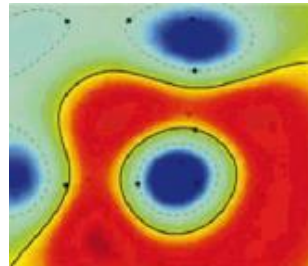


# Kernel Methods

Introduction to SVMs, KPCA, RDE

---



---

Lecture by Klaus-Robert Müller, TUB 2024

# Basic ideas in learning theory

Three scenarios: regression, classification & density estimation.

Learn  $f$  from examples

$(\mathbf{x}_1, y_1), \dots, (\mathbf{x}_N, y_N) \in \mathbb{R}^n \times \mathbb{R}^m$  or  $\{\pm 1\}$ , generated from  $P(\mathbf{x}, y)$ ,

such that expected number of errors on test set (drawn from  $P(\mathbf{x}, y)$ ),

$$R[f] = \int \frac{1}{2} |f(\mathbf{x}) - y|^2 dP(\mathbf{x}, y),$$

is minimal (*Risk Minimization (RM)*).

**Problem:**  $P$  is unknown.  $\longrightarrow$  need an *induction principle*.

*Empirical risk minimization (ERM)*: replace the average over  $P(\mathbf{x}, y)$  by an average over the training sample, i.e. minimize the training error

$$R_{emp}[f] = \frac{1}{N} \sum_{i=1}^N \frac{1}{2} |f(\mathbf{x}_i) - y_i|^2$$

# Basic ideas in learning theory II

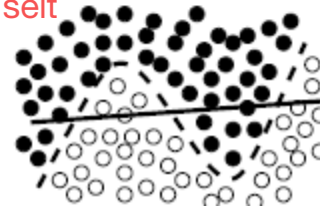
---

- Law of large numbers:  $R_{emp}[f] \rightarrow R[f]$  as  $N \rightarrow \infty$ .  
“consistency” of ERM: for  $N \rightarrow \infty$ , ERM should lead to the same result as RM?
- **No:** *uniform* convergence needed (Vapnik)  $\rightarrow$  **VC theory**.  
Thm. [classification] (Vapnik 95): with a probability of at least  $1 - \eta$ ,

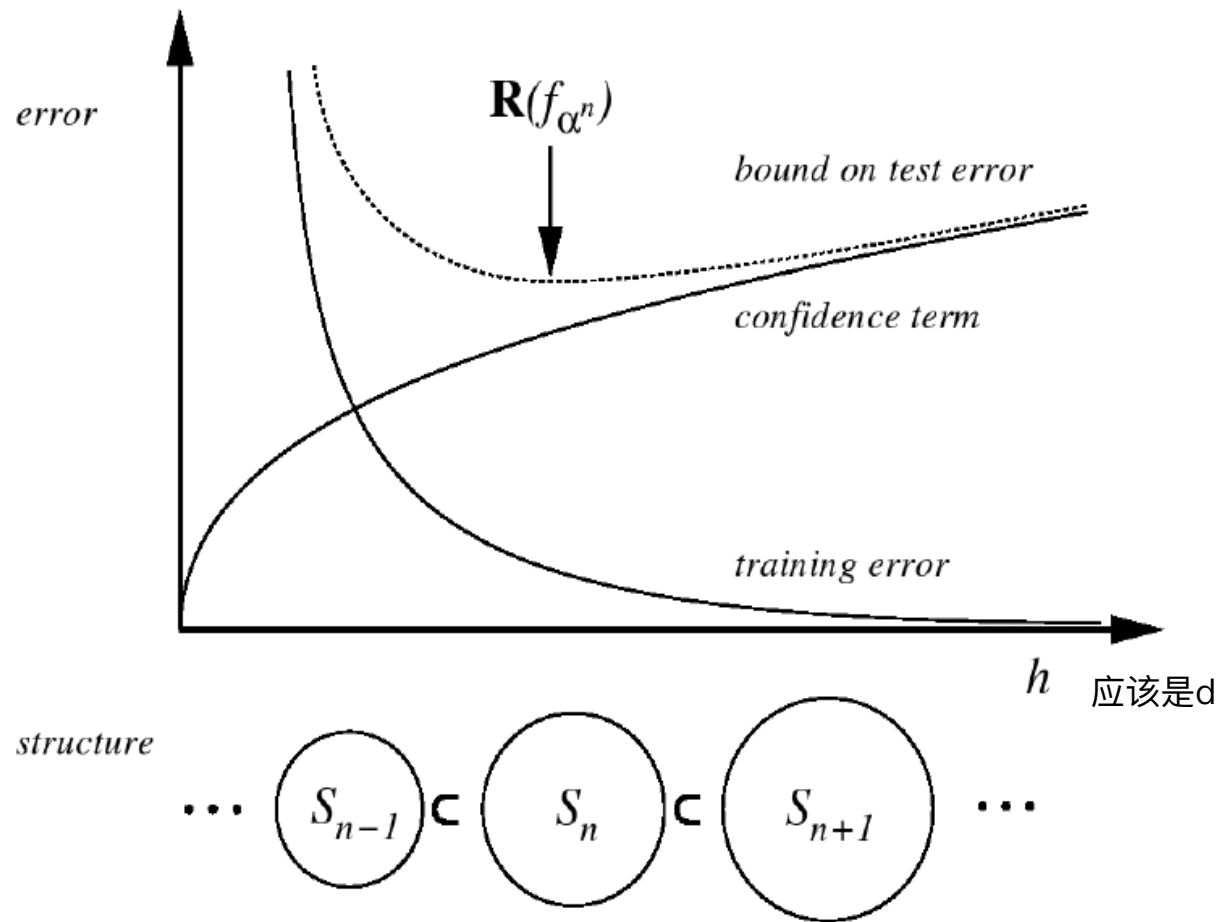
$$R[f] \leq R_{emp}[f] + \sqrt{\frac{d \left( \log \frac{2N}{d} + 1 \right) - \log(\eta/4)}{N}}.$$

- **Structural risk minimization (SRM)**: introduce structure on set of functions  $\{f_\alpha\}$  & minimize RHS to get low risk! (Vapnik 95)
- $d$  is VC dimension, measuring complexity of function class

right hand seit



## Structural Risk Minimization: the picture



Learning  $f$  requires small training error *and* small complexity of the set  $\{f_{\alpha}\}$ .

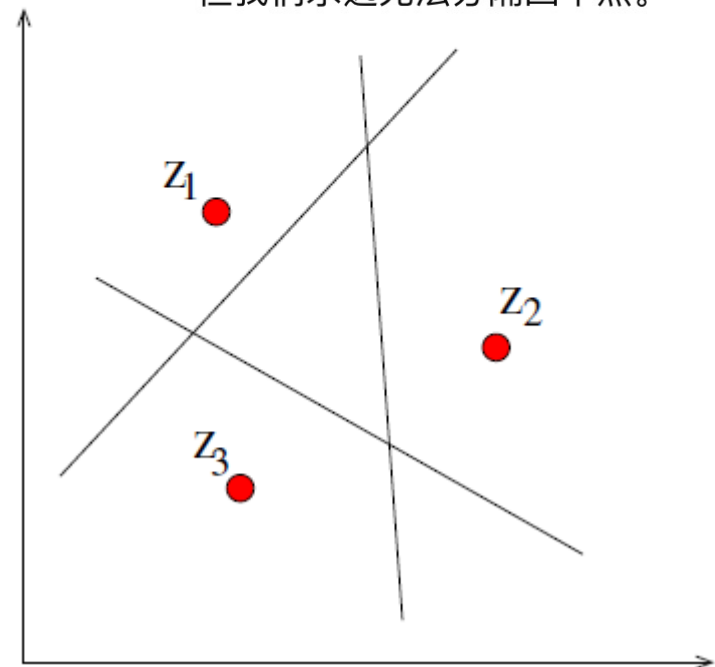
# VC Dimensions: an examples

Half-spaces in  $\mathbf{R}^2$ :

