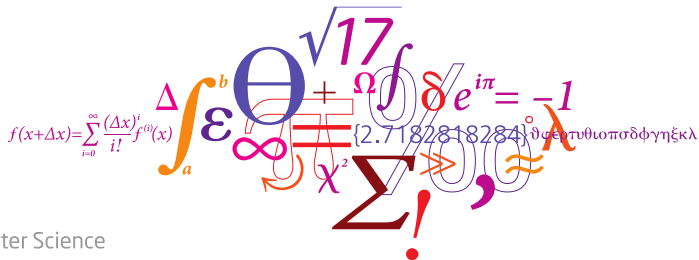


02471 Machine Learning for Signal Processing

Kernel methods – Support vector regression

Tommy Sonne Alstrøm

Cognitive Systems section



Outline

- Course admin
- Last week review
- Anomaly detection and de-noising
- Support vector regression

Material: 11.8.

So far:

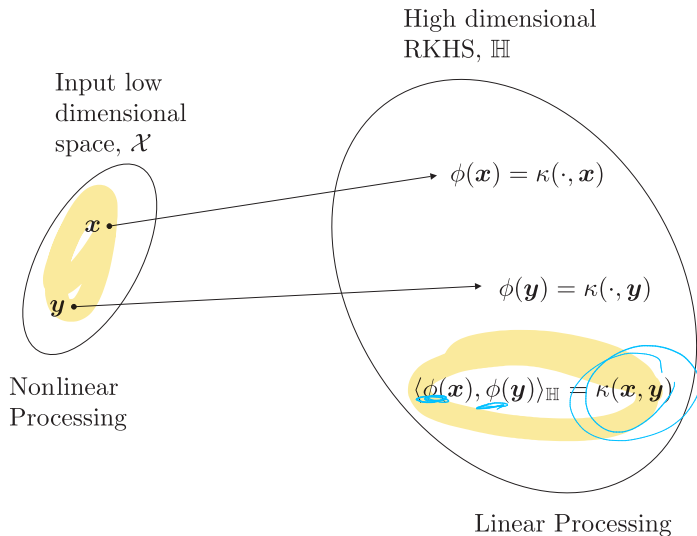
- Parameter estimation [Regularization, biased estimation, mean squared error minimization].
- Signal representations [Time frequency analysis, sparsity aware learning, factor models].
- Filtering signals [Linear regression, adaptive filtering using stochastic gradient decent (LMS, NLMS, RLS), adaptive filtering using regularization].
- Sequential models [hidden markov models, linear dynamical systems, kalman filter].

Next weeks

- Kernel methods [kernel ridge regression, support vector regression].

