

# Learning with kernels

60' Tutorial

Justin Dauwels

dauwels@isi.ee.ethz.ch

Signal and Information Processing Laboratory  
Department of Information Technology and Electrical Engineering  
ETH Zurich, Switzerland

# Overview

- Introduction
- From a simple kernel algorithm ...
- ... to general kernel algorithms
- Key aspects of kernels
- Kernels from graphical models
- Support vector machines
- Applications
- Kernel machines: Pros and Cons
- Conclusions
- Where to look for more information?

# Introduction

## History of learning

- **1960s**: Algorithms for detecting **linear** patterns  
e.g., the perceptron, PCA, LMS, RLS
- **mid 1980s**: “**nonlinear** revolution”  
e.g., neural networks, decision trees, graphical models
- **mid 1990s**: **Kernel machines**  
e.g., support vector machines, kernel PCA/LMS/RLS

# Introduction (2)

## Tasks solved by kernel machines

- Density estimation
- Clustering
- Compression
- Regression
- **Classification**

## Input data

- Vectors in  $\mathbb{R}^n$
- Strings (e.g., DNA-sequences)
- Documents (e.g., webpages)
- Graphical models (e.g., HMMs)
- ...

# Introduction (3)

Kernel methods consist of **two** parts

- **Mapping** of the data into suitable **high-dimensional dot-product** space (“feature space”)
- **Learning algorithm** (based on the dot product) designed to discover **linear** patterns in that space

**Good** idea, since

- Increasing dimensionality makes problem often **easier**
- Detection of **linear** patterns is **well-understood**
- **Kernel trick** for computing dot product in feature space

# Intermezzo: Dot product spaces

Let  $K \triangleq \mathbb{R}$  or  $K \triangleq \mathbb{C}$ . A **linear space** (vector space) over  $K$  is a set  $V$  with the following properties:

- An operator “+”:  $V \times V \rightarrow V$  is defined such that  $V, +$  is a **commutative group**
- A “**multiplication**”  $K \times V \rightarrow V$  is defined with properties
  - $a(u + v) = au + av$
  - $(a + b)v = av + bv$
  - $(ab)v = a(bv)$
  - $\mathbf{1} \cdot v = v$

# Intermezzo: Dot product spaces (2)

Let  $V$  be a linear space over  $K$ . An **inner product** (dot product, scalar product) is a function  $V \times V \rightarrow K : (v, w) \rightarrow \langle v, w \rangle$ , such that

- $\langle u + v, w \rangle = \langle u, w \rangle + \langle v, w \rangle$
- $\langle v, w \rangle = \overline{\langle w, v \rangle}$
- $\langle av, w \rangle = a\langle v, w \rangle$
- $\langle v, v \rangle \in \mathbb{R}$  and  $\langle v, v \rangle \geq 0$
- $\langle v, v \rangle = 0 \iff v = \mathbf{0}$

The **norm** of a vector  $v \in V$  is defined as  $\|v\| \triangleq \sqrt{\langle v, v \rangle}$ .

# Introduction (4)

## Supervised binary classification

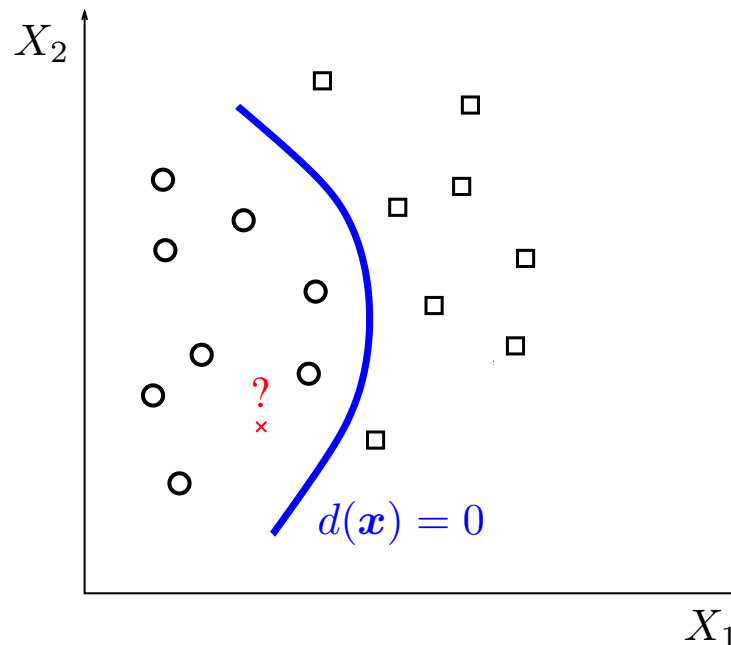
Set  $S = \{(\mathbf{x}_1, y_1), \dots, (\mathbf{x}_n, y_n)\}$  of points  $\mathbf{x}_i \in \mathbb{R}^n$  with labels  $y_i \in \{-1, 1\}$ .

Find **prediction function**  $g(\mathbf{x}) = \text{sign}(d(\mathbf{x}))$  such that  $P(g(\mathbf{x}) \neq y)$  is small.

$d(\mathbf{x})$  is **discriminant function**,  $D = \{\mathbf{x} \in \mathbb{R}^n : d(\mathbf{x}) = 0\}$  is **decision boundary**.

MAP-estimator:  $g_{\text{MAP}}(\mathbf{x}) = \text{sign}[p(y = 1|\mathbf{x}) - p(y = -1|\mathbf{x})]$

However:  $p(\mathbf{x}, y)$  **unknown!**

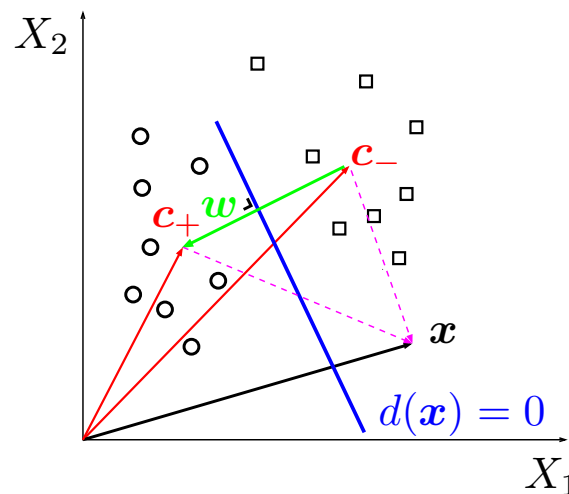


# A simple kernel-based classification algorithm

**Minimum distance classifier:** assign pattern  $\mathbf{x}$  to class with **nearest mean**

$$\begin{aligned}g(\mathbf{x}) &= \text{sign}(\|\mathbf{x} - \mathbf{c}_-\| - \|\mathbf{x} - \mathbf{c}_+\|) \\&= \text{sign}\left(\frac{1}{n_+} \sum_{i|y_i=+1} \langle \mathbf{x}, \mathbf{x}_i \rangle - \frac{1}{n_-} \sum_{i|y_i=-1} \langle \mathbf{x}, \mathbf{x}_i \rangle + b\right) \\&= \text{sign}(\langle \mathbf{w}, \mathbf{x} \rangle + b)\end{aligned}$$

$$\mathbf{c}_\pm = \frac{1}{n_\pm} \sum_{i|y_i=\pm 1} \phi(\mathbf{x}_i) \quad b = \frac{1}{2}(\|\mathbf{c}_-\|^2 - \|\mathbf{c}_+\|^2) \quad \mathbf{w} = \mathbf{c}_+ - \mathbf{c}_-$$



# General kernel-based classification algorithms

## Minimum distance classifier

$$g(\mathbf{x}) = \text{sign} \left( \frac{1}{n_+} \sum_{i|y_i=+1} \langle \mathbf{x}, \mathbf{x}_i \rangle - \frac{1}{n_-} \sum_{i|y_i=-1} \langle \mathbf{x}, \mathbf{x}_i \rangle + b \right)$$

