Approximation Properties of Reproducing Kernels

M. Eberts (Meister), S. Fischer, N. Schmid, P. Thomann, & I. Steinwart

Institute of Stochastics and Applications University of Stuttgart

> Guangzhou May 20th, 2017

- Approximation theoretic questions related to kernel-based learning
- More flexible kernels: spatial decompositions
- More flexible kernels: deeper compositions

- X space of input samples Y space of labels, usually $Y \subset \mathbb{R}$.
- Already observed samples

$$D = ((x_1, y_1), \ldots, (x_n, y_n)) \in (X \times Y)^n$$

- X space of input samples Y space of labels, usually $Y \subset \mathbb{R}$.
- Already observed samples

$$D = ((x_1, y_1), \ldots, (x_n, y_n)) \in (X \times Y)^n$$

Goal:

With the help of D find a function $f_D : X \to \mathbb{R}$ such that $f_D(x)$ is a good prediction of the label y for new, unseen x.

• Learning method:

Assigns to every training set *D* a predictor $f_D : X \to \mathbb{R}$.

Problem:

The labels y are \mathbb{R} -valued.

Goal:

Estimate label y for new data x as accurate as possible.

Example:



Assumptions

- We have bounded labels Y = [-1, 1].
- *P* is an unknown probability measure on $X \times [-1, 1]$.
- $D = ((x_1, y_1), \dots, (x_n, y_n)) \in (X \times Y)^n$ is sampled from P^n .
- Future samples (x, y) will also be sampled from *P*.
- (For this talk) we mostly use the least squares loss

$$L(y,t) := (y-t)^2$$

to assess quality of a prediction t for y.

• The risk of a predictor $f: X \to \mathbb{R}$ is the average loss

$$\mathcal{R}_{L,P}(f) := \int_{X \times Y} L(y, f(x)) dP(x, y) .$$

• The Bayes risk is the smallest possible risk

$$\mathcal{R}^*_{L,P}:= \infig\{ \, \mathcal{R}_{L,P}(f) \mid f: X o \mathbb{R} \; (\mathsf{measurable}) \; ig\} \; .$$

• The Bayes predictor for the least squares loss is $f_{L,P}^*(x) := \mathbb{E}(Y|x)$, i.e.

$$\mathcal{R}_{L,P}(f_{L,P}^*) = \mathcal{R}_{L,P}^*$$

• The excess risk satisfies

$$\mathcal{R}_{L,P}(f) - \mathcal{R}^*_{L,P} = \|f - f^*_{L,P}\|^2_{L_2(P_X)}.$$

Kernel-based learning methods

- Let *H* be a reproducing kernel Hilbert space, here with bounded kernel
- Let $L: Y \times \mathbb{R} \to [0,\infty)$ be a *convex* loss

Kernel-based learning methods

- Let H be a reproducing kernel Hilbert space, here with bounded kernel
- Let $L: Y \times \mathbb{R} \to [0,\infty)$ be a *convex* loss
- Kernel-based learning methods (e.g. SVMs) solve the problem

$$f_{D,\lambda} = \arg\min_{f \in H} \lambda \|f\|_{H}^{2} + \frac{1}{n} \sum_{i=1}^{n} L(y_{i}, f(x_{i})) , \qquad (1)$$

where $\lambda > 0$ is a free regularization parameter. Solution is of the form

$$f_{D,\lambda} = \sum_{i=1}^n \alpha_i k(x_i, \cdot).$$

Historical Notes

- G. Wahba (1971 –): Least squares loss
- V. Vapnik et al. (1992 –): Hinge loss
- Other losses in the last decade or so.

