

Chapter 3: Kernel Methods

Kernel Regression and Support Vector Machines

Maschinelles Lernen 1 -
Grundverfahren WS20/21

Prof. Gerhard Neumann
KIT, Institut für Anthropomatik und Robotik

Learning Outcomes

- What are kernels and how are they useful?
- What do we mean by the “Kernel trick”?
- How to use kernels in regression (using Kernel Regression)?
- How to use kernels in classification (using SVMs)?
- Understand how to obtain dual optimization problems from the primal
- ... and its relation to kernel methods

Today's Agenda!

ML Algorithms

Kernels:

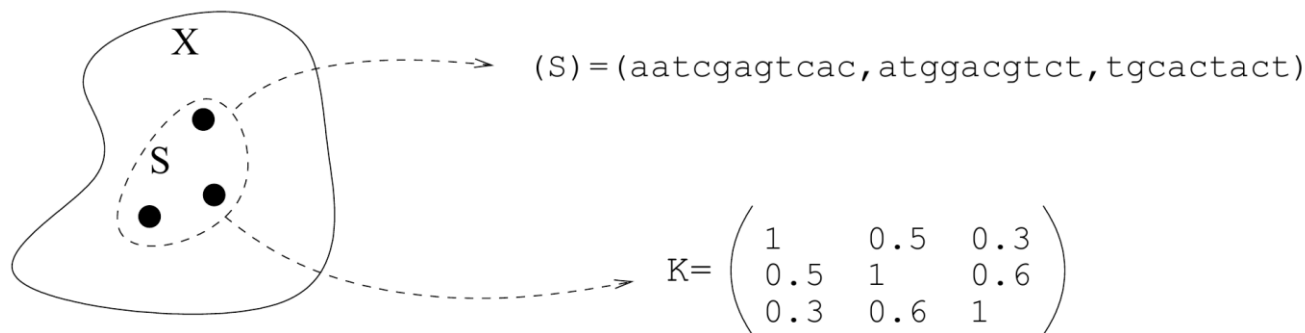
- Definition and properties
- Kernel trick

Kernel Regression:

- Kernel trick for Ridge Regression
- Analytical Solution

What is a kernel?

Representation by point-wise comparisons



- Define a “comparison function” $k : \mathcal{X} \times \mathcal{X} \rightarrow \mathbb{R}$
- Represent a set of points $\mathcal{S} = \{x_1, \dots, x_n\}$ by **the $n \times n$ matrix** $[K]_{ij} = k(x_i, x_j)$

Kernel Matrix

Properties:

- \mathbf{K} is always an $n \times n$ matrix, whatever the nature of data: the same algorithm will work for any type of data (vectors, strings, ...).
- Total modularity between the choice of function k and the choice of the algorithm.
- Poor scalability with respect to the dataset size (n^2 to compute and store \mathbf{K})...
- We will restrict ourselves to a particular class of pairwise comparison functions.

Positive definite kernels

A **positive definite kernel function** k is a function $k : \mathcal{X} \times \mathcal{X} \rightarrow \mathbb{R}$ that is:

(i) **Symmetric:** $\forall \mathbf{x}, \mathbf{x}' : k(\mathbf{x}, \mathbf{x}') = k(\mathbf{x}', \mathbf{x})$

(ii) **Similarity matrix is always positive definite**

$$\mathbf{a}^T \mathbf{K} \mathbf{a} = \sum_{i=1}^n \sum_{j=1}^n a_i a_j k(\mathbf{x}_i, \mathbf{x}_j) \geq 0, \quad \forall \mathbf{a}, \forall S = \{\mathbf{x}_1, \dots, \mathbf{x}_n\}$$

Kernel methods are algorithms that take such matrices as input.