## General kernel classifier (“support vector machine”)

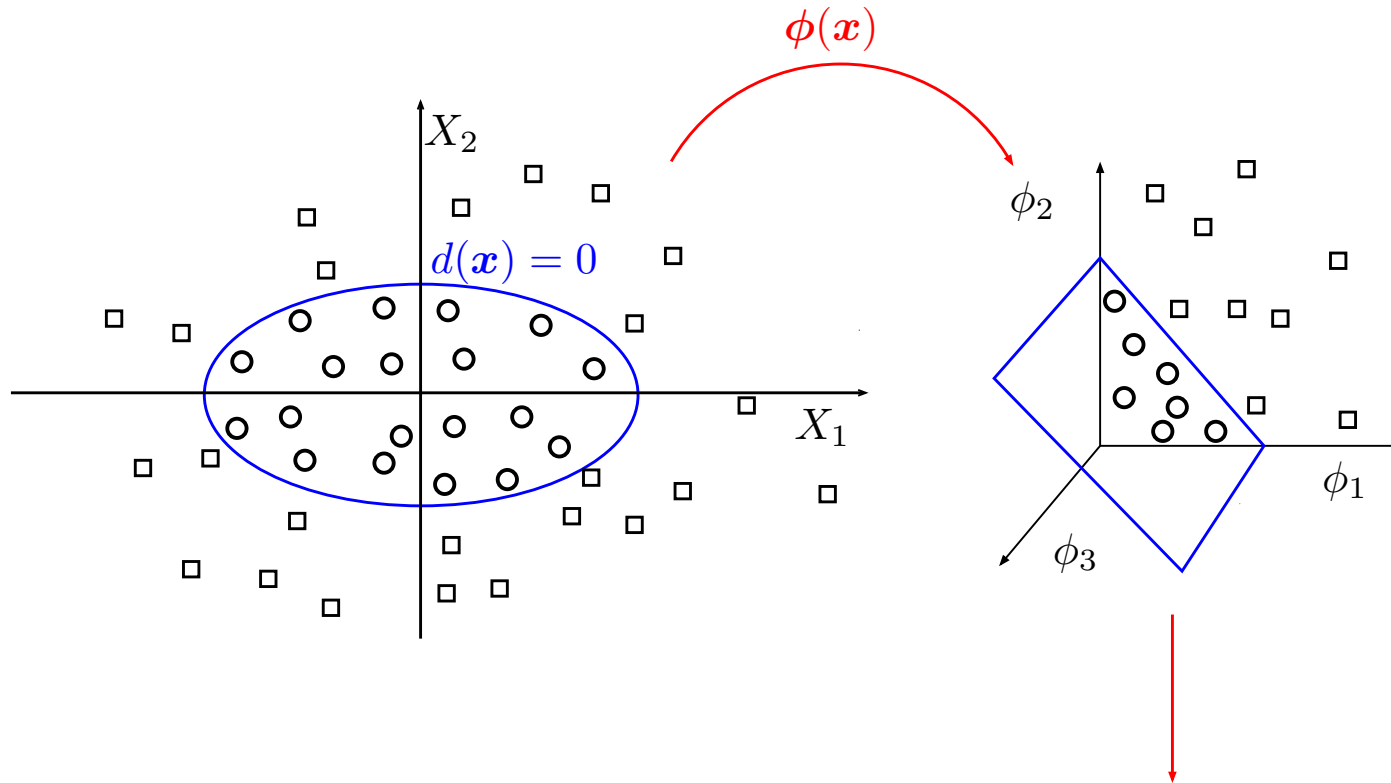
$$g(\mathbf{x}) = \text{sign} \left( \sum_{i=1}^n \beta_i \kappa(\mathbf{x}, \mathbf{x}_i) + b \right)$$

Remarks:

- **General** kernels  $\kappa(\cdot, \cdot)$
- Coefficients  $\beta_i$  **not uniform** within classes
- Point  $\mathbf{x}$  classified by **comparing** to all **input patterns**  $\mathbf{x}_i$  with **nonzero**  $\beta_i$  (“support vectors”)
- **Non-parametric** classifiers, kernels centered on input patterns

# Mapping

Step 1



Step 2

Algorithms in feature space based on dot product

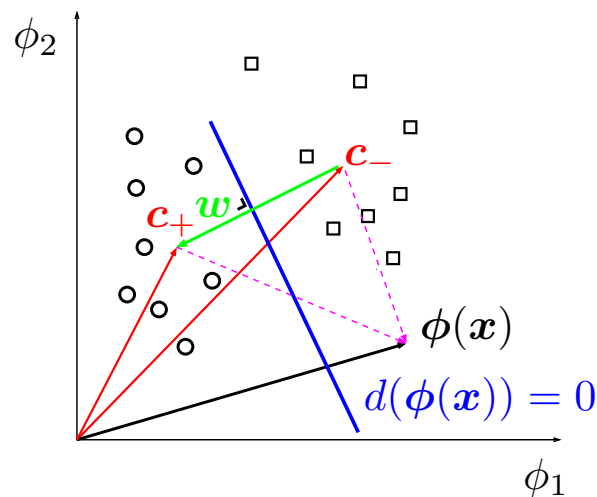
$$\phi : \mathbf{x} = (x_1, x_2) \mapsto \phi(\mathbf{x}) = (x_1^2, x_2^2, \sqrt{2}x_1x_2) \in \mathbb{R}^3$$

# A simple kernel-based classification algorithm

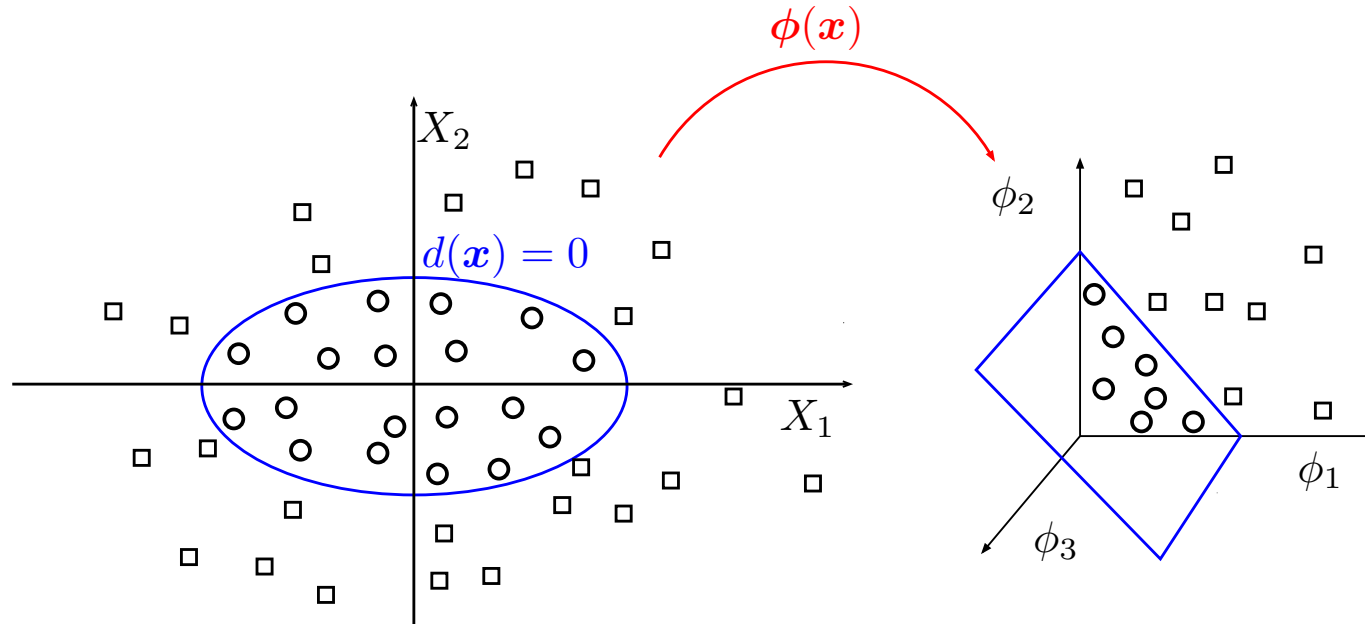
## Minimum distance classifier in *feature space*

$$\begin{aligned}g(\mathbf{x}) &= \text{sign}(\|\phi(\mathbf{x}) - \mathbf{c}_-\| - \|\phi(\mathbf{x}) - \mathbf{c}_+\|) \\&= \text{sign}\left(\frac{1}{n_+} \sum_{i|y_i=+1} \langle \phi(\mathbf{x}), \phi(\mathbf{x}_i) \rangle - \frac{1}{n_-} \sum_{i|y_i=-1} \langle \phi(\mathbf{x}), \phi(\mathbf{x}_i) \rangle + b\right) \\&= \text{sign}(\langle \mathbf{w}, \phi(\mathbf{x}) \rangle + b)\end{aligned}$$

$$\mathbf{c}_\pm = \frac{1}{n_\pm} \sum_{i|y_i=\pm 1} \phi(\mathbf{x}_i) \quad b = \frac{1}{2}(\|\mathbf{c}_-\|^2 - \|\mathbf{c}_+\|^2) \quad \mathbf{w} = \mathbf{c}_+ - \mathbf{c}_-$$



# Kernel trick



$$\phi : \mathbf{x} = (x_1, x_2) \mapsto \phi(\mathbf{x}) = (x_1^2, x_2^2, \sqrt{2}x_1x_2) \in \mathbb{R}^3$$

