

Cognitive Algorithms Lecture 4

Kernel Methods

Klaus-Robert Müller, **Ali Hashemi, Lorenz Vaitl,**
Augustin Krause, Joanina Obersdorff, Ken Schreiber

Berlin Institute of Technology
Dept. Machine Learning

Linear Regression

The most popular loss function to optimize \mathbf{w}
is the **least-square error** [Gauß, 1809; Legendre, 1805]

$$\mathcal{E}_{lsq}(\mathbf{w}) = \sum_{i=1}^N (y_i - \mathbf{w}^\top \mathbf{X}_i)^2$$



C.F. Gauß (1777-1855)



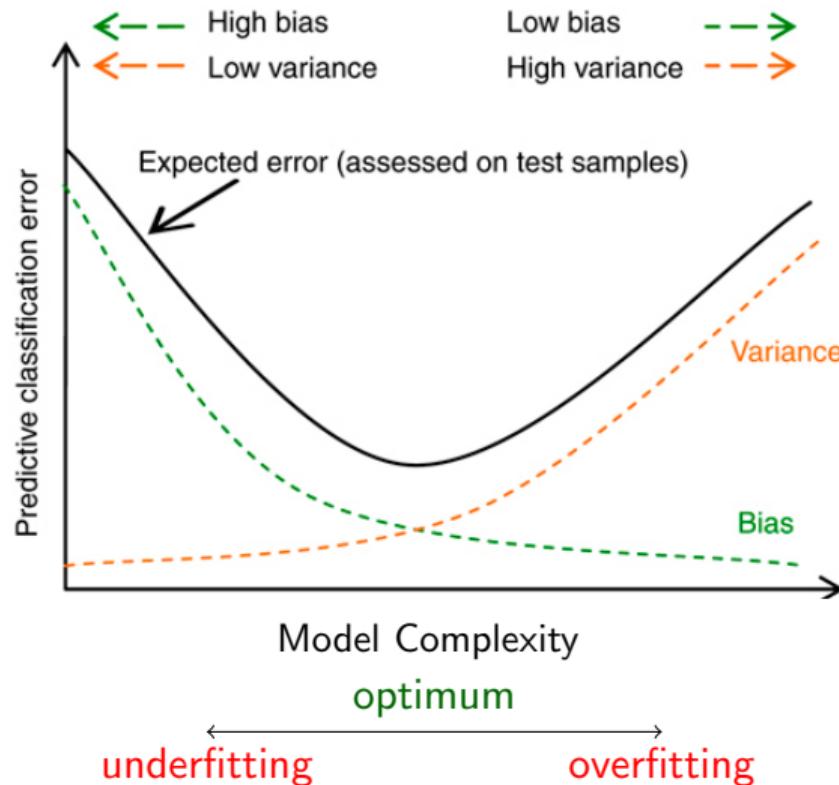
A.M. Legendre (1752-1833)

Gauss-Markov-Theorem

Under the model assumption $y = \mathbf{w}^\top \cdot \mathbf{x} + \epsilon$ with uncorrelated noise ϵ , our ordinary least squares estimator $\hat{\mathbf{w}} = (X X^\top)^{-1} X y$ is the Best Linear Unbiased Estimator (BLUE), i.e. the minimum variance unbiased estimator that is linear in the y .

But: in some cases biased estimators with lower variance might be more suitable

The Bias-Variance Tradeoff

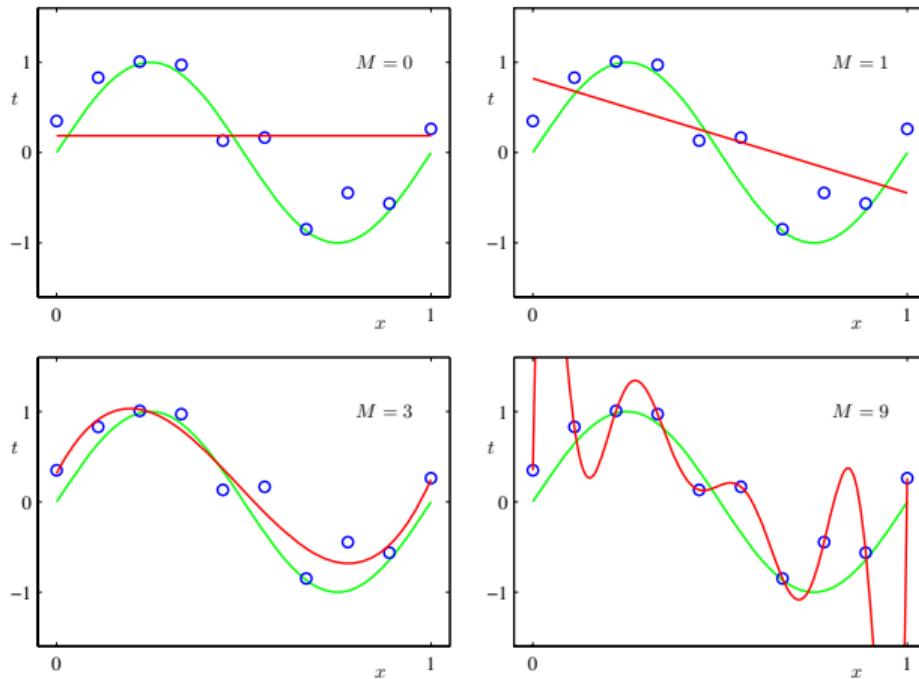


Example: Linear regression for a polynomial function

$$\hat{y}(x) = w_0 + w_1 \cdot x^1 + \dots + w_M \cdot x^M$$

→ Here we are using the basis function

$\phi_M(x) = (x^0, x^1, \dots, x^M)$,
which has a $(M + 1)$ -dimensional feature space \mathcal{F}



