , there exist a constant  $0 < c_1 < \infty$ , and two other constants  $0 < c_2 < c_3 < \infty$ , for which following holds:  $\lambda_{\min}(\mathbf{S}_n) \geq c_1 n^{-\alpha_1}$  and  $c_2 n^{-\alpha_2} \leq \lambda_{\min}(\mathbf{E}(\mathbf{L}_n)) < \lambda_{\max}(\mathbf{E}(\mathbf{L}_n)) \leq c_3 n^{\alpha_3}$ .

(C.5)  $\max_{i \in \{1, \dots, n\}} \|\mathbf{x}_i\|_{\infty} = O(1)$ .

(C.6)  $\|\boldsymbol{\beta}\|_{\infty} = O(1)$ .

(C.7)  $n^{-1} \sum_{i=1}^n \mathbf{E}|y_i|^6 = O(1)$ .

(C.8)  $n^{-1} \sup_{\{\|\boldsymbol{\alpha}\|, \|\boldsymbol{\kappa}\|=1\}} \sum_{i=1}^n |\boldsymbol{\alpha}^{\top} \mathbf{x}_i|^8 |\boldsymbol{\kappa}^{\top} \mathbf{x}_i|^8 = O(1)$ .

(C.9.i)  $\max_{i \in \{1, \dots, n\}} |h'(\mathbf{x}_i^{\top} \boldsymbol{\beta})| = O(1)$ .

(C.9.ii)  $\max_{i \in \{1, \dots, n\}} |h''(\mathbf{x}_i^{\top} \boldsymbol{\beta})| = O(1)$ .

(C.9.iii)  $\max_{i \in \{1, \dots, n\}} |(g^{-1})'(\mathbf{x}_i^{\top} \boldsymbol{\beta})| = O(1)$ .

(C.9.iv)  $n^{-1} \sum_{i=1}^n \sup_{\{|z_i - \mathbf{x}_i^{\top} \boldsymbol{\beta}| < \delta\}} |(g^{-1})''(z_i)|^8 = O(1)$ .

(C.9.v)  $n^{-1} \sum_{i=1}^n \sup_{\{|z_i - \mathbf{x}_i^{\top} \boldsymbol{\beta}| < \delta\}} |h'''(z_i)|^{12} = O(1)$ .

$$E(y_i | \mathbf{x}_i) = \mu_i$$

$$g(\mu_i) = \mathbf{x}_i^{\top} \boldsymbol{\beta}$$

$$g^{-1}(\mathbf{x}_i^{\top} \boldsymbol{\beta}) = h(\mathbf{x}_i^{\top} \boldsymbol{\beta})$$

$$= g^{-1}(\mu_i)$$

1. The support of the density for  $X$ , i.e., the set  $\{x : f(x; \theta) > 0\}$  does not depend on  $\theta$ .
2. The true value  $\theta_*$  is contained in an open subset of  $\Theta$ , and for all  $\theta$  in this open subset, the density  $f(x; \theta)$  is differentiable with respect to  $\theta$  for all  $x$  in the support of the density.
3.  $f(x; \theta)$  is three times continuously differentiable with respect to  $\theta$  for all  $x$  in the support of the density.
4.  $E_\theta\{u(\theta; X)\} = \mathbf{0}$ , and  $E_\theta\{u(\theta; X)u^\top(\theta; X)\} = E_\theta\{-\partial u(\theta; X)/\partial\theta^\top\}$ , and these expectations exist and are finite for all  $\theta$  in the open subset defined in 2.
5. The matrix  $I_1(\theta) = E_\theta\{u(\theta; X)u^\top(\theta; X)\}$  is positive definite for all  $\theta$  in the open subset defined in 2.
6. There exist functions  $M_{abc}(\cdot)$  such that

$$\left| \frac{\partial^3}{\partial\theta_a\partial\theta_b\partial\theta_c} \ell(\theta; \mathbf{x}) \right| \leq M_{abc}(\mathbf{x}), \quad \|\theta - \theta_*\| \leq \delta \quad \text{and} \quad E_{\theta_*}\{M_{abc}(X)\} < \infty.$$

$X_1, \dots, X_n$  i.i.d.  $f(x; \theta), \theta \in \Theta$

# Bayesian inference

Under regularity conditions on the model, and the prior, the posterior is asymptotically normal.

More precisely

$$\int_{a_n}^{b_n} \pi(\theta | \mathbf{x}) d\theta \xrightarrow[n \rightarrow \infty]{p} \int_a^b \phi(t) dt$$

$\mathbf{x} \sim$  under the model for  $X$

obs = iid  $X$   
 $\mathbf{x} = (x_1, \dots, x_n)$

$$a_n = \hat{\theta} + a I_n(\hat{\theta})^{-1/2}, \quad b_n = \hat{\theta} + b I_n(\hat{\theta})^{-1/2}$$

center + scaling

$I_n(\theta)$  is the expected Fisher information in the sample, and can be replaced (to this order of approximation) by the observed Fisher information  $J(\hat{\theta})$ .

$J(\hat{\theta})$  in my notation; text defines  $J = I^{-1}$  and uses  $H$  for  $-\ell''$

Informally,  $\pi(\theta | \mathbf{x}) \sim N\{\hat{\theta}, I_n^{-1}(\hat{\theta})\}$ , and an approximate 95% credible interval for  $\theta$  is  $\hat{\theta} \pm 1.96 \hat{se}$

same as 95% CI

$$I_n^{1/2}(\hat{\theta}) \text{ or } J^{1/2}(\hat{\theta})$$

Laplace approximation

Laplace approx<sup>17</sup>

equivalent by,  $\pi(\theta|x) \approx N(\hat{\theta}_\pi, J_{\pi\pi}^{-1}(\hat{\theta}))$  (\*)

$$\hat{\theta}_\pi = \underset{\theta}{\operatorname{argsup}} \pi(\theta|x)$$

under ref. on  $f(v, \theta)$  &  $\pi(\theta)$

$$J_{\pi\pi} = - \frac{\partial^2 \log \pi(\theta|x)}{\partial \theta^2} \Big|_{\theta = \hat{\theta}_\pi}$$

$$\pi(\theta|x) = \frac{L(\theta; x) \pi(\theta)}{m(x)} \leftarrow \max_{\theta}$$

