Nonparametric Inference in Cosmology and Astrophysics: Biases and Variants

Christopher R. Genovese Department of Statistics Carnegie Mellon University http://www.stat.cmu.edu/~genovese/

Collaborators: Larry Wasserman, Peter Freeman, Chad Schafer, Alex Rojas, and the InCA group (www.incagroup.org) Department of Statistics Carnegie Mellon University

This work partially supported by NIH Grant 1 R01 NS047493-01 and NSF Grants ACI 0121671.

Example: Dark Energy

Gold SNe Sample



Redshift z

Example: Dark Energy (cont'd)

$$Y_i = r(z_i) + \epsilon_i, \qquad i = 1, \ldots, n.$$

 Y_i : distance measure derived from supernova luminosity z_i : redshift

Want to make inferences about the Dark Energy Equation of State

$$w(z) = T(r, r', r'')$$

=
$$\frac{H_0^2 \Omega_M (1+z)^3 + \frac{2}{3} \frac{r''(z)}{(r'(z))^3}}{H_0^2 \Omega_M (1+z)^3 - \frac{1}{(r'(z))^2}}$$

Example: Dark Energy (cont'd)



Example: CMB







Example: CMB Power Spectrum (cont'd)

$$Y_i = f(x_i) + \epsilon_i, \qquad i = 1, \dots, n$$

- Y_i : estimated power spectrum
- x_i : multipole index

Want to make inferences about features of f such as location and relative heights of peaks.



Example: Galaxy Star Formation Rate

A galaxy's evolution is affected by its local environment.

SFR $S_i = f(D_i) + \epsilon_i$ Density $D_i = \hat{\rho}(x_i) \approx \rho(x_i) + \delta_i$.





Estimated Conditional Density Functions

Why Nonparametric?

- 1. When we don't have a well-justified parametric (finite-dimensional) model for the object of interest.
- 2. When we have a well-justified parametric model but have enough data to go after even more detail.
- 3. When we can do as well (or better) more simply.
- 4. As a way of assessing sensitivity to model assumptions.

Goal: make sharp inferences about unknown functions with a minimum of assumptions.

This involves estimation procedures, but we also need an accurate assessment of uncertainty.

Road Map

- 1. Smoothing
- 2. The Six Biases
- 3. Confidence Sets

Road Map

- 1. Smoothing
- 2. The Six Biases
- 3. Confidence Sets

The Nonparametric Regression Problem

Observe data (X_i, Y_i) for $i = 1, \ldots, n$ where

 $Y_i = f(X_i) + \epsilon_i,$

where $E(\epsilon_i) = 0$ and the X_i s can be fixed (x_i) or random. Leading cases: 1. $x_i = i/n$ and $Cov(\epsilon) \equiv \Sigma = \sigma^2 I$.

2.
$$X_i \text{ IID } g \text{ and } \text{Cov}(\epsilon) \equiv \Sigma = \sigma^2 I.$$

Key Assumption: $f \in \mathcal{F}$ for some infinite dimensional space \mathcal{F} . Examples

1. Sobolev: $\mathcal{F} \equiv \mathcal{W}_p(C) = \{f: \int |f|^2 < \infty \text{ and } \int |f^{(p)}|^2 \leq C^2\}$ 2. Lipschitz: $\mathcal{F} \equiv \mathcal{H}(A) = \{f: |f(x) - f(y)| \leq A|x - y|, \text{ for all } x, y\}$

Goal: Make inferences about f or about specific features of f.

Variants of the Problem

- \bullet Inference for Derivatives of f
- Estimating Variance functions
- Regression in High dimensions
- Inferences about specific functionals of f

Related Problems:

- Density Estimation
- Spectral Density Estimation

 $\log \sigma^2(x)$



Rate-Optimal Estimators

Choose a performance measure, or risk function, e.g., $R(\hat{f}, f) = \mathsf{E} \int (\hat{f} - f)^2$ or $R(\hat{f}, f) = \mathsf{E} |\hat{f}(x_0) - f(x_0)|^2$)

Want \hat{f} that minimizes worst-case risk over \mathcal{F} (minimax). But typically must settle for achieving the optimal minimax rate of convergence r_n :

 $\inf_{\widehat{f}_n} \sup_{f \in \mathcal{F}} R(\widehat{f}_n, f) \asymp r_n$

In infinite-dimensional problems, $r_n\sqrt{n} \rightarrow \infty$.