A Typical Oracle Inequality

• Consider the approximation (regularization) error

$$A(\lambda) := \inf_{f \in H} \lambda \|f\|_{H}^{2} + \mathcal{R}_{L,P}(f) - \mathcal{R}_{L,P}^{*}(f)$$

• Assume an (dyadic) entropy number behavior

$$e_i(I:H\to L_2(P_X)) \preceq i^{-1/(2p)}$$

Then with probability P^n not smaller than $1 - e^{-\tau}$ we have

$$\mathcal{R}_{L,P}(f_{D,\lambda}) - \mathcal{R}^*_{L,P} \leq K\Big(A(\lambda) + \frac{1}{\lambda^p n} + \frac{\tau A(\lambda)}{\lambda n}\Big)$$

Remarks:

- If rate $A(\lambda) \rightarrow 0$ for $\lambda \rightarrow 0$ known, we obtain learning rates.
- Entropy behaviour is equivalent to a similar eigenvalue behaviour of

$$egin{aligned} T_k &: L_2(P_X) o L_2(P_X) \ &T_k f &:= \int_X k(x,\cdot) f(x) \, dP_X(x) \end{aligned}$$

• For Banach spaces $F \hookrightarrow E$ and $x \in E$, the K-functional is

$$K(x,t) := \inf_{y \in F} ||x - y||_E + t ||y||_F, \qquad t > 0.$$

For 0 < β < 1, 1 ≤ r ≤ ∞, the interpolation space [E, F]_{β,r} consists of those x ∈ E with finite ||x||_{β,r}, where

$$\|x\|_{\beta,r} := \begin{cases} \left(\int_{0}^{\infty} (t^{-\beta} K(x,t))^{r} t^{-1} dt\right)^{1/r} & \text{if } 1 \leq r < \infty \\ \sup_{t>0} t^{-\beta} K(x,t) & \text{if } r = \infty \,. \end{cases}$$

• We are interested in the spaces $[L_2(P_X), [H]_{\sim}]_{\beta,r}$.

Interpolation Spaces vs. Approximation Properties

Smale & Zhou, 2003 $A(\lambda) \preceq \lambda^{\beta}$ if and only if $f_{L,P}^* \in [L_2(P_X), [H]_{\sim}]_{\beta,\infty}$.

Operator techniques (Caponnetto and De Vito, 2007, ...) Rates for $\mathcal{R}_{L,P}(f_{D,\lambda}) \rightarrow \mathcal{R}^*_{L,P}$ are obtained if

$$f^*_{L,P} \in \operatorname{im} T^{\beta/2}_k$$

Smale & Zhou, 2003 If X is compact, supp $P_X = X$ and k continuous, then

$$\left[L_2(P_X),[H]_{\sim}\right]_{\beta+\varepsilon,\infty}\,\subset\,\operatorname{im}\,T_k^{\beta/2}\,\subset\,\left[L_2(P_X),[H]_{\sim}\right]_{\beta,\infty},$$

Both approximation assumptions are almost the same, since

$$[L_2(P_X), [H]_{\sim}]_{\beta+\varepsilon,\infty} \hookrightarrow [L_2(P_X), [H]_{\sim}]_{\beta,1} \hookrightarrow [L_2(P_X), [H]_{\sim}]_{\beta,\infty}$$

Spectral Theorem, Revisited

- Let k be a reproducing kernel with compact $I_{k,\nu}: H \to L_2(\nu)$.
- Then $T_{k,\nu} = I_{k,\nu} \circ I_{k,\nu}^*$ is selfadjoint, positive, and compact.
- Let (µ_i)_{i∈I} be the family of non-zero eigenvalues of T_{k,ν} and ([ẽ_i]_∼) be a corresponding ONS of eigenfunctions in L₂(ν).

Then $e_i := \mu_i^{-1} I_{k,\nu}^* [\tilde{e}_i]_{\sim} \in H$ satisfies $[e_i]_{\sim} = [\tilde{e}_i]_{\sim}$ and we have:

Spectral Theorem, Revisited

- Let k be a reproducing kernel with compact $I_{k,\nu}: H \to L_2(\nu)$.
- Then $T_{k,\nu} = I_{k,\nu} \circ I_{k,\nu}^*$ is selfadjoint, positive, and compact.
- Let (µ_i)_{i∈I} be the family of non-zero eigenvalues of T_{k,ν} and ([ẽ_i]_∼) be a corresponding ONS of eigenfunctions in L₂(ν).