# Laplace approximation

$$\pi(\theta|x) = \frac{L(\theta;x)\pi(\theta)}{\int L(\theta;x)\pi(\theta)d\theta} \quad \theta \in \mathbb{R}$$

$$\begin{aligned} \text{den} &= \int e^{l(\theta;x)} \pi(\theta) d\theta \xrightarrow{\text{green arrow}} \int e^{l(\theta;x) + \ln \pi(\theta)} d\theta \\ &\stackrel{\text{green arrow}}{=} \int e^{l(\hat{\theta}) + (\theta - \hat{\theta})l'(\hat{\theta}) + \frac{1}{2}(\theta - \hat{\theta})^2 l''(\hat{\theta})} \left\{ \pi(\hat{\theta}) + (\theta - \hat{\theta})\pi'(\hat{\theta}) \right\} d\theta \\ &= \pi(\hat{\theta}) e^{l(\hat{\theta})} \cdot \int e^{-\frac{1}{2}(\theta - \hat{\theta})^2 \{ -l''(\hat{\theta}) \}} \left\{ \frac{1}{\pi} + (\theta - \hat{\theta}) \frac{\pi'(\hat{\theta})}{\pi(\hat{\theta})} \right\} d\theta \\ &= \underbrace{\sqrt{2\pi} \pi(\hat{\theta}) e^{l(\hat{\theta})} j(\hat{\theta})^{-1/2}}_{\text{green box}} \cdot \int \underbrace{j(\hat{\theta})^{1/2}}_{\text{green box}} e^{-\frac{1}{2}j(\hat{\theta})(\theta - \hat{\theta})^2} \left\{ 1 + (\theta - \hat{\theta}) \frac{j'(\hat{\theta})}{j(\hat{\theta})} \right\} d\theta \end{aligned}$$

$$= \sqrt{2\pi} \pi(\hat{\theta}) e^{\ell(\hat{\theta})} \mathcal{J}(\hat{\theta})^{-1/2}$$

$$\pi(\theta | \underline{x}) \doteq \frac{1}{\sqrt{2\pi}} e^{\frac{\ell(\theta) - \ell(\hat{\theta})}{\sqrt{2\pi}}} \cdot \pi(\hat{\theta}) \mathcal{J}(\hat{\theta})^{-1/2}$$

$$\theta \in \mathbb{R}^d \doteq \frac{1}{\sqrt{2\pi}} d e^{\ell(\underline{\theta}) - \ell(\hat{\underline{\theta}})} \cdot \pi(\hat{\underline{\theta}}) |\mathcal{J}(\hat{\underline{\theta}})|^{-1/2}$$

$d_n \uparrow n$

( $\hat{\theta} = \text{MLE}$ )  
 $\mathcal{J}(\hat{\theta}) = \text{obs'd}$

INLA integrated nested Laplace approx<sup>n</sup>  
 $\uparrow$   
 (normal version)

- conjugate priors

• posterior is in the same class as prior same form

$$f(x; \theta) = e^{\eta(\theta) \cdot t(x) - \kappa(\theta) - d(x)} \quad \pi(\theta; \alpha) =$$

- non-informative priors

$$x \sim N(\mu, \sigma^2) \quad \text{flat, "ignorance"}$$

- convenience priors

$$\mu \sim N(a, b) \quad \beta_j \text{ iid } N(b, \sigma)$$

$$\sigma^2 \sim \text{Inv. Gamma}(\cdot, \cdot)$$

- minimally/weakly informative priors Gelman

⊥

- hierarchical priors

$$\beta \sim N(\underline{0}, \tau^2 \cdot I); \quad \tau^2 \sim \text{Inv. G}(\nu)$$

- empirical Bayes

HW 2 : estimate parameters of the prior

$$\pi(\theta | \underline{x}) = \frac{L(\theta; \underline{x}) \pi(\theta)^{\alpha, \beta}}{\int L(\theta; \underline{x}) \pi(\theta)^{\alpha, \beta} d\theta} \equiv m(\underline{x}^{\alpha, \beta})$$

could use  $m(\underline{x}; \alpha, \beta)$  +  $x_1, \dots, x_n$  to  
estimate  $\alpha$  &  $\beta$

non-informative

•  $\pi(\theta) \propto I^{1/2}(\theta)$

$I^{1/2}(\theta) = \pi(\theta)$        $I^{1/2}(\tau) = \pi(\tau)$   
 $\tau = \tau(\theta)$

• Example:  $X \sim \text{Bin}(n, \theta)$        $I(\theta) = n/\{\theta(1-\theta)\}$ ,       $0 < \theta < 1$

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

$\mathcal{N}(\mu, \sigma^2)$  (ntbc)

$\Rightarrow \pi(\mu, \sigma) = \frac{1}{\sigma^2}$

• for regression-scale models Jeffreys recommends  $\pi(\beta, \sigma) \propto 1/\sigma$  ← leads to t-p dist

• For normal theory linear regression, the conjugate prior is Normal for  $\beta$  and Inverse gamma for  $\sigma$ , with  $\beta \perp \sigma$ . These are often the defaults.

one Bayes estimator  $E(\theta | \underline{x})$

• Posterior means are typically weighted averages of MLE and prior mean, with weight on MLE  $\rightarrow 1$  with  $n$

sample estimate  $\rightarrow$  prior mean

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

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

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

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

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

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

$$\underline{m} \sim \mathcal{N}(\underline{\mu}_m, \sigma_m^2)$$

$$\underline{v} \sim \mathcal{N}(\underline{\mu}_v, \sigma_v^2)$$

$$\text{corr}(\underline{m}, \underline{v}) \text{ is } \theta$$

$$\text{Cov. } \theta \sigma_v \sigma_m$$

$$L(\mu_m, \mu_v, \sigma_v^2, \sigma_m^2, \theta; \underline{m}, \underline{v})$$

