

Asymptotic Normality of Support Vector Machines

Robert Hable
Department of Mathematics
University of Bayreuth

Nonparametric Regression

$$Y = f_0(X) + \varepsilon$$

with

- ▶ Y : output variable (observable)
- ▶ X : input variable (observable)
- ▶ f_0 : regression function (unknown)
- ▶ ε : error term (not observable)

Goal: Estimation of the unknown regression function f_0

Support Vector Machines

$$Y_i = f_0(X_i) + \varepsilon_i, \quad (X_i, Y_i) \sim P \quad \text{i.i.d.}, \quad i \in \{1, \dots, n\}$$

Goal: Estimation of $f_0 : \mathcal{X} \rightarrow \mathcal{Y} \subset \mathbb{R}$

Support Vector Machines

$$Y_i = f_0(X_i) + \varepsilon_i, \quad (X_i, Y_i) \sim P \quad \text{i.i.d.}, \quad i \in \{1, \dots, n\}$$

Goal: Estimation of $f_0 : \mathcal{X} \rightarrow \mathcal{Y} \subset \mathbb{R}$

- ▶ Loss function

$$L : \mathcal{Y} \times \mathbb{R} \rightarrow [0, \infty)$$

$L(y, t)$: loss caused by estimation $t = \hat{f}_n(x)$ if y is true

Support Vector Machines

$$Y_i = f_0(X_i) + \varepsilon_i, \quad (X_i, Y_i) \sim P \quad \text{i.i.d.}, \quad i \in \{1, \dots, n\}$$

Goal: Estimation of $f_0 : \mathcal{X} \rightarrow \mathcal{Y} \subset \mathbb{R}$

- ▶ Loss function

$$L : \mathcal{Y} \times \mathbb{R} \rightarrow [0, \infty)$$

$L(y, t)$: loss caused by estimation $t = \hat{f}_n(x)$ if y is true

- ▶ Risk of an estimate $\hat{f}_n : \mathcal{X} \rightarrow \mathbb{R}$

$$\int L(y, \hat{f}_n(x)) P(d(x, y))$$

Support Vector Machines

$$Y_i = f_0(X_i) + \varepsilon_i, \quad (X_i, Y_i) \sim P \quad \text{i.i.d.}, \quad i \in \{1, \dots, n\}$$

Goal: Estimation of $f_0 : \mathcal{X} \rightarrow \mathcal{Y} \subset \mathbb{R}$

- ▶ Loss function

$$L : \mathcal{Y} \times \mathbb{R} \rightarrow [0, \infty)$$

$L(y, t)$: loss caused by estimation $t = \hat{f}_n(x)$ if y is true

- ▶ empirical Risk of an estimate $\hat{f}_n : \mathcal{X} \rightarrow \mathbb{R}$

$$\frac{1}{n} \sum_{i=1}^n L(y_i, \hat{f}_n(x_i))$$

Support Vector Machines

$$Y_i = f_0(X_i) + \varepsilon_i, \quad (X_i, Y_i) \sim P \quad \text{i.i.d.}, \quad i \in \{1, \dots, n\}$$

Goal: Estimation of $f_0 : \mathcal{X} \rightarrow \mathcal{Y} \subset \mathbb{R}$

- ▶ Loss function

$$L : \mathcal{Y} \times \mathbb{R} \rightarrow [0, \infty)$$

$L(y, t)$: loss caused by estimation $t = \hat{f}_n(x)$ if y is true

- ▶ empirical Risk of an estimate $\hat{f}_n : \mathcal{X} \rightarrow \mathbb{R}$

$$\frac{1}{n} \sum_{i=1}^n L(y_i, \hat{f}_n(x_i))$$

- ▶ RKHS H (certain Hilbert space of functions $f : \mathcal{X} \rightarrow \mathbb{R}$)

Support Vector Machines

$$Y_i = f_0(X_i) + \varepsilon_i, \quad (X_i, Y_i) \sim P \quad \text{i.i.d.}, \quad i \in \{1, \dots, n\}$$

Goal: Estimation of $f_0 : \mathcal{X} \rightarrow \mathcal{Y} \subset \mathbb{R}$

- ▶ Loss function

$$L : \mathcal{Y} \times \mathbb{R} \rightarrow [0, \infty)$$

$L(y, t)$: loss caused by estimation $t = \hat{f}_n(x)$ if y is true

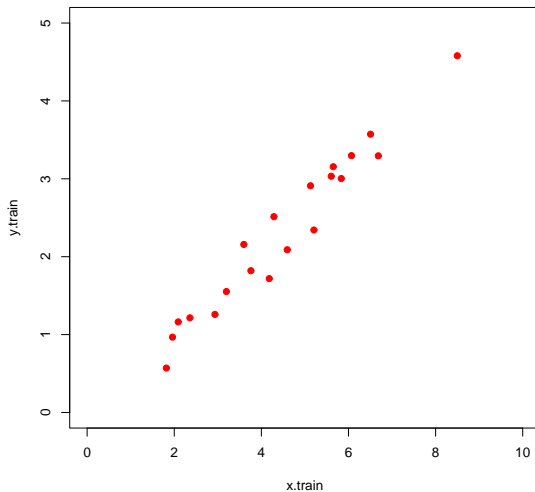
- ▶ empirical Risk of an estimate $\hat{f}_n : \mathcal{X} \rightarrow \mathbb{R}$

