Computationally Efficient Nonparametric Testing

Guang Cheng¹ Dept of Statistics Purdue University www.science.purdue.edu/bigdata/

> Statistics@Cambridge May 26, 2017

¹Acknowledge NSF, ONR and Simons Foundation.

Acknowledgment



Zuofeng Shang @ Binghamton



Meimei Liu @ Purdue

Million Song Dataset

- 1,000,000 music tracks released from the year 1922 to 2011
- Predictor variables: 12 timbre averages
- Goal: identify predictors associated with the year of release
- Data source: free, public, supported in part by NSF



• Consider a nonparametric model

$$y_i = f(x_i) + \epsilon_i$$
, for $i = 1, \ldots, n$

• The standard KRR estimate f by

$$\widehat{f}_n := \operatorname*{arg\,min}_{f \in \mathcal{H}} \frac{1}{n} \sum_{i=1}^n (y_i - f(x_i))^2 + \underbrace{\lambda \|f\|_{\mathcal{H}}^2}_{\text{roughness penalty}}$$

reproducing kernel Hilbert space (RKHS)

Computational challenges

• The KRR theory says that

$$\widehat{f}_n(\cdot) = \sum_{i=1}^n \widehat{w}_i K(\cdot, x_i),$$

$$\downarrow_{\text{kernel function}}$$

where

$$\widehat{\mathbf{w}} = n^{-1} \underbrace{(\mathbf{K} + \lambda \mathbf{I}_{n})^{-1}}_{n \times n} \underbrace{\mathbf{y}}_{n \times 1},$$

 ${\bf K}$ is $n\times n$ kernel matrix and ${\bf y}$ is n-dimensional response

Computational challenges

• The KRR theory says that

$$\widehat{f}_n(\cdot) = \sum_{i=1}^n \widehat{w}_i K(\cdot, x_i),$$

where

$$\widehat{\mathbf{w}} = n^{-1} \underbrace{(\mathbf{K} + \lambda \mathbf{I}_{n})^{-1}}_{n \times n} \underbrace{\mathbf{y}}_{n \times 1},$$

 ${\bf K}$ is $n\times n$ kernel matrix and ${\bf y}$ is n-dimensional response

• Time (in seconds) taken to invert an *n*-dimensional matrix:

| \overline{n} | 10^{3} | 10^{4} | 10^{5} | 10^{6} |
|----------------|----------|----------|----------|--------------|
| Time | 0.78 | 148.8 | 28269.5 | > two months |

How to address "curse of sample size?"



Our contributions

- We propose two computationally efficient testing methods in nonparametric settings:
 - one is based on randomized sketches
 - another employs parallel computing
- Characterize computational limits, i.e., the minimal computational cost to preserve statistical optimality
- Existing work only focus on estimation [Zhang et al, 2015; Yang et al, 2016]

1 Nonparametric testing based on randomized sketches

2 Nonparametric testing based on parallel computing

3 Numerical results

• Consider a hypothesis testing problem

$$H_0: f = f_0$$
 vs. $H_1: f \neq f_0$

• A Wald type test statistic is

$$T_{n,\lambda} = \|\widehat{f}_R - f_0\|_{L^2}^2$$

$$\widehat{\square}_{\text{an estimator of } f}$$

• Reject H_0 if $T_{n,\lambda}$ is large

An algorithm based on randomized sketches

• The standard KRR is equivalent to the following

$$\widehat{\mathbf{w}} = \operatorname*{arg\,min}_{\mathbf{w}\in\mathbb{R}^n} \mathbf{w}^T \mathop{\mathbf{K}^2}_{\substack{\uparrow\\n\times n}} \mathbf{w} - \frac{2}{n} \mathbf{y}^T \mathop{\mathbf{K}}_{\mathbf{w}} \mathbf{w} + \lambda \mathbf{w}^T \mathop{\mathbf{K}}_{\substack{\uparrow\\n\times n}} \mathbf{w} \qquad (1)$$