Then $e_i := \mu_i^{-1} I_{k,\nu}^* [\tilde{e}_i]_{\sim} \in H$ satisfies $[e_i]_{\sim} = [\tilde{e}_i]_{\sim}$ and we have:

- $([e_i]_{\sim})$ is an ONS in $L_2(\nu)$.
- $(\sqrt{\mu_i}e_i)$ is an ONS in *H*.

$$(\ker I_{k,\nu})^{\perp} = \overline{\operatorname{im} I_{k,\nu}^*} = \overline{\operatorname{span}}\{\sqrt{\mu_i}e_i : i \in I\}$$
$$(\ker T_{k,\nu})^{\perp} = (\ker I_{k,\nu}^*)^{\perp} = \overline{\operatorname{im} I_{k,\nu}} = \overline{\operatorname{span}}\{[e_i]_{\sim} : i \in I\}$$

Consequence.

 $L_2(\nu)$ and H "share a subspace" described by (e_i) .

Isometric Copy of *H* **in** $L_2(\nu)$

$$[H]_{\sim} = \left\{ \sum_{i \in I} a_i \mu_i^{1/2} [e_i]_{\sim} : (a_i) \in \ell_2(I) \right\}$$

Closure of *H* in $L_2(\nu)$

$$\overline{[H]}_{\sim}^{L_2(\nu)} = \left\{ \sum_{i \in I} a_i [e_i]_{\sim} : (a_i) \in \ell_2(I) \right\}$$

Question What is in between?

Power Spaces

• For $\beta \in [0,1]$ we can consider the following subspace of $L_2(\nu)$:

$$\begin{split} [H]^{\beta}_{\sim} &:= \left\{ \sum_{i \in I} a_i \mu_i^{\beta/2} [e_i]_{\sim} : (a_i) \in \ell_2(I) \right\} \\ &= \left\{ \sum_{i \in I} b_i [e_i]_{\sim} : (b_i) \in \ell_2(\mu^{-\beta}) \right\}, \end{split}$$

where $\ell_2(\mu^{-\beta})$ is a weighted sequence space with norm: $\|(b_i)\|_{\ell_2(\mu^{-\beta})}^2 := \sum_{i \in I} b_i^2 \mu_i^{-\beta}$

Power Spaces

• For $\beta \in [0,1]$ we can consider the following subspace of $L_2(\nu)$:

$$\begin{split} [H]^{\beta}_{\sim} &:= \left\{ \sum_{i \in I} a_i \mu_i^{\beta/2} [e_i]_{\sim} : (a_i) \in \ell_2(I) \right\} \\ &= \left\{ \sum_{i \in I} b_i [e_i]_{\sim} : (b_i) \in \ell_2(\mu^{-\beta}) \right\}, \end{split}$$

where $\ell_2(\mu^{-\beta})$ is a weighted sequence space with norm:

$$\|(b_i)\|^2_{\ell_2(\mu^{-\beta})} := \sum_{i \in I} b_i^2 \mu_i^{-\beta}$$

• By construction, $(\mu_i^{\beta/2}[e_i]_{\sim})_{i\in I}$ is an ONB of $[H]_{\sim}^{\beta}$ and

$$[H]^{0}_{\sim} = \overline{[H]}^{L}_{\sim 2}(\nu)$$
$$[H]^{1}_{\sim} = [H]_{\sim}$$
$$[H]^{\beta}_{\sim} = \operatorname{im} T^{\beta/2}_{k,\nu}$$

S. & Scovel, 2012

If $I_{k,\nu}: H \to L_2(\nu)$ is compact, then, for $\beta \in (0,1)$, we have

im
$$T_{k,\nu}^{\beta/2} = [H]_{\sim}^{\beta} \cong [L_2(\nu), [H]_{\sim}]_{\beta,2}$$

S. & Scovel, 2012 If $I_{k,\nu}: H \to L_2(\nu)$ is compact, then, for $\beta \in (0,1)$, we have

im
$$T_{k,\nu}^{\beta/2} = [H]_{\sim}^{\beta} \cong [L_2(\nu), [H]_{\sim}]_{\beta,2}.$$

Idea of the Proof.