Example: Linear kernel

The **linear kernel** is the simplest kernel for vectors

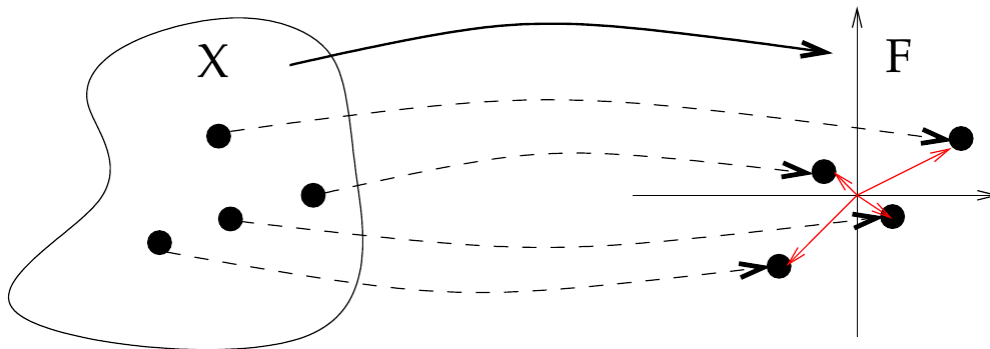
- Its defined by the scalar product:

$$k(\mathbf{x}, \mathbf{x}') = \langle \mathbf{x}, \mathbf{x}' \rangle, \quad \text{where } \langle \cdot, \cdot \rangle \text{ denotes the inner product}$$

- It is **always positive definite**:

$$\sum_{i=1}^n \sum_{j=1}^n a_i a_j \langle \mathbf{x}_i, \mathbf{x}_j \rangle = \left\| \sum_i a_i \mathbf{x}_i \right\|^2 \geq 0$$

Kernels in Feature Spaces



Let $\phi : \mathcal{X} \rightarrow \mathbb{R}^d$ be an arbitrary **feature function**, then $k(\mathbf{x}, \mathbf{x}') = \langle \phi(\mathbf{x}), \phi(\mathbf{x}') \rangle$ defines a **positive definite kernel**.

Proof:
$$\sum_{i=1}^n \sum_{j=1}^n a_i a_j \langle \phi(\mathbf{x}_i), \phi(\mathbf{x}_j) \rangle = \left\| \sum_i a_i \phi(\mathbf{x}_i) \right\|^2 \geq 0$$

Kernels as inner products

Theorem (Aransjan 1950):

k is a **positive definite kernel** on the set \mathcal{X} if and only if there exists a **feature space** \mathcal{H} and a feature mapping

$$\phi : \mathcal{X} \rightarrow \mathcal{H}$$

such that for any $x, x' \in \mathcal{X}$:

$$k(x, x') = \langle \phi(x), \phi(x') \rangle$$

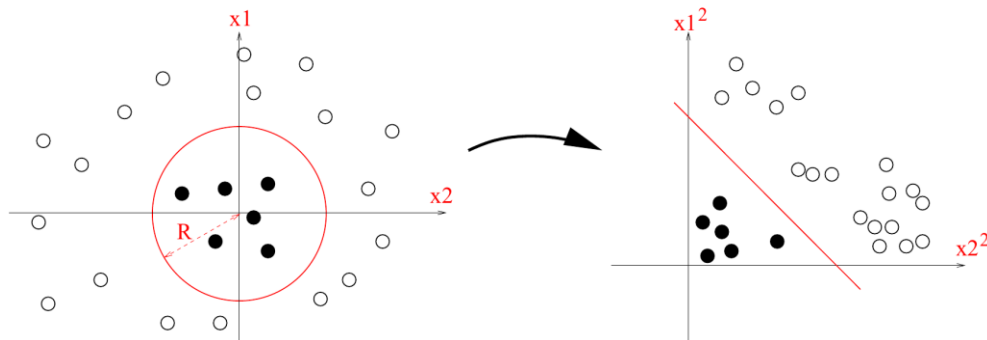
➤ **Every p.d. kernel comes with an associated feature space!**

Example: polynomial kernel

For $\mathbf{x} = [x_1, x_2]^T$, let $\phi(\mathbf{x}) = [x_1^2, \sqrt{2}x_1x_2, x_2^2]$