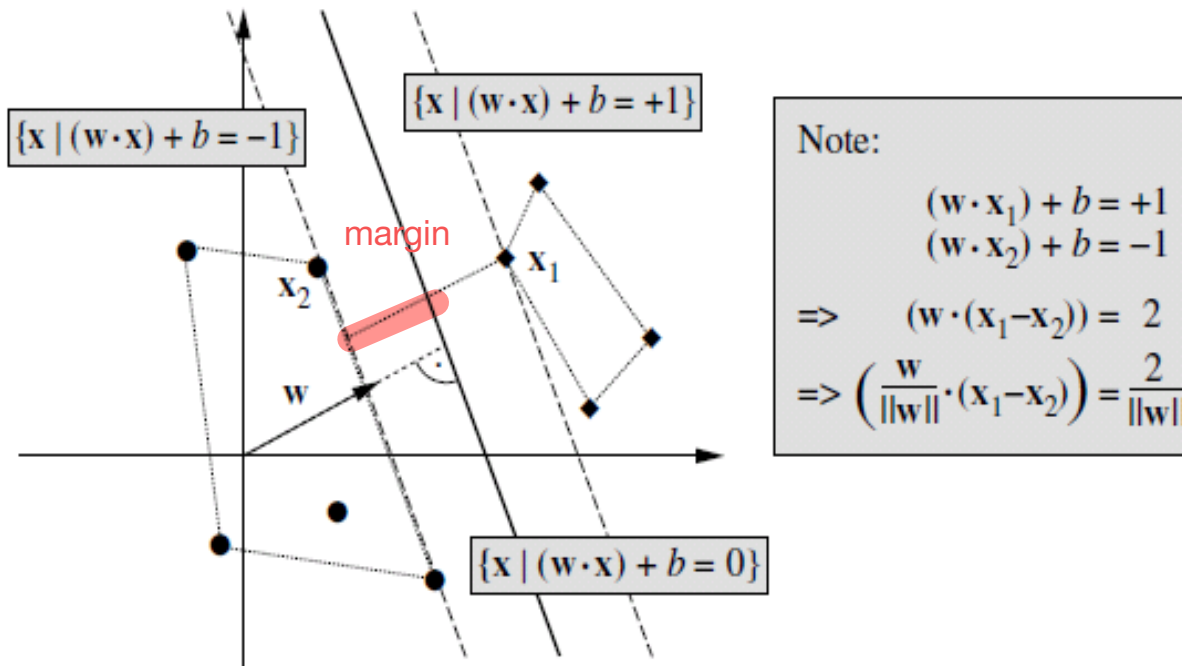
$$f(x, y) = \text{sgn}(a + bx + cy), \quad \text{with parameters } a, b, c \in \mathbf{R}$$

- Clearly, we can shatter three non-collinear points.
- But we can never shatter four points.
- Hence the VC dimension is  $d = 3$
- in  $n$  dimensions: VC dimension is  $d = n + 1$

- 很明显，我们可以分隔三个不共线的点。
- 但我们永远无法分隔四个点。



# Linear Hyperplane Classifier



- hyperplane  $y = \text{sgn}(w \cdot x + b)$  in canonical form if  $\min_{x_i \in X} |(w \cdot x_i) + b| = 1$ , i.e. scaling freedom removed.
- larger margin  $\sim 1/\|w\|$  is giving better generalization  $\rightarrow$  LMC!

# VC Theory applied to hyperplane classifiers

---

- Theorem (Vapnik 95): For hyperplanes in canonical form  
VC-dimension satisfying

$$d \leq \min\{[R^2 \|\mathbf{w}\|^2] + 1, n + 1\}.$$

Here,  $R$  is the radius of the smallest sphere containing data.  
Use  $d$  in SRM bound

$$R[f] \leq R_{emp}[f] + \sqrt{\frac{d \left( \log \frac{2N}{d} + 1 \right) - \log(\eta/4)}{N}}.$$

- maximal margin = minimum  $\|\mathbf{w}\|^2 \rightarrow$  good generalization, i.e.  
low risk, i.e. optimize

$$\min \|\mathbf{w}\|^2$$

- independent of the dimensionality of the space!

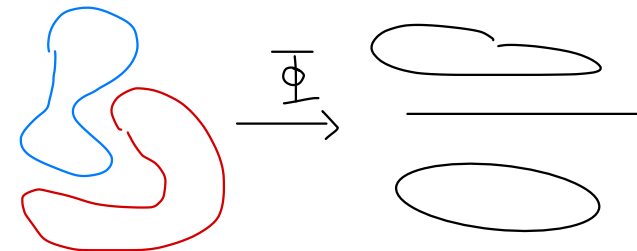


# Feature Spaces & curse of dimensionality

---

The **Support Vector (SV)** approach: *preprocess* the data with

$$\begin{aligned}\Phi : \mathbf{R}^N &\rightarrow F \\ \mathbf{x} &\mapsto \Phi(\mathbf{x}) \\ \text{where } N &\ll \dim(F).\end{aligned}$$



to get data  $(\Phi(\mathbf{x}_1), y_1), \dots, (\Phi(\mathbf{x}_N), y_N) \in F \times \mathbf{R}^M$  or  $\{\pm 1\}$ .

Learn  $\tilde{f}$  to construct  $f = \tilde{f} \circ \Phi$

- classical statistics: **harder**, as the data are high-dimensional
- SV-Learning: (in some cases) **simpler**:

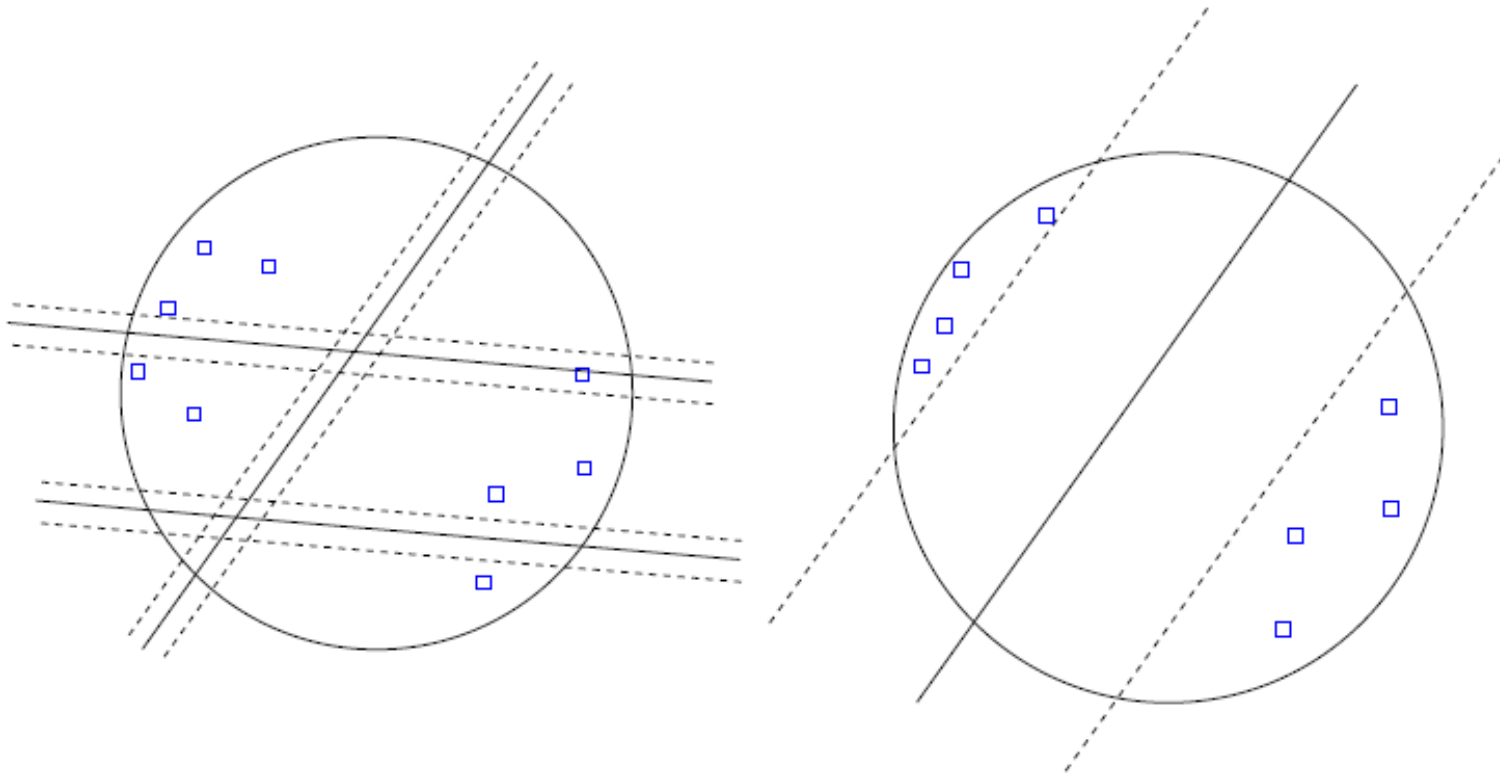
If  $\Phi$  is chosen such that  $\{\tilde{f}\}$  allows small training error *and* has low complexity, then we can guarantee good generalization.

The *complexity* matters, not the *dimensionality* of the space.



# Margin Distributions – large margin hyperplanes

---



# Feature Spaces & curse of dimensionality

---

The **Support Vector (SV)** approach: *preprocess* the data with

$$\Phi : \mathbf{R}^N \rightarrow F$$

$$\mathbf{x} \mapsto \Phi(\mathbf{x})$$

where  $N \ll \dim(F)$ .

to get data  $(\Phi(\mathbf{x}_1), y_1), \dots, (\Phi(\mathbf{x}_N), y_N) \in F \times \mathbf{R}^M$  or  $\{\pm 1\}$ .

Learn  $\tilde{f}$  to construct  $f = \tilde{f} \circ \Phi$

- classical statistics: **harder**, as the data are high-dimensional
- SV-Learning: (in some cases) **simpler**:

If  $\Phi$  is chosen such that  $\{\tilde{f}\}$  allows small training error *and* has low complexity, then we can guarantee good generalization.

