

Information Consistency of Nonparametric Gaussian Process Methods

Matthias W. Seeger

Max Planck Institute for Biological Cybernetics
Tübingen, Germany

www.kyb.tuebingen.mpg.de/bs/people/seeger/

with Sham M. Kakade, Dean P. Foster



MAX-PLANCK-GESELLSCHAFT



BIOLOGISCHE KYBERNETIK

The Prediction Game

- **Conditional** sequence prediction under **log loss**:

$(\mathbf{x}_1, \dots, \mathbf{x}_n) \rightarrow (y_1, \dots, y_n)$

- Predict $Q(\cdot | \mathbf{y}_{<i}, \mathbf{x}_{\leq i})$
- Incur loss $-\log Q(y_i | \mathbf{y}_{<i}, \mathbf{x}_{\leq i})$

Cumulative loss: $\sum_{i=1}^n -\log Q(y_i | \mathbf{y}_{<i})$

- Fixed $P(y|u)$, function space \mathcal{F} . **Bayesian strategy**:

$$P_{bs}(y_i | \mathbf{y}_{<i}) = \mathbb{E}_{P_{bs}(f | \mathbf{y}_{<i})} [P(y_i | f(\mathbf{x}_i))],$$

$$\underbrace{dP_{bs}(f | \mathbf{y}_{<i})}_{\text{Posterior}} \propto \prod_{j < i} P(y_j | f(\mathbf{x}_j)) \underbrace{dP_{bs}(f)}_{\text{Prior}}$$

Cumulative loss: $L_{bs}(\mathbf{y}_{\leq n}) = -\log P_{bs}(\mathbf{y}_{\leq n})$

- **Expert strategy**: For fixed $f \in \mathcal{F}$ predict $P(y_i | f(\mathbf{x}_i))$.

Cumulative loss: $L_f(\mathbf{y}_{\leq n}) = -\log P(\mathbf{y}_{\leq n} | f)$

The Prediction Game

- **Conditional** sequence prediction under **log loss**:

$(\mathbf{x}_1, \dots, \mathbf{x}_n) \rightarrow (y_1, \dots, y_n)$

- Predict $Q(\cdot | \mathbf{y}_{<i}, \mathbf{x}_{\leq i})$
- Incur loss $-\log Q(y_i | \mathbf{y}_{<i}, \mathbf{x}_{\leq i})$

Cumulative loss: $\sum_{i=1}^n -\log Q(y_i | \mathbf{y}_{<i})$

- Fixed $P(y|u)$, function space \mathcal{F} . **Bayesian strategy**:

$$P_{bs}(y_i | \mathbf{y}_{<i}) = E_{P_{bs}(f | \mathbf{y}_{<i})} [P(y_i | f(\mathbf{x}_i))],$$

$$\underbrace{dP_{bs}(f | \mathbf{y}_{<i})}_{\text{Posterior}} \propto \prod_{j < i} P(y_j | f(\mathbf{x}_j)) \underbrace{dP_{bs}(f)}_{\text{Prior}}$$

Cumulative loss: $L_{bs}(\mathbf{y}_{\leq n}) = -\log P_{bs}(\mathbf{y}_{\leq n})$

- **Expert strategy**: For fixed $f \in \mathcal{F}$ predict $P(y_i | f(\mathbf{x}_i))$.

Cumulative loss: $L_f(\mathbf{y}_{\leq n}) = -\log P(\mathbf{y}_{\leq n} | f)$

The Prediction Game

- **Conditional** sequence prediction under **log loss**:

$(\mathbf{x}_1, \dots, \mathbf{x}_n) \rightarrow (y_1, \dots, y_n)$

- Predict $Q(\cdot | \mathbf{y}_{<i}, \mathbf{x}_{\leq i})$
- Incur loss $-\log Q(y_i | \mathbf{y}_{<i}, \mathbf{x}_{\leq i})$

Cumulative loss: $\sum_{i=1}^n -\log Q(y_i | \mathbf{y}_{<i})$

- Fixed $P(y|u)$, function space \mathcal{F} . **Bayesian strategy**:

$$P_{bs}(y_i | \mathbf{y}_{<i}) = E_{P_{bs}(f | \mathbf{y}_{<i})} [P(y_i | f(\mathbf{x}_i))],$$

$$\underbrace{dP_{bs}(f | \mathbf{y}_{<i})}_{\text{Posterior}} \propto \prod_{j < i} P(y_j | f(\mathbf{x}_j)) \underbrace{dP_{bs}(f)}_{\text{Prior}}$$

Cumulative loss: $L_{bs}(\mathbf{y}_{\leq n}) = -\log P_{bs}(\mathbf{y}_{\leq n})$

- **Expert strategy**: For fixed $f \in \mathcal{F}$ predict $P(y_i | f(\mathbf{x}_i))$.