The kernel is defined by:

$$\begin{aligned}k(\mathbf{x}, \mathbf{x}') &= x_1^2x_1'^2 + 2x_1x_2x_1'x_2' + x_2^2x_2'^2 \\&= (x_1x_1' + x_2x_2')^2 \\&= \langle \mathbf{x}, \mathbf{x}' \rangle^2\end{aligned}$$



Kernel for polynomials of degree d :

$$k(\mathbf{x}, \mathbf{x}') = \langle \mathbf{x}, \mathbf{x}' \rangle^d$$

Example: Gaussian Kernel

The Gaussian kernel is defined by:

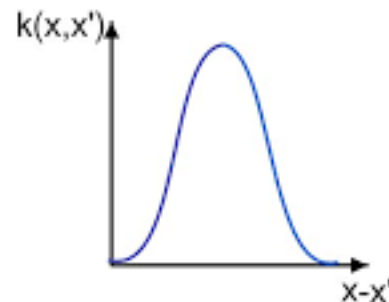
$$k(\mathbf{x}, \mathbf{y}) = \exp\left(-\frac{\|\mathbf{x} - \mathbf{y}\|^2}{2\sigma^2}\right)$$

- where σ is the bandwidth parameter

Often also called:

- Radial basis function kernel (RBF)
- Squared exponential kernel

It is the **most used kernel for kernel methods**



Is the Gaussian kernel a valid p.d. kernel?

Remember: If we can show that the kernel is a valid product of feature vectors, then it is p.d.

- Consider the following feature function:

$$\phi_{\mu}(\mathbf{x}) = 1/Z \exp\left(-\frac{\|\mathbf{x} - \mu\|^2}{4\sigma^2}\right), \quad \forall \mu \in \mathbb{R}^d$$

- I.e. we have an **infinite amount** of features (for every possible center μ)
- Z is a normalization constant (which we will ignore)

Inner product:

- Inner product becomes an integral** due to infinite amount of dimensions

$$\langle \phi_{\mu}(\mathbf{x}), \phi_{\mu}(\mathbf{y}) \rangle = \int \phi_{\mu}(\mathbf{x}) \phi_{\mu}(\mathbf{y}) d\mu$$

Is the Gaussian kernel a valid p.d. kernel?

Inner product:

$$\begin{aligned}\langle \phi_{\mu}(\mathbf{x}), \phi_{\mu}(\mathbf{y}) \rangle &= \int \phi_{\mu}(\mathbf{x}) \phi_{\mu}(\mathbf{y}) d\mu \\ &\propto \int \exp\left(-\frac{\|\mathbf{x} - \mu\|^2}{4\sigma^2}\right) \exp\left(-\frac{\|\mathbf{y} - \mu\|^2}{4\sigma^2}\right) d\mu \dots \text{ignore normalization constants} \\ &\propto \int \mathcal{N}(\mu|\mathbf{x}, \sigma^2/2\mathbf{I}) \mathcal{N}(\mu|\mathbf{y}, \sigma^2/2\mathbf{I}) d\mu \dots \text{product of 2 Gaussians (see Gaussian identities)} \\ &= \mathcal{N}(\mathbf{x}|\mathbf{y}, \sigma^2\mathbf{I}) \underbrace{\int \mathcal{N}(\mu|\dots, \dots) d\mu}_{=1} \\ &\propto \exp\left(-\frac{\|\mathbf{x} - \mathbf{y}\|^2}{2\sigma^2}\right) = k(\mathbf{x}, \mathbf{y})\end{aligned}$$

I.e. the Gaussian kernel is the inner product of 2 infinite dimensional feature vectors!

Kernel Trick

So why do we do this?