An algorithm based on randomized sketches

• The standard KRR is equivalent to the following

$$\widehat{\mathbf{w}} = \operatorname*{arg\,min}_{\mathbf{w}\in\mathbb{R}^n} \mathbf{w}^T \mathop{\mathbf{K}^2}_{\substack{\uparrow\\n\times n}} \mathbf{w} - \frac{2}{n} \mathbf{y}^T \mathop{\mathbf{K}}_{\mathbf{w}} \mathbf{w} + \lambda \mathbf{w}^T \mathop{\mathbf{K}}_{\substack{\uparrow\\n\times n}} \mathbf{w} \qquad (1)$$

• Let **S** be an $s \times n$ random matrix. Replacing $\mathbf{w} = \mathbf{S}^T \boldsymbol{\alpha}$, (1) becomes

$$\widehat{\boldsymbol{\alpha}} = \underset{\boldsymbol{\alpha} \in \mathbb{R}^{s}}{\arg\min} \alpha^{T} (\underbrace{\mathbf{S}\mathbf{K}^{2}\mathbf{S}^{\mathrm{T}}}_{s \times s}) \boldsymbol{\alpha} - \frac{2}{n} \mathbf{y}^{T} (\mathbf{K}\mathbf{S}^{T}) \boldsymbol{\alpha} + \lambda \boldsymbol{\alpha}^{T} (\underbrace{\mathbf{S}\mathbf{K}\mathbf{S}^{\mathrm{T}}}_{s \times s}) \boldsymbol{\alpha}$$
Solution:
$$\widehat{\boldsymbol{\alpha}} = n^{-1} (\underbrace{\mathbf{S}\mathbf{K}^{2}\mathbf{S}^{\mathrm{T}}}_{s \times s} + \lambda \operatorname{\mathbf{S}KS}^{\mathrm{T}}_{s})^{-1} (\mathbf{S}\mathbf{K}) \mathbf{y}$$

An algorithm based on randomized sketches

• The standard KRR is equivalent to the following

$$\widehat{\mathbf{w}} = \operatorname*{arg\,min}_{\mathbf{w}\in\mathbb{R}^n} \mathbf{w}^T \mathop{\mathbf{K}^2}_{\substack{\uparrow\\n\times n}} \mathbf{w} - \frac{2}{n} \mathbf{y}^T \mathop{\mathbf{K}} \mathbf{w} + \lambda \mathbf{w}^T \mathop{\mathbf{K}}_{\substack{\uparrow\\n\times n}} \mathbf{w} \qquad (1)$$

• Let **S** be an $s \times n$ random matrix. Replacing $\mathbf{w} = \mathbf{S}^T \boldsymbol{\alpha}$, (1) becomes

$$\widehat{\boldsymbol{\alpha}} = \underset{\boldsymbol{\alpha} \in \mathbb{R}^{s}}{\arg\min} \alpha^{T} (\underbrace{\mathbf{S}\mathbf{K}^{2}\mathbf{S}^{\mathrm{T}}}_{s \times s}) \boldsymbol{\alpha} - \frac{2}{n} \mathbf{y}^{T} (\mathbf{K}\mathbf{S}^{T}) \boldsymbol{\alpha} + \lambda \boldsymbol{\alpha}^{T} (\underbrace{\mathbf{S}\mathbf{K}\mathbf{S}^{\mathrm{T}}}_{s \times s}) \boldsymbol{\alpha}$$

Solution:
$$\widehat{\boldsymbol{\alpha}} = n^{-1} (\underbrace{\mathbf{S}\mathbf{K}^{2}\mathbf{S}^{\mathrm{T}}}_{s \times s} + \lambda \operatorname{\mathbf{S}KS}^{\mathrm{T}}_{s \times s})^{-1} (\mathbf{S}\mathbf{K}) \mathbf{y}$$

• Estimate f by the following sketched KRR (SKRR):

$$\widehat{f}_R(\cdot) = \sum_{i=1}^n (\mathbf{S}^T \widehat{\boldsymbol{\alpha}})_i K(\cdot, x_i)$$

Independent sub-Gaussian entries, e.g., Gaussian or Bernoulli



Computing time vs. projection dimension



Figure 1: Computing time of $T_{n,\lambda}$

Existence of minimal projection dimension



Figure 2: $n = 2^{12}$. Significance level 0.05. Run 1,000 replications

How to characterize minimal projection dimension?



