

# Statistical Theory for Data Science

STA2212H S LEC9101

Week 4

January 27 2026

arXiv:2411.17395v2 [math.ST] 7 Apr 2025

Asymptotics for estimating a diverging number of parameters —  
with and without sparsity

Jana Gauss and Thomas Nagler

*Department of Statistics, LMU Munich  
Munich Center for Machine Learning (MCML)*

April 8, 2025

## Abstract

We consider high-dimensional estimation problems where the number of parameters diverges with the sample size. General conditions are established for consistency, uniqueness, and asymptotic normality in both unpenalized and penalized estimation settings. The conditions are weak and accommodate a broad class of estimation problems, including ones with non-convex and group structured penalties. The wide applicability of the results is illustrated through diverse examples, including generalized linear models, multi-sample inference, and stepwise estimation procedures.

## 1 Introduction

In modern applications, statisticians are facing increasingly complex and high-dimensional problems. Many data sets have a huge number of variables, calling for similarly many parameters  $p$ . In other scenarios, the number of variables is moderate, but adequately modeling the data requires highly complex, non-linear models with many parameters. The traditional fixed- $p$ -large- $n$  paradigm is inadequate in such situations.

This article adopts an asymptotic perspective, allowing both the sample size  $n$  and the number of parameters  $p_n$  to diverge. We consider general parametric problems where the estimator  $\hat{\theta}$  solves an estimating equation

$$\frac{1}{n} \sum_{i=1}^n \phi(\mathbf{X}_i; \hat{\theta}) = \mathbf{0} \in \mathbb{R}^{p_n}, \quad (1)$$

with some function  $\phi: \mathbb{R}^{p_n} \rightarrow \mathbb{R}^{p_n}$ . A classical example are risk minimization problems, where  $\phi$  is the gradient of a loss function. The estimating equation framework is also suited for more complex

# Asymptotics for estimating a diverging number of parameters — with and without sparsity

Jana Gauss and Thomas Nagler

*Department of Statistics, LMU Munich  
Munich Center for Machine Learning (MCML)*

April 8, 2025

## Abstract

We consider high-dimensional estimation problems where the number of parameters diverges with the sample size. General conditions are established for consistency, uniqueness, and asymptotic normality in both unpenalized and penalized estimation settings. The conditions are weak and accommodate a broad class of estimation problems, including ones with non-convex and group structured penalties. The wide applicability of the results is illustrated through diverse examples, including generalized linear models, multi-sample inference, and stepwise estimation procedures.

## 1 Introduction

In modern applications, statisticians are facing increasingly complex and high-dimensional

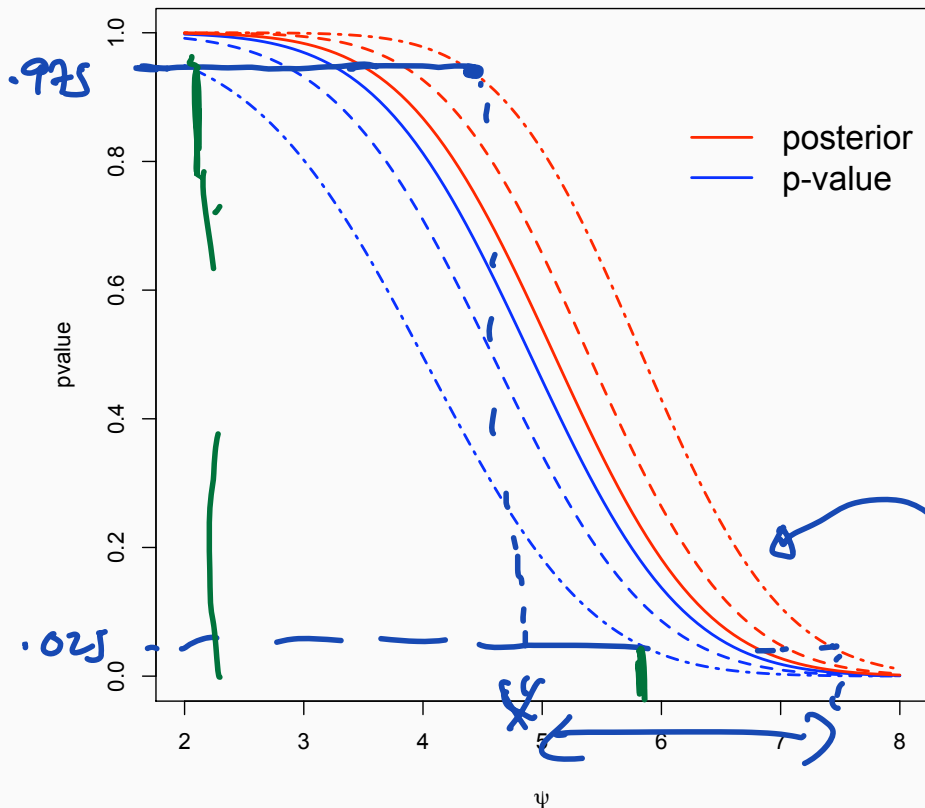
SEMIPARAMETRIC EFFICIENT EMPIRICAL HIGHER ORDER INFLUENCE  
FUNCTION ESTIMATORSBY LIN LIU<sup>\*</sup>, RAJARSHI MUKHERJEE<sup>†</sup>, WHITNEY K. NEWEY<sup>‡</sup>, JAMES M. ROBINS<sup>§</sup>

Robins et al. (2008) applied the theory of higher order influence functions (HOIFs) to derive an estimator of the mean  $\psi$  of an outcome  $Y$  in a missing data model with  $Y$  missing at random conditional on a vector  $X$  of continuous covariates; their estimator, in contrast to other existing estimators but ours, is semiparametric efficient under the minimal Hölder smoothness conditions derived in Robins et al. (2009b), together with an additional (non-minimal) Hölder smoothness condition on the density  $g$  of  $X$ , because that particular estimator depends on a non-parametric estimate of  $g$ . In this paper, we introduce a new HOIF estimator that has the same asymptotic properties as the previous one, but imposes no smoothness requirement on  $g$ . This improvement is significant for two reasons. First, one rarely has the knowledge about the smoothness properties of  $g$ . Second, even when  $g$  is smooth, and even if  $X$  is just multivariate with fixed dimensions, accurate nonparametric estimation of its density is generally not feasible at the sample sizes often encountered in practice. In fact, to our knowledge, this new HOIF estimator to be studied here remains the *only* semiparametric efficient estimator of  $\psi$  under minimal Hölder smoothness conditions, despite the rapidly growing literature on causal effect estimation. We also show that our estimator can be generalized to the entire class of functionals considered by Robins et al. (2008) which includes the average effect of a treatment on a response  $Y$  when a vector  $X$  suffices to control for confounding and the expected conditional variance of  $Y$  given  $X$ . Simulation experiments are also conducted, which demonstrate that our new estimator outperforms previous ones proposed in earlier works on HOIFs in finite samples, when  $g$  is not very smooth.