- Kernels can be used for all feature based algorithms that can be rewritten such that they contain inner products of feature vectors
 - This is true for almost all feature based algorithms (Linear regression, SVMs, ...)
 - This is called the **Kernel Trick**
- Kernels can be used to map the data x in an **infinite dimensional feature** space (i.e., a function space)
 - The feature vector **never** has to be represented **explicitly**
 - As long as we can **evaluate the inner product** of two feature vectors
- Hence, we obtain a **more powerful representation** than standard linear feature models

A few kernel identities

Let $\Phi_X = \begin{bmatrix} \phi(\mathbf{x}_1)^T \\ \vdots \\ \phi(\mathbf{x}_N)^T \end{bmatrix} \in \mathbb{R}^{N \times d}$ then the following identities hold:

- **Kernel matrix:** $\mathbf{K} = \Phi_X \Phi_X^T$
 - Check: $[\mathbf{K}]_{ij} = \phi(\mathbf{x}_i)^T \phi(\mathbf{x}_j) = k(\mathbf{x}_i, \mathbf{x}_j)$
- **Kernel vector:** $\mathbf{k}(\mathbf{x}^*) = \begin{bmatrix} k(\mathbf{x}_1, \mathbf{x}^*) \\ \vdots \\ k(\mathbf{x}_N, \mathbf{x}^*) \end{bmatrix} = \begin{bmatrix} \phi(\mathbf{x}_1)^T \phi(\mathbf{x}^*) \\ \vdots \\ \phi(\mathbf{x}_N)^T \phi(\mathbf{x}^*) \end{bmatrix} = \Phi_X \phi(\mathbf{x}^*)$

Today's Agenda!

ML Algorithms

Kernels:

- Definition and properties
- Kernel trick

Kernel Regression:

- Kernel trick for Ridge Regression
- Analytical Solution

Kernel ridge Regression

Recap: Ridge Regression

- Squared error function + L2 regularization
- Linear feature space
- Not directly applicable in infinite dimensional feature spaces

Objective:

$$L_{\text{ridge}} = \underbrace{(\mathbf{y} - \Phi \mathbf{w})^T (\mathbf{y} - \Phi \mathbf{w})}_{\text{sum of squared errors}} + \lambda \underbrace{\mathbf{w}^T \mathbf{w}}_{L_2 \text{ regularization}}$$

Solution:

$$\mathbf{w}_{\text{ridge}}^* = \underbrace{(\Phi^T \Phi + \lambda \mathbf{I})^{-1}}_{d \times d \text{ matrix inversion}} \Phi^T \mathbf{y} \quad \text{Matrix inversion infeasible in infinite dimensions}$$

Kernel Ridge regression

We can apply the “kernel trick”:

- Rewrite solution as inner products of the feature space!
- We can do this by using the following matrix identity

$$(I + AB)^{-1}A = A(I + BA)^{-1}$$

- “Searle set of identities”, The Matrix Cookbook

$$w^* = \underbrace{(\Phi^T \Phi + \lambda I)^{-1}}_{d \times d \text{ matrix inversion}} \Phi^T y = \Phi^T \underbrace{(\Phi \Phi^T + \lambda I)^{-1}}_{N \times N \text{ matrix inversion}} y$$

- With $A = \Phi^T$ and $B = \Phi$

Kernel ridge regression

The “kernelized” solution is given by:

$$\mathbf{w}^* = \Phi^T \underbrace{(\Phi\Phi^T + \lambda I)^{-1}}_{N \times N \text{ matrix inversion}} \mathbf{y} = \Phi^T \underbrace{(K + \lambda I)^{-1} \mathbf{y}}_{\boldsymbol{\alpha}} = \Phi^T \boldsymbol{\alpha}$$

- Instead of inverting a $d \times d$ matrix, we can now invert an $N \times N$ matrix
- Is beneficial for $d \gg N$ (e.g., infinite)
- Still, $\mathbf{w}^* \in \mathbb{R}^d$ is potentially infinite dimensional and can not be represented

Yet, we can evaluate the function f that is specified by \mathbf{w}^* :

$$f(\mathbf{x}) = \phi(\mathbf{x})^T \mathbf{w}^* = \phi(\mathbf{x})^T \Phi^T \boldsymbol{\alpha} = \mathbf{k}(\mathbf{x})^T \boldsymbol{\alpha} = \sum_i \alpha_i k(\mathbf{x}_i, \mathbf{x})$$