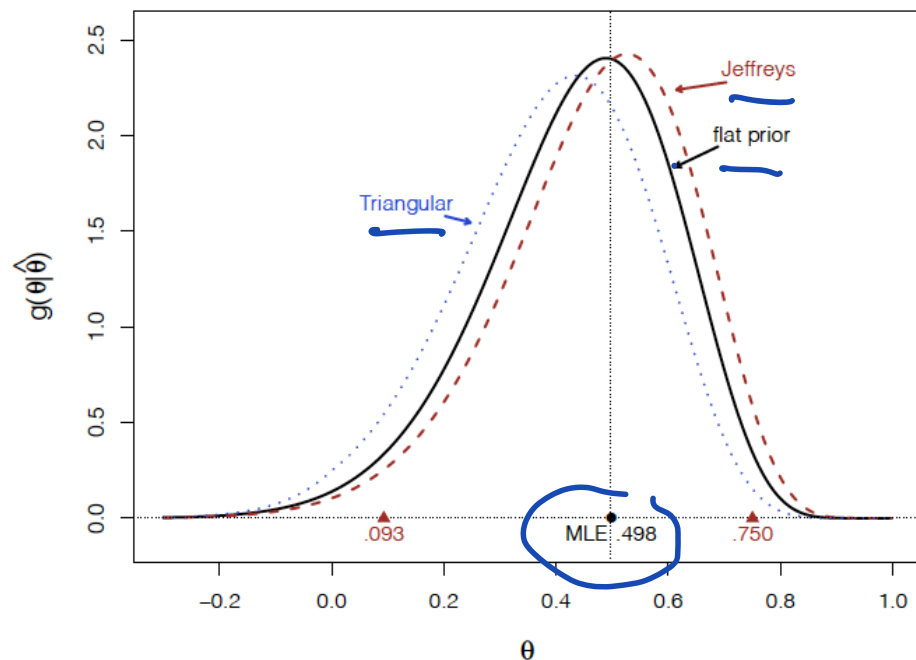
could find a  $f$  of

$\hat{\theta}, \theta$  only to use as a Lik.

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

↑

$$\pi(\theta|\hat{\theta}) \propto f(\hat{\theta}|\theta)\pi(\theta)$$



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

$$\pi_J = \frac{1}{1-\theta^2} \quad -1 \leq \theta \leq 1$$

$$\pi_{\text{flat}} = 1$$

$$\pi_{\text{triang}} \quad \begin{array}{c} \uparrow \\ \text{triangle} \\ \text{on } [-1, 1] \end{array}$$

11.2 · Inference

579

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

$x/n$

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

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

$$\hat{P}(\text{dying in A}) = 0$$

prior for hospital A  $Beta(1, 1)$

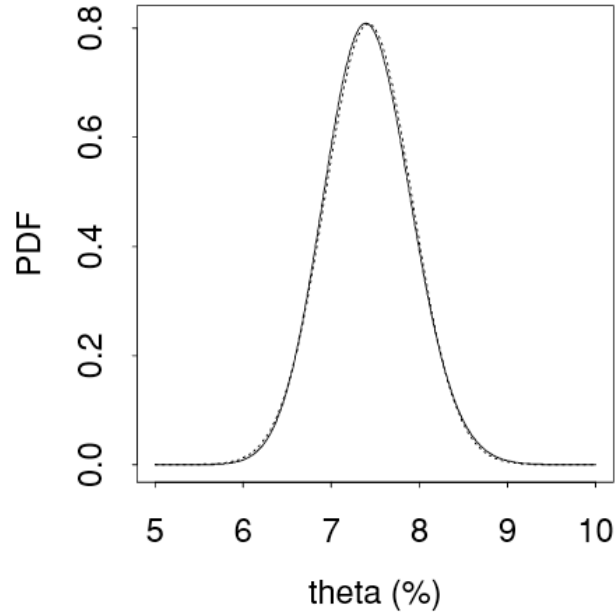
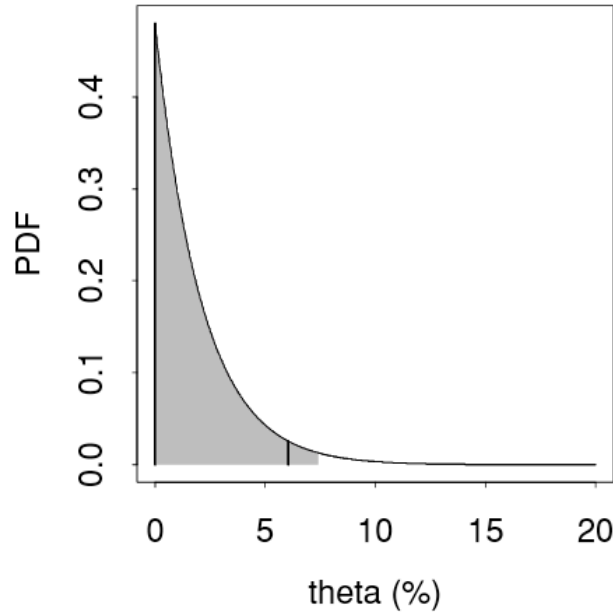
$$\frac{1}{49}$$

posterior mean  $\hat{\theta}_A$

posterior mean

580

11 · Bayesian Models



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

↓ same  $\theta$

put all hospitals together; 208 failures ‘

(or HW 2)

# Marginalization

- Bayes posterior carries all the information about  $\theta$ , given  $\mathbf{x}$  by definition

- probabilities for any set  $A$  computed using the posterior distribution

- $\text{pr}(\Theta \in A | \mathbf{x}) = \int_A \pi(\underline{\theta} | \underline{x}) d\underline{\theta}$

- if  $\theta = (\psi, \lambda), \dots$   
     $\uparrow \quad \uparrow$   
     $\leftarrow \int \pi(\psi, \underline{\lambda} | \underline{x}) d\underline{\lambda} = \pi_m(\psi | \underline{x})$   
    marginal post.