**1. Introduction.** Robins et al. (2008), together with a companion technical report Robins et al. (2016) containing more details, introduced novel U-statistic based estimators of a class of nonlinear functionals in semi- and non-parametric models. Construction of these estimators was based on the theory of Higher Order Influence Functions (henceforth referred to as HOIFs) (Robins et al., 2008). HOIFs are U-statistics that represent higher order derivatives of a functional. The

1. Midterm 2 shift from March 3 to March 10? ✓
2. Recap: Stein's problem (notes corrected), sandwich variance,  $M$ -estimators (recent arxiv), interval estimation
3. Nonparametric estimation
4. Hypothesis and significance testing
5. Questions/Review re Midterm Feb 3 7.10 – 8.00 pm
6. Papers re project

Normal Circle,  $k=2, 5, 10$



$$X_i \sim N(\mu_i, 1) \quad \text{ind't} \\ i = 1, \dots, k$$

$$\pi(\underline{\mu}) = 1 \quad \psi = \sum_{i=1}^k \mu_i^2$$

$$E(\psi | \underline{x}) \text{ post.mean} = k + \sum x_i^2$$

$$E_{\text{model}}(E(\psi | \underline{x})) = k + k + \psi$$

$$\psi | \underline{x} \sim \chi_k^2(\sum x_i^2) \quad \text{exact}$$

$$P_n(\Psi \geq \psi | \underline{x}) \quad \text{Bayes surv. f.}$$

$$p\text{-value function} \\ p(\psi) = P_n \left\{ \sum X_i^2 \leq \sum x_i^2 \right\} \\ \underline{x} | \underline{\mu} \quad \text{ntbc.}$$

- model assumption  $X_1, \dots, X_n$  i.i.d.  $f(x; \theta), \theta \in \Theta$
- maximum likelihood estimator based on model:  $\sum_{i=1}^n \ell'(\hat{\theta}_n; X_i) = 0$
- 

$$\sqrt{n}\{\hat{\theta}_n - \theta(F)\} \xrightarrow{d} N(0, \sigma^2), \quad \sigma^2 = \frac{\int [\ell'\{\theta(F); x\}]^2 dF(x)}{(\int [\ell''\{\theta(F); x\}]^2 dF(x))^2} = \frac{\text{var}_F\{\ell'(\theta; X)\}}{E_F\{-\ell''(\theta; X)\}^2}$$

true  
 $X_i \sim F(\cdot)$

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

needs to be est'd

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

$X_i \sim F(\cdot)$

• maximum likelihood estimator based on model:  $\sum_{i=1}^n \ell'(\hat{\theta}_n; X_i) = 0$

$$\sqrt{n}\{\hat{\theta}_n - \theta(F)\} \xrightarrow{d} N(0, \sigma^2), \quad \sigma^2 = \frac{\int [\ell'\{\theta(F); x\}]^2 dF(x)}{(\int [\ell''\{\theta(F); x\}]^2 dF(x))^2} = \frac{\text{var}_{F_n}\{\ell'(\theta; X)\} \frac{1}{n} \sum \ell'(\theta; x_i)^2}{E_n\{-\ell''(\theta; X)\}^2}$$

$\text{var}(\ell') \triangleq \int \ell'(\theta; x)^2 dF_n(x) \quad \frac{1}{n} @ \text{ each } x_i$

$\left\{ \frac{1}{n} \sum_{i=1}^n -\ell''(\theta; x_i) \right\}_{\tilde{\theta}_n}$

→ • model is correct; estimator is root of  $\sum_{i=1}^n \psi(\theta; x_i) = 0$

• assume  $\psi(\theta; X_i)$  has expected value 0, finite variance, and is differentiable w.r.t  $\theta$

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

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

...cond's...

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

- upper, lower, central

$$C_n(x) = (-\infty, a_n] \text{ or } [b_n, \infty)$$

$$\text{or } [a_n(x), b_n(x)]$$

- length and coverage

↑  
 $1-\alpha$

$$\frac{b_n(x) - a_n(x)}{\text{or } E\{b_n(x) - a_n(x)\}} \rightarrow P_{\theta} \{ a_n(x) \leq \theta \leq b_n(x) \} \geq 1-\alpha \quad \forall \theta$$

- exact or approximate

- inversion of pivotal quantity

$S(x; \theta)$  with known dist: " $= 1-\alpha$  same"

$$G_S(s)$$

- posterior credible intervals

$$\int_a^b \pi(\theta|x) d\theta = 1-\alpha$$

e.g.  $\int_{-\infty}^a \pi(\theta|x) d\theta = \alpha/2$

$$P_n \{ S(x; \theta) \leq s^\alpha \} = \alpha$$

↑  
find  $s^\alpha$

invert that =  $G_S^{-1}(\alpha)$   
to get a bd for  $\theta$

# Confidence regions; HPD regions

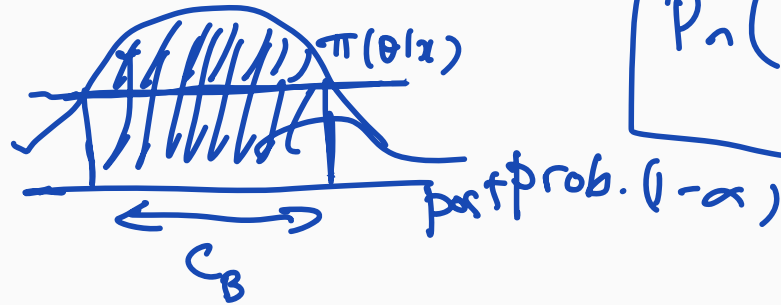
$$\theta \in \mathbb{R}^p \quad \mathbb{R}^k \quad \mathbb{R}^d$$

$$R(\underline{x}) \quad C_n(\underline{x}) = \{ \underline{\theta} : P_{n, \underline{\theta}} \{ R(\underline{x}) \ni \underline{\theta} \} \geq 1 - \alpha$$

gen'n  
Conf. int.

