

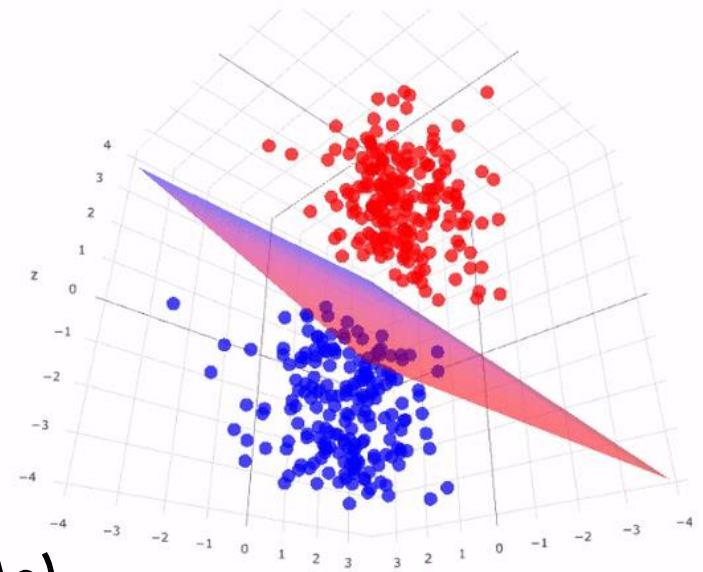
CHAPTER 13:

KERNEL MACHINES

Stella Grasshof

Overview

- Revisit
- Kernel Machines:
Support Vector Machines (SVMs)
- Supervised classification of:
(non)linearly separable classes
- Regression
- Unsupervised learning: one-class classification
- Kernel PCA



Reminder: Multivariate Derivatives

- A function with vector input and scalar output:

$$f : \mathbb{R}^2 \rightarrow \mathbb{R}, \mathbf{v} = (v_1, v_2)^T$$

$$f(\mathbf{v}) = \mathbf{v}^T \mathbf{v} + 2 = v_1^2 + v_2^2 + 2$$

- Its gradient is defined by its partial derivatives as

$$\nabla f : \mathbb{R}^2 \rightarrow \mathbb{R}^2$$

$$\nabla f(\mathbf{v}) = \begin{pmatrix} \frac{\partial}{\partial v_1} f(\mathbf{v}) \\ \frac{\partial}{\partial v_2} f(\mathbf{v}) \end{pmatrix} = \begin{pmatrix} 2v_1 \\ 2v_2 \end{pmatrix} = 2\mathbf{v}$$

Reminder: Multivariate Derivatives

- Compute the gradient of a function

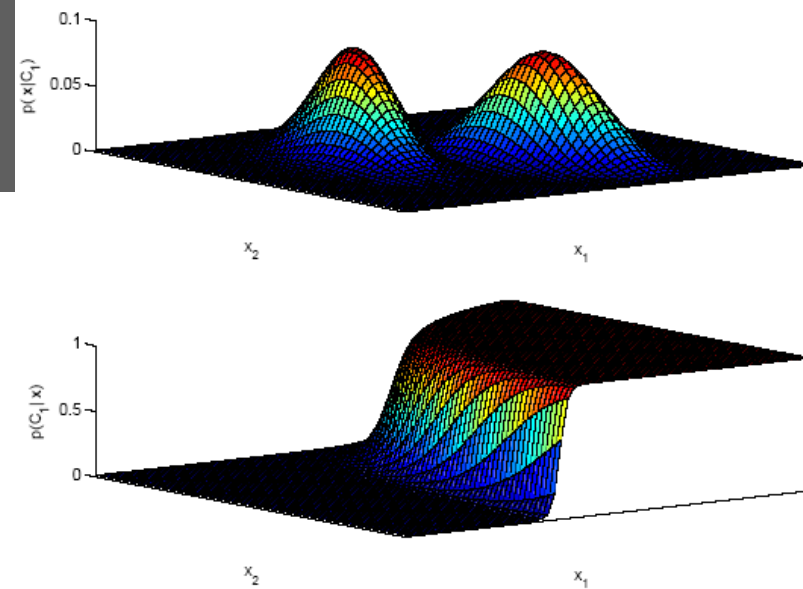
$$g(\mathbf{v}) = \sum_{i=1}^2 \sum_{j=1}^2 v_i v_j = v_1^2 + v_2^2 + 2v_1 v_2$$

$$\begin{aligned}\nabla g(\mathbf{v}) &= \begin{pmatrix} 2v_1 + 2v_2 \\ 2v_2 + 2v_1 \end{pmatrix} = 2 \begin{pmatrix} v_1 + v_2 \\ v_1 + v_2 \end{pmatrix} \\ &= 2 \begin{pmatrix} 1 & 1 \\ 1 & 1 \end{pmatrix} \begin{pmatrix} v_1 \\ v_2 \end{pmatrix}\end{aligned}$$

- Example of rewriting (helpful for exercise)

$$\begin{pmatrix} v_1 + 4v_2 \\ 2v_1 + 3v_2 \end{pmatrix} = \begin{pmatrix} 1 & 4 \\ 2 & 3 \end{pmatrix} \begin{pmatrix} v_1 \\ v_2 \end{pmatrix}$$

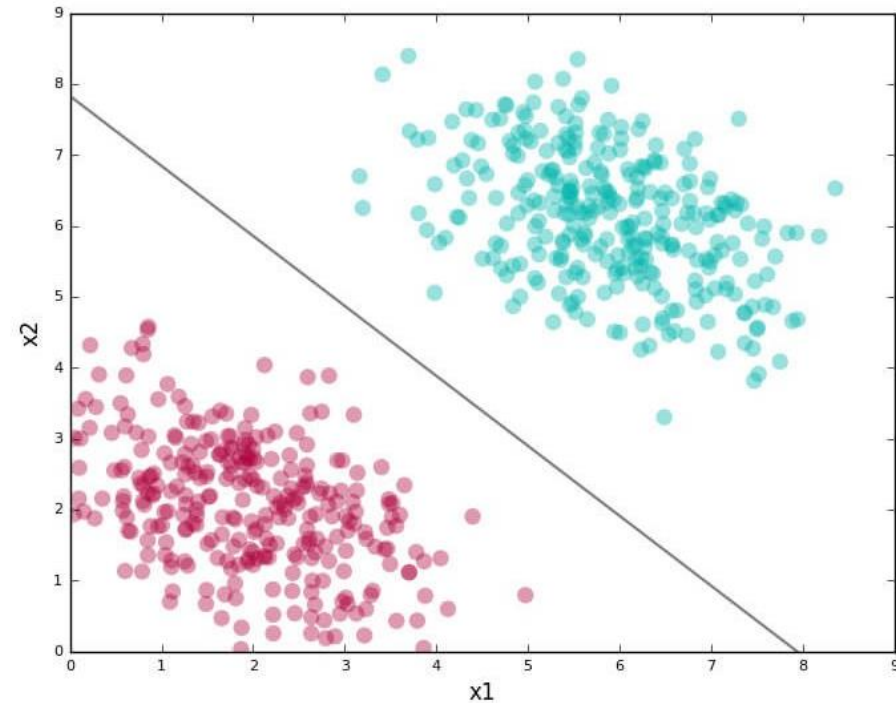
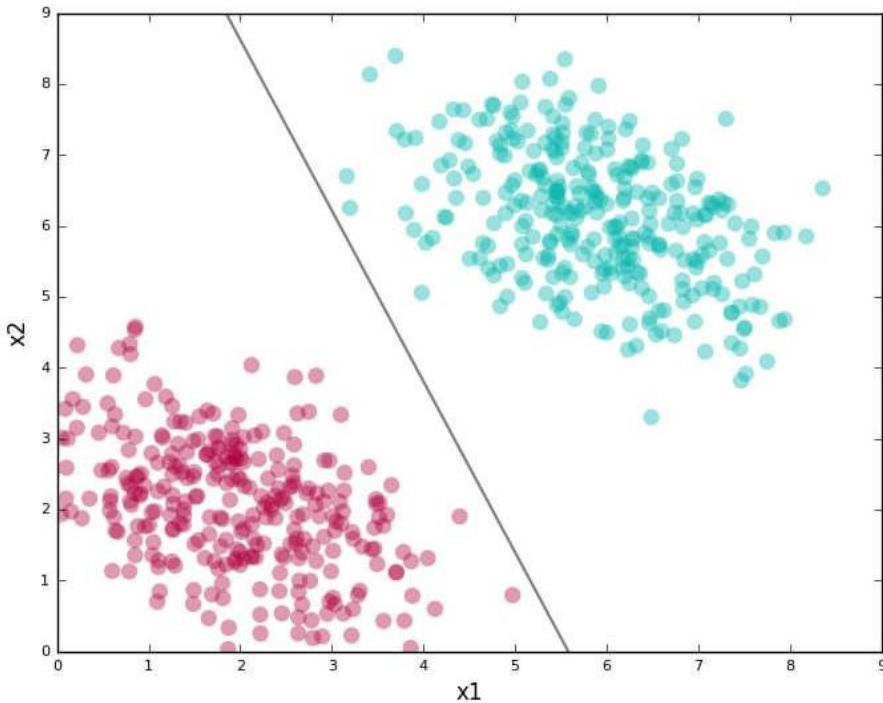
Kernel Machines



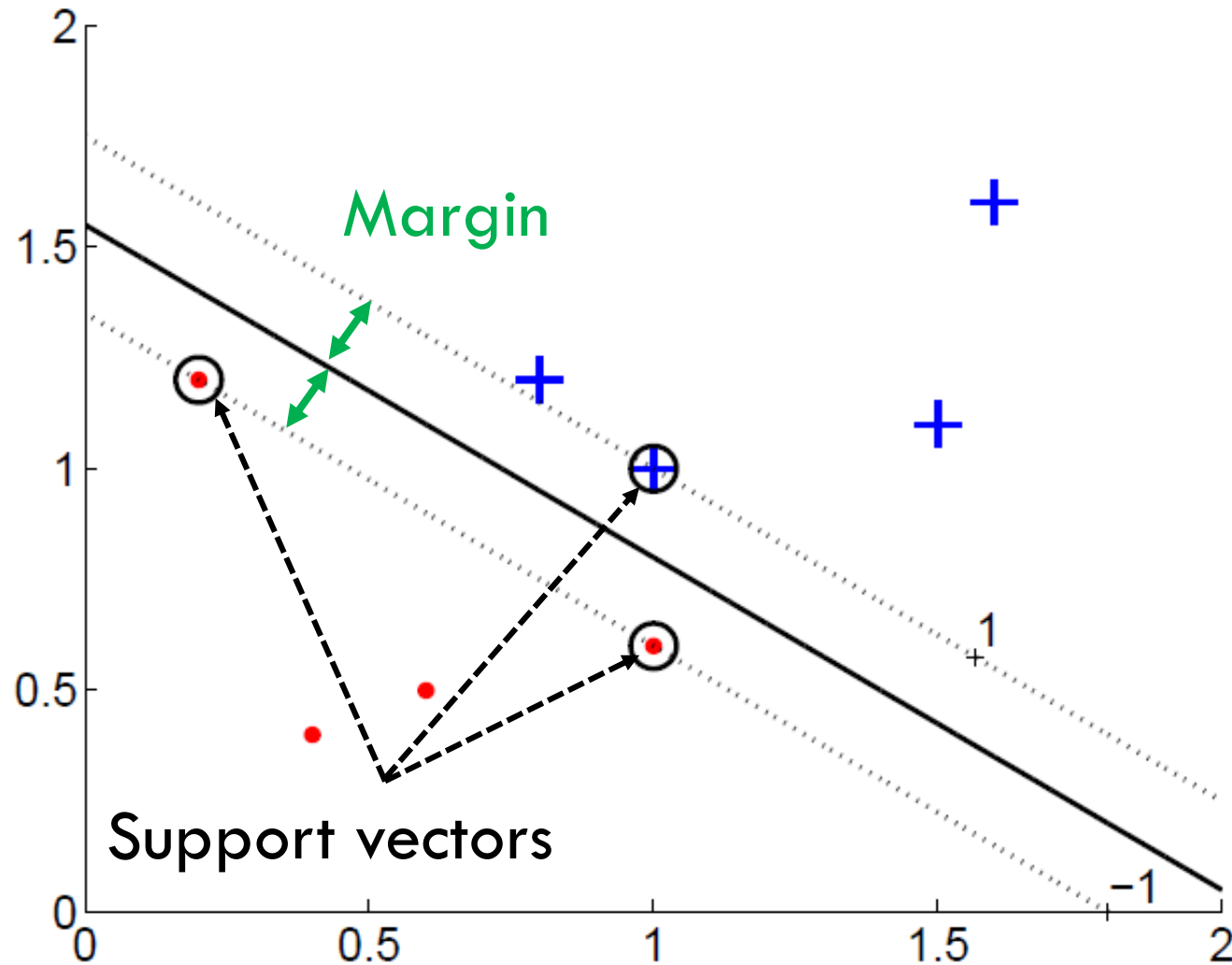
- Discriminant-based:
previously defined by densities,
here:
we want few training points aka **support vectors**
- Kernel functions:
application-specific measures of similarity
- No need to represent instances as vectors
- Convex optimization problems with a unique solution

Separating two classes by one line

Which is better?



Support Vector Machines (SVMs)



Support Vector Machines (SVMs)

Support Vector Machines (SVMs)

find the **optimal separating hyperplane** which:

- Separates the two classes
- Maximizes the margin between them, i.e.
maximize distance between plane and the points
closest to it
- Hyperplane is defined by few training samples

Repeat: discriminant function

Classification:

- Discriminant function require the probability density function (pdf):

$$g_i(\mathbf{x}) = \log p(\mathbf{x}|C_i) + \log P(C_i)$$

$$g_k(\mathbf{x}) = \max_i g_i(\mathbf{x})$$

- Separates the two classes

$$g_1(\mathbf{x}_i) = \mathbf{w}_1^T \mathbf{x}_i + w_{10}$$

$$g_2(\mathbf{x}_i) = \mathbf{w}_2^T \mathbf{x}_i + w_{20}$$

$$\Rightarrow g(\mathbf{x}_i) = g_1(\mathbf{x}_i) - g_2(\mathbf{x}_i) = (\mathbf{w}_1 - \mathbf{w}_2)^T \mathbf{x}_i + (w_{10} - w_{20})$$

$$g(\mathbf{x}_i) = \mathbf{w}^T \mathbf{x}_i + w_0 \begin{cases} > 0 & , \text{ class } C_1 \\ \leq 0 & , \text{ class } C_2 \end{cases}$$

Optimal Separating Hyperplane

- Given points $x_i \in \mathbb{R}^D$, $i = 1, \dots, N$ of two classes with labels:

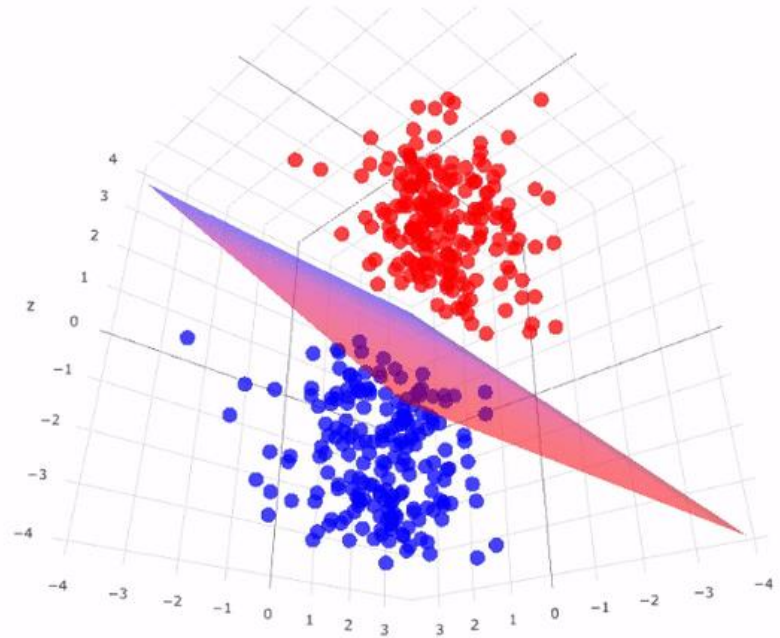
$$r_i = \begin{cases} +1, & \text{if } x_i \in C_1 \\ -1, & \text{if } x_i \in C_2 \end{cases}$$

- Points on the hyperplane:

$$w^T x + w_0 = 0$$

- Class decision as previously

$$w^T x_i + w_0 \begin{cases} > 0, & \text{class } C_1 : r_i = +1 \\ \leq 0, & \text{class } C_2 : r_i = -1 \end{cases}$$



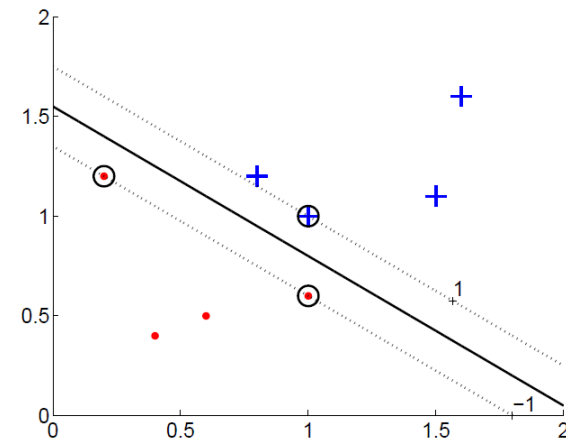
Optimal Separating Hyperplane

- Previous class decision

$$\mathbf{w}^T \mathbf{x}_i + w_0 \begin{cases} > 0, & \text{class } C_1 : r_i = +1 \\ \leq 0, & \text{class } C_2 : r_i = -1 \end{cases}$$

- Now: we want points to be some distance away of the hyperplane

$$\mathbf{w}^T \mathbf{x}_i + w_0 \begin{cases} \geq 1, & r_i = +1 \\ \leq -1, & r_i = -1 \end{cases}$$



This can be rewritten

$$r_i(\mathbf{w}^T \mathbf{x}_i + w_0) \geq 1$$

Optimal Separating Hyperplane

- Data points $x_i \in \mathbb{R}^D$ with labels $r_i \in \{-1, 1\}$

then decision boundary is

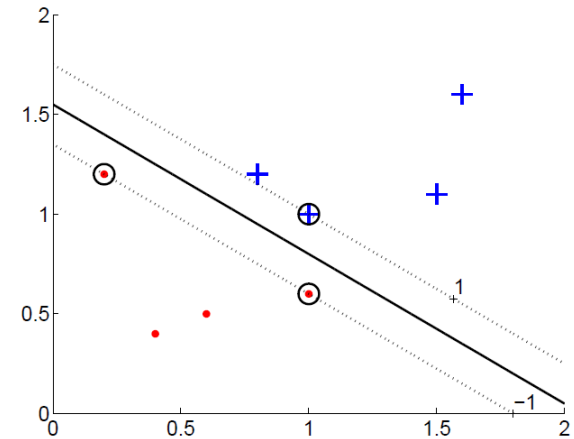
$$r_i(\mathbf{w}^T \mathbf{x}_i + w_0) \geq 1$$

unknown: \mathbf{w} , w_0

- Distance of points to the hyperplane should be large:

$$\max \frac{|\mathbf{w}^T \mathbf{x}_i + w_0|}{\|\mathbf{w}\|_2} \geq \rho, \quad \forall i$$

$$\max \frac{r_i(\mathbf{w}^T \mathbf{x}_i + w_0)}{\|\mathbf{w}\|_2} \geq \rho, \quad \forall i$$



Optimal Separating Hyperplane

- Data points $x_i \in \mathbb{R}^D$ with labels $r_i \in \{-1, 1\}$

then decision boundary is

$$r_i(w^T x_i + w_0) \geq 1$$

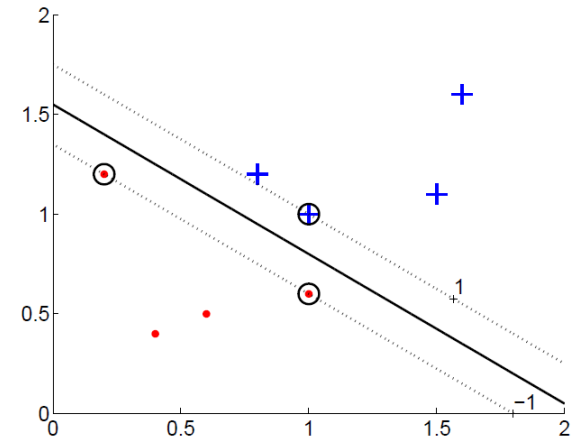
unknown: w , w_0

- Distance of points to the hyperplane should be large:

$$\max \frac{\overset{\text{maximize}}{r_i(w^T x_i + w_0)}}{\underset{\text{minimize}}{\|w\|_2}} \geq \rho, \quad \forall i$$

- Fix $\|w\|_2 \rho = 1$ then:

$$\min \|w\|_2^2, \text{ subject to } r_i(w^T x_i + w_0) \geq 1, \quad \forall i$$



Optimal Separating Hyperplane

- Data points $x_i \in \mathbb{R}^D$ with labels $r_i \in \{-1, 1\}$

$$r_i(\mathbf{w}^T \mathbf{x}_i + w_0) \geq 1$$

unknown: \mathbf{w} , w_0

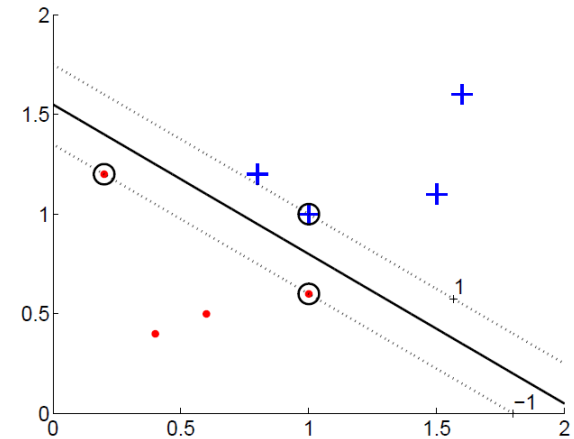
- Distance of points to the hyperplane should be large:

$$\max \frac{\overset{\text{maximize}}{r_i(\mathbf{w}^T \mathbf{x}_i + w_0)}}{\underset{\text{minimize}}{\|\mathbf{w}\|_2}} \geq \rho, \quad \forall i$$

- Lagrange function

$$L(\mathbf{w}, w_0, \alpha_i) = \frac{1}{2} \|\mathbf{w}\|_2^2 - \sum_{i=1}^N \alpha_i (r_i(\mathbf{w}^T \mathbf{x}_i + w_0) - 1)$$

$$= \frac{1}{2} \|\mathbf{w}\|_2^2 - \sum_{i=1}^N \alpha_i r_i(\mathbf{w}^T \mathbf{x}_i + w_0) + \sum_{i=1}^N \alpha_i, \quad \alpha_i \geq 0$$



Optimal Separating Hyperplane

- Data points $\mathbf{x}_i \in \mathbb{R}^D$ with labels $r_i \in \{-1, 1\}$

$$r_i(\mathbf{w}^T \mathbf{x}_i + w_0) \geq 1$$

unknown: \mathbf{w} , w_0

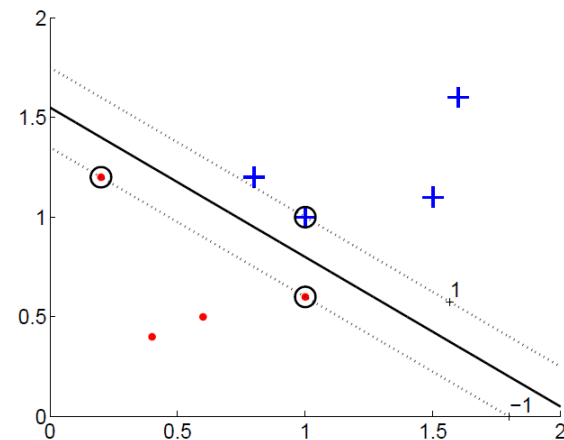
- Lagrange function

$$L(\mathbf{w}, w_0, \alpha_i) = \frac{1}{2} \|\mathbf{w}\|_2^2 - \sum_{i=1}^N \alpha_i r_i (\mathbf{w}^T \mathbf{x}_i + w_0) + \sum_{i=1}^N \alpha_i$$

- Minimize L for \mathbf{w} , w_0

$$\nabla_{\mathbf{w}} L(\mathbf{w}, w_0, \alpha_i) = \mathbf{w} - \sum_{i=1}^N \alpha_i r_i \mathbf{x}_i \stackrel{!}{=} 0$$

$$\frac{\partial}{\partial w} L(\mathbf{w}, w_0, \alpha_i) = - \sum_{i=1}^N \alpha_i r_i \stackrel{!}{=} 0$$



Optimal Separating Hyperplane

- Data points $x_i \in \mathbb{R}^D$ with labels $r_i \in \{-1, 1\}$

$$r_i(w^T x_i + w_0) \geq 1$$

unknown: w, w_0

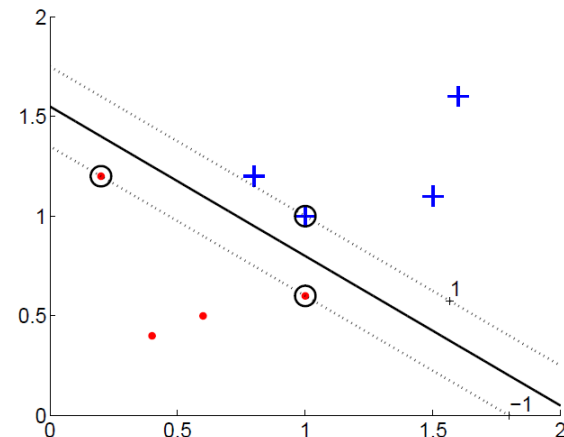
- Lagrange function

$$L(w, w_0, \alpha_i) = \frac{1}{2} \|w\|_2^2 - \sum_{i=1}^N \alpha_i r_i (w^T x_i + w_0) + \sum_{i=1}^N \alpha_i$$

- Minimize L for w, w_0

$$w = \sum_{i=1}^N \alpha_i r_i x_i$$

$$\sum_{i=1}^N \alpha_i r_i = 0$$



Optimal Separating Hyperplane

- Data points $\mathbf{x}_i \in \mathbb{R}^D$ with labels $r_i \in \{-1, 1\}$

$$r_i(\mathbf{w}^T \mathbf{x}_i + w_0) \geq 1$$

unknown: \mathbf{w} , w_0

- Lagrange function

$$\mathbf{w} = \sum_{i=1}^N \alpha_i r_i \mathbf{x}_i$$

$$\sum_{i=1}^N \alpha_i r_i = 0$$

$$\begin{aligned} L(\mathbf{w}, w_0, \alpha_i) &= \frac{1}{2} \|\mathbf{w}\|_2^2 - \sum_{i=1}^N \alpha_i r_i (\mathbf{w}^T \mathbf{x}_i + w_0) + \sum_{i=1}^N \alpha_i \\ &= \frac{1}{2} \mathbf{w}^T \mathbf{w} - \mathbf{w}^T \sum_{i=1}^N \alpha_i r_i \mathbf{x}_i - w_0 \sum_{i=1}^N \alpha_i r_i + \sum_{i=1}^N \alpha_i \\ &= -\frac{1}{2} \sum_{i=1}^N \sum_{j=1}^N \alpha_i \alpha_j r_i r_j \mathbf{x}_i^T \mathbf{x}_j + \sum_{i=1}^N \alpha_i \end{aligned}$$

Optimal Separating Hyperplane

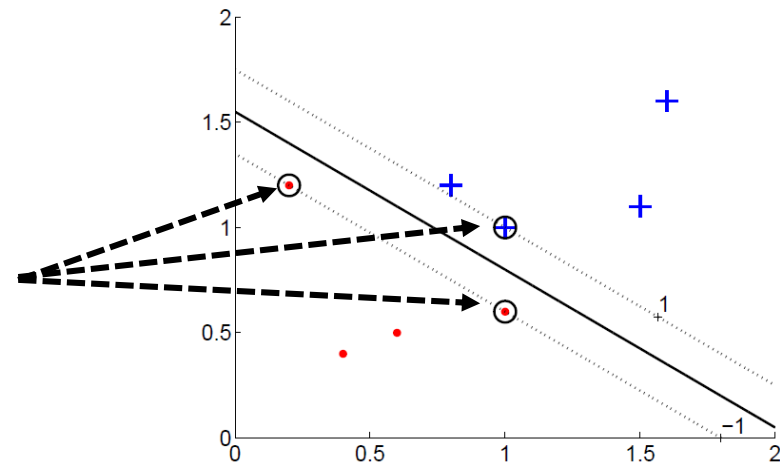
- Data points $\mathbf{x}_i \in \mathbb{R}^D$ with labels $r_i \in \{-1, 1\}$

$$r_i(\mathbf{w}^T \mathbf{x}_i + w_0) \geq 1$$

- **Known:** $\mathbf{w} = \sum_{i=1}^N \alpha_i r_i \mathbf{x}_i$

and if \mathbf{x}_k **support vector**, then

$$w_0 = r_k - \mathbf{w}^T \mathbf{x}_k$$



- **unknown:** $\alpha_i > 0$ are **support vectors**

- Max Lagrange function for α_i , with: $\sum_{i=1}^N \alpha_i r_i = 0, \quad \alpha_i \geq 0$

$$L(\mathbf{w}, w_0, \alpha_i) = -\frac{1}{2} \sum_{i=1}^N \sum_{j=1}^N \alpha_i \alpha_j r_i r_j \mathbf{x}_i^T \mathbf{x}_j + \sum_{i=1}^N \alpha_i$$

Sequential Minimal Optimization (SMO)

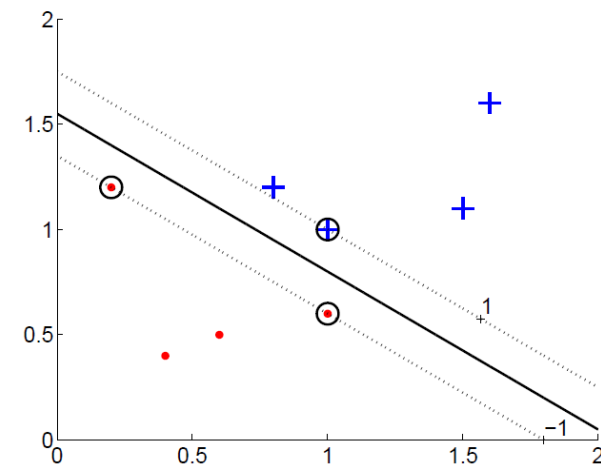
$$\max_{\alpha_i} L(\mathbf{w}, w_0, \alpha_i) = -\frac{1}{2} \sum_{i=1}^N \sum_{j=1}^N \alpha_i \alpha_j r_i r_j \mathbf{x}_i^T \mathbf{x}_j + \sum_{i=1}^N \alpha_i$$
$$\sum_{i=1}^N \alpha_i r_i = 0, \quad 0 \leq \alpha_i \leq C$$

How to estimate the Lagrangian Multipliers α_i ?

By: **SMO**

Idea: estimate two: α_k, α_l , keep the rest fixed

- 1) Find α_k that violates conditions
- 2) Pick a second multiplier α_l and optimize the pair α_k, α_l
- 3) Repeat 1) and 2)

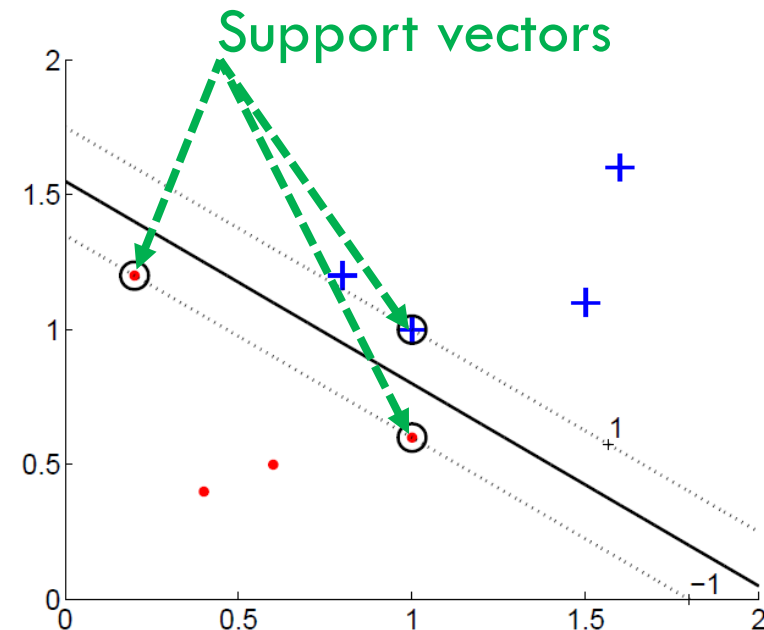
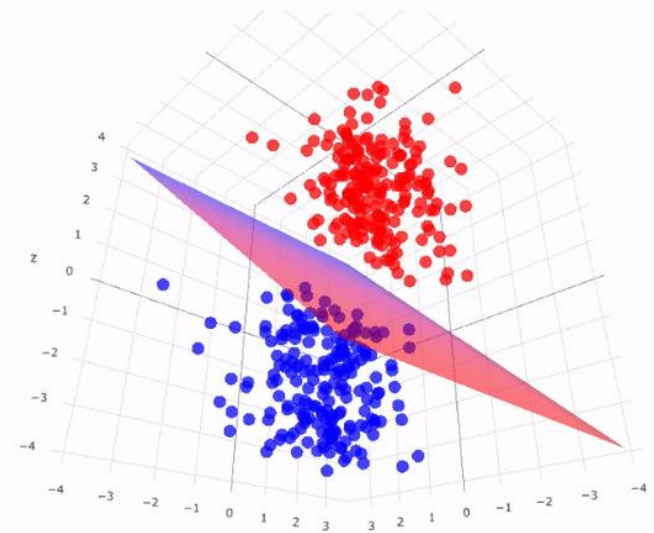


Summary SVM

- Data points $x_i \in \mathbb{R}^D$ with labels $r_i \in \{-1, 1\}$
- SVM:
 - (1) find support vectors, defined by $\alpha_i > 0$ which then define w, w_0
 - (2) the parameters needed to classify new points

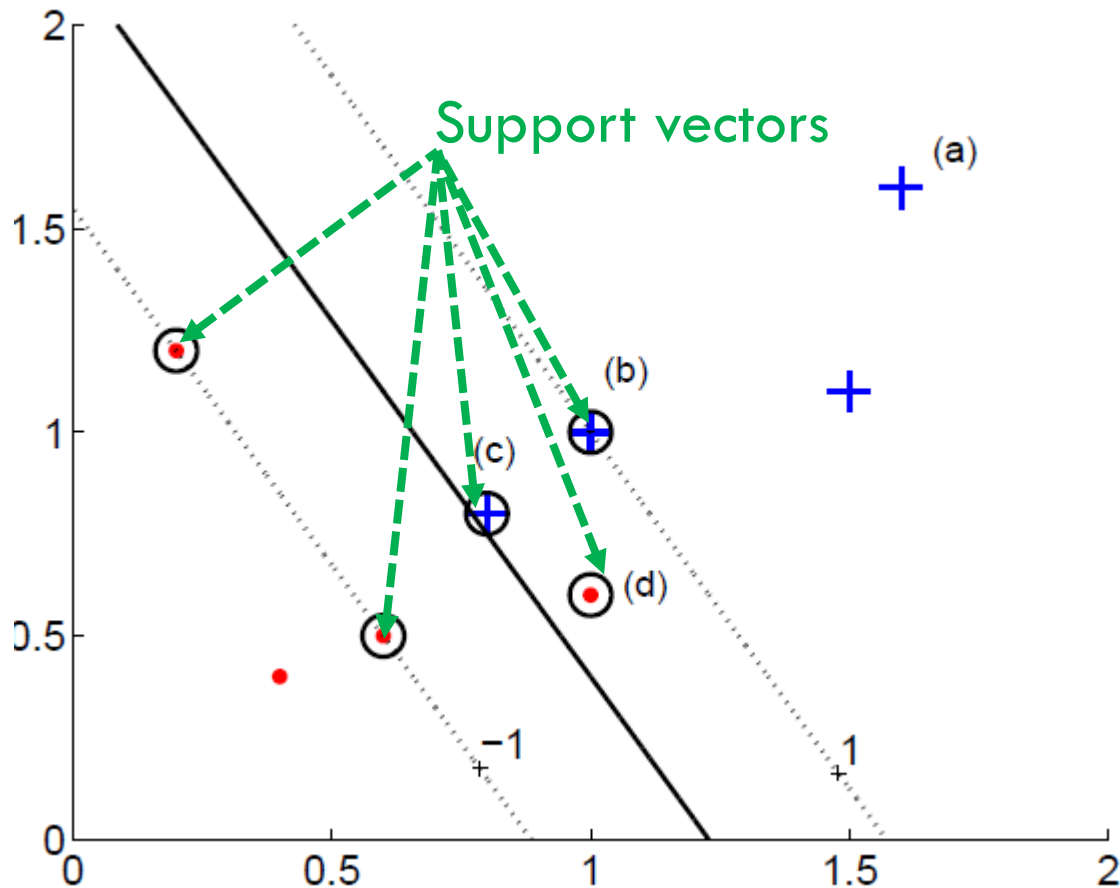
$$w^T x + w_0 \begin{cases} \geq 1, & r_i = +1 \\ \leq -1, & r_i = -1 \end{cases}$$

- Problem:
 - Assumption linearly separable
 - Two classes only



Soft Margin Hyperplane

There is no Separating Hyperplane



Soft Margin Hyperplane

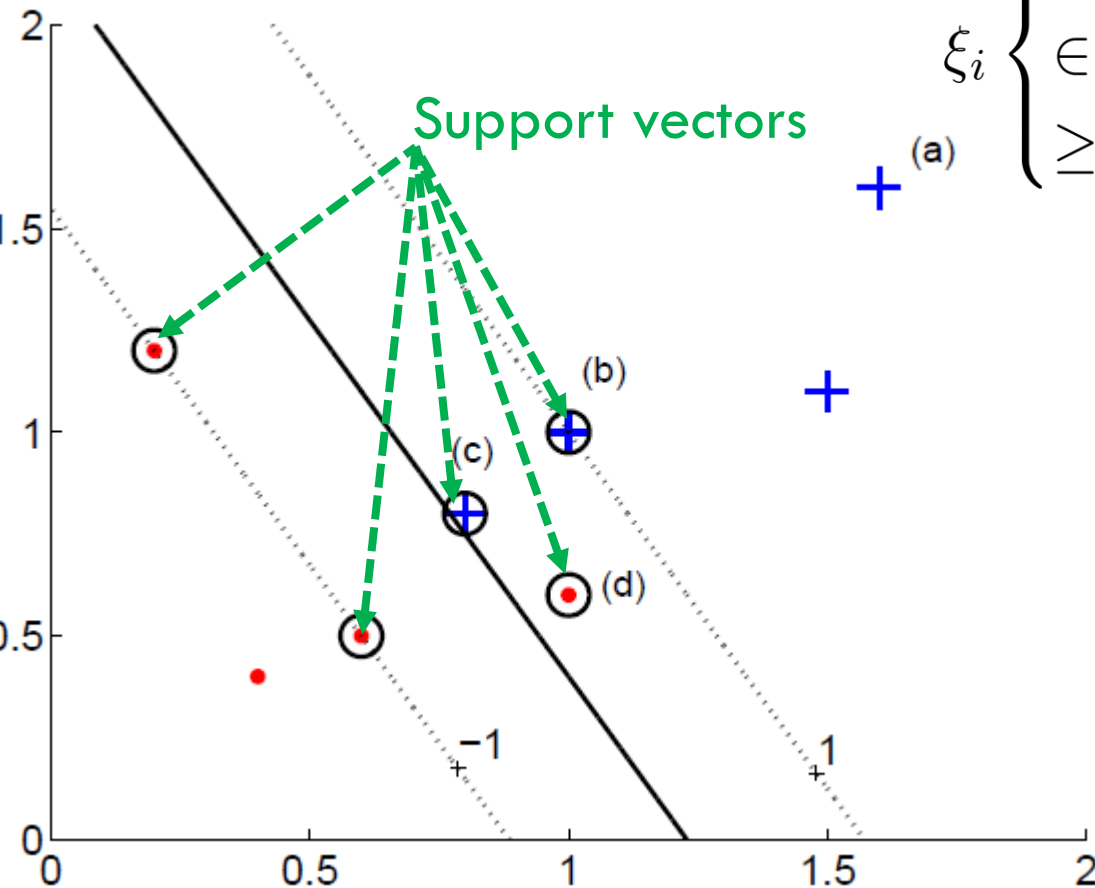
□ Linearly separable

$$r_i(\mathbf{w}^T \mathbf{x} + w_0) \geq 1$$

□ Not linearly separable

$$r_i(\mathbf{w}^T \mathbf{x} + w_0) \geq 1 - \xi$$

$$\xi_i \begin{cases} = 0 & , \mathbf{x}_i \text{ ok (a), (b)} \\ \in (0, 1) & , \mathbf{x}_i \text{ in margin (c)} \\ \geq 1 & , \mathbf{x}_i \text{ falsely classified (d)} \end{cases}$$



Soft Margin Hyperplane

- Linearly separable $r_i(\mathbf{w}^T \mathbf{x}_i + w_0) \geq 1$
- Not linearly separable $r_i(\mathbf{w}^T \mathbf{x}_i + w_0) \geq 1 - \xi_i$
- Soft error $\sum_{i=1}^N \xi_i$ $\xi_i \begin{cases} = 0 & , \mathbf{x}_i \text{ ok} \\ \in (0, 1) & , \mathbf{x}_i \text{ in margin} \\ \geq 1 & , \mathbf{x}_i \text{ falsely classified} \end{cases}$

- New Optimization function, where C is must be set

$$L(\mathbf{w}, w_0, \alpha_i, \beta_i, \xi_i) = \frac{1}{2} \|\mathbf{w}\|_2^2 - \sum_{i=1}^N \alpha_i r_i ((\mathbf{w}^T \mathbf{x}_i + w_0) - 1 + \xi_i) \\ + C \sum_i \xi_i - \sum_{i=1}^N \beta_i \xi_i$$

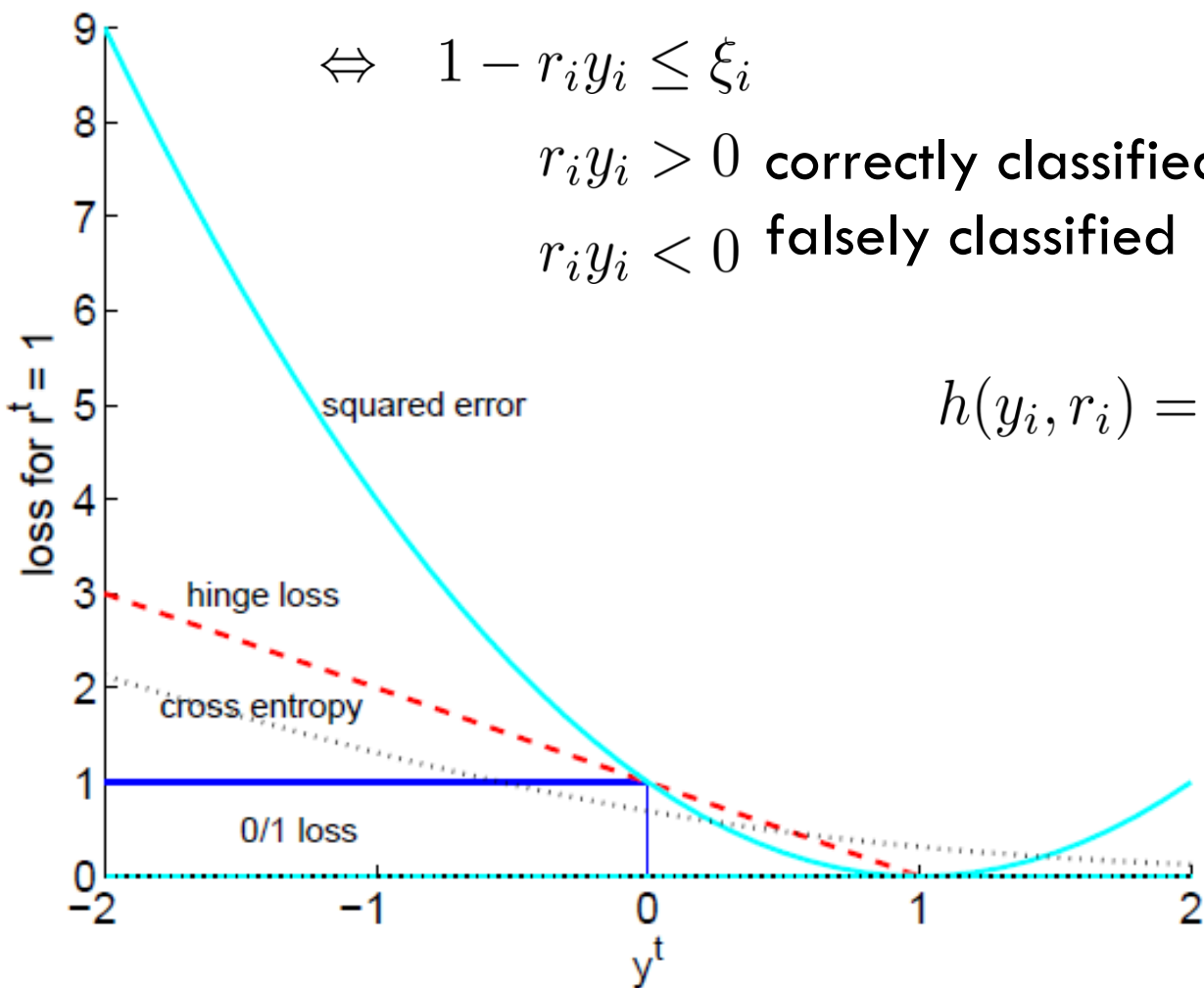
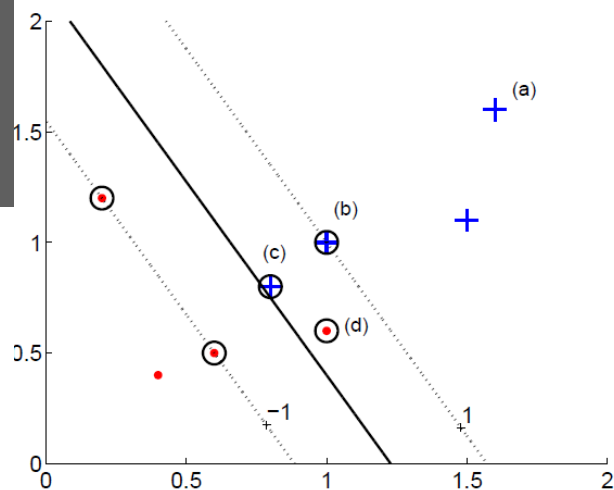
Hinge Loss

$$r_i \underbrace{(w^T x + w_0)}_{=y_i} \geq 1 - \xi_i$$

$$\Leftrightarrow 1 - r_i y_i \leq \xi_i$$

$r_i y_i > 0$ correctly classified

$r_i y_i < 0$ falsely classified



$$h(y_i, r_i) = \begin{cases} 0, & \text{if } y_i r_i \geq 1 \\ 1 - y_i r_i, & \text{otherwise} \end{cases}$$

Hinge loss is more robust than square error

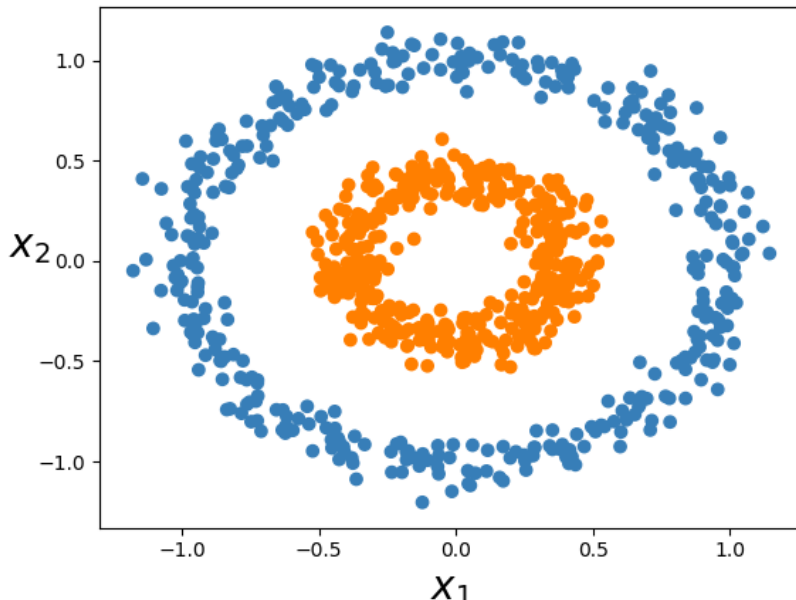
Kernel Trick

For non-linearly separable data:
transfer in higher dimensional space

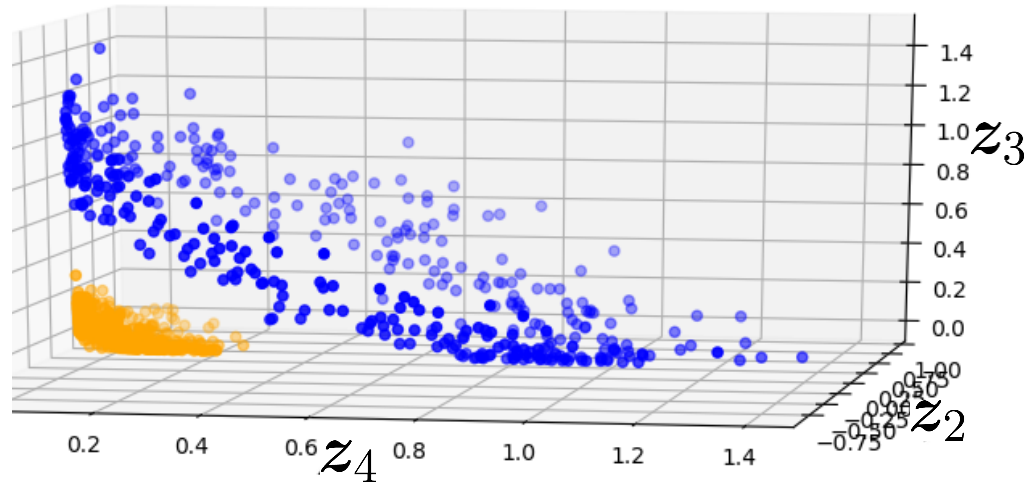
$$\Phi : \mathbb{R}^2 \mapsto \mathbb{R}^4$$

$$z_i = \Phi(x_i) = (1, x_1^2, x_2^2, \sqrt{2}x_1x_2)$$

before



after



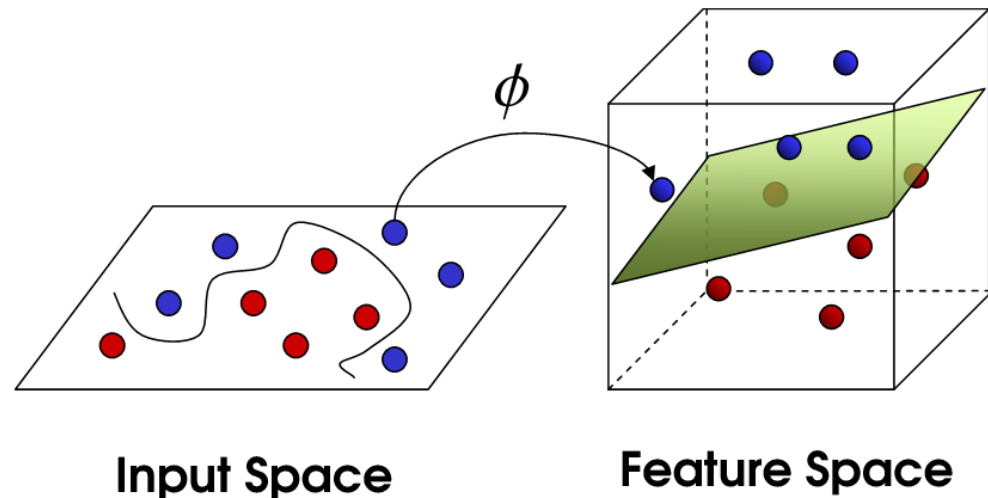
Kernel Trick

- Basis function as preprocessing

$$\Phi : \mathbb{R}^D \rightarrow \mathbb{R}^L, D \ll L$$

$$z_i = \Phi(x_i)$$

$$g(x) = w^T \Phi(x)$$



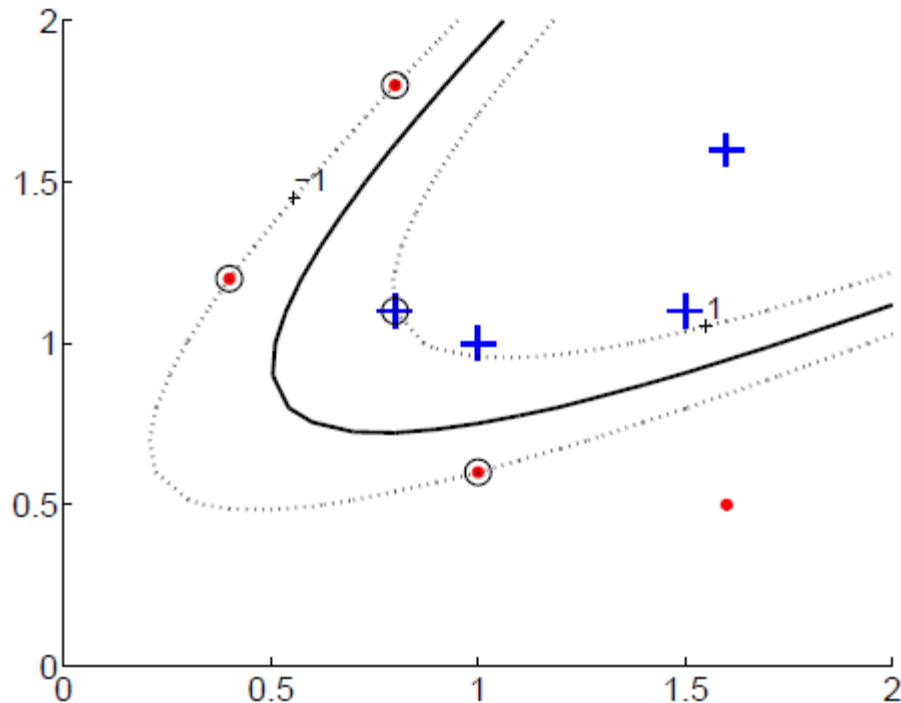
- Kernel function $K(.,.)$

$$g(x) = w^T \Phi(x) = \sum_i \alpha_i r_i \Phi(x_i)^T \Phi(x)$$

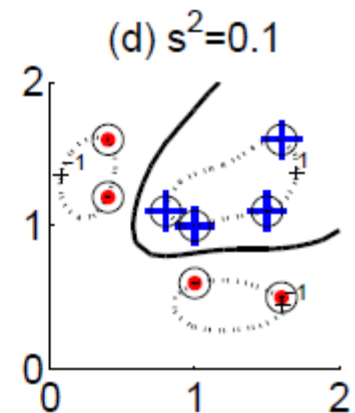
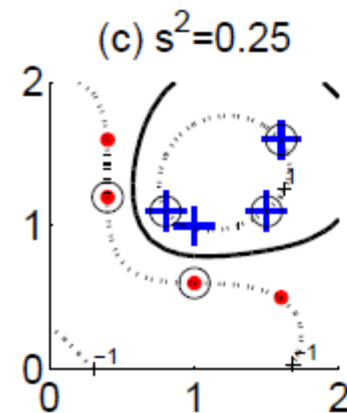
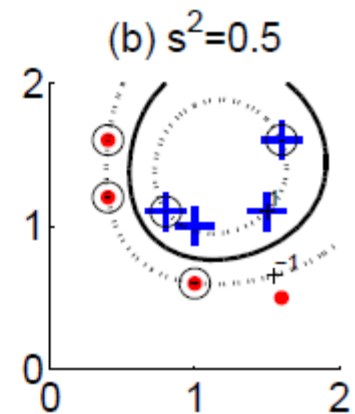
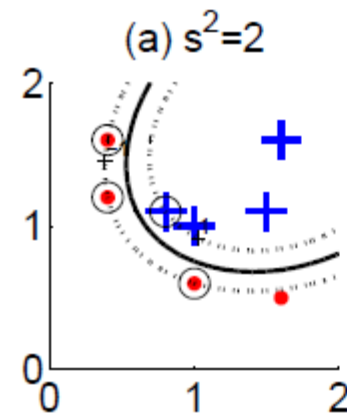
$$= \sum_i \alpha_i r_i K(x_i, x)$$

Kernel: examples

Polynomial



Radial-basis functions

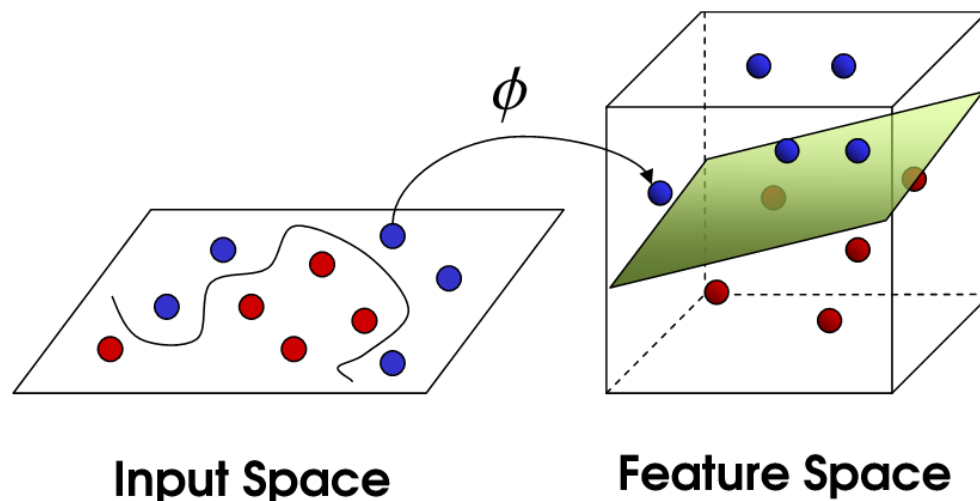


How to define a Kernel

$$\Phi : \mathbb{R}^D \rightarrow \mathbb{R}^L, D \ll L$$

$$g(\mathbf{x}) = \mathbf{w}^T \Phi(\mathbf{x}) = \sum_i \alpha_i r_i \Phi(\mathbf{x}_i)^T \Phi(\mathbf{x}) = \sum_i \alpha_i r_i K(\mathbf{x}_i, \mathbf{x})$$

- Application dependent!
- Kernel function is a measure of similarity:
the more similar the input, the higher:
(e.g. p.d.f.)



How to define a Kernel

$$\Phi : \mathbb{R}^D \rightarrow \mathbb{R}^L, D \ll L$$

$$g(\mathbf{x}) = \mathbf{w}^T \Phi(\mathbf{x}) = \sum_i \alpha_i r_i \Phi(\mathbf{x}_i)^T \Phi(\mathbf{x}) = \sum_i \alpha_i r_i K(\mathbf{x}_i, \mathbf{x})$$

- Kernel function = measure of similarity

- Empirical kernel map:

Define a set of templates \mathbf{m}_i and score function s

$$\Phi(\mathbf{x}) = [s(\mathbf{x}, \mathbf{m}_1), \dots, s(\mathbf{x}, \mathbf{m}_M)]$$

- Combine kernels

$$K(\mathbf{x}, \mathbf{x}_i) = \sum_n \lambda_n K_n(\mathbf{x}, \mathbf{x}_i)$$

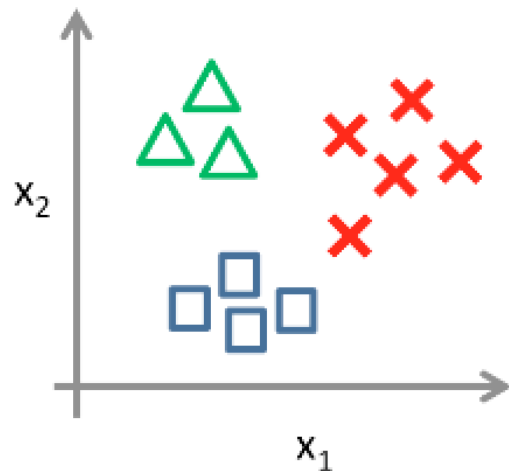
- Combination of two different representations,
e.g. image and sound




Multiclass Kernel Machines

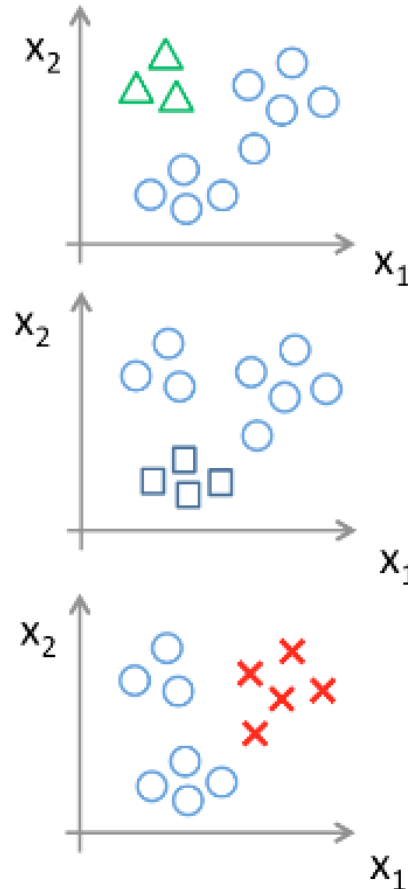
Consider K classes

- One-vs-all: K subproblems

One-vs-all (one-vs-rest):



Class 1: 
Class 2: 
Class 3: 

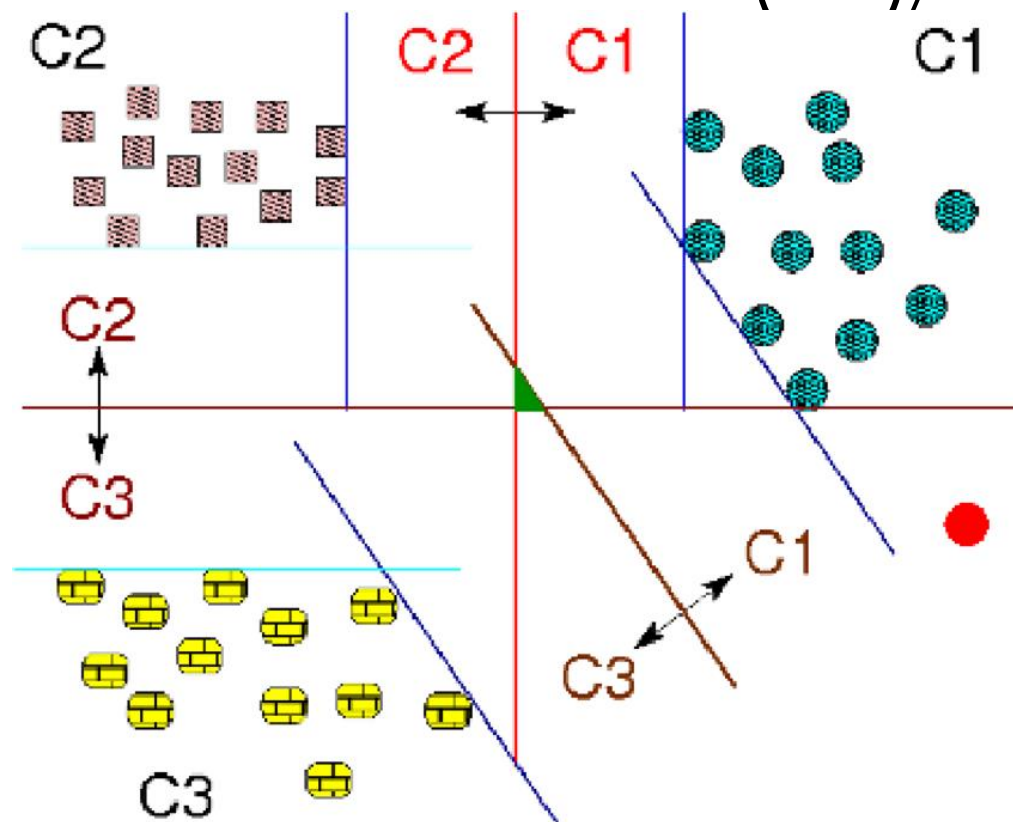


Multiclass Kernel Machines

Consider K classes

□ Pairwise separation (10.4.):

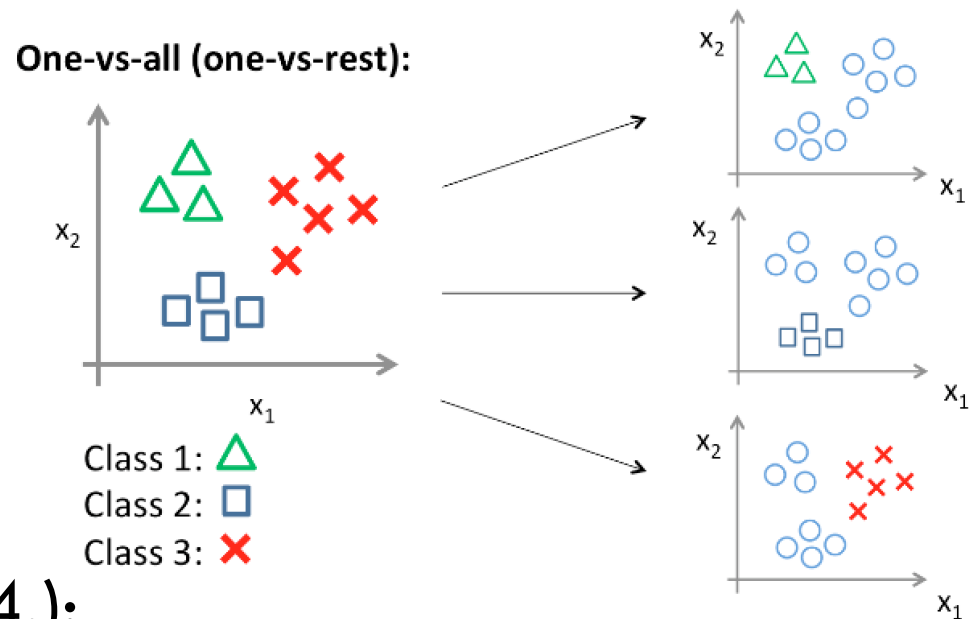
focus on two classes at a time: $K(K-1)/2$



Multiclass Kernel Machines

Consider K classes

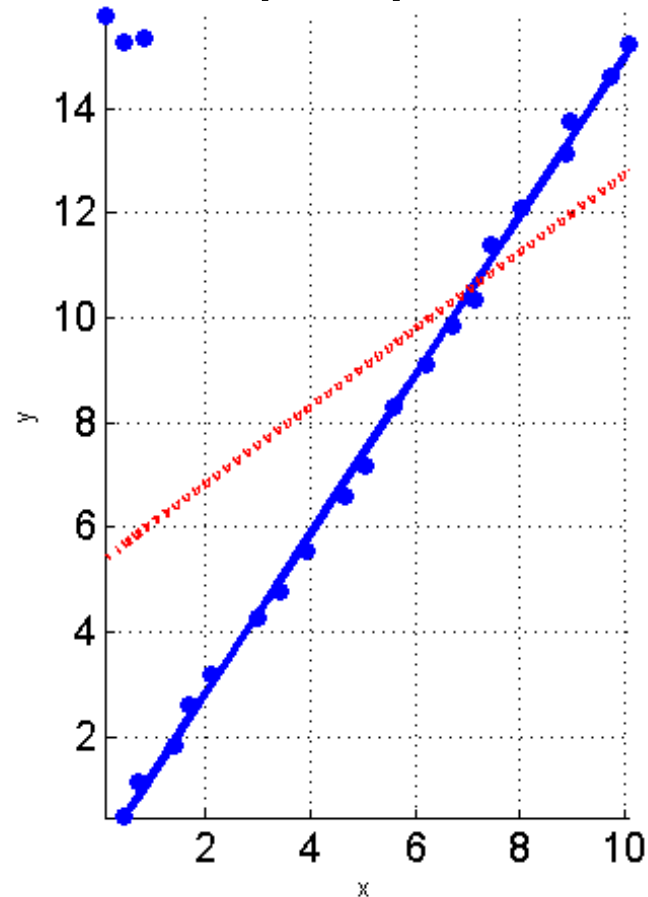
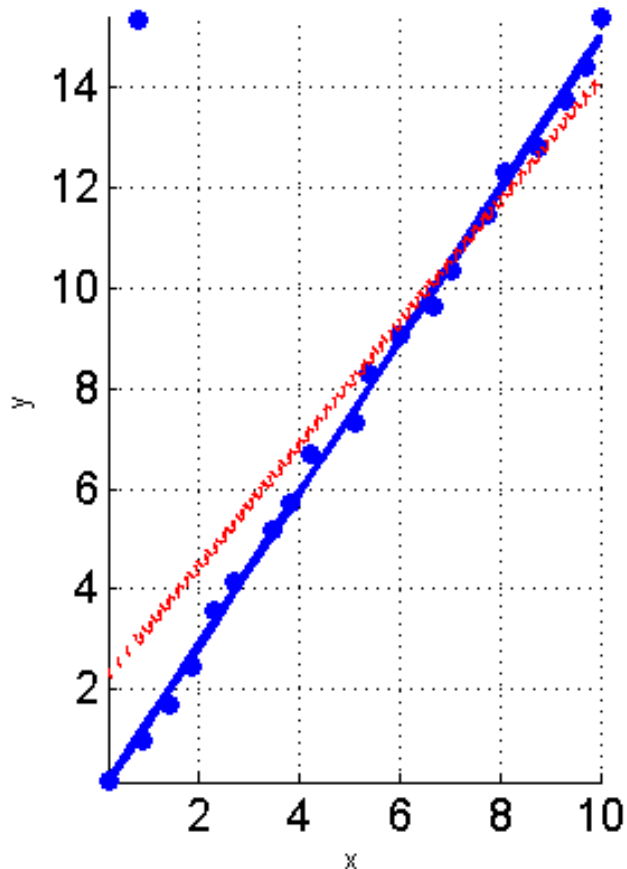
- 1-vs-all
define K problems
- Pairwise separation (10.4.):
focus on two classes at a time: $K(K-1)/2$
- Error-Correcting Output Codes (Sec. 17.5)
consider sets of 2 classes
- Single multiclass optimization
involves all classes



Regression

Fit a line / hyperplane through the data
previous problem:

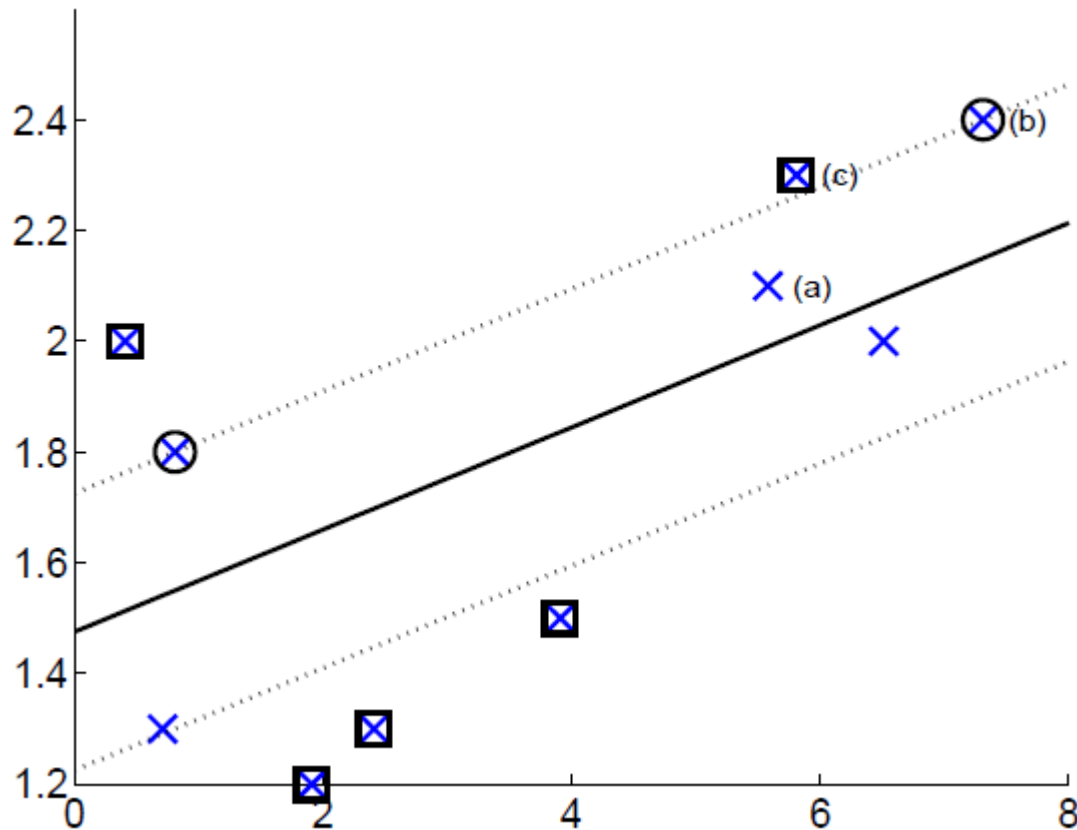
outliers, because all points contribute equally



SVM for Regression

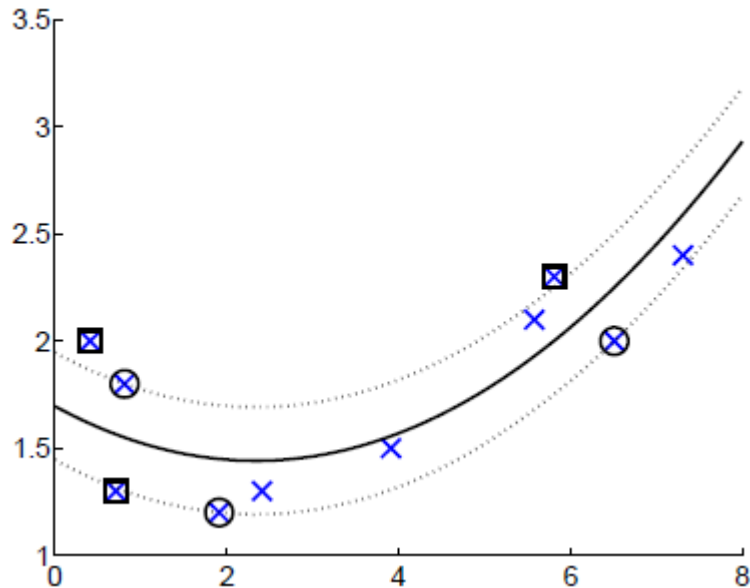
Fitted line is weighted sum of support vectors: $\square \circ$

$$f(\mathbf{x}) = \mathbf{w}^T \mathbf{x} + w_0 = \sum_i (\alpha_{i+} - \alpha_{i-}) \mathbf{x}_i^T \mathbf{x} + w_0$$



Kernel Regression

□ Polynomial kernel