- or, if  $\psi = \psi(\theta) \rightarrow \int_{\{\underline{\theta} : \psi(\underline{\theta}) = \psi\}} \pi(\underline{\theta} | \underline{x}) d\underline{\theta}$

- in this context, 'flat' priors can have a large influence on the marginal posterior

$$\begin{aligned}
 X_i &\sim N(\mu_i, \frac{1}{n}) & X_i &\sim N(\mu_i, 1/n) & \pi(\mu_i) &= 1 & i=1, \dots, k \\
 \sqrt{n}X_i &\sim N(\mu_i\sqrt{n}, 1) & \pi(\mu_i | x_i) &\sim N(x_i, 1/n) & & & E(\psi | \underline{x}) = n + 2\sum x_i^2 \\
 nX_i^2 &\sim \chi^2_1(\mu_i^2 n) & & & & & E[\quad] = n + n + \sum \mu_i^2
 \end{aligned}$$

$$\pi(\underline{\mu} | \underline{x}) \propto \prod_{i=1}^k e^{-\frac{n}{2}(x_i - \mu_i)^2}$$

$$\psi = \sum_{i=1}^k \mu_i^2 \quad \int_{\{\underline{\mu} : \sum \mu_i^2 = \psi\}} \pi(\underline{\mu} | \underline{x}) d\underline{\mu} = \pi(\psi | \underline{x})$$

$$\text{posterior of } n \sum_{i=1}^k \mu_i^2 \sim \chi^2_k \left( n \sum_{i=1}^k x_i^2 \right) \quad n \sum_{i=1}^k x_i^2 \sim \chi^2_k \left( n \sum_{i=1}^k \mu_i^2 \right)$$

$$E(\tilde{\psi}_B | \underline{x}) = \frac{k}{n} + \sum x_i^2$$

$$E(E(\tilde{\psi}_B | \underline{x})) = \frac{2k}{n} + \sum \mu_i^2$$

$$= \frac{2k}{n} + \psi$$

- 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:

notation

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

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

$\hat{\theta}_n \rightarrow ??$  if it converges  
~~??~~

e.g.  $X_i \sim N(\mu, 1)$

$$\bar{X} \rightarrow E_F(X_i)$$

$$= \int x_i dF(x_i)$$

let  $\theta(F)$  be defined as

$$E_F \{ \ell'(\theta; X_i) \} = 0$$

$$\theta = \bar{\theta}(F)$$

- 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:

notation

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

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

- define the parameter  $\theta(F)$  by *(estmand)*

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

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

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

$$0 = \sum l'(\hat{\theta}; x_i) = \sum l'(\theta(F); x_i) + \left\{ \hat{\theta} - \theta(F) \right\}^x \sum l''(\theta(F); x_i)$$

$$\sqrt{n} \{ \hat{\theta}_n - \theta(F) \} \doteq \frac{\frac{1}{\sqrt{n}} \sum l'(\theta(F); x_i) + \dots}{-\frac{1}{n} \sum l''(\theta(F); x_i)} + R_n$$

$$\frac{1}{\sqrt{n}} \sum l'(\theta(F); x_i) \xrightarrow{d} N(0, \text{var}_F\{l'(\theta(F); x_i)\})$$

$$-\frac{1}{n} \sum l''(\quad) \xrightarrow{P} E_F\{l''(\theta(F); x_i)\}$$

no Bartlett identity

- $$\sqrt{n}\{\hat{\theta}_n - \theta(F)\} \xrightarrow{d} N(0, \sigma^2)$$
- $$\sigma^2 = \frac{\int [l'\{\theta(F); x\}]^2 dF(x)}{(\int [l''\{\theta(F); x\}]^2 dF(x))^2}$$
- more generally, for  $\theta \in \mathbb{R}^p$ ,
 
$$\sqrt{n}\{\hat{\theta}_n - \theta(F)\} \xrightarrow{d} N_p\{0, \underline{G}^{-1}(F)\}$$

$\sqrt{n}(\hat{\theta}_n - \theta) \rightarrow N(0, I^{-1}(\theta))$   
 $\uparrow$   
 Fisher info  
 $\swarrow$   
 Godambe info
- $$G(F) = J(F)I^{-1}(F)J(F), \quad G^{-1}(F) = J^{-1} I J^{-1}$$
- $$J(F) = \int -l''\{\theta(F); x_i\} dF(x_i), \quad I(F) = \int \{l'\{\theta(F); x_i\}\} \{l'\{\theta(F); x_i\}\}^T dF(x_i)$$

"Sandwich ~~est.~~ variance formula"  
 Godambe information sandwich variance

- Consider maximizing a function other than the log-likelihood function
- VdV notation  $M_n(\theta; \mathbf{x}) \equiv \frac{1}{n} \sum_{i=1}^n m(\theta; x_i)$  ——— *M est*
- Typically found by solving  $M'_n(\theta; \mathbf{x}) = 0$ , or

VdV Z-estimators

$$\Psi_n(\theta; \mathbf{x}) = \sum_{i=1}^n \Psi(\theta; x_i) = 0 \quad \text{Z-estimators}$$

- SM notation:  $g(\mathbf{Y}; \theta) = \sum_{i=1}^n g(Y_i; \theta)$  is a (set of) **estimating equations**
- $g(y; \theta)$  is an **unbiased estimating function** if

$$E\{g(Y; \theta)\} = 0 = \int g(y; \theta) f(y; \theta) dy$$

- the score function is an unbiased estimating function
- the solution of  $g(\mathbf{y}; \theta) = 0$  is an M-estimator

*model  $f(y; \theta)$   
not using  
score eq<sup>n</sup> to  
est.  $\theta$*

- Examples: moment estimators
- Example: median
- Example: Huber estimator

- Examples: moment estimators  $g(x; \theta) = (x^k - \mu_k)$   $\frac{1}{n} \sum X_i^k = \hat{\mu}_k = E(X)^k$  SM Ex.7.15

- Example: median  $g(x; \theta) = \text{sign}(x - \theta)$   $\frac{1}{n} \sum \text{sign}(X_i - \theta) = 0$  VdV Ex.5.4  
 $\hat{\theta}_n = \text{median}(X_1, \dots, X_n)$

- Example: Huber estimator  $g(x; \theta) = \text{sign}(x - \theta)$ , where SM Ex 7.19

$$g(x; \theta) = \begin{cases} -k, & x \leq \theta - k, \\ x, & -k < x - \theta < k \\ k, & x \geq \theta + k \end{cases}$$

$\bar{X}_n$  "truncated at  $\pm k$ "

