

Kernel-Based Learning & Multivariate Modeling

MIRI Master

Lluís A. Belanche

`belanche@cs.upc.edu`

Soft Computing Research Group

Universitat Politècnica de Catalunya

2021-2022

Kernel-Based Learning & Multivariate Modeling

Syllabus

Part 1

Sep 15 Introduction to kernel-based learning

Sep 22 The SVM for classification, regression & novelty detection (I)

Sep 29 The SVM for classification, regression & novelty detection (II)

Oct 06 Kernel design (I): theoretical issues

Oct 13 Kernel design (II): practical issues

Oct 20 Kernelizing ML & statistical algorithms

Oct 27 Advanced topics

Kernel-Based Learning

Introduction

Desiderata for satisfactory learning methods (my particular view):

Robustness to outliers, errors and/or wrong model assumptions
tolerance

Efficiency in the computational sense (necessary to handle large datasets)
complexity, fast languages, use of supercomputers etc

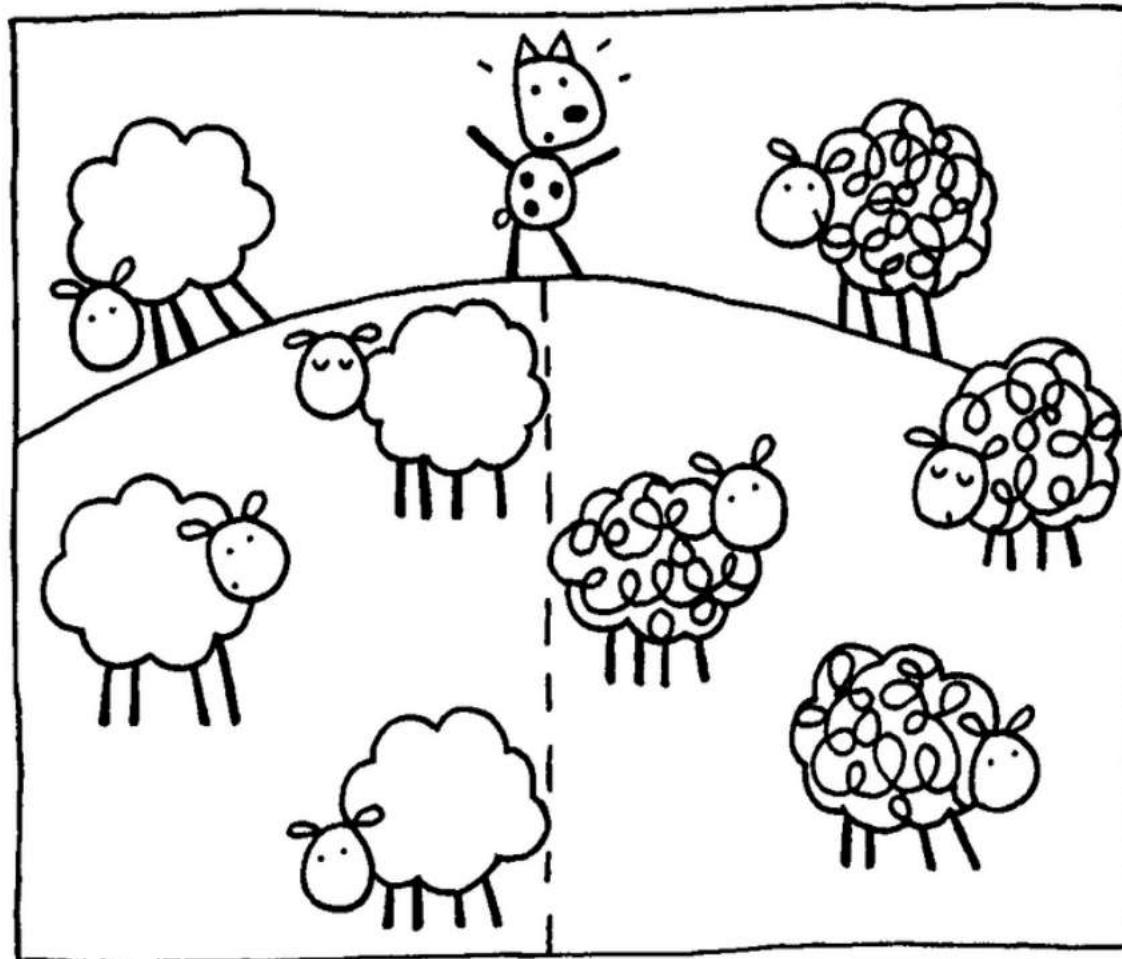
Flexibility to perform different tasks

Controlable non-linearity to deliver complexity surplus and accept explicit complexity control

Versatility to accept different data types and incorporate prior knowledge

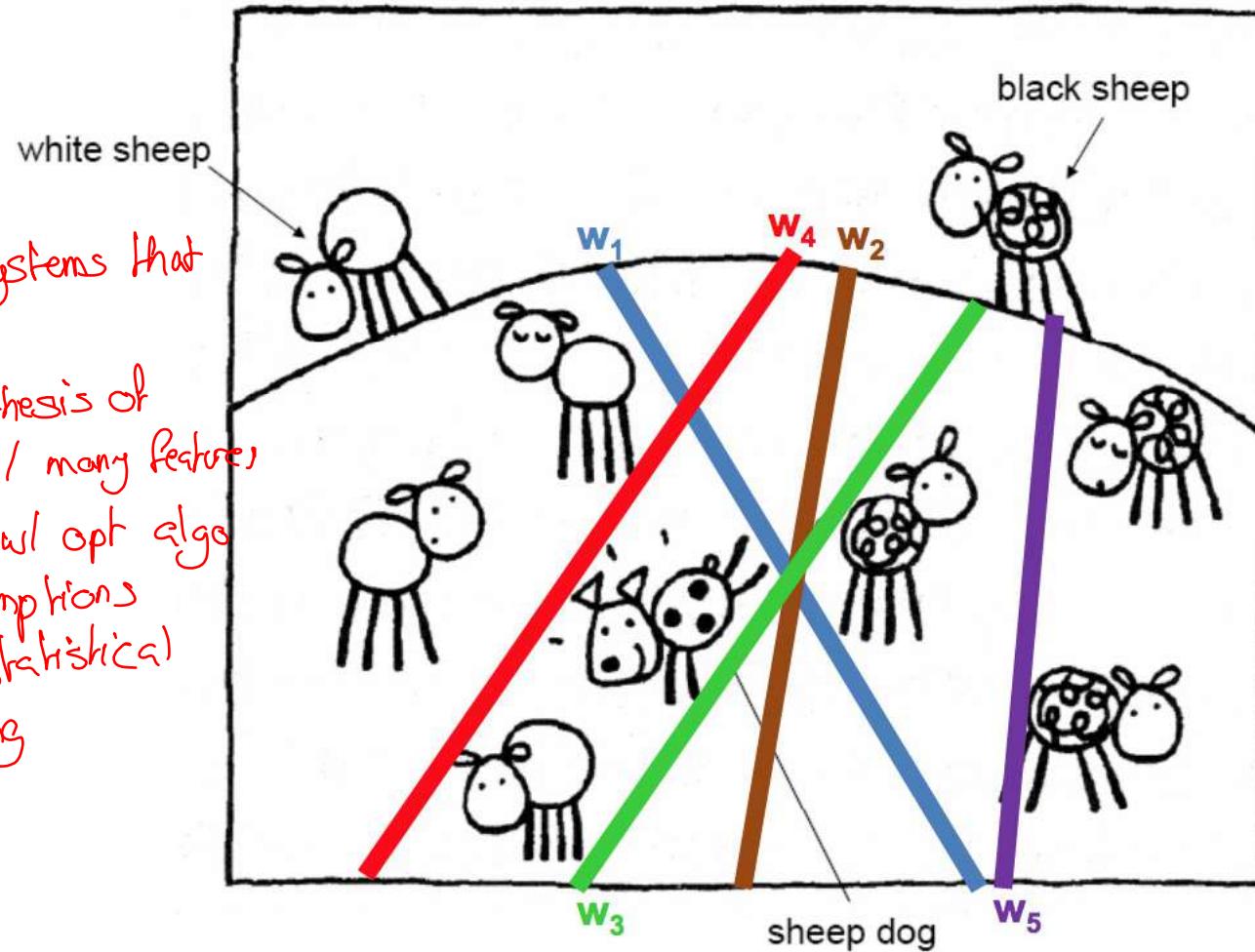
⇒ **Generalize** well to unseen data (as well as possible)

Kernel-Based Learning



Kernel-Based Learning

calculation is only estimation and classifier is random because it depends on sample



SVMs are systems that add info

↳ use hypothesis of function w/ many features

↳ trained w/ opt algo w/ assumptions from statistical learning

odd info to solve

"SVMs are regularized learning systems that leverage a hypothesis space of linear functions in a high dimensional feature space, and trained with a learning algorithm from optimization theory that implements a learning bias derived from statistical learning theory."

n-dim space w/ all aff of data

set of all hypo

assumptions

Kernel-Based Learning

Introduction

Kernel methods constitute a very powerful paradigm:

- provide a principled way to perform non linear, nonparametric learning;
- rely on solid functional analytic foundations and enjoy optimal statistical properties;
- they have been successfully applied to a wide range of data analysis problems (e.g. in computational biology);
- are specially well suited for high-dimensional data, noisy environments or heterogenous information (unlike other approaches like neural networks) but small data sets

Kernel-Based Learning

Introduction

- the number of hyper-parameters is usually very small compared to other approaches (again unlike neural networks);
- their main limitations are:
 - their difficult applicability in large data scenarios because of stringent computational requirements in terms of time and (especially) memory;
 - the inability to extract increasingly abstract features from the data (“single hidden layer”, unlike, yes, neural networks!)

→ *Kernel methods are a central tool in machine learning*

Kernel-Based Learning

Key aspects of kernel methods

1. Input data is embedded (mapped) into a **vector space**
2. **Linear relations** are sought among the elements of this vector space
3. The coordinates of the images (mapped data) are *not* needed: only their **pairwise inner products**
dot product: $a_1 \cdot b_1 + \dots + a_n \cdot b_n$
4. Often these inner products can be computed (efficiently and implicitly) in the input space (via a **kernel function**) → the PROMISE

Kernel-Based Learning

Introduction

Kernel-based methods consist of two ingredients:

1. The (right) **kernel function** to convert input data into a **similarity matrix**

2. The (chosen) **learning algorithm** that uses the kernel matrix and the data to produce a model that performs a prescribed task

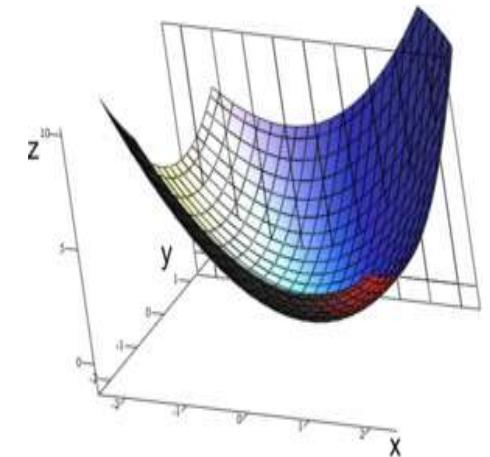
$$K = \begin{bmatrix} k(x_1, x_1) & k(x_1, x_2) & \cdots & k(x_1, x_n) \\ k(x_2, x_1) & k(x_2, x_2) & \cdots & k(x_2, x_n) \\ \vdots & \vdots & \ddots & \vdots \\ k(x_n, x_1) & k(x_n, x_2) & \cdots & k(x_n, x_n) \end{bmatrix}$$

Kernel-Based Learning

Introduction

A kernel function is:

1. a special case of similarity measure;
keep in mind: not all similarity measures are kernel functions
2. a positive semi-definite function; $\text{all eigenvalues} \geq 0$
 $K(x,x) \geq 0$
3. a general inner product in a (special) Hilbert space



It must be psd so that optimization problems can be convex and, therefore, uniquely solvable by standard algorithms (moreover, there is a global optimum)

Kernel-Based Learning

Introduction

- kernel functions can be used to analyze virtually any kind of data (and to express many forms of previous knowledge) without preprocessing the data
- kernel functions allow the obtention of more general (non-linear) versions of learning algorithms as long as they are base on inner products or Euclidean distances

Kernel-Based Learning

Introduction

Examples of **kernel functions**:

- RBF kernels
- Polynomial kernels
- String kernels
- Anova kernels
- Fisher kernels

Examples of (kernel) **learning algorithms**:

- Support Vector Machine (SVM) and Relevance Vector Machine (RVM)
- kernel Fisher Discriminant Analysis (FDA)
- kernel Principal Components analysis (PCA)
- kernel Canonical Correspondence analysis (CCA)
- kernel (regularized) regression (logistic, linear)
- kernel k -means

Kernel-Based Learning

Introduction

Designing/choosing an appropriate kernel function is not easy in general:

- Consider the **RBF kernel**:
$$k(x, x') = \exp\left(-\frac{1}{2} \frac{\|x - x'\|^2}{\sigma^2}\right), \quad \sigma^2 \in \mathbb{R}$$

Radial Basis Func

eucl. norm

$k \in (0, 1]$
 $k(x, x') = 1 \Leftrightarrow x = x'$

$\exists H, \exists \phi: \mathcal{X} \rightarrow H$ s.t. $\langle \phi(x), \phi(x') \rangle_{\text{inner prod}} = k(x, x')$
- All information of a problem (besides the target values) is tunneled through the kernel matrix $\mathbf{K} = (k_{ij})$, where $k_{ij} := k(x_i, x_j)$
- if σ^2 is very small, then $\mathbf{K} \approx \mathbf{I}$ (all data are dissimilar): over-fitting
too detailed
- if σ^2 is very large, then $\mathbf{K} \approx \mathbf{J}$ (all data are similar): under-fitting

Kernel-Based Learning

Notation

We have a data sample $S = \{(\bar{x}_1, t_1), \dots, (\bar{x}_n, t_n)\}$

- (\bar{x}_i, t_i) are drawn i.i.d. from some unknown (joint) prob. distribution
independently & identically distributed
- $\bar{x}_i \in \mathcal{X}, t_i \in \mathcal{T}$ (input space, output space)

Kernel-Based Learning

Linear regression

Problem: We wish to find a function $f(\bar{x}) = \bar{w}^\top \bar{x}$ which best models a data set S for the case $\mathcal{X} = \mathbb{R}^d$ and $\mathcal{T} = \mathbb{R}$

d = dimension

↖ transposed

- If $f(\bar{x}) = w^\top \bar{x} + w_0$, we rewrite to add one dimension as $x := (1, \bar{x}^\top)^\top$;
 $w := (w_0, \bar{w}^\top)^\top$ $w_0 \cdot 1 + \bar{w}^\top \bar{x}$
- Call $\mathbf{X}_{n \times (d+1)}$ the matrix of the x_i^\top and $t = (t_1, \dots, t_n)^\top$
- If S is generated as $(x, f(x))$ for some f , the x_i vectors are linearly independent and $n = d+1$, then there is a unique solution for $\mathbf{X}w = t$ given by $\hat{w} = \mathbf{X}^{-1}t$; in any other case, the problem is ill-posed

data sample $S = \{(x, f(x))\}$ w. x_i lin. independent $\& n = d+1$
 $\Rightarrow \exists$ unique sol $\mathbf{X}w = t$ by $\hat{w} = \mathbf{X}^{-1} \cdot t$

Kernel-Based Learning

III-posed & well-posedness

A problem is **ill-posed** if the solution may not always exist, is not uniquely determined or is unstable (small variations in the initial conditions of the problem cause the solution to change quite a lot).



Jacques
Hadamard
(1865-1963)

- A basic example would be the solution of linear systems of equations
- Learning problems are in general ill-posed
- The solution is to use an **inductive principle**; one of the most popular is the **regularization** principle

$$t_i = f(x_i) + \varepsilon_i$$

$$t_i = \omega^T x + x_0 + \varepsilon_i$$

$$\Leftrightarrow \varepsilon_i = t_i - \omega^T x - x_0$$

not σ^2 from kernel!

This is only
for regression

Kernel-Based Learning

Ridge regression

Under the standard assumptions $t_i = f(x_i) + \varepsilon_i$ and $\varepsilon_i \sim N(0, \sigma^2)$, a maximum likelihood argument leads to the minimization of the regularized empirical error (a.k.a penalized likelihood):

$$E_\lambda(\omega) := \sum_{i=1}^n (t_i - \omega^T x_i)^2 + \lambda \sum_{j=1}^d w_j^2$$

bc Gaussian *weil $E(\cdot)^2 = (\cdot)^T \cdot (\cdot)$*

$$= \underbrace{\|t - X\omega\|^2}_{\|t - X\omega\|^2} + \underbrace{\lambda \|\bar{w}\|^2}_{\lambda \|\bar{w}\|^2}$$

1. Note $E_\lambda(\omega) = \|t - X\omega\|^2 + \lambda \|\bar{w}\|^2$

2. The parameter $\lambda > 0$ defines a trade-off between the fit to the data and the complexity of the model (2-norm of \bar{w} in this case)

Kernel-Based Learning

Ridge regression: primal representation

Grows w/ param

Setting $\nabla_{\mathbf{w}} E_{\lambda} = 0$, we obtain the (regularized) normal equations

$$-2\mathbf{X}^T(\mathbf{t} - \mathbf{X}\mathbf{w}) + 2\lambda\mathbf{w} = 0$$

with solution $\hat{\mathbf{w}} = (\mathbf{X}^T\mathbf{X} + \lambda\mathbf{I}_{d+1})^{-1}\mathbf{X}^T\mathbf{t}$

and therefore $f(x) = \hat{\mathbf{w}}^T x = \mathbf{t}^T \mathbf{X} (\mathbf{X}^T\mathbf{X} + \lambda\mathbf{I}_{d+1})^{-1} x.$

1. Since \mathbf{X} is $n \times (d + 1)$, the matrix $\mathbf{X}^T\mathbf{X}$ is $(d + 1) \times (d + 1)$

2. $\mathbf{X}^T\mathbf{X} + \lambda\mathbf{I}_{d+1}$ always has an inverse, for all $\lambda > 0$

3. The “model size” does not grow with data size (a **parametric model**)

Kernel-Based Learning

Dual representation grows w/ data

It turns out that the regularized solution can also be written as:

$$\hat{w} = \sum_{i=1}^n \hat{\alpha}_i x_i$$

→ use this to go
from primal to dual

In consequence,

$$f(x) = \hat{w}^T x = \sum_{i=1}^n \hat{\alpha}_i x_i^T x$$

very diff
compared to $x^T x$

1. The vector of parameters $\hat{\alpha} = (\hat{\alpha}_1, \dots, \hat{\alpha}_n)^T$ is $\hat{\alpha} = (\mathbf{X}\mathbf{X}^T + \lambda \mathbf{I}_n)^{-1} \mathbf{t}$
2. The **Gram matrix** $\mathbf{X}\mathbf{X}^T$ ("matrix of inner products") is $n \times n$

how grows
w/ data

Kernel-Based Learning

Primal and dual

So we have the **primal** and the **dual** forms for $f(x)$:

$$\text{primal} \quad f(x) = \hat{w}^\top x = \sum_{j=0}^d \hat{w}_j x_j \quad \text{and} \quad \text{dual} \quad f(x) = \sum_{i=1}^n \hat{\alpha}_i x_i^\top x$$

The **dual** form is more convenient when $d \gg n$:

param data size

- The primal requires the computation and inversion of $\mathbf{X}^\top \mathbf{X} + \lambda \mathbf{I}_{(d+1)}$, requiring $O(nd^2 + d^3)$ operations param model
- The dual requires the computation and inversion of $\mathbf{X} \mathbf{X}^\top + \lambda \mathbf{I}_n$, requiring $O(dn^2 + n^3)$ operations

Kernel-Based Learning

How can we perform non-linear regression?