- Interpolating $L_2(\nu)$ and $[H]_{\sim}$ is the same as interpolating $\ell_2(I)$ and $\ell_2(\mu^{-1})$.
- We have $[\ell_2(I), \ell_2(\mu^{-1})]_{\beta,2} \cong \ell_2(\mu^{-\beta}).$

Rates for Fixed Kernel

Generic Setting (S., Scovel, & Hush, 2009 + S. & Scovel, 2012)

- Assume $\mu_i \preceq i^{-1/p}$
- Assume $f_{L,P}^* \in [L_2(P_X), H]_{\beta,\infty}$ for some $\beta \in (0, 1]$.
- Assume $[L_2(P_X), H]_{s,1} \hookrightarrow L_{\infty}(P_X)$ for $s = \min\{1, p/(1-\beta)\}$. This is equivalent to

$$\|f\|_{\infty} \leq c \|f\|_{H}^{s} \|f\|_{L_{2}(P_{X})}^{1-s}, \qquad f \in H$$

Then kernel method can learn with the optimal rate $n^{-\frac{\beta}{\beta+p}}$.

Special Case: Sobolev Setting (e.g. Kohler)

- X ball in ℝ^d and H := W^m(X) Sobolev space with m > d/2.
 → Least squares with splines.
- P_X uniform distribution and $f^*_{L,P} \in B^s_{2,2}(X)$ for some $s \in (d/2, m]$.

The kernel method can learn with the optimal rate $n^{-\frac{2s}{2s+d}}$

Improved Convergence

Fischer & S., 2017

- Assume $\mu_i \preceq i^{-1/p}$
- Assume $f_{L,P}^* \in [L_2(P_X), H]_{\beta,2}$ for some $\beta \in (0, 1]$.
- Assume $[L_2(P_X), H]_{\alpha,2} \hookrightarrow L_{\infty}(P_X)$ for some $\alpha \in (0, 1)$.

Then, for a suitable sequence (λ_n) the decision functions f_{D,λ_n} converges to $f_{L,P}^*$ in the norm of $[L_2(P_X), H]_{\gamma,2}$ for $\gamma \in [0, \beta]$ with rate n^{-r} , where

$$r = rac{eta - \gamma}{\max\{lpha, eta\} + p}$$

Example.

Let $H = W^m(X)$ and $f_{L,P}^* \in B_{2,2}^s(X)$ for some $s \in (d/2, m]$. For $t \in (0, s)$, the rate in $B_{2,2}^t(X)$ is n^{-r} , where

$$r = \frac{2s - 2t}{2s + d}$$

This improves and generalizes results by Smale & Zhou (2007), Capaonetto & de Vito (2007), S. et al (2009), and Blanchard & Mücke (2016)

Smale & Zhou, 2003 Consider Gaussian RKHS $H_{\gamma}(X)$ with kernel

$$k_{\gamma}(x,x') := \exp(-\gamma^{-2}||x-x'||_2^2), \qquad x,x' \in X.$$

Then $A_{\gamma}(\lambda) \preceq \lambda^{\beta}$ for some $\beta \in (0,1]$ implies $f_{L,P}^* \in C^{\infty}(X)$.

Solution

Consider width γ as a free parameter. \rightsquigarrow Theory presented so far does not work anymore.

Eberts & S., 2011/3

- X ball in \mathbb{R}^d and H_γ is RKHS of Gaussian kernel k_γ .
- P_X has bounded Lebesgue density.
- Pick λ and γ by a training/validation approach.

Then, for $s \ge 1$, every $f_{L,P}^* \in W_2^s(X)$ is learned with the rate $n^{-\frac{2s}{2s+d}+\varepsilon}$ without knowing s.

The extra factor n^{ε} can be replaced by a logarithmic factor.

Key idea of the proof

Bound approximation error by convoluting $f_{L,P}^*$ with weighted sum of kernels $k_{\gamma_1}, \ldots k_{\gamma_m}$.