- Example: quantile regression

$\rightarrow \hat{\beta}_\tau = \arg \min_{\beta} \rho_\tau(y_i - x_i^T \beta), \quad \rho_\tau(u) = u\{\tau - I(u \leq 0)\}$

Koenker<sup>1</sup>

<sup>1</sup>Koenker, Roger, and Kevin F. Hallock (2001). Quantile Regression. *J. Econ. Persp.* **15**, 143–56.  $\leftarrow$

- $\mathbf{0} = g(\mathbf{x}; \tilde{\theta}) =$

$$\hat{\theta}_n \quad \frac{1}{n} \sum g(x_i; \theta) = \underline{\underline{0}} = g(\mathbf{x}; \tilde{\theta}) = g(\mathbf{x}; \theta) + (\tilde{\theta} - \theta) \dot{g}(\mathbf{x}; \theta) + \dots$$

$$\tilde{\theta} - \theta \doteq \frac{\frac{1}{n} \sum g(x_i; \theta)}{-\frac{1}{n} \dot{g}(x_i; \theta)} \doteq \frac{\frac{1}{n} \sum g(x_i; \theta)}{E\{-\dot{g}(X; \theta)\}}$$

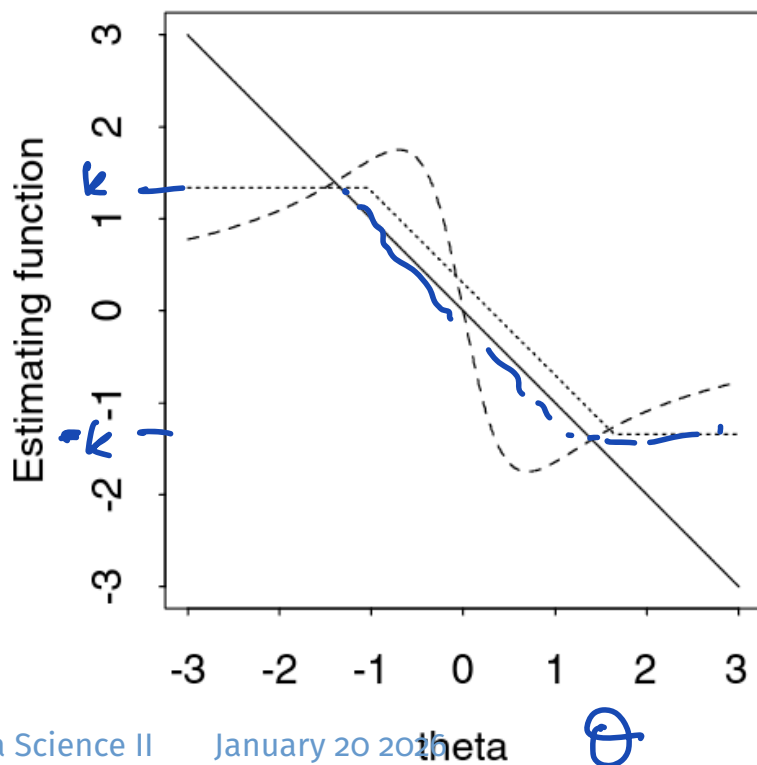
$$\dot{g}(x; \theta) = \frac{\partial}{\partial \theta} g(x; \theta)$$

$$E(\tilde{\theta} - \theta) \doteq \mathbf{0}, \quad \text{var}(\tilde{\theta} - \theta) \doteq \frac{1}{n} \text{var}\{g(X; \theta)\} / E^2\{\dot{g}(X; \theta)\} \quad (1)$$

$$\sqrt{n}(\tilde{\theta} - \theta) \xrightarrow{d} N\{\mathbf{0}, G^{-1}(\theta)\}, \quad (2)$$

$$G(\theta) = E\left\{ \underbrace{-\frac{\partial}{\partial \theta^T} g(X; \theta)} \right\} [\text{var}\{g(X; \theta)\}]^{-1} E\left\{ \underbrace{-\frac{\partial}{\partial \theta} g(X; \theta)} \right\} \quad (3)$$

7 · Estimation and Hypothesis Testing



**Figure 7.3** Estimating functions. Left: construction of  $g(y; \theta)$  (heavy) as the sum of  $g(y_j; \theta)$  for a sample of size  $n = 3$  shown by the rug. The lines  $g = 0$  (dots) and  $\theta = \tilde{\theta}$  (dashes) are also shown. Right: estimating functions for the mean (solid), the Huber estimator (dots) and a redescending M-estimator (dashes), slightly offset to avoid overplotting.

$g(x, \theta)$

theta  $\theta$

- a  $(1 - \alpha)$ -level confidence interval for  $\theta \in \mathbb{R}$  is  $\{L(\mathbf{x}), U(\mathbf{x})\}$ , with

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

$$\Pr_{\theta} \{L(\mathbf{X}) \leq \theta \leq U(\mathbf{X})\} \geq 1 - \alpha \quad \text{for all } \theta,$$

"with = for some  $\theta$ " (MS)

equality?

- similarly upper and lower confidence bounds

$$\int P_{\theta} \{L(\underline{x}) \leq \theta\} \geq 1 - \alpha$$

$$\int_{-\infty}^{\theta} l(\underline{x}) f(\underline{x}; \theta) d\underline{x} \geq 1 - \alpha$$

$$P_{\theta} \{u(\underline{x}) \geq \theta\} \geq 1 - \alpha$$

$$\int_{\theta}^{\infty} u(\underline{x}) f(\underline{x}; \theta) d\underline{x} \geq 1 - \alpha$$

- a  $(1 - \alpha)$ -level confidence interval for  $\theta \in \mathbb{R}$  is  $\{L(\mathbf{x}), U(\mathbf{x})\}$ , with

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

equality?

- similarly upper and lower confidence bounds
- a  $(1 - \alpha)$ -level confidence **region** for  $\theta \in \Theta$ , is a set  $R(\mathbf{X}) \subset \Theta$ , with

$$\text{pr}\{\theta \in R(\mathbf{X})\} \geq 1 - \alpha$$

$\theta \in \mathbb{R}^p$   
(say)

- a  $(1 - \alpha)$ -level confidence interval for  $\theta \in \mathbb{R}$  is  $\{L(\mathbf{x}), U(\mathbf{x})\}$ , with

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

equality?