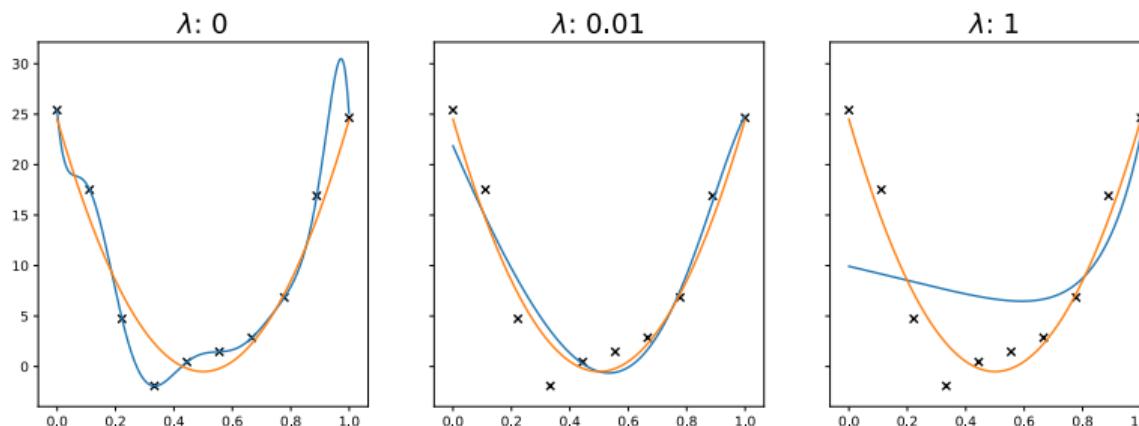
[Bishop, 2007]

Ridge Regression

Control the complexity of the solution \mathbf{w} .

This is done by constraining the norm of \mathbf{w} ,

$$\mathcal{E}_{RR}(\mathbf{w}) = \|\mathbf{y} - \mathbf{w}^\top \mathbf{X}\|^2 + \lambda \|\mathbf{w}\|^2$$



Ridge Regression

Computing the derivative w.r.t. w yields

$$\frac{\partial \mathcal{E}_{RR}(\mathbf{w})}{\partial \mathbf{w}} = -2X\mathbf{y}^\top + 2XX^\top \mathbf{w} + \lambda 2\mathbf{w}.$$

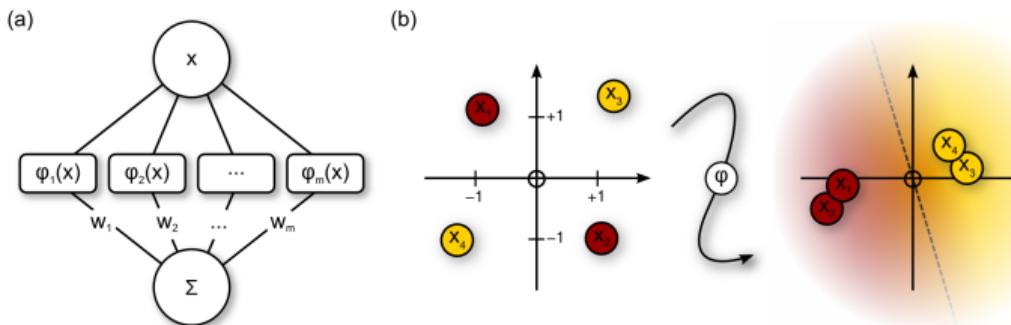
Setting the gradient to zero and rearranging terms the optimal \mathbf{w} is

$$\begin{aligned} 2XX^\top \mathbf{w} + \lambda 2\mathbf{w} &= 2X\mathbf{y}^\top \\ (XX^\top + \lambda I)\mathbf{w} &= X\mathbf{y}^\top \\ \mathbf{w} &= (XX^\top + \lambda I)^{-1}X\mathbf{y}^\top \end{aligned}$$

⇒ Biased estimator, but smaller variance

[Hoerl and Kennar, 1970; Tychonoff, 1943]

Kernelizing linear methods



[Jäkel et al., 2009]

1. Map the data into a (high dimensional) feature space, $\mathbf{x} \mapsto \varphi(\mathbf{x})$
2. Look for linear relations in the feature space

Definition Kernel

Definition

We define a Kernel function k as a dot product in feature space \mathcal{F} where

$$\begin{aligned}\varphi : \mathcal{X} &\rightarrow \mathcal{F} \\ x &\mapsto \varphi(x)\end{aligned}$$

and so

$$\begin{aligned}k : \mathcal{X} \times \mathcal{X} &\rightarrow \mathbb{R} \\ k(x_i, x_j) &= \varphi(x_i)^T \cdot \varphi(x_j)\end{aligned}$$

Kernel Trick (Kernel Substitution)

For any algorithm that can be formulated such that the input vector \mathbf{x} enters only in terms of scalar products $\mathbf{x}^T \cdot \mathbf{x}'$, we can replace each scalar product by a kernel $k(\mathbf{x}, \mathbf{x}') = \varphi(\mathbf{x})^T \cdot \varphi(\mathbf{x}')$.

Why do it?