The *complexity* matters, not the *dimensionality* of the space.

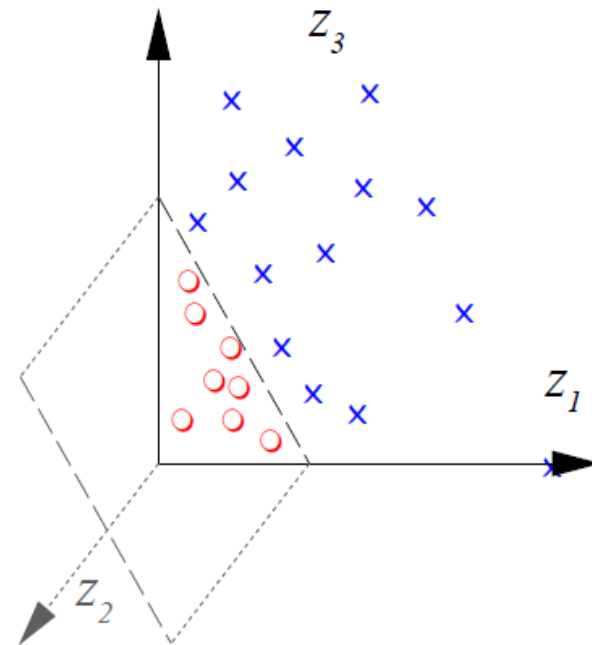
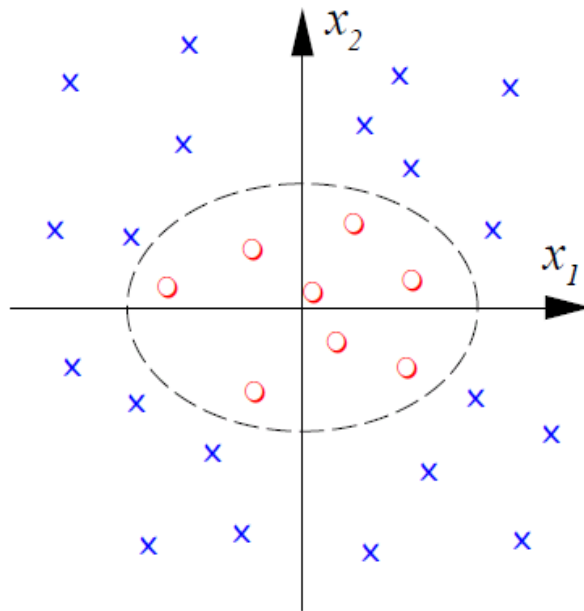


# Nonlinear Algorithms in Feature Space

Example: all second order monomials

$$\Phi : \mathbb{R}^2 \rightarrow \mathbb{R}^3$$

$$(x_1, x_2) \mapsto (z_1, z_2, z_3) := (x_1^2, \sqrt{2} x_1 x_2, x_2^2)$$



# The kernel trick: an example

---

(cf. Boser, Guyon & Vapnik 1992)

$$\begin{aligned}(\Phi(\mathbf{x}) \cdot \Phi(\mathbf{y})) &= (x_1^2, \sqrt{2} x_1 x_2, x_2^2)(y_1^2, \sqrt{2} y_1 y_2, y_2^2)^\top \\ &= (\mathbf{x} \cdot \mathbf{y})^2 \\ &=: k(\mathbf{x}, \mathbf{y})\end{aligned}$$

- Scalar product in (**high dimensional**) feature space can be computed in  $\mathbf{R}^2$ !
- works only for Mercer Kernels  $k(\mathbf{x}, \mathbf{y})$

# Kernology

---

[Mercer] If  $k$  is a continuous kernel of a positive integral operator on  $L_2(\mathcal{D})$  (where  $\mathcal{D}$  is some compact space),

$$\int f(\mathbf{x})k(\mathbf{x}, \mathbf{y})f(\mathbf{y}) \, d\mathbf{x} \, d\mathbf{y} \geq 0, \quad \text{for } f \neq 0$$

it can be expanded as

$$k(\mathbf{x}, \mathbf{y}) = \sum_{i=1}^{N_F} \lambda_i \psi_i(\mathbf{x}) \psi_i(\mathbf{y})$$

with  $\lambda_i > 0$ , and  $N_F \in \mathbf{N}$  or  $N_F = \infty$ . In that case

$$\Phi(\mathbf{x}) := \begin{pmatrix} \sqrt{\lambda_1} \psi_1(\mathbf{x}) \\ \sqrt{\lambda_2} \psi_2(\mathbf{x}) \\ \vdots \end{pmatrix}$$

satisfies  $(\Phi(\mathbf{x}) \cdot \Phi(\mathbf{y})) = k(\mathbf{x}, \mathbf{y})$ .



## Kernology II

Examples of common kernels:

$$\text{Polynomial} \quad k(\mathbf{x}, \mathbf{y}) = (\mathbf{x} \cdot \mathbf{y} + c)^d$$

~~$$\text{Sigmoid} \quad k(\mathbf{x}, \mathbf{y}) = \tanh(\kappa(\mathbf{x} \cdot \mathbf{y}) + \theta)$$~~

$$\text{RBF} \quad k(\mathbf{x}, \mathbf{y}) = \exp(-\|\mathbf{x} - \mathbf{y}\|^2 / (2\sigma^2))$$

$$\text{inverse multiquadric} \quad k(\mathbf{x}, \mathbf{y}) = \frac{1}{\sqrt{\|\mathbf{x} - \mathbf{y}\|^2 + c^2}}$$

**Note:** kernels correspond to regularization operators (a la Tichonov) with regularization properties that can be conveniently expressed in Fourier space, e.g. Gaussian kernel corresponds to general smoothness assumption (Smola et al 98 )!

# A RKHS representation of $\mathcal{F}$

---

$$\tilde{\Phi} : \mathbf{R}^N \longrightarrow \mathcal{H}, \quad \mathbf{x} \mapsto k(\mathbf{x}, \cdot)$$

Need a dot product  $\langle \cdot, \cdot \rangle$  for  $\mathcal{H}$  such that

$$\langle \tilde{\Phi}(\mathbf{x}), \tilde{\Phi}(\mathbf{y}) \rangle = k(\mathbf{x}, \mathbf{y}), \quad \text{i.e. require} \quad \langle k(\mathbf{x}, \cdot), k(\mathbf{y}, \cdot) \rangle = k(\mathbf{x}, \mathbf{y}).$$

For a Mercer kernel  $k(\mathbf{x}, \mathbf{y}) = \sum_j \lambda_j \psi_j(\mathbf{x}) \psi_j(\mathbf{y})$ , with  $\lambda_i > 0$  for all  $i$ , and  $(\psi_i \cdot \psi_j)_{L_2(\mathcal{C})} = \delta_{ij}$ , this can be achieved by choosing  $\langle \cdot, \cdot \rangle$  such that

$$\langle \psi_i, \psi_j \rangle = \delta_{ij} / \lambda_i.$$

$\mathcal{H}$ , the closure of the space of all functions

$$f(\mathbf{x}) = \sum_i a_i k(\mathbf{x}, \mathbf{x}_i),$$

with dot product  $\langle \cdot, \cdot \rangle$ , is called **reproducing kernel Hilbert space**



# Hyperplane $y = \text{sgn}(\mathbf{w} \cdot \Phi(x) + b)$ in $\mathcal{F}$

---

$$\begin{array}{ll} \min & \|\mathbf{w}\|^2 \\ \text{subject to} & y_i \cdot [(\mathbf{w} \cdot \Phi(\mathbf{x}_i)) + b] \geq 1 \quad \text{for } i = 1 \dots N \end{array}$$

(i.e. training data separated correctly, otherwise introduce slack variables).

$$L(\mathbf{w}, b, \alpha) = \frac{1}{2} \|\mathbf{w}\|^2 - \sum_{i=1}^N \alpha_i (y_i \cdot ((\mathbf{w} \cdot \Phi(\mathbf{x}_i)) + b) - 1).$$

obtain unique  $\alpha_i$  by QP (no local minima!): **dual problem**

$$\frac{\partial}{\partial b} L(\mathbf{w}, b, \alpha) = 0, \quad \frac{\partial}{\partial \mathbf{w}} L(\mathbf{w}, b, \alpha) = 0,$$

$$\text{i.e.} \quad \sum_{i=1}^N \alpha_i y_i = 0 \quad \text{and} \quad \mathbf{w} = \sum_{i=1}^N \alpha_i y_i \Phi(\mathbf{x}_i).$$

Substitute both into  $L$  to get the **dual problem**





# Hyperplane in $\mathcal{F}$ with slack variables: SVM

---

$$\min \quad \|\mathbf{w}\|^2 + C \sum_{i=1}^N \xi_i^p$$

subject to  $y_i \cdot [(\mathbf{w} \cdot \Phi(\mathbf{x}_i)) + b] \geq 1 - \xi_i$  and  $\xi_i \geq 0$  for  $i = 1 \dots N$