$$\frac{1}{n} \sum_{i=1}^n L(y_i, \hat{f}_n(x_i))$$

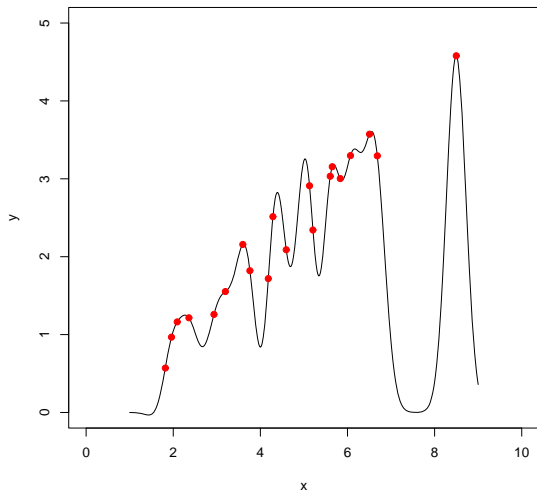
- ▶ RKHS H (certain Hilbert space of functions $f : \mathcal{X} \rightarrow \mathbb{R}$)
- ▶ Support vector machine

$$S_n((x_1, y_1), \dots, (x_n, y_n)) = \arg \inf_{f \in H} \frac{1}{n} \sum_{i=1}^n L(y_i, f(x_i))$$

Overfitting



Overfitting



Support Vector Machines

$$Y_i = f_0(X_i) + \varepsilon_i, \quad (X_i, Y_i) \sim P \quad \text{i.i.d.}, \quad i \in \{1, \dots, n\}$$

Goal: Estimation of $f_0 : \mathcal{X} \rightarrow \mathcal{Y} \subset \mathbb{R}$

- ▶ Loss function

$$L : \mathcal{Y} \times \mathbb{R} \rightarrow [0, \infty)$$

$L(y, t)$: loss caused by prediction t if y is the true value

- ▶ empirical Risk of an estimate $f : \mathcal{X} \rightarrow \mathbb{R}$

$$\frac{1}{n} \sum_{i=1}^n L(y_i, f(x_i))$$

- ▶ RKHS H (certain Hilbert space of functions $f : \mathcal{X} \rightarrow \mathbb{R}$)
- ▶ Support vector machine

$$S_n((x_1, y_1), \dots, (x_n, y_n)) = \arg \inf_{f \in H} \frac{1}{n} \sum_{i=1}^n L(y_i, f(x_i))$$

Support Vector Machines

$$Y_i = f_0(X_i) + \varepsilon_i, \quad (X_i, Y_i) \sim P \quad \text{i.i.d.}, \quad i \in \{1, \dots, n\}$$

Goal: Estimation of $f_0 : \mathcal{X} \rightarrow \mathcal{Y} \subset \mathbb{R}$

- ▶ Loss function

$$L : \mathcal{Y} \times \mathbb{R} \rightarrow [0, \infty)$$

$L(y, t)$: loss caused by prediction t if y is the true value

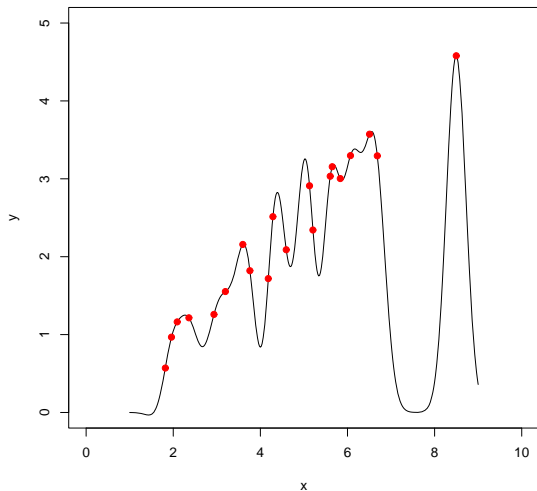
- ▶ empirical Risk of an estimate $f : \mathcal{X} \rightarrow \mathbb{R}$

$$\frac{1}{n} \sum_{i=1}^n L(y_i, f(x_i))$$

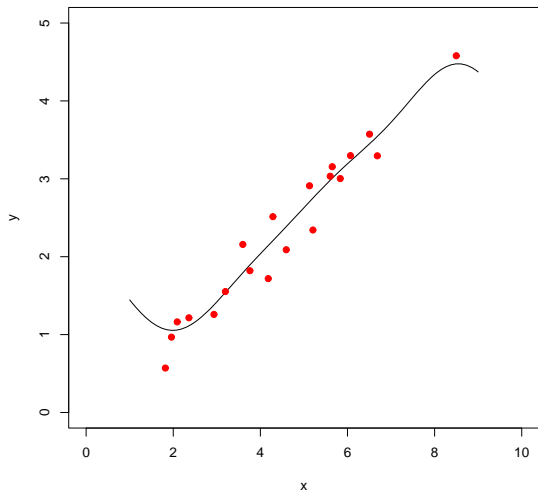
- ▶ RKHS H (certain Hilbert space of functions $f : \mathcal{X} \rightarrow \mathbb{R}$)
- ▶ Support vector machine

$$S_n((x_1, y_1), \dots, (x_n, y_n)) = \arg \inf_{f \in H} \frac{1}{n} \sum_{i=1}^n L(y_i, f(x_i)) + \lambda \|f\|_H^2$$