Cumulative loss: $L_f(\mathbf{y}_{\leq n}) = -\log P(\mathbf{y}_{\leq n} | f)$

Information Consistency

- Likelihood $P(y|u)$. Strategy Q . Competitor space $\mathcal{F}_{comp} \subset \mathcal{F}$. Input distribution $d\mu(\mathbf{x})$.
- Q **information consistent** over \mathcal{F}_{comp} w.r.t. μ iff $\forall f \in \mathcal{F}_{comp}$:

$$\frac{\mathbb{E}_{\mathbf{x}_{\leq n} \sim \mu^n} [D [P(\mathbf{y}_{\leq n}|f, \mathbf{x}_{\leq n}) \parallel Q(\mathbf{y}_{\leq n}|\mathbf{x}_{\leq n})]]}{n} \rightarrow 0 \quad (n \rightarrow \infty)$$

- Link to prediction game:

$$D[P(\mathbf{y}_{\leq n}|f) \parallel P_{bs}(\mathbf{y}_{\leq n})] = \mathbb{E}[L_{bs}(\mathbf{y}_{\leq n}) - L_f(\mathbf{y}_{\leq n})]$$

Uniform (over $\mathbf{y}_{\leq n}$) regret bound implies information consistency

Information Consistency

- Likelihood $P(y|u)$. Strategy Q . Competitor space $\mathcal{F}_{comp} \subset \mathcal{F}$. Input distribution $d\mu(\mathbf{x})$.
- Q **information consistent** over \mathcal{F}_{comp} w.r.t. μ iff $\forall f \in \mathcal{F}_{comp}$:

$$\frac{\mathbb{E}_{\mathbf{x}_{\leq n} \sim \mu^n} [\mathbb{D} [P(\mathbf{y}_{\leq n}|f, \mathbf{x}_{\leq n}) \parallel Q(\mathbf{y}_{\leq n}|\mathbf{x}_{\leq n})]]}{n} \rightarrow 0 \quad (n \rightarrow \infty)$$

- Link to prediction game:

$$\mathbb{D}[P(\mathbf{y}_{\leq n}|f) \parallel P_{bs}(\mathbf{y}_{\leq n})] = \mathbb{E}[L_{bs}(\mathbf{y}_{\leq n}) - L_f(\mathbf{y}_{\leq n})]$$

Uniform (over $\mathbf{y}_{\leq n}$) regret bound implies information consistency

- Likelihood $P(y|u)$. Strategy Q . Competitor space $\mathcal{F}_{comp} \subset \mathcal{F}$. Input distribution $d\mu(\mathbf{x})$.
- Q **information consistent** over \mathcal{F}_{comp} w.r.t. μ iff $\forall f \in \mathcal{F}_{comp}$:

$$\frac{E_{\mathbf{x}_{\leq n} \sim \mu^n} [D [P(\mathbf{y}_{\leq n}|f, \mathbf{x}_{\leq n}) \parallel Q(\mathbf{y}_{\leq n}|\mathbf{x}_{\leq n})]]}{n} \rightarrow 0 \quad (n \rightarrow \infty)$$

- Link to prediction game:

$$D[P(\mathbf{y}_{\leq n}|f) \parallel P_{bs}(\mathbf{y}_{\leq n})] = E[L_{bs}(\mathbf{y}_{\leq n}) - L_f(\mathbf{y}_{\leq n})]$$

Uniform (over $\mathbf{y}_{\leq n}$) regret bound implies information consistency

Information Consistency (II)

- Our result implies:
 - **Gaussian process prediction** (dP_{bs} GP prior with kernel K) is information consistent over reproducing kernel Hilbert space (RKHS) \mathcal{H}_K w.r.t. any μ , for any stationary K with $K(\mathbf{x}, \mathbf{x}) < \infty$
 - \mathcal{H}_K is (usually) infinite dimensional (nonparametric prediction)
 - Our approach allows simple computation of information consistency rate bounds from spectrum of kernel operator
 - **Main result**

$$\underbrace{-\log P_{bs}(\mathbf{y}_{\leq n})}_{L_{bs}(\mathbf{y}_{\leq n})} \leq \underbrace{-\log P(\mathbf{y}_{\leq n}|f)}_{L_f(\mathbf{y}_{\leq n})} + \frac{1}{2} \|f\|_K^2 + \frac{1}{2} \underbrace{\log |I + cK|}_{\mathbb{E}[\cdot]: \text{Regret term}}$$

for any $\mathbf{y}_{\leq n} \in \mathbb{R}^n$ and any $f \in \mathcal{H}_K$.