For $\mathbf{x}, \mathbf{x}' \in \mathbb{R}^{d \times 1} \Rightarrow \mathbf{x}^T \cdot \mathbf{x}' \in \mathbb{R}^{1 \times 1}$ regardless of d

For $\varphi(\mathbf{x}) \in \mathbb{R}^{\tilde{d} \times 1}$ instead of \mathbf{x} with typically $\tilde{d} \gg d$

\Rightarrow we get more complex (powerful) models

By using kernel $k(\mathbf{x}, \mathbf{x}') \in \mathbb{R}^{1 \times 1}$

We do not need to explicitly calculate high-dimensional $\varphi(\mathbf{x})$

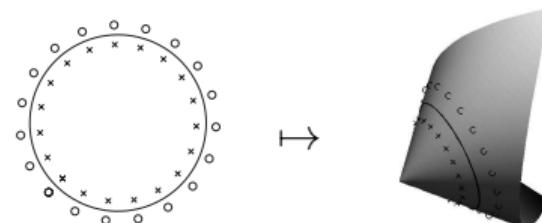
Example kernel

$$\varphi : \mathbf{x} = (x_1, x_2)^\top \mapsto \varphi(\mathbf{x}) = \begin{pmatrix} x_1^2 \\ \sqrt{2}x_1x_2 \\ x_2^2 \end{pmatrix}$$

The corresponding kernel:

$$\begin{aligned} k(\mathbf{x}, \mathbf{y}) &= \varphi(\mathbf{x})^\top \cdot \varphi(\mathbf{y}) \\ &= (x_1^2, \sqrt{2}x_1x_2, x_2^2) \cdot (y_1^2, \sqrt{2}y_1y_2, y_2^2)^\top \\ &= x_1^2y_1^2 + 2x_1x_2y_1y_2 + x_2^2y_2^2 \\ &= (x_1y_1 + x_2y_2)^2 = (\mathbf{x}^\top \mathbf{y})^2 \end{aligned}$$

Visualizing $\varphi(\mathbf{x})$



→ With $k(\cdot, \cdot)$ we implicitly work in \mathbb{R}^3 , but only operate in \mathbb{R}^2

Definition (Positive semi-definite symmetric kernels)

A kernel

$$k \nabla : \mathcal{X} \times \mathcal{X} \rightarrow \mathbb{R}$$

is said to be positive semi-definite symmetric

if for any $\{x_1, \dots, x_n\} \subseteq \mathcal{X}$, the matrix $K = [k(x_i, x_j)]_{ij} \in \mathbb{R}^{n \times n}$ is symmetric positive semidefinite.¹

Symmetric Matrix:

$$A = A^T$$

Positive semi-definite Matrix:

$$x^T A x \geq 0 \quad \forall x$$

alternatively, all eigenvalues λ of A non-negative.

¹For ease of notation we may use $K(X, X') = [K(x_i, x'_j)]_{ij}$ for describing the matrix of the kernel-function evaluated on all sample-pairs. $K(X, X)$ is often called the *Gram matrix of X*.

Mercer's Theorem [Mercer, 1909]

Non-technical formulation:

If kernel is positive semi-definite symmetric then one can find a map φ such that
 $k(x, x') = \varphi(x)^T \varphi(x')$ for almost all $x \in \mathcal{X}$

construct a feature space F with a scalar product
↙

- To see if your kernel is valid, show it is positive semi-definite symmetric!
- You can construct kernels from other kernels, e.g. by sum, product or exponentiating

Alternatively show $k(x, x') = \varphi(x)^T \varphi(x')$

Some Popular Kernel Functions

Linear Kernel

$$k(\mathbf{x}_i, \mathbf{x}_j) = \mathbf{x}_i^\top \mathbf{x}_j$$

Polynomial Kernel

$$k(\mathbf{x}_i, \mathbf{x}_j) = (\mathbf{x}_i^\top \mathbf{x}_j + c)^p$$

Gaussian Kernel/ RBF (more on this later)

$$k(\mathbf{x}_i, \mathbf{x}_j) = e^{\|\mathbf{x}_i - \mathbf{x}_j\|^2 / -2\sigma^2}$$

Note that we are never directly operating in the feature space \mathcal{F} !

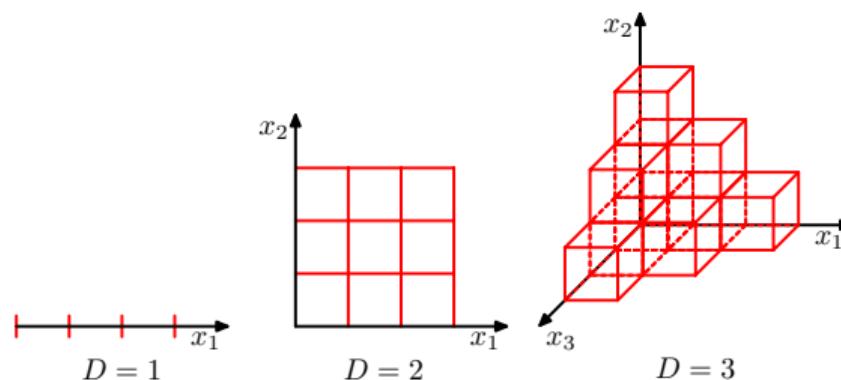
The curse of dimensionality

A big problem with high dimensional feature spaces

When the dimensionality increases, the volume of the space increases so fast that the available data becomes sparse

稀少

The amount of data needed for a reliable result often grows exponentially with the dimensionality



[Bishop, 2007]

How do we find the optimal weight

$$\mathbf{w} \in \mathbb{R}^{\tilde{d}} ?$$

Seems impossible with $\tilde{d} \gg n$



Representer Theorem [Kimeldorf and Wahba, 1971]

The minimizing function f^* of a (regularized) error function on some training data \mathbf{x}_i can be written in terms of the kernel k as

$$f^*(\mathbf{x}) = \sum_{i=1}^N \alpha_i k(\mathbf{x}, \mathbf{x}_i).$$

i.e. only scalar products between $\varphi(\mathbf{x})$ and the training data $\varphi(\mathbf{x}_i)$ are relevant.
For the model $f(\mathbf{x}) = \mathbf{w}^\top \varphi(\mathbf{x})$ this means that \mathbf{w} can be written as

$$\mathbf{w} = \sum_{i=1}^N \alpha_i \varphi(\mathbf{x}_i).$$

Kernelizing Algorithms

From before:

Kernel Trick (Kernel Substitution)

For any algorithm that can be formulated such that the input vector \mathbf{x} enters only in terms of scalar products $\mathbf{x}^T \cdot \mathbf{x}'$, we can replace each scalar product by a kernel $k(\mathbf{x}, \mathbf{x}') = \varphi(\mathbf{x})^T \cdot \varphi(\mathbf{x}')$.