Overfitting



Overfitting



Reproducing Kernel Hilbert Space (RKHS)

Support Vector Machines

$$S_n : (\mathcal{X} \times \mathcal{Y})^n \rightarrow H,$$

$$((x_1, y_1), \dots, (x_n, y_n)) \mapsto \arg \inf_{f \in H} \frac{1}{n} \sum_{i=1}^n L(y_i, f(x_i)) + \lambda \|f\|_H^2$$

with H a reproducing kernel Hilbert space (RKHS)

Reproducing Kernel Hilbert Space (RKHS)

Support Vector Machines

$$S_n : (\mathcal{X} \times \mathcal{Y})^n \rightarrow H,$$

$$((x_1, y_1), \dots, (x_n, y_n)) \mapsto \arg \inf_{f \in H} \frac{1}{n} \sum_{i=1}^n L(y_i, f(x_i)) + \lambda \|f\|_H^2$$

with H a reproducing kernel Hilbert space (RKHS)

Reproducing kernel Hilbert space H

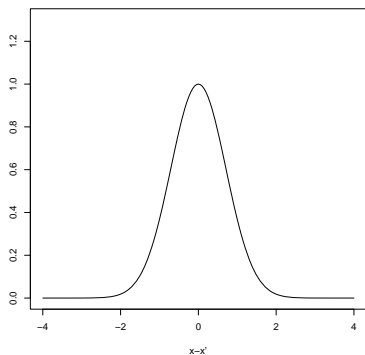
- ▶ a Hilbert space of functions $f : \mathcal{X} \rightarrow \mathbb{R}$
- ▶ generated by a kernel function $k : \mathcal{X} \times \mathcal{X} \rightarrow \mathbb{R}$
- ▶ reproducing property

$$\langle f, k(x, \cdot) \rangle_H = f(x) \quad \forall x \in \mathcal{X}, \quad \forall f \in H$$

Example: Gaussian Kernel

Gaussian Kernel $\mathcal{X} = \mathbb{R}$

$$k : \mathbb{R} \times \mathbb{R} \rightarrow \mathbb{R}, \quad (x, x') \mapsto \exp\left(-\frac{1}{\gamma^2}|x - x'|^2\right)$$

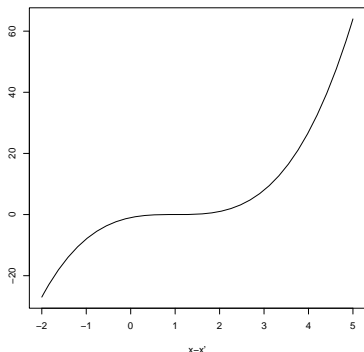


$H \subset L_p(P)$ dense

Example: Polynomial Kernel

Polynomial Kernel $\mathcal{X} = \mathbb{R}$

$$k : \mathbb{R} \times \mathbb{R} \rightarrow \mathbb{R}, \quad (x, x') \mapsto (x \cdot x' + c)^m$$



$$H = \{f : \mathbb{R} \rightarrow \mathbb{R} \mid f \text{ a polynomial with degree } \leq m\} \cong \mathbb{R}^{m+1}$$

Representer Theorem

How to calculate the SVM?

$$D_n = ((x_1, y_1), \dots, (x_n, y_n))$$

$$\text{SVM: } f_{D_n, \lambda} = \arg \inf_{f \in H} \frac{1}{n} \sum_{i=1}^n L(y_i, f(x_i)) + \lambda \|f\|_H^2$$

Representer Theorem

How to calculate the SVM?

$$D_n = ((x_1, y_1), \dots, (x_n, y_n))$$

$$\text{SVM: } f_{D_n, \lambda} = \arg \inf_{f \in H} \frac{1}{n} \sum_{i=1}^n L(y_i, f(x_i)) + \lambda \|f\|_H^2$$

Representer Theorem

There are $\alpha_{D_n, 1}, \dots, \alpha_{D_n, n} \in \mathbb{R}$ such that

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

Representer Theorem

How to calculate the SVM?

$$D_n = ((x_1, y_1), \dots, (x_n, y_n))$$

$$\text{SVM: } f_{D_n, \lambda} = \arg \inf_{f \in H} \frac{1}{n} \sum_{i=1}^n L(y_i, f(x_i)) + \lambda \|f\|_H^2$$

Representer Theorem

There are $\alpha_{D_n, 1}, \dots, \alpha_{D_n, n} \in \mathbb{R}$ such that

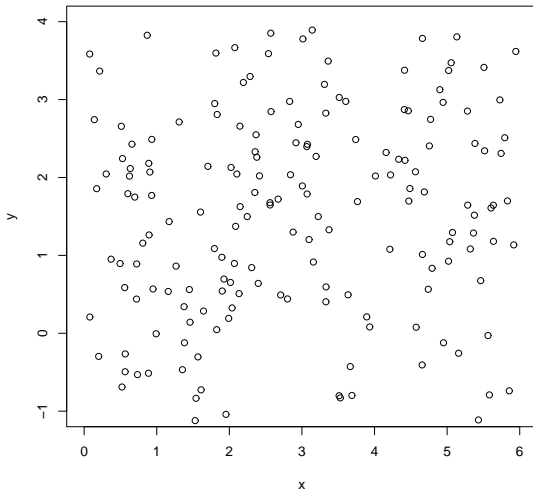
$$f_{D_n, \lambda} = \sum_{i=1}^n \alpha_{D_n, i} k(x_i, \cdot).$$

→ **just solve a finite convex optimization problem**

... and this really works?