First create a *feature map*, a function $\phi : \mathbb{R}^d \rightarrow \mathbb{R}^D$, with $D \in \mathbb{N} \cup \{\infty\}$

$$x \mapsto \phi(x) = (\phi_1(x), \phi_2(x), \dots, \phi_D(x))^\top$$

- $\phi(x)$ is called a *feature vector*
- $\{\phi(x) : x \in \mathbb{R}^d\}$ is the *feature space* (possibly part of a larger vector space)
- as a technicality, we could easily add $\phi_0(x) = 1$ as before, if desired

Kernel-Based Learning

Hilbert spaces

The right structure for this vector space \mathcal{H} is that of a **Hilbert space**:

- A vector space endowed with an **inner product** $\langle \cdot, \cdot \rangle$ whose associated norm defines a complete metric

1. define inner prod $\langle \cdot, \cdot \rangle : \mathcal{H} \times \mathcal{H} \rightarrow \mathbb{R}$
2. define norm out of inner prod $\|f\| := \sqrt{\langle f, f \rangle}$
3. check completeness

- Completeness means that all **Cauchy sequences** defined in \mathcal{H} converge to an element of \mathcal{H} $\forall \epsilon > 0, \exists N : \|x_n - x_m\| < \epsilon, \forall n, m > N$

- Example: the l_2 space of square-summable sequences
$$l_2 := \left\{ \text{sequences } \{x_n\}_{n \in \mathbb{N}} \mid \sum_{n=0}^{\infty} (x_n)^2 < \infty \right\}$$
 Complete bc \mathbb{R} is complete

(informally) Generalizes the notion of **Euclidean space**: an **infinite-dimensional space** with the **structure of \mathbb{R}^d** (distances, lengths and angles are well-defined)



David
Hilbert
(1862-1943)

Kernel-Based Learning

(non-linear)

The regression function has now the primal representation:

$$f(x) = \langle w, \phi(x) \rangle = \sum_{j=1}^D w_j \phi_j(x)$$

compute dual rep as
Exercise

- This feature space has the structure of \mathbb{R}^D (a vector space)
- In consequence, there is also the dual representation:

$$f(x) = \sum_{i=1}^n \hat{\alpha}_i \langle \phi(x_i), \phi(x) \rangle$$

$$w = \sum \hat{\alpha}_i x_i$$

(BTW, how general are all these results?)

Kernel-Based Learning

Feature maps and kernels (1)

Given a feature map $\phi : \mathbb{R}^d \rightarrow \mathcal{H}$, being \mathcal{H} a Hilbert space, we define its associated **kernel function** $k : \mathbb{R}^d \times \mathbb{R}^d \rightarrow \mathbb{R}$ as:

$$k(x, x') = \langle \phi(x), \phi(x') \rangle, \quad x, x' \in \mathbb{R}^d$$

One key point is that, for some feature maps, computing $k(x, x')$ is independent of D (the dimension of \mathcal{H})

if $k : X \times X \rightarrow \mathbb{R}$ is symm PSD Func,

→ the PROMISE!!! then $\exists \phi : X \rightarrow \mathcal{H}$ s.t. $\langle \phi(x), \phi(x') \rangle_{\mathcal{H}} = k(x, x') \quad \forall x, x' \in X$
⇒ reverse also applies

Kernel-Based Learning

Feature maps and kernels (2)

Our regression function has now the dual representation:

$$f(x) = \sum_{i=1}^n \hat{\alpha}_i \langle \phi(x_i), \phi(x) \rangle = \sum_{i=1}^n \hat{\alpha}_i k(x_i, x)$$

often $O(sparse)$

This dual representation:

- ... is a **non-linear model** (in the **input space**)
- ... is a **linear model** (in the **feature space**)

We then have a **non-parametric model** (complexity grows with data size)

Kernel-Based Learning

Back to ridge regression ...

The new vector of parameters $\hat{\alpha} = (\hat{\alpha}_1, \dots, \hat{\alpha}_n)^\top$ is now given by

$$\hat{\alpha} = (\mathbf{K} + \lambda \mathbf{I}_n)^{-1} \mathbf{t},$$

where (as introduced earlier) $\mathbf{K} = (k_{ij})$, with $k_{ij} = k(x_i, x_j)$.

So we can do ridge regression based only on \mathbf{K} (and throw away \mathbf{X})

→ this is called **Kernel ridge regression (KRR)**

Kernel-Based Learning

Kernel ridge regression

What if we take the (simplest) choice $\phi(x) = x$? In this case $d = D$ and $k(x, x') = \langle x, x' \rangle = x^\top x'$. The regularized solution reads:

$$f(x) = \sum_{i=1}^n \hat{\alpha}_i x_i^\top x$$
$$= (X X^\top + \lambda I_n)^{-1} t$$

$X X^\top$
 $n \times n$

where

$$\hat{\alpha} = (X X^\top + \lambda I_n)^{-1} t,$$

$\hat{\alpha} \in \mathbb{R}^n$

(so $K = X X^\top$ in this particularly simple case)

This means we have **generalized** (the **dual** of) standard ridge regression via a kernel function (we have **kernelized** it!)
use smaller input

Kernel-Based Learning

Kernelizing ...

Many (classical and new) learning algorithms can be kernelized:

1. They require solving a problem where the **data appear** in the form of **pairwise inner products** (or **pairwise Euclidean distances**)
2. The solution is expressed as a **linear combination** of the kernel function centered at the data (ideally we only use *some* of the data: **sparsity**)
3. Examples include SVMs, ridge regression, perceptrons, FDA, PLS [supervised], as well as PCA, k-means, Parzen Windows [unsupervised]

Kernel-Based Learning

General feature maps

A feature map is of the general form $\phi : \mathcal{X} \rightarrow \mathcal{H}$

The associated kernel function $k : \mathcal{X} \times \mathcal{X} \rightarrow \mathbb{R}$ is:

$$k(x, x') = \langle \phi(x), \phi(x') \rangle, \quad x, x' \in \mathcal{X}$$

\mathcal{X} can be any space, \mathcal{H} is always a (special case of) Hilbert space

In our starting development, $\mathcal{X} = \mathbb{R}^d$

Kernel-Based Learning

Regularization-based learning algorithms

The key for this is the **regularization framework**; consider again the regularized empirical error for mapped data:

$$E_\lambda(\mathbf{w}) = \sum_{i=1}^n (t_i - \langle \mathbf{w}, \phi(x_i) \rangle)^2 + \lambda \|\mathbf{w}\|^2, \quad \lambda > 0$$

which we now generalize to arbitrary **loss** functions $L : \mathbb{R} \times \mathbb{R} \rightarrow \mathbb{R}$:
Very general function

$$E_\lambda(\mathbf{w}) = \sum_{i=1}^n \underbrace{L(t_i, \langle \mathbf{w}, \phi(x_i) \rangle)}_{(t_i - \langle \mathbf{w}, \phi(x_i) \rangle)^2} + \underbrace{\lambda \|\mathbf{w}\|^2}_{\text{regularizer}}, \quad \lambda > 0$$

Kernel-Based Learning

Regularization-based learning algorithms

Theorem (informal). If L is differentiable w.r.t. its second argument and \hat{w} is a minimizer of $E_\lambda(w)$, then we can represent:

$$\hat{w} = \sum_{i=1}^n \hat{\alpha}_i \phi(x_i)$$

and therefore

$$f(x) = \langle \hat{w}, \phi(x) \rangle = \sum_{i=1}^n \hat{\alpha}_i k(x_i, x)$$

This result is usually called the **Representer Theorem**. In the case of KRR, $\hat{\alpha} = (K + \lambda I_n)^{-1} t$.

Organization of the course

Theory

Practice Lab

In-class exercises

Practical work

Partial exam: 08/11/21 10:00-11:30h

Kernel-Based Learning

Suggested theoretical exercise

1. Derive the (primal) solution of ridge regression by minimization of the regularized empirical error, and derive a formula for the model
2. Prove that the obtained (regularized) solution can be also written as:

$$\mathbf{w}^* = \sum_{n=1}^N \alpha_n \mathbf{x}_n$$

In consequence,

$$f(\mathbf{x}) = \sum_{n=1}^N \alpha_n \langle \mathbf{x}_n, \mathbf{x} \rangle, \text{ where } \boldsymbol{\alpha} = (X X^T + \lambda I_N)^{-1} \mathbf{t}$$

Kernel-Based Learning

Suggested practical exercise

Consider the function

$$f(x) = 0.5 \frac{\sin(x - a)}{x - a} + 0.8 \frac{\sin(x - b)}{x - b} + 0.3 \frac{\sin(x - c)}{x - c}, \quad x \in \mathbb{R}$$

where $a = 10, b = 50, c = 80$.

1. Generate a dataset of $N = 1052$ examples where the x are equally-spaced in $[0.1, 100]$ and the targets for regression are obtained as $t = f(x) + N(0, 0.05^2)$
2. Fit standard polynomial regression with some degrees of your choice
3. Fit kernel ridge regression with the RBF kernel and some σ and λ of your (careful) choice, until you are satisfied with the fit.

Kernel-Based Learning & Multivariate Modeling

MIRI Master

Lluís A. Belanche

`belanche@cs.upc.edu`

Soft Computing Research Group

Universitat Politècnica de Catalunya

2021-2022

Kernel-Based Learning & Multivariate Modeling

Syllabus

Part 1

Sep 15 Introduction to kernel-based learning

Sep 22 The SVM for classification, regression & novelty detection (I)

Sep 29 The SVM for classification, regression & novelty detection (II)

Oct 06 Kernel design (I): theoretical issues

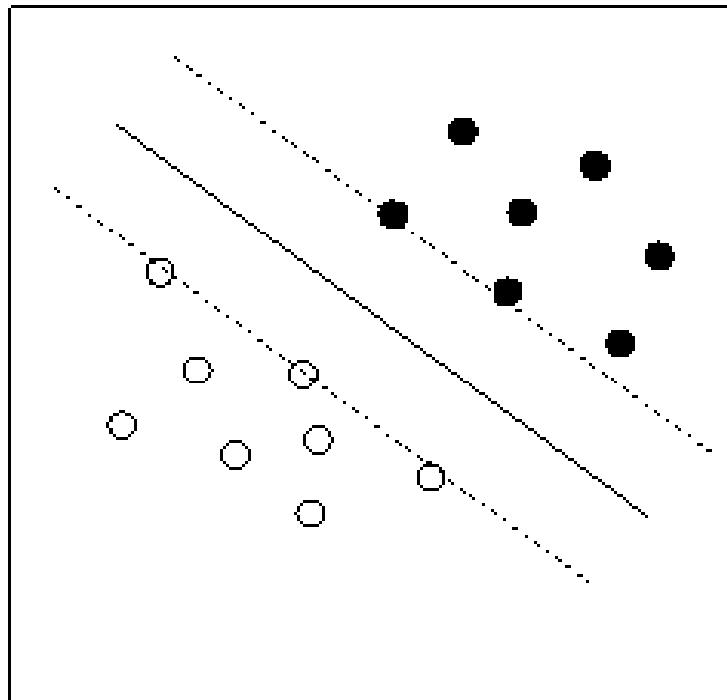
Oct 13 Kernel design (II): practical issues

Oct 20 Kernelizing ML & statistical algorithms

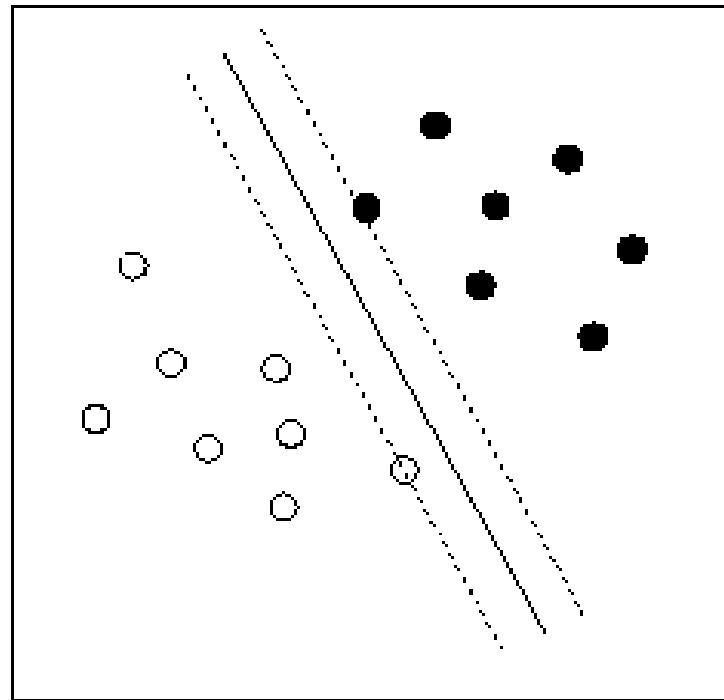
Oct 27 Advanced topics

Support Vector Machines

Preliminaries



(a) Larger margin



(b) Smaller margin

Which solution is more likely to lead to better **generalization**?

works for all models & data sets

Support Vector Machines

Preliminaries

- Criterion for building a two-class classifier:

Maximize the **margin = width of the separation** between the classes, defined by the **distance to the nearest training examples**

- Working Hypotheses:**

1. The data are **linearly separable** ("linsep") –very unlikely, but see later
2. The **larger the margin**, the better the **generalization** (a first intuition)

Goal: find the separating hyperplane with the **largest margin**

Support Vector Machines

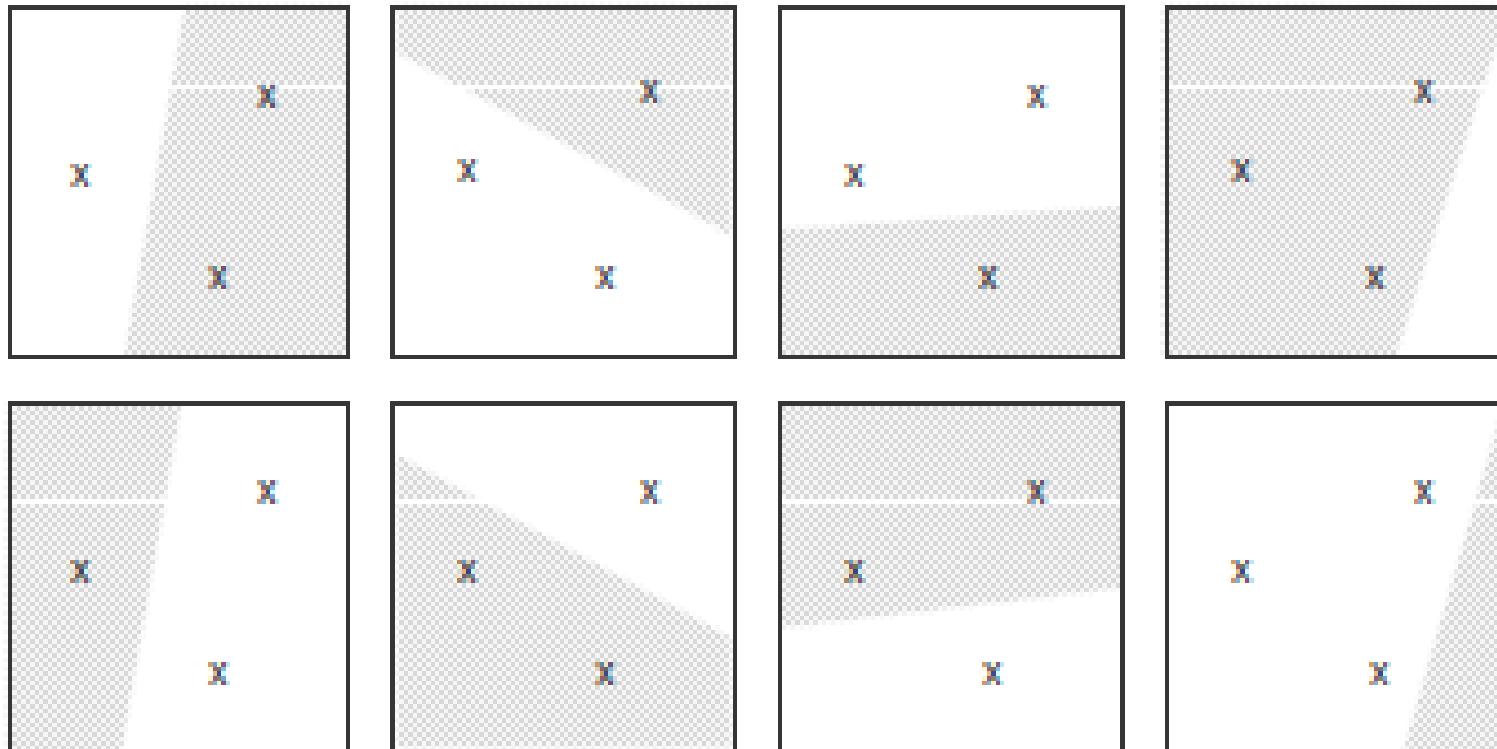
Preliminaries

For a two-class classifier, the **VC dimension** ϑ is the maximum number of points that can be separated in all possible 2^ϑ ways (**shattered**) by using functions representable by the classifier.

- Note it is *sufficient* that one set of ϑ points exists that can be shattered for the VC dimension to be at least ϑ
- If the VC dimension of a class is ϑ , this means there is at least one set of ϑ points that can be shattered by members of the class. It does not mean that every set of ϑ points can be shattered
- If no set of $\vartheta + 1$ points can be shattered by members of the class, then the VC dimension of the class is at most ϑ

Support Vector Machines

An example



$$v(\mathbb{R}^2) = 3$$

- In \mathbb{R}^2 we can shatter these three points (VC dim is ≥ 3)
- No set of four or more points can be shattered (VC dim is < 4)

Support Vector Machines

Why is the VC dimension relevant?

Theorem (Vapnik and Chervonenkis, 1974). Let D be an i.i.d data sample of size n and $\mathcal{Y}^{\{0,1\}}$ a class of parametric binary classifiers. Let ϑ denote the VC dimension of \mathcal{Y} . Take $y \in \mathcal{Y}$ with empirical error $R_n(y)$ on D . For all $\eta > 0$ it holds true that, with probability at least $1 - \eta$, the true error of y is bounded by:

$$R(y) \leq R_n(y) + \underbrace{H(n, \vartheta, \eta)}_{\text{"error term"}}$$

where

$$H(n, \vartheta, \eta) = \sqrt{\frac{\vartheta(\ln(2n/\vartheta) + 1) - \ln(\eta/4)}{n}}$$

Support Vector Machines

Formalisation

We have a data set $D = \{(x_1, t_1), \dots, (x_n, t_n)\}$, with $x_i \in \mathbb{R}^d$ and $t_i \in \{-1, +1\}$, describing a two-class problem.