$\mathbf{K} = (K(\mathbf{x}_i, \mathbf{x}_j))$ kernel matrix

c bounds $-d^2 \log P(y|u)/du^2$

Information Consistency (II)

- Our result implies:
 - **Gaussian process prediction** (dP_{bs} GP prior with kernel K) is information consistent over reproducing kernel Hilbert space (RKHS) \mathcal{H}_K w.r.t. any μ , for any stationary K with $K(\mathbf{x}, \mathbf{x}) < \infty$
 - \mathcal{H}_K is (usually) infinite dimensional (nonparametric prediction)
 - Our approach allows simple computation of information consistency rate bounds from spectrum of kernel operator
 - **Main result**

$$\underbrace{-\log P_{bs}(\mathbf{y}_{\leq n})}_{L_{bs}(\mathbf{y}_{\leq n})} \leq \underbrace{-\log P(\mathbf{y}_{\leq n}|f)}_{L_f(\mathbf{y}_{\leq n})} + \frac{1}{2} \|f\|_K^2 + \frac{1}{2} \underbrace{\log |I + c\mathbf{K}|}_{\mathbb{E}[\cdot]: \text{Regret term}}$$

for any $\mathbf{y}_{\leq n} \in \mathbb{R}^n$ and any $f \in \mathcal{H}_K$.

$\mathbf{K} = (K(\mathbf{x}_i, \mathbf{x}_j))$ kernel matrix

c bounds $-d^2 \log P(y|u)/du^2$

Information Consistency (II)

- Our result implies:
 - **Gaussian process prediction** (dP_{bs} GP prior with kernel K) is information consistent over reproducing kernel Hilbert space (RKHS) \mathcal{H}_K w.r.t. any μ , for any stationary K with $K(\mathbf{x}, \mathbf{x}) < \infty$
 - \mathcal{H}_K is (usually) infinite dimensional (nonparametric prediction)
 - Our approach allows simple computation of information consistency rate bounds from spectrum of kernel operator
 - **Main result**

$$\underbrace{-\log P_{bs}(\mathbf{y}_{\leq n})}_{L_{bs}(\mathbf{y}_{\leq n})} \leq \underbrace{-\log P(\mathbf{y}_{\leq n}|f)}_{L_f(\mathbf{y}_{\leq n})} + \frac{1}{2} \|f\|_K^2 + \frac{1}{2} \underbrace{\log |I + c\mathbf{K}|}_{\mathbb{E}[\cdot]: \text{Regret term}}$$

for any $\mathbf{y}_{\leq n} \in \mathbb{R}^n$ and any $f \in \mathcal{H}_K$.

$\mathbf{K} = (K(\mathbf{x}_i, \mathbf{x}_j))$ kernel matrix

c bounds $-d^2 \log P(y|u)/du^2$

Information Consistency (II)

- Our result implies:
 - **Gaussian process prediction** (dP_{bs} GP prior with kernel K) is information consistent over reproducing kernel Hilbert space (RKHS) \mathcal{H}_K w.r.t. any μ , for any stationary K with $K(\mathbf{x}, \mathbf{x}) < \infty$
 - \mathcal{H}_K is (usually) infinite dimensional (nonparametric prediction)
 - Our approach allows simple computation of information consistency rate bounds from spectrum of kernel operator
 - **Main result**

$$\underbrace{-\log P_{bs}(\mathbf{y}_{\leq n})}_{L_{bs}(\mathbf{y}_{\leq n})} \leq \underbrace{-\log P(\mathbf{y}_{\leq n}|f)}_{L_f(\mathbf{y}_{\leq n})} + \frac{1}{2} \|f\|_K^2 + \frac{1}{2} \underbrace{\log |I + c\mathbf{K}|}_{\mathbb{E}[\cdot]: \text{Regret term}}$$

for any $\mathbf{y}_{\leq n} \in \mathbb{R}^n$ and any $f \in \mathcal{H}_K$.

$\mathbf{K} = (K(\mathbf{x}_i, \mathbf{x}_j))$ kernel matrix

c bounds $-d^2 \log P(y|u)/du^2$

Proof Idea

- First $f \in \mathcal{H}_n := \text{span}\{K(\cdot, \mathbf{x}_i)\}$. **Variational inequality:**

$$\underbrace{-\log P_{bs}(\mathbf{y}_{\leq n})}_{\text{"Free Energy"}} \leq \underbrace{E_Q[-\log P(\mathbf{y}_{\leq n}|u(\cdot))]}_{\text{"Energy"}} + \underbrace{D[Q \| P_{bs}]}_{\text{"Entropy"}} \quad \forall dQ(f)$$