$$\rightarrow C_B(\underline{x}) = \{ \underline{\theta} : \pi(\underline{\theta} | \underline{x}) \geq k_\alpha \}$$

credible region  
for  $\underline{\theta}$



$$k_\alpha \text{ chosen s.t.} \\ P_n(\underline{\theta}_n \in C_B | \underline{x}) = 1 - \alpha$$

approx.  $C_n(\underline{x}) \quad (\hat{\underline{\theta}} - \underline{\theta})^T I_n(\underline{\theta})(\hat{\underline{\theta}} - \underline{\theta}) \leq \chi_p^2(1-\alpha)$

$C_B \quad \{ \underline{\theta} - \hat{\underline{\theta}}_n \}^T J_\pi(\underline{\theta})(\underline{\theta} - \hat{\underline{\theta}}_n) \leq \chi_p^2(1-\alpha) \quad \text{e.g.}$

• model  $Y \sim f(y; \underline{\psi}, \underline{\lambda})$ ,  $\psi \in \mathbb{R}^d, \lambda \in \mathbb{R}^{p-d}, \theta = (\psi, \lambda)$   $y = (y_1, \dots, y_n)$

• profile log-likelihood function  $\underline{l_p(\psi)} = l(\psi, \hat{\lambda}_\psi)$   $\dot{j}_p(\psi) = -l''_{\psi\psi}(\psi)$  maximize over  $\lambda$

• approximate pivotal quantities MS Thm 7.4,5

$\{ \psi : (\hat{\psi} - \psi)^T \dot{j}_p(\hat{\psi})(\hat{\psi} - \psi) \leq \chi^2_{d, 1-\alpha} \}$   $\rightarrow$  c.R. for  $\psi$

$\hat{\psi} \pm 1.96 \dot{j}_p^{-1/2}(\hat{\psi})$   $(\hat{\psi} - \psi)^T \dot{j}_p(\hat{\psi})(\hat{\psi} - \psi) \sim \chi^2_p$  under  $f(y; \psi, \lambda)$

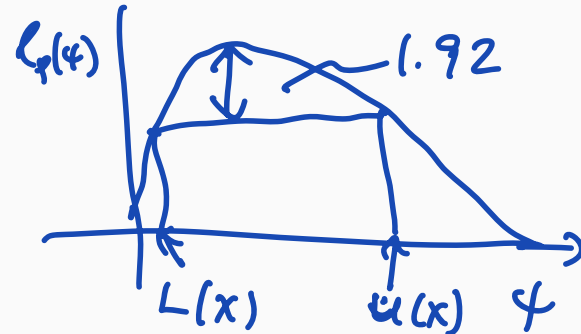
$2\{l_p(\hat{\psi}) - l_p(\psi)\} \sim \chi^2_p$

• approximate confidence regions



$\{ \psi : (\hat{\psi} - \psi)^T \dot{j}_p(\hat{\psi})(\hat{\psi} - \psi) \leq \chi^2_{d, 1-\alpha} \}$

$\{ \psi : 2\{l_p(\hat{\psi}) - l_p(\psi)\} \leq \chi^2_{d, 1-\alpha} \}$



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

plug-in est'rs

$$X_{(1)} \leq \dots \leq X_{(n)}$$

- empirical cdf

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

- properties:

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

Binomial

any fixed  $x$

- pointwise approximate confidence limits

$$\underline{\hat{F}_n(x) \pm z_{1-\alpha/2} [\hat{F}_n(x)\{1 - \hat{F}_n(x)\}]^{1/2}}$$

for any fixed  $x$   
has coverage  $\approx 1 - \alpha$

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

$\{F(x); x \in \mathcal{X}\}$  stoch. proc.

Glivenko-Cantelli Theorem:

$$\sup_x |\hat{F}_n(x) - F(x)| \xrightarrow{p} 0$$

$\{\hat{F}_n(x); x \in \mathcal{X}\}$   
empirical process

Dvoretzky-Kiefer-Wolfowitz Inequality

$$\Pr\{\sup_x |\hat{F}_n(x) - F(x)| > \epsilon\} \leq 2 \exp(-2n\epsilon^2)$$

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

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

98 7. Estimating the CDF and Statistical Functionals

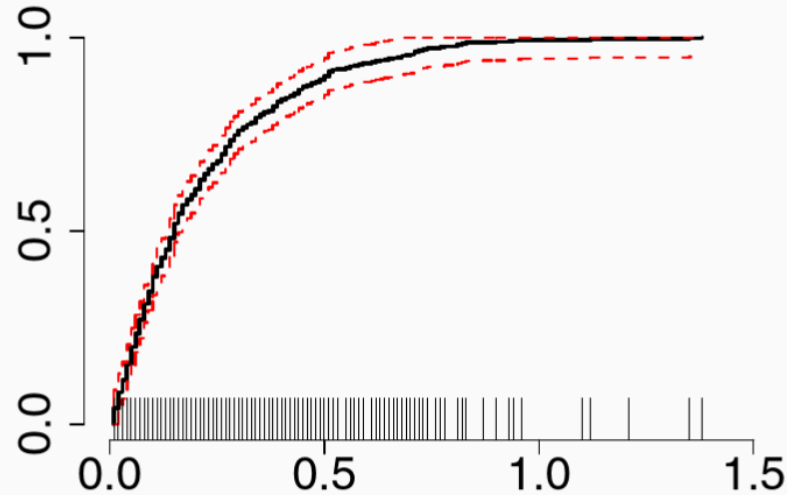


FIGURE 7.1. Nerve data. Each vertical line represents one data point. The solid line is the empirical distribution function. The lines above and below the middle line are a 95 percent confidence band.

**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  $\hat{F}_n(.6) - \hat{F}_n(.4) = .93 - .84 = .09$ . ■

$$\hat{F}_n(x) = \text{e.c.d.f.}$$

$$\hat{F}'_n(x) = \frac{1}{n} \text{ at each } x_i \quad i=1, \dots, n$$

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

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

$k(\cdot)$  kernel  $f$

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

↑ scale

e.g.  $N(0, 1)$

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

biased est.