We wish to find a linear function f which best models D :

1. • Set up an **affine function** $g(x) = \underbrace{\langle w, x \rangle + b}_{\text{sign(Vorzeichen} = 1, 0, -1\})}$
2. • Obtain a **linear discriminant** as $f(x) = \underbrace{\text{sgn}(g(x))}_{\text{sign(Vorzeichen} = 1, 0, -1\})}$
3. • We would like to find w, b such that:
 $\underbrace{\langle w, x_i \rangle + b}_{\text{when } t_i = +1} > 0$, when $t_i = +1$ $\underbrace{\langle w, x_i \rangle + b}_{\text{when } t_i = -1} < 0$, when $t_i = -1$ ← error
that is $t_i(\langle w, x_i \rangle + b) > 0$ or simply $t_i(g(x_i)) > 0$ $(1 \leq i \leq n)$

Support Vector Machines

Formalisation

- The quantity $t_i g(x_i)$ is called the **functional margin** of x_i (there will be an “error” whenever $t_i g(x_i) < 0$)
- Define the **loss** $L(t_i, \langle \mathbf{w}, \mathbf{x}_i \rangle) := \underbrace{\max(1 - t_i g(x_i), 1)}_{\substack{\text{has error} \\ t_i > 0}}$
- Given the plane $\pi : g(\mathbf{x}) = \langle \mathbf{w}, \mathbf{x} \rangle + b = 0$, the distance
$$d(\mathbf{x}, \pi) = \frac{|g(\mathbf{x})|}{\|\mathbf{w}\|}$$
 is called the **geometrical margin** of \mathbf{x} .
$$= \frac{|\langle \mathbf{w}, \mathbf{x} \rangle + b|}{\|\mathbf{w}\|}$$
- The **optimal separating hyperplane** (OSH) for linsep data is the one that maximizes the geometrical margin:

$$\max_{\mathbf{w}, b} \left\{ \min_{1 \leq i \leq n} d(\mathbf{x}_i, \pi) \right\} \quad \text{subject to } t_i (\langle \mathbf{w}, \mathbf{x}_i \rangle + b) > 0 \quad (1 \leq i \leq n)$$

max dist to nearest x

Support Vector Machines

Formalisation

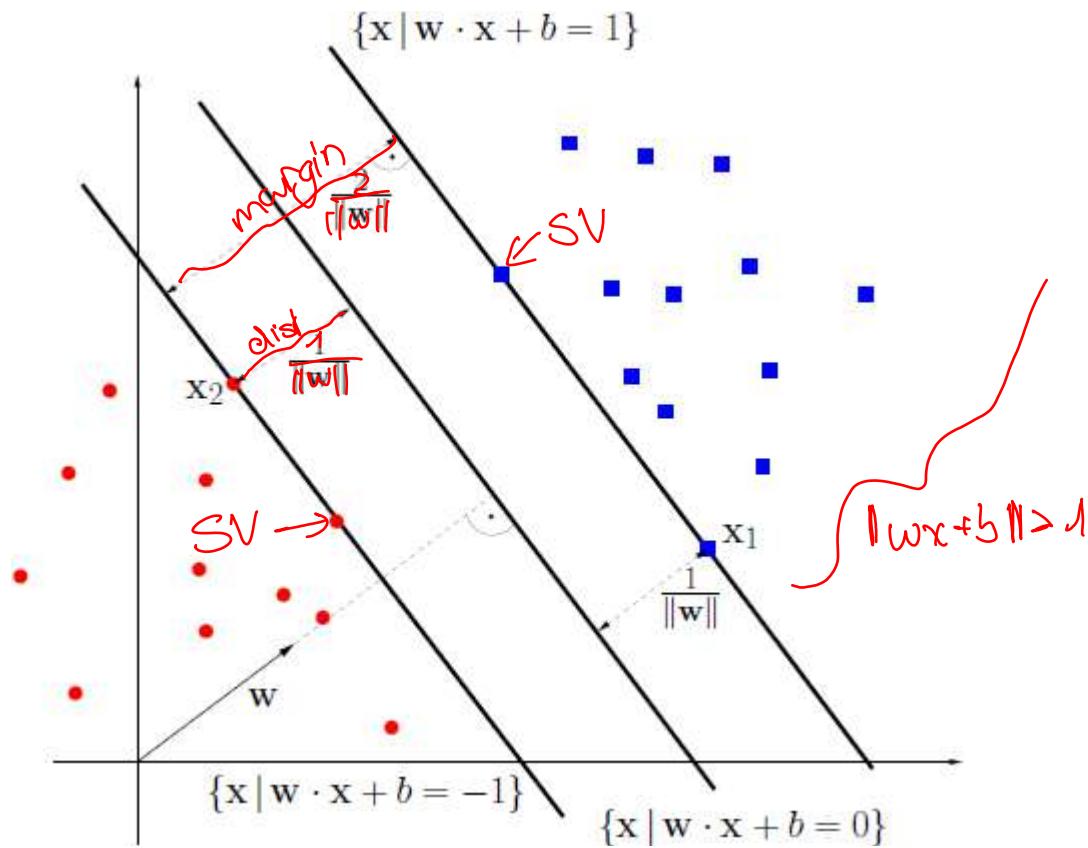
max Zähler von dist

- Rescaling w, b such that $|\langle w, x \rangle + b| = 1$ for the points closest to the hyperplane, we obtain $|\langle w, x \rangle + b| \geq 1$. The **support vectors** (SVs) are those $\{x_i / |\langle w, x_i \rangle + b| = 1\}$. x_i s.t. $1 - 1 = 1 \Leftarrow$ directly on line
- The new loss is $\max(1 - t_i g(x_i), 0) =: (1 - t_i g(x_i))_+$ (**hinge loss**)
- The **margin** is twice the distance of any SV to the plane π :
 $2 d(x_{SV}, \pi) = 2/\|w\|$, since $|g(x_{SV})| = 1$
 $d(x, \pi) = \frac{|g(x)|}{\|w\|}$
- Therefore we find the **canonical OSH** by solving $\max \text{ margin}$

$$\max_{w,b} \left\{ \frac{2}{\|w\|} \quad / \quad t_i (\langle w, x_i \rangle + b) \geq 1, \quad 1 \leq i \leq n \right\}$$

Support Vector Machines

Geometrical view of the OSH



Support Vector Machines

A look on what's to come

1. The solution for w can be expressed as $w = \sum_{i=1}^n t_i \alpha_i x_i, \alpha_i \geq 0.$
(as a consequence of the **Representer theorem**) $P(x) = \langle \hat{w}, \phi(x) \rangle$
2. A fraction of the training data vectors will have $\alpha_i = 0$ (**sparsity**, as a consequence of the chosen error function)
3. The x_i for which $\alpha_i > 0$ will coincide with the **support vectors**
4. The **discriminant function** (classifier) is written $f_{\text{SVM}}(x) = \text{sgn}(\langle w, x \rangle + b) = \text{sgn} \left(\sum_{i=1}^n t_i \alpha_i \langle x, x_i \rangle + b \right)$

$$f_{\text{SVM}}(x) = \text{sgn}(\langle w, x \rangle + b) = \text{sgn} \left(\sum_{i=1}^n t_i \alpha_i \langle x, x_i \rangle + b \right)$$

Support Vector Machines

More than an intuition

- Separating hyperplanes in \mathbb{R}^d have VC dimension $d + 1$
- When we use a feature map into a very high dimension $D \in (\mathbb{N} \cup \{\infty\})$, VC dimension will grow accordingly
- If we bound the margin of the hyperplanes, we limit VC dimension (therefore, we have an explicit control on complexity)

Support Vector Machines

More than an intuition

Theorem. Consider canonical hyperplanes $f(\mathbf{x}) = \text{sgn}(\langle \mathbf{w}, \mathbf{x} \rangle + b)$ and a data set $D = \{(\mathbf{x}_1, t_1), \dots, (\mathbf{x}_n, t_n)\}$, with $\mathbf{x}_i \in \mathbb{R}^d$ and $t_i \in \{-1, +1\}$. The **subclass** of linear classifiers with margin $m \geq m_0$ has VC dimension ϑ bounded by

VC dim bounded by dist/radius

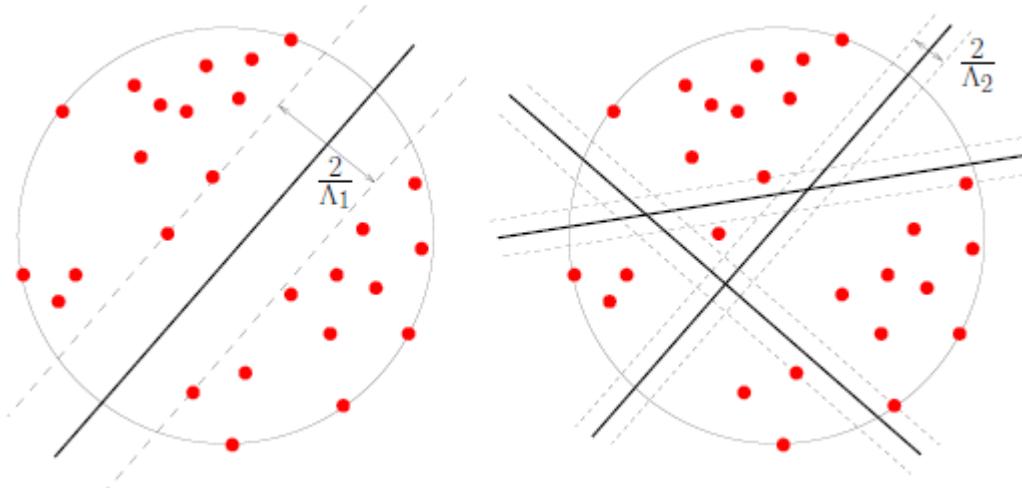
$$\vartheta \leq \min \left(\left\lceil \frac{R^2}{m_0^2} \right\rceil, d \right) + 1$$

$$\vartheta = \min \left(\frac{(R/m)^2}{m_0^2} \right) \text{oder } \dim$$

where R is the radius of the smallest sphere centered at the origin containing the \mathbf{x}_i .

Support Vector Machines

More than an intuition



smaller margin: more ways to shatter (2^v possib)

- Left: hyperplanes with a large margin have reduced chances to separate the data (the VC dimension is small)
- Right: smaller margins allow more separating hyperplanes (the VC dimension is large)

Support Vector Machines

Formulation

\max (weit Kehrwest) Margin = $\frac{z}{\|w\|}$

$$\underset{\substack{w,b \\ (\text{loss})}}{\text{minimize}} \quad \frac{1}{2} \|w\|^2 \leftarrow \frac{\|w\|^2}{z}$$

Square loss

$$\text{subject to } t_i (\langle w, x_i \rangle + b) \geq 1, \quad 1 \leq i \leq n$$

This is solved (numerically) by QP techniques:

- Quadratic (therefore convex) function subject to linear constraints
- Unique solution (or set of equivalent ones)
- Therefore, no local minima

Support Vector Machines

Formulation

For the set of constraints to be satisfied, the data set must be linsep; this is a very unrealistic requirement in practice

- We could aim at minimizing the number of violated constraints $|\{n \mid t_i(\langle w, x_i \rangle + b) < 1\}|$, but this turns out to be NP-hard ...
- Instead, we can minimize a convex function of w :

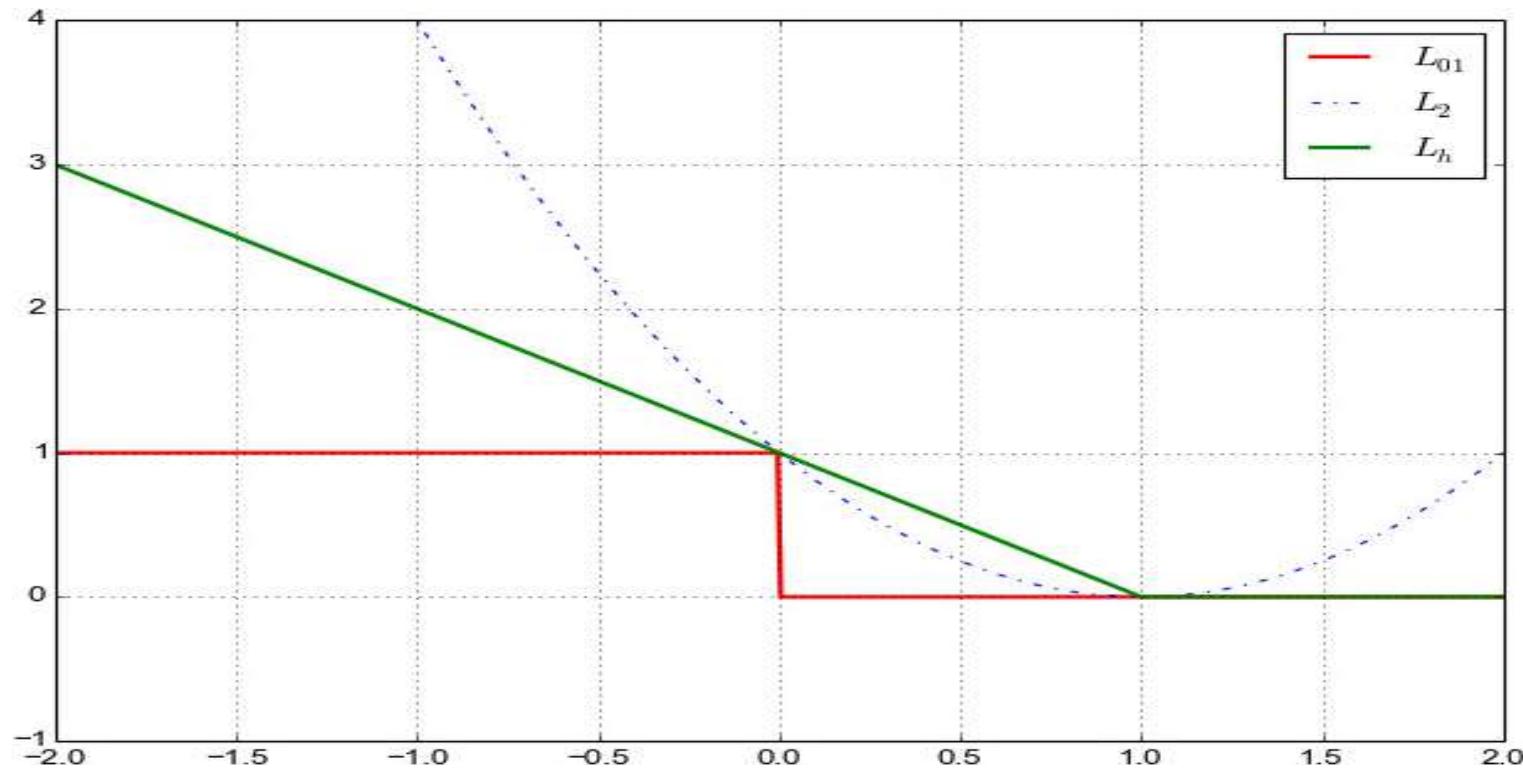
$$\underset{w,b}{\text{minimize}} \quad \frac{1}{2} \|w\|^2 + C \sum_{i=1}^n (1 - t_i g(x_i))^+, \quad C > 0$$

new point
error term

- Yes, the new term is the total hinge loss!
- cost param regularizing sol → doesn't directly affect margin value*

Support Vector Machines

Formulation



L_{01} is the 0/1 loss; L_2 is the square loss; L_h is the hinge loss

$$\text{Combi} \quad \frac{\|w\|^2}{2} + C \sum_{i=1}^n \ell(g(x_i))$$

Support Vector Machines

Margin violations

- quad prob
• This problem is rewritten as another QP, by introducing a set of margin violations ε_i –called **slack variables**–, for each x_i :

$$\underset{\mathbf{w}, b, \{\varepsilon_i\}}{\text{minimize}} \quad \frac{1}{2} \|\mathbf{w}\|^2 + C \sum_{i=1}^n \varepsilon_i$$

subject to $t_i (\underbrace{\langle \mathbf{w}, \mathbf{x}_i \rangle + b}_{\text{error term}}) \geq 1 - \varepsilon_i$ and $\varepsilon_i \geq 0$ ($1 \leq i \leq n$)
 $\Rightarrow \varepsilon_i \geq 1 - t_i g(\mathbf{x}) \geq 0$

- This is a **soft** margin ($\varepsilon_i > 0$ implying \mathbf{x}_i would violate the original constraint)
not \geq
- For a training error to occur $\varepsilon_i > 1$ and so $\sum_{i=1}^n \varepsilon_i$ is an **upper bound** on the number of training errors
- The optimal slacks satisfy $\varepsilon_i = (1 - t_i g(\mathbf{x}))_+$

Support Vector Machines

Excursion: Lagrange multipliers

The famous method of **Lagrange multipliers** allows the optimization of smooth functions subject to equality constraints. (local optimum)

The Karush, Kuhn and Tucker (KKT) theory extends Lagrange's method to include **inequality constraints**.

Consider the problem of minimizing $f(\mathbf{x})$ in a convex $\Omega \subset \mathbb{R}^d$, subject to:

- $g_j(\mathbf{x}) \leq 0$ affine functions, $1 \leq j \leq k$
- $h_j(\mathbf{x}) = 0$ affine functions, $1 \leq j \leq l$

Support Vector Machines

Excursion: Lagrange multipliers

Define the **Lagrangian** as:

$$\mathcal{L}(x, \alpha, \beta) = f(x) + \sum_{j=1}^k \alpha_j g_j(x) + \sum_{j=1}^l \beta_j h_j(x)$$

Annotations in blue:

- A blue arrow points from the word "Funktion" to the term $f(x)$.
- A blue arrow points from the word "equality constraint" to the term $\sum_{j=1}^k \alpha_j g_j(x)$.
- A blue arrow points from the word "multiplier" to the term $\sum_{j=1}^l \beta_j h_j(x)$.

where f, g_j, h_j are continuously differentiable functions.

Support Vector Machines

Excursion: Lagrange multipliers

Theorem. Necessary and sufficient conditions for a point x^* to be an optimum are the existence of $\alpha^* \in \mathbb{R}^k$ and $\beta^* \in \mathbb{R}^l$ such that:

ableiten

1. $\frac{\partial \mathcal{L}(x^*, \alpha^*, \beta^*)}{\partial x} = 0$

2. $\frac{\partial \mathcal{L}(x^*, \alpha^*, \beta^*)}{\partial \beta} = 0$

nach x & β ableiten
↳ beides = 0

3. $\alpha_j^* g_j(x^*) = 0, 1 \leq j \leq k$ (KKT complementarity conditions)

For opt $\alpha^* \rightarrow \alpha^* \cdot g(x^*) = 0$

4. $g_j(x^*) \leq 0, 1 \leq j \leq k$

≥ 0 ≤ 0

5. $\alpha_j^* \geq 0, 1 \leq j \leq k$

eins von beiden ist 0

Support Vector Machines

SVM Lagrangian (primal)

We construct the **Lagrangian**:

$$\mathcal{L} = \frac{1}{2} \|\mathbf{w}\|^2 - \sum_{i=1}^n \alpha_i \left\{ t_i (\langle \mathbf{w}, \mathbf{x}_i \rangle + b) - 1 + \varepsilon_i \right\} + C \sum_{i=1}^n \varepsilon_i - \sum_{i=1}^n \mu_i \varepsilon_i$$

only for error term

cost param *multiplier*
s.t. $\varepsilon \geq 0$

- The $\alpha_i, \mu_i \geq 0$ are the **Lagrange multipliers**; the μ_i ensure that $\varepsilon_i \geq 0$
- The solution is a **saddle point** of \mathcal{L} : minimum w.r.t. \mathbf{w}, b and the ε_i and maximum w.r.t. the α_i and μ_i

Support Vector Machines

Lagrangian form

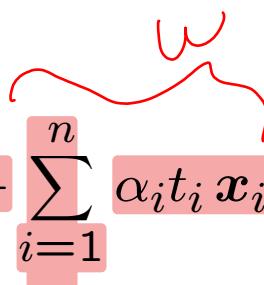
The gradient of \mathcal{L} with respect to w, b and ε_i must vanish:

s.f. only $\alpha_i t_i x_i$ left
 \Rightarrow Kernel

$$\frac{\partial \mathcal{L}}{\partial b} = \sum_{i=1}^n \alpha_i t_i = 0,$$

$$\frac{\partial \mathcal{L}}{\partial w} = w - \sum_{i=1}^n \alpha_i t_i x_i = 0,$$

$$\frac{\partial \mathcal{L}}{\partial \varepsilon_i} = C - \alpha_i - \mu_i = 0$$



In addition, the KKT complementarity conditions must hold:

$$\alpha_i \left(t_i (\langle w, x_i \rangle + b) - 1 + \varepsilon_i \right) = 0$$

$$\begin{aligned}
 L(w, b, a) &= \frac{1}{2} \|w\|^2 - \sum_{i=1}^n \underbrace{\alpha_i [t_i (w^\top x_i + b) - 1]}_{t_i (\langle w, x_i \rangle + b)} \\
 &= \frac{1}{2} \|\sum_{i=1}^n \alpha_i t_i x_i\|^2 - \sum_{i=1}^n \alpha_i \left[t_i \left(\left(\sum_j \alpha_j t_j x_j \right)^\top x_i + b \right) - 1 \right] \\
 &\quad - b \sum_{i=1}^n \alpha_i t_i - (\sum_{i=1}^n \alpha_i t_i x_i)^\top (\sum_{i=1}^n \alpha_i t_i x_i) \\
 &= \sum_{i=1}^n \alpha_i - \frac{1}{2} (\sum_{i=1}^n \alpha_i t_i x_i)^\top (\sum_{i=1}^n \alpha_i t_i x_i) \\
 &= \sum_{i=1}^n \alpha_i - \frac{1}{2} \sum_j \sum_j \alpha_i \alpha_j t_i t_j (x_i^\top x_j) \\
 &= \sum_{i=1}^n \alpha_i - \frac{1}{2} \sum_j \sum_j \alpha_i \alpha_j t_i t_j (x_i \cdot x_j)
 \end{aligned}$$

$$|w = \sum \alpha_i t_i x_i$$

Support Vector Machines

Dual formulation

The Lagrangian \mathcal{L} is convex; its optimization is equivalent to the maximization of its concave **dual problem** \mathcal{L}_D :

(over n)