For example, $r_n = n^{-\frac{2p}{2p+1}}$ on \mathcal{W}_p .

Rate-optimal estimators exist for a wide variety of spaces and risk functions.

Smoothing Methods

- (Quasi-) Linear Methods
 - Kernels and Local Polynomial Regression
 - Roughness Penalty Regularization
 - $(\widehat{f}_n = \arg\min_{\xi} \sum_{i=1}^n (Y_i \xi(X_i))^2 + \lambda Q(f)$, e.g., smoothing splines)
 - Basis Decomposition
- Nonlinear Methods
 - Wavelet Shrinkage
 - Variable Bandwidth Kernels

(Attempt to adapt spatially by changing smoothing over domain; appealing but hard.)

- Others
 - Scale-Space Methods (e.g., SiZer)

(Consider all levels of smoothing simultaneously.)

(Quasi-) Linear Smoothers

An estimator \widehat{f}_n is a *linear smoother* if there exists functions $s(x) = (s_1(x), \ldots, s_n(x))$ such that

$$\widehat{f}_n(x) = \sum_{j=1}^n s_j(x) Y_j.$$

Writing
$$\widehat{f}_n = (\widehat{f}_n(x_1), \dots, \widehat{f}_n(x_n))$$
 and $S_{ij} = s_j(x_i)$, we have
 $\widehat{f}_n = SY$.

For nonparametric smoothers, the function $s(x) \equiv s_h(x)$ depends on a free *smoothing parameter* h that governs complexity of \hat{f} .

The smoothing parameter must be selected, usually from the data. A quasi-linear smoother is linear conditional on h.

(Quasi-) Linear Smoothers: Kernels

The Nadaraya-Watson Kernel estimator, for suitable function K and bandwidth h, is

$$\widehat{f}_n(x) = \frac{\sum_{i=1}^n K\left(\frac{X_i - x}{h}\right) Y_i}{\sum_{j=1}^n K\left(\frac{X_j - x}{h}\right)} = \sum_{i=1}^n s_i(x) Y_i.$$

This weighted average is defined by a *local* optimization problem: for each x, find $a_0 \equiv a_0(x)$ to minimize

$$\sum_{i=1} K\left(\frac{X_i - x}{h}\right) (Y_i - a_0)^2.$$

This is fitting a local constant.

(Quasi-) Linear Smoothers: Local Polynomials

Consider higher order approximation for u near x:

$$f(u) \approx a_0 + a_1(u - x) + \dots + a_p \frac{(u - x)^p}{p!} \equiv p_x(u; a).$$

Let $\hat{a} \equiv \hat{a}(x)$ minimize

$$\sum_{i=1} K\left(\frac{X_i - x}{h}\right) (Y_i - p_x(X_i; a))^2.$$

Then,

$$\widehat{f}_n(x) = p_x(x; \widehat{a}) = \widehat{a}_0(x).$$

The p = 0 case gives back the kernel estimator.

Remarks on Local Polynomials

- Although $\widehat{f}_n(x) = \widehat{a}_0(x)$, this is not simply fitting a local constant.
- Performance is insensitive to the choice of kernel K but highly sensitive to the choice of smoothing parameter h.
- Local polynomial estimators for p > 0 automatically correct some of the biases inherent in the kernel estimator.
- Work well for estimating derivatives. For estimating the vth derivative, prefer p v odd.