Fenchel Duality. Similar to PAC-Bayesian technique

- Use specific GP Q with $E_Q[u(\cdot)] = f(\cdot)$:

$$\inf_{f \in \mathcal{H}_n} -\log P(\mathbf{y}_{\leq n}|f) + \frac{1}{2} \|f\|_K^2 \geq \underbrace{-\log P_{bs}(\mathbf{y}_{\leq n}) - \frac{1}{2} \log |I + cK|}_{\text{Equality for } P(y|u)=N(y|u, c^{-1})}$$

- $-\log P(\mathbf{y}_{\leq n}|f)$ depends on $f(\mathbf{x}_1), \dots, f(\mathbf{x}_n)$ only
 $\Rightarrow \inf_{\mathcal{H}_n}[\cdot] = \inf_{\mathcal{H}}[\cdot]$ by **representer theorem**

Proof Idea

- First $f \in \mathcal{H}_n := \text{span}\{K(\cdot, \mathbf{x}_i)\}$. **Variational inequality:**

$$\underbrace{-\log P_{bs}(\mathbf{y}_{\leq n})}_{\text{"Free Energy"}} \leq \underbrace{E_Q[-\log P(\mathbf{y}_{\leq n}|u(\cdot))]}_{\text{"Energy"}} + \underbrace{D[Q \| P_{bs}]}_{\text{"Entropy"}} \quad \forall dQ(f)$$

Fenchel Duality. Similar to PAC-Bayesian technique

- Use specific GP Q with $E_Q[u(\cdot)] = f(\cdot)$:

$$\inf_{f \in \mathcal{H}_n} -\log P(\mathbf{y}_{\leq n}|f) + \frac{1}{2} \|f\|_K^2 \geq \underbrace{-\log P_{bs}(\mathbf{y}_{\leq n}) - \frac{1}{2} \log |I + cK|}_{\text{Equality for } P(y|u)=N(y|u, c^{-1})}$$

- $-\log P(\mathbf{y}_{\leq n}|f)$ depends on $f(\mathbf{x}_1), \dots, f(\mathbf{x}_n)$ only
 $\Rightarrow \inf_{\mathcal{H}_n}[\cdot] = \inf_{\mathcal{H}}[\cdot]$ by **representer theorem**

Proof Idea

- First $f \in \mathcal{H}_n := \text{span}\{K(\cdot, \mathbf{x}_i)\}$. **Variational inequality:**

$$\underbrace{-\log P_{bs}(\mathbf{y}_{\leq n})}_{\text{"Free Energy"}} \leq \underbrace{E_Q[-\log P(\mathbf{y}_{\leq n}|u(\cdot))]}_{\text{"Energy"}} + \underbrace{D[Q \| P_{bs}]}_{\text{"Entropy"}} \quad \forall dQ(f)$$

Fenchel Duality. Similar to PAC-Bayesian technique

- Use specific GP Q with $E_Q[u(\cdot)] = f(\cdot)$:

$$\inf_{f \in \mathcal{H}_n} -\log P(\mathbf{y}_{\leq n}|f) + \frac{1}{2} \|f\|_K^2 \geq \underbrace{-\log P_{bs}(\mathbf{y}_{\leq n}) - \frac{1}{2} \log |I + cK|}_{\text{Equality for } P(y|u)=N(y|u, c^{-1})}$$

- $-\log P(\mathbf{y}_{\leq n}|f)$ depends on $f(\mathbf{x}_1), \dots, f(\mathbf{x}_n)$ only
 $\Rightarrow \inf_{\mathcal{H}_n}[\cdot] = \inf_{\mathcal{H}}[\cdot]$ by **representer theorem**

The Regret Term

- $R = \log |\mathbf{I} + c\mathbf{K}|$. Information consistency if $E[R]/n \rightarrow 0$

$$R = \log \left| \mathbf{I} + cn\mathbf{U}\hat{\Lambda}\mathbf{U}^T \right| = \sum_{i=1}^n \log(1 + cn\hat{\lambda}_i)$$

$$\hat{\lambda}_1 \geq \hat{\lambda}_2 \geq \dots \geq \hat{\lambda}_n. \hat{\lambda}_j = 0 \text{ for } j > n$$

- Stationary kernel $K(\mathbf{x}_i, \mathbf{x}_j) = K(\mathbf{x}_i - \mathbf{x}_j)$, probability $d\mu$:
Kernel operator on $L_2(d\mu)$, discrete spectrum $\lambda_1 \geq \lambda_2 \geq \dots$
Summable: $\sum_{i \geq 1} \lambda_i = \int K(\mathbf{x}, \mathbf{x}) d\mu(\mathbf{x}) = K(\mathbf{0}) < \infty$
- Split sum into two parts:

$$E[R] = S_1 + S_2, \quad S_1 = E \left[\sum_{i=1}^{s_0} \log(1 + cn\hat{\lambda}_i) \right] = O(s_0 \log n),$$

$$S_2 = E \left[\sum_{i > s_0} \log(1 + cn\hat{\lambda}_i) \right] \leq \underbrace{cn E \left[\sum_{i > s_0} \hat{\lambda}_i \right]}_{\log(1+x) \leq x}$$