$j \subseteq k,l$

minimize
 $w, b, \{\alpha_i\}$

$$\mathcal{L}_D = \sum_{i=1}^n \alpha_i - \frac{1}{2} \sum_{i=1}^n \sum_{j=1}^n \alpha_i \alpha_j t_i t_j \langle x_i, x_j \rangle$$

subject to $0 \leq \alpha_i \leq C$ ($1 \leq i \leq n$), and $\sum_{i=1}^n \alpha_i t_i = 0$

- Neither $\mu_i, \varepsilon_i, w, b$ appear in the dual form; maximization is only w.r.t. the α_i
- This optimization problem is expressed only in terms of inner products of the data points: the dual lends itself to kernelisation
- How many free parameters? n (independent of data dimension)

Support Vector Machines

Dual formulation

A closer look at the KKT complementarity conditions:

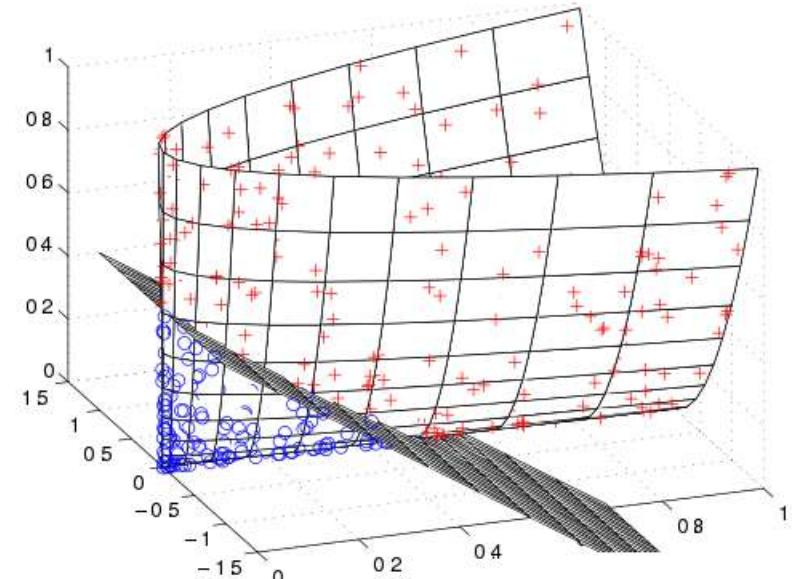
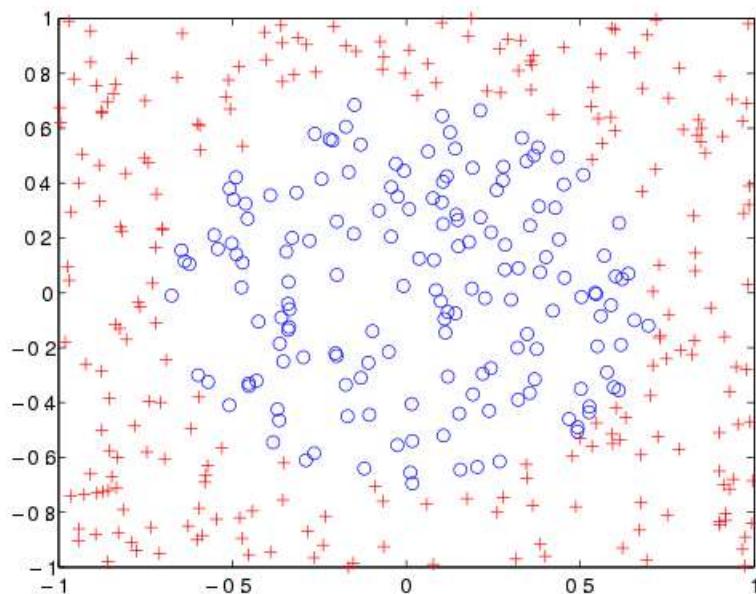
- $\alpha_i = 0$ implies $t_i g(x_i) > 1$ and $\varepsilon_i = 0$ (x_i is **not a SV**) *Liegt außerhalb Margin*
- $\alpha_i \in (0, C)$ implies $t_i g(x_i) = 1$ and $\varepsilon_i = 0$ (x_i is a **non-bound SV**) *Kann auf Margin liegen, aber nicht immer SV*
- $\alpha_i = C$ implies $t_i g(x_i) < 1$ and $\varepsilon_i > 0$ (x_i is a **bound SV**)
(in particular, $\varepsilon_i > 1$ implies x_i is a **training error**)

Support Vector Machines

The SVM goes non-linear

Recall the idea of mapping input data into some Hilbert space (called the **feature space**) via a non-linear mapping $\phi : \mathcal{X} \rightarrow \mathcal{H}$

The associated kernel function is $k(x, x') = \langle \phi(x), \phi(x') \rangle_{\mathcal{H}}$, $x, x' \in \mathcal{X}$



Support Vector Machines

SVM kernelization

- We now substitute x_i by $\phi(x_i)$, then build the OSH in \mathcal{H}
- The dual of the new QP problem is formulated exactly as before, replacing $\langle x_i, x_j \rangle$ with $\langle \phi(x_i), \phi(x_j) \rangle_{\mathcal{H}} = k(x_i, x_j)$
- The discriminant function becomes:

$$= \text{sgn}(\underbrace{\omega^\top \phi(x)}_{k(x, x_i)} + b) = \text{sgn}\left(\sum_{i=1}^n \alpha_i t_i \underbrace{\phi(x_i)^\top \phi(x)}_{k(x, x_i)} + b\right)$$

$$f_{\text{SVM}}(x) = \text{sgn}\left(\sum_{i=1}^n \alpha_i t_i k(x, x_i) + b\right)$$

Support Vector Machines

LOOCV bounds (I)

A rough but simple bound on LOOCV (leave-one-out CV) error can be computed as:

$$\text{LOOCV}(n) \leq \frac{1}{n} \mathbb{E}(n_{SV})$$

n_{SV} is the number of SVs for a given sample of size n

The $\mathbb{E}()$ is taken over all such samples

Support Vector Machines

LOO bounds (II)

Theorem. The LOOCV error of a stable SVM^(*) on a set of training patterns x_i is bounded by $|\{i \mid (2\alpha_i R^2 + \varepsilon_i) \geq 1\}|/n$, where R^2 is an upper bound on $k(x, x)$ and $k(x, x') \geq 0$.

$$\leq \frac{|\{i \mid 2\alpha_i R^2 + \varepsilon_i \geq 1\}|}{n}$$

- This quantity can be extracted easily from the solution
- This LOOCV error is an unbiased estimate of true error

(*) A SVM is stable if there is at least one non-bound SV (see *Estimating the Generalization Performance of a SVM Efficiently*. T. Joachims; In ICML, 2000)

Support Vector Machines

Final remarks (I)

- The fact that the **OSH** is determined only by the support vectors is most remarkable, since usually this number will be small
- The **support vectors** (SVs) are:
 1. the only training examples that define the solution
 2. the most difficult examples to classify
- This means all the **relevant information** in the data set is summarized by the **SVs**: we would have obtained the same result by using *only* the SVs from the outset

Support Vector Machines

Final remarks (II)

primal

- The SVM is specially well suited for “large d , low n ” problems, because:
many param, less data
 1. complexity grows with n (non-parametric model)
 2. space requirements (the kernel matrix) also grows with n
 3. generalization error does not depend on d
- The “architecture” is determined automatically by the method (not by experimentation, as in neural networks)

Support Vector Machines

Hot topics

- Choice of the **best kernel** is an open issue; **kernel design** is an active area of research
- More efficient algorithms for solving **big QP** problems are being developed
- Sometimes the **fraction of SVs** is very high (indicating a poor fit); it is possible to control this fraction directly (ν -SVMs)
- Performance usually depends on a careful choice of the external parameters: C and those of the kernel function; we need principled ways for **hyper-parameter** selection

Support Vector Machines

Where to look for more ...

- *An Introduction to Kernel-based Learning Algorithms.* K.-R. Mueller, S. Mika, G. Raetsch, K. Tsuda, and B. Schoelkopf, IEEE Neural Networks, 12(2):181-201, 2001.
- *A Tutorial on Support Vector Machines for Pattern Recognition.* Christopher Burges.
<https://research.microsoft.com/en-us/um/people/cburges/papers/svmtutorial.pdf>
- *Learning with Kernels: Support Vector Machines, Regularization, Optimization, and Beyond.* Bernhard Schoelkopf and Alexander J. Smola, MIT Press, 2001.
- *An Introduction to Support Vector Machines and Other Kernel-based Learning Methods.* Nello Cristianini and John Shawe-Taylor, Cambridge University Press, 2000.
- *Kernel Methods for Pattern Analysis.* John Shawe-Taylor and Nello Cristianini, Cambridge University Press, 2004.
- *The Nature of Statistical Learning Theory.* V. Vapnik, Springer, 2nd ed., 1999

Kernel-Based Learning & Multivariate Modeling

MIRI Master

Lluís A. Belanche

`belanche@cs.upc.edu`

Soft Computing Research Group

Universitat Politècnica de Catalunya

2021-2022

Kernel-Based Learning & Multivariate Modeling

Syllabus

Part 1

Sep 15 Introduction to kernel-based learning

Sep 22 The SVM for classification, regression & novelty detection (I)

Sep 29 The SVM for classification, regression & novelty detection (II)

Oct 06 Kernel design (I): theoretical issues

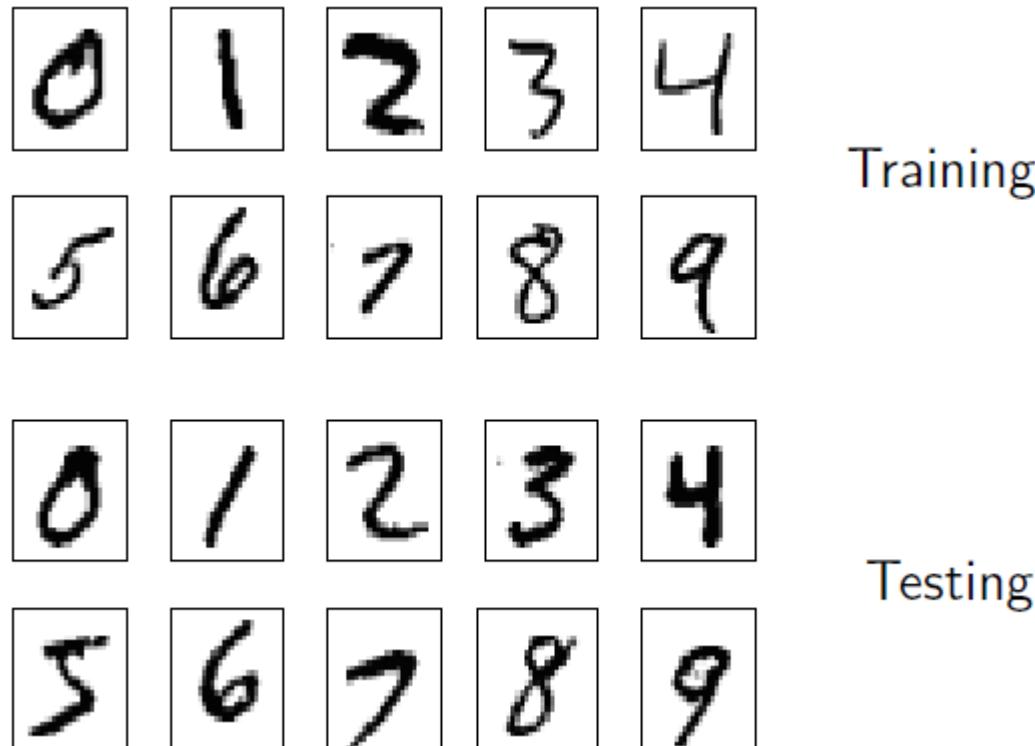
Oct 13 Kernel design (II): practical issues

Oct 20 Kernelizing ML & statistical algorithms

Oct 27 Advanced topics

Support Vector Machines

Application (digit recognition)



- Handwritten zip code recognition traces back to the 1960's

Support Vector Machines

Application (digit recognition)

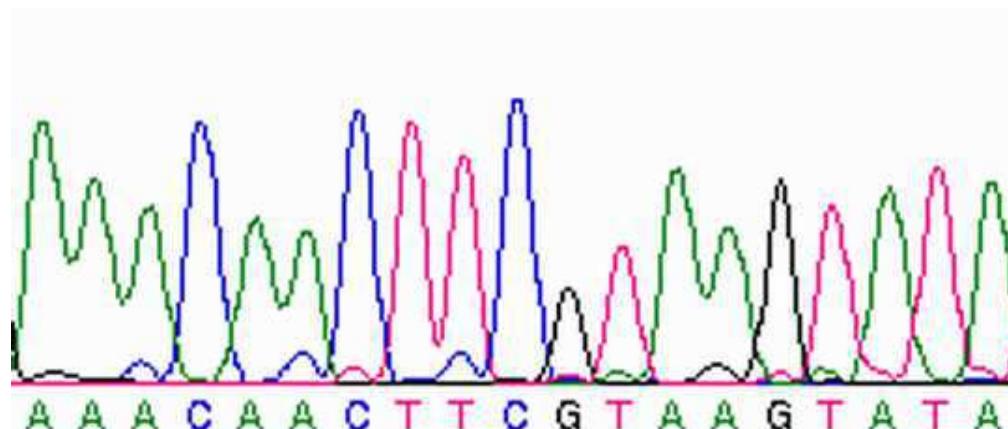
- MNIST handwritten zip code recognition
- 60.000 training, 10.000 test examples (28×28 pixels)

Classifier	test error
Linear classifier	8.4 %
3-nearest-neighbor	2.4 %
SVM	1.4 %
Tangent distance	1.1 %
LeNet4	1.1 %
Boosted LeNet4	0.7 %
Translation invariant SVM	0.56 %

Support Vector Machines

Application: Classification of DNA sequences

- A *promoter* is a region of DNA that initiates or facilitates transcription of a particular gene
- The dataset consists of 106 DNA sequences described by 57 categorical variables, represented as the nucleotide at each position: [A]denine, [C]ytosine, [G]uanine, [T]hymine
- The response is a binary class: “+” for a promoter gene and “-” for a non-promoter gene



Support Vector Machines

Application: Classification of DNA sequences

The similarity between two multivariate categorical vectors is the fraction of the number of matching values.

Overlap/Dirac kernel]

$$k_0(x, x') = \frac{1}{d} \sum_{i=1}^d \mathbb{I}_{\{x_i = x'_i\}}$$

no matching values

Another kernel that can be used is the RBF kernel:

Gaussian Radial Basis Function kernel:

$$k_{\text{RBF}}(x, x') = \exp(-\gamma \|x - x'\|^2), \gamma > 0$$

In order to use this kernel, categorical variables with m modalities are coded using a binary expansion representation.

Support Vector Machines

Application: Classification of DNA sequences

Univariate kernel $k_1^{(U)}$:

$$k_1^{(U)}(x_i, x'_i) = \begin{cases} h_\alpha(P_Z(x_i)) & \text{if } x_i = x'_i \\ 0 & \text{if } x_i \neq x'_i \end{cases}$$

?

where

$$h_\alpha(z) = (1 - z^\alpha)^{1/\alpha}, \quad \alpha > 0$$

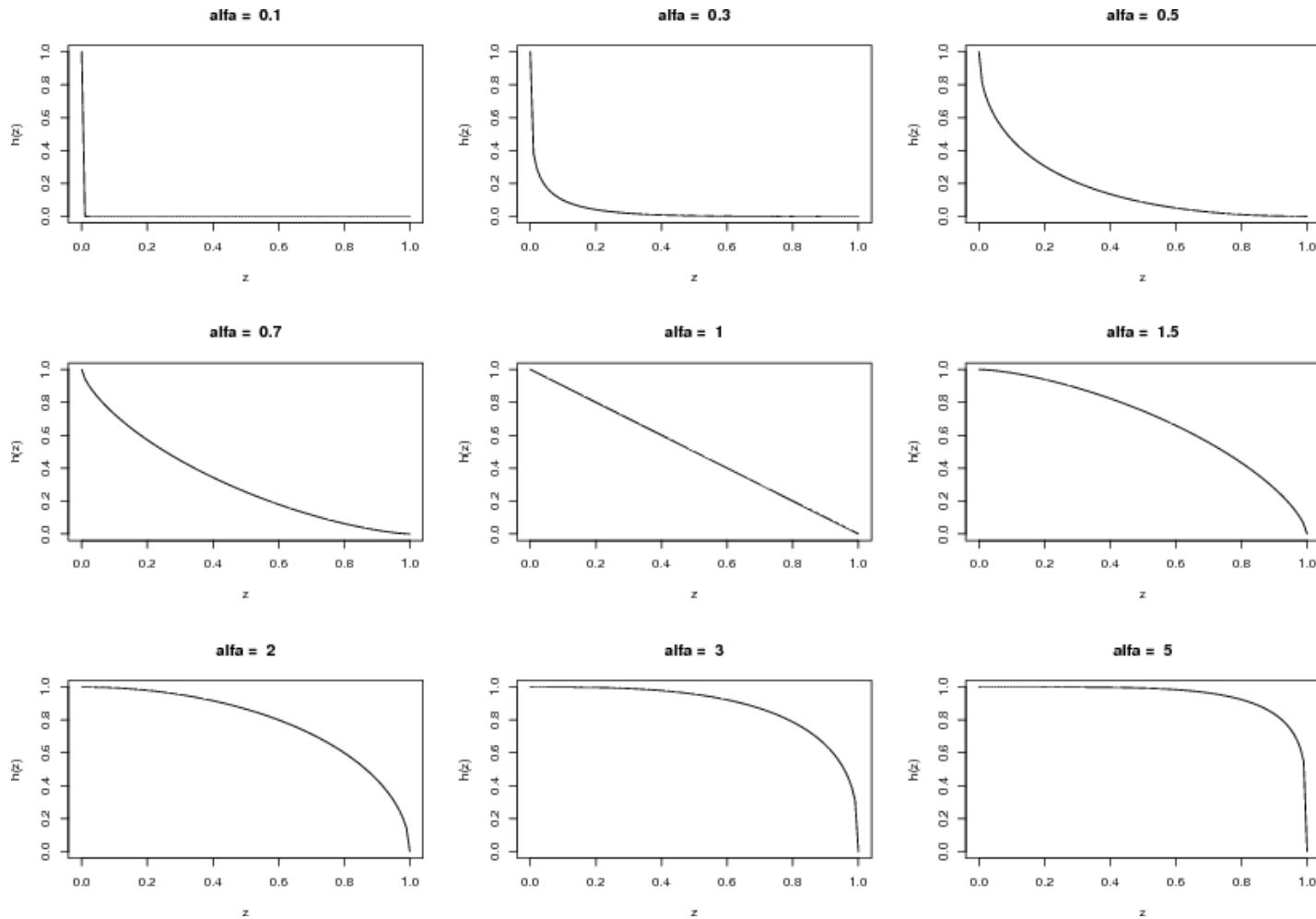
Multivariate kernel k_1 :

$$k_1(x, x') = \exp\left(\frac{\gamma}{d} \sum_{i=1}^d k_1^{(U)}(x_i, x'_i)\right), \quad \gamma > 0$$

The kernel matrices generated by k_1 are p.s.d.