Spatial Decompositions

Optimization Problem

$$f_{D,\lambda} = \arg\min_{f \in H} \lambda \|f\|_{H}^{2} + \frac{1}{n} \sum_{i=1}^{n} L(y_{i}, f(x_{i}))$$

Example: Dual Problem for Hinge Loss

$$\alpha^* \in \arg \max_{\alpha \in [0, \frac{1}{2\lambda_n}]^d} \sum_{i=1}^n \alpha_i - \frac{1}{2} \sum_{i,j=1}^n y_i y_j \alpha_i \alpha_j k(x_i, x_j)$$

Re-substitution

$$f_{D,\lambda} = \sum_{i=1}^{n} y_i \alpha_i^* k(\cdot, x_i),$$

Computational Requirements

- The size of the optimization problem grows linearly in *n*.
- The kernel matrix $(k(x_i, x_j))$ grows quadratically in n.
- Computing the decision functions grows linearly in *n*.
- Solving the optimization problem costs between $O(n^2)$ and $O(n^3)$

Computational Requirements

- The size of the optimization problem grows linearly in *n*.
- The kernel matrix $(k(x_i, x_j))$ grows quadratically in n.
- Computing the decision functions grows linearly in *n*.
- Solving the optimization problem costs between $O(n^2)$ and $O(n^3)$

Consequences

- For 64GB machines, kernel matrices for n > 100.000 cannot be stored.
- Training for such sample sizes, even if only a fixed parameter pair (λ, σ) is considered, may take up to hours.
- Using kernel methods without tricks is impossible for data sizes ranging in the millions.

- Nyström Method (Williams & Seeger, 2001): Approximate matrix by low-rank matrix. Theoretical analysis by Rudi, Camariano, & Rosasco (2015)
- Random Fourier Features (Rahimi & Recht, 2007): Approximate kernel by finite dimensional kernel. Bounds for approximation by Sriperumbudur & Szabo (2015)
- Chunking: Devide data into smaller subsets. Analysis for random splits: Zhang, Duchi & Wainwright (2014), Zhou, Guo & Lin (2015). next talk! Spatial splits: now

- Split bounded $X \subset \mathbb{R}^d$ into cells A_1, \ldots, A_m of diameter $\leq r$.
- On each cell A_j train a kernel method with Gaussian kernel and the data in A_j, i.e.

$$D_j := \{(x_i, y_i) \in D : x_i \in A_j\}$$
.

- The hyper-parameters λ and σ are found by training/validation on each cell separately.
- To predict y for some test sample x, only take the decision function that is constructed on the cell A_j with x ∈ A_j.

Rates for Localized kernel methods (Meister & S., 2016)

• Pick some
$$eta > 0$$
 and $r_n \sim n^{-1/eta}$

• Assume that $f_{L,P}^* \in W_2^s(X)$ for some $s < \frac{\beta-d}{2}$.

Then the localized kernel method learns with rate $n^{-\frac{2s}{2s+d}+\varepsilon}$.

Rates for Localized kernel methods (Meister & S., 2016)

• Pick some
$$eta > 0$$
 and $r_{\it n} \sim {\it n}^{-1/eta}$

• Assume that $f_{L,P}^* \in W_2^s(X)$ for some $s < \frac{\beta-d}{2}$.

Then the localized kernel method learns with rate $n^{-\frac{2s}{2s+d}+\varepsilon}$.

Remarks

- Good adaptivity requires large β .
- Large β leads to large cells.

 → Trade-off between statistics and computational complexity.
- Similar results for quantile regression.

• The split kernel method can be viewed as an ordinary kernel method using the RKHS

$$\mathcal{H} = igoplus_{j=1}^m \sqrt{\lambda_j} \mathcal{H}_{\mathcal{A}_j,\sigma_j}$$

- Investigate how properties of the local RKHS influence properties of the global *H* in view of *P*.
- Again we are facing a kernel more complex than usual.