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

fixed  $x$

- mean-squared error

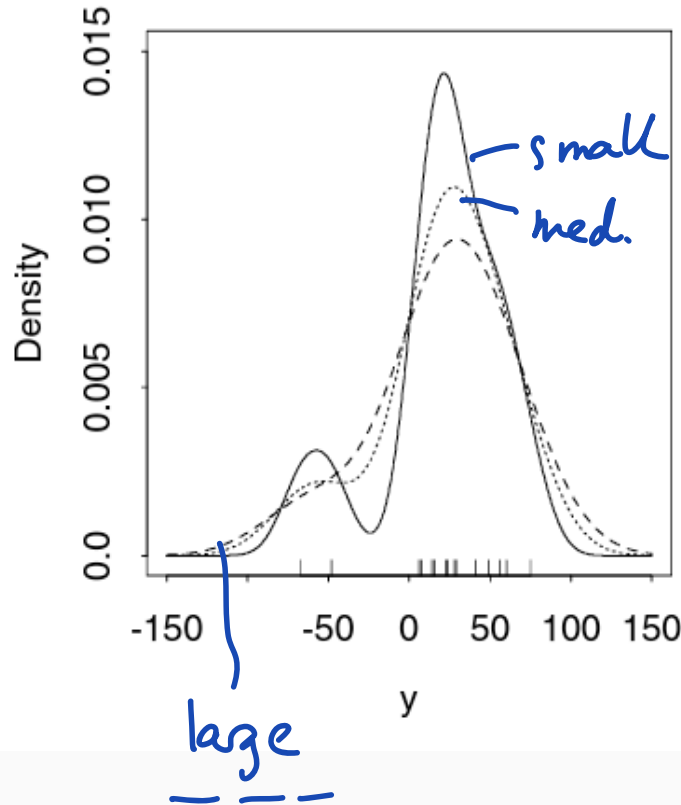
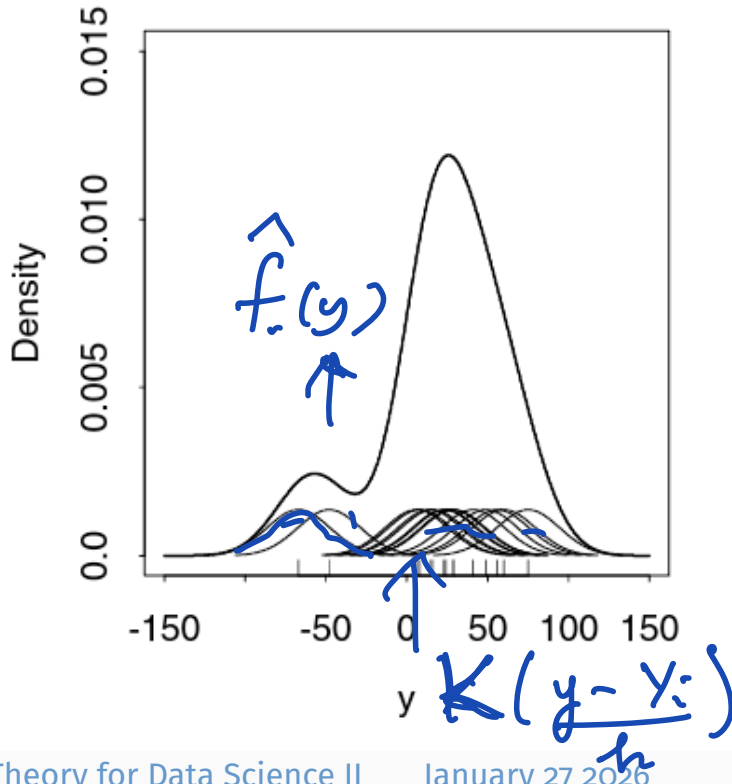
$$\text{MSE} = \frac{h^4}{4} (f''(x))^2 + \frac{1}{nh} f(x) \int K^2(u) du \left( + O(h) ? \right)$$

$$\frac{\partial}{\partial h} \dots = 0 \Rightarrow h^* = h^*(x)$$

$$h \propto n^{-1/5} \text{ gives } \text{MSE} \sim O(n^{-4/5})$$

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).

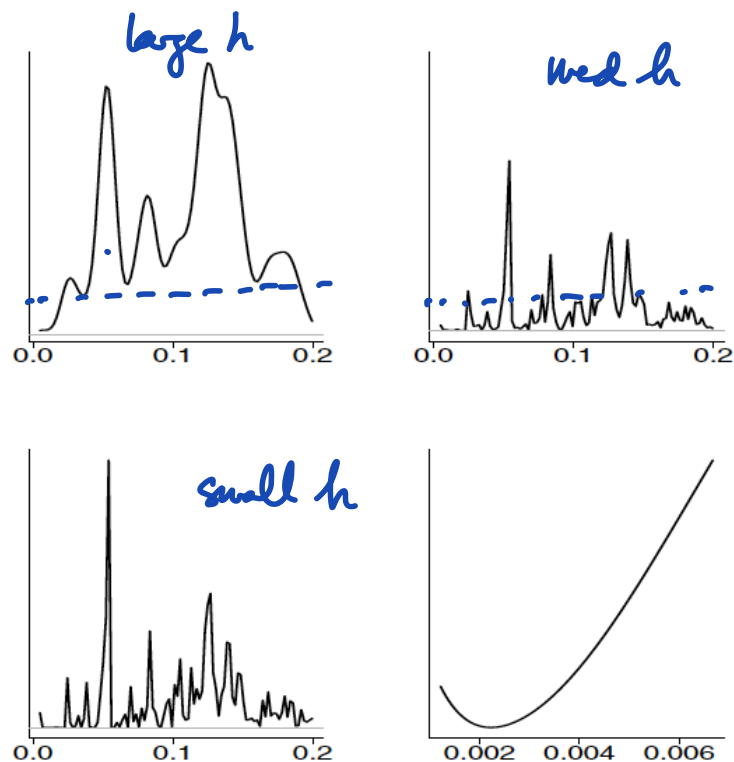


FIGURE 20.6. Kernel density estimators and estimated risk for the astronomy data. Top left: oversmoothed. Top right: just right (bandwidth chosen by cross-validation). Bottom left: undersmoothed. Bottom right: cross-validation curve as a function of bandwidth  $h$ . The bandwidth was chosen to be the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

## 318 20. Nonparametric Curve Estimation

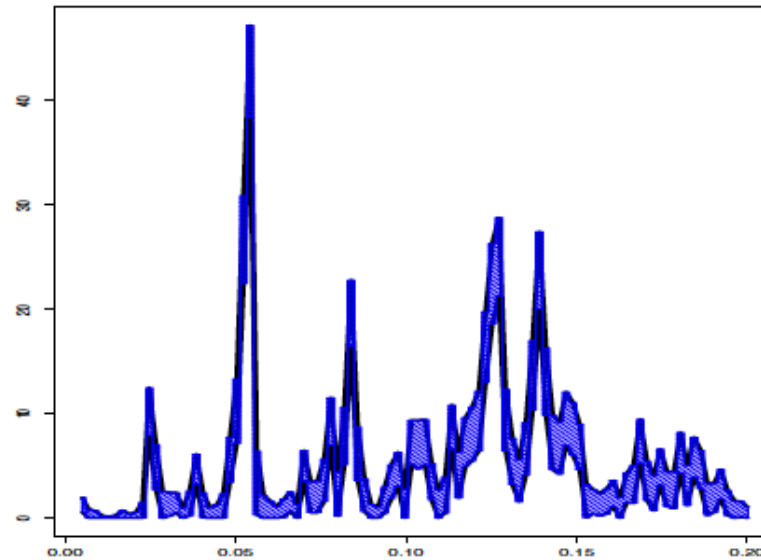


FIGURE 20.7. 95 percent confidence bands for kernel density estimate for the astronomy data.