Support Vector Machines

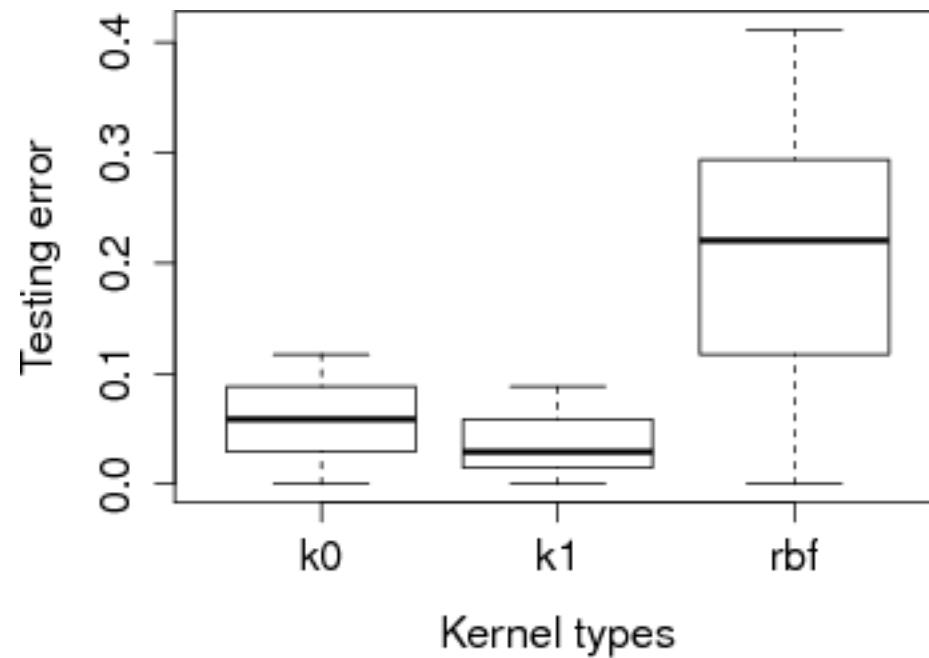
Application: Classification of DNA sequences



The family of inverting functions $h_\alpha(z)$ for different values of α

Support Vector Machines

Application: Classification of DNA sequences



Test error distributions on the PromoterGene problem

(joint work with M. Villegas)

Support Vector Machines

The role of the C parameter

Increasing the value of C ...

(only recommended)

can affect but not always

- penalizes margin errors more \Rightarrow narrower margin \Rightarrow larger VC-dimension
- allows the $\alpha_i \leq C$ to be larger (so more opportunities for outliers)
- increases training times

Support Vector Machines

The role of the C parameter

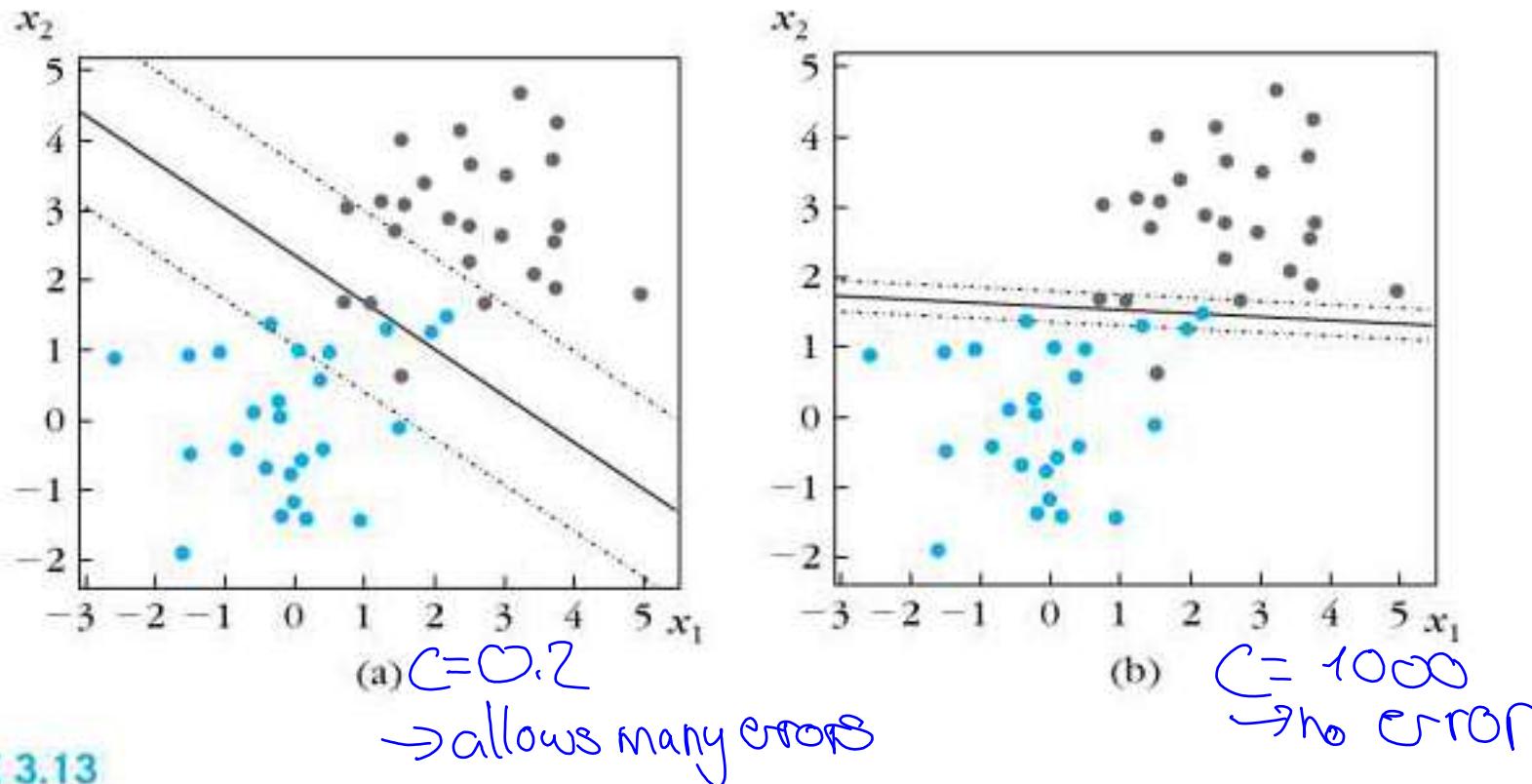


FIGURE 3.13

An example of two nonseparable classes and the resulting SVM linear classifier (full line) with the associated margin (dotted lines) for the values (a) $C = 0.2$ and (b) $C = 1000$. In the latter case, the location and direction of the classifier as well as the width of the margin have changed in order to include a smaller number of points inside the margin.

Support Vector Machines

ν -SVMs

There are two commonly used versions of the SVM for classification:

'C-SVC': original SVM formulation, uses a parameter $C \in (0, \infty)$ to apply a **penalty to the optimization** for those data points not entirely separated by the OSH (violating the margins) **penalty for opt**

'nu-SVC': C is replaced by $\nu \in (0, 1)$: **lower & upper bound**

- **upper bound** on the fraction of examples which are **training errors** (missclassified)
- **lower bound** on the fraction of **points** which are **SVs**.

Support Vector Machines

SVMs for regression

“The Support Vector method can also be applied to the case of regression, maintaining all the main features that characterise the maximal margin algorithm: a non-linear function is learned by a linear learning machine in a kernel-induced feature space while the capacity of the system is controlled by a parameter that does not depend on the dimensionality of the space.”

—from N. Cristianini and J. Shawe-Taylor, *An introduction to Support Vector Machines* (2000)

Support Vector Machines

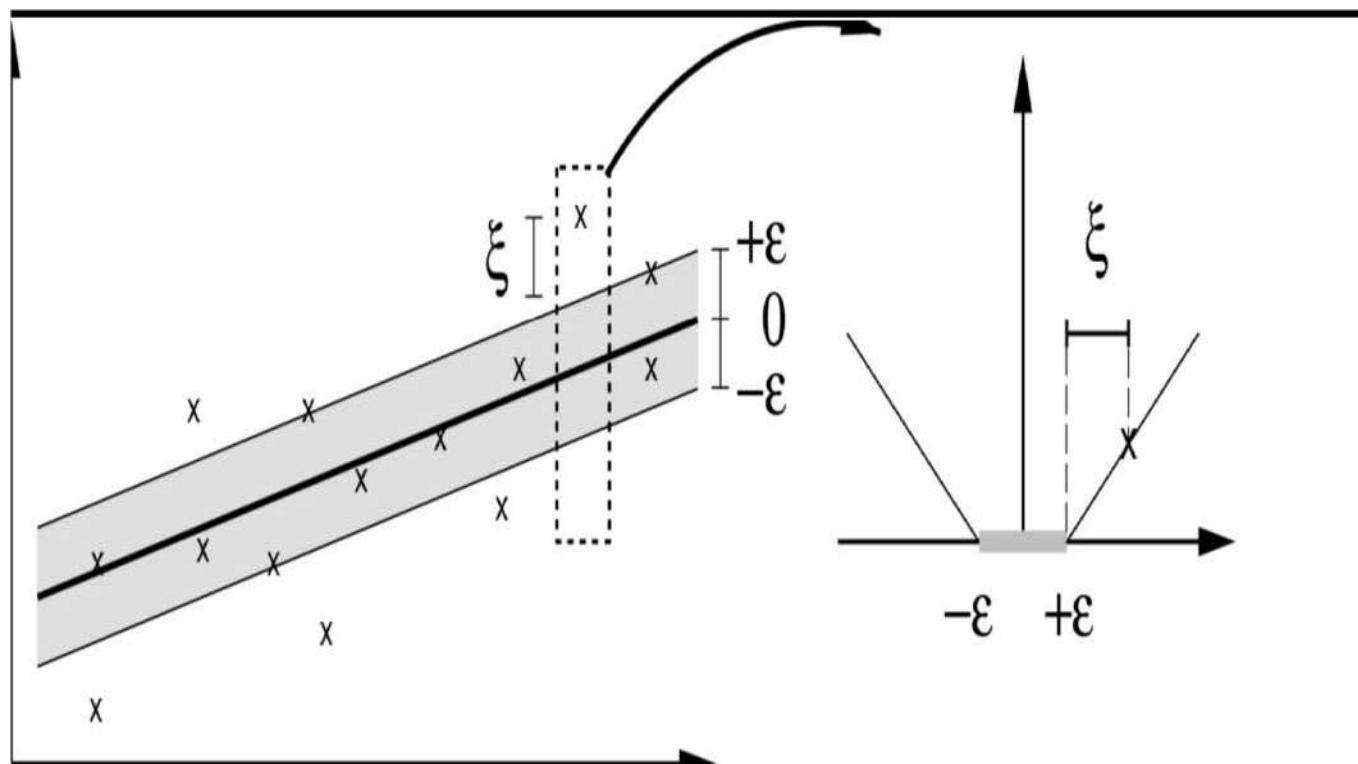
SVMs for regression

Here we choose the **ε -insensitive** loss:

$$L(t_i, \langle \mathbf{w}, \mathbf{x}_i \rangle) = |t_i - g(\mathbf{x}_i)|_\varepsilon = \max(|t_i - g(\mathbf{x}_i)| - \varepsilon, 0)$$

algebraic error
dist to line

where $g(\mathbf{x}) = \langle \mathbf{w}, \mathbf{x} \rangle + b$



Support Vector Machines

SVMs for regression

minimize $\frac{1}{2} \|w\|^2 + C \sum_{i=1}^n (\xi_i + \xi_i^*)$

↑
cost param

unter oben

subject to $\langle w, x_i \rangle + b - t_i \leq \varepsilon + \xi_i,$
 $t_i - \langle w, x_i \rangle + b \leq \varepsilon + \xi_i^*,$
 $\xi_i, \xi_i^* \geq 0$

where the ξ_i, ξ_i^* are again slack variables controlling the “violations”

Support Vector Machines

SVMs for regression

Then feature maps $\phi(\cdot)$ are introduced, the primal optimization problem is transformed into the dual, and kernelised, to give:

$$y_{\text{SVM}}(x) = \sum_{i=1}^n \beta_i k(x, x_i)$$

if $\xi = 0 \Rightarrow \beta = 0$, (not SV)
otherwise $\neq 0$ (is SV)
 \Leftrightarrow only use SVs

with $0 \leq \alpha_i, \alpha_i^* \leq C$

$$= \sum_{i=1}^n \beta_i \langle \phi(x), \phi(x_i) \rangle = \sum \beta_i x^T x_i$$

For convenience, we have defined $\beta_i := \alpha_i - \alpha_i^*$.

Support Vector Machines

SVMs for regression

A closer look at the structure of the solution:

- Data points that end up **within** the ε -tube have **inactive slacks** (i.e., $\xi_i = \xi_i^* = 0$) and therefore $\beta_i = 0$ (**not SVs**) *between lines*
- Data points that end up **not within** the ε -tube have exactly one active slack (i.e., either $\xi_i > 0$ and $\xi_i^* = 0$, or vice versa) and therefore $\beta_i \neq 0$ (**non-bound SVs**) *One on line, other not*
- Data points that end up **outside** the ε -tube have exactly one bound slack (i.e., $\xi_i = C$ and $\xi_i^* = 0$, or vice versa) and therefore $\beta_i \neq 0$ (**bound SVs**) *one on line, one outside*
(0,∞) penalizes

Support Vector Machines

SVMs for regression

In comparison to ridge regression, the only difference is in the choice of the loss (since both are **regularized machines** and both are amenable to **kernelisation**) and its consequences:

- Deviations lower than ε are ignored
 $\text{dev} < \varepsilon$ ignored
 \leftarrow ridge
- The loss grows linearly (and not quadratically) in the residual, making it more robust against outliers
- The solution is **sparse** (the number of **basis functions** $\phi(x_i)$ is automatically adapted)

Support Vector Machines

C versus ε

- C determines the trade off between model complexity (flatness) and tolerance to deviations larger than ε
- ε controls the width of the ε -insensitive tube

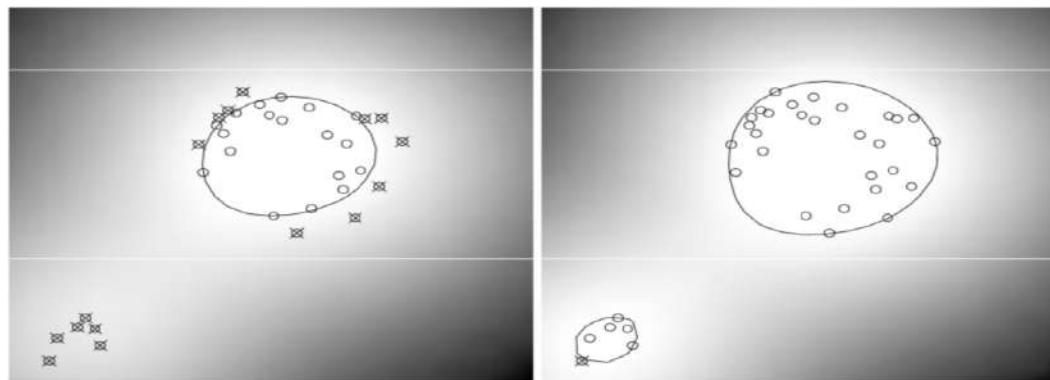
Larger ε or C implies less SVs (while smaller ε or C implies more SVs); but larger ε gives flatter models while larger C implies more complex models

Hence, both parameters affect model complexity and number of SVs (but in a different way).

Support Vector Machines

SVMs for novelty detection (I)

- You are given a dataset drawn from a pdf $p(x)$; the x can be handwritten digits (recognizable/strange), process status (normal/faulty), credit card transactions (normal/fraudulent), ...
- The goal is to estimate a “simple” subset S of input space s.t. the *probability* that a test point drawn from p lies *outside* S equals some a priori specified $\rho \in (0, 1)$:



—from Alex Smola: Hilbert Space Methods: Basics, Applications, Open Problems
<http://alex.smola.org/talks/rsisesvm.pdf>

Support Vector Machines

SVMs for novelty detection (II)

USPS dataset of handwritten digits: 9,298 digit images of size 16×16

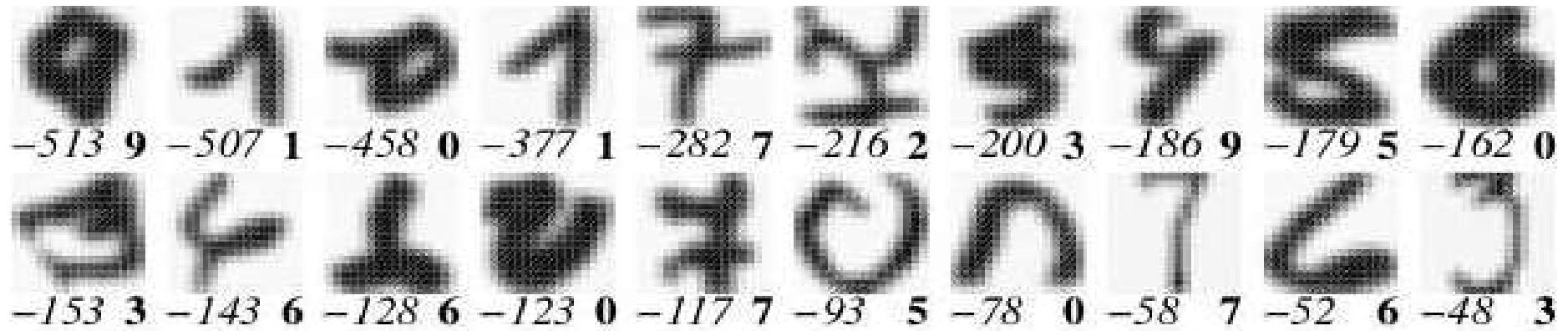


Figure 2: Outliers identified by the proposed algorithm, ranked by the negative output of the SVM (the argument of the sgn in the decision function). The outputs (for convenience in units of 10^{-5}) are written underneath each image in italicics, the (alleged) class labels are given in bold face. Note that most of the examples are “difficult” in that they are either atypical or even mislabelled.

The 20 worst outliers for the USPS test set (here $\rho = 0.05$)

—from Schölkopf et al, *Support Vector Method for Novelty Detection*, NIPS'2000

Support Vector Machines

Tricks of the trade for the kernel

1. Standardizing the variables is in general good (assumed numerical)

$$x_S = \frac{x - \mu_x}{\sigma_x}$$

2. Kernel matrices close to the identity or close to the “all-ones” matrix are also an indication of bad kernel parameter: avoid these situations

3. Kernel matrix values should not be very large or very small (or both): if so, normalize the kernel matrix

$$K_n(x, x') = \frac{k(x, x')}{\sqrt{K(x, x')} \sqrt{K(x, x')}} = \frac{\langle \phi(x), \phi(x') \rangle}{\|\phi(x)\| \cdot \|\phi(x')\|} = \cos \langle \phi(x), \phi(x') \rangle$$

$\in [-1, 1]$

Kernel-Based Learning & Multivariate Modeling

MIRI Master

Lluís A. Belanche

`belanche@cs.upc.edu`

Soft Computing Research Group

Universitat Politècnica de Catalunya

2021-2022

Kernel-Based Learning & Multivariate Modeling

Syllabus

Part 1

Sep 15 Introduction to kernel-based learning

Sep 22 The SVM for classification, regression & novelty detection (I)

Sep 29 The SVM for classification, regression & novelty detection (II)

Oct 06 Kernel design (I): theoretical issues

Oct 13 Kernel design (II): practical issues

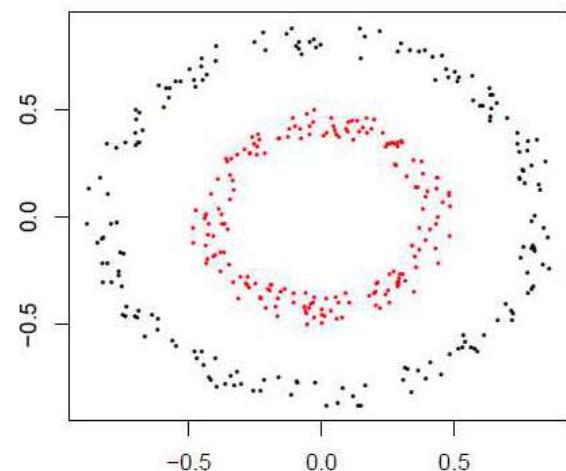
Oct 20 Kernelizing ML & statistical algorithms

Oct 27 Advanced topics

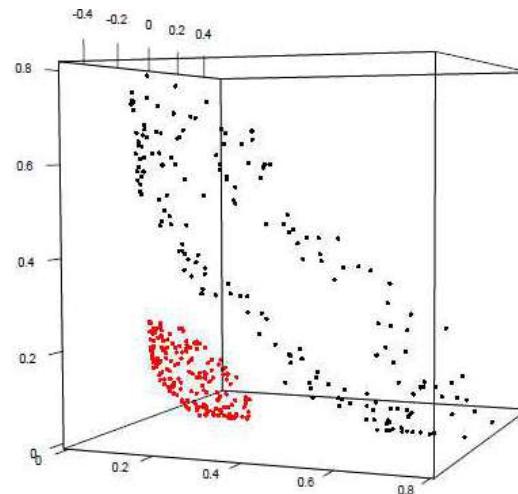
Kernel design (I): theoretical issues

General feature maps

Recall the idea of mapping input data into some Hilbert space (called the *feature space*) via a non-linear mapping $\phi : \mathcal{X} \rightarrow \mathcal{H}$



(a) Input Space (data not linearly separable)



(b) Feature Space (data linearly separable)

Kernel design (I): theoretical issues

Hilbert spaces

An abstract complete **vector space** endowed with an inner product:

Inner product requires symmetry, bilinearity and PSD-ness

Completeness means all Cauchy sequences converge to an element within the space (w.r.t. the norm induced by the inner product)

Show by adding & multiplying elem in $H \rightarrow$ result in H

Kernel design (I): theoretical issues

Characterization of Kernels

Given a function $k : \mathcal{X} \times \mathcal{X} \rightarrow \mathbb{R}$, which properties make it a valid kernel function for ML?