The Regret Term

- $R = \log |\mathbf{I} + c\mathbf{K}|$. Information consistency if $E[R]/n \rightarrow 0$

$$R = \log \left| \mathbf{I} + cn\mathbf{U}\hat{\Lambda}\mathbf{U}^T \right| = \sum_{i=1}^n \log(1 + cn\hat{\lambda}_i)$$

$$\hat{\lambda}_1 \geq \hat{\lambda}_2 \geq \dots \geq \hat{\lambda}_n. \hat{\lambda}_j = 0 \text{ for } j > n$$

- Stationary kernel $K(\mathbf{x}_i, \mathbf{x}_j) = K(\mathbf{x}_i - \mathbf{x}_j)$, probability $d\mu$:
Kernel operator on $L_2(d\mu)$, discrete spectrum $\lambda_1 \geq \lambda_2 \geq \dots$
Summable: $\sum_{i \geq 1} \lambda_i = \int K(\mathbf{x}, \mathbf{x}) d\mu(\mathbf{x}) = K(\mathbf{0}) < \infty$
- Split sum into two parts:

$$E[R] = S_1 + S_2, \quad S_1 = E \left[\sum_{i=1}^{s_0} \log(1 + cn\hat{\lambda}_i) \right] = O(s_0 \log n),$$

$$S_2 = E \left[\sum_{i > s_0} \log(1 + cn\hat{\lambda}_i) \right] \leq \underbrace{cn E \left[\sum_{i > s_0} \hat{\lambda}_i \right]}_{\log(1+x) \leq x}$$

The Regret Term

- $R = \log |\mathbf{I} + c\mathbf{K}|$. Information consistency if $E[R]/n \rightarrow 0$

$$R = \log \left| \mathbf{I} + cn\mathbf{U}\hat{\Lambda}\mathbf{U}^T \right| = \sum_{i=1}^n \log(1 + cn\hat{\lambda}_i)$$

$$\hat{\lambda}_1 \geq \hat{\lambda}_2 \geq \dots \geq \hat{\lambda}_n. \quad \hat{\lambda}_j = 0 \text{ for } j > n$$

- Stationary kernel $K(\mathbf{x}_i, \mathbf{x}_j) = K(\mathbf{x}_i - \mathbf{x}_j)$, probability $d\mu$:
Kernel operator on $L_2(d\mu)$, discrete spectrum $\lambda_1 \geq \lambda_2 \geq \dots$
Summable: $\sum_{i \geq 1} \lambda_i = \int K(\mathbf{x}, \mathbf{x}) d\mu(\mathbf{x}) = K(\mathbf{0}) < \infty$
- Split sum into two parts:

$$E[R] = S_1 + S_2, \quad S_1 = E \left[\sum_{i=1}^{s_0} \log(1 + cn\hat{\lambda}_i) \right] = O(s_0 \log n),$$

$$S_2 = E \left[\sum_{i > s_0} \log(1 + cn\hat{\lambda}_i) \right] \leq \underbrace{cn E \left[\sum_{i > s_0} \hat{\lambda}_i \right]}_{\log(1+x) \leq x}$$

The Regret Term (II)

- Empirical eigenvalues are overestimates

Shawe-Taylor *et.al.*, IEEE IT 05

$$\sum_{i \geq 1} \hat{\lambda}_i = n^{-1} \text{tr} \mathbf{K} = K(\mathbf{0}) = \sum_{i \geq 1} \lambda_i,$$

$$\sum_{i=1}^{s_0} \hat{\lambda}_i \geq \sum_{i=1}^{s_0} \lambda_i$$

- For any s_0 :

$$S_2 \leq cn \sum_{i > s_0} \hat{\lambda}_i \leq cn \sum_{i > s_0} \lambda_i$$

- All in all: For some $C \approx 1$:

$$E[R]/n = (S_1 + S_2)/n \leq Cs_0(\log n)/n + c \sum_{i > s_0} \lambda_i$$

Choose $s_0(n) \rightarrow \infty$, $s_0(n)(\log n)/n \rightarrow 0$.