$$\frac{\hat{\theta} \pm 1.96 \hat{SE}}{f \theta}$$

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

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

$\bar{f}_n(x)$  is a smoothed version of the density estimator <sup>and</sup>

re-defines what we're est'g

$$f(x) \rightarrow \int_{x-\varepsilon}^{x+\varepsilon} f(z) p(z) dz = \bar{f}$$

$K(\cdot)$

To construct confidence bands, we use something similar to histograms. Again, the confidence band is for the smoothed version,

$$\bar{f}_n = \mathbb{E}(\hat{f}_n(x)) = \int \frac{1}{h} K\left(\frac{x-u}{h}\right) f(u) du,$$

of the true density  $f$ .<sup>4</sup> Assume the density is on an interval  $(a, b)$ . The band is

$$\ell_n(x) = \hat{f}_n(x) - q \widehat{\text{se}}(x), \quad u_n(x) = \hat{f}_n(x) + q \widehat{\text{se}}(x) \quad (20.27)$$

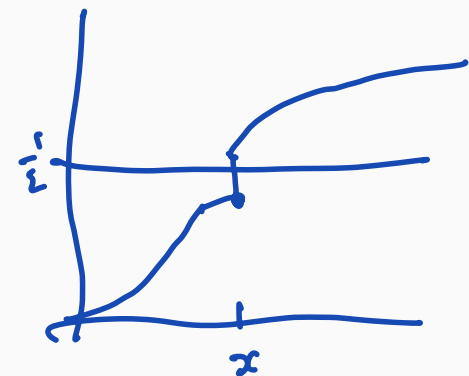
where

$$\begin{aligned} \widehat{\text{se}}(x) &= \frac{s(x)}{\sqrt{n}}, \\ s^2(x) &= \frac{1}{n-1} \sum_{i=1}^n (Y_i(x) - \bar{Y}_n(x))^2, \\ Y_i(x) &= \frac{1}{h} K\left(\frac{x-X_i}{h}\right), \\ \underline{q} &= \Phi^{-1}\left(\frac{1 + (1-\alpha)^{1/m}}{2}\right), \\ m &= \frac{b-a}{\omega} \end{aligned}$$

where  $\omega$  is the width of the kernel. In case the kernel does not have finite width then we take  $\omega$  to be the effective width, that is, the range over which the kernel is non-negligible. In particular, we take  $\omega = 3h$  for the Normal

- Define an **estimand**  $\theta = T(F)$  or  $\theta(F)$   $T: \mathcal{F} \rightarrow \Theta$
- Substitution (plug-in) estimator  $\hat{\theta}_n = T(\hat{F}_n)$   $\hat{F}_n$  e.c.d.f.
- Examples of functionals: moments, quantiles, solution of  $\int \psi\{x; T(F)\} dF(x) = 0$

$$E(x^3) = \int x^3 dF(x) = T_3(F) \quad \frac{1}{n} \sum X_i^3 = T_3(\hat{F}_n)$$



$$T_m = \text{median}(F) = F^{-1}\left(\frac{1}{2}\right) = \inf \left\{ x : F(x) \geq \frac{1}{2} \right\}$$

$$\Rightarrow T_m(\hat{F}_n) = \hat{F}_n^{-1}\left(\frac{1}{2}\right)$$

$$\int \psi(x, T(\hat{F}_n)) d\hat{F}_n(x) = 0 \quad \Rightarrow \quad \frac{1}{n} \sum \psi(x_i, T(\hat{F}_n)) = 0$$

- Define an **estimand**  $\theta = T(F)$  or  $\theta(F)$   $T : \mathcal{F} \rightarrow \Theta$
- Substitution (plug-in) estimator  $\hat{\theta}_n = T(F_n)$
- Examples of functionals: moments, quantiles, solution of  $\int \psi\{x; T(F)\}dF(x) = 0$
- Lorenz curve; Gini index MS Ex 4.21;  $F$  on  $[0, \infty)$

$$q_F(t) = \frac{\int_0^t \hat{F}^{-1}(s) ds}{\int_0^1 \hat{F}^{-1}(s) ds}, \quad \theta(F) = 2 \int_0^1 \{t - \hat{q}_F(t)\} dt$$

- Define an **estimand**  $\theta = T(F)$  or  $\theta(F)$   $T : \mathcal{F} \rightarrow \Theta$
- Substitution (plug-in) estimator  $\hat{\theta}_n = T(F_n)$
- Examples of functionals: moments, quantiles, solution of  $\int \psi\{x; T(F)\}dF(x) = 0$
- Lorenz curve; Gini index MS Ex 4.21;  $F$  on  $[0, \infty)$

$$q_F(t) = \frac{\int_0^t F^{-1}(s)ds}{\int_0^1 F^{-1}(s)ds}, \quad \theta(F) = 2 \int_0^1 \{t - q_F(t)\}dt$$

- **Estimator:**  $\hat{\theta}_n$  solves  $\frac{1}{n} \sum_{i=1}^n \psi(X_i; \theta) = 0$   $= \int \psi\{x; T(\hat{F}_n)\}d\hat{F}_n(x)$
- $F^{-1}(t) = \inf\{x : F(x) \geq t\}$ ,  $\hat{\theta}_n = 2 \int_0^1 \{t - q_{\hat{F}_n}(t)\}dt$

- Define an **estimand**  $\theta = T(F)$ , **estimator**  $\hat{\theta}_n = T(\hat{F}_n)$
- Influence function

$$T : \mathcal{F} \rightarrow \Theta$$

if limit exists

(efficient " ")

$$\phi(x; F) = \frac{d}{dt} \theta\{(1-t)F + t\delta(x)\} \Big|_{t=0}$$

$$\delta(x) = \begin{cases} 0 & x \leq 0 \\ 1 & x > 0 \end{cases}$$

- If we have  $\hat{\theta}_n - \theta(F) = \frac{1}{n} \sum_{i=1}^n \phi(X_i; F) + R_n$ ,  $\sqrt{n}R_n \xrightarrow{p} 0$ , then

$$E\{\phi(X_i; F)\} = 0$$

estimator  $\theta(\hat{F}_n)$

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

$$\sigma^2(F) = \int \phi^2(x; F) dF(x)$$

MS Ex.4.35

$$\left\{ \sqrt{n} \{ \hat{F}_n(t) - F(t) \} ; t \in \mathcal{R} \right\}$$

Classic reference:

Van der Vaart & Wellner (1996). *Weak Convergence and Empirical Processes*, Springer.