(introduce slack variables if training data **not** separated correctly)

$$L(\mathbf{w}, b, \alpha) = \frac{1}{2} \|\mathbf{w}\|^2 - \sum_{i=1}^N \alpha_i (y_i \cdot ((\mathbf{w} \cdot \Phi(\mathbf{x}_i)) + b) - 1).$$

obtain unique  $\alpha_i$  by QP (no local minima!): **dual problem**

$$\frac{\partial}{\partial b} L(\mathbf{w}, b, \alpha) = 0, \quad \frac{\partial}{\partial \mathbf{w}} L(\mathbf{w}, b, \alpha) = 0,$$

$$\text{i.e.} \quad \sum_{i=1}^N \alpha_i y_i = 0 \quad \text{and} \quad \mathbf{w} = \sum_{i=1}^N \alpha_i y_i \Phi(\mathbf{x}_i).$$

Substitute both into  $L$  to get the **dual problem**



$$\frac{1}{2} \|\mathbf{w}\|^2 = \sum_{i=1}^N \alpha_i \left( y_i \cdot \left( \sum_{j=1}^N \alpha_j y_j \Phi(\mathbf{x}_j) \cdot \Phi(\mathbf{x}_i) + b \right) - 1 \right)$$

## Dual Problem

$$= \frac{1}{2} \sum_{i=1}^N \left( \alpha_i y_i \Phi(\mathbf{x}_i) \right)^2 - \sum_{i=1}^N \sum_{j=1}^N \alpha_i \alpha_j y_i y_j \Phi(\mathbf{x}_i) \cdot \Phi(\mathbf{x}_j) - 0 + \sum_{i=1}^N \alpha_i$$


---

maximize

$$W(\alpha) = \sum_{i=1}^N \alpha_i - \frac{1}{2} \sum_{i,j=1}^N \alpha_i \alpha_j y_i y_j k(\mathbf{x}_i, \mathbf{x}_j)$$

subject to

$$C \geq \alpha_i \geq 0, \quad i = 1, \dots, N, \quad \text{and} \quad \sum_{i=1}^N \alpha_i y_i = 0.$$

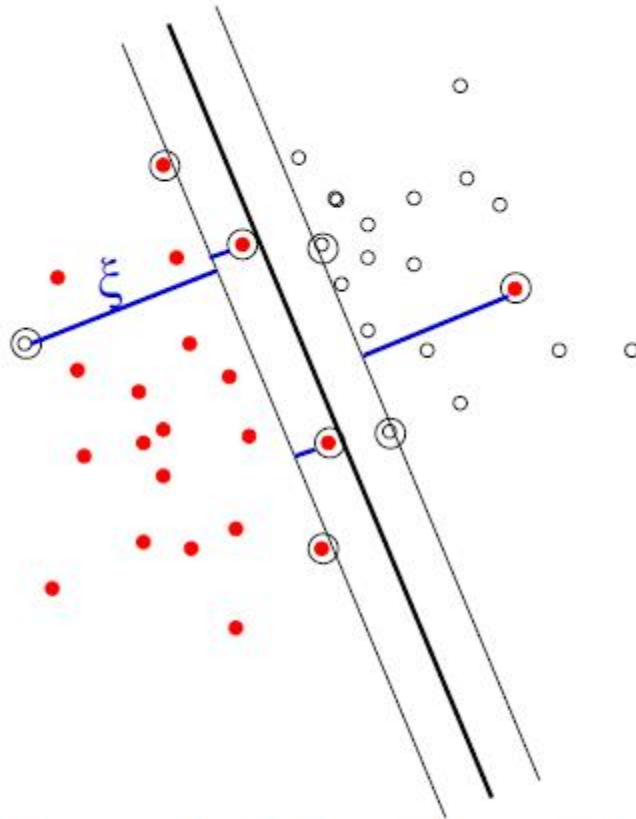
Note: solution determined by training examples (SVs) on /in the margin. Remark: duality gap.

$$y_i \cdot [(\mathbf{w} \cdot \Phi(\mathbf{x}_i)) + b] > 1 \quad \implies \alpha_i = 0 \longrightarrow \mathbf{x}_i \text{ irrelevant or}$$

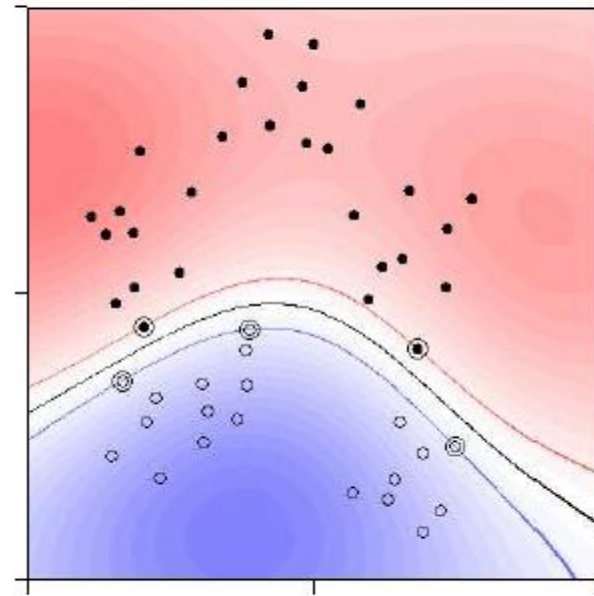
$$y_i \cdot [(\mathbf{w} \cdot \Phi(\mathbf{x}_i)) + b] = 1 \quad (\text{on /in margin}) \longrightarrow \mathbf{x}_i \text{ Support Vector}$$

# A Toy Example: $k(\mathbf{x}, \mathbf{y}) = \exp(-\|\mathbf{x} - \mathbf{y}\|^2)$

---



linear SV with slack variables



nonlinear SVM, Domain:  $[-1, 1]^2$

# Kernel Trick

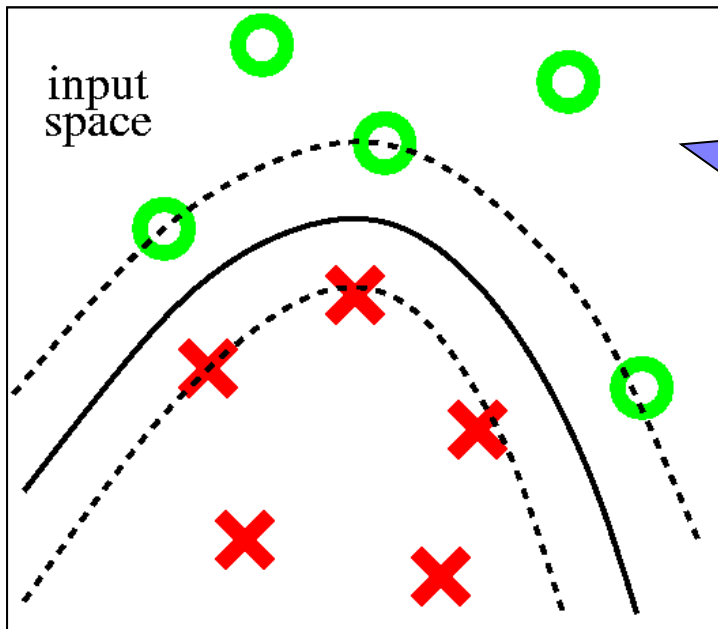
---

- Saddle Point:  $\mathbf{w} = \sum_{i=1}^N \alpha_i y_i \Phi(\mathbf{x}_i)$ .
- Hyperplane in  $\mathcal{F}$ :  $y = \text{sgn}(\mathbf{w} \cdot \Phi(x) + b)$
- putting things together “kernel trick”

$$\begin{aligned} f(\mathbf{x}) &= \text{sgn}(\mathbf{w} \cdot \Phi(\mathbf{x}) + b) \\ &= \text{sgn}\left(\sum_{i=1}^N \alpha_i y_i \Phi(\mathbf{x}_i) \cdot \Phi(\mathbf{x}) + b\right) \\ &= \text{sgn}\left(\sum_{i \in \#SV_S} \alpha_i y_i k(\mathbf{x}, \mathbf{x}_i) + b\right) \quad \text{sparse!} \end{aligned}$$

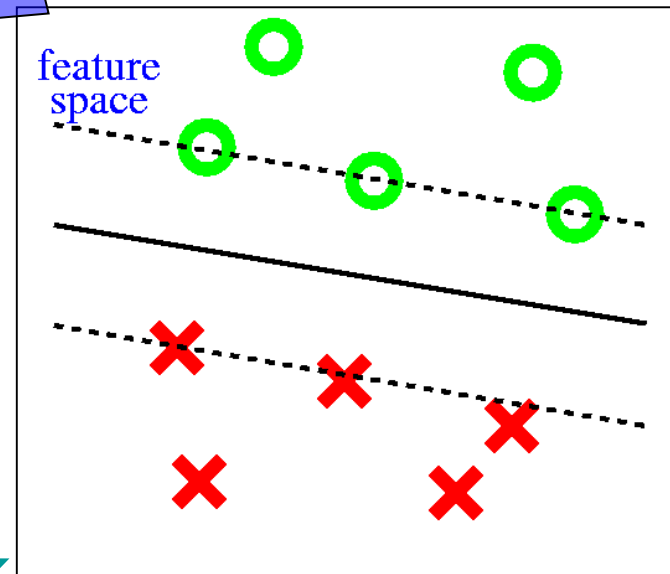
- trick:  $k(\mathbf{x}, \mathbf{y}) = \Phi(\mathbf{x}) \cdot \Phi(\mathbf{y})$ , i.e. **never use  $\Phi$ : only  $k$ !!!**