\Rightarrow Information consistency

Thanks to Ingo Steinwart

The Regret Term (II)

- Empirical eigenvalues are overestimates

Shawe-Taylor *et.al.*, IEEE IT 05

$$\sum_{i \geq 1} \hat{\lambda}_i = n^{-1} \text{tr} \mathbf{K} = K(\mathbf{0}) = \sum_{i \geq 1} \lambda_i,$$

$$\sum_{i=1}^{s_0} \hat{\lambda}_i \geq \sum_{i=1}^{s_0} \lambda_i$$

- For any s_0 :

$$S_2 \leq cn \sum_{i > s_0} \hat{\lambda}_i \stackrel{!}{\leq} cn \sum_{i > s_0} \lambda_i$$

- All in all: For some $C \approx 1$:

$$E[R]/n = (S_1 + S_2)/n \leq Cs_0(\log n)/n + c \sum_{i > s_0} \lambda_i$$

Choose $s_0(n) \rightarrow \infty$, $s_0(n)(\log n)/n \rightarrow 0$.

\Rightarrow Information consistency

Thanks to Ingo Steinwart

The Regret Term (II)

- Empirical eigenvalues are overestimates

Shawe-Taylor *et.al.*, IEEE IT 05

$$\sum_{i \geq 1} \hat{\lambda}_i = n^{-1} \text{tr} \mathbf{K} = K(\mathbf{0}) = \sum_{i \geq 1} \lambda_i,$$

$$\sum_{i=1}^{s_0} \hat{\lambda}_i \geq \sum_{i=1}^{s_0} \lambda_i$$

- For any s_0 :

$$\mathcal{S}_2 \leq cn \sum_{i > s_0} \hat{\lambda}_i \stackrel{!}{\leq} cn \sum_{i > s_0} \lambda_i$$

- All in all: For some $C \approx 1$:

$$E[R]/n = (\mathcal{S}_1 + \mathcal{S}_2)/n \leq Cs_0(\log n)/n + c \sum_{i > s_0} \lambda_i$$

Choose $s_0(n) \rightarrow \infty$, $s_0(n)(\log n)/n \rightarrow 0$.

\Rightarrow **Information consistency**

Thanks to Ingo Steinwart

Information Consistency Rates

- Is R **concentrated** ($\text{Var}[R] \leq \text{E}[R]$)?
 - $\text{E}[R]$ can be small. Standard global approach (McDiarmid) not strong enough
 - Local results on $\sum_{i>s_0} \hat{\lambda}_i$ could be useful Zwald *et.al.*, Mach. Learn. 07
- For fixed stationary K, μ :
 Bounds on $\sum_{i>s_0} \lambda_i \Rightarrow$ Information consistency rate bounds
- **Gaussian** $K(\mathbf{r}) = \exp(-b\|\mathbf{r}\|^2)$, **Gaussian** $\mu(\mathbf{x}) = N(\mathbf{x}|\mathbf{0}, (4a)^{-1}I)$,
 $\mathbf{r}, \mathbf{x} \in \mathbb{R}^d$: Spectrum known explicitly Zhu *et.al.*, NATO ASI 98

$$\text{E}[R] = O\left((\log n)^{d+1}\right)$$

Very small (infinite dimensional \mathcal{H}_K):
 Smoothness through kernel for nonparametric statistic \leftrightarrow
 dimensionality for parametric statistics.
 \Rightarrow Gaussian kernel **extremely smoothing**

Information Consistency Rates

- Is R **concentrated** ($\text{Var}[R] \leq \text{E}[R]$)?
 - $\text{E}[R]$ can be small. Standard global approach (McDiarmid) not strong enough
 - Local results on $\sum_{i>s_0} \hat{\lambda}_i$ could be useful Zwald *et.al.*, Mach. Learn. 07
- For fixed stationary K, μ :
 Bounds on $\sum_{i>s_0} \lambda_i \Rightarrow$ Information consistency rate bounds
- **Gaussian** $K(\mathbf{r}) = \exp(-b\|\mathbf{r}\|^2)$, **Gaussian** $\mu(\mathbf{x}) = N(\mathbf{x}|\mathbf{0}, (4a)^{-1}I)$,
 $\mathbf{r}, \mathbf{x} \in \mathbb{R}^d$: Spectrum known explicitly Zhu *et.al.*, NATO ASI 98

$$\text{E}[R] = O\left((\log n)^{d+1}\right)$$

Very small (infinite dimensional \mathcal{H}_K):
 Smoothness through kernel for nonparametric statistic \leftrightarrow
 dimensionality for parametric statistics.
 \Rightarrow Gaussian kernel **extremely smoothing**

Information Consistency Rates

- Is R **concentrated** ($\text{Var}[R] \leq \text{E}[R]$)?
 - $\text{E}[R]$ can be small. Standard global approach (McDiarmid) not strong enough
 - Local results on $\sum_{i>s_0} \hat{\lambda}_i$ could be useful Zwald *et.al.*, Mach. Learn. 07
- For fixed stationary K, μ :
 Bounds on $\sum_{i>s_0} \lambda_i \Rightarrow$ Information consistency rate bounds
- **Gaussian** $K(\mathbf{r}) = \exp(-b\|\mathbf{r}\|^2)$, **Gaussian** $\mu(\mathbf{x}) = N(\mathbf{x}|\mathbf{0}, (4a)^{-1}\mathbf{I})$,
 $\mathbf{r}, \mathbf{x} \in \mathbb{R}^d$: Spectrum known explicitly Zhu *et.al.*, NATO ASI 98

$$\text{E}[R] = O\left((\log n)^{d+1}\right)$$

Very small (infinite dimensional \mathcal{H}_K):

Smoothness through kernel for nonparametric statistic \leftrightarrow
 dimensionality for parametric statistics.