Examples and comparison to RBF regression

For a Gaussian kernel, the prediction corresponds to

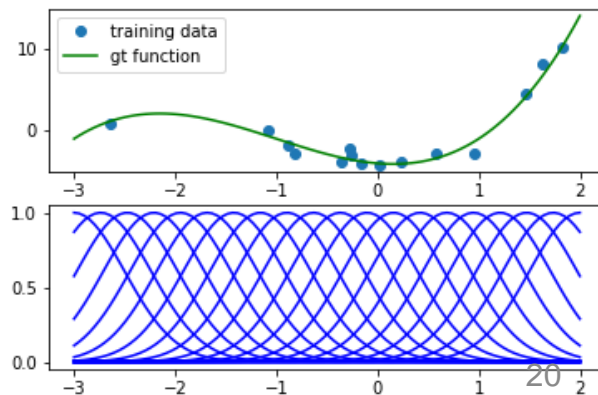
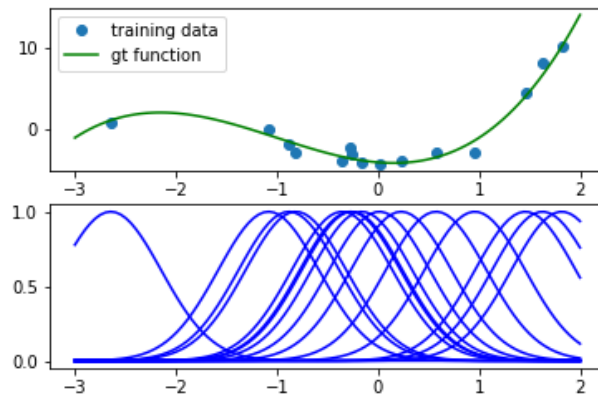
$$f(\mathbf{x}) = \sum_i \alpha_i k(\mathbf{x}_i, \mathbf{x}) = \sum_i \alpha_i \exp\left(-\frac{\|\mathbf{x} - \mathbf{x}_i\|^2}{2\sigma^2}\right)$$

- The kernel allows setting the **centres adaptively to the available data!**
- One centre per data-point

Comparison: Linear regression with radial basis function (RBF) features

$$f(\mathbf{x}) = \sum_i w_i \phi_i(\mathbf{x}_i) = \sum_i w_i \exp\left(-\frac{\|\mathbf{x} - \boldsymbol{\mu}_i\|^2}{2\sigma^2}\right)$$

$\boldsymbol{\mu}_i \dots i^{\text{th}}$ center location (fixed)



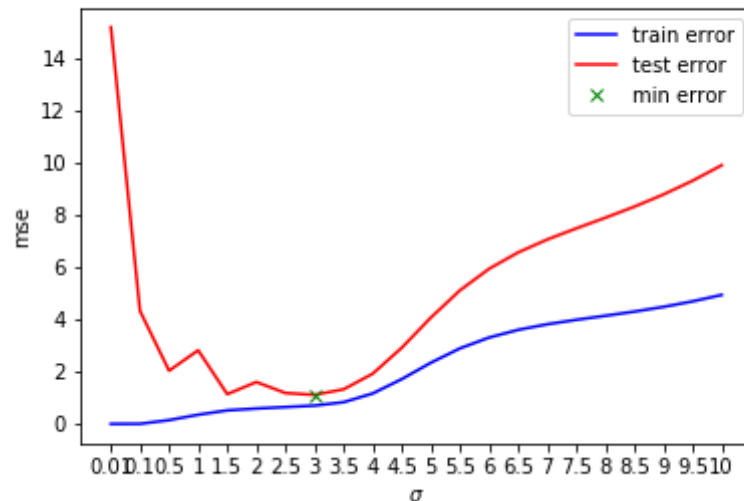
Selecting the hyper-parameters

- The parameters of the kernel, e.g., sigma in

$$k(\mathbf{x}, \mathbf{y}) = \exp\left(-\frac{\|\mathbf{x} - \mathbf{y}\|^2}{2\sigma^2}\right)$$