\Rightarrow existence of a map $\phi : \mathcal{X} \rightarrow \mathcal{H}$ s.t.

1. \mathcal{H} is a Hilbert space and
2. $k(x, x') = \langle \phi(x), \phi(x') \rangle_{\mathcal{H}}$ holds?

Kernel design (I): theoretical issues

Characterization of Kernels

A symmetric function k is called **positive semi-definite (PSD)** in \mathcal{X} if:

for every $n \in \mathbb{N}$, and every choice $x_1, \dots, x_n \in \mathcal{X}$,

the Gram matrix $\mathbf{K} = (k_{ij})$, where $k_{ij} = k(x_i, x_j)$, is PSD.

Theorem. A function $k : \mathcal{X} \times \mathcal{X} \rightarrow \mathbb{R}$ admits the existence of a map $\phi : \mathcal{X} \rightarrow \mathcal{H}$ s.t. \mathcal{H} is a Hilbert space and $k(x, x') = \langle \phi(x), \phi(x') \rangle_{\mathcal{H}}$ if and only if k is a **symmetric** and **PSD** function in \mathcal{X} .

Prove example on next slide

\Rightarrow We depart from a $\phi: \mathcal{X} \rightarrow \mathcal{H}$. We define $k: \mathcal{X} \times \mathcal{X} \rightarrow \mathbb{R}$

$$(x, x') \mapsto k(x, x') = \langle \phi(x), \phi(x') \rangle_{\mathcal{H}}$$

• Symmetry ✓

• $\forall m \in \mathbb{N}, \forall x_1 - x_n, x_i \in \mathcal{X},$
 $\forall c \in \mathbb{R}^n, c^T K c = \sum_{i=1}^n \sum_{j=1}^n c_i k_{ij} c_j$

$$= \sum_{i=1}^n \sum_{j=1}^n c_i \underbrace{\langle \phi(x_i), \phi(x_j) \rangle_{\mathcal{H}}}_{c_j} c_j$$

$$= \langle \sum_{i=1}^n c_i \phi(x_i), \sum_{j=1}^n c_j \phi(x_j) \rangle_{\mathcal{H}}$$

$$= \left\| \sum_{i=1}^n c_i \phi(x_i) \right\|^2 \geq 0$$

alles in $\|\cdot\|^2$ packen, weil ≥ 0

PSD

Kernel design (I): theoretical issues

On positive semi-definiteness

There are many **equivalent characterizations** of the PSD property for real symmetric matrices. Here are some: $A_{n \times n}$ is PSD if and only if ...

1. all of its eigenvalues are non-negative

$$A \cdot \vec{v} = \lambda \cdot \vec{v}$$

eigenvalue

2. the determinants of all of its leading principal minors are non-negative

* ↘ ← oder Dreiecksmatrix
Diagonale mult of all matrixes $i \times i$ w/ $i=1 \dots n$

3. there is a PSD matrix B such that $BB^T = A$ (this matrix is unique, denoted with $B = A^{1/2}$, and called the **principal square root** of A)

4. $\forall c \in \mathbb{R}^n, c^T A c \geq 0$ the def used by prop on prev slide
↳ zum Quadrat in Erg ideal

Kernel design (I): theoretical issues

Generating the inner product

Given a kernel k symmetric and PSD, consider the space of functions:

$$\phi : \mathcal{X} \rightarrow \mathbb{R}^{\mathcal{X}}$$
$$x \mapsto \phi(x) : k(x, \cdot)$$

stand-in für beliebigen Wert

Define the (soon-to-be) vector space

$$\mathcal{H}_{\text{pre}} := \text{span}\{\phi(x) / x \in \mathcal{X}\}$$

set of all linear
combi's
intersection of subspaces

all possible
finite combos

$$= \left\{ f(\cdot) = \sum_{i=1}^n \alpha_i k(x_i, \cdot) / n \in \mathbb{N}, x_i \in \mathcal{X}, \alpha_i \in \mathbb{R} \right\}$$

linear combi's of $\phi(x) = k(x, \cdot)$

Kernel design (I): theoretical issues

Generating the inner product

Let $f, g \in \mathcal{H}_{\text{pre}}$; define an **inner product** in \mathcal{H}_{pre} as

$$\langle f, g \rangle = \left\langle \sum_{i=1}^n \alpha_i k(x_i, \cdot), \sum_{j=1}^m \beta_j k(x'_j, \cdot) \right\rangle = \sum_{i=1}^n \sum_{j=1}^m \alpha_i \beta_j k(x_i, x'_j)$$

Note that $\langle f, k(x, \cdot) \rangle = \sum_{i=1}^n \alpha_i k(x_i, x) = f(x)$

$$f(\cdot) = \sum \alpha_i k(x_i, \cdot)$$

This is called the **reproducing property** of the kernel

Kernel design (I): theoretical issues

Generating the inner product

Let's check we have a valid inner product space:

1. $\langle f, g \rangle = \langle g, f \rangle$ (symmetry)

2. $\langle f, g \rangle = \sum_{i=1}^n \alpha_i g(x_i) = \sum_{j=1}^m \beta_j f(x'_j)$ (bilinearity)

Funktionen sind gleich
mit bestimmter Gewichtung

3. $\langle f, f \rangle \geq 0$ with equality iff f is the zero function (PSD-ness)
almost always $f(s)=0$

This inner product satisfies the Cauchy-Schwartz inequality:

$$|\langle f, g \rangle| \leq \sqrt{\langle f, f \rangle} \cdot \sqrt{\langle g, g \rangle}, \quad \forall f, g \in \mathcal{H}_{\text{pre}}$$

Kernel design (I): theoretical issues

Generating the inner product

- Once we have an inner product, we have a **norm** $\|f\| := \sqrt{\langle f, f \rangle}$
- Moreover, we have a **metric** $d(f, g) := \|f - g\| = \sqrt{\langle f-g, f-g \rangle}$
- For any metric space, one can construct a **complete** metric space which contains the former as a **dense subspace***; if completion is applied to an inner product space, the result is a **Hilbert space \mathcal{H}**

(*): Let (X, d) be a metric space, and $X_0 \subset X$. Then X_0 is **dense** in X if and only if $\forall x \in X$ there is a sequence of points $x_n \in X_0$ that has **limit** x .

Kernel design (I): theoretical issues

The Kernel Trick

Such a space is called a **Reproducing Kernel Hilbert Space (RKHS)**

Given the mapping $\phi : \mathcal{X} \rightarrow \mathcal{H}$, the **kernel trick** consists in performing the mapping and the inner product simultaneously by defining its associated kernel function:

$$k(x, x') = \langle \phi(x), \phi(x') \rangle_{\mathcal{H}}, \quad x, x' \in \mathcal{X}$$

inner prod w/o map info

This way it is possible to compute inner products in \mathcal{H} without explicitly performing/knowing the map (e.g. Gram matrices, the OSH)

We depart from kernel function: symmetric & PSD
 Let's build associated Hilbert space.

1) Let $\mathcal{H} := \emptyset$

2) $\forall x \in X$, add the function $k_x: X \rightarrow \mathbb{R}^X$
 $x \mapsto k(x, \cdot)$

3) Make \mathcal{H} the span of the set in 2), that is

$$\mathcal{H} := \text{span}\{k_x / x \in X\} = \{f(\cdot) = \sum_{i=1}^n \alpha_i k_{x_i}(\cdot) / \alpha_i \in \mathbb{R}, x_i \in X, i=1, \dots, n\}$$

4) Define inner prod in this space

$$\text{Fig as } f = \sum_{i=1}^n \alpha_i k_{x_i}, g = \sum_{j=1}^m \beta_j k_{x_j} \rightarrow \left\langle \sum_{i=1}^n \alpha_i k_{x_i}, \sum_{j=1}^m \beta_j k_{x_j} \right\rangle_H := \sum_{i=1}^n \sum_{j=1}^m \alpha_i \beta_j k(x_i, x_j)$$

5) We define a norm & a distance

- $\forall f \in \mathcal{H}, \|f\|_{\mathcal{H}} := \sqrt{\langle f, f \rangle_H}$

- $\forall f, g, d_H(f, g) := \|f - g\|_H$

- add fc H Hr (limits of all converging seq (complete))

↳ HS of entirely generated by $k \rightarrow H_k$

Properties of HS generated this way

$\forall f \in \mathcal{H}_K, \langle f, k_x \rangle_{H_K} = f(x)$, the reproducing property of the kernel

\Rightarrow Reproducing kernel Hilbert Space (RKHS)

Take $x, x' \in X$, define $\phi: X \rightarrow \mathbb{R}, x \mapsto \phi(x) := k_x$

$$\langle \phi(x), \phi(x') \rangle_{H_K} = \langle k_x, k_{x'} \rangle_{H_K} = \langle k(x, \cdot), k(x', \cdot) \rangle_{H_K} = k(x, x') = k_x(x')$$

Build Hilbert space
 (same as slides but sum on 1 page)

Norm

- 1) $\|z\| = 0 \Leftrightarrow z = 0$
- 2) $\|a \cdot z\| = |a| \cdot \|z\|$
- 3) $\|z + t\| \leq \|z\| + \|t\|$

Kernel design (I): theoretical issues

The Kernel Trick: an example

Take $k(x, x') = \langle x, x' \rangle^q$, for $x, x' \in \mathbb{R}^d$. What is the underlying feature map ϕ ?

⇒ Answer: the space spanned by all products of exactly q dimensions of \mathbb{R}^d .

Example: $x, x' \in \mathbb{R}^3$, and $q = 2$:

$$\begin{aligned} k(x, x') &= \langle x, x' \rangle^2 = \left\langle \begin{pmatrix} x_1 \\ x_2 \\ x_3 \end{pmatrix}, \begin{pmatrix} x'_1 \\ x'_2 \\ x'_3 \end{pmatrix} \right\rangle^2 \text{ write how inner prod is calc} \\ &= (x_1 x'_1 + x_2 x'_2 + x_3 x'_3)^2 = (x_1 x'_1 + x_2 x'_2)^2 + 2(x_1 x'_1 + x_2 x'_2)x_3 x'_3 + (x_3 x'_3)^2 \\ &= \left\langle \begin{pmatrix} x_1^2 \\ \sqrt{2}x_1 x_2 \\ \sqrt{2}x_1 x_3 \\ \sqrt{2}x_2 x_3 \\ x_2^2 \\ x_3^2 \end{pmatrix}, \begin{pmatrix} (x'_1)^2 \\ \sqrt{2}x'_1 x'_2 \\ \sqrt{2}x'_1 x'_3 \\ \sqrt{2}x'_2 x'_3 \\ (x'_2)^2 \\ (x'_3)^2 \end{pmatrix} \right\rangle \text{ d lang w/ jewels q} \\ &= \langle \phi(x), \phi(x') \rangle \text{ d.q lang} \end{aligned}$$

bring it in form
 $\langle \phi(x), \phi(x') \rangle$

Kernel design (I): theoretical issues

Popular choices for the Kernel

Polynomial kernels (relation to GLDs)

$$k(\mathbf{x}, \mathbf{x}') = (a \langle \mathbf{x}, \mathbf{x}' \rangle + c)^q, \quad q \in \mathbb{N}, a > 0, c \geq 0 \in \mathbb{R}$$

Gaussian kernels (relation to RBFNNs)

$$k(\mathbf{x}, \mathbf{x}') = \exp(-\gamma \|\mathbf{x} - \mathbf{x}'\|^2), \quad \gamma > 0 \in \mathbb{R}$$

Laplacian kernels (relation to ???)

$$k(\mathbf{x}, \mathbf{x}') = \exp(-\gamma \|\mathbf{x} - \mathbf{x}'\|), \quad \gamma > 0 \in \mathbb{R}$$

Sigmoidal kernels (relation to MLPs)

$$k(\mathbf{x}, \mathbf{x}') = g(a \langle \mathbf{x}, \mathbf{x}' \rangle + c)$$

with g a sigmoidal (e.g., logistic, tanh, ...) and particular choices for a, c

Kernel design (I): theoretical issues

Kernel construction

Which **operations** (e.g., products, sums, composition, etc) on kernels produce new kernels? (*closure properties*)

Example:

Consider functions $p : \mathbb{R} \rightarrow \mathbb{R}$.

If k is a kernel, when is $p \circ k$ a kernel?

$$p(x) = \sum_{q=1}^g c_q x^q, c_q \geq 0$$

$$p(z) = (z+1)^q$$

$$k(x, x') = (x^T x + 1)^q$$

Kernel design (I): theoretical issues

Closure properties

- Inner products: finite (sums), infinite countable (series) or infinite uncountable (integrals)
- Scalar operations, sums and direct sums
- Products and tensor products
- Limits of point-wise convergent sequences
- Composition with certain analytic functions
- Normalization

Kernel design (I): theoretical issues

Inner products

1. Let $f_1, \dots, f_n : \mathcal{X} \rightarrow \mathbb{R}$ be a finite collection of functions:

$$k(x, x') = \sum_{i=1}^n f_i(x) \cdot f_i(x')$$

2. Let $\{f_n\}_n$ be a sequence of functions $\mathcal{X} \rightarrow \mathbb{R}$; if the series is convergent:

$$k(x, x') = \sum_{n=1}^{\infty} f_n(x) \cdot f_n(x')$$

3. Let $f : \mathcal{X} \times W \rightarrow \mathbb{R}$ be a parameterized (indexed) set of functions; if the integral is well-defined:

$$k(x, x') = \int_W f(x; w) \cdot f(x'; w) dw$$

Kernel design (I): theoretical issues

Scalar operations, sums and direct sums

Take $k_1, k_2 : \mathcal{X} \times \mathcal{X} \rightarrow \mathbb{R}$ and $k' : \mathcal{Y} \times \mathcal{Y} \rightarrow \mathbb{R}$ kernels

- $a \cdot k_1(x, x') + b$, $a > 0, b \geq 0$

$$\begin{aligned} k_1, k_2 : \mathcal{X} \times \mathcal{X} \rightarrow \mathbb{R} &\Rightarrow k, (x, x') \mapsto k_1(x, x') + k_2(x, x') \quad \text{sum} \\ k_1 : \mathcal{X} \times \mathcal{X} \rightarrow \mathbb{R} &\rightarrow k((x, y), (x', y')) = k_1(x, x') + k_2(y, y') \\ k_2 : \mathcal{Y} \times \mathcal{Y} \rightarrow \mathbb{R} & \end{aligned}$$

direct sum

- k_+ : $\mathcal{X} \times \mathcal{X} \rightarrow \mathbb{R}$ defined as

$$k_+(x, x') = k_1(x, x') + k_2(x, x')$$

- k_{\oplus} : $(\mathcal{X} \times \mathcal{Y}) \times (\mathcal{X} \times \mathcal{Y}) \rightarrow \mathbb{R}$ defined as

$$k_{\oplus}((x, y), (x', y')) = k_1(x, x') + k'(y, y')$$

Kernel design (I): theoretical issues

Products and tensor products

Take $k_1, k_2 : \mathcal{X} \times \mathcal{X} \rightarrow \mathbb{R}$ and $k' : \mathcal{Y} \times \mathcal{Y} \rightarrow \mathbb{R}$ kernels

- $k_{\cdot} : \mathcal{X} \times \mathcal{X} \rightarrow \mathbb{R}$ defined as

$$k_{\cdot}(x, x') = k_1(x, x') \cdot k_2(x, x')$$

- $k_{\odot} : (\mathcal{X} \times \mathcal{Y}) \times (\mathcal{X} \times \mathcal{Y}) \rightarrow \mathbb{R}$ defined as

$$k_{\odot}((x, y), (x', y')) = k_1(x, x') \cdot k'(y, y')$$

Kernel design (I): theoretical issues

Limits of sequences

Let $\{k_n\}_n : \mathcal{X} \times \mathcal{X} \rightarrow \mathbb{R}$ be a sequence of kernels; if, for all $x, x' \in \mathcal{X}$, the limit exists,

then $k_\infty : \mathcal{X} \times \mathcal{X} \rightarrow \mathbb{R}$ defined as

$$k_\infty(x, x') := \lim_{n \rightarrow \infty} k_n(x, x'), \quad \forall x, x' \in \mathcal{X}$$

is a valid kernel.

Kernel design (I): theoretical issues

Composition with analytic functions

Theorem. Let f be a real analytic function with radius of convergence $R > 0$ s.t. all the coefficients in its power series expansion are non-negative. Let $k : \mathcal{X} \times \mathcal{X} \rightarrow \mathbb{R}$ be a kernel fulfilling $|k(x, x')| < R$.

Then $k_f : \mathcal{X} \times \mathcal{X} \rightarrow \mathbb{R}$ given by $k_f(x, x') := f(k(x, x'))$ is a valid kernel.

Example: $f(z) = \exp(z)$

A real function f is *analytic* in an open set $\Omega \subset \mathbb{R}$ iff for every $x_0 \in \Omega$ there is a neighborhood of x_0 for which the Taylor series expansion of f in x_0 coincides with $f(x)$.

Kernel design (I): theoretical issues

Operations in feature space

Norms in feature space:

$$\|\phi(x)\|_{\mathcal{H}} = \sqrt{\langle \phi(x), \phi(x) \rangle_{\mathcal{H}}} = \sqrt{k(x, x)}$$

Norms of linear combinations in feature space:

$$\left\| \sum_i \alpha_i \phi(x_i) \right\|_{\mathcal{H}}^2 = \langle K\alpha, \alpha \rangle = \alpha^T K \alpha$$

$$\langle \alpha_i \phi(x_i), \alpha_i \phi(x_i) \rangle \simeq k(x, x) \alpha_i^2 = \alpha_i^T K(x, x) \alpha_i$$

Kernel design (I): theoretical issues

Operations in feature space

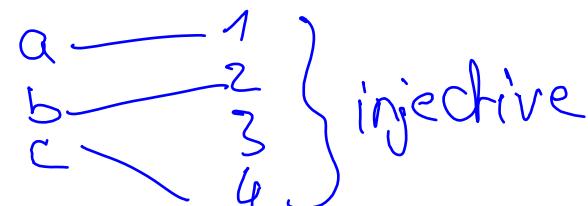
Distances in feature space:

$$\begin{aligned} \|\phi(x) - \phi(x')\|^2 &= \|\phi(x)\|^2 + \|\phi(x')\|^2 - 2\langle \phi(x), \phi(x') \rangle \\ &= \langle \phi(x), \phi(x) \rangle + \langle \phi(x'), \phi(x') \rangle - 2\langle \phi(x), \phi(x') \rangle = k(x, x) + k(x', x') - 2k(x, x') \end{aligned}$$

$$\|\phi(x) - \phi(x')\|_{\mathcal{H}} = \sqrt{\langle \phi(x), \phi(x) \rangle_{\mathcal{H}} + \langle \phi(x'), \phi(x') \rangle_{\mathcal{H}} - 2 \langle \phi(x), \phi(x') \rangle_{\mathcal{H}}}$$

and then $d_{\mathcal{H}}(x, x') := \sqrt{k(x, x) + k(x', x') - 2k(x, x')}$ is an Euclidean metric (distance) when ϕ is **injective** (otherwise it would be a pseudo-metric).

Suppose ϕ not injective & I have a kernel k s.t. $k(x, x) = 1 \quad \forall x, x' \in X$
 $\Rightarrow d_{\mathcal{H}}(x, x') = \sqrt{1 + 1 - 2k(x, x')} = \sqrt{2} \cdot \sqrt{1 - k(x, x')} \propto \sqrt{1 - k(x, x')}$



Kernel design (I): theoretical issues

Normalizing kernels

If k is a kernel, then so is:

$$k_n(x, x') := \frac{k(x, x')}{\sqrt{k(x, x)} \cdot \sqrt{k(x', x')}}$$

Moreover, $|k_n(x, x')| \leq 1$ and $k_n(x, x) = 1$.