Commonly used kernels

Mercer's Theorem:

$$K(x, x') = \sum_{i=1}^{\infty} \mu_i \phi_i(x) \phi_i(x'),$$

eigenvalue
eigenvalue

- Polynomial kernel: $\mu_k \simeq k^{-2m}$ for m > 1/2, e.g., m = 2 (cubic spline)
- Exponential kernel: $\mu_k \simeq \exp(-\alpha k^p)$ for $\alpha, p > 0$
- Finite rank kernel: $\mu_k = 0 \text{ for } k > r$

| | Polynomial kernel | Exponential kernel |
|-------|---------------------|--------------------------|
| s^* | $n^{rac{2}{4m+1}}$ | $(\log n)^{rac{1}{2p}}$ |

Practical choice of projection dimension



Theorem 1

Suppose $\lambda \to 0$. Then we have under H_0 ,

$$\frac{T_{n,\lambda} - \mu_{n,\lambda}}{\sigma_{n,\lambda}} \xrightarrow{d} N(0,1),$$

where $\mu_{n,\lambda} = E_{H_0}\{T_{n,\lambda}\}$ and $\sigma_{n,\lambda}^2 = Var_{H_0}\{T_{n,\lambda}\}.$

Theorem 1

Suppose $\lambda \to 0$. Then we have under H_0 ,

$$\frac{T_{n,\lambda} - \mu_{n,\lambda}}{\sigma_{n,\lambda}} \xrightarrow{d} N(0,1),$$

where $\mu_{n,\lambda} = E_{H_0}\{T_{n,\lambda}\}$ and $\sigma_{n,\lambda}^2 = Var_{H_0}\{T_{n,\lambda}\}.$

• Testing rule is

$$\phi_{n,\lambda} = I(|T_{n,\lambda} - \mu_{n,\lambda}| > z_{1-\alpha/2}\sigma_{n,\lambda}),$$

$$\phi_{n,\lambda} = 1 \iff \text{reject } H_0$$

Power

Let

$$s^* = \begin{cases} \frac{2}{4m+1}, & \text{polynomial kernel} \\ (\log n)^{\frac{1}{p}}, & \text{exponential kernel} \end{cases}$$

Theorem 2

Suppose $s \gtrsim s^*$ and $\lambda \to 0$. Then for any $\epsilon > 0$, there exist $c_{\epsilon}, n_{\epsilon} > 0$, s.t., for any $n \ge n_{\epsilon}$,

$$\inf_{f \in \mathcal{B}, \|f - f_0\|_{L^2} \ge c_{\epsilon} d_{n,\lambda}} P_f(\phi_{n,\lambda} = 1) \ge 1 - \epsilon, \quad \text{(high power)}$$

where $\mathcal{B} = \{f \in \mathcal{H} : ||f||_{\mathcal{H}} \leq C\}$ for a positive constant C and the separation rate $d_{n,\lambda} = \sqrt{\lambda + \sigma_{n,\lambda}}$. Bias²

• Small separation rate \iff powerful test

Minimal separation rate: Bias² vs. SD tradeoff



 $d_{n,\lambda}$ attains $d_{n,\lambda}^*$ at $\lambda = \lambda^*$. Rate of λ^* : Rademacher complexity

Minimal projection dimension

Theorem 3

Suppose $s \ll s^*$ and $\lambda \to 0$. Then there exists a positive sequence $\beta_{n,\lambda}$ with $\lim_{n\to\infty} \beta_{n,\lambda} = \infty$, s.t.