So how can we formulate an algorithm like this?

- In the following we will kernelize Ridge Regression as an example on how to kernelize linear models
- We'll recast the problem into an equivalent dual representation (can be done with many linear models)

Recap: Linear Regression

The Linear Regression model in matrix notation then becomes

$$\hat{\mathbf{y}} = \mathbf{w}^\top \mathbf{X}.$$

Linear Regression minimizes the least-squares loss function

$$\mathcal{E}_{lsq}(\mathbf{w}) = \sum_{i=1}^n (y_i - \mathbf{w}^\top \mathbf{x}_i)^2 = \|\mathbf{y} - \mathbf{w}^\top \mathbf{X}\|^2$$

$$\begin{array}{c} N \rightarrow \\ \textcolor{red}{\hat{\mathbf{y}}} \end{array} = \begin{array}{c} D \rightarrow \\ \mathbf{w}^\top \end{array} \cdot \begin{array}{c} N \rightarrow \\ D \downarrow \\ \mathbf{X} \end{array}$$

Beware: \mathbf{y} is a row vector, while \mathbf{w} is a column vector.

Recap: Ridge Regression

Linear ridge regression looks for a linear combination of features \mathbf{w} that minimizes the prediction error and has a small norm

We can write this term in several equivalent ways:

$$\begin{aligned}\mathcal{E}_{RR}(\mathbf{w}) &= \sum_n (y_n - \mathbf{w}^\top \mathbf{x}_n)^2 + \lambda \sum_d w_d^2 \\ &= \|\mathbf{y} - \mathbf{w}^\top \mathbf{X}\|^2 + \lambda \|\mathbf{w}\|^2 \\ &= \mathbf{y}^\top \mathbf{y} - 2\mathbf{w}^\top \mathbf{X} \mathbf{y}^\top + \mathbf{w}^\top \mathbf{X} \mathbf{X}^\top \mathbf{w} + \lambda \mathbf{w}^\top \mathbf{w}\end{aligned}$$

Kernel Trick: We can use kernels if algorithm depends on training data X and test sample \mathbf{x} only through scalar products.

From Linear to Kernel Ridge Regression

$$\mathcal{E}_{RR}(\mathbf{w}) = \mathbf{y}\mathbf{y}^T - 2\mathbf{w}^T X \mathbf{y}^T + \mathbf{w}^T X X^T \mathbf{w} + \lambda \mathbf{w}^T \mathbf{w}$$

Computing the derivative w.r.t. \mathbf{w} yields

$$\frac{\partial \mathcal{E}_{RR}(\mathbf{w})}{\partial \mathbf{w}} = -2X\mathbf{y}^T + 2X X^T \mathbf{w} + 2\lambda \mathbf{w}$$

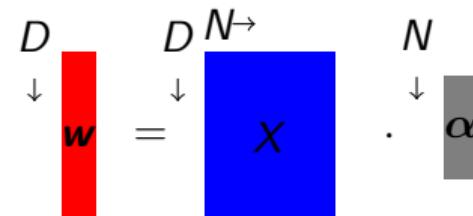
Setting the gradient to 0 and rearranging terms the optimal \mathbf{w} satisfies

$$\mathbf{w} = X \underbrace{\frac{1}{\lambda} (\mathbf{y}^T - X^T \mathbf{w})}_{:= \boldsymbol{\alpha} (\in \mathbb{R}^{n \times 1})} = \sum_i^n \alpha_i \mathbf{x}_i$$

→ we showed that the Representer Theorem is valid for RR!

From Linear to Kernel Ridge Regression

$$\mathcal{E}_{RR}(\mathbf{w}) = \mathbf{y}\mathbf{y}^\top - 2\mathbf{w}^\top \mathbf{X}\mathbf{y}^\top + \mathbf{w}^\top \mathbf{X}\mathbf{X}^\top \mathbf{w} + \lambda \mathbf{w}^\top \mathbf{w}$$
$$\mathbf{w} = \mathbf{X}\boldsymbol{\alpha}$$



Let's look at the error function \mathcal{E}_{RR} w.r.t. $\boldsymbol{\alpha}$

$$\mathcal{E}_{RR}(\boldsymbol{\alpha}) = \mathbf{y}\mathbf{y}^\top - 2\boldsymbol{\alpha}^\top \underbrace{\mathbf{X}^\top \mathbf{X}}_K \mathbf{y}^\top + \boldsymbol{\alpha}^\top \underbrace{\mathbf{X}^\top \mathbf{X} \mathbf{X}^\top}_K \mathbf{X} \boldsymbol{\alpha} + \lambda \boldsymbol{\alpha}^\top \underbrace{\mathbf{X}^\top \mathbf{X}}_K \mathbf{X} \boldsymbol{\alpha}$$

We call this form **dual representation**

Only scalar products \rightarrow we can put in kernels:

$$\mathbf{x}_i^\top \mathbf{x}_j \rightarrow \varphi(\mathbf{x}_i)^\top \varphi(\mathbf{x}_j) = k(\mathbf{x}_i, \mathbf{x}_j) = K_{ij}$$

We write $k(\mathbf{X}, \mathbf{X})$ as K

Kernel Ridge Regression:

$$\alpha = \frac{1}{\lambda}(\mathbf{y}^\top - \varphi(X_{train})^\top \mathbf{w})$$

$$\lambda \alpha = \mathbf{y}^\top - \varphi(X_{train})^\top \varphi(X_{train}) \alpha$$

$$\mathbf{y}^\top = (\varphi(X_{train})^\top \varphi(X_{train}) + \lambda I) \alpha$$

$$\alpha = (\varphi(X_{train})^\top \varphi(X_{train}) + \lambda I)^{-1} \mathbf{y}^\top$$

$$\alpha = (K + \lambda I)^{-1} \mathbf{y}^\top$$

This α minimizes $\mathcal{E}_{RR}(\alpha)$.