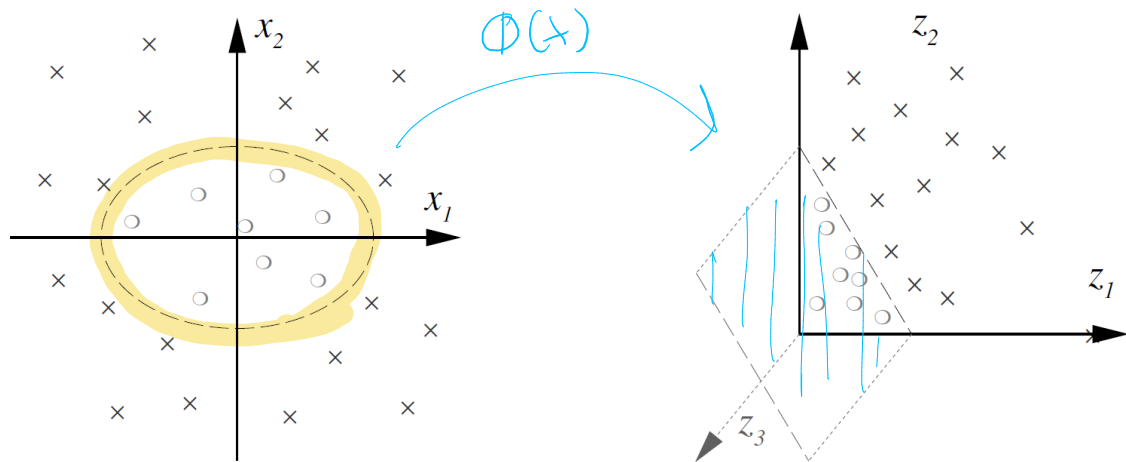
Last week review

Consequences – linear processing



Last week review

This particular mapping leads to linear estimation



Source: Learning With Kernels: Support Vector Machines, Regularization, Optimization, and Beyond — 2002, by Schölkopf, Bernhard; Smola, Alexander J.

Last week review

Algorithms can generally be kernelized, if they can be expressed as inner products

- ① Map (implicitly) the input training data to an RKHS.
- ② Solve the linear estimation task in \mathbb{H} .
- ③ Cast the algorithm in terms of inner products $\langle \mathbf{x}_n, \mathbf{x}_m \rangle$ (\mathbb{R}^l is a Hilbert space).
- ④ Replace each inner product by a kernel evaluation, that is, $\langle \phi(\mathbf{x}_n), \phi(\mathbf{x}_m) \rangle = \kappa(\mathbf{x}_n, \mathbf{x}_m)$.

Simply example:

$$\phi(\mathbf{x}) = \begin{bmatrix} x_1^2 & \sqrt{2}x_1x_2 & x_2^2 \end{bmatrix}^T$$

The corresponding kernel would be

$$\kappa(\mathbf{x}, \mathbf{y}) = (\mathbf{x}^T \mathbf{y})^2, \quad \text{Polynomial kernel}$$

Some kernel properties

Representer theorem in a nutshell

Linearity comes from representer theorem, stating a solution can (under certain circumstances) be written as

$$\hat{f}(\cdot) = \sum_{n=1}^N \theta_n \kappa(\cdot, \mathbf{x}_n)$$

fix today

Kernel properties, assuming $\kappa(\cdot, \mathbf{x}) \in \mathbb{H}$, and has the reproducing property

$$\langle \kappa(\cdot, \mathbf{x}), \kappa(\cdot, \mathbf{y}) \rangle = \kappa(\mathbf{x}, \mathbf{y}) = \kappa(\mathbf{y}, \mathbf{x})$$

$$\langle \phi(\mathbf{x}), \phi(\mathbf{y}) \rangle = \kappa(\mathbf{x}, \mathbf{y}), \quad \text{Kernel Trick}$$

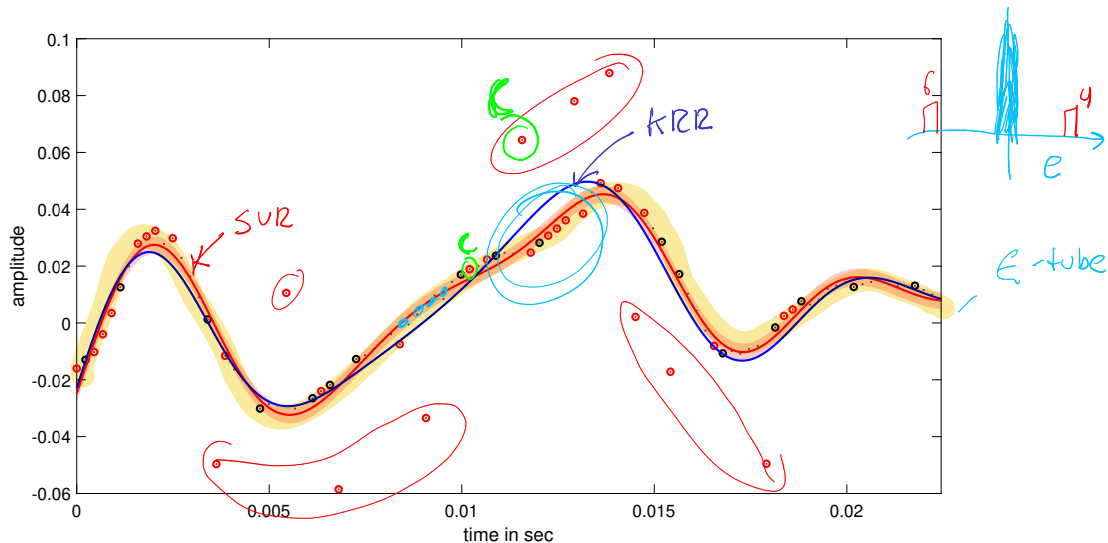
$$\mathcal{K} = \begin{bmatrix} \kappa(\mathbf{x}_1, \mathbf{x}_1) & \cdots & \kappa(\mathbf{x}_1, \mathbf{x}_N) \\ \vdots & \vdots & \vdots \\ \kappa(\mathbf{x}_N, \mathbf{x}_1) & \cdots & \kappa(\mathbf{x}_N, \mathbf{x}_N) \end{bmatrix}$$