 $\limsup_{n \to \infty} \inf_{f \in \mathcal{B}, \|f - f_0\|_{L^2} \ge \beta_{n,\lambda} d_{n,\lambda}^*} P_f(\phi_{n,\lambda} = 1) \le \alpha. \quad (\text{low power})$

Minimal projection dimension

Theorem 3

Suppose $s \ll s^*$ and $\lambda \to 0$. Then there exists a positive sequence $\beta_{n,\lambda}$ with $\lim_{n\to\infty} \beta_{n,\lambda} = \infty$, s.t.

 $\limsup_{n \to \infty} \inf_{f \in \mathcal{B}, \|f - f_0\|_{L^2} \ge \beta_{n,\lambda} d^*_{n,\lambda}} P_f(\phi_{n,\lambda} = 1) \le \alpha. \quad (\text{low power})$

• s^* is a *sharp* lower bound for projection dimension

| | Polynomial Kernel | Exponential Kernel |
|-------------------|-----------------------|------------------------------------|
| s^* | $n^{rac{2}{4m+1}}$ | $(\log n)^{\frac{1}{p}}$ |
| $d^*_{n,\lambda}$ | $n^{-rac{2m}{4m+1}}$ | $(\log n)^{\frac{1}{4p}} n^{-1/2}$ |
| λ^* | $n^{-rac{4m}{4m+1}}$ | $(\log n)^{\frac{1}{2p}} n^{-1}$ |

Can we perform hypothesis testing when kernel smoothness is unknown?



Consider a polynomial kernel of order m with m unknown

An adaptive test

Consider a polynomial kernel of order m with m unknown

• For any integer m, find

$$\tau_m = \frac{T_{n,\lambda}(m) - \mu_{n,\lambda}(m)}{\sigma_{n,\lambda}(m)}$$

An adaptive test

Consider a polynomial kernel of order m with m unknown

• For any integer m, find

$$\tau_m = \frac{T_{n,\lambda}(m) - \mu_{n,\lambda}(m)}{\sigma_{n,\lambda}(m)}$$



$$\tau_n^* = \max_{1 \le m \le m_n} \tau_m$$
$$\uparrow_{m_n \asymp (\log n)^{d_0}, d_0 \in (0, 1/2)}$$

An adaptive test

Consider a polynomial kernel of order m with m unknown

• For any integer m, find

$$\tau_m = \frac{T_{n,\lambda}(m) - \mu_{n,\lambda}(m)}{\sigma_{n,\lambda}(m)}$$



$$\tau_n^* = \max_{1 \le m \le m_n} \tau_m$$
$$\uparrow_{m_n \asymp (\log n)^{d_0}, d_0 \in (0, 1/2)}$$

3 Let

Theorem 4

Suppose $m_n \asymp (\log n)^{d_0}$ for $d_0 \in (0, 1/2)$. Then for any $\alpha \in (0, 1)$, under H_0 , $P(\tau_{n,m_n} \le c_\alpha) \to 1 - \alpha$, as $n \to \infty$, where $c_\alpha = -\log(-\log(1 - \alpha))$.

Proof is based on Stein's leave-one-out method (Stein, 1986) Testing rule is

$$\phi_{n,\lambda}^* = I(\tau_{n,m_n} > c_{\alpha}),$$

$$\phi_{n,\lambda}^* = 1 \iff \text{reject } H_0$$

To achieve high power, choose

$$s(m) \gtrsim n^{\frac{2}{4m+1}} (\log \log n)^{-\frac{1}{4m+1}}$$
, for $m = 1, \dots, m_n$

Minimal separation rate:

 $\delta(n,m)$ is optimal (Spokoiny, 1996)

1 Nonparametric testing based on randomized sketches



3 Numerical results

What if we can use parallel computing?