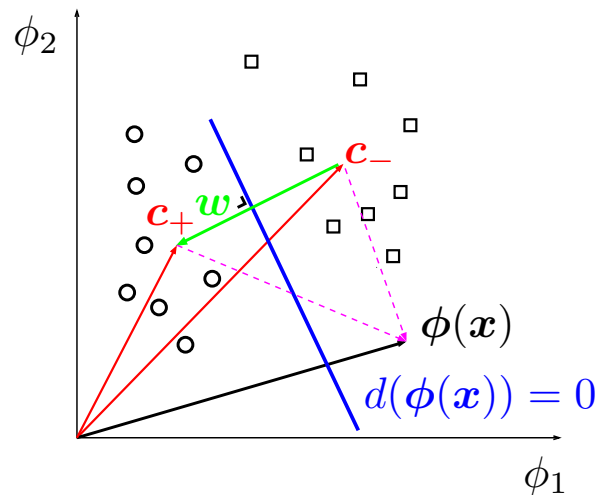
$$\begin{aligned} \langle \phi(\mathbf{x}), \phi(\mathbf{z}) \rangle &= \langle (x_1^2, x_2^2, \sqrt{2}x_1x_2), (z_1^2, z_2^2, \sqrt{2}z_1z_2) \rangle \\ &= x_1^2z_1^2 + x_2^2z_2^2 + 2x_1x_2z_1z_2 \\ &= (x_1z_1 + x_2z_2)^2 \\ &= \langle \mathbf{x}, \mathbf{z} \rangle^2 = \kappa(\mathbf{x}, \mathbf{z}) \end{aligned}$$

# A simple kernel-based classification algorithm

## Minimum distance classifier in *feature space*

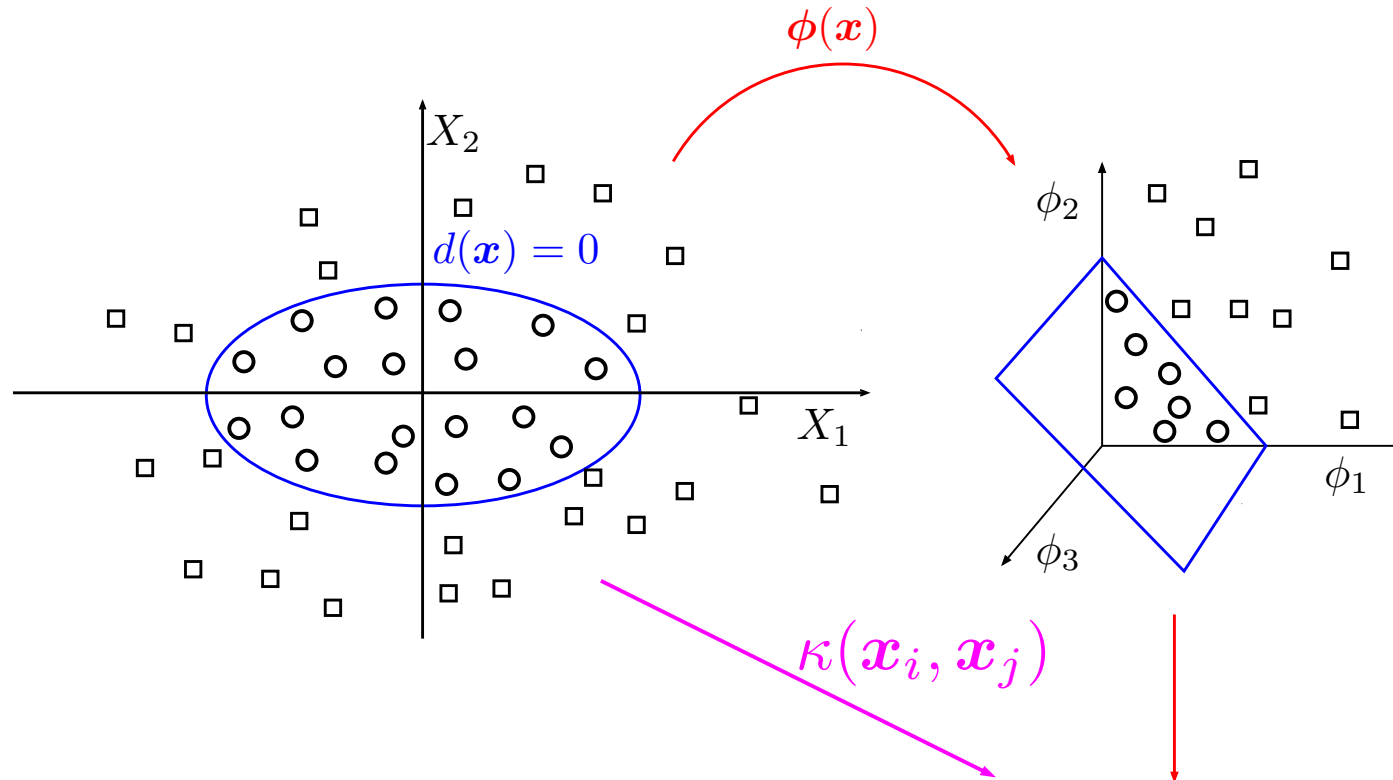
$$\begin{aligned}g(\mathbf{x}) &= \text{sign}(\|\phi(\mathbf{x}) - \mathbf{c}_-\| - \|\phi(\mathbf{x}) - \mathbf{c}_+\|) \\&= \text{sign}\left(\frac{1}{n_+} \sum_{i|y_i=+1} \kappa(\mathbf{x}, \mathbf{x}_i) - \frac{1}{n_-} \sum_{i|y_i=-1} \kappa(\mathbf{x}, \mathbf{x}_i) + b\right) \\&= \text{sign}(\langle \mathbf{w}, \phi(\mathbf{x}) \rangle + b)\end{aligned}$$

$$\mathbf{c}_\pm = \frac{1}{n_\pm} \sum_{i|y_i=\pm 1} \phi(\mathbf{x}_i) \quad b = \frac{1}{2}(\|\mathbf{c}_-\|^2 - \|\mathbf{c}_+\|^2) \quad \mathbf{w} = \mathbf{c}_+ - \mathbf{c}_-$$



# Kernel trick (2)

Step 1



Step 2

Algorithms in feature space based on dot product

$$\phi : \mathbf{x} = (x_1, x_2) \mapsto \phi(\mathbf{x}) = (x_1^2, x_2^2, \sqrt{2}x_1x_2) \in \mathbb{R}^3$$
$$\kappa(\mathbf{x}, \mathbf{z}) = \langle \mathbf{x}, \mathbf{z} \rangle^2$$

# General kernel-based classification algorithms

Minimum distance classifier in input space

$$g(\mathbf{x}) = \text{sign} \left( \frac{1}{n_+} \sum_{i|y_i=+1} \langle \mathbf{x}, \mathbf{x}_i \rangle - \frac{1}{n_-} \sum_{i|y_i=-1} \langle \mathbf{x}, \mathbf{x}_i \rangle + b \right)$$

Minimum distance classifier in feature space

$$g(\mathbf{x}) = \text{sign} \left( \frac{1}{n_+} \sum_{i|y_i=+1} \kappa(\mathbf{x}, \mathbf{x}_i) - \frac{1}{n_-} \sum_{i|y_i=-1} \kappa(\mathbf{x}, \mathbf{x}_i) + b \right)$$

Support vector machine

$$g(\mathbf{x}) = \text{sign} \left( \sum_{i=1}^n \beta_i \kappa(\mathbf{x}, \mathbf{x}_i) + b \right)$$

- General kernels  $\kappa(., .)$
- Point  $\mathbf{x}$  classified by comparing to all input patterns  $\mathbf{x}_i$  with nonzero  $\beta_i$  (“support vectors”)

# Key aspects of kernels

## Definitions

A **kernel** is a function  $\kappa : \mathcal{X} \times \mathcal{X} \rightarrow \mathcal{X}$  that for all  $\mathbf{x}, \mathbf{z} \in \mathcal{X}$  satisfies

$\kappa(\mathbf{x}, \mathbf{z}) = \langle \phi(\mathbf{x}), \phi(\mathbf{z}) \rangle$ , where  $\phi$  is a mapping from  $\mathcal{X}$  to an inner product space  $F$ ,  $\phi : \mathbf{x} \mapsto \phi(\mathbf{x}) \in F$ .

e.g.,  $\kappa(\mathbf{x}, \mathbf{z}) = \langle \mathbf{x}, \mathbf{z} \rangle^2$

Given a set  $S = \{\mathbf{x}_1, \dots, \mathbf{x}_l\}$ , a **Gram matrix** is defined as the  $l \times l$  matrix  $\mathbf{G}$  whose entries are  $\mathbf{G}_{ij} = \langle \mathbf{x}_i, \mathbf{x}_j \rangle$ .