Ridge Regression

Train \mathbf{w} which minimizes $\mathcal{E}_{RR}(\mathbf{w})$:

$$\mathbf{w} = (\varphi(X) \varphi(X)^\top + \lambda I)^{-1} \varphi(X) \mathbf{y}^\top$$

Training KRR

$$\text{Elements for RR: } \begin{pmatrix} \sum_i \varphi_1(x_i) \varphi_1(x_i) & \cdots \\ \varphi_1(x) \cdot \varphi_1(x)^\top & \begin{pmatrix} \cdots \\ \sum_i \varphi_n(x_i) \varphi_n(x_i) & \cdots \\ \vdots \end{pmatrix} \end{pmatrix}$$

For KRR we write $\varphi(\mathbf{x}) \in \mathbb{R}^{\tilde{D}}$ instead of \mathbf{x}

$$\text{We defined } \alpha_i = \frac{1}{\lambda}(y_i - \varphi(\mathbf{x}_i)^\top \mathbf{w})$$

$$\text{Optimal } \mathbf{w} = \varphi(X_{train}) \alpha$$

KRR trains $\alpha \in \mathbb{R}^n$, RR trains $\mathbf{w} \in \mathbb{R}^{\tilde{D}}$

For $\varphi(\mathbf{x})^\top \varphi(\mathbf{x}') = k(\mathbf{x}, \mathbf{x}')$ the two models are equivalent (but not the runtime complexity)

$$\begin{matrix} \tilde{D} \\ \downarrow \\ \mathbf{w} \end{matrix} = \begin{matrix} \tilde{D} \\ \downarrow \\ \varphi(X) \end{matrix} \xrightarrow{N} \begin{matrix} N \\ \downarrow \\ \alpha \end{matrix}$$

Predictions for new data \mathbf{x}_{new}

$$\begin{aligned}y_{new} &= \mathbf{w}^T \varphi(\mathbf{x}_{new}) \\&= (\varphi(\mathcal{X}_{train})\boldsymbol{\alpha})^T \varphi(\mathbf{x}_{new}) \\&= \boldsymbol{\alpha}^T \varphi(\mathcal{X}_{train})^T \varphi(\mathbf{x}_{new}) \\&= \boldsymbol{\alpha}^T k(\mathcal{X}_{train}, \mathbf{x}_{new}) \\&= \mathbf{y}_{train}(K + \lambda I)^{-1} k(\mathcal{X}_{train}, \mathbf{x}_{new})\end{aligned}$$

Optimal $\boldsymbol{\alpha} = (K + \lambda I)^{-1} \mathbf{y}^T$

$\mathbf{w} = \varphi(\mathcal{X}_{train})\boldsymbol{\alpha}$

$K(a, b) = K(b, a)^T$

We call $K = k(\mathcal{X}_{train}, \mathcal{X}_{train})$
the Gram matrix

Summary Kernel Ridge Regression

- Input: kernel function $k(\mathbf{x}, \mathbf{x}') = \varphi(\mathbf{x})^T \varphi(\mathbf{x}')$, regularization hyperparameter λ , training dataset $(X_{train}, \mathbf{y}_{train})$ and test datapoints X_{test}
- Training ²:

$$\boldsymbol{\alpha} = (k(X_{train}, X_{train}) + \lambda I)^{-1} \mathbf{y}_{train}^T$$

- Predicting:

$$\hat{\mathbf{y}}_{test} = \boldsymbol{\alpha}^T k(X_{train}, X_{test})$$

²Since we need X_{train} again for predictions, we cannot forget the data like for RR

Kernels as Similarity Measures

$$f^*(\mathbf{x}_{\text{new}}) = \sum_{i=1}^N \alpha_i k(\mathbf{x}_{\text{new}}, \mathbf{x}_i)$$

Kernel methods are *memory-based* methods:

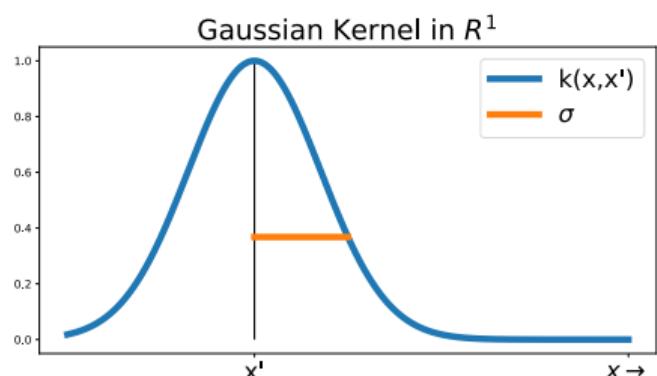
- store the entire training set
- define similarity of data points by kernel function

$$k(\cdot, \cdot) : \mathbb{R}^d \times \mathbb{R}^d \rightarrow \mathbb{R}$$

- new predictions require comparison with previously learned examples

Law of Universal Generalization

- Roger Shepard argued that the **Perceptual Similarity** of new data x decays exponentially with distance from prototype x' [Shepard, 1987]
 - motivation for using the Gaussian Kernel



Gaussian Kernel (RBF Kernel)

$$k(x', x) = e^{-\frac{(x' - x)^2}{2\sigma^2}}$$

Because $\exp(x) = \sum_{k=0}^{\infty} \frac{x^k}{k!}$, the basis function $\varphi(x)$ is infinite dimensional