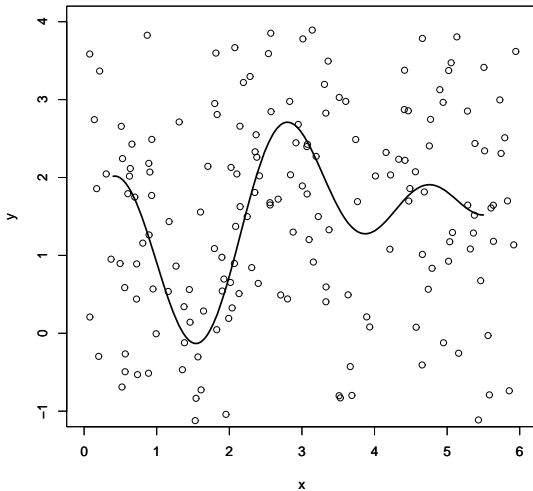
... and this really works?

Yes, quite good.



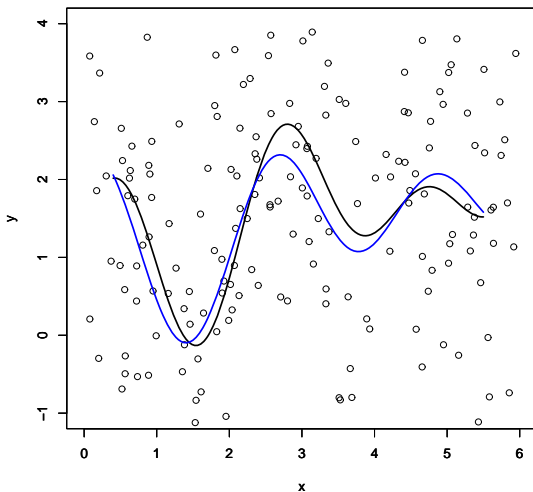
... and this really works?

Yes, quite good.



... and this really works?

Yes, quite good.



Risk-Consistency

Risk of a predictor $f : \mathcal{X} \rightarrow \mathbb{R}$

$$\mathcal{R}_P(f) = \int L(y, f(x)) P(d(x, y)) \quad \hat{=} \quad \text{Quality of } f$$

$$\mathbf{D}_n = ((X_1, Y_1), \dots, (X_n, Y_n))$$

$$\text{SVM: } f_{\mathbf{D}_n, \lambda_n} = \arg \inf_{f \in H} \frac{1}{n} \sum_{i=1}^n L(Y_i, f(X_i)) + \lambda_n \|f\|_H^2$$

Risk-Consistency

Risk of a predictor $f : \mathcal{X} \rightarrow \mathbb{R}$

$$\mathcal{R}_P(f) = \int L(y, f(x)) P(d(x, y)) \quad \hat{=} \quad \text{Quality of } f$$

$$\mathbf{D}_n = ((X_1, Y_1), \dots, (X_n, Y_n))$$

$$\text{SVM: } f_{\mathbf{D}_n, \lambda_n} = \arg \inf_{f \in H} \frac{1}{n} \sum_{i=1}^n L(Y_i, f(X_i)) + \lambda_n \|f\|_H^2$$

Risk-consistency

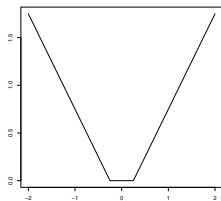
$$\mathcal{R}_P(f_{\mathbf{D}_n, \lambda_n}) \xrightarrow{n \rightarrow \infty} \inf_{f: \mathcal{X} \rightarrow \mathbb{R}} \mathcal{R}_P(f) \quad \text{in probability}$$

essentially if

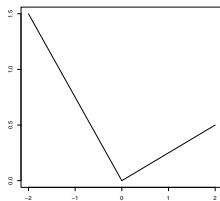
- ▶ $H \subset L_P(P)$ dense (e.g. Gaussian kernel)
- ▶ $\lambda_n \rightarrow 0$ not too fast (!)

Robustness

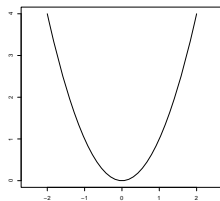
Loss function L



ϵ -insensitive



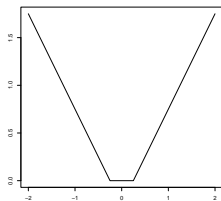
pinball



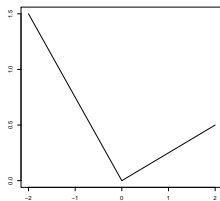
least squares

Robustness

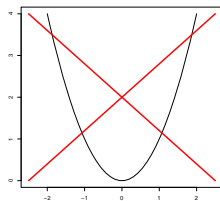
Loss function L should be Lipschitz continuous



ϵ -insensitive



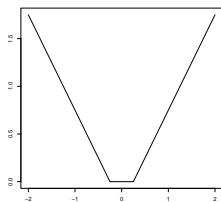
pinball



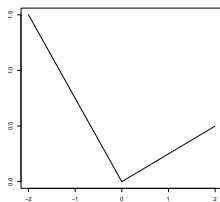
least squares

Robustness

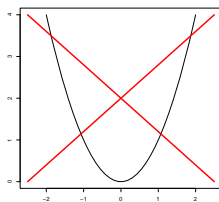
Loss function L should be Lipschitz continuous