# Hypothesis and significance testing

- model  $Y \sim f(y; \psi, \lambda)$ ,  $\psi \in \mathbb{R}^d, \lambda \in \mathbb{R}^{p-d}, \theta = (\psi, \lambda)$

$$y = (y_1, \dots, y_n)$$

- approximate pivotal quantities

MS Thm 7.4,5

likelihood-based

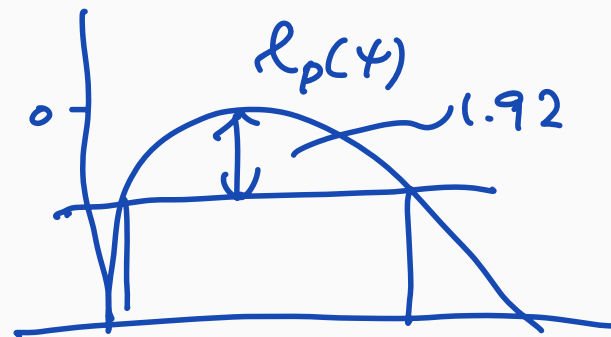
$$\begin{aligned} & (\hat{\psi} - \psi)^T j_p(\hat{\psi})(\hat{\psi} - \psi) \\ & 2\{\ell_p(\hat{\psi}) - \ell_p(\psi)\} \end{aligned}$$



- approximate confidence regions

$$\begin{aligned} & \rightarrow \{\psi : (\hat{\psi} - \psi)^T j_p(\hat{\psi})(\hat{\psi} - \psi) \leq \chi_{d,1-\alpha}^2\} \\ & \rightarrow \{\psi : 2\{\ell_p(\hat{\psi}) - \ell_p(\psi)\} \leq \chi_{d,1-\alpha}^2\} \end{aligned}$$

- generalized likelihood ratio test *pivot ? check?*



- Textbook version

$$\Lambda = \frac{\sup_{\theta \in \Theta} f(\mathbf{X}; \theta)}{\sup_{\theta \in \Theta_0} f(\mathbf{X}; \theta)} = \frac{\sup_{\theta \in \Theta} L(\theta)}{\sup_{\theta \in \Theta_0} L(\theta)}$$

Test of  $H_0: \theta \in \Theta_0 \subset \Theta$

- Theorem 7.5  
(MS)

$$2 \log \Lambda_n \xrightarrow{d} V \sim \chi_d^2$$

$$d = \dim \Theta_0$$

- constraint  $\theta \in \Theta_0$  has  $d$  restrictions on  $\theta = (\theta_1, \dots, \theta_p)$

- version using profile constrains  $\psi = (\theta_1, \dots, \theta_d)$ ,  $d < p$

$$H_0: \theta_1, \dots, \theta_d = \psi_1, \dots, \psi_d \text{ fixed}$$

$\rightarrow$  free to vary  
 $\Theta_0 \subset \Theta$

- motivated as a test of  $H_0: \theta \in \Theta_0$  vs  $H_A: \theta \notin \Theta_0$

1 value of  $\psi_0$



$$X_1, \dots, X_n \sim f(\mathbf{x}; \theta)$$

AoS range  $\mathbf{x} \in \mathcal{X}$

- **Null** and **alternative** hypothesis:  $H_0 : \theta \in \Theta_0, H_1 : \theta \in \Theta_1, \Theta_0 \cup \Theta_1 = \Theta$

- Rejection region:  $R \subset \mathcal{X}$ ; if  $\mathbf{x} \in R$  “reject”  $H_0$

AoS 10.0

$$R = \{\mathbf{x} : \phi(\mathbf{x}) = 1\}$$

- Test statistic and critical value:  $R = \{\mathbf{x} \in \mathcal{X} : t(\mathbf{x}) > \underline{c}\}$

$\underline{c}$  to be chosen

- Test (decision) function:  $\phi : \mathcal{X} \rightarrow \{0, 1\}$

MS 7.3

$$\phi(\mathbf{X}) = 1 \text{ decide } \theta \in \Theta_1, \text{ else decide } \theta \in \Theta_0$$

- **Type I** and **Type II** error:  $\text{pr}\{\mathbf{X} \in R\}, \theta \in \Theta_0,$   $\text{pr}\{\mathbf{X} \notin R\}, \theta \in \Theta_1$   
 $\text{pr}\{t(\mathbf{X}) \geq c\},$   $\text{pr}\{t(\mathbf{X}) < c\}$

- goal is to identify  $R$ ,  $\phi(\cdot)$ , or  $T = t(\mathbf{X})$  with small Type I and Type II errors

- goal is to identify  $R$ ,  $\phi(\cdot)$ , or  $T = t(\mathbf{X})$  with small Type I and Type II errors
- can't reduce both errors at once see text following Ex. 7.10
- classical solution: require **Type 1 error**  $\leq \alpha$  size  $\alpha$
- subject to this constraint, minimize **type 2 error**

- goal is to identify  $R$ ,  $\phi(\cdot)$ , or  $T = t(\mathbf{X})$  with small Type I and Type II errors

- can't reduce both errors at once

see text following Ex. 7.10

- classical solution: require **Type 1 error**  $\leq \alpha$

size  $\alpha$

- subject to this constraint, minimize **type 2 error**

- find a **test statistic**,  $T = t(\mathbf{X})$ , and define  $R = \{\mathbf{x} : t(\mathbf{x}) \geq c\}$  s.t.

$$\underbrace{\text{pr}_{\theta \in \Theta_0} \{t(\mathbf{X}) \geq c\} \leq \alpha,}_{\text{then}} \underbrace{\max_{\theta \in \Theta_1} \text{pr}_{\theta \in \Theta_1} \{t(\mathbf{X}) \geq c\}}$$

- power function

$$\beta(\theta) = \text{pr}_{\theta} \{t(X) \geq c\}$$

(implicit ass<sup>n</sup> =  
large  $t(\underline{x})$   
is evidence  
against  $H_0$ )