# Support Vector Machines in a nutshell



$$\Phi \text{ rsp. } K(x,y) = \Phi(x) \cdot \Phi(y)$$

$\Phi$



**good theory**

non-linear decision by  
implicitly **mapping** the data

into feature space by SV **kernel** function **K**

# Kernels ...

---

- kernels hold key to learning problem.
- **choosing kernels ...**
  - Mercer condition ( $\ell_2$  integrability & positivity)
  - kernel reflects prior (Smola, Schölkopf & Müller 98, Girosi 98)
  - approximating LOO bounds give good model selection results (Tsuda et al. 2001, Vapnik & Chapelle 2000)
- So: **engineer** an appropriate kernel from prior knowledge! (Jaakola and Haussler 1998, Watkins 2000, Zien et al 2000, recently a large body of interesting work)
- And: use **careful** model selection to find appropriate kernel parameters, i.e. chose appropriate degree of polynomial or bandwidth of Gaussian kernel

# Digestion: Use of kernels

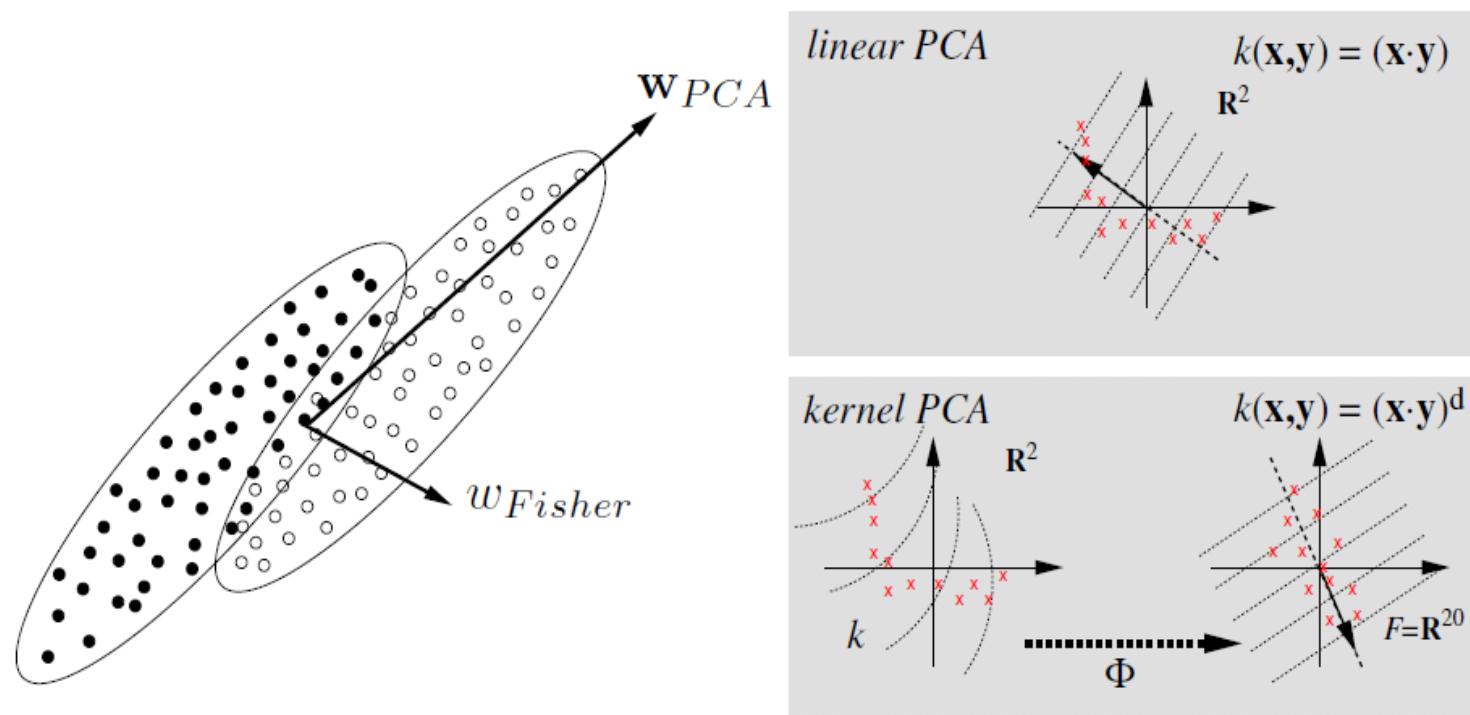
---

- **Question:** What makes kernel methods (e.g. SVM) perform well?
- **Answer:**
  - In the first place: a good idea/theory.
  - But also: **The kernel**
- Using kernels, we work explicitly in extremely high dimensional spaces (RKHS) with interesting features for themselves (depending on the kernel) [SSM et al. 98]
- Common choices: Gaussian kernel  $\exp(-\|x - y\|^2/c)$  or polynomial kernel  $(x \cdot y)^d$ .
- Almost any linear algorithm can be transformed to feature space. [SSM et al. 98]
- With suitable regularization it outperforms its linear counterpart. [Mika et al. 02]  
[Zien et al. 00, Tsuda et al. 02, Sonnenburg et al. 05]
- **The kernel can be adopted to specific tasks, e.g. using prior knowledge**



Kernels for graphs,  
trees, strings etc.

## Remark: Kernelizing linear algorithms



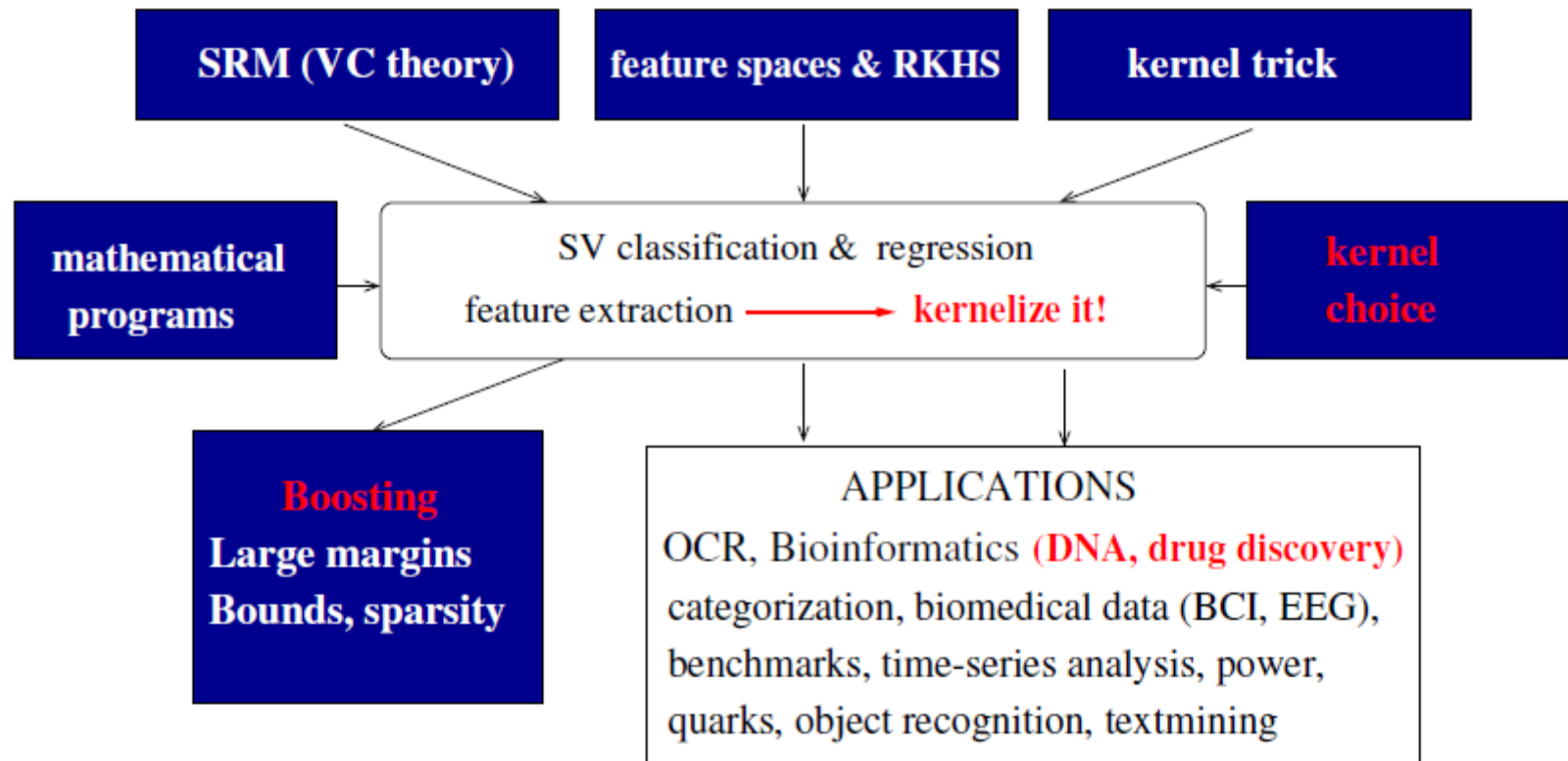
(cf. Schölkopf, Smola and Müller 1996, 1998, Schölkopf et al 1999, Mika et al, 1999, 2000, 2001, Müller et al 2001, Harmeling et al 2003, ...)