ϵ -insensitive



pinball



least squares

- ▶ bounded influence function (Christmann & Van Messem, 2008)
- ▶ bounds for maxbias (Christmann & Van Messem, 2008)
- ▶ qualitative robustness (Hable & Christmann, 2010)

Rates of Convergence

Risk-consistency

$$\mathcal{R}_P(f_{\mathbf{D}_n, \lambda_n}) \xrightarrow{n \rightarrow \infty} \inf_{f: \mathcal{X} \rightarrow \mathbb{R}} \mathcal{R}_P(f) \quad \text{in probability}$$

Rates of Convergence

Risk-consistency

$$\mathcal{R}_P(f_{\mathbf{D}_n, \lambda_n}) \xrightarrow{n \rightarrow \infty} \inf_{f: \mathcal{X} \rightarrow \mathbb{R}} \mathcal{R}_P(f) \quad \text{in probability}$$

How fast is this convergence?

Is there a **uniform** rate r_n such that

$$r_n \left(\mathcal{R}_P(f_{\mathbf{D}_n, \lambda_n}) - \inf_{f: \mathcal{X} \rightarrow \mathbb{R}} \mathcal{R}_P(f) \right) \xrightarrow{n \rightarrow \infty} 0 \quad \text{in probability}$$

for **every** P ?

Rates of Convergence

Risk-consistency

$$\mathcal{R}_P(f_{\mathbf{D}_n, \lambda_n}) \xrightarrow{n \rightarrow \infty} \inf_{f: \mathcal{X} \rightarrow \mathbb{R}} \mathcal{R}_P(f) \quad \text{in probability}$$

How fast is this convergence?

Is there a **uniform** rate r_n such that

$$r_n \left(\mathcal{R}_P(f_{\mathbf{D}_n, \lambda_n}) - \inf_{f: \mathcal{X} \rightarrow \mathbb{R}} \mathcal{R}_P(f) \right) \xrightarrow{n \rightarrow \infty} 0 \quad \text{in probability}$$

for **every** P ? \longrightarrow **No!** (no-free-lunch theorem)

Rates of Convergence

Risk-consistency

$$\mathcal{R}_P(f_{\mathbf{D}_n, \lambda_n}) \xrightarrow{n \rightarrow \infty} \inf_{f: \mathcal{X} \rightarrow \mathbb{R}} \mathcal{R}_P(f) \quad \text{in probability}$$

How fast is this convergence?

Is there a **uniform** rate r_n such that

$$r_n \left(\mathcal{R}_P(f_{\mathbf{D}_n, \lambda_n}) - \inf_{f: \mathcal{X} \rightarrow \mathbb{R}} \mathcal{R}_P(f) \right) \xrightarrow{n \rightarrow \infty} 0 \quad \text{in probability}$$

for **every** P ? \longrightarrow **No!** (no-free-lunch theorem)

Instead,

rates r_n of convergence under assumptions on P

e.g. Steinwart and Scovel (2007), Caponnetto and De Vito (2007), Blanchard et al. (2008), Steinwart et al. (2009), Mendelson and Neeman (2010)

Smooth Approximation of the Regression Function

Goal: estimate a solution $f^* : \mathcal{X} \rightarrow \mathbb{R}$ of

$$\mathcal{R}_P(f) = \min_{f : \mathcal{X} \rightarrow \mathbb{R}}$$

or

$$\inf_{f \in H} \mathcal{R}_P(f) = \min_{f \in H}$$

Smooth Approximation of the Regression Function

Goal: estimate a solution $f^* : \mathcal{X} \rightarrow \mathbb{R}$ of

$$\mathcal{R}_P(f) = \min_{f : \mathcal{X} \rightarrow \mathbb{R}}$$

or

$$\inf_{f \in H} \mathcal{R}_P(f) = \min_{f \in H}$$

However, these optimization problems

- ▶ are ill-posed and
- ▶ there is no uniform rate of convergence to the solution (without substantial assumptions on P)

Smooth Approximation of the Regression Function

Goal: estimate a solution $f^* : \mathcal{X} \rightarrow \mathbb{R}$ of

$$\mathcal{R}_P(f) = \min_{f : \mathcal{X} \rightarrow \mathbb{R}}$$

or

$$\inf_{f \in H} \mathcal{R}_P(f) = \min_{f \in H}$$

However, these optimization problems

- ▶ are ill-posed and
- ▶ there is no uniform rate of convergence to the solution
(without substantial assumptions on P)

Instead, consider the regularized problem

$$\mathcal{R}_P(f) + \lambda_0 \|f\|_H^2 = \min_{f \in H}$$

Smooth Approximation of the Regression Function

- ▶ Instead of estimating a solution $f^* : \mathcal{X} \rightarrow \mathbb{R}$ of

$$\mathcal{R}_P(f) = \min! \quad f : \mathcal{X} \rightarrow \mathbb{R}$$

we may estimate the solution f_{P,λ_0} of the regularized problem

$$\mathcal{R}_P(f) + \lambda_0 \|f\|_H^2 = \min! \quad f \in H.$$

f_{P,λ_0} serves as a “smoother approximation” of f^* .