\Rightarrow Gaussian kernel **extremely smoothing**

Information Consistency Rates

- Is R **concentrated** ($\text{Var}[R] \leq \text{E}[R]$)?
 - $\text{E}[R]$ can be small. Standard global approach (McDiarmid) not strong enough
 - Local results on $\sum_{i>s_0} \hat{\lambda}_i$ could be useful Zwald *et.al.*, Mach. Learn. 07
- For fixed stationary K, μ :
 Bounds on $\sum_{i>s_0} \lambda_i \Rightarrow$ Information consistency rate bounds
- **Gaussian** $K(\mathbf{r}) = \exp(-b\|\mathbf{r}\|^2)$, **Gaussian** $\mu(\mathbf{x}) = N(\mathbf{x}|\mathbf{0}, (4a)^{-1}\mathbf{I})$,
 $\mathbf{r}, \mathbf{x} \in \mathbb{R}^d$: Spectrum known explicitly Zhu *et.al.*, NATO ASI 98

$$\text{E}[R] = O\left((\log n)^{d+1}\right)$$

Very small (infinite dimensional \mathcal{H}_K):

Smoothness through kernel for nonparametric statistic \leftrightarrow
 dimensionality for parametric statistics.

\Rightarrow Gaussian kernel **extremely smoothing**

Information Consistency Rates (II)

- Widom (Trans. AMS 63): Asymptotic expressions for λ_i if
 - Spectral density (Fourier transform) of K decays polynomially
 - $d\mu$ has density with bounded support
- **Matérn class**: Kernels with Student's t spectral densities: Ornstein-Uhlenbeck ($\nu = 1/2$) \leftrightarrow Gaussian ($\nu \rightarrow \infty$)
- K **Matérn** (d.o.f. $\nu > 0$), $\mu(\mathbf{x})$ **bounded support**, $\mathbf{x} \in \mathbb{R}^d$:

$$E[R] = O\left(n^{d/(2\nu+d)}(\log n)^{2\nu/(2\nu+d)}\right)$$

Drawback: Constant as large as $|\text{supp}\mu|^{2\nu+d}$

- K **Matérn** (d.o.f. ν), $\tilde{\mu}$ **Student's t** (d.o.f. $\nu_2 > \nu$), $\mu(\mathbf{x}) \propto \tilde{\mu}(\mathbf{x})\mathbf{I}_{\{\|\mathbf{x}\| \leq T\}}$:

$$E[R] = O\left(n^{d/(2\nu+d)}(\log n)^{2\nu/(2\nu+d)}\right)$$

Constant **independent** of T .

But speed of convergence may depend on T

Information Consistency Rates (II)

- Widom (Trans. AMS 63): Asymptotic expressions for λ_i if
 - Spectral density (Fourier transform) of K decays polynomially
 - $d\mu$ has density with bounded support
- Matérn class**: Kernels with Student's t spectral densities: Ornstein-Uhlenbeck ($\nu = 1/2$) \leftrightarrow Gaussian ($\nu \rightarrow \infty$)
- K **Matérn** (d.o.f. $\nu > 0$), $\mu(\mathbf{x})$ **bounded support**, $\mathbf{x} \in \mathbb{R}^d$:

$$E[R] = O\left(n^{d/(2\nu+d)}(\log n)^{2\nu/(2\nu+d)}\right)$$

Drawback: Constant as large as $|\text{supp}\mu|^{2\nu+d}$

- K **Matérn** (d.o.f. ν), $\tilde{\mu}$ **Student's t** (d.o.f. $\nu_2 > \nu$), $\mu(\mathbf{x}) \propto \tilde{\mu}(\mathbf{x})\mathbf{I}_{\{\|\mathbf{x}\| \leq T\}}$:

$$E[R] = O\left(n^{d/(2\nu+d)}(\log n)^{2\nu/(2\nu+d)}\right)$$

Constant **independent** of T .