AoS

assume

$\mathbb{H} = \mathbb{R}$ ;

$\text{pr} \uparrow$  in  $\theta > \theta_0$

- for testing simple  $H_0$  against simple  $H_1$

$$f(x; \theta_0)$$

$$f(x; \theta_1)$$

- test statistic

$$T = \frac{L(\theta_1; \mathbf{x})}{L(\theta_0; \mathbf{x})} = \frac{f(\mathbf{x}; \theta_1)}{f(\mathbf{x}; \theta_0)}$$

$$\left[ \frac{L(\hat{\theta}_1; \mathbf{x})}{L(\hat{\theta}_0; \mathbf{x})} \right]$$

↑  
gen'd

- critical region

$$\{\mathbf{x} : t(\mathbf{x}) \geq k\}$$

- Choose  $k = k_\alpha$  to satisfy

$$\underline{\text{pr}_{H_0}(T \geq k_\alpha) = \alpha}$$

- This test is a most powerful test of  $H_0$  against  $H_1$  at level  $\alpha$

smallest  $T_2$  error  
largest  $1 - T_2$  error

# A neatly-typed proof (from SM 7.3)

Let  $R$  be the rejection region for the test based on

$$R = \{\mathbf{x} : T(\mathbf{x}) \geq k_\alpha\}$$

Let  $R'$  be some other rejection region also of size  $\alpha$

$$\begin{aligned} \alpha &= \int_R f_0(\mathbf{x}) d\mathbf{x} = \int_{R'} f_0(\mathbf{x}) d\mathbf{x} \\ \int_{R-R'} f_0(\mathbf{x}) d\mathbf{x} &= \int_{R'-R} f_0(\mathbf{x}) d\mathbf{x} \end{aligned}$$

$\underbrace{\hspace{10em}}_{\subseteq R} \quad \downarrow \quad \underbrace{\hspace{10em}}_{\subseteq R^c}$

$$\int_{R-R'} f_1(\mathbf{x}) d\mathbf{x} \geq \int_{R'-R} f_1(\mathbf{x}) d\mathbf{x}$$

On LHS  $f_1(\mathbf{x}) \geq k_\alpha f_0(\mathbf{x})$ .

On RHS  $f_1(\mathbf{x}) < k_\alpha f_0(\mathbf{x})$ .

Add integral over intersection  $\underline{\underline{R \cap R'}}$

$$T = f_1(\mathbf{x})/f_0(\mathbf{s})$$

$$\leq \alpha$$

$$\begin{aligned} A - B \\ \equiv A \cap B^c \end{aligned}$$

$$\begin{aligned} R - R' &\subset R \\ R' - R &\subset R^c \end{aligned}$$

## A neatly-typed proof (from MS)

Let  $\phi(\mathbf{x})$  be the test function for the test based on  $T$ .

Let  $\psi(\mathbf{x})$  be any other function that maps  $\mathbf{x}$  to  $[0, 1]$ .

If

$$E_{H_0}\{\psi(\mathbf{X})\} \leq E_{H_0}\{\phi(\mathbf{X})\} = \alpha$$

then it must follow that

$$E_{H_1}\{\psi(\mathbf{X})\} \leq E_{H_1}\{\phi(\mathbf{X})\}$$

Proof:  $\forall \mathbf{x}$ ,

$$\psi(\mathbf{x})\{f_1(\mathbf{x}) - kf_0(\mathbf{x})\} \leq \phi(\mathbf{x})\{f_1(\mathbf{x}) - kf_0(\mathbf{x})\}$$

Integrate and re-arrange terms to get the result

- The formal theory of testing imagines a decision to “reject  $H_0$ ” or not, according as  $X \in R$  or  $X \notin R$ , for some defined region  $R \subset \mathcal{X}$  e.g.  $|Z| > 1.96$
- This is useful for deriving the form of optimal tests, but not useful in practice.
- Doesn't distinguish between  $Z = 1.97$  and  $Z = 19.7$ , for example.
- $P$ -values give more precise information about the null hypothesis

- The formal theory of testing imagines a decision to “reject  $H_0$ ” or not, according as  $X \in R$  or  $X \notin R$ , for some defined region  $R \subset \mathcal{X}$  e.g.  $|Z| > 1.96$
- This is useful for deriving the form of optimal tests, but not useful in practice.
- Doesn't distinguish between  $Z = 1.97$  and  $Z = 19.7$ , for example.
- P-values give more precise information about the null hypothesis

• AoS definition: p-value =  $\inf\{\alpha : T(X_n) \in R_\alpha\}$

$|Z|$  - score 1.97 Def 10.11

• MS definition:  $p(\mathbf{x}) = \inf\{\alpha : \phi_\alpha(\mathbf{x}) = 1\}$

$p < 0.05$  7.5

• SM definition  $p_{obs} = \Pr_{H_0}\{T(X_n) \geq t_{obs}\}$

$p = 0.052$

$t_{obs} = T(\underline{x})$

$|Z| = 4.5$   $p = 10^{-2}$

- **Hypothesis tests** typically means:

- $H_0, H_1$
- critical/rejection region  $R \subset \mathcal{X}$ ,
- level  $\alpha$ , power  $1 - \beta$
- conclusion: “reject  $H_0$  at level  $\alpha$ ” or “do not reject  $H_0$  at level  $\alpha$ ”
- planning: maximize power for some relevant alternative

minimize type II error

# Hypothesis tests and significance tests

- **Hypothesis tests** typically means:

- $H_0, H_1$
- critical/rejection region  $R \subset \mathcal{X}$ ,
- level  $\alpha$ , power  $1 - \beta$
- conclusion: “reject  $H_0$  at level  $\alpha$ ” or “do not reject  $H_0$  at level  $\alpha$ ”
- planning: maximize power for some relevant alternative

minimize type II error

- **Significance tests** typically means:

- $H_0$ ,
- test statistic  $T$
- observed value  $t^{obs}$ ,
- $p$ -value  $p^{obs} = \Pr(T \geq t^{obs}; H_0)$
- alternative hypothesis often only implicit

large  $T$  points to alternative

- $X_1, \dots, X_n$  i.i.d.  $F(\cdot)$
- $H_0 : \mu = \mu_0, \mu = F^{-1}(1/2)$  median of distribution
- $H_1 : \mu > \mu_0$
- test statistic

both  $H$  composite

$$T = \sum_{i=1}^n 1\{X_i > \mu_0\}$$

- under  $H_0$ ,

$$T \sim \text{Binom}(n, 1/2)$$

- $p$ -value

$$p_{obs} = \text{pr}_{H_0}(T \geq t_{obs}) = \sum_{r=t_{obs}}^n \binom{n}{r} \frac{1}{2^n} \doteq 1 - \Phi \left\{ \frac{2(t_{obs} - n/2)}{\sqrt{n}} \right\}.$$