Smooth Approximation of the Regression Function

- ▶ Instead of estimating a solution $f^* : \mathcal{X} \rightarrow \mathbb{R}$ of

$$\mathcal{R}_P(f) = \min! \quad f : \mathcal{X} \rightarrow \mathbb{R}$$

we may estimate the solution f_{P,λ_0} of the regularized problem

$$\mathcal{R}_P(f) + \lambda_0 \|f\|_H^2 = \min! \quad f \in H.$$

f_{P,λ_0} serves as a “smoother approximation” of f^* .

- ▶ The regularized problem is equivalent to

$$\mathcal{R}_P(f) = \min! \quad f \in H, \quad \|f\|_H \leq r_0.$$

r_0 : bound on complexity of “smoother approximation”

Asymptotic Normality of Regularized Problem

Under some

- ▶ assumptions on \mathcal{X} , L , and k ($\leftrightarrow H$)
- ▶ but (essentially) no assumptions on P ,

we have

$$\sqrt{n} \left(\mathcal{R}(f_{\mathbf{D}_n, \lambda_0}) - \mathcal{R}(f_{P, \lambda_0}) \right) \rightsquigarrow \mathcal{N}(0, \sigma^2)$$

Asymptotic Normality of Regularized Problem

Under some

- ▶ assumptions on \mathcal{X} , L , and k ($\leftrightarrow H$)
- ▶ but (essentially) no assumptions on P ,

we have

$$\sqrt{n} \left(\mathcal{R}(f_{\mathbf{D}_n, \lambda_0}) - \mathcal{R}(f_{P, \lambda_0}) \right) \rightsquigarrow \mathcal{N}(0, \sigma^2)$$

and, even more,

$$\sqrt{n} (f_{\mathbf{D}_n, \lambda_0} - f_{P, \lambda_0}) \rightsquigarrow \text{Gaussian process in } H$$

Asymptotic Normality of Regularized Problem

Under some

- ▶ assumptions on \mathcal{X} , L , and k ($\leftrightarrow H$)
- ▶ but (essentially) no assumptions on P ,

we have

$$\sqrt{n} \left(\mathcal{R}(f_{\mathbf{D}_n, \lambda_0}) - \mathcal{R}(f_{P, \lambda_0}) \right) \rightsquigarrow \mathcal{N}(0, \sigma^2)$$

and, even more,

$$\sqrt{n} (f_{\mathbf{D}_n, \lambda_0} - f_{P, \lambda_0}) \rightsquigarrow \text{Gaussian process in } H$$

Asymptotic Normality of Regularized Problem

Under some

- ▶ assumptions on \mathcal{X} , L , k ($\leftrightarrow H$), and $\lambda_n \xrightarrow[n \rightarrow \infty]{} \lambda_0$
- ▶ but (essentially) no assumptions on P ,

we have

$$\sqrt{n} \left(\mathcal{R}(f_{\mathbf{D}_n, \lambda_n}) - \mathcal{R}(f_{P, \lambda_0}) \right) \rightsquigarrow \mathcal{N}(0, \sigma^2)$$

and, even more,

$$\sqrt{n} (f_{\mathbf{D}_n, \lambda_n} - f_{P, \lambda_0}) \rightsquigarrow \text{Gaussian process in } H$$

Asymptotic Normality of Regularized Problem

Under some

- ▶ assumptions on \mathcal{X} , L , k ($\leftrightarrow H$), and $\lambda_{\mathbf{D}_n} \xrightarrow{n \rightarrow \infty} \lambda_0$
- ▶ but (essentially) no assumptions on P ,

we have

$$\sqrt{n} \left(\mathcal{R}(f_{\mathbf{D}_n, \lambda_{\mathbf{D}_n}}) - \mathcal{R}(f_{P, \lambda_0}) \right) \rightsquigarrow \mathcal{N}(0, \sigma^2)$$

and, even more,

$$\sqrt{n} (f_{\mathbf{D}_n, \lambda_{\mathbf{D}_n}} - f_{P, \lambda_0}) \rightsquigarrow \text{Gaussian process in } H$$

Asymptotic Normality of Regularized Problem

Corollary

In particular, we also have for every $x_1, \dots, x_m \in \mathcal{X}$

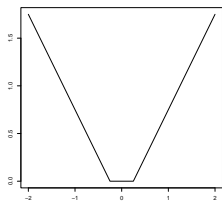
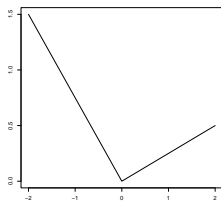
$$\sqrt{n} \begin{pmatrix} f_{\mathbf{D}_n, \lambda_{\mathbf{D}_n}}(x_1) - f_{P, \lambda_0}(x_1) \\ \vdots \\ f_{\mathbf{D}_n, \lambda_{\mathbf{D}_n}}(x_m) - f_{P, \lambda_0}(x_m) \end{pmatrix} \rightsquigarrow \mathcal{N}_m(0, \Sigma)$$

where Σ is a covariance matrix.

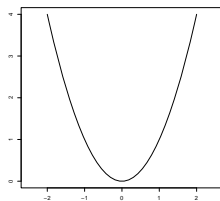
(follows from the reproducing property of the kernel k)

under some assumptions ...