Divide-and-Conquer (DC)

Consider the same nonparametric regression model:

 $y = f(x) + \epsilon$

$$\bar{f}_N = \frac{1}{s} \sum_{j=1}^s \widehat{f}_j$$

• Consider the same hypothesis testing problem

$$H_0: f = f_0$$
 vs. $H_1: f \neq f_0$

• Consider a Wald type test statistic

$$T_{N,\lambda} = \|\bar{f}_N - f_0\|_{L^2}^2$$

• Prefer a large number of divisions to reduce computational cost

Computing time vs. number of divisions



Figure 3: Computing time is decreasing with s

Existence of maximal number of divisions



Figure 4: N=10,000. Significance level 0.05. Run 1,000 replicates

Phase transition diagram

Consider smoothing spline regression with a smoothing parameter λ and smoothness $m \geq 1.$



Figure 5: Locations of (λ, s) achieving testing optimality

Main theorem: testing consistency

Theorem 6

Suppose $\lambda \to 0$, $n \to \infty$ when $N \to \infty$, and $\lim_{N\to\infty} n\lambda^{1/2m}$ exists (which could be infinity). Then, we have under H_0 ,

$$\frac{T_{N,\lambda} - \mu_{N,\lambda}}{\sigma_{N,\lambda}} \xrightarrow{d} N(0,1), \text{ as } N \to \infty,$$

where
$$\mu_{N,\lambda} = \mathbb{E}_{H_0}\{T_{N,\lambda}\}$$
 and $\sigma_{N,\lambda}^2 = \operatorname{Var}_{H_0}\{T_{N,\lambda}\}.$

Testing rule is

$$\phi_{N,\lambda} = I(|T_{N,\lambda} - \mu_{N,\lambda}| \ge z_{1-\alpha/2}\sigma_{N,\lambda}),$$

$$\phi_{N,\lambda} = 1 \iff \text{reject } H_0$$

Main theorem: power

Consider m-order smoothing splines

Let

$$d_{N,\lambda} = \sqrt{\lambda + n^{-2m} + \sigma_{N,\lambda}}$$

Theorem 7

Suppose $s \leq N^{\frac{4m-1}{4m+1}}$, $\lambda \to 0$, $n \to \infty$ when $N \to \infty$, and $\lim_{N\to\infty} n\lambda^{1/2m}$ exists (which could be infinity). Then for any $\varepsilon > 0$, there exist $C_{\varepsilon}, N_{\varepsilon} > 0$ s.t. for any $N \geq N_{\varepsilon}$,

$$\inf_{\substack{f \in \mathcal{B} \\ \|f - f_0\|_2 \ge C_{\varepsilon} d_{N,\lambda}}} P_f\left(\phi_{N,\lambda} = 1\right) \ge 1 - \varepsilon, \quad \text{(high power)}$$

where $\mathcal{B} = \{f \in S^m(\mathbb{I}) : ||f||_{\mathcal{H}} \leq C\}$ for a positive constant C.

•
$$\lambda^* = N^{-\frac{4m}{4m+1}}$$

Theorem 8

Suppose $s \gg N^{\frac{4m-1}{4m+1}}$, $\lambda \to 0$, $n \to \infty$ when $N \to \infty$, and $\lim_{N\to\infty} n\lambda^{1/2m}$ exists (possibly infinity). Then there exists a positive sequence $\beta_{N,\lambda}$ with $\lim_{N\to\infty} \beta_{N,\lambda} = \infty$ s.t.

$$\limsup_{N \to \infty} \inf_{\substack{f \in \mathcal{B} \\ \|f - f_0\|_2 \ge \beta_{N,\lambda} d_{N,\lambda}^*}} P_f(\phi_{N,\lambda} = 1) \le \alpha.$$
 (low power)

• $s^{**} = N^{\frac{4m-1}{4m+1}}$ is a *sharp* upper bound for number of divisions to maintain testing optimality