Kernel Ridge Regression:

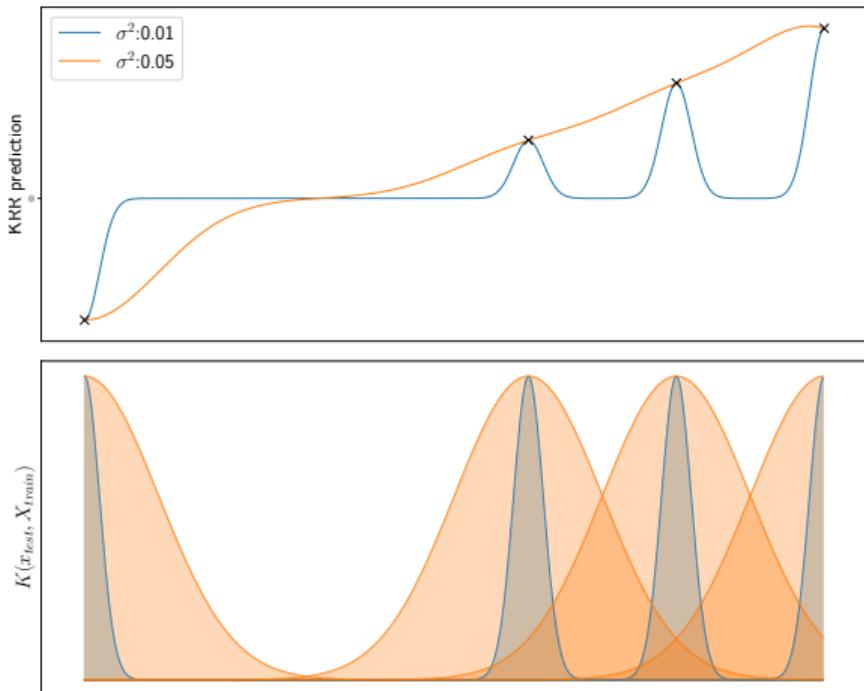
Example

KRR with Gaussian Kernels:

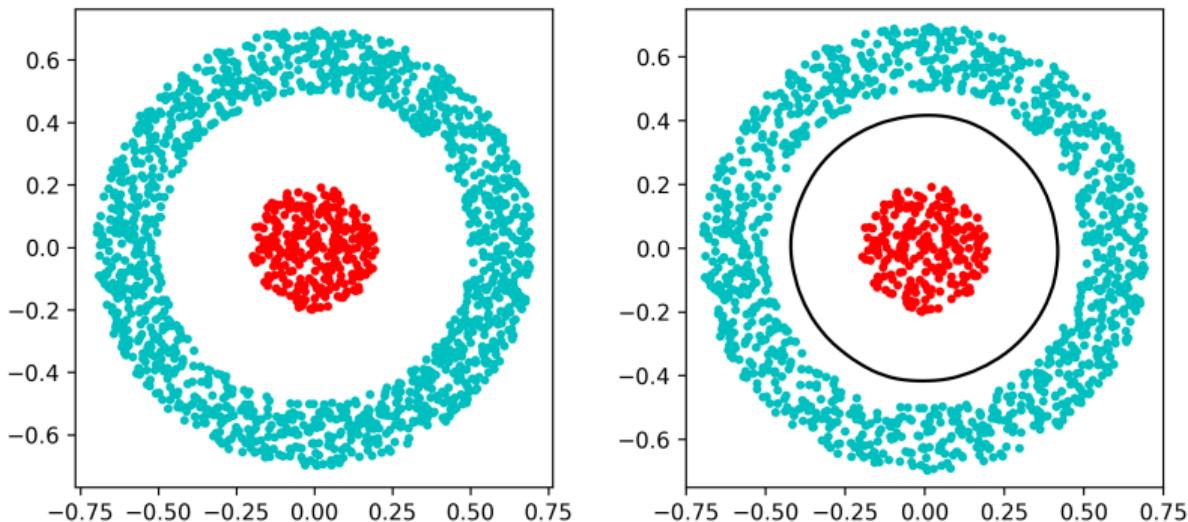
$$k(x, x') = \exp \left\{ -\frac{\|x - x'\|^2}{2\sigma^2} \right\}$$

Predictions:

$$y_{new} = y_{train}(K + \lambda I)^{-1} k(X_{train}, x_{new})$$



We can also kernelize LDA with e.g. a Gaussian Kernel ($\sigma = 0.1$)



→ We are able to do non-linear classification!

Kernel Methods - Pros and Cons

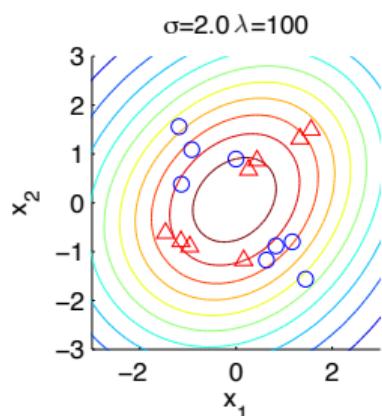
- + Powerful Modeling Tool
 - (non-linear problems become linear in kernel space)
 - + Omnipurpose Kernels
 - (Gaussian works well in many cases)
 - + Kernel methods can handle symbolic objects
 - + When you have less data points than your data has dimensions kernel methods can offer a dramatic speedup

 - Difficult to understand what's happening in kernel space
 - Model complexity increases with number of data points
 - If you have too much data, kernel methods can be slow
- + 强大的建模工具
(非线性问题在核空间中变为线性)
 - + 通用内核
(高斯在许多情况下效果很好)
 - + 内核方法可以处理符号对象
 - + 当您的数据点少于数据的维度时，内核方法可以提供显著的加速
 - 难以理解内核空间中发生了什么
 - 模型复杂性随着数据点数量的增加而增加
→ 如果数据过多，内核方法可能会很慢