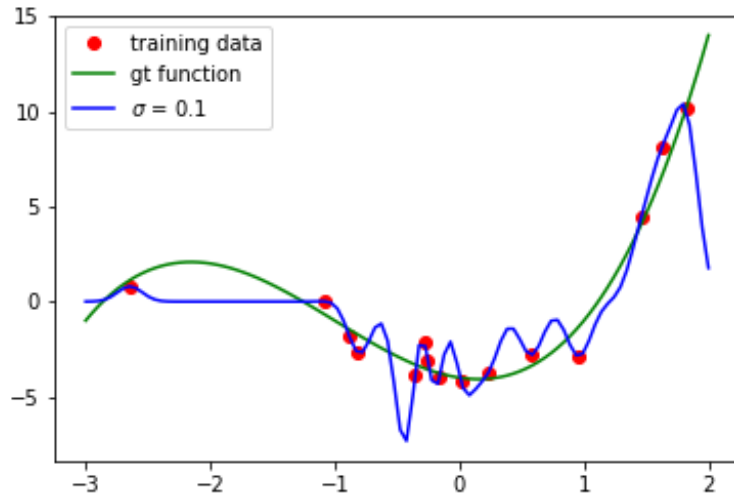
are called **hyper-parameters**.

- Choosing them is again a model-selection problem that can be solved via cross-validation.

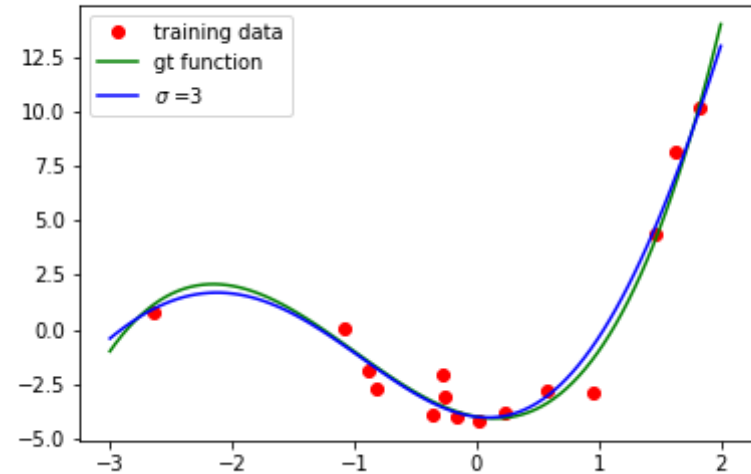


Different bandwidth factors

Overfitting



Good fit



Summary: Kernel ridge regression

The solution for kernel ridge regression is given by

$$f^*(x) = k(x)^T (\mathbf{K} + \lambda \mathbf{I})^{-1} \mathbf{y}$$

- No evaluations of the feature vectors needed
- Only pair-wise scalar products (evaluated by the kernel)
- Need to invert a $N \times N$ matrix (can be costly)

Note:

- We have to **store all samples** in kernel-based methods
 - **Computationally expensive** (matrix inverse is $O(n^{2.376})$) !
- Hyper-parameters of the method are given by the kernel-parameters
 - Can be optimized on validation-set
- Very **flexible function representation**, only few hyper-parameters

Takeaway messages

What have we learned today?

- Kernels estimate the similarity between samples
- They represent an **inner product** in a feature space
 - Allows to use potentially infinite dimensional
 - That's ok due to the kernel trick and regularization
- Many standard ML algorithms can be **"kernelized"**
 - I.e. rewritten in terms of inner products
 - **Regression:** Kernel Ridge regression, Gaussian Processes (to be covered), Support Vector Regression (not covered)
 - **Classification:** SVMs, Kernel Logistic Regression (not covered)
- ✓ Very flexible representation that adapts to the complexity of the data
- ✓ Works well with small data sets
- ✗ Hard to scale to more complex problems



Self-test questions

You should know now:

- What is the definition of a kernel and its relation to an underlying feature space?
- Why are kernels more powerful than traditional feature-based methods?
- What do we mean by the kernel trick?
- How do we apply the kernel trick to ridge regression?