Nonlinear Smoothers: Wavelet Shrinkage

Wavelets provide a sparse representation for a wide class of functions.

$$\begin{split} f &= \sum_{k} \alpha_{J_0,k} \varphi_j + \sum_{j=J_0}^{\infty} \sum_{k} \beta_{jk} \psi_{jk} \\ &= \text{Coarse Part} + \text{Successively Refined Details} \end{split}$$

Schematic:

- 1. Fast Discrete Wavelet Transform (DWT): $\begin{pmatrix} \widetilde{\alpha} \\ \widetilde{\beta} \end{pmatrix} = W\mathbf{Y}$.
- 2. Nonlinear shrinkage of detail coefficients: $\hat{\beta} = \eta_{\lambda}(\tilde{\beta})$. For example,

$$\widehat{eta} = \mathsf{sgn}(\widetilde{eta}) \, \left(|\widetilde{eta}| - \lambda
ight)_+.$$

3. Invert Transform: $\widehat{\mathbf{f}} = W^{-1}(\widetilde{\alpha}, \widehat{\beta})^T$.

Nonlinear Smoothers: Wavelet Shrinkage (cont'd)

Result: Adaptive Estimators

The same procedure is nearly optimal (asymptotic minimax) across a range of assumed spaces.

Current State of the Art: Johnstone and Silverman (2005)

- Put mixture prior on each β
- Posterior median yields threshold
- Works with translation invariant transform
- Handles boundaries, estimates derivatives, excellent performance.

The Bias-Variance Tradeoff

All of these methods have a smoothing parameter that determines the complexity of the fit.

Key goal: Choose the correct level of smoothing.

 $R(\widehat{f},f) = \mathsf{E} \int (f(x) - \widehat{f}(x))^2 \, dx = \int \mathsf{bias}^2(x) \, dx + \int \mathsf{variance}(x) \, dx$



The Bias-Variance Tradeoff (cont'd)

Example: CMB Spectrum





400

Choosing the Smoothing Parameter

1. Plug-in

Asymptotically optimal bandwidths are often of the form $h_* = C(f)r_n$ for $r_n \to 0$. Then, $\hat{h}_{\rm PI} = C(\tilde{f})r(n)$ for a pilot estimate \tilde{f} . Do not perform well in general (Loader 1999).

2. Cross-Validation

Divide $\{1, \ldots, n\}$ into m disjoint subsets S_1, \ldots, S_m . Let $\widehat{f}^{[-\ell]}$ be the estimate obtained with *the same procedure* but omitting the subset of data (X_i, Y_i) with $i \in S_\ell$.

Choosing the Smoothing Parameter (cont'd)

Then
$$\widehat{\mathsf{PE}}(h) = \frac{1}{m} \sum_{\ell=1}^{m} \frac{1}{\#(S_{\ell})} \sum_{i \in S_{\ell}} (Y_i - \widehat{f}^{[-\ell]}(X_i))^2.$$

When m = n and $S_{\ell} = \{\ell\}$, this becomes

$$\widehat{\mathsf{PE}}(h) = \frac{1}{n} \sum_{i=1}^{n} \left(\frac{Y_i - \widehat{f}(X_i)}{1 - S_{ii}(h)} \right)^2,$$

where S(h) is the smoothing matrix. Choose \hat{h} to minimize $\widehat{PE}(h)$. (There are several variants.) Then, $PE(\hat{h}) \leq c_1 \inf_h PE(h) + o(1)$.

Choosing the Smoothing Parameter (cont'd)

3. Risk Estimation

Would like h to minimize the risk $R(\hat{f}_h, f) = \int (\hat{f}_h - f)^2$. Find $\hat{R}(h)$ such that $\mathbb{E}\hat{R}(h) = R(\hat{f}_h, f)$. Choose \hat{h} to minimize $\hat{R}(h)$.