The effect is to project each point onto the unit sphere, since

$$1 = k_n(x, x) = \langle \phi_n(x), \phi_n(x) \rangle = \|\phi_n(x)\|^2$$

Kernel design (I): theoretical issues

General linear kernel

Theorem. If $A_{d \times d}$ is a PSD matrix, then the function $k : \mathbb{R}^d \times \mathbb{R}^d \rightarrow \mathbb{R}$ given by $k(x, x') = x^\top A x'$ is a kernel.

Proof. Since A is PSD we can write it in the form $A = BB^\top$. For every $n \in \mathbb{N}$, and every choice $x_1, \dots, x_n \in \mathbb{R}^d$, we form the matrix $\mathbf{K} = (k_{ij})$, where $k_{ij} = k(x_i, x_j) = x_i^\top A x_j$. Then for every $c \in \mathbb{R}^n$:

$$= x_i^\top B B^\top x_j = B^\top x_i \cdot B^\top x_j$$

$$\sum_{i=1}^n \sum_{j=1}^n c_i c_j k_{ij} = \sum_{i=1}^n \sum_{j=1}^n c_i c_j x_i^\top A x_j = \sum_{i=1}^n \sum_{j=1}^n c_i c_j (B^\top x_i)^\top (B^\top x_j)$$

$$= \left\| \sum_{i=1}^n c_i (B^\top x_i) \right\|^2 \geq 0. \quad \text{Note that } \phi(x) = B^\top x$$

Kernel design (I): theoretical issues

Polynomial kernels

1. If k is a kernel and p is a (non-zero) polynomial of degree q with non-negative coefficients, then the function

$$k_p(x, x') := p(k(x, x'))$$

is also a kernel.

2. The special case where k is linear and $p(z) = (az + c)^q, a > 0, c \geq 0 \in \mathbb{R}$ leads to the so-called **polynomial kernel**

Kernel design (I): theoretical issues

Translation invariant and radial kernels

We say that a kernel $k : \mathbb{R}^d \times \mathbb{R}^d \rightarrow \mathbb{R}$ is:

Translation invariant if it has the form $k(x, x') = T(x - x')$, where $T : \mathbb{R}^d \rightarrow \mathbb{R}$ is a differentiable function

Radial if it has the form $k(x, x') = t(\|x - x'\|)$, where $t : [0, \infty) \rightarrow [0, \infty)$ is a differentiable function

Radial kernels fulfill $k(x, x) = t(0)$.

Kernel design (I): theoretical issues

The Gaussian kernel

Consider the function $t(z) = \exp(-\gamma z^2)$, $\gamma > 0$. The resulting radial kernel is known as the **Gaussian RBF kernel**:

$$k(x, x') = \exp(-\gamma \|x - x'\|^2)$$

Many people refer to it simply as “the RBF kernel”

You can also find it as:

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

Kernel design (I): theoretical issues

Using the exponential

1. If k is a kernel and $\gamma > 0$, then the function

$$k(x, x') = \exp(\gamma k(x, x'))$$

is also a kernel.

2. If k is a kernel and $\gamma > 0$, then the function

$$k(x, x') = \exp \left(-\gamma [k(x, x) + k(x', x') - 2k(x, x')] \right)$$

is also a kernel.

Kernel design (I): theoretical issues

Characterization of Kernels

A symmetric function k is called **conditionally positive semi-definite** (CPSD) in \mathcal{X} if for every $n \in \mathbb{N}$, and every choice $x_1, \dots, x_n \in \mathcal{X}$, the matrix $\mathbf{K} = (k_{ij})$, where $k_{ij} = k(x_i, x_j)$, is CPSD.

A real symmetric matrix $A_{n \times n}$ is CPSD if and only if $\forall c \in \mathbb{R}^n$ such that $c^T 1 = 0$, $c^T A c \geq 0$.

new part

It turns out that it suffices for a kernel to be CPSD! Since the class of CPSD kernels is larger than that of PSD kernels:

1. a larger set of kernel functions are usable by kernel machines
2. a larger set of learning algorithms are prone to kernelization

Kernel-Based Learning & Multivariate Modeling

MIRI Master

Lluís A. Belanche

`belanche@cs.upc.edu`

Soft Computing Research Group

Universitat Politècnica de Catalunya

2021-2022

Kernel-Based Learning & Multivariate Modeling

Syllabus

Part 1

Sep 15 Introduction to kernel-based learning

Sep 22 The SVM for classification, regression & novelty detection (I)

Sep 29 The SVM for classification, regression & novelty detection (II)

Oct 06 Kernel design (I): theoretical issues

Oct 13 Kernel design (II): practical issues

Oct 20 Kernelizing ML & statistical algorithms

Oct 27 Advanced topics

Kernel design (II): practical issues

Euclidean space \mathbb{R}^d , but not only ...

- Kernels on real vectors (whole families)
- Kernels on binary vectors (bitstrings = sets)
- General structured kernels:
 - All-subsets kernel
 - Convolution kernels
- Kernels on discrete structures:
 - Tree kernels
 - Graph kernels
- Kernels on distributions (generative kernels):
 - P-kernels
 - Marginalized kernels
- String kernels (text)

... and many others (functional data, categorical data, ...)

Kernel design (II): practical issues

All-subsets kernel

Consider a feature space with one feature for every subset $A \subseteq \{1, \dots, d\}$ of the input variables:

For $x \in \mathbb{R}^d$, feature A is given by $\phi_A(x) = \prod_{i \in A} x_i$ (note $\phi_\emptyset(x) = 1$)

The kernel is defined by the mapping $\phi : x \rightarrow (\phi_A(x))_{A \subseteq \{1, \dots, d\}}$

$$\begin{aligned} k(x, x') &= \langle \phi(x), \phi(x') \rangle = \sum_{A \subseteq \{1, \dots, d\}} \phi_A(x) \phi_A(x') \\ &= \sum_{A \subseteq \{1, \dots, d\}} \prod_{i \in A} x_i x'_i = \prod_{i=1}^d (1 + x_i x'_i) \end{aligned}$$

The last step is obtained by expanding $(1 + x_1 x'_1)(1 + x_2 x'_2) \dots (1 + x_d x'_d)$

Kernel design (II): practical issues

All-subsets kernel

We have the freedom to downplay some features (and thus emphasize others) by introducing weighting factors $w_i \geq 0$ for each feature i :

$$\phi_A(x) = \prod_{i \in A} \sqrt{w_i} x_i$$

therefore

$$k_w(x, x') = \prod_{i=1}^d (1 + w_i x_i x'_i)$$

Kernel design (II): practical issues

Bitstring/Binary variables/Sets

Let $x, x' \in \{0, 1\}^d$, representing absence/presence of a binary trait:

1. The Simple Matching Coefficient (SMC) is the fraction of 1 – 1 matches, and it is a kernel on $\{0, 1\}^d$.
$$\frac{\text{\# matching values}}{\text{\# all values}}$$

Proof. For every $n \in \mathbb{N}$, and every choice $x_1, \dots, x_n \in \{0, 1\}^d$, we form the matrix $K = (k_{ij})$, where $k_{ij} = k(x_i, x_j) = \frac{1}{d}x_i^\top x_j$.

2. The Jaccard Coefficient is the fraction of 1 – 1 matches among the traits present in either data vector, and it is a kernel on $\{0, 1\}^d$.

$$\frac{|A \cap B|}{|A \cup B|}$$

Kernel design (II): practical issues

Bitstring/Binary variables/Sets

- Given two sets $A, B \subset U$, where U is finite, consider

$$k(A, B) = \frac{1}{|U|} \sum_{a \in A} \sum_{b \in B} k_{\text{base}}(a, b)$$

- If k_{base} is the overlap kernel $k(a, b) = \begin{cases} 1 & \text{if } a = b; \\ 0 & \text{otherwise.} \end{cases}$

we get $k(A, B) = \frac{|A \cap B|}{|U|}$, the equivalent of the SMC.

- The equivalent of the Jaccard kernel would be $k(A, B) = \frac{|A \cap B|}{|A \cup B|}$.

Kernel design (II): practical issues

Generative kernels

Given a probability distribution on $\mathcal{X} \times \mathcal{Z}$, we can compare data points by assigning a high value if both have **high conditional probability**:

$$k(x, x') = \sum_{z \in \mathcal{Z}} p(x|z)p(x'|z)P(z) \quad \text{discrete case}$$

$$k(x, x') = \int_{\mathcal{Z}} p(x|z)p(x'|z)p(z) dz \quad \text{continuous case}$$

The feature maps are $(\phi(x))_z = p(x|z)\sqrt{p(z)}$

Idea: $p(x, x', z) = p(x, x'|z)p(z) = p(x|z)p(x'|z)p(z)$

Kernel design (II): practical issues

Generative kernels

Given a probability distribution on $\mathcal{X} \times \mathcal{Z}$, and a kernel on $\mathcal{X} \times \mathcal{Z}$ pairs, we can define:

$$k(\mathbf{x}, \mathbf{x}') = \sum_z \sum_{z'} k((\mathbf{x}, z), (\mathbf{x}', z')) p(z|\mathbf{x}) p(z'|\mathbf{x}')$$

Typical applications of generative kernels are found in **graphical models**:

- \mathcal{X} are the **observed** variables and \mathcal{Z} are the **hidden** (latent) variables
- A kernel for the observed ones is obtained by taking the expectation w.r.t. the hidden ones (*marginalizing* them away)
- Examples: HMMs for sequences or stochastic context-free grammars for RNA sequences –see “Kernel methods in genomics and computational biology” by J.P. Vert

Kernel design (II): practical issues

The Spectrum (aka n -Gram) kernel

- Let Σ be a finite alphabet: an n -Gram is a block of n adjacent characters in Σ

Define $k(x, x') := \sum_{s \in \Sigma^n} |s \in x| \cdot |s \in x'|$ *X is word*

Example: Word aababc in alphabet $\Sigma = \{a, b, c\}$, $n = 2$:

aa	ab	ac	ba	bb	bc	...
1	2	0	1	0	1	...

- While the feature space is large, the feature vectors are sparse; this kernel can be computed in $O(|x| + |x'|)$ time and memory (the actual number of distinct n -Grams in a text is very small)

Kernel design (II): practical issues

Kernels from graphs

- Consider a graph $G = (V, E)$, where the set of vertices (nodes) V are the data points and E is the set of edges. Call $N = |V|$, the number of nodes
- The idea is to compute a (base) matrix $S_{N \times N}$ whose entries are the weights of the edges and consider $S^2 = SS$ (S need not be symmetric)
- Typical use: **connectivity matrix** of G : the (i, j) element of S^2 is the number of paths of length exactly 2 between i and j

Examples:

1. protein-protein interactions
2. people-to-people interactions

In 2, the (i, j) element of S^2 is the number of common friends between data points i and j (it can be thought of as a measure of their similarity)

Kernel design (II): practical issues

Kernels from graphs

Notes:

- The entries of S may be real-valued numbers (e.g., symmetric bounded similarities)
- Higher powers of S measure higher-order similarities
- Only the even powers are guaranteed to be PSD

Consider, for a given $\lambda \in (0, 1)$:

$$\sum_{k=0}^{\infty} \frac{1}{k!} \lambda^k S^k = \exp(\lambda S)$$

1. If S is symmetric, then $S = U \Lambda U^T$ (spectral decomposition), so $S^2 = (U \Lambda U^T)(U \Lambda U^T) = U \Lambda^2 U^T$.
2. In general, we have $S^k = U \Lambda^k U^T$ and therefore:

$$K := \exp(\lambda S) = U \exp(\lambda \Lambda) U^T$$

is an example of a **diffusion** kernel.

Example in a real application domain

Handling missing values in microbiology

- Modern modelling problems are difficult for a number of reasons, including the challenge of dealing with a significant amount of missing information
- **Missing values** almost always represent a serious problem because they force to **preprocess** the dataset and a good deal of **effort** is normally put in this part of the modelling
- In order to process such datasets with kernel methods, an imputation procedure is then deemed a necessary but demanding step

Example in a real application domain

Handling missing values in microbiology

- The study of fecal source pollution in waterbodies is a major problem in ensuring the welfare of human populations
- **Microbial source tracking** (MST) methods attempt to identify the source of contamination, allowing for improved risk analysis and better water management
- The available dataset includes 148 observations about 10 chemical, microbial, and eukaryotic markers of fecal pollution in water
- All variables (except the class variable) are binary, i.e., they signal the presence or absence of a particular marker

Example in a real application domain

Handling missing values in microbiology

Origin	HF183	HF134	CF128	Humito	Pomito	Bomito	ADO	DEN
Human :50	0 :68	0 :81	0 :104	0 :35	0 :83	0 :78	0 :56	0 :80
Cow :26	1 :40	1 :26	1 : 5	1 :79	1 :32	1 :32	1 :59	1 :34
Poultry:31	?:31	?:32	?:30	?:25	?:24	?:29	?:24	?:25
Pig :32								

Summary (counts) table for the full dataset. The first column is the target class. The symbol ? denotes a missing value.

The percentage of missing values is around 19.8 %, and all the predictive variables have percentages between 17 % and 23 %

Example in a real application domain

Handling missing values in microbiology

Let the symbol ? denote a missing element, for which only equality is defined. Let $k : X \times X \rightarrow \mathbb{R}$ be a symmetric kernel in X and P a probability mass function (PMF) in X . Then the function $k^\text{?}(x, y)$ given by

$$k^\text{?}(x, y) = \begin{cases} k(x, y), & \text{if } x, y \neq \text{?}; \\ g(x) = \sum_{y' \in X} P(y') k(x, y'), & \text{if } x \neq \text{?} \text{ and } y = \text{?}; \\ g(y) = \sum_{x' \in X} P(x') k(x', y), & \text{if } x = \text{?} \text{ and } y \neq \text{?}; \\ G = \sum_{x' \in X} P(x') \sum_{y' \in X} P(y') k(x', y'), & \text{if } x = y = \text{?} \end{cases}$$

is a kernel in $X \cup \{\text{?}\}$.

Example in a real application domain

Handling missing values in microbiology

For the particular case of binary variables $x, y \in \{v_1, v_2\}$, a convenient approach is to define the kernel:

$$k_{0/1}(x, y) = \mathbb{I}_{\{x=y\}}$$

where

$$\mathbb{I}_{\{z\}} = \begin{cases} 1 & \text{if } z \text{ is true} \\ 0 & \text{if } z \text{ is false} \end{cases}$$

Example in a real application domain

Handling missing values in microbiology

Consider now $x, y \in \{0, 1\}^d$. When we apply the Theorem to this kernel, we obtain an extended multivariate kernel:

$$\mathcal{K}_1(x, y) = \frac{1}{d} \sum_{i=1}^d \begin{cases} 1 & \text{if } x_i = y_i = 1 ; \\ P_i(x_i), & \text{if } x_i \neq ? \text{ and } y_i = ?; \\ P_i(y_i), & \text{if } x_i = ? \text{ and } y_i \neq ?; \\ (P_i(0))^2 + (P_i(1))^2, & \text{if } x_i = y_i = ?; \\ 0, & \text{otherwise} \end{cases}$$

This kernel is a **generalization** of the classical ***Simple Matching Coefficient***, proposed by Sokal and Michener for numerical taxonomy

Example in a real application domain

Handling missing values in microbiology

Let the symbol $?$ denote a missing element, for which only equality is defined. Let $k : X \times X \rightarrow \mathbb{R}$ be a symmetric kernel in $X = \{0, 1\}^d$. Let $c(x)$ be the set of completions of x . Given two vectors $x, y \in X$, the function

$$\mathcal{K}_2(x, y) = \frac{1}{|c(x)||c(y)|} \sum_{x' \in c(x)} \sum_{y' \in c(y)} k(x', y') \quad (1)$$

is a kernel in $X \cup \{?\}$.

Example in a real application domain

Handling missing values in microbiology

Approach	C	10x10cv	10x10cv for each class			
			Human	Cow	Poultry	Swine
\mathcal{K}_1	2.0	79.3	95.4	64.5	75.2	69.4
\mathcal{K}_2	1.6	78.2	92.6	62.8	71.8	74.2
MI-1	1.0	79.9	92.7	66.4	69.4	80.2
MI-2	1.0	79.0	94.5	57.5	70.8	78.8

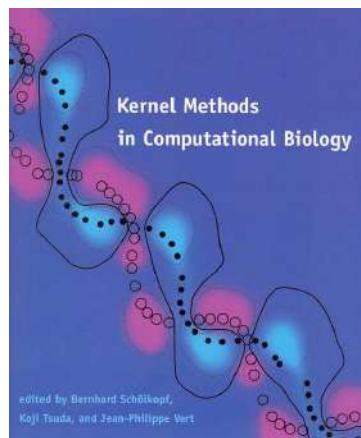
Mean 10x10cv accuracies for the four approaches to handle missing values. Also shown are best cost parameter C and detailed class performance.

(joint work with G. Nebot, T. Aluja and V. Kobayashi)

Kernel design (II): practical issues

More Kernels!

Kernels abound in computational biology and computational chemistry
(e.g., phylogenetic profiles, protein 3D structures)



Example: the prediction of **interacting proteins** to reconstruct an interaction network can be posed as a binary classification problem: given a pair of proteins, do they interact or not?

→ we need kernel between *pairs* of proteins!

Kernel design (II): practical issues

More Kernels!

The available data is about each single protein; it is then natural to derive kernels for **pairs** of proteins k_{pair} from any kernel k for **single** proteins:

$$k_{\text{pair}}((A, B), (C, D)) := k(A, C)k(B, D) + k(A, D)k(B, C)$$

(there is usually no order in a protein pair, so we try both matches)

-
- Using Product Kernels to Predict Protein Interactions. *Advances in Biochemical Engineering/Biotechnology* (110), pp 215-245 (2007)
 - Kernel methods for predicting protein-protein interactions. *Bioinformatics*. 2005

Kernel design (II): practical issues

Conclusions

- The power of kernel methods partly relies in the ability to process virtually any sort of data as soon as a valid kernel is defined
- Importance of designing kernels that do not constitute explicit inner products between objects, and therefore fully exploit the kernel trick
- Possibility of learning the kernel function (or the kernel matrix) from the training data
- Theoretical analyses are needed on the implications of the kernel choice for the success of specific kernel-based methods

Kernel-Based Learning & Multivariate Modeling

MIRI Master

Lluís A. Belanche

`belanche@cs.upc.edu`

Soft Computing Research Group

Universitat Politècnica de Catalunya

2021-2022

Kernel-Based Learning & Multivariate Modeling

Syllabus

Part 1

Sep 15 Introduction to kernel-based learning

Sep 22 The SVM for classification, regression & novelty detection (I)

Sep 29 The SVM for classification, regression & novelty detection (II)

Oct 06 Kernel design (I): theoretical issues

Oct 13 Kernel design (II): practical issues

Oct 20 Kernelizing ML & statistical algorithms

Oct 27 Advanced topics

Kernelizing ML & stats algorithms

[Q] What is “kernelizing”?

[A] do the inner product/distance in feature space:

- If an algorithm is **inner product-based**, the idea is:

1. substitute $\langle x_i, x_j \rangle$ by $\langle \phi(x_i), \phi(x_j) \rangle_{\mathcal{H}}$
2. replace this by $k(x_i, x_j)$

- If an algorithm is (Euclidean) **distance-based**, the idea is:

1. substitute $\|x_i - x_j\|$ by $\|\phi(x_i) - \phi(x_j)\|_{\mathcal{H}}$
2. replace this by $\sqrt{k(x_i, x_i) + k(x_j, x_j) - 2k(x_i, x_j)}$

Kernelizing ML & stats algorithms

A problem with the inner product

- If two data points x_i and x_j are translated as:

$$\begin{cases} x_i \leftarrow x_i - x_0 \\ x_j \leftarrow x_j - x_0 \end{cases}$$

then the inner product $\langle x_i, x_j \rangle$ changes substantially

- This is not suitable for algorithms that should be translation-invariant (e.g. PCA), unless we center the data in feature space

Kernelizing ML & stats algorithms

Conditional positive semi-definiteness

A symmetric function k is called **conditionally positive semi-definite** (CPSD) in \mathcal{X} if for every $n \in \mathbb{N}$, and every choice $x_1, \dots, x_n \in \mathcal{X}$, the matrix $\mathbf{K} = (k_{ij})$, where $k_{ij} = k(x_i, x_j)$ is CPSD.