# Digestion

$$R[f] \leq R_{emp}[f] + \sqrt{\frac{d\left(\log \frac{2N}{d} + 1\right) - \log(\eta/4)}{N}}$$

$$K(\mathbf{x}, \mathbf{y}) = \Phi(\mathbf{x}) \cdot \Phi(\mathbf{y})$$



# Part II

---



# Implementation Issues: working set methods

---

matrix notation: Let  $\alpha = (\alpha_1, \dots, \alpha_M)^\top$ , let  $\mathbf{y} = (y_1, \dots, y_M)^\top$ , let  $H$  be the matrix with the entries  $H_{ij} = y_i y_j k(\mathbf{x}_i, \mathbf{x}_j)$ , and let  $\mathbf{1}$  denote the vector of length  $M$  consisting of all 1s.

dual SVM Problem becomes:

$$\max_{\alpha} \quad \mathbf{1}^\top \alpha - \frac{1}{2} \alpha^\top H \alpha, \quad (1)$$

$$\text{subject to} \quad \mathbf{y}^\top \alpha = 0, \quad (2)$$

$$\alpha - C\mathbf{1} \leq 0, \quad (3)$$

$$\alpha \geq 0. \quad (4)$$

## Implementation Issues: working set methods II

---

$\alpha_B$  of the variables in the working set at a current iteration and  $\alpha_N$  remaining variables.  $H$  is thus partitioned as  $H = \begin{bmatrix} H_{BB} & H_{BN} \\ H_{NB} & H_{NN} \end{bmatrix}$ ,

at each iteration, is obtained:

$$\max_{\alpha} \quad (\mathbf{1}_B^\top - \alpha_N^\top H_{NB}) \alpha_B - \frac{1}{2} \alpha_B^\top H_{BB} \alpha_B, \quad (5)$$

$$\text{subject to} \quad \mathbf{y}_B^\top \alpha_B = -\mathbf{y}_N^\top \alpha_N, \quad (6)$$

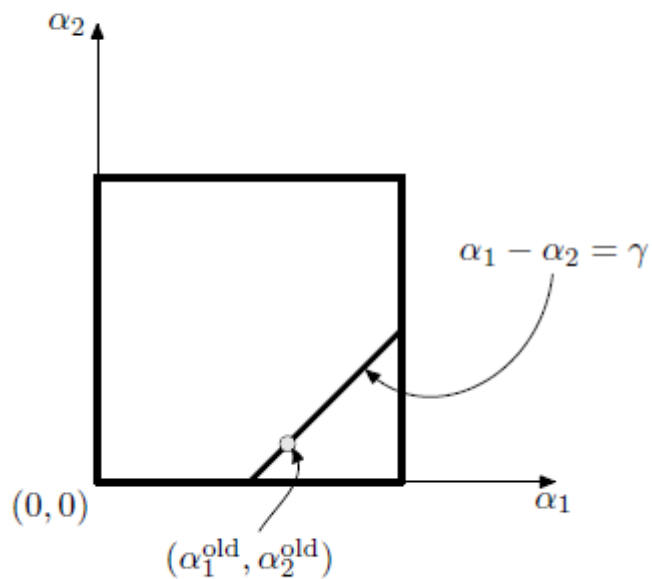
$$\alpha_B - C \mathbf{1}_B \leq 0, \quad (7)$$

$$\alpha_B \geq 0. \quad (8)$$

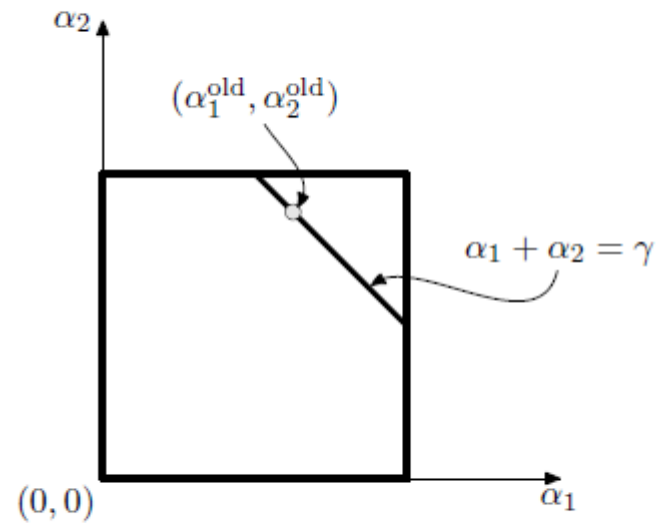
Usual small working set, iteration is carried out until KKT conditions, are satisfied to the required precision  $\epsilon$ . monitor gap.

# John Platt's SMO

- extreme: use only two points in working set and compute optimal solution *analytically*



$$y_1 \neq y_2$$



$$y_1 = y_2$$

## SMO continued

---

eliminating  $\alpha_1$  yields update rule for  $\alpha_2$ :

$$\alpha_2^{\text{new}} = \alpha_2^{\text{old}} - \frac{y_2(E_1 - E_2)}{\eta}, \quad (9)$$

where

$$E_1 = \sum_{j=1}^M y_j \alpha_j k(\mathbf{x}_1, \mathbf{x}_j) + b - y_1, \quad (10)$$

$$E_2 = \sum_{j=1}^M y_j \alpha_j k(\mathbf{x}_2, \mathbf{x}_j) + b - y_2, \quad (11)$$

$$\eta = 2 k(\mathbf{x}_1, \mathbf{x}_2) - k(\mathbf{x}_1, \mathbf{x}_1) - k(\mathbf{x}_2, \mathbf{x}_2). \quad (12)$$

Next, the bound constraints should be taken care of. Depending on the geometry, one computes the following lower and upper bounds on the

value of the variable  $\alpha_2$ :

$$L = \begin{cases} \max(0, \alpha_2^{\text{old}} - \alpha_1^{\text{old}}), & \text{if } y_1 \neq y_2, \\ \max(0, \alpha_2^{\text{old}} + \alpha_1^{\text{old}} - C), & \text{if } y_1 = y_2, \end{cases}$$

$$H = \begin{cases} \min(C, C + \alpha_2^{\text{old}} - \alpha_1^{\text{old}}), & \text{if } y_1 \neq y_2, \\ \min(C, \alpha_2^{\text{old}} + \alpha_1^{\text{old}}), & \text{if } y_1 = y_2. \end{cases}$$

The constrained optimum is then found by clipping the unconstrained optimum to the ends of the line segment:

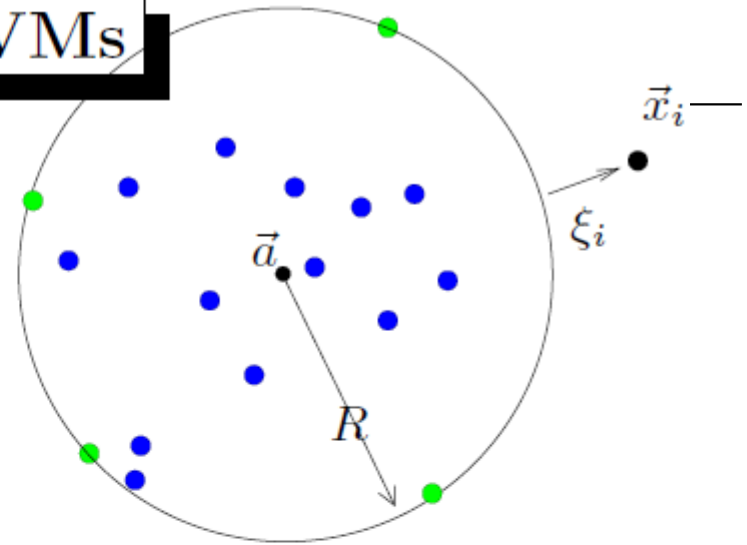
$$\alpha_2^{\text{new}} := \begin{cases} H, & \text{if } \alpha_2^{\text{new}} \geq H, \\ L, & \text{if } \alpha_2^{\text{new}} \leq L, \\ \alpha_2^{\text{new}}, & \text{otherwise.} \end{cases}$$