Generalization and Model Selection

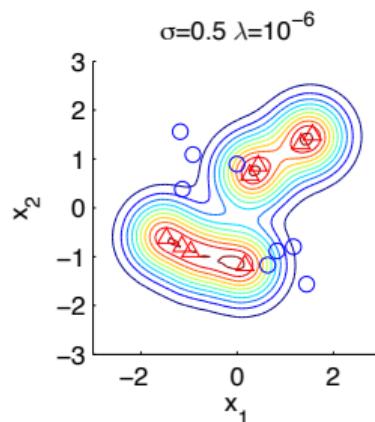
The best model is the model that *generalizes best*

Underfitting



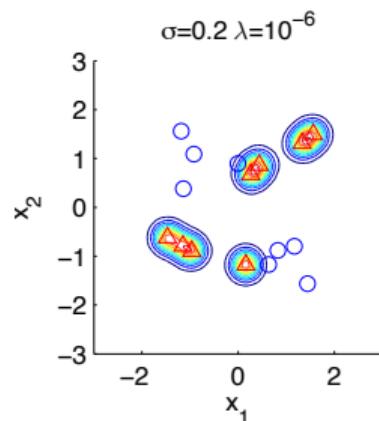
Model is too simple
→ Bad generalization

Better fit



Appropriate complexity
→ Good generalization

Overfitting



Model is too complex

Recap: Cross-Validation

Split data set in F different **training** and **test** data

$$\text{fold 1} [\underbrace{x_1, x_2, x_3, x_4}_{\mathcal{F}_1^{\text{train}}}, \underbrace{x_5, x_6}_{\mathcal{F}_1^{\text{test}}}]$$

$$\text{fold 2} [\underbrace{x_1, x_2}_{\mathcal{F}_1^{\text{test}}}, \underbrace{x_3, x_4, x_5, x_6}_{\mathcal{F}_1^{\text{train}}}]$$

fold 3 ...

For each fold:

Train your model on the training data

Test your model on the test data

How to Achieve Good Generalization?

When using powerful algorithms (MLPs, KRR, ...)
every data set can be modeled perfectly! (overfitting)

But we want to model new data well (generalization)

Cross-validation can be used for **either**:

Model selection

Optimize hyper-parameters of a model for generalization performance

Model evaluation

Test how good an algorithm with fixed parameters actually is

Cross-Validation: Model Selection or Model Evaluation

Model Evaluation

Report **mean evaluation score** – e.g.
accuracy – across folds

Model Selection

Take hyper-parameter with the highest
mean score across folds

Want to estimate the performance model which we optimize on unseen data:

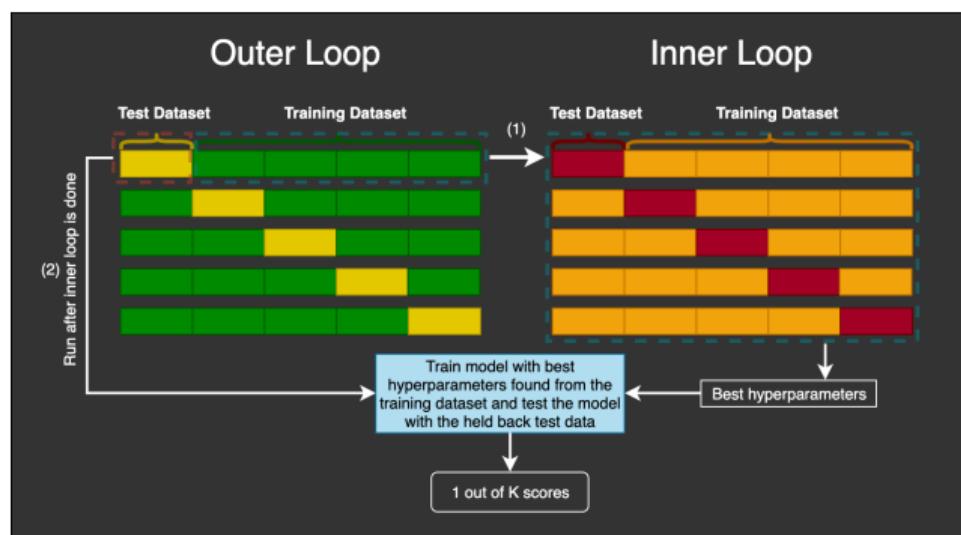
- If we did selection and evaluation on the same test fold:
 - We would be too optimistic
 - because we use same set for optimizing and evaluating ³
 - would be similar to reporting train error
- Solution: **Nested Cross-Validation**

³e.g. see this sklearn example

Nested Cross-Validation

It is simply
CV for model selection
nested inside
CV for model evaluation

- Two loops, two hold out folds
- Only gives you an estimator for performance
- **Does not replace parameter tuning**