Data: covertype in binary classification from LIBSVM site, d = 54
 Method: Hinge loss and number of samples in cells are controlled



Deeper Compositions

Structure of a Neural Network



At each non-input node, we perform the operation

$$x \mapsto \sigma(\langle w, x \rangle + b)$$

- Do we need the network structure on the right to classify?
- Can we replace the feature modification on the left by something else?

A simple Network with One Hidden Layer

• Input space X = [0, 1]

• One hidden layer with *m* ReLU-units each performing

$$x \mapsto \Phi_j(x) := |w_j x + b_j|_+, \qquad j = 1, \ldots, m.$$

• Output layer creates a function

$$egin{aligned} &x\mapsto \langle v,\Phi(x)
angle_{\ell_2^d} = \sum_{j=1}^m v_j ig| w_j x + b_j ig|_+ \end{aligned}$$

Thus it realizes an element in the RKHS with FM $\Phi := (\Phi_1, \dots, \Phi_m)$.

A simple Network with One Hidden Layer

• Input space X = [0, 1]

• One hidden layer with *m* ReLU-units each performing

$$x \mapsto \Phi_j(x) := |w_j x + b_j|_+, \qquad j = 1, \ldots, m.$$

• Output layer creates a function

$$egin{aligned} &x\mapsto \langle v,\Phi(x)
angle_{\ell_2^d} = \sum_{j=1}^m v_j ig| w_j x + b_j ig|_+ \end{aligned}$$

Thus it realizes an element in the RKHS with FM $\Phi := (\Phi_1, \dots, \Phi_m)$.

• For fixed $w, b \in \mathbb{R}^m$ this RKHS is a set of piecewise linear functions with kinks at

$$-\frac{b_1}{w_1},\ldots,-\frac{b_m}{w_m}$$

• The NN represents all piecewise linear functions with at most m-1 kinks and most with m kinks .

 \rightsquigarrow nonlinear structure, parametric method for each fixed design

Observation

Each layer performs a non-linear transformation

 $\mathbb{R}^{m_i} \to \mathbb{R}^{m_{i+1}}$ $x \mapsto \Phi_{w_i, b_i}(x)$

Entire feature map is $\Phi := \Phi_{w_L, b_L} \circ \cdots \circ \Phi_{w_1, b_1}$

Idea for Rest of the Talk

Replace finite-dimensional spaces by infinite dimensional Hilbert spaces

$$H_i
ightarrow H_{i+1}$$

 $x \mapsto \Phi_{w_i}(x)$

Use the kernel of the resulting feature map $\Phi := \Phi_{w_l} \circ \cdots \circ \Phi_{w_1}$

Bach, Lanckriet, and Jordan 2004 L = 2, linear kernel in second layer \rightsquigarrow Multiple kernel learning

Cho and Saul, 2009 General setup and some examples

Zhuang, Tsang, and Hoi, 2011 L = 2, sum of kernels in composition step, pseudo-dimension bound

Strobl and Visweswaran, 2013 Sums of kernels in each composition step, VC-bounds

Tang, 2013

 $\Phi_{L-1} \circ \cdots \circ \Phi_1$ is a neural net with *M* output nodes, Φ_L is linear "SVM".

Wilson, Hu, Salakhutdinov, and Xing, 2016

 $\Phi_{L-1} \circ \cdots \circ \Phi_1$ is a neural net with *M* output nodes, Φ_L is non-linear.

Observations

Let *H* be a Hilbert space and $\Phi : X \to H$.

• We obtain a new kernel on X by

$$k_{\gamma,H,X}(x,x') := \expigl(-\gamma^{-2} \|\Phi(x) - \Phi(x')\|_H^2igr), \qquad x,x' \in X\,,$$

• If $k(x,x') := \langle \Phi(x), \Phi(x') \rangle$ with $k(x,x) \equiv c$, then

$$k_{\gamma,X,H}(x,x') = \exp\left(-2\gamma^{-2}(c-k(x,x'))\right)$$

• If $\Phi_{\gamma,H}: H o H_{\gamma,H}$ is a feature map of $k_{\gamma,H}$ on H, then

 $\Phi_{\gamma,H} \circ \Phi$

is a feature map of $k_{\gamma,H,X}$.