- ▶ $\mathcal{X} \subset \mathbb{R}^d$ compact
- ▶ k more than $d/2$ -times continuously differentiable
- ▶ L smooth (2-times differentiable) and integrable

 ε -insensitive

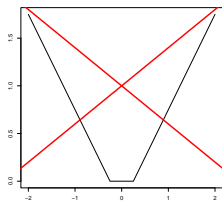
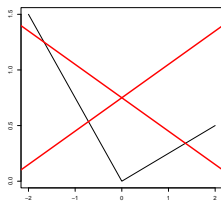
pinball



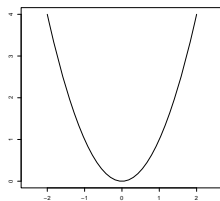
least squares

under some assumptions ...

- ▶ $\mathcal{X} \subset \mathbb{R}^d$ compact
- ▶ k more than $d/2$ -times continuously differentiable
- ▶ L smooth (2-times differentiable) and integrable

 ϵ -insensitive

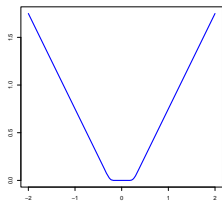
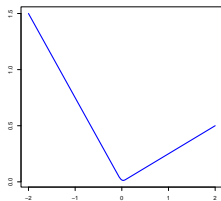
pinball



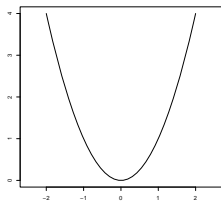
least squares

under some assumptions ...

- ▶ $\mathcal{X} \subset \mathbb{R}^d$ compact
- ▶ k more than $d/2$ -times continuously differentiable
- ▶ L smooth (2-times differentiable) and integrable

smoothed ε -insensitive

smoothed pinball



least squares

under some assumptions ...

and

▶ $\sqrt{n}(\lambda_{\mathbf{D}_n} - \lambda_0) \xrightarrow[n \rightarrow \infty]{} 0$ in probability.

under some assumptions ...

and

▶ $\sqrt{n}(\lambda_{\mathbf{D}_n} - \lambda_0) \xrightarrow[n \rightarrow \infty]{} 0$ in probability.

For example:

Choose a (large) constant $c > 0$ and do a cross-validation in

$$\left[\lambda_0, \lambda_0 + \frac{c}{\sqrt{n \ln(n)}} \right].$$

Sketch of the Proof: Functional Delta-Method

Consider the SVM-functional

$$S : \mathcal{M}_1 \rightarrow H, \quad P \mapsto f_{P, \lambda_0}$$

SVM-functional represents SVM-estimator:

$$f_{\mathbf{D}_n, \lambda_0} = S(\mathbb{P}_{\mathbf{D}_n})$$

Sketch of the Proof: Functional Delta-Method

Consider the SVM-functional

$$S : \mathcal{M}_1 \rightarrow H, \quad P \mapsto f_{P, \lambda_0}$$

SVM-functional represents SVM-estimator:

$$f_{\mathbf{D}_n, \lambda_0} = S(\mathbb{P}_{\mathbf{D}_n})$$

1. Show that $\sqrt{n}(\mathbb{P}_{\mathbf{D}_n} - P)$ converges weakly to a Gaussian process in a suitable space $\ell_\infty(\mathcal{G})$.

Sketch of the Proof: Functional Delta-Method

Consider the SVM-functional

$$S : \mathcal{M}_1 \rightarrow H, \quad P \mapsto f_{P, \lambda_0}$$

SVM-functional represents SVM-estimator:

$$f_{\mathbf{D}_n, \lambda_0} = S(\mathbb{P}_{\mathbf{D}_n})$$

1. Show that $\sqrt{n}(\mathbb{P}_{\mathbf{D}_n} - P)$ converges weakly to a Gaussian process in a suitable space $\ell_\infty(\mathcal{G})$.
2. Show that S is Hadamard-differentiable.

Sketch of the Proof: Functional Delta-Method

Consider the SVM-functional

$$S : \mathcal{M}_1 \rightarrow H, \quad P \mapsto f_{P, \lambda_0}$$

SVM-functional represents SVM-estimator:

$$f_{\mathbf{D}_n, \lambda_0} = S(\mathbb{P}_{\mathbf{D}_n})$$

1. Show that $\sqrt{n}(\mathbb{P}_{\mathbf{D}_n} - P)$ converges weakly to a Gaussian process in a suitable space $\ell_\infty(\mathcal{G})$.
2. Show that S is Hadamard-differentiable.
3. Then, it follows from the functional delta-method that

$$\sqrt{n}(f_{\mathbf{D}_n, \lambda_0} - f_{P, \lambda_0}) = \sqrt{n}(S(\mathbb{P}_{\mathbf{D}_n}) - S(P))$$

converges weakly to a Gaussian process.

Sketch of the Proof: Functional Delta-Method

Consider the SVM-functional

$$S : \mathcal{M}_1 \rightarrow H, \quad P \mapsto f_{P, \lambda_0}$$

SVM-functional represents SVM-estimator:

$$f_{\mathbf{D}_n, \lambda_0} = S(\mathbb{P}_{\mathbf{D}_n})$$

1. Show that $\sqrt{n}(\mathbb{P}_{\mathbf{D}_n} - P)$ converges weakly to a Gaussian process in a suitable space $\ell_\infty(\mathcal{G})$.
2. Show that S is Hadamard-differentiable:
3. Then, it follows from the functional delta-method that

$$\sqrt{n}(f_{\mathbf{D}_n, \lambda_0} - f_{P, \lambda_0}) = \sqrt{n}(S(\mathbb{P}_{\mathbf{D}_n}) - S(P))$$

converges weakly to a Gaussian process.