In many cases, can show that this approximates true minimum. Example: Basis coefficients $\hat{f}_h = \sum_{k=0}^h \beta_k \phi_k$, h = 0, ... n.

$$R(\widehat{f}_h, f) = h \frac{\sigma^2}{n} + \sum_{k=h+1}^n \beta_k^2 \iff \widehat{R}(h) = h \frac{\sigma^2}{n} + \sum_{k=h+1}^n (\beta_k^2 - \frac{\sigma^2}{n}).$$

Whatever the method, $\frac{\widehat{h}_n - h_n}{h_n} \to 0$ slowly. Estimating the tuning parameter is an intrinsically hard problem.

Road Map

- 1. Smoothing
- 2. The Six Biases
- 3. Confidence Sets

The Six Biases

- 1. Model Bias
- 2. Design Bias
- 3. Boundary Bias
- 4. Derivative Bias
- 5. Measurement-Error Bias
- 6. Coverage Bias

Most of these can be dealt with well-understood statistical techniques.

Model Bias

Without assumptions, there can be no conclusions. – John Tukey

The performance of nonparametric procedures can depend on the space \mathcal{F} assumed to contain f.

Choosing \mathcal{F} too restrictively incurs a bias from the un-modeled component of f; choosing \mathcal{F} too expansively can lead the procedure's to be driven by implausible cases.

And how can we justify the abstract choice of a particular space or ball size etc. in terms of concrete data and prior information?

Model Bias (cont'd): Adaptive Estimators

It's unsatisfying to depend too strongly on intangible assumptions such as whether $f \in \mathcal{W}_p(C)$ or $f \in \mathcal{H}(A)$.

Instead, we want procedures to *adapt* to the unknown smoothness.

For example, \hat{f}_n is a *(rate) adaptive procedure* over the \mathcal{W}_p spaces if when $f \in \mathcal{W}_p$

 $\widehat{f}_n \to f$ at rate $n^{-2p/2p+1}$

without knowing p.

Rate adaptive estimators exist over a variety of function families and over a range of norms (or semi-norms).

Adaptive confidence sets? (later)

Design Bias

Some procedures affected by clustering or positioning of X_i s. Example: Kernel and local linear estimators have

variance
$$(x) = \frac{1}{nh_n} \frac{\sigma^2(x) \int K^2(u) du}{g(x)} + o_P\left(\frac{1}{nh_n}\right).$$

But, the kernel estimator p = 0 has bias

$$h_n^2 \left(\frac{1}{2} f''(x) + \frac{f'(x)g'(x)}{\underbrace{g(x)}_{\text{design bias}}} \right) \int u^2 K(u) du + o_P(h^2)$$

whereas the local linear estimator p = 1 has asymptotic bias

$$h_n^2\left(\frac{1}{2}f''(x)\right)\int u^2K(u)du+o_P(h^2).$$

Design Bias (cont'd)

Another Example: Wavelets and choice of origin



from Coifman and Donoho (1995)

One solution: Translation Invariant wavelet transform

 $\widehat{f} = \frac{1}{n} \sum_{\Delta=1}^{n} \text{Shift}_{\Delta}^{-1} \circ \text{DWT}^{-1} \circ \text{Threshold} \circ \text{DWT} \circ \text{Shift}_{\Delta},$

where $Shift_{\Delta}$ is a circular shift of an n vector by Δ components.

Boundary Bias

Many methods exhibit nontrivial biases near data boundaries.



This is potentially important in problems such as the CMB and Dark Energy because boundary behavior is of scientific interest.

Boundary biases get worse in high-dimensions because a greater proportion of points lie near boundaries.

Boundary Bias (cont'd)

Example: Kernel Smoothing

The kernel "hits" fewer points near the boundary.

Local polynomial for $p \ge 1$ is an easy fix, called "automatic boundary carpentry".

Example: Periodic Wavelets

Standard orthonormal wavelets on an interval are periodic, but can lead to significant bias when the function is not periodic. Cohen, Daubechies, Jawerth and Vial (1993) construction: a preconditioning step that creates boundary wavelets.