↙ ↘  
↑  $\hat{se}$

$$\left( n \hat{p}_0 (1 - \hat{p}_0) \right)^{\frac{1}{2}}$$

- $H_0 : \mu = \mu_0$       $H_1 : \mu > \mu_0$
- Test statistic  $T = \sum_{i=1}^n \mathbf{1}\{X_i > \mu_0\}$
- Rejection region  $R = \{T \geq c_\alpha\}$
- $c_\alpha \approx n/2 - n^{1/2}z_\alpha/2$

$$\mu = F^{-1}(1/2)$$

Normal approx

- Power =  $\text{pr}_{H_1}(\text{reject } H_0) = \text{pr}_{H_1}(T \geq c_\alpha)$
- to calculate power we need values for  $\mu$  and for  $F$

Need distribution of  $T$  under  $H_1$

- $H_0 : \mu = \mu_0 \quad H_1 : \mu > \mu_0$
- Test statistic  $T = \sum_{i=1}^n 1\{X_i > \mu_0\}$
- Rejection region  $R = \{T \geq c_\alpha\}$
- $c_\alpha \approx n/2 - n^{1/2}z_\alpha/2$

no Normal ass<sup>n</sup>  
needed

$$\mu = F^{-1}(1/2)$$

Normal approx

- Power =  $\text{pr}_{H_1}(\text{reject } H_0) = \text{pr}_{H_1}(T \geq c_\alpha)$
- to calculate power we need values for  $\mu$  and for  $F$

Need distribution of  $T$  under  $H_1$

- SM assumes  $F$  is  $N(\mu, \sigma^2)$ , so

$$\delta = n^{1/2}(\mu_1 - \mu_0)/\sigma$$

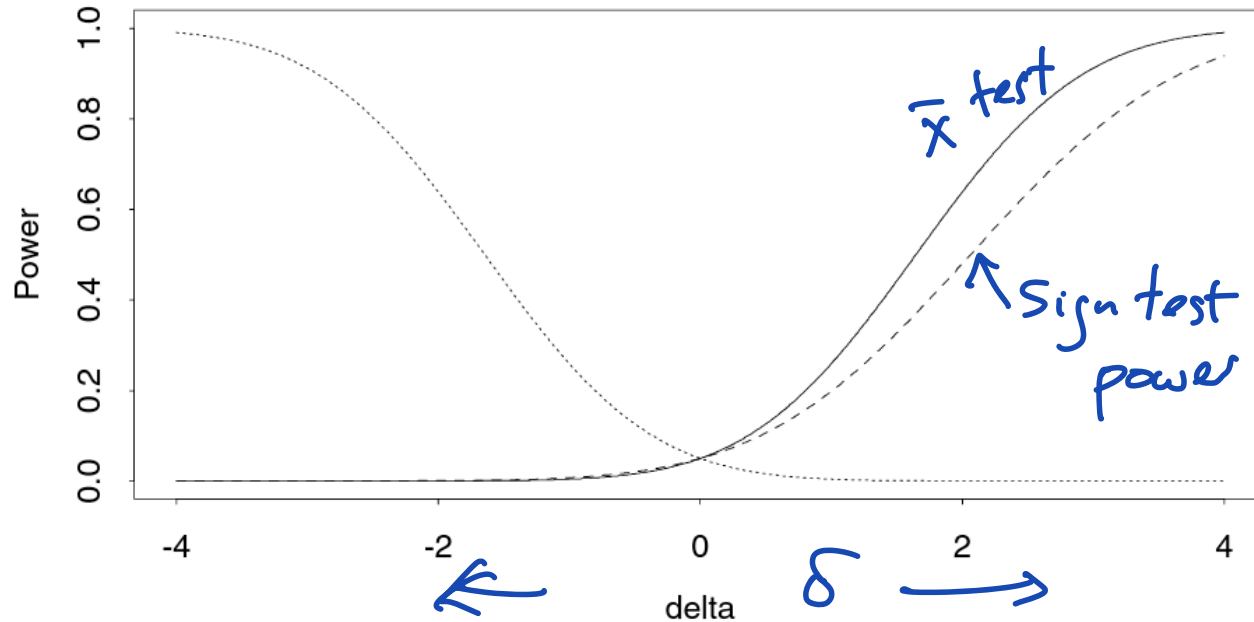
$$\text{pr}_{\mu_1}(T \geq c_\alpha) = \text{pr}_{\mu_1}(T \geq \underbrace{n/2 - n^{1/2}z_\alpha/2}_{c_n}) \doteq \Phi \left\{ \frac{n\Phi(n^{-1/2}\delta) - n/2 + n^{1/2}z_\alpha}{[n\Phi(n^{-1/2}\delta)\{1 - \Phi(n^{-1/2})\}]} \right\}$$

$$\doteq \Phi\{z_\alpha + \delta(2/\pi)^{1/2}\}$$

- test based on  $\bar{X}$  has power  $\Phi(z_\alpha + \delta)$

334

## 7 · Estimation and Hypothesis Testing



**Figure 7.6** Power functions for a test of whether the mean of a  $N(\mu, \sigma^2)$  random sample of size  $n$  equals  $\mu_0$  against the alternative  $\mu = \mu_1$ , as a function of  $\delta = n^{1/2}(\mu_1 - \mu_0)/\sigma$ . The test size is  $\alpha = 0.05$ . The solid curve is the power function for a test of  $\mu_1 > \mu_0$  based on  $\bar{y}$ , and the dashed line is the power function for the sign test. Both critical regions are of form  $\bar{y} > t_\alpha$ . The dotted curve is the power function for  $\bar{y}$  when the critical region is  $\bar{y} < t_\alpha$ .

$$\bar{X}_n \sim N(\mu, \sigma^2/n)$$

$$\hat{\sigma}^2 \xrightarrow{p} ?$$

$$E(\hat{\sigma}^2) = \dots$$

$$? P(|\bar{x} - \mu| > \frac{\varepsilon}{\sqrt{n}}) < \delta \text{ var}(\bar{X}) ?$$

$$\text{WLLN } \hat{\sigma}^2 \xrightarrow{p} \dots$$
$$\text{var}(\hat{\sigma}^2) \rightarrow 0 \uparrow$$