Finally, the value of  $\alpha_1^{\text{new}}$  is computed:

$$\alpha_1^{\text{new}} = \alpha_1^{\text{old}} + y_1 y_2 (\alpha_2^{\text{old}} - \alpha_2^{\text{new}}). \quad (13)$$

- Use heuristics to choose examples

# One-Class SVMs



Fitting a hypersphere around the data

$$\max_{\alpha} \quad \sum_{i=1}^M \alpha_i k(\mathbf{x}_i, \mathbf{x}_i) - \frac{1}{2} \sum_{i,j=1}^M \alpha_i \alpha_j k(\mathbf{x}_i, \mathbf{x}_j), \quad (14)$$

subject to  $0 \leq \alpha_i \leq C, \quad i = 1, \dots, M,$

$$\sum_{i=1}^M \alpha_i = 1.$$

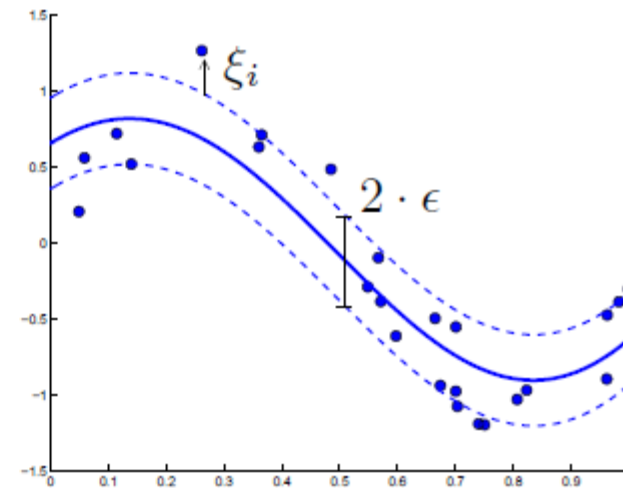
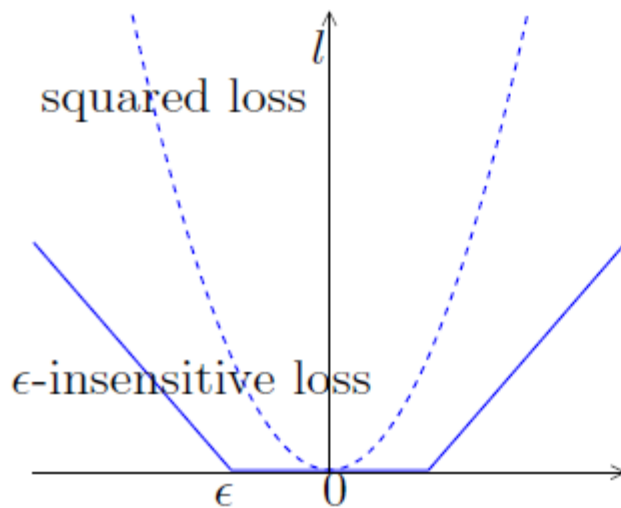
new object belongs to target class? (cf. Tax 01, Schölkopf et al. 01)

$$f(\mathbf{x}) = \text{sign}(R^2 - k(\mathbf{x}, \mathbf{x}) + 2 \sum_i \alpha_i k(\mathbf{x}, \mathbf{x}_i) - \sum_{i,j} \alpha_i \alpha_j k(\mathbf{x}_i, \mathbf{x}_j)). \quad (15)$$



# SVMs for Regression

$$\ell(f(\mathbf{x}), y) = (f(\mathbf{x}) - y)^2,$$
$$\ell(f(\mathbf{x}), y) = \begin{cases} |f(\mathbf{x}) - y| - \epsilon, & \text{if } |f(\mathbf{x}) - y| > \epsilon, \\ 0, & \text{otherwise.} \end{cases}$$

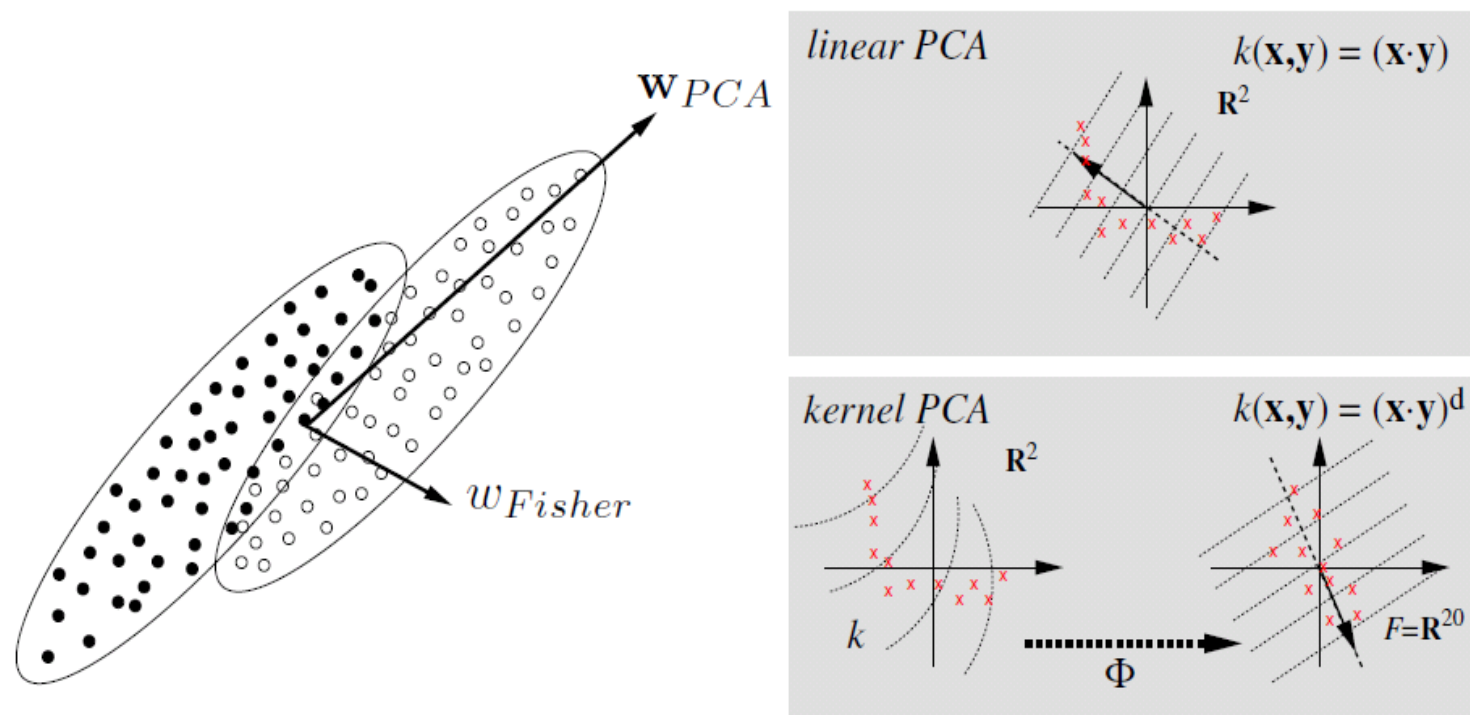


(cf. Vapnik 95, Smola and Schölkopf 02)

The primal formulation for the SVR

$$\begin{aligned} \min_{\mathbf{w}, b, \boldsymbol{\xi}^{(*)}} \quad & \frac{1}{2} \|\mathbf{w}\|^2 + C \sum_{i=1}^M (\xi_i + \xi_i^*), \\ \text{subject to} \quad & ((\mathbf{w}^\top \mathbf{x}_i) + b) - y_i \leq \epsilon + \xi_i, \\ & y_i - ((\mathbf{w}^\top \mathbf{x}_i) + b) \leq \epsilon + \xi_i^*, \\ & \xi_i^{(*)} \geq 0, \quad i = 1, \dots, M. \end{aligned}$$

## Remark: Kernelizing linear algorithms

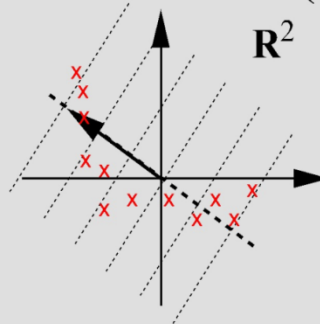


(cf. Schölkopf, Smola and Müller 1996, 1998, Schölkopf et al 1999, Mika et al, 1999, 2000, 2001, Müller et al 2001, Harmeling et al 2003, ...)