Idea

So far we have

$$k_{\gamma,X,H}(x,x') = \exp(-2\gamma^{-2}(c-k(x,x')))$$
(2)

• For
$$I \subset \{1, \ldots, d\}$$
 we write $x_I := (x_i)_{i \in I}$.

- For $I_1, \ldots I_m \subset \{1, \ldots, d\}$, let k_1, \ldots, k_m be kernels on $\mathbb{R}^{|I_1|}, \ldots, \mathbb{R}^{|I_1|}$.
- Assume that $k_i(x,x) \equiv 1$.

For $I := I_1 \cup \cdots \cup I_m$ consider the kernel

$$k(x,x') := \sum_{i=1}^{m} w_i^2 k_i(x_{I_i},x'_{I_i}), \qquad x,x' \in X_I.$$

in (2). This kernel is denoted by k_w . This can be iterated!

Definition

Let H be the RKHS of the kernel

$$k(x, x') := \sum_{i=1}^{m} w_i^2 k_i(x_{I_i}, x'_{I_i}), \qquad x, x' \in X_I.$$

Then the resulting hierarchical Gaussian kernel $k_{\gamma,X_I,H}$, that is

$$k_{\gamma,X,H}(x,x') = \exp\left(-2\gamma^{-2}(c-k(x,x'))\right)$$

is said to be:

- of depth 1, if all kernels k_1, \ldots, k_m are linear kernels.
- of depth L > 1, if all k₁,..., k_m are hierarchical Gaussian kernels of depth L − 1.

Construction III

Example 1

Hierarchical Gaussian kernels of depth L = 1 are of the form

$$k_{\mathbf{w}}(x,x') := \exp\left(-\sum_{i\in I} w_i^2 (x_i - x_i')^2
ight), \qquad x,x'\in X,$$

ARD kernel

Example 2

Hierarchical Gaussian kernels of depth L = 2 are of the form

$$k_{\mathbf{W}^{(1)},\mathbf{w},\gamma}(x,x') = \exp\left(-2\gamma^{-2}\sum_{i=1}^{m} w_i^2 (1 - k_{\mathbf{w}_i}(x_{l_i},x'_{l_i}))\right)$$
$$= \exp\left(-2\gamma^{-2}\sum_{i=1}^{m} w_i^2 \left(1 - \exp\left(-\sum_{j \in I_i} w_{j,i}^2 (x_j - x'_j)^2\right)\right)\right)$$

Structure of a Hierarchical Gaussian Kernel



Example of a hierarchical Gaussian kernels of depth L = 3.

Definition

A continuous kernel on a compact metric space X is universal, if its RKHS is dense in C(X).

Theorem (Christmann & S., 2010)

A kernel of the form

$$k_{\gamma,H,X}(x,x') := \exp\left(-\gamma^{-2} \|\Phi(x) - \Phi(x')\|_H^2\right), \qquad x,x' \in X,$$

is universal, if $\boldsymbol{\Phi}$ is continuous and injective.

A Bit Theory II

Theorem (S. & Thomann, 2016)

A hierarchical Gaussian kernel of depth $L \ge 1$ is universal, if it does not ignore coordinates.

Corollary (S. & Thomann, 2016)

Every SVM using a fixed hierarchical Gaussian kernel of depth $L \ge 1$ that does not ignore coordinates is universally consistent.

Remarks

- Learning rates for weights changing with sample size *n*?
- For which distributions do hierarchical Gaussian kernels help?
- Learning the kernel can be, in principle, decoupled from learning a classifier/regressor.

A Bit Theory III

A few words on the proof ...

- Induction over L
- At the highest level we have

$$k_{\gamma,X_I,H}(x,x') = \prod_{i=1}^{I} k_{\gamma/w_i,X,H_i}(x_{I_i},x'_{I_i}), \qquad x,x' \in X_I.$$

 If k₁ and k₂ are universal kernels on X₁ and X₂, then k₁ ⊗ k₂ defined by

$$k_I \otimes k_J(x,x') := k_I(x_I,x_I') \cdot k_J(x_J,x_J'), \qquad x,x' \in X_{I \cup J}$$

is a universal kernel on $X_{I\cup J}$. Use Stone-Weierstraß.