- similarly upper and lower confidence bounds
- a  $(1 - \alpha)$ -level confidence **region** for  $\theta \in \Theta$ , is a set  $R(\mathbf{X}) \subset \Theta$ , with

$$\text{pr}\{\theta \in R(\mathbf{X})\} \geq 1 - \alpha$$

- a **pivotal quantity**  $h(\mathbf{x}; \theta)$  is a function of  $\mathbf{X}$  and  $\theta$  with a known distribution:  
we can find  $a, b$  s.t.

$$\text{pr}\{a \leq h(\mathbf{X}; \theta) \leq b\} \stackrel{\geq}{=} 1 - \alpha$$

continuity  
if  $X \sim \text{cont}^s$

- inversion of this gives a  $(1 - \alpha)$  confidence region for  $\theta$

$$X \sim \text{Binom}(n, \theta), \quad 0 \leq \theta \leq 1, \quad \hat{\theta}(x) = x/n$$

$$\hat{\theta} \sim N\{\theta, (1-\theta)\theta/n\}$$

approximate  $1 - \alpha$  CI version 1

$$\hat{\theta} \pm 1.96 \sqrt{\hat{\theta}(1-\hat{\theta})/n}$$

↖ se  
↙ 95%

Wald

approximate  $1 - \alpha$  CI version 2

$$P_n \left\{ |\hat{\theta} - \theta| \leq z_{\alpha/2} \left\{ \frac{(1-\theta)\theta}{n} \right\}^{1/2} \right\} = 1 - \alpha$$

Wilson

guaranteed  $1 - \alpha$  CI

$$-z \sqrt{\frac{\theta(1-\theta)}{n}} \leq (\hat{\theta} - \theta) \leq z \sqrt{(1-\theta)\theta/n}$$

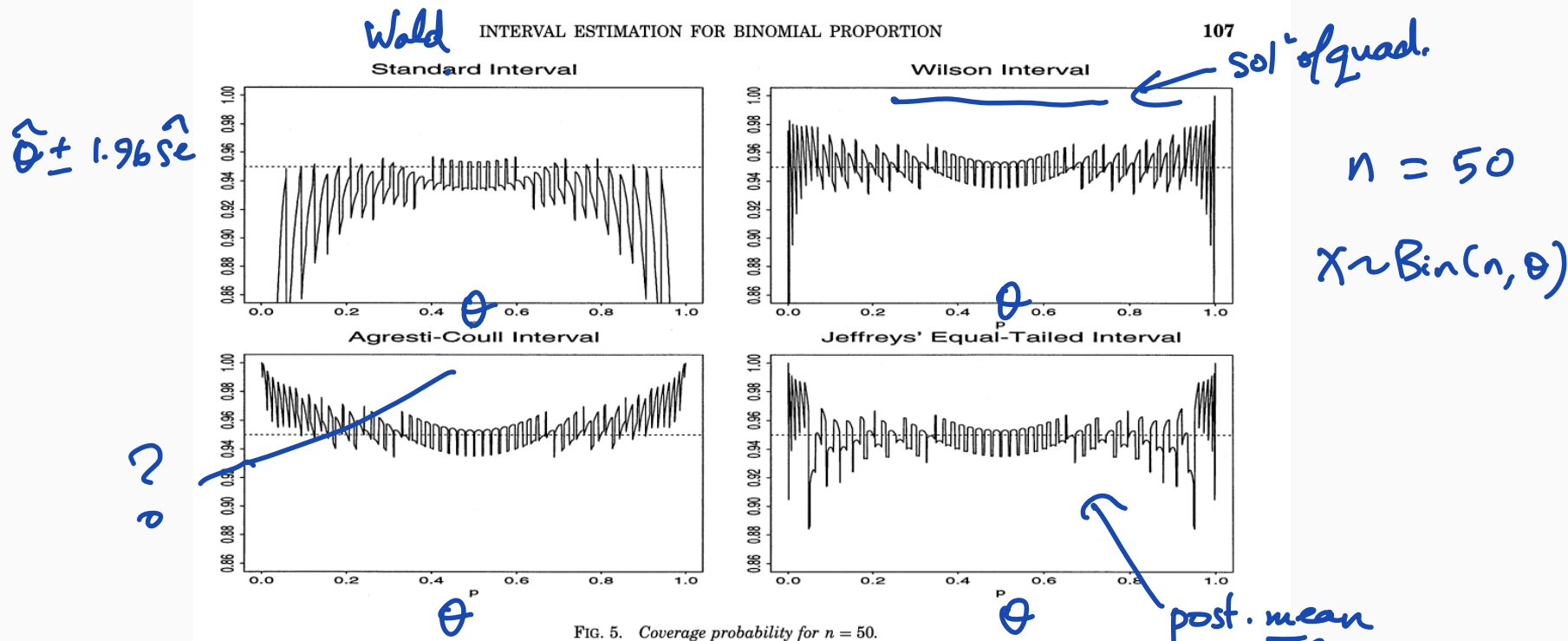
nonparametric

Hoeffding's inequality:

$$(\hat{\theta} - \epsilon_n, \hat{\theta} + \epsilon_n), \quad \epsilon_n^2 = \log(2/\alpha)/(2n)$$

$$\text{pr}(\theta \in \hat{\theta} - \epsilon_n, \hat{\theta} + \epsilon_n) \geq 1 - \alpha$$

no normal approx



<sup>2</sup>Brown, L.D., Cai, T., DasGupta, A. (2001). Interval Estimation for a Binomial Proportion, *Statistical Science* **16**, 101–133. [link](#)

- recall  $X_1, \dots, X_n$ , i.i.d.  $F(\cdot)$

plug-in est'rs

- empirical cdf

$$\hat{F}_n(t) = \frac{1}{n} \sum_{i=1}^n 1\{X_{(i)} \leq t\}$$

- properties:

$$\mathbb{E}\{\hat{F}_n(t)\} = F(t), \quad \text{var}\{\hat{F}_n(t)\} = \frac{1}{n}F(t)\{1 - F(t)\}$$

any fixed  $t$

- pointwise approximate confidence limits  $\hat{F}_n(t) \pm z_{1-\alpha/2}[\hat{F}_n(t)\{1 - \hat{F}_n(t)\}]^{1/2}$