Sketch of the Proof: Functional Delta-Method

How to deal with random parameters λ_{D_n} ?

Sketch of the Proof: Functional Delta-Method

How to deal with random parameters λ_{D_n} ?

Problem

$$\int L(y, f(x)) P(d(x, y)) + \lambda \|f\|_H^2 = \min! \quad f \in H.$$

is equivalent to

$$\frac{\lambda_0}{\lambda} \left(\int L(y, f(x)) P(d(x, y)) + \lambda \|f\|_H^2 \right) = \min! \quad f \in H$$

Sketch of the Proof: Functional Delta-Method

How to deal with random parameters λ_{D_n} ?

Problem

$$\int L(y, f(x)) P(d(x, y)) + \lambda \|f\|_H^2 = \min! \quad f \in H.$$

is equivalent to

$$\frac{\lambda_0}{\lambda} \left(\int L(y, f(x)) P(d(x, y)) + \lambda \|f\|_H^2 \right) = \min! \quad f \in H$$

and

$$\int L(y, f(x)) \left(\frac{\lambda_0}{\lambda} P \right) (d(x, y)) + \lambda_0 \|f\|_H^2 = \min! \quad f \in H.$$

Hence, $f_{P, \lambda} = f_{\frac{\lambda_0}{\lambda} P, \lambda_0} = S\left(\frac{\lambda_0}{\lambda} P\right)$

Sketch of the Proof: Functional Delta-Method

Consider the SVM-functional

$$S : \mathcal{M}_1 \rightarrow H, \quad P \mapsto f_{P, \lambda_0}$$

SVM-functional represents SVM-estimator:

$$f_{\mathbf{D}_n, \lambda_{\mathbf{D}_n}} = S\left(\frac{\lambda_0}{\lambda_{\mathbf{D}_n}} \mathbb{P}_{\mathbf{D}_n}\right)$$

1. Show that $\sqrt{n}(\mathbb{P}_{\mathbf{D}_n} - P)$ converges weakly to a Gaussian process in a suitable space $\ell_\infty(\mathcal{G})$.
2. Show that S is Hadamard-differentiable:
3. Then, it follows from the functional delta-method that

$$\sqrt{n}(f_{\mathbf{D}_n, \lambda_0} - f_{P, \lambda_0}) = \sqrt{n}(S(\mathbb{P}_{\mathbf{D}_n}) - S(P))$$

converges weakly to a Gaussian process.

Sketch of the Proof: Functional Delta-Method

Consider the SVM-functional

$$S : \mathcal{M}_1 \rightarrow H, \quad P \mapsto f_{P, \lambda_0}$$

SVM-functional represents SVM-estimator:

$$f_{\mathbf{D}_n, \lambda_{\mathbf{D}_n}} = S\left(\frac{\lambda_0}{\lambda_{\mathbf{D}_n}} \mathbb{P}_{\mathbf{D}_n}\right)$$

1. Show that $\sqrt{n}(\mathbb{P}_{\mathbf{D}_n} - P)$ converges weakly to a Gaussian process in a suitable space $\ell_\infty(\mathcal{G})$.
2. Show that S is Hadamard-differentiable:
3. Then, it follows from the functional delta-method that

$$\sqrt{n}(f_{\mathbf{D}_n, \lambda_{\mathbf{D}_n}} - f_{P, \lambda_0}) = \sqrt{n}\left(S\left(\frac{\lambda_0}{\lambda_{\mathbf{D}_n}} \mathbb{P}_{\mathbf{D}_n}\right) - S(P)\right)$$

converges weakly to a Gaussian process.

Sketch of the Proof: Functional Delta-Method

Consider the SVM-functional

$$S : \mathcal{M}_f \rightarrow H, \quad \mu \mapsto f_{\mu, \lambda_0}$$

SVM-functional represents SVM-estimator:

$$f_{\mathbf{D}_n, \lambda_{\mathbf{D}_n}} = S\left(\frac{\lambda_0}{\lambda_{\mathbf{D}_n}} \mathbb{P}_{\mathbf{D}_n}\right)$$

1. Show that $\sqrt{n}(\mathbb{P}_{\mathbf{D}_n} - P)$ converges weakly to a Gaussian process in a suitable space $\ell_\infty(\mathcal{G})$.
2. Show that S is Hadamard-differentiable:
3. Then, it follows from the functional delta-method that

$$\sqrt{n}(f_{\mathbf{D}_n, \lambda_{\mathbf{D}_n}} - f_{P, \lambda_0}) = \sqrt{n}\left(S\left(\frac{\lambda_0}{\lambda_{\mathbf{D}_n}} \mathbb{P}_{\mathbf{D}_n}\right) - S(P)\right)$$

converges weakly to a Gaussian process.

References

- ▶ **R. Hable (2010)**: Asymptotic Normality of Support Vector Machines for Classification and Regression. *arXiv:1010.0535v2*
- ▶ **I. Steinwart and A. Christmann (2008)**: *Support vector machines*. Springer, New York.
- ▶ **V. N. Vapnik (1998)**: *Statistical learning theory*. John Wiley & Sons, New York.

The handout to this talk is also available on my homepage

<http://www.staff.uni-bayreuth.de/~btms04/index.html>