A "theoretical" suggestion



1 Nonparametric testing based on randomized sketches

2 Nonparametric testing based on parallel computing



Consider the hypothesis testing problem

$$H_0: f = 0$$
 vs. $H_1: f \neq 0$

•
$$y_i = f(x_i) + \epsilon_i, i = 1, \dots, n, n = 2^9, 2^{10}, 2^{11}, 2^{12}$$

•
$$\epsilon_i \stackrel{iid}{\sim} N(0,1)$$
 and $x_i \stackrel{iid}{\sim} Unif[0,1]$

•
$$f(x) = c(3\beta_{30,17}(x) + 2\beta_{3,11}(x))$$
 with $c = 0, 0.01, 0.02, 0.03$

Beta density

• Standard KRR

- Randomized sketches (RS):
 - projection dimension $2n^{2/9}$

- Adaptive randomized sketches (Adaptive RS):
 - projection dimension $2n^{2/9}(\log \log n)^{-1/9}$
 - number of multiple tests $\sqrt{\log n}$

- Divide-and-Conquer (DC):
 - number of divisions $0.5n^{7/9}$

| Method | $n = 2^9$ | $n = 2^{10}$ | $n = 2^{11}$ | $n = 2^{12}$ |
|---------------|-----------|--------------|--------------|--------------|
| Standard KRR | 1.44 | 6.67 | 42.61 | 332.08 |
| \mathbf{RS} | 0.31 | 0.40 | 0.79 | 2.98 |
| Adaptive RS | 0.53 | 1.55 | 4.33 | 15.93 |
| DC | 0.28 | 0.35 | 0.54 | 2.13 |

Computing time (in seconds) of one trial on a Windows computer with 2GB of memory and a single-threaded 2.70Ghz CPU

Size and power



Figure 6: Significance level 0.05. Averages over 1,000 replications

41 / 46

Real data analysis

- Million Song Dataset:
 - 1,000,000 observations (music tracks)
 - year of release \boldsymbol{y}
 - 12 timbre averages x_1, \ldots, x_{12}
- Goal: identify significant timbre averages
- Fit $y = f(x_j)$ +error for $1 \le j \le 12$. Test $H_0: f$ is constant

• Gaussian kernel $K(x, y) = \exp\left(-\frac{(x-y)^2}{\phi}\right)$; unknown $\phi > 0$

- Estimate $\hat{\phi} \approx 3$ based on distributed GCV (Xu, Shang, C., 2016) and 463,715 training samples
- Wald test on 536,285 testing samples:
 - projection dimension $536, 285^{\frac{2}{9}} \approx 19$
 - parallel processors $536, 285^{\frac{7}{9}}/100 \approx 286$

| | x_1 | x_2 | x_3 | x_4 | x_5 | x_6 |
|----|----------------------|----------------------|----------------------|-------------------------|-------------------------|-------------------|
| RS | 0.0318 | 0.0943 | 0.7740 | 0.4069 | 0.7586 | 0.3699 |
| DC | 0.0205 | 0.0319 | 0.8212 | 0.3211 | 0.4729 | 0.4103 |
| | | | | | | |
| | x_7 | x_8 | x_9 | x_{10} | x_{11} | x_{12} |
| RS | $\frac{x_7}{0.0433}$ | $\frac{x_8}{0.2305}$ | $\frac{x_9}{0.1642}$ | $\frac{x_{10}}{0.5591}$ | $\frac{x_{11}}{0.8979}$ | $x_{12} = 0.0201$ |

Table 1: P-values for testing marginal association by two methods

- Matrix multiplication for RS: about 380 seconds
- Wald test: about 10–20 seconds
- One node in a cluster with 20 cores



Figure 7: x_1, x_7, x_{12} have significant patterns, with p-values < 0.05 for both RS and DC methods.