Given a set  $S = \{\mathbf{x}_1, \dots, \mathbf{x}_l\}$ , the **kernel matrix** associated to kernel  $\kappa$  is defined as the  $l \times l$  matrix  $\mathbf{K}$  whose entries are  $\mathbf{K}_{ij} = \langle \phi(\mathbf{x}_i), \phi(\mathbf{x}_j) \rangle = \kappa(\mathbf{x}_i, \mathbf{x}_j)$ .

## Properties of Gram and kernel matrices

Gram matrices are **symmetric** and **positive semi-definite**, i.e.,  $\forall \mathbf{x} \in \mathbb{R}^l : \mathbf{x}^T \mathbf{G} \mathbf{x} \geq 0$ .

# Key aspects of kernels

## Definition

A function  $\kappa : \mathcal{X} \times \mathcal{X} \rightarrow \mathbb{R}$  satisfies the **finitely positive property** if it is a **symmetric** function for which the **matrices**  $\mathbf{K}_{ij} = \kappa(\mathbf{x}_i, \mathbf{x}_j)$  formed by restriction to any finite subset of the space  $\mathcal{X}$  are **positive semi-definite**.

Remarks:

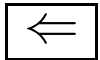
- Every **kernel** satisfies the **finitely positive property**.
- But is **every function** that satisfies this condition a kernel?

# Key aspects of kernels (2)

## Theorem (Generating an inner product space from a kernel)

A function  $\kappa : \mathcal{X} \times \mathcal{X} \rightarrow \mathbb{R}$  which is either continuous or has a finite domain, is a **kernel**, i.e., it can be decomposed as  $\kappa(\mathbf{x}, \mathbf{z}) = \langle \phi(\mathbf{x}), \phi(\mathbf{z}) \rangle$ , where  $\phi$  is a mapping into an inner product space  $F$ , **if and only** if it satisfies the **finitely positive semi-definite property**.

## Sketch of the proof



Mapping  $\phi : \mathbf{x} \in \mathcal{X} \rightarrow \phi(\mathbf{x}) = \kappa(\mathbf{x}, \cdot) \in F_k$

Linear space  $F = \{ \sum_{i=1}^l \alpha_i \kappa(\mathbf{x}_i, \cdot) : l \in \mathbb{N}, \mathbf{x}_i \in \mathcal{X}, \alpha_i \in \mathbb{R} \}$

Inner product:

$$f(\mathbf{x}) = \sum_{i=1}^l \alpha_i \kappa(\mathbf{x}_i, \mathbf{x}) \text{ and } g(\mathbf{x}) = \sum_{i=1}^n \beta_i \kappa(\mathbf{z}_i, \mathbf{x})$$

$$\langle f, g \rangle = \sum_{i=1}^l \sum_{j=1}^n \alpha_i \beta_j \kappa(\mathbf{x}_i, \mathbf{z}_j) = \sum_{i=1}^l \alpha_i g(\mathbf{x}_i) = \sum_{j=1}^n \beta_j f(\mathbf{z}_j)$$

$$\|f\|^2 = \sum_{i=1}^l \sum_{j=1}^l \alpha_i \alpha_j \kappa(\mathbf{x}_i, \mathbf{x}_j) = \boldsymbol{\alpha}' \mathbf{K} \boldsymbol{\alpha} \geq 0 \quad \forall f \in F$$

$$\|f\|^2 = 0 \Leftrightarrow f(\mathbf{x}) = 0, \text{ since } f(\mathbf{x}) = \langle f, \phi(\mathbf{x}) \rangle \leq \|f\| \|\phi(\mathbf{x})\| = 0$$

# Key aspects of kernels (2)

## Theorem (Construction of kernels)

Let  $\kappa_1$  and  $\kappa_2$  be kernels over  $\mathcal{X} \times \mathcal{X}$ ,  $a \in \mathbb{R}^+$  and  $f$  a real valued function on  $\mathcal{X}$ .

Then the following functions are kernels:

(i)  $\kappa(\mathbf{x}, \mathbf{z}) = \kappa_1(\mathbf{x}, \mathbf{z}) + \kappa_2(\mathbf{x}, \mathbf{z})$ ,

(ii)  $\kappa(\mathbf{x}, \mathbf{z}) = \kappa_1(\mathbf{x}, \mathbf{z})\kappa_2(\mathbf{x}, \mathbf{z})$ ,

(iii)  $\kappa(\mathbf{x}, \mathbf{z}) = a\kappa_1(\mathbf{x}, \mathbf{z})$ ,

(iv)  $\kappa(\mathbf{x}, \mathbf{z}) = f(\mathbf{x})f(\mathbf{z})$ .

## Popular kernels on $\mathbb{R}^n$

$\langle \mathbf{x}, \mathbf{z} \rangle^d$ ,  $\exp(-\frac{\|\mathbf{x}-\mathbf{z}\|^2}{2\sigma^2})$  and  $\tanh(\sigma\langle \mathbf{x}, \mathbf{z} \rangle) + \tau$ , where  $d \in \mathbb{N}$ ,  $\sigma$  and  $\tau \in \mathbb{R}$

# From graphical models to kernels

## “Message passing” kernel (P-kernel)

$$\begin{aligned}\kappa(\mathbf{x}, \mathbf{z}) &= p(\mathbf{x}, \mathbf{z}) \\ &= \sum_m p(\mathbf{x}, \mathbf{z} | m) p_M(m) \\ &= \sum_m p(\mathbf{x} | m) p(\mathbf{z} | m) p_M(m) \\ &= \sum_m \sum_{\boldsymbol{\theta}} p(\mathbf{x} | m, \boldsymbol{\theta}) p(\mathbf{z} | m, \boldsymbol{\theta}) p_{\Theta}(\boldsymbol{\theta} | m) p_M(m)\end{aligned}$$

- It is a **kernel!** (even using **approximate** inference algorithms!)
- Kernel matrix is computed by **message passing** on graph.
- **Applications**
  - **DNA sequence analysis:** phylogenetic trees, HMMs
  - **Information retrieval:** hierarchical model for documents
  - **Signal processing:** e.g., EMG/speech signals

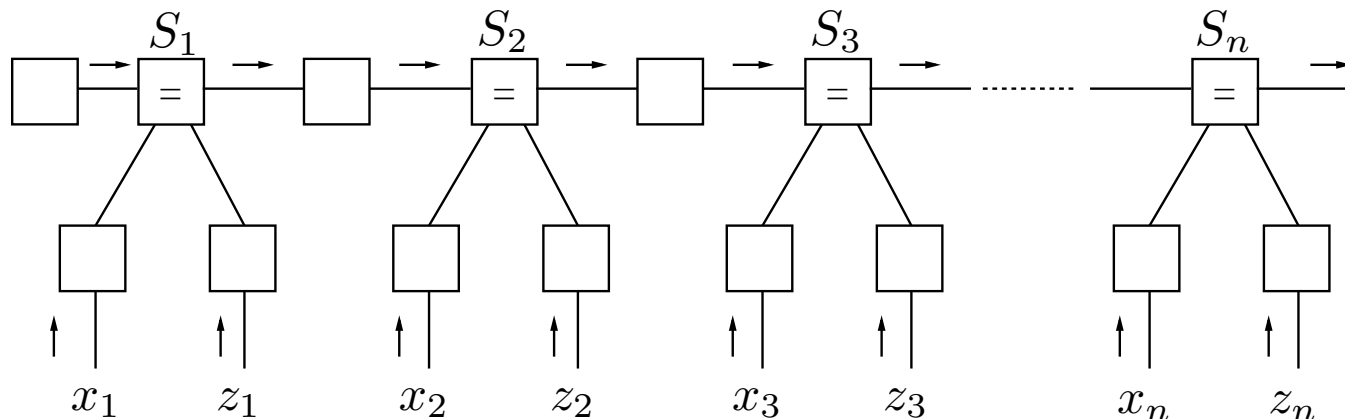
# From graphical models to kernels (2)

“Message passing” kernel (P-kernel)

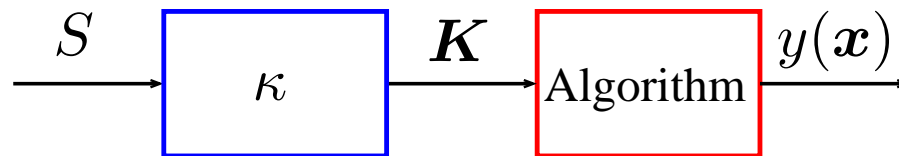
$$\kappa(\mathbf{x}, \mathbf{z}) = \sum_{\boldsymbol{\theta}} p(\mathbf{x}|\boldsymbol{\theta})p(\mathbf{z}|\boldsymbol{\theta})p_{\Theta}(\boldsymbol{\theta})$$

Example: Hidden Markov Model

$$\kappa(\mathbf{x}, \mathbf{z}) = \sum_{s_1} \cdots \sum_{s_n} p(s_1) \prod_{k=2}^n p(x_k | s_k) p(z_k | s_k) p(s_k | s_{k-1})$$



# Overall picture



## Kernel $\kappa$

- Depends on data structure of  $\mathbf{x}$  (strings, documents, graphs)
- Corresponds to embedding in (high dimensional) space
- The sole nonlinear element in the system

## Kernel matrix $\mathbf{K}$ = information bottleneck

- Input of the learning algorithm
- Encodes the input data  $\mathbf{x}$

## Learning algorithm

- Operates on the kernel matrix
- Does NOT depend on the data type of the input!