Derivative Bias

When estimating derivatives f', f'', ..., $f^{(d)}$, note that $(\widehat{f}_n)^{(d)}$ is NOT a good estimator of $f^{(d)}$. Illustration: Let $\phi_0(x) = 1$ and $\phi_k(x) = \sqrt{2} \cos(\pi k x)$, for $k \ge 1$. If $f = \sum \beta_k \phi_k$ on [0, 1], then $f'' = \sum -\pi^2 k^2 \beta_k \phi_k$. $k \ge 0$ k > 0Let $\hat{\beta}_k = \frac{1}{n} \sum_{i=0}^n Y_i \phi_k(x_i) \approx N(\beta_k, \sigma^2/n).$ Then, $\hat{\beta}_k^{(d)} \equiv -\pi^2 k^2 \hat{\beta}_k \approx N(-\pi^2 k^2 \beta_k, \pi^4 k^4 \frac{\sigma^2}{\sigma})$. Ouch. Inference for derivatives requires different levels of smoothing.

Derivative Bias (cont'd)

With local polynomial estimation, this takes a manageable form.

Recall $p_x(u; a) = a_0 + a_1(u - x) + \dots + a_p \frac{(u - x)^p}{p!}$ and $\hat{f}_n(x) = \hat{a}_0(x)$. Then, $\hat{f}_n^{(d)}(x) = \hat{a}_d(x)$. Note $\hat{a}_0^{(d)} \neq \hat{a}_d$. Choose p - d odd (e.g., p = d + 1) to avoid boundary and design bias.

The (asymptotically) optimal bandwidth for d derivatives is then

$$h_d(k) = C(k, p) h_0$$

where C(k, p) is a known constant.

What about confidence sets? There it's not so easy. Example: Dark Energy.

Measurement-Error Bias

Errors on the covariates can lead to counter-intuitive effects. Simplified Example: Galaxy Star Formation Rates versus Local Density

SFR
$$S_i = f(D_i) + \epsilon_i$$

Density $D_i = \hat{\rho}(x_i) \approx \rho(x_i) + \delta_i$.

Common result is *attenuation bias*, but the opposite is possible (Carroll, Ruppert, and Stefanski 1995).

Bias depends on the nature and extent of the errors.

Measurement-Error Bias

Observe:

$$Y_i = f(X_i) + \epsilon_i$$
$$\widetilde{X}_i = X_i + U_i$$

If we use the (\widetilde{X}_i, Y_i) to estimate f, the estimator will be inconsistent $\int (\widehat{f}_n(x) - f)^2 \neq 0.$

Not enough to simply increase the error bars on Y.

The extra bias
$$\sigma_U^2\left(\frac{g'(x)}{g(x)}f'(x) + \frac{f''(x)}{2}\right) \not\to 0.$$

Measurement-Error Bias (cont'd)

Toy Illustration of Attenuation Bias:

$$Y_i = \beta X_i + \epsilon_i$$
, but we observe (W_i, Y_i) where $W_i = X_i + U_i$.

Then, least squares estimate consistent for $\tilde{\beta} = \beta \frac{\sigma_X^2}{\sigma_V^2 + \sigma_T^2} < \beta$.

Also
$$\operatorname{Var}(Y \mid W) = \sigma_{\epsilon}^2 + \beta^2 \frac{\sigma_X^2 \sigma_U^2}{\sigma_X^2 + \sigma_U^2}.$$

This is not just an issue of variance – *both* variance and bias are affected.

Measurement-Error Bias (cont'd)

One solution: SIMEX (Cook and Stefanski 1994)

Key idea: determine effect of error empirically via simulation.

Toy Illustration Revisited: Given data sets with measurement error variances by factors of $(1 + \lambda_m)$ for $0 = \lambda_1 < \lambda_1 < \cdots < \lambda_m$

Get a regression for $r(\lambda) = \frac{\beta^2 \sigma_X^2}{\sigma_X^2 + (1+\lambda)\sigma_U^2}$. Want to estimate r(-1).

In general, new data sets made from old by adding noise.