Source: mlfromscratch.com/nested-cross-validation-python-code

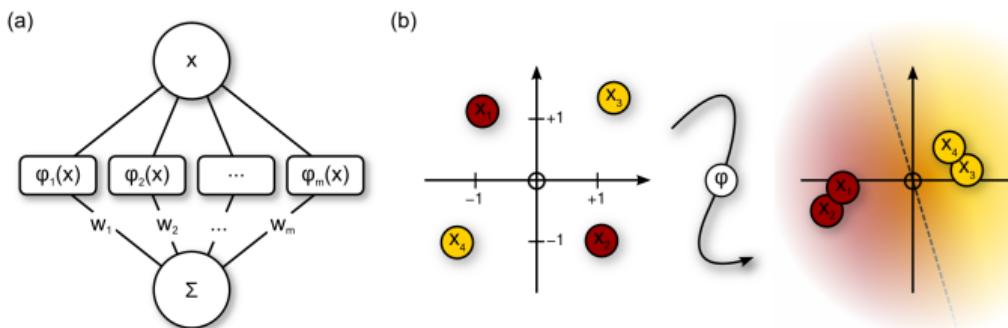
Nested Cross-Validation

Algorithm 1: Cross-Validation for Model Selection and Evaluation

Require: Data $(x_1, y_1), \dots, (x_N, y_N)$, parameters $\sigma_1, \dots, \sigma_S$, Number of CV folds F

```
1: Split data in  $F$  disjunct folds
2: for Outer folds  $f_{\text{outer}} = 1, \dots, F$  do
3:   Pick folds  $\{1, \dots, F\} \setminus f_{\text{outer}}$  for Model Selection
4:   Model Selection
5:   for Fold  $f_{\text{inner}} = 1, \dots, F - 1$  do
6:     for Parameter  $s = 1, \dots, S$  do
7:       Train model on folds  $\{1, \dots, F\} \setminus \{f_{\text{outer}}, f_{\text{inner}}\}$  with parameter  $\sigma_s$ 
8:       Compute prediction on fold  $f_{\text{inner}}$ 
9:     end for
10:   end for
11:   Pick best parameter  $\sigma_s$  for all  $f_{\text{inner}}$ 
12:   Model Evaluation
13:   Train model on folds  $\{1, \dots, F\} \setminus f_{\text{outer}}$  with parameter  $\sigma_s$ 
14:    $\text{Performance}_{\text{outer}} \leftarrow \text{Test model on fold } f_{\text{outer}}$ 
15: end for
16: return Average of  $\text{Performance}_{\text{outer}}$ 
```

Kernelizing linear methods



[Jäkel et al., 2009]

1. Map the data into a (high dimensional) feature space, $\mathbf{x} \mapsto \varphi(\mathbf{x})$
2. Look for linear relations in the feature space
 - Work in that space by considering scalar product of data points,
 $k(\mathbf{x}_i, \mathbf{x}_j) = \langle \varphi(\mathbf{x}_i), \varphi(\mathbf{x}_j) \rangle$
 - Many linear models have a *dual representation* that only uses scalar products between the data points
 - k is called the *kernel function*

Summary

Kernel Ridge Regression

Non-linear regression

Predictions involve comparison of new and old data

Predictions based on linear combination of (non-linear) similarity measures

Optimization requires inversion of kernel matrix

$(N \times N, \rightarrow \mathcal{O}(N^3))$

→ Difficult for very large data sets

Generalization and Model Selection

Good prediction on new data is called generalization

Cross-Validation is a simple and powerful framework for model selection

Nested Cross-Validation gives you a valid estimate for generalization of your model class (*without giving you parameters or hyperparameters*)

References

- C. M. Bishop. *Pattern Recognition and Machine Learning (Information Science and Statistics)*. 2007.
- C. F. Gauß. *Theoria motus corporum coelestium in sectionibus conicis solem ambientium*. Göttingen, 1809.
- A. E. Hoerl and R. W. Kennar. Ridge regression: Applications to nonorthogonal problems. *Technometrics*, 12(1):69–82, 1970.
- F. Jäkel, B. Schölkopf, and F. A. Wichmann. Does cognitive science need kernels? *Trends Cogn Sci*, 13(9):381–388, 2009. doi: 016/j.tics.2009.06.002.
- G. Kimeldorf and G. Wahba. Some results on tchebycheffian spline functions. *Journal of mathematical analysis and applications*, 33(1):82–95, 1971.
- A.-M. Legendre. *Nouvelles méthodes pour la détermination des orbites des comètes*, chapter Sur la méthode des moindres quarrés. Firmin Didot, <http://imgbase-scd-ulb.u-strasbg.fr/displayimage.php?pos=-141297>, 1805.
- J. Mercer. Xvi. functions of positive and negative type, and their connection the theory of integral equations. *Philosophical transactions of the royal society of London. Series A, containing papers of a mathematical or physical character*, 209(441-458):415–446, 1909.
- R. N. Shepard. Toward a universal law of generalization for psychological science. *Science*, 237(4820):1317–23, 1987.
- A. N. Tychonoff. On the stability of inverse problems. *Doklady Akademii Nauk SSSR*, 39(5):195–198, 1943.

References

- C. M. Bishop. *Pattern Recognition and Machine Learning (Information Science and Statistics)*. 2007.
- C. F. Gauß. *Theoria motus corporum coelestium in sectionibus conicis solem ambientium*. Göttingen, 1809.
- A. E. Hoerl and R. W. Kennar. Ridge regression: Applications to nonorthogonal problems. *Technometrics*, 12(1):69–82, 1970.
- F. Jäkel, B. Schölkopf, and F. A. Wichmann. Does cognitive science need kernels? *Trends Cogn Sci*, 13(9):381–388, 2009. doi: 016/j.tics.2009.06.002.
- G. Kimeldorf and G. Wahba. Some results on tchebycheffian spline functions. *Journal of mathematical analysis and applications*, 33(1):82–95, 1971.
- A.-M. Legendre. *Nouvelles méthodes pour la détermination des orbites des comètes*, chapter Sur la méthode des moindres quarrés. Firmin Didot, <http://imgbase-scd-ulb.u-strasbg.fr/displayimage.php?pos=-141297>, 1805.
- J. Mercer. Xvi. functions of positive and negative type, and their connection the theory of integral equations. *Philosophical transactions of the royal society of London. Series A, containing papers of a mathematical or physical character*, 209(441-458):415–446, 1909.
- R. N. Shepard. Toward a universal law of generalization for psychological science. *Science*, 237(4820):1317–23, 1987.
- A. N. Tychonoff. On the stability of inverse problems. *Doklady Akademii Nauk SSSR*, 39(5):195–198, 1943.