# General kernel-based classification algorithms

Minimum distance classifier in input space

$$g(\mathbf{x}) = \text{sign} \left( \frac{1}{n_+} \sum_{i|y_i=+1} \langle \mathbf{x}, \mathbf{x}_i \rangle - \frac{1}{n_-} \sum_{i|y_i=-1} \langle \mathbf{x}, \mathbf{x}_i \rangle + b \right)$$

Minimum distance classifier in feature space

$$g(\mathbf{x}) = \text{sign} \left( \frac{1}{n_+} \sum_{i|y_i=+1} \kappa(\mathbf{x}, \mathbf{x}_i) - \frac{1}{n_-} \sum_{i|y_i=-1} \kappa(\mathbf{x}, \mathbf{x}_i) + b \right)$$

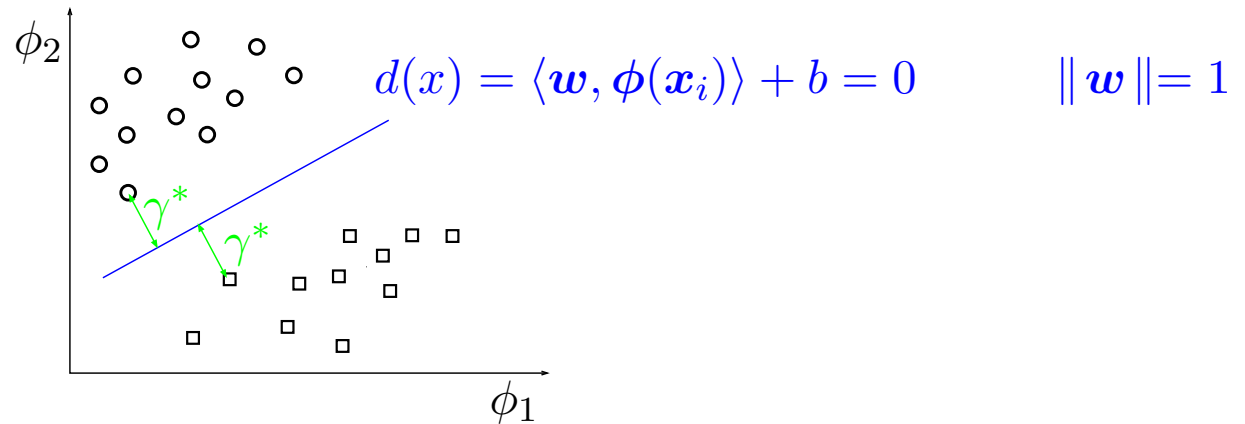
Support vector machine

$$g(\mathbf{x}) = \text{sign} \left( \sum_{i=1}^n \beta_i \kappa(\mathbf{x}, \mathbf{x}_i) + b \right)$$

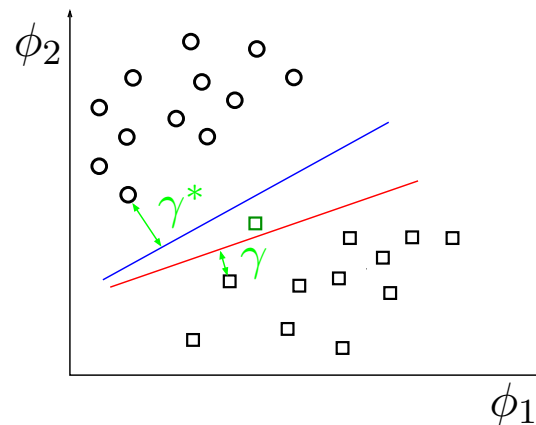
- General kernels  $\kappa(.,.)$
- Point  $\mathbf{x}$  classified by comparing to all input patterns  $\mathbf{x}_i$  with nonzero  $\beta_i$  (“support vectors”)

# Support Vector Machines

## Maximal margin classifier

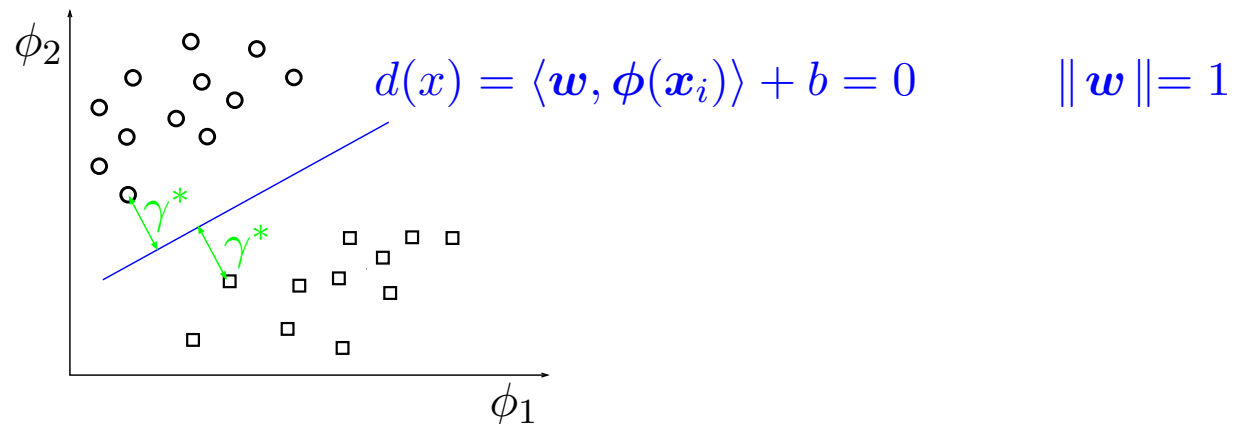


Maximize margin  $\gamma$  = distance between hyperplane and closest data point  
 $\Rightarrow$  Good generalization



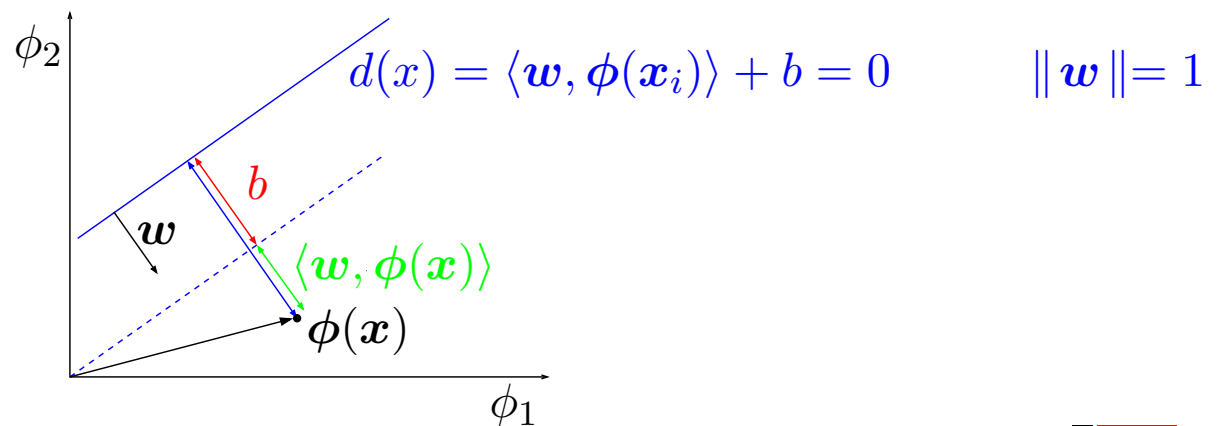
# Support Vector Machines (2)