# Kernel PCA

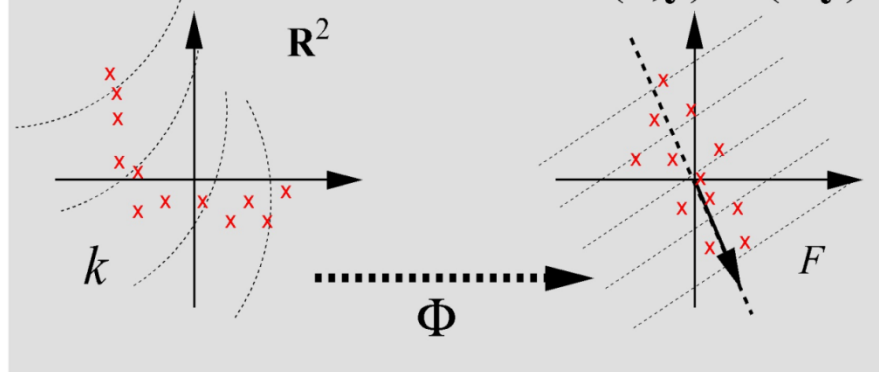
*linear PCA*

$$k(\mathbf{x}, \mathbf{y}) = (\mathbf{x} \cdot \mathbf{y})$$



*kernel PCA*

$$k(\mathbf{x}, \mathbf{y}) = (\mathbf{x} \cdot \mathbf{y})^d$$



## PCA in high dimensional feature spaces

$$\mathbf{x}_1, \dots, \mathbf{x}_N, \quad \Phi : \mathbb{R}^D \rightarrow F, \quad \mathbf{C} = \frac{1}{N} \sum_{j=1}^N \Phi(\mathbf{x}_j) \Phi(\mathbf{x}_j)^\top$$

Eigenvalue problem

$$\lambda \mathbf{V} = \mathbf{C} \mathbf{V} = \frac{1}{N} \sum_{j=1}^N (\Phi(\mathbf{x}_j) \cdot \mathbf{V}) \Phi(\mathbf{x}_j).$$

For  $\lambda \neq 0$ ,  $\mathbf{V} \in \text{span}\{\Phi(\mathbf{x}_1), \dots, \Phi(\mathbf{x}_N)\}$ , thus  $\mathbf{V} = \sum_{i=1}^N \alpha_i \Phi(\mathbf{x}_i)$ .

Multiplying with  $\Phi(\mathbf{x}_k)$  from the left yields

$$\mathbf{N} \lambda (\Phi(\mathbf{x}_k) \cdot \mathbf{V}) = (\Phi(\mathbf{x}_k) \cdot \mathbf{C} \mathbf{V}) \text{ for all } k = 1, \dots, N$$

# Nonlinear PCA as an Eigenvalue problem

---

Define an  $N \times N$  matrix

$$K_{ij} := (\Phi(\mathbf{x}_i) \cdot \Phi(\mathbf{x}_j)) = k(\mathbf{x}_i, \mathbf{x}_j)$$

to get

$$N\lambda K\alpha = K^2\alpha$$

where  $\alpha = (\alpha_1, \dots, \alpha_N)^\top$ .

Solve

$$N\lambda\alpha = K\alpha$$

$$\longrightarrow (\lambda_k, \alpha^k)$$

$$(\mathbf{V}^k \cdot \mathbf{V}^k) = 1 \iff N\lambda_k(\alpha^k \cdot \alpha^k) = 1$$

# Feature Extraction

Compute projections on the Eigenvectors

$$\mathbf{v}^k = \sum_{i=1}^M \alpha_i^k \Phi(\mathbf{x}_i)$$

in  $F$ :

for a test point  $\mathbf{x}$  with image  $\Phi(\mathbf{x})$  in  $F$  we get the features  
("kernel PCA components")

$$\begin{aligned} f_k(x) = (\mathbf{v}^k \cdot \Phi(\mathbf{x})) &= \sum_{i=1}^M \alpha_i^k (\Phi(\mathbf{x}_i) \cdot \Phi(\mathbf{x})) \\ &= \sum_{i=1}^M \alpha_i^k k(\mathbf{x}_i, \mathbf{x}) \end{aligned}$$

# Centering in Feature Space

---

Center the data in  $F$ :

$$\tilde{\Phi}(\mathbf{x}_i) := \Phi(\mathbf{x}_i) - \frac{1}{N} \sum_{i=1}^N \Phi(\mathbf{x}_i)$$

For  $\tilde{\Phi}(\mathbf{x}_i)$ , everything works fine.

Express  $\tilde{K}$  in terms of  $K$ , using  $(1_N)_{ij} := 1/N$ :

$$\tilde{K}_{ij} = K - 1_N K - K 1_N + 1_N K 1_N.$$

Compute  $\tilde{K}$  and solve the Eigenvalue problem.

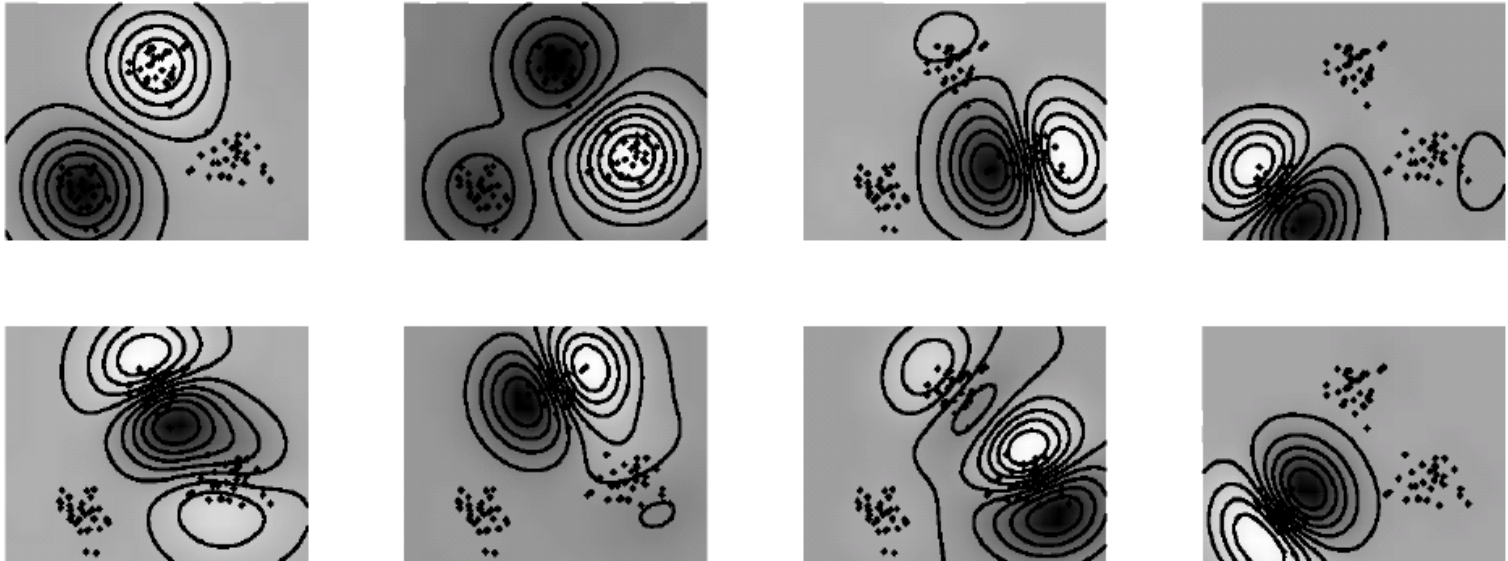
Similar for feature extraction.






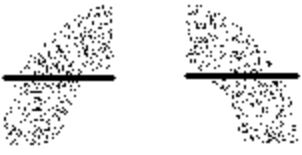
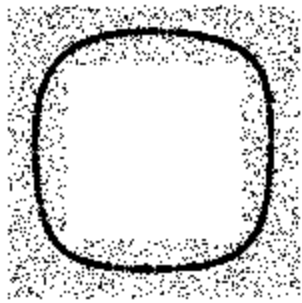
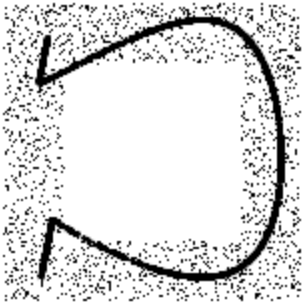
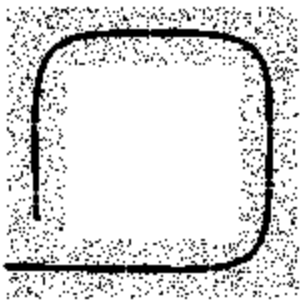
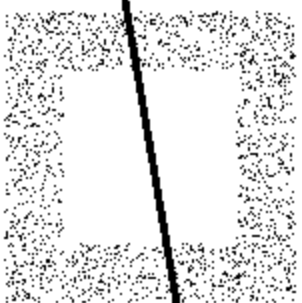
# Example: 8 kPCA components with RBF kernel

---

$$k(\mathbf{x}, \mathbf{y}) = \exp \left( -\frac{\|\mathbf{x} - \mathbf{y}\|^2}{0.1} \right)$$

























# Denoising

kernel PCA (4 PCs)	nonlinear autoencoder	Principal Curves	linear PCA (1 PC)
			
			

Principal curves: Hastie & Stützle, 1989

Nonlinear autoencoder: e.g. Kramer, 1991

## Denoising II

	Gaussian noise	'speckle' noise
orig.		
noisy		
$n = 1$		
4		
16		
64		
256		
$n = 1$		
4		
16		
64		
256	