$$\mathcal{K} = \mathcal{K}^T$$

$$\mathbf{a} \mathcal{K} \mathbf{a} \geq 0, \quad \mathbf{a} \in \mathbb{R}^l, \quad \mathcal{K} \text{ is positive semi-definite}$$

- Reproducing kernel Hilbert space enables linear processing while obtaining nonlinear signal processing.
- If you limit yourself to apply already proven reproducing kernels, you do not need to understand the theory behind RKHS to apply kernel methods. Use list of kernels.
- Choice of kernel and kernel parameters is critical for performance.
- No restrictions on x , only on the kernels.
- Requires tuning of parameters, but is not “that” sensitive.

Anomaly detection and de-noising
Results - support vector regression



One solution – the Huber loss

Assume the regression task (η_n is i.i.d. noise)

$$y_n = g(\mathbf{x}_n) + \eta_n, \quad n = 1, 2, \dots, N$$

Huber showed that, assuming the noise follows a symmetric pdf, the optimal loss is



Huber loss

DESIGN

$$\mathcal{L}(y, f(\mathbf{x})) = \begin{cases} \epsilon |y - f(\mathbf{x})| - \frac{\epsilon^2}{2}, & \text{if } |y - f(\mathbf{x})| > \epsilon \\ \frac{1}{2} |y - f(\mathbf{x})|^2, & \text{if } |y - f(\mathbf{x})| \leq \epsilon \end{cases}$$

Handwritten notes: $\epsilon |y - f(\mathbf{x})| - \frac{\epsilon^2}{2} = \frac{\epsilon^2}{2}$ and $\frac{1}{2} \epsilon^2 = \frac{\epsilon^2}{2}$

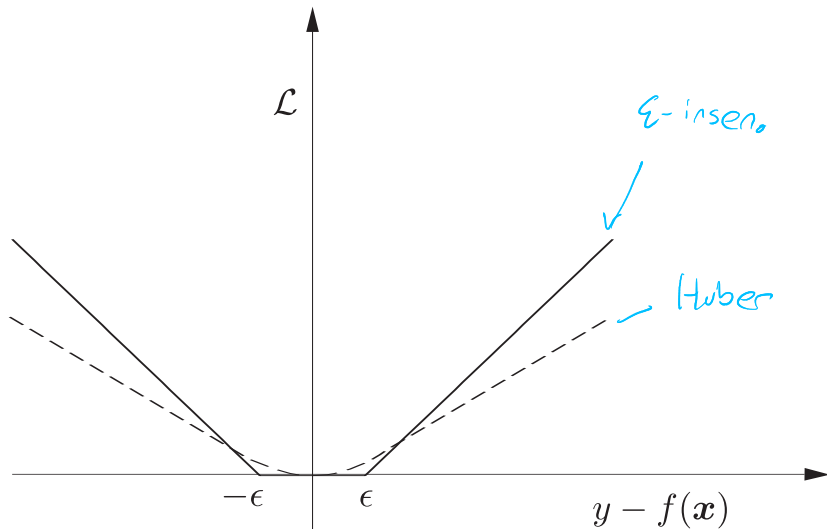
An alternative, and more computationally efficient loss is the ϵ -insensitive loss