## Maximal margin classifier



Maximize margin  $\gamma =$  distance between hyperplane and closest data point

$$\gamma = \min_{1 \leq i \leq l} y_i d(\mathbf{x}_i) = \min_{1 \leq i \leq l} y_i (\langle \mathbf{w}, \phi(\mathbf{x}_i) \rangle + b)$$



# Intermezzo: Constrained optimization

## Primal problem

Minimize  $f_0(\mathbf{x})$ ,  $\mathbf{x} \in \mathbb{R}^d$ , subject to  $f_i(\mathbf{x}) \leq 0$  ( $1 \leq i \leq m$ ) and  $h_i(\mathbf{x}) = 0$  ( $1 \leq i \leq p$ )

Define **Lagrangian**  $L : \mathbb{R}^d \times \mathbb{R}^m \times \mathbb{R}^p \rightarrow \mathbb{R}$

$$L(\mathbf{x}, \boldsymbol{\lambda}, \boldsymbol{\nu}) \triangleq f_0(\mathbf{x}) + \sum_{i=1}^m \lambda_i f_i(\mathbf{x}) + \sum_{i=1}^p \nu_i h_i(\mathbf{x})$$

and **Lagrange dual function**  $g : \mathbb{R}^m \times \mathbb{R}^p \rightarrow \mathbb{R}$  as  $g(\boldsymbol{\lambda}, \boldsymbol{\nu}) = \inf_{\mathbf{x} \in D} L(\mathbf{x}, \boldsymbol{\lambda}, \boldsymbol{\nu})$

$(\boldsymbol{\lambda}, \boldsymbol{\nu})$  **Lagrange multipliers** or **dual variables**

## Property

For any  $\boldsymbol{\nu}$  and  $\boldsymbol{\lambda} \geq 0$ :  $g(\boldsymbol{\lambda}, \boldsymbol{\nu}) \leq p^*$

## Dual problem

Maximize  $g(\boldsymbol{\lambda}, \boldsymbol{\nu})$  subject to  $\boldsymbol{\lambda} \geq 0$

$(\boldsymbol{\lambda}^*, \boldsymbol{\nu}^*)$  **optimal Lagrange multipliers**

## Weak duality

$d^* \leq p^*$  (even if primal problem is **not convex**);  $p^* - d^*$  is **optimal duality gap**

## Strong duality

$d^* = p^*$  (“constraint qualifications”)

# Intermezzo: Constrained optimization (2)

## Karush-Kuhn-Tucker conditions

If  $f_i$  are **convex** and  $h_i$  are **affine** and if  $(\tilde{\mathbf{x}}, \tilde{\boldsymbol{\lambda}}, \tilde{\boldsymbol{\nu}})$  satisfy the KKT-conditions

$$h_i(\tilde{\mathbf{x}}) = 0, \quad i = 1, \dots, p$$

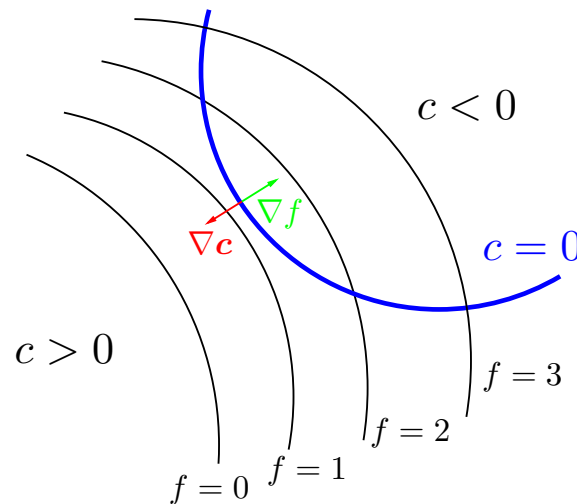
$$f_i(\tilde{\mathbf{x}}) \leq 0, \quad i = 1, \dots, m$$

$$\tilde{\lambda}_i \geq 0, \quad i = 1, \dots, m$$

$$f_i(\tilde{\mathbf{x}})\tilde{\lambda}_i = 0, \quad i = 1, \dots, m$$

$$\nabla f_0(\mathbf{x}) + \sum_{i=1}^m \lambda_i \nabla f_i(\mathbf{x}) + \sum_{i=1}^p \nu_i \nabla h_i(\mathbf{x}) = 0,$$

then  $\tilde{\mathbf{x}}$  and  $(\tilde{\boldsymbol{\lambda}}, \tilde{\boldsymbol{\nu}})$  are **primal and dual optimal** with **zero duality gap**.



# Support Vector Machines (3)

## Maximal margin classifier

(Constrained) Optimization problem

$$\begin{aligned} & \max_{\mathbf{w}, b, \gamma} \quad \gamma \\ & \text{subject to} \quad y_i(\langle \mathbf{w}, \phi(\mathbf{x}_i) \rangle + b) \geq \gamma, \quad i = 1, \dots, \ell, \text{ and } \|\mathbf{w}\|^2 = 1 \end{aligned}$$

## Primal Lagrangian

$$L(\mathbf{w}, b, \gamma, \boldsymbol{\alpha}, \lambda) = -\gamma - \sum_{i=1}^{\ell} \alpha_i [y_i(\langle \mathbf{w}, \phi(\mathbf{x}_i) \rangle + b) - \gamma] + \lambda (\|\mathbf{w}\|^2 - 1) \quad (\alpha_i \succeq 0)$$

Zero gradient equations lead to the equalities

$$\sum_{i=1}^{\ell} \alpha_i y_i = 0; \quad \sum_{i=1}^{\ell} \alpha_i = 1; \quad \mathbf{w} = \lambda \sum_{i=1}^{\ell} y_i \alpha_i \phi(\mathbf{x}_i); \quad d(\mathbf{x}) = \lambda \sum_{i=1}^{\ell} \alpha_i y_i \kappa(\mathbf{x}, \mathbf{x}_i) + b$$

Substitution in the primal Lagrangian leads to dual Lagrangian

$$L(\boldsymbol{\alpha}) = - \left( \sum_{i,j=1}^{\ell} \alpha_i \alpha_j y_i y_j \kappa(\mathbf{x}_i, \mathbf{x}_j) \right)^{1/2} \quad (\alpha_i \succeq 0, \forall i)$$

= **Convex opt. problem**, iff  $\kappa(\mathbf{x}_i, \mathbf{x}_j)$  is **positive semi-definite**, i.e., iff  $\kappa$  is a **kernel**!

# Support Vector Machines (4)

Maximal margin classifier

Dual Lagrangian

$$L(\boldsymbol{\alpha}) = - \left( \sum_{i,j=1}^{\ell} \alpha_i \alpha_j y_i y_j \kappa(\mathbf{x}_i, \mathbf{x}_j) \right)^{1/2} \quad (\alpha_i \succeq 0, \forall i)$$

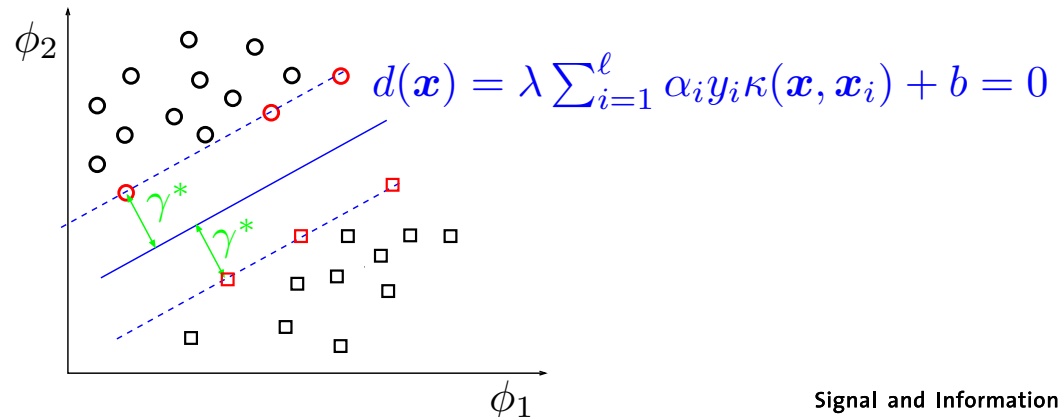
Discriminant function

$$d(\mathbf{x}) = \lambda \sum_{i=1}^{\ell} \alpha_i y_i \kappa(\mathbf{x}, \mathbf{x}_i) + b$$

Karush-Kuhn-Tucker (KKT) conditions

$$\alpha_i^* [y_i (\langle \mathbf{w}^*, \boldsymbol{\phi}(\mathbf{x}_i) \rangle + b^*) - \gamma^*] = 0, \quad i = 1, \dots, \ell.$$

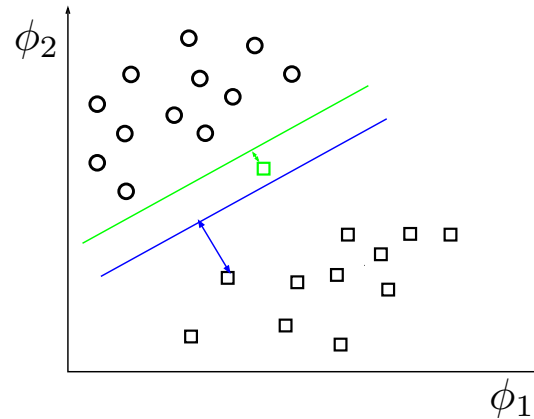
In words:  $\alpha_i^*$  is nonzero iff  $\boldsymbol{\phi}(\mathbf{x}_i)$  has geometric margin  $\gamma^*$ !