A real symmetric matrix $A_{n \times n}$ is CPSD if and only if $\forall c \in \mathbb{R}^n$ such that $c^T 1 = 0$, $c^T A c \geq 0$.

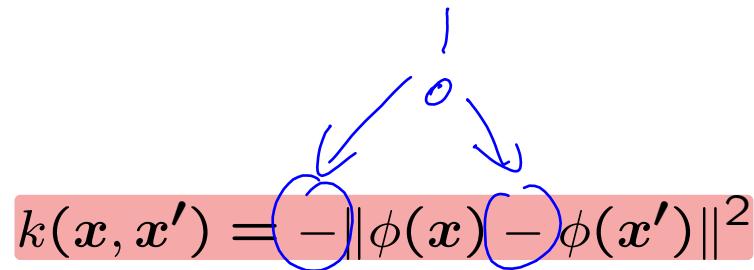
Example: in $\mathcal{X} = \mathbb{R}^d$, $k(x, x') = -\|x - x'\|^2$ is CPSD (but not PSD).

Exercise: For which values of β is $-\|x - x'\|^\beta$ CPSD? ¶

Kernelizing ML & stats algorithms

The kernel trick for distances

Theorem. Let k be a CPSD kernel on \mathcal{X} , satisfying $k(x, x) = 0$ for all $x \in \mathcal{X}$. Then there exists a Hilbert space \mathcal{H} and a mapping $\phi : \mathcal{X} \rightarrow \mathcal{H}$ such that

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


It turns out that it suffices for a kernel to be CPSD! Since the class of CPSD kernels is larger than that of PSD kernels, a larger set of learning algorithms are prone to kernelization.

Kernelizing ML & stats algorithms

Example: k -nearest neighbours

(kNN classifies new examples by finding the k closest examples in the sample and taking a majority vote)

$$\begin{aligned} k(x, x) + k(x', x') - 2k(x, x') = \\ \langle \phi(x), \phi(x) \rangle + \langle \phi(x'), \phi(x') \rangle - 2 \langle \phi(x), \phi(x') \rangle = \\ \|\phi(x)\|^2 + \|\phi(x')\|^2 - 2 \langle \phi(x), \phi(x') \rangle = \\ \|\phi(x) - \phi(x')\|^2 =: d_{\mathcal{H}}^2(x, x') \end{aligned}$$

- This (square) Euclidean distance in *feature space* can be calculated using 3 calls to the kernel function (or 1 if k is normalized), for k PSD
- Note that $\sqrt{-k(x, x')}$ is also a Euclidean distance, for k CPSD

Kernelizing ML & stats algorithms

Standard PCA

- We are given a data set $X = \{x_1, \dots, x_n\}$, $x_i \in \mathbb{R}^d$ for $i = 1, \dots, n$ which is centered around the origin, i.e. $\sum_{i=1}^n x_i = 0$
- The sample covariance matrix C of the data is defined as

$$C = \frac{1}{n} \sum_{i=1}^n x_i x_i^\top$$

Kernelizing ML & stats algorithms

Standard PCA

What are the PCs?

reduce dimensions -> simplify data

1. The goal of PCA is to replace the original axes with the principal components (PC) of the data
2. We choose the first PC as a projection direction such that the projections of the data onto it have maximum variance
3. All subsequent components can be defined as orthogonal projection directions with the next largest variance

Kernelizing ML & stats algorithms

Kernel PCA from Standard PCA

- If we could define our principal components to be arbitrary **manifolds** we could get **higher variance** and **better separability**
- If we first (non-linearly) **map our data onto a higher-dimensional space** where the data falls neatly onto some **hyperplane** we can perform Standard PCA in that space
- These PCs will map back onto the arbitrary manifolds that we need in our lower-dimensional space

Kernelizing ML & stats algorithms

Kernel PCA from Standard PCA

Kernel PCA allows us to perform PCA in this higher dimension using the kernel trick, doing all our calculations in a lower dimension.

Recall the idea of mapping input data into some Hilbert space (called the *feature space*) via a non-linear mapping $\phi : \mathcal{X} \rightarrow \mathcal{H}$:

→ the new sample covariance matrix C of the data is given by

$$C = \frac{1}{n} \sum_{i=1}^n \phi(\mathbf{x}_i) \phi(\mathbf{x}_i)^\top$$

Kernelizing ML & stats algorithms

Kernel PCA from Standard PCA

To find the first PC, again we need to solve $Cv = \lambda v$. Therefore:

$$\frac{1}{n} \sum_{i=1}^n \phi(x_i) (\phi(x_i)^\top v) = \lambda v$$

Since $\lambda \neq 0$, v must be in the span of the set of vectors $\phi(x_i)$, i.e. it can be written as some linear combination thereof:

→ there must exist some α_i such that $v = \sum_{i=1}^n \alpha_i \phi(x_i)$.

Kernelizing ML & stats algorithms

Kernel PCA from Standard PCA

$$\lambda \sum_{j=1}^n \alpha_j \phi(x_j) = \frac{1}{n} \sum_{i=1}^n \phi(x_i) \phi(x_i)^\top \sum_{j=1}^n \alpha_j \phi(x_j)$$

switch to our alternative notation ...

$$\lambda \sum_{j=1}^n \alpha_j \phi(x_j) = \frac{1}{n} \sum_{i=1}^n \phi(x_i) \left\langle \phi(x_i), \sum_{j=1}^n \alpha_j \phi(x_j) \right\rangle$$

Kernelizing ML & stats algorithms

Kernel PCA from Standard PCA

and now a trick ... introduce an inner product with $\phi(\mathbf{x}_l)$ for an arbitrary l on both sides ($1 \leq l \leq n$):

$$\begin{aligned} & \left\langle \lambda \sum_{j=1}^n \alpha_j \phi(\mathbf{x}_j), \phi(\mathbf{x}_l) \right\rangle = \lambda \sum_{j=1}^n \alpha_j \langle \phi(\mathbf{x}_j), \phi(\mathbf{x}_l) \rangle \\ & \left\langle \frac{1}{n} \sum_{i=1}^n \phi(\mathbf{x}_i) \left\langle \phi(\mathbf{x}_i), \sum_{j=1}^n \alpha_j \phi(\mathbf{x}_j) \right\rangle, \phi(\mathbf{x}_l) \right\rangle = \left\langle \frac{1}{n} \sum_{i=1}^n \phi(\mathbf{x}_i) \sum_{j=1}^n \alpha_j \langle \phi(\mathbf{x}_i), \phi(\mathbf{x}_j) \rangle, \phi(\mathbf{x}_l) \right\rangle \\ & = \left\langle \frac{1}{n} \sum_{j=1}^n \alpha_j \sum_{i=1}^n \phi(\mathbf{x}_i) \langle \phi(\mathbf{x}_i), \phi(\mathbf{x}_j) \rangle, \phi(\mathbf{x}_l) \right\rangle = \frac{1}{n} \sum_{j=1}^n \alpha_j \sum_{i=1}^n \langle \phi(\mathbf{x}_i), \phi(\mathbf{x}_l) \rangle \langle \phi(\mathbf{x}_i), \phi(\mathbf{x}_j) \rangle \end{aligned}$$

siehe slide vorher

Kernelizing ML & stats algorithms

Kernel PCA from Standard PCA

Now we can introduce the kernel; since

$$k(\mathbf{x}, \mathbf{x}') = \phi(\mathbf{x})^\top \phi(\mathbf{x}') = \langle \phi(\mathbf{x}), \phi(\mathbf{x}') \rangle_{\mathcal{H}}, \quad \mathbf{x}, \mathbf{x}' \in \mathcal{X}$$

$$\lambda \sum_{j=1}^n \alpha_j \langle \phi(\mathbf{x}_j), \phi(\mathbf{x}_l) \rangle = \frac{1}{n} \sum_{j=1}^n \alpha_j \sum_{i=1}^n \langle \phi(\mathbf{x}_i), \phi(\mathbf{x}_l) \rangle \langle \phi(\mathbf{x}_i), \phi(\mathbf{x}_j) \rangle$$

is rewritten as:

$$\lambda \sum_{j=1}^n \alpha_j k(\mathbf{x}_j, \mathbf{x}_l) = \frac{1}{n} \sum_{j=1}^n \alpha_j \sum_{i=1}^n k(\mathbf{x}_i, \mathbf{x}_l) k(\mathbf{x}_i, \mathbf{x}_j)$$

Kernelizing ML & stats algorithms

Kernel PCA from Standard PCA

Now let $\alpha = (\alpha_1, \dots, \alpha_n)^\top$ and $\mathbf{K} = (k_{ij})$, where $k_{ij} = k(x_i, x_j)$:

$$\lambda \mathbf{K}\alpha = \frac{1}{n} \mathbf{K}^2 \alpha$$

Again an eigenvalue problem! put $p := \mathbf{K}\alpha$ and rewrite:

$$\frac{1}{n} \mathbf{K}p = \lambda p$$

eigenvector
↓
eigenvalue

where p, λ are the eigenvectors and eigenvalues of $\frac{1}{n} \mathbf{K}$ (we get rid of α)

$\frac{1}{n} \mathbf{K}$ has the same eigenvectors/values as \mathbf{K} , with eigenvalues scaled by n

Kernelizing ML & stats algorithms

Kernel PCA from Standard PCA

- In Standard PCA, the matrix XX^T grows with the **dimension** of the data points d
- In Kernel PCA, the matrix K grows with the **number** of data points n ; in consequence, we get n non-linear PCs
- When the kernel function is standard dot product, Kernel PCA solution reduces to Standard PCA

Kernelizing ML & stats algorithms

Kernel PCA from Standard PCA

- Centered data is required to perform an effective PCA
- Even though the $\{x_i\}$ were centered, the $\{\phi(x_i)\}$ are *not* guaranteed to be centered in feature space!
- We have to “center” the kernel matrix before doing Kernel PCA:

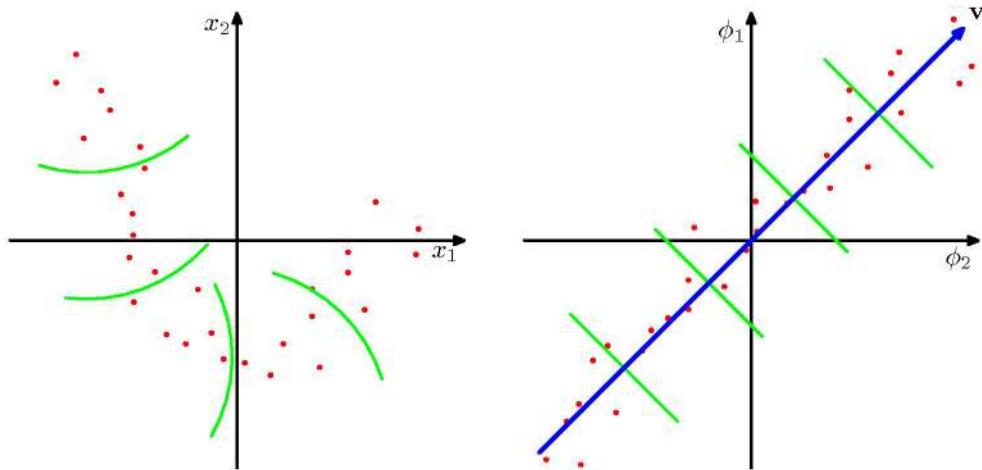
$$K := K - \frac{1}{n} \mathbf{1} \mathbf{K} - \frac{1}{n} \mathbf{K} \mathbf{1} + \frac{1}{n^2} \mathbf{1} \mathbf{K} \mathbf{1}$$

where $\mathbf{1}$ is a $n \times n$ matrix of ones

only here for
correct dims

Kernelizing ML & stats algorithms

An illustration of KPCA



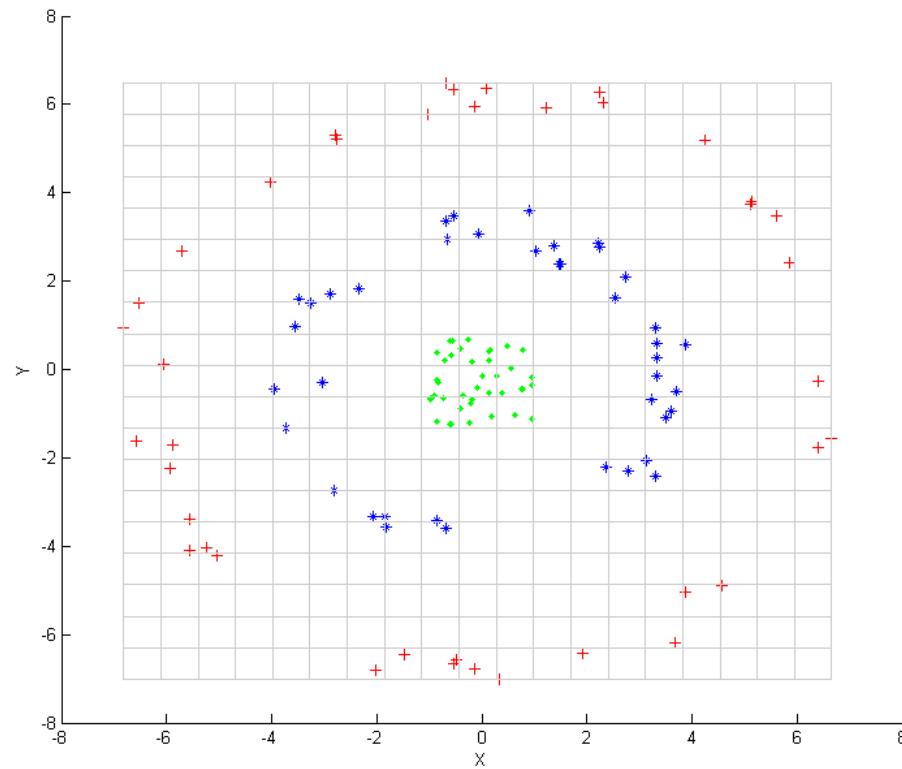
Left (input space): Green lines are desired non-linear projections

Right (feature space): Green lines are linear projections onto the first non-linear PC

(from *Pattern Recognition and Machine Learning*, C. M. Bishop, Springer, 2006)

Kernelizing ML & stats algorithms

Example (I)

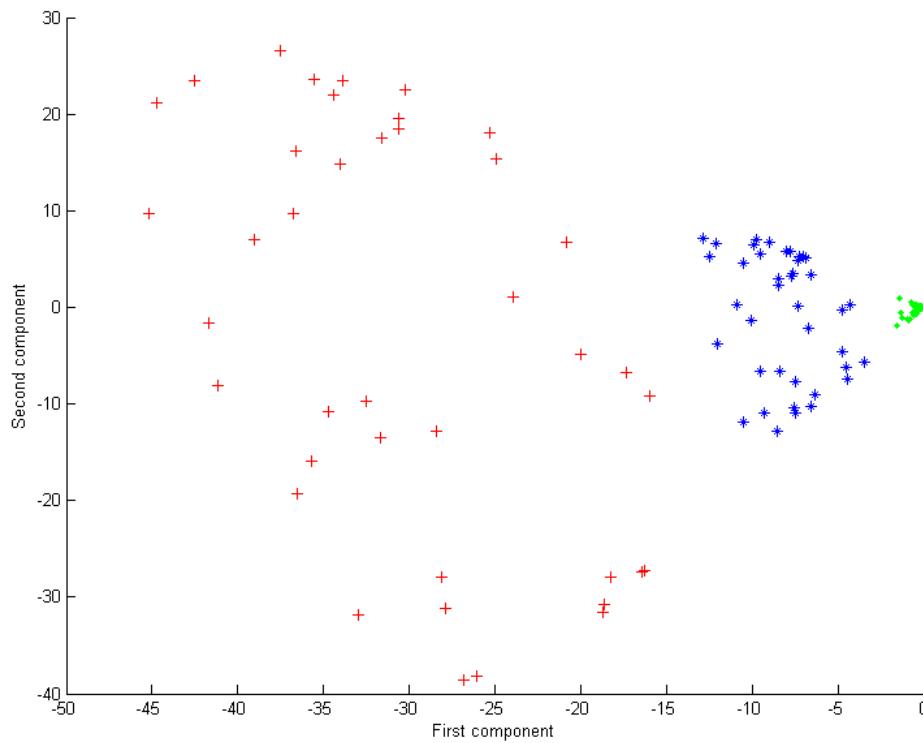


Three concentric clouds of points

(Source: Wikipedia, Petter Strandmark)

Kernelizing ML & stats algorithms

Example (II)



$$k(x, x') = (x^\top x' + 1)^2$$

(the color of the points is not part of the algorithm; (Wikipedia, P. Strandmark)

Kernelizing ML & stats algorithms

Bad and good news

In Standard PCA, it is common to:

1. use the eigenvalues to rank the eigenvectors based on how much of the data variation is captured by each PC
2. estimate a “convenient” number of PCs and perform dimensionality reduction by projecting the original data onto the reduced set of eigenvectors

We **cannot compute the kPCs** (since they reside in the high-dimensional feature space), only the projections of our data onto the kPCs; however, we can **perform new projections on test data** (as in Standard PCA).

Kernelizing ML & stats algorithms

Summary of Standard PCA

1. We are given a data set of d -dimensional vectors $X = \{x_i\}$ for $i = 1, \dots, n$ which we **center** around the origin, as $x_i := x_i - \frac{1}{n} \sum_{j=1}^n x_j$
2. Compute the **sample covariance matrix** of the data as $C = \frac{1}{n} \sum_{i=1}^n x_i x_i^\top$
3. Compute the **eigenvalues** and **eigenvectors** (λ_i, v_i) of C ; **decide** a number $l < d$. Let $\hat{x}_i = Vx_i$, where $V = [v_1; \dots; v_l]$
4. **Return** the l -dimensional data sample $\hat{X} = \{\hat{x}_i\}$ for $i = 1, \dots, n$

Kernelizing ML & stats algorithms

Summary of Kernel PCA

1. We are given a data set of objects $X = \{x_i\}$ for $i = 1, \dots, n$
2. Choose a **kernel** function k and compute the **kernel matrix K**
3. **Center** the kernel matrix as: $K := K - \frac{1}{n} \mathbf{1} \mathbf{K} - \frac{1}{n} \mathbf{K} \mathbf{1} + \frac{1}{n^2} \mathbf{1} \mathbf{K} \mathbf{1}$
4. Compute the **eigenvalues** and **eigenvectors** (λ_i, v_i) of K ; **decide** a number $l < n$. Set $\alpha_j = v_j / \sqrt{\lambda_j}$ and let $\hat{x}_i = \left(\sum_{k=1}^n \alpha_{jk} k(x_i, x_k) \right)_{j=1}^l$
5. **Return** the l -dimensional data sample $\hat{X} = \{\hat{x}_i\}$ for $i = 1, \dots, n$

Kernelizing ML & stats algorithms

Key aspects of kernel methods

Kernel-based methods consist of two ingredients:

1. The **kernel function** (this is non-trivial)

2. The **algorithm** taking kernels as input

- Data items are embedded into a **vector space where linear relations** are computed
- The coordinates of these images are not needed (and are usually unknown): only their **pairwise inner products** matter
- These inner products can sometimes be **computed** efficiently and implicitly in the input space (**kernel function**)
- This requires expressing a problem wherein the data appear in the **form of inner products or distances**
- The **solution vector** is expressed as a **linear combination of kernel evaluations** centered at the data

Kernelizing ML & stats algorithms

Work so far ...

Many (classical and new) learning algorithms can be “kernelized”:

- Support Vector Machine (SVM) and Relevance Vector Machine (SVM)
- Fisher Discriminant Analysis (kFDA), Principal Components Analysis (kPCA) and Canonical Correlation Analysis (kCCA), Independent Component Analysis (kICA)
- Kernel (regularized) linear and logistic regression
- kernel kNN (k Nearest Neighbours)
- Clustering (Spectral Clustering, Kernel k-means)
- PLS, Parzen Windows, Vector Quantization, ...
- Statistical indexes (e.g. Kendall's tau correlation coefficient to measure similarity between permutations)