ϵ -insensitive loss

$$\mathcal{L}(y, f(\mathbf{x})) = \begin{cases} |y - f(\mathbf{x})| - \epsilon, & \text{if } |y - f(\mathbf{x})| > \epsilon \\ 0, & \text{if } |y - f(\mathbf{x})| \leq \epsilon \end{cases} \Rightarrow \text{SVR}$$

Anomaly detection and de-noising

The loss functions



Support vector regression

Support vector regression

Introduce slack variables

This loss leads to support vector regression through a couple of steps.

First, rewrite the loss by introducing slack variables ξ_n and $\tilde{\xi}_n$:

$$-0.2 \leq 0.1 + 0$$

$$y_n - f(x_n) \leq \epsilon + \xi_n$$

Min ξ_n

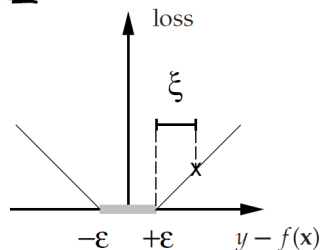
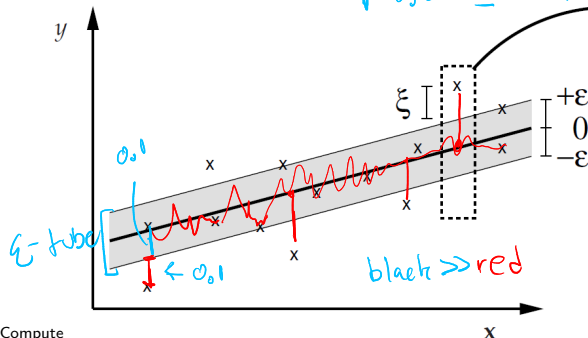
$$\ominus (y_n - f(x_n)) \leq \epsilon + \tilde{\xi}_n$$

$$\xi \geq 0$$

$$\xi_n \cdot \tilde{\xi}_n = 0$$

$$- \text{error} \leq$$

$$+ 0.2 \leq 0.1 + 0.1$$



Regularized minimization of slack variables

13.1.1

DESIGN

func. smoother

model misfit

$$\begin{aligned}
 \arg \min_{\theta, \xi, \tilde{\xi}} J(\theta, \xi, \tilde{\xi}) &:= \frac{1}{2} \|\theta\|^2 + \left(\sum_{n=1}^N \xi_n + \sum_{n=1}^N \tilde{\xi}_n \right) \\
 \text{s.t.} \quad y_n - f(\mathbf{x}_n) &\leq \epsilon + \xi_n \\
 -(y_n - f(\mathbf{x}_n)) &\leq \epsilon + \tilde{\xi}_n \\
 \xi_n &\geq 0 \\
 \tilde{\xi}_n &\geq 0
 \end{aligned}$$

Next step ??

$$\Rightarrow \xi \cdot \tilde{\xi} = 0$$

Lagrangian multipliers

From appendix C.2. If we have the problem:

prob. statement

$$\begin{aligned} \min_{\boldsymbol{\theta}} \quad & J(\boldsymbol{\theta}) \\ \text{s.t.} \quad & f_i(\boldsymbol{\theta}) \geq 0, \quad i = 1, 2, \dots, m \end{aligned}$$

25, errors..

then the Lagrangian is

$$\mathcal{L}(\boldsymbol{\theta}, \boldsymbol{\lambda}) = J(\boldsymbol{\theta}) - \sum_{i=1}^m \lambda_i f_i(\boldsymbol{\theta})$$

And is solved by:

$$\begin{aligned} \frac{\partial}{\partial \boldsymbol{\theta}} \mathcal{L}(\boldsymbol{\theta}, \boldsymbol{\lambda}) \Big|_{\boldsymbol{\theta}=\boldsymbol{\theta}_*} &= \mathbf{0} \\ \lambda_i &\geq 0 \quad i = 1, 2, \dots, m \\ \lambda_i f_i(\boldsymbol{\theta}_*) &= 0 \quad i = 1, 2, \dots, m \end{aligned}$$