• Universal kernels have injective feature maps. $\rightsquigarrow k_{\gamma/w_i, X, H_i}$ are universal by induction assumption

Data Set	SVM	HKL	Ours	RF	DNN
BANK	.2978 ±.0024	$.2939 \pm .0028$.2596 ±.0039	.2687 ±.0027	$.2931 \pm .0025$
CADATA	$.0538 \pm .0016$	$.0625 \pm .0014$.0525 ±.0019	.0509 ±.0015	$.0550 \pm .0015$
COD	$.1574 \pm .0023$	$.1734 \pm .0013$.1309 ±.0050	$.1725 \pm .0020$.1154 ±.0013
COVTYPE	$.5205 \pm .0043$	$.6100 \pm .0042$.3995 ±.0148	.4878 ±.0041	$.5027 \pm .0063$
CPUSMALL	$.0036 \pm .0002$	$.0046 \pm .0004$.0034 ±.0002	.0032 ±.0002	$.0038 \pm .0001$
CYCLE	$.0105 \pm .0003$	$.0122 \pm .0003$.0098 ±.0005	.0084 ±.0003	$.0121 \pm .0003$
HIGGS	$.9021 \pm .0017$	$.8178 \pm .0074$	$.8023 \pm .0175$.7770 ±.0024	$.9162 \pm .0024$
LETTER	.0451 ±.0015	$.1151 \pm .0018$.0339 ±.0014	$.0577 \pm .0015$.0448 ±.0018
MAGIC	.4007 ±.0083	$.4282 \pm .0082$	$.3900 \pm .0093$.3772 ±.0079	.3783 ±.0085
PENDIGITS	.0079 ±.0007	$.0243 \pm .0012$.0070 ±.0007	$.0127 \pm .0012$.0079 ±.0010
SATIMAGE	.0488 ±.0029	$.1078 \pm .0059$.0467 ±.0030	$.0525 \pm .0026$	$.0525 \pm .0033$
SEISMIC	$.3113 \pm .0013$	$.3189 \pm .0022$.2981 ±.0016	.2955 ±.0012	.2975 ±.0014
SHUTTLE	$.0046 \pm .0003$	$.0129 \pm .0007$.0042 ±.0004	.0008 ±.0002	$.0059 \pm .0004$
THYROID	$.1750 \pm .0081$	$.1637 \pm .0083$.1538 ±.0080	.0251 ±.0031	.1522 ±.0080
UPDRS	$.0537 \pm .0052$	$.1774 \pm .0090$.0059 ±.0021	.0305 ±.0016	$.0531 \pm .0042$

Detailed Comparison of 3 Best Methods



Paper

M. Eberts and I. Steinwart, *Optimal regression rates for SVMs using Gaussian kernels*, Electron. J. Stat. 7, 1-42, 2013.

M. Meister and I. Steinwart, *Optimal learning rates for localized kernel methods,* J. Mach. Learn. Res. 17, 1-44, 2016.

S. Fischer and I. Steinwart, *Sobolev norm learning rates for regularized least-squares algorithm*, https://arxiv.org/abs/1702.07254

I. Steinwart, P. Thomann, and N. Schmid, *Learning with hierarchical Gaussian kernels*, http://arxiv.org/abs/1612.00824, 2016

I. Steinwart, D. Hush, and C. Scovel, *Optimal rates for regularized least squares regression*, 22nd Annual Conference on Learning Theory (COLT), 79-93, 2009.

I. Steinwart and C. Scovel, Mercer's theorem on general domains: on the interaction between measures, kernels, and RKHSs, Constr. Approx. 35, 363-417, 2012.

Software

I. Steinwart, LiquidSVM, http://www.isa.uni-stuttgart.de/Steinwart/software, R, Java, and Python interface by P. Thomann, Matlab interface by N. Schmid Paper at https://arxiv.org/abs/1702.06899