Measurement-Error Bias (cont'd)

Another Solution: Special kernels (Fan and Truong 1990)

$$\widehat{f}_n(x) = \frac{\sum_{i=1}^n K_n\left(\frac{x - \widetilde{X}_i}{h_n}\right) Y_i}{\sum_{i=1}^n K_n\left(\frac{x - \widetilde{X}_i}{h_n}\right)}$$

where

$$K_n(x) = \frac{1}{2\pi} \int e^{-itx} \frac{\phi_K(t)}{\phi_U(t/h_n)} dt,$$

where ϕ_K is the Fourier transform of a kernel K and ϕ_U is the characteristic function of U.

This is a standard kernel estimator except for the unusual kernel K_n , chosen to reduce bias.

How to do confidence bands/spheres for this case? We are currently working on that.

Coverage Bias

Using a rate-optimal smoothing parameter gives

 $\mathsf{bias}^2\approx\mathsf{var}.$

Loosely, if
$$\tilde{f} = \mathsf{E}\hat{f}$$
 and $s = \sqrt{\mathsf{Var}\,\hat{f}}$, then
$$\frac{\hat{f} - f}{s} = \frac{\hat{f} - \tilde{f}}{s} + \frac{\tilde{f} - f}{s} \approx \mathsf{N}(0, 1) + \frac{\mathsf{bias}}{\sqrt{\mathsf{var}}}.$$

So, " $\widehat{f} \pm 2s$ " undercovers.

Two common solutions in the literature:

- Bias Correction: Shift confidence set by estimated bias.
- Undersmoothing: Smooth so that var dominates bias².

Road Map

- 1. Smoothing
- 2. The Six Biases
- 3. Confidence Sets

Confidence Sets

In practice, we usually need more than \widehat{f} .

We want to make *inferences* about features of f: shape, magnitude, peaks, inclusion, derivatives.

One approach: construct a $1 - \alpha$ confidence set for f, a random set C such that $P\{C \ni f\} = 1 - \alpha$.

Usually C is a ball or a band around some \widehat{f} .

Three challenges:

- 1. Bias
- 2. Simultaneity
- 3. Relevance

Relevance

In small samples, confidence balls and bands need not constrain all features of interest.

For example, number of peaks:



Alternative: confidence intervals for *specific functionals* of f

Two practical problems:

- 1. Many relevant functionals (e.g., peak locations) hard to work with.
- 2. One often ends up choosing functionals post-hoc.

Better to obtain construct a confidence set for the whole object with post-hoc protection for inferences about many functionals.

Remark: What We Want

1. For asymptotic confidence procedures, prefer uniform coverage:

$$\sup_{f\in\mathcal{F}} \left| \mathsf{P} \Big\{ \mathcal{C}_n \ni f \Big\} - (1-\alpha) \right| \to 0.$$

This ensures that the coverage error depends only on n, not on f.

2. We would also like adaptive confidence procedures.

That is, maintain coverage on \mathcal{F} but can use the data to tailor the set's diameter to the unknown f.

Confidence Bands Cannot Adapt

Confidence bands are of the form $C = \{f: L \leq f \leq U\}$ for some random functions L and U.

Unfortunately, confidence bands whose width is determined from the data cannot do better than fixed-width bands.

Let $\ensuremath{\mathcal{D}}$ denote fixed-diameter confidence band. Then,

$$\liminf_{n\to\infty} \frac{\inf_{\mathcal{C}} \inf_{f\in\mathcal{F}} \mathsf{E}_f(s_n(\mathcal{C}))}{\inf_{\mathcal{D}} \inf_{f\in\mathcal{F}} \mathsf{E}_f(s_n(\mathcal{D}))} > 0.$$

This continues to hold (Low 1997, Genovese and Wasserman 2005) even when smoothness constraints are imposed.

Bottom line: For commonly used smoothers, neither the width nor the tuning parameter of the optimal confidence bands depends on the data.