But speed of convergence may depend on T

Information Consistency Rates (II)

- Widom (Trans. AMS 63): Asymptotic expressions for λ_i if
 - Spectral density (Fourier transform) of K decays polynomially
 - $d\mu$ has density with bounded support
- Matérn class**: Kernels with Student's t spectral densities: Ornstein-Uhlenbeck ($\nu = 1/2$) \leftrightarrow Gaussian ($\nu \rightarrow \infty$)
- K **Matérn** (d.o.f. $\nu > 0$), $\mu(\mathbf{x})$ **bounded support**, $\mathbf{x} \in \mathbb{R}^d$:

$$E[R] = O\left(n^{d/(2\nu+d)}(\log n)^{2\nu/(2\nu+d)}\right)$$

Drawback: Constant as large as $|\text{supp}\mu|^{2\nu+d}$

- K **Matérn** (d.o.f. ν), $\tilde{\mu}$ **Student's t** (d.o.f. $\nu_2 > \nu$),
 $\mu(\mathbf{x}) \propto \tilde{\mu}(\mathbf{x})\mathbf{I}_{\{\|\mathbf{x}\| \leq T\}}$:

$$E[R] = O\left(n^{d/(2\nu+d)}(\log n)^{2\nu/(2\nu+d)}\right)$$

Constant **independent** of T .

But speed of convergence may depend on T

Information Consistency Rates (II)

- Widom (Trans. AMS 63): Asymptotic expressions for λ_i if
 - Spectral density (Fourier transform) of K decays polynomially
 - $d\mu$ has density with bounded support
- Matérn class**: Kernels with Student's t spectral densities: Ornstein-Uhlenbeck ($\nu = 1/2$) \leftrightarrow Gaussian ($\nu \rightarrow \infty$)
- K **Matérn** (d.o.f. $\nu > 0$), $\mu(\mathbf{x})$ **bounded support**, $\mathbf{x} \in \mathbb{R}^d$:

$$E[R] = O\left(n^{d/(2\nu+d)}(\log n)^{2\nu/(2\nu+d)}\right)$$

Drawback: Constant as large as $|\text{supp}\mu|^{2\nu+d}$

- K **Matérn** (d.o.f. ν), $\tilde{\mu}$ **Student's t** (d.o.f. $\nu_2 > \nu$), $\mu(\mathbf{x}) \propto \tilde{\mu}(\mathbf{x})\mathbf{I}_{\{\|\mathbf{x}\| \leq T\}}$:

$$E[R] = O\left(n^{d/(2\nu+d)}(\log n)^{2\nu/(2\nu+d)}\right)$$

Constant **independent** of T .

But speed of convergence may depend on T

Conclusions

- Nonparametric information consistency of GP prediction for wide range of stationary kernels
- Information consistency rates reduced to kernel operator spectrum (entropy and covering numbers can be used; I. Steinwart)
- Nonparametrics: Smoothness more important than dimensionality
- Nonparametric MDL beyond Shtarkov Kakade *et.al.*, NIPS 06; Grünwald 07, Sect. 13.5
- Open problems
 - Concentration of regret term
 - Regret term for unbounded support $d\mu$

Conclusions

- Nonparametric information consistency of GP prediction for wide range of stationary kernels
- Information consistency rates reduced to kernel operator spectrum (entropy and covering numbers can be used; I. Steinwart)
- Nonparametrics: Smoothness more important than dimensionality
- Nonparametric MDL beyond Shtarkov Kakade *et.al.*, NIPS 06; Grünwald 07, Sect. 13.5
- Open problems
 - Concentration of regret term
 - Regret term for unbounded support $d\mu$

Conclusions

- Nonparametric information consistency of GP prediction for wide range of stationary kernels
- Information consistency rates reduced to kernel operator spectrum (entropy and covering numbers can be used; I. Steinwart)
- Nonparametrics: Smoothness more important than dimensionality
- Nonparametric MDL beyond Shtarkov Kakade *et.al.*, NIPS 06; Grünwald 07, Sect. 13.5
- Open problems
 - Concentration of regret term
 - Regret term for unbounded support $d\mu$

Conclusions

- Nonparametric information consistency of GP prediction for wide range of stationary kernels
- Information consistency rates reduced to kernel operator spectrum (entropy and covering numbers can be used; I. Steinwart)
- Nonparametrics: Smoothness more important than dimensionality
- Nonparametric MDL beyond Shtarkov Kakade *et.al.*, NIPS 06; Grünwald 07, Sect. 13.5
- Open problems
 - Concentration of regret term
 - Regret term for unbounded support $d\mu$

Conclusions

- Nonparametric information consistency of GP prediction for wide range of stationary kernels
- Information consistency rates reduced to kernel operator spectrum (entropy and covering numbers can be used; I. Steinwart)
- Nonparametrics: Smoothness more important than dimensionality
- Nonparametric MDL beyond Shtarkov Kakade *et.al.*, NIPS 06; Grünwald 07, Sect. 13.5
- Open problems
 - Concentration of regret term
 - Regret term for unbounded support $d\mu$