The solution

 ξ_i and $\tilde{\xi}_i$ LM. for ξ andAssume $f(\mathbf{x}) = \boldsymbol{\theta}^T \mathbf{x} + \theta_0$, this leads to the solution

$$\hat{\boldsymbol{\theta}}^T = \sum_{n=1}^N (\tilde{\lambda}_n - \lambda_n) \mathbf{x}_n^T$$

which leads to the solution

SVR

$$\hat{y} = f(\mathbf{x}) = \sum_{n=1}^N (\tilde{\lambda}_n - \lambda_n) \mathbf{x}_n^T \mathbf{x} + \hat{\theta}_0$$

 $k(\mathbf{x}_n, \mathbf{x})$

This is an unusual solution! Can we leverage on this?

Do you see any weakness in this solution?

 $\sum_{n \in \text{outside or on the } \epsilon\text{-tube}}$

Support vector regression

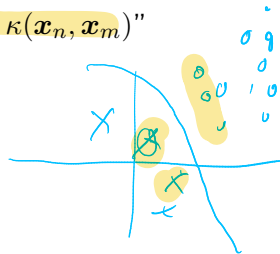
Remember our kernelization recipe

Step 1–4 basically:

“Replace each inner product by a kernel evaluation, that is, $\langle \mathbf{x}_n, \mathbf{x}_m \rangle = \kappa(\mathbf{x}_n, \mathbf{x}_m)$ ”

The prediction equation was

$$\hat{y} = f(\mathbf{x}) = \sum_{n=1}^N (\tilde{\lambda}_n - \lambda_n) \mathbf{x}_n^T \mathbf{x} + \hat{\theta}_0$$



The kernel trick results in

Support vector regression prediction

$$\hat{y} = \sum_{n=1}^N (\tilde{\lambda}_n - \lambda_n) \kappa(\mathbf{x}_n, \mathbf{x}) + \hat{\theta}_0$$

The points for which either λ_n or $\tilde{\lambda}_n$ are non-zero are called the **support vectors**. ex 13.64

Additionally, it can be shown that; $\lambda_n \tilde{\lambda}_n = 0$, and $0 \leq \tilde{\lambda}_n \leq C$, $0 \leq \lambda_n \leq C$

Support vector regression prediction

Sol:

$$\hat{y} = \sum_{n=1}^N (\tilde{\lambda}_n - \lambda_n) \kappa(\mathbf{x}_n, \mathbf{x}) + \hat{\theta}_0$$

Support vector regression minimization

LM \Rightarrow

$$\begin{aligned} \arg \max_{\lambda, \tilde{\lambda}} \quad & \sum_{n=1}^N (\tilde{\lambda}_n - \lambda_n) y_n - \epsilon (\tilde{\lambda}_n - \lambda_n) \\ & - \frac{1}{2} \sum_{n=1}^N \sum_{m=1}^N (\tilde{\lambda}_n - \lambda_n) (\tilde{\lambda}_m - \lambda_m) \kappa(\mathbf{x}_n, \mathbf{x}_m) \\ \text{s.t.} \quad & 0 \leq \tilde{\lambda}_n \leq C, \quad 0 \leq \lambda_n \leq C, \quad n = 1, 2, \dots, N \\ & \sum_{n=1}^N \tilde{\lambda}_n = \sum_{n=1}^N \lambda_n \end{aligned}$$

- Huber loss is the optimal robust loss if the noise follows a symmetric pdf.
- A tractable loss that leads to support vector regression is the ϵ -insensitive loss.
- Support vector regression:
 - seems to not have "training", but store points instead.
 - is solved by Lagrangian multipliers.
 - is very useful for de-noising and for anomaly detection.

Winter is ~~coming~~ here... exam...

Once more, thanks for all who evaluated. We use the feedback!

The course usually lacks TAs, please consider being TA next year.

Thank you for attending.