# Support Vector Machines (5)

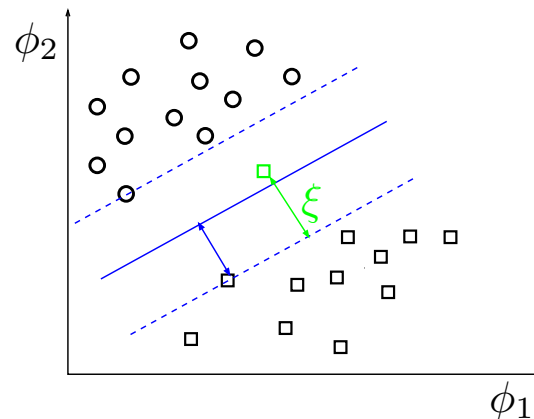
## Hard margin classifier

**Sensitive** to noise/outliers!



## Soft margin classifier

Tolerates some **misclassification** (“slack variables”  $\xi$ )



# Support Vector Machines (6)

## Soft margin classifier

### (Constrained) Optimization problem

$$\min_{\mathbf{w}, b, \gamma, \xi} -\gamma + C \sum_{i=1}^{\ell} \xi_i$$

$$\text{subject to } y_i(\langle \mathbf{w}, \phi(\mathbf{x}_i) \rangle + b) \geq \gamma - \xi_i, \quad \xi_i \geq 0, \quad i = 1, \dots, \ell, \text{ and } \|\mathbf{w}\|^2 = 1$$

### Primal Lagrangian

$$L(\mathbf{w}, b, \gamma, \boldsymbol{\alpha}, \boldsymbol{\lambda}) = -\gamma + C \sum_{i=1}^{\ell} \xi_i - \sum_{i=1}^{\ell} \alpha_i [y_i(\langle \mathbf{w}, \phi(\mathbf{x}_i) \rangle + b) - \gamma + \xi_i] \\ - \sum_{i=1}^{\ell} \beta_i \xi_i + \lambda (\|\mathbf{w}\|^2 - 1) \quad (\alpha_i \succeq 0, \beta_i \succeq 0, \forall i)$$

### Zero gradient equations lead to the equalities

$$\sum_{i=1}^{\ell} \alpha_i y_i = 0; \quad \sum_{i=1}^{\ell} \alpha_i = 1; \quad \mathbf{w} = \sum_{i=1}^{\ell} y_i \alpha_i \phi(\mathbf{x}_i); \quad C = \alpha_i + \beta_i; \quad d(\mathbf{x}) = \sum_{i=1}^{\ell} \alpha_i y_i \kappa(\mathbf{x}, \mathbf{x}_i)$$

### Dual Lagrangian

$$L(\boldsymbol{\alpha}) = - \left( \sum_{i,j=1}^{\ell} \alpha_i \alpha_j y_i y_j \kappa(\mathbf{x}_i, \mathbf{x}_j) \right)^{1/2} \quad (C \succeq \alpha_i \succeq 0, \forall i)$$

= **Convex opt. problem**, iff  $\kappa(\mathbf{x}_i, \mathbf{x}_j)$  is **positive semi-definite**, i.e., iff  $\kappa$  is a **kernel**!

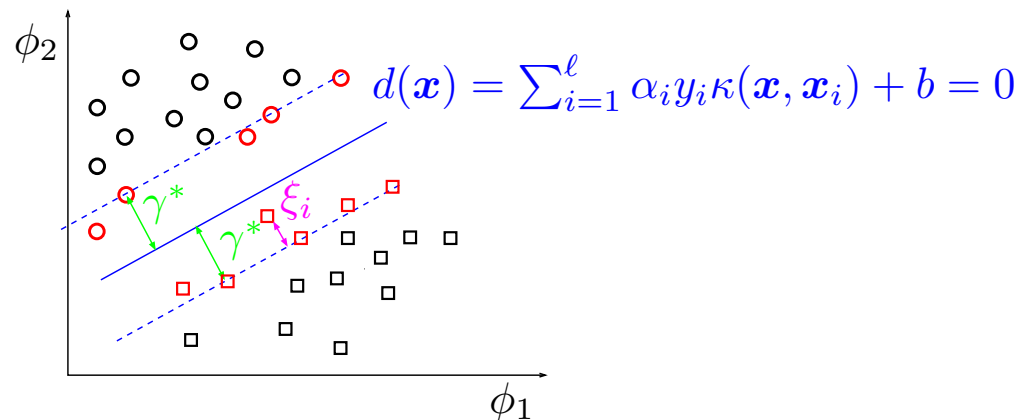
# Support Vector Machines (7)

## Soft margin classifier

### Karush-Kuhn-Tucker (KKT) conditions

$$\alpha_i^* [y_i (\langle \mathbf{w}^*, \phi(\mathbf{x}_i) \rangle + b^*) - \gamma^* + \xi_i^*] = 0; \quad \xi_i^* (\alpha_i^* - C) = 0, \quad i = 1, \dots, \ell.$$

$\alpha_i^*$  is nonzero iff input  $\mathbf{x}_i$  has geometric margin  $\gamma^*$  or  $\gamma^* - \xi_i^*$  (then  $\alpha_i^* = C$ )



# General kernel-based classification algorithms

Minimum distance classifier in input space

$$g(\mathbf{x}) = \text{sign} \left( \frac{1}{n_+} \sum_{i|y_i=+1} \langle \mathbf{x}, \mathbf{x}_i \rangle - \frac{1}{n_-} \sum_{i|y_i=-1} \langle \mathbf{x}, \mathbf{x}_i \rangle + b \right)$$

Minimum distance classifier in feature space

$$g(\mathbf{x}) = \text{sign} \left( \frac{1}{n_+} \sum_{i|y_i=+1} \kappa(\mathbf{x}, \mathbf{x}_i) - \frac{1}{n_-} \sum_{i|y_i=-1} \kappa(\mathbf{x}, \mathbf{x}_i) + b \right)$$

Support vector machine

$$g(\mathbf{x}) = \text{sign} \left( \sum_{i=1}^n \beta_i \kappa(\mathbf{x}, \mathbf{x}_i) + b \right)$$

- General kernels  $\kappa(., .)$
- Point  $\mathbf{x}$  classified by comparing to all input patterns  $\mathbf{x}_i$  with nonzero  $\beta_i$  (“support vectors”)

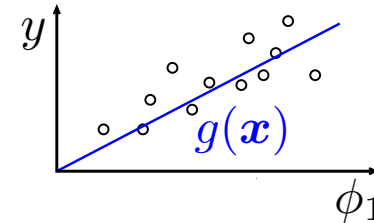
# Other kernel-based learning algorithms

## Least mean squares regression

Set  $S = \{(\mathbf{x}_1, y_1), \dots, (\mathbf{x}_n, y_n)\}$  of points  $\mathbf{x}_i \in \mathbb{R}^n$  with labels  $y_i \in \mathbb{R}$ .

Find **regression function**  $g(\mathbf{x}) = \langle \mathbf{w}, \phi(\mathbf{x}) \rangle = \sum_{i=1}^n w_i \phi(x_i)$  such that  $L(\mathbf{w}, S, \lambda) = \sum_{i=1}^l (y_i - g(\mathbf{x}_i))^2 + \lambda \|\mathbf{w}\|^2$  with  $\lambda \geq 0$  is **minimal**.