- recall  $X_1, \dots, X_n$ , i.i.d.  $F(\cdot)$

plug-in est'rs

- empirical cdf

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

- properties:

$$\mathbb{E}\{\hat{F}_n(t)\} = F(t), \quad \text{var}\{\hat{F}_n(t)\} = \frac{1}{n}F(t)\{1 - F(t)\}$$

any fixed  $t$

- pointwise approximate confidence limits  $\hat{F}_n(t) \pm z_{1-\alpha/2}[\hat{F}_n(t)\{1 - \hat{F}_n(t)\}]^{1/2}$

- simultaneous confidence band** :  $\text{pr}\{L(t) \leq F(t) \leq U(t) \text{ for all } t\} \geq 1 - \alpha$ :

$$L(t) = \max\{\hat{F}_n(t) - \epsilon_n, 0\}, \quad U(t) = \min\{\hat{F}_n(t) + \epsilon_n, 1\}, \quad \epsilon_n = \left\{ \frac{1}{2n} \log \left( \frac{2}{\alpha} \right) \right\}^{1/2}$$

98 7. Estimating the CDF and Statistical Functionals

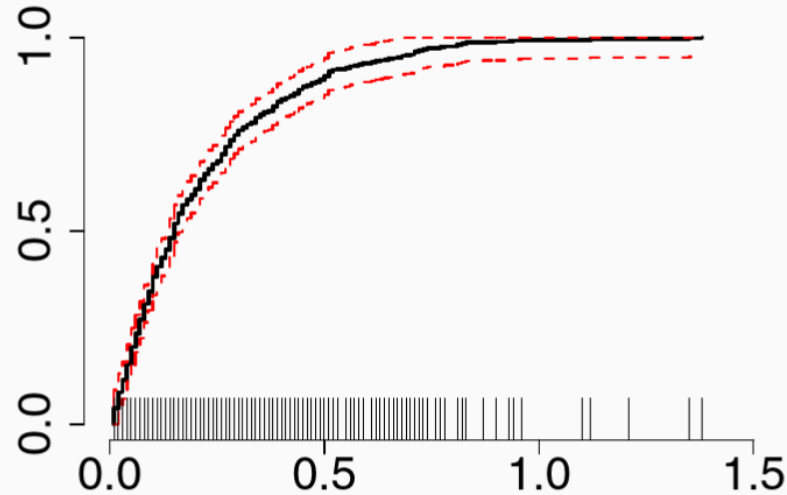


FIGURE 7.1. Nerve data. Each vertical line represents one data point. The solid line is the empirical distribution function. The lines above and below the middle line are a 95 percent confidence band.

**7.2 Example (Nerve Data).** Cox and Lewis (1966) reported 799 waiting times between successive pulses along a nerve fiber. Figure 7.1 shows the empirical CDF  $\hat{F}_n$ . The data points are shown as small vertical lines at the bottom of the plot. Suppose we want to estimate the fraction of waiting times between .4 and .6 seconds. The estimate is  $\hat{F}_n(.6) - \hat{F}_n(.4) = .93 - .84 = .09$ . ■

- $X_1, \dots, X_n$  i.i.d.,  $X_i \sim f(\cdot)$
- kernel density estimator

$$\hat{f}(x) = \frac{1}{n} \sum_{i=1}^n \frac{1}{h} K\left(\frac{x - X_i}{h}\right)$$

- with a symmetric kernel function, for small  $h$ ,

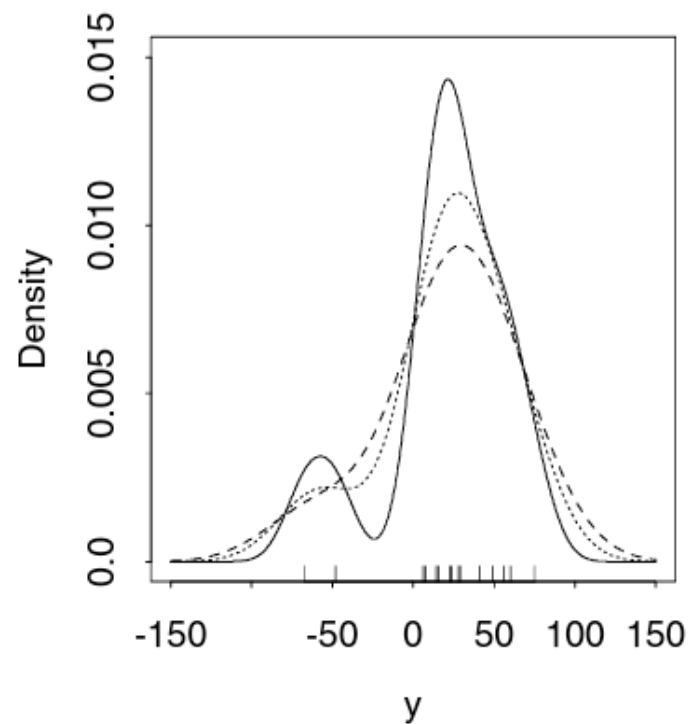
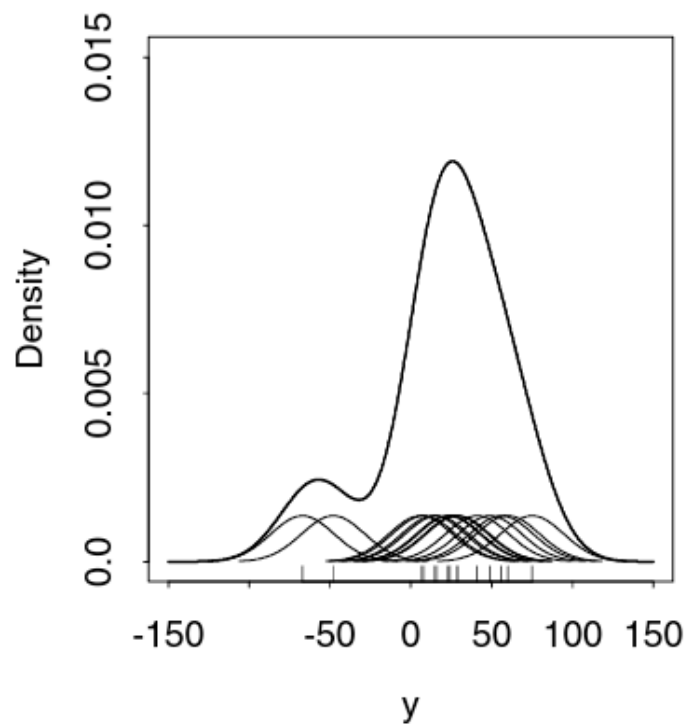
$$\mathbb{E}\{\hat{f}(x)\} = f(x) + \frac{1}{2}h^2f''(x) + O(h^4),$$