• Nonparametric testing based on RS and DC

• Equally powerful as standard method with less computing cost

• Computational limits are characterized

Thanks! and Questions?

Backup slides

A By-product: minimax optimal estimation

Recall that $T_{n,\lambda} = \|\widehat{f}_R - f_0\|_{L^2}^2$ is the squared estimation error

Corollary (Optimal Estimation Rate)

Suppose that $\lambda \to 0$, $n\lambda \to \infty$. With probability greater than $1 - 2 \exp\{-cn\lambda\}$, we have

$$\|\widehat{f}_R - f_0\|_{L^2}^2 \le c'(\lambda + \frac{s_\lambda}{n})$$

with $s_{\lambda} = \min\{j : \widehat{\mu}_j \leq \lambda\}.$

- Polynomial kernel: $\lambda = n^{-\frac{2m}{2m+1}}$ and $s_{\lambda} = n^{\frac{1}{2m+1}}$
- Exponential kernel: $\lambda = (\log n)^{1/p} n^{-1}$ and $s_{\lambda} = (\log n)^{1/p}$
- Recover the minimax optimal estimation results in Yang, Pilanci and Wainwright (2016, AoS) under the same set of conditions

Define $\Delta = \mathbf{KS}^T (\mathbf{SK}^2 \mathbf{S}^T + \lambda \mathbf{SKS}^T)^{-1} \mathbf{SK}$. Under $H_0: f = 0$, we have $y_i = \epsilon_i$ for i = 1, ..., n, and hence

$$T_{n,\lambda} = \|\widehat{f}_R\|_{L^2}^2 \approx \|\widehat{f}_R\|_n^2 = \frac{1}{n} \epsilon^T \Delta^2 \epsilon \quad \text{(for large } n\text{)}$$

So we have the following approximations for practical use:

$$\mu_{n,\lambda} \approx \frac{1}{n} \operatorname{Tr}(\Delta^2) = \frac{1}{n} \operatorname{Tr}(\Gamma^2)$$

 and

$$\sigma_{n,\lambda}^2\approx \frac{2}{n^2}{\rm Tr}(\Delta^4)=\frac{2}{n^2}{\rm Tr}(\Gamma^4)$$

where

$$\boldsymbol{\Gamma} = \mathbf{S}\mathbf{K}^2\mathbf{S}^T(\mathbf{S}\mathbf{K}^2\mathbf{S}^{\mathrm{T}} + \lambda\,\mathbf{S}\mathbf{K}\mathbf{S}^{\mathrm{T}})^{-1} \hspace{0.1in} (s \times s)$$

Projection dimension for adaptive testing

Theorem 5

Suppose

$$s(m) \gtrsim n^{\frac{2}{4m+1}} (\log \log n)^{-\frac{1}{4m+1}}$$

for each $m \in \{1, \ldots, m_n \asymp (\log n)^{d_0}\}$. Then, for any $\varepsilon > 0$, there exist positive constants $\tilde{c}_{\varepsilon}, \tilde{n}_{\varepsilon}$ such that for any $n \ge \tilde{n}_{\varepsilon}$

$$\inf_{\substack{f \in \mathcal{B}_{n,m} \\ \|f - f_0\|_{L^2} \ge \tilde{c}_{\varepsilon} \delta(n,m)}} P_f(\phi_{n,\lambda}^* = 1) \ge 1 - \varepsilon, \text{(high power)}$$

where $\mathcal{B}_{n,m} = \{f \in \mathcal{H}(m) : (\mathbf{f})^T \mathbf{K}^{-1} \mathbf{f} \leq 1\}$ with $\mathbf{f} = (f(x_1), \cdots, f(x_n))$, and \mathbf{K} is the kernel matrix, and

$$\delta(n,m) \equiv n^{-2m/(4m+1)} (\log \log n)^{m/(4m+1)}.$$
price for adaptivity

• $\delta(n,m)$ is optimal (Spokoiny, 1996)