Notation:  $\mathbf{X} \triangleq (\phi(\mathbf{x}_1), \dots, \phi(\mathbf{x}_n))^T$



**Solution:**

$$(\mathbf{X}^T \mathbf{X} + \lambda \mathbf{I}_n) \mathbf{w} = \mathbf{X}^T \mathbf{Y} \quad (\text{normal equations})$$

$$g(\mathbf{x}) = \mathbf{Y}^T \mathbf{X} (\mathbf{X}^T \mathbf{X} + \lambda \mathbf{I}_n)^{-1} \phi(\mathbf{x})$$

$$= \sum_{i=1}^l \alpha_i \kappa(\mathbf{x}_i, \mathbf{x})$$

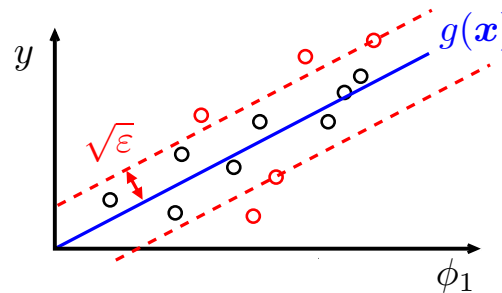
$$\boldsymbol{\alpha} = (\mathbf{K} + \lambda \mathbf{I}_n)^{-1} \mathbf{Y}$$

**BUT:  $\boldsymbol{\alpha}$  is not sparse!**

# Other kernel-based learning algorithms

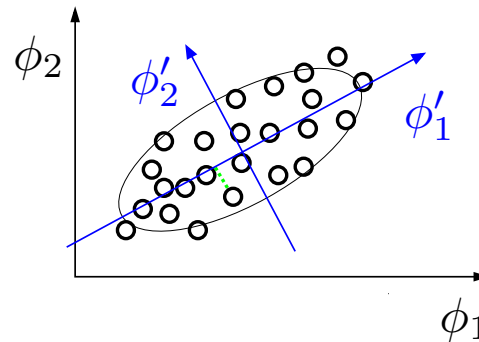
## Support vector regression

$\varepsilon$ -sensitive loss function  $L_\varepsilon(\mathbf{w}, S, \lambda) = \sum_{i=1}^l \max(0, (y_i - g(\mathbf{x}_i))^2 - \varepsilon) + \lambda \|\mathbf{w}\|^2$



# Other kernel-based learning algorithms (2)

## Principal Component analysis (PCA)



IDEA: Projection of data on  $k$ -dimensional linear subspace

Center of mass  $\boldsymbol{\mu} = 1/\ell \sum_{i=1}^{\ell} \boldsymbol{\phi}(\mathbf{x}_i)$

Empirical covariance matrix  $\mathbf{C} = 1/\ell \sum_{i=1}^{\ell} (\boldsymbol{\phi}(\mathbf{x}_i) - \boldsymbol{\mu})(\boldsymbol{\phi}(\mathbf{x}_i) - \boldsymbol{\mu})^T$

Eigenvalue problem

$$[\mathbf{U}, \boldsymbol{\Lambda}] = \text{eig}(\ell \mathbf{C})$$

Projection

$$\tilde{\mathbf{x}}_i = (\mathbf{u}_1^T \boldsymbol{\phi}(\mathbf{x}_i), \dots, \mathbf{u}_p^T \boldsymbol{\phi}(\mathbf{x}_i))^T \text{ with } \lambda_1 \geq \lambda_2 \geq \dots \geq \lambda_p$$

# Other kernel-based learning algorithms (3)

## Kernel PCA

### Key insight

If  $\mathbf{v}$  is eigenvector of  $\mathbf{K}$  with eigenvalue  $\lambda$ , then  $\mathbf{u} = \mathbf{X}^T \mathbf{v}$  is an eigenvector of  $\ell\mathbf{C}$ .

Indeed:

$$\ell\mathbf{C}\mathbf{u} = \ell\mathbf{C}\mathbf{X}^T \mathbf{v} = \mathbf{X}^T \mathbf{X} \mathbf{X}^T \mathbf{v} = \mathbf{X}^T \mathbf{K} \mathbf{v} = \lambda \mathbf{X}^T \mathbf{v} = \lambda \mathbf{u} \text{ and } \|\mathbf{u}\|^2 = \lambda$$

As a consequence:

$$\mathbf{u}_j = \lambda_j^{-1/2} \sum_{i=1}^{\ell} (v_j)_i \phi(\mathbf{x}_i) = \sum_{i=1}^{\ell} (\alpha_j)_i \phi(\mathbf{x}_i) \text{ with } \alpha_j = \lambda_j^{-1/2} v_j$$
$$\mathbf{u}_j^T \phi(\mathbf{x}) = \sum_{i=1}^{\ell} (\alpha_j)_i \kappa(\mathbf{x}_i, \mathbf{x})$$

### Algorithm

$$\tilde{\mathbf{K}} = \mathbf{K} - 1/\ell \mathbf{j} \mathbf{j}^T \mathbf{K} - 1/\ell \mathbf{K} \mathbf{j} \mathbf{j}^T + 1/\ell^2 (\mathbf{j}^T \mathbf{K} \mathbf{j}) \mathbf{j} \mathbf{j}^T \text{ (centering)}$$

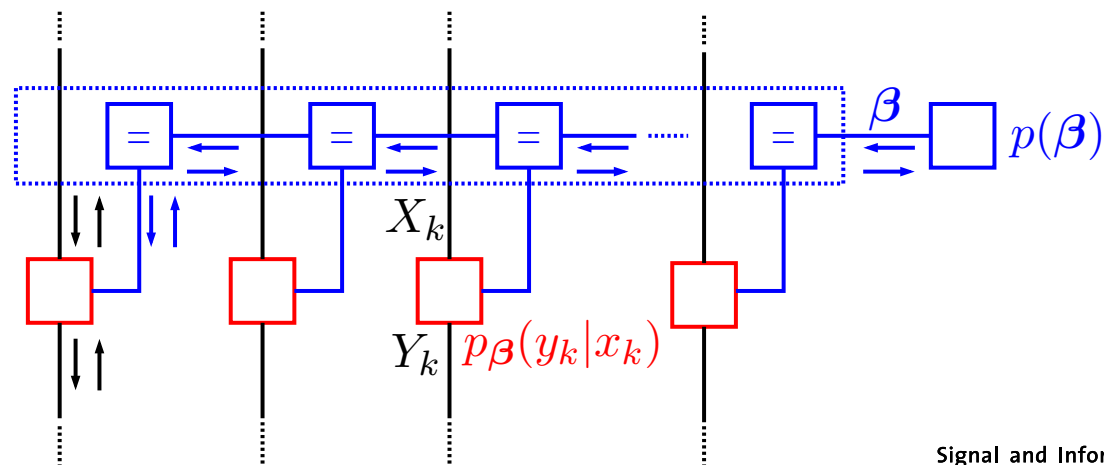
$$[\mathbf{V}, \mathbf{\Lambda}] = \text{eig}(\tilde{\mathbf{K}})$$

$$\alpha_j = \frac{1}{\sqrt{\lambda_j}} v_j, \quad j = 1, \dots, k \text{ with } \lambda_1 \geq \lambda_2 \geq \dots \geq \lambda_p$$

$$\tilde{\mathbf{x}}_i = \left( \sum_{j=1}^{\ell} (\alpha_1)_j \tilde{\kappa}(\mathbf{x}_i, \mathbf{x}), \dots, \sum_{j=1}^{\ell} (\alpha_k)_j \tilde{\kappa}(\mathbf{x}_i, \mathbf{x}) \right)$$

# Applications

- Handwriting recognition
- Speaker identification
- Face detection
- Text categorization
- Computational biology, e.g, analysis of microarray data
- Applications @ ISI
  - Message passing kernels
  - “Learning node functions” (density estimation)
  - Combinations



# Kernel machines: Pros and Cons

## Pros

- Learning algos in feature space are **well understood**
- **Modular** architecture
- **Good performance** (after some tweaking)

## Cons

- How to choose an **appropriate kernel**?
- How to determine the kernel **parameters**?

# Conclusion

- Kernel machines are **non-parametric** learning algorithms.
- **Two** main ideas
  - **(Implicit) mapping** of the input data into suitable high-dimensional **dot-product** space (“feature space”)
  - **Learning algorithm** (based on the dot product) designed to discover **linear** patterns in feature space.
- **Kernel trick** for computing inner products in feature space.
- **Kernel matrix** is interface between input data and (kernel) algorithm.
- Kernels can be derived from **graphical models**.
- **Sparseness**: only “support vectors” are relevant.

# Where to look for more information?

- Books (available @ ISIbib !)
  - *Kernel Methods for Pattern Analysis*,  
J. Shawe-Taylor and N. Cristianini
  - *Learning with Kernels*,  
B. Schölkopf and A. Smola
- Web  
<http://www.kernel-machines.org>
- ML Summer School 2005 in Canberra (Jan)  
<http://canberra05.mlss.cc/>

Thank you for your attention!

# Learning with kernels

60' Tutorial

Justin Dauwels

dauwels@isi.ee.ethz.ch

Signal and Information Processing Laboratory  
Department of Information Technology and Electrical Engineering  
ETH Zurich, Switzerland