Building Confidence Bands: Volume of Tubes

If
$$\widehat{f}(x) = \sum_{i=1}^{n} \ell_i(x) Y_i$$
, for weights $\ell(x) = (\ell_1(x), \dots, \ell_n(x))$, then

$$\inf_{f\in\mathcal{F}} \mathsf{P}\left\{\widehat{f}(x) - c\widehat{\sigma} \|\ell(x)\| \le f(x) \le \widehat{f}(x) + c\widehat{\sigma} \|\ell(x)\|, \forall x\right\} = 1 - \alpha,$$

for suitable class \mathcal{F} (Sun and Loader 1994). The constant c solves the equation $\alpha = K_{\ell}\phi(c) + 2(1 - \Phi(c))$.

Special case:
$$f(x) = \sum_{i=1}^{n} \ell_i(x)\theta_i$$
.
Then, $|\widehat{f}(x) - f(x)| = \left|\sum_{i=1}^{n} \ell_i(x)\epsilon_i\right| = |\langle \ell(x), \epsilon \rangle|$, so
 $\alpha = \mathsf{P}\left\{\sup_{x} \left|\frac{\widehat{f}(x) - f(x)}{\|\ell(x)\|}\right| > c\sigma\right\} = \mathsf{P}\left\{\sup_{x} \left|\langle \frac{\ell(x)}{\|\ell(x)\|}, \frac{\epsilon}{\|\epsilon\|}\rangle\right| > \frac{c\sigma}{\|\epsilon\|}\right\}$

Reduces to finding the volume of a tube on the sphere S^{n-1} .

Building Confidence Bands: Parametric

Approximately minimum expected size parametric confidence bands (Schafer 2004)



Confidence Balls

1. Expand $f = \sum_k \beta_k \phi_k$ in orthonormal basis (e.g., cosine basis). 2. Shrink naive estimators by $\hat{\beta}_k = \lambda_k \tilde{\beta}_k$, $1 \ge \lambda_1 \ge \cdots \ge \lambda_n$. 3. Choose λ by minizing estimated risk $\hat{R}(\lambda)$.

4.
$$C_n = \left\{ \beta : \sum_k (\hat{\beta}_k - \beta_k)^2 \leq \frac{k_\alpha}{\sqrt{n}} + \hat{R}(\hat{\lambda}) \right\}.$$

Confidence balls can adapt to unknown smoothness.

Can impose constraints post-hoc and get valid inferences.

Confidence Ball Center vs Concordance Model



- Concordance model is an MLE based on WMAP and four other data sets.
- Confidence ball center based on WMAP data only.

Eyes on the Ball I: Parametric Probes

Simultaneous 95% CIs on Peak Heights, Locations, and Height Ratios



Multipole /

Eyes on the Ball I: Parametric Probes (cont'd)

Varied baryon fraction $(\Omega_b h^2)$ keeping $\Omega_{ ext{total}} \equiv 1$



Extended search over millions of spectra (Bryan et al. 2006).

Eyes on the Ball II: Model Checking

Inclusion in the confidence ball provides simultaneous goodness-of-fit tests for parametric (or other) models.



Comment: Coverage and Posterior Probability

 \bullet Frequentist Confidence Set ${\cal C}$

$$\min_{f} \mathsf{P}\big\{\mathcal{C} \ni f\big\} \ge 1 - \alpha. \tag{1}$$

 \bullet Bayesian Posterior Region ${\cal B}$

$$\mathsf{P}\left\{f\in\mathcal{B}\mid\mathsf{Data}\right\}\geq1-lpha.$$
 (2)

• In nonparametric problems, can have (2) hold and yet have

$$\min_{f} \mathsf{P}\big\{\mathcal{B} \ni f\big\} \approx 0. \tag{3}$$

Road Map

- 1. Smoothing
- 2. The Six Biases
- 3. Confidence Sets

Take-Home Points

- The crucial decision in nonparametric estimation is choosing the correct amount of smoothing.
- We must avoid the six biases, but fortunately methods exist to deal with most of them.
- Estimates alone are not enough, also need an assessment of uncertainty. Various types of confidence sets can be constructed, but their dependence on the assumptions can be delicate.