$$\text{var}\{\hat{f}(x)\} = \frac{1}{nh}f(x) \int K^2(u)du$$

- mean-squared error

306

7 · Estimation and Hypothesis Testing



**Figure 7.2** Kernel density estimates for maize data. Left: construction of kernel estimate (heavy) as sum of 15 scaled normal densities centred at the  $y_j$ , with  $h = 19.5$ . Right: density estimates with  $h = 13.3$  (solid),  $h = 23.2$  (dots) and  $h = 30$  (dashes).

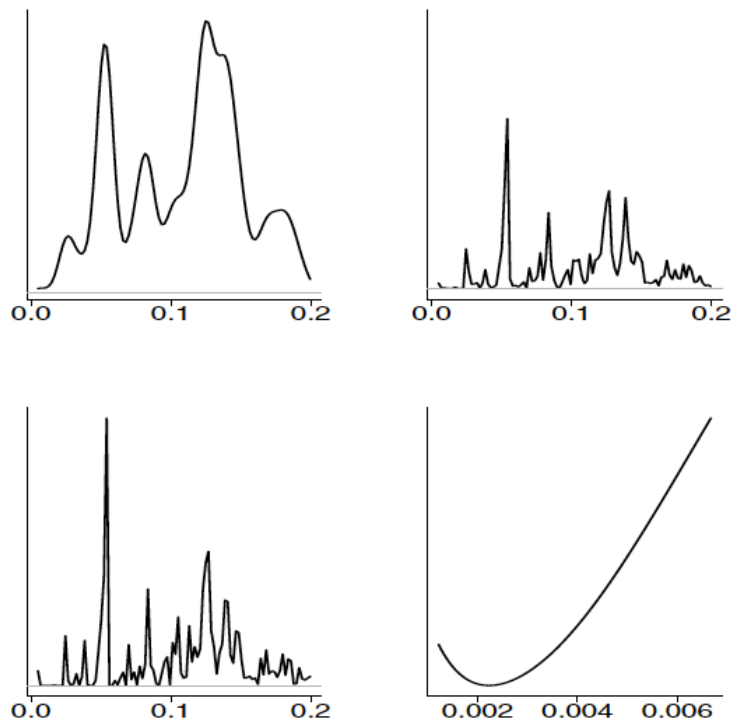


FIGURE 20.6. Kernel density estimators and estimated risk for the astronomy data. Top left: oversmoothed. Top right: just right (bandwidth chosen by cross-validation). Bottom left: undersmoothed. Bottom right: cross-validation curve as a function of bandwidth  $h$ . The bandwidth was chosen to be the value of  $h$  where the curve is a minimum.

## 318 20. Nonparametric Curve Estimation

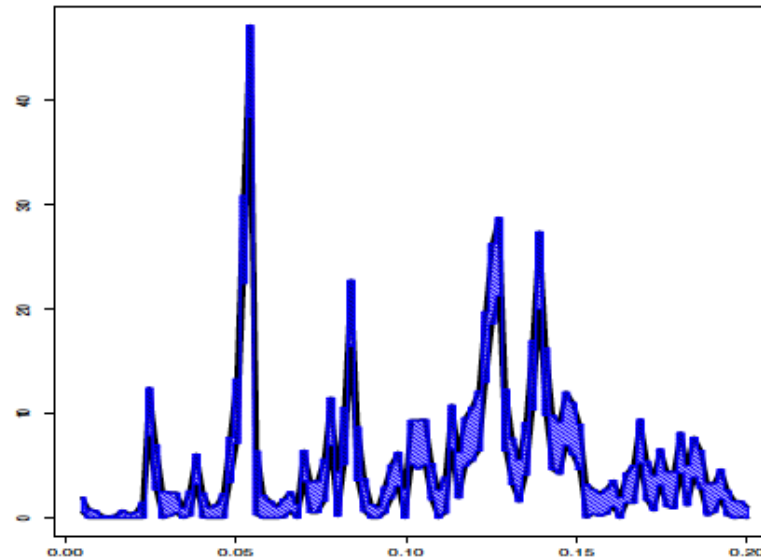
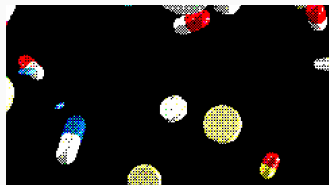


FIGURE 20.7. 95 percent confidence bands for kernel density estimate for the astronomy data.

**20.9 Definition.** A pair of functions  $(\ell_n(x), u_n(x))$  is a  $1 - \alpha$  confidence band (or confidence envelope) if

$$\mathbb{P}\left(\ell(x) \leq \bar{f}_n(x) \leq u(x) \text{ for all } x\right) \geq 1 - \alpha. \quad (20.16)$$



Link

My trail:

“Published in 2025”

[https://www2.unlearn.ai/l/1055293/2025-04-29/4d7f2/1055293/1757437116rfdemBsS/PD\\_C](https://www2.unlearn.ai/l/1055293/2025-04-29/4d7f2/1055293/1757437116rfdemBsS/PD_C)

ASA Substack Dec 2025

<https://asabiopreport.substack.com/p/unleashing-ai-generated-digital-twins>

arxiv paper <https://arxiv.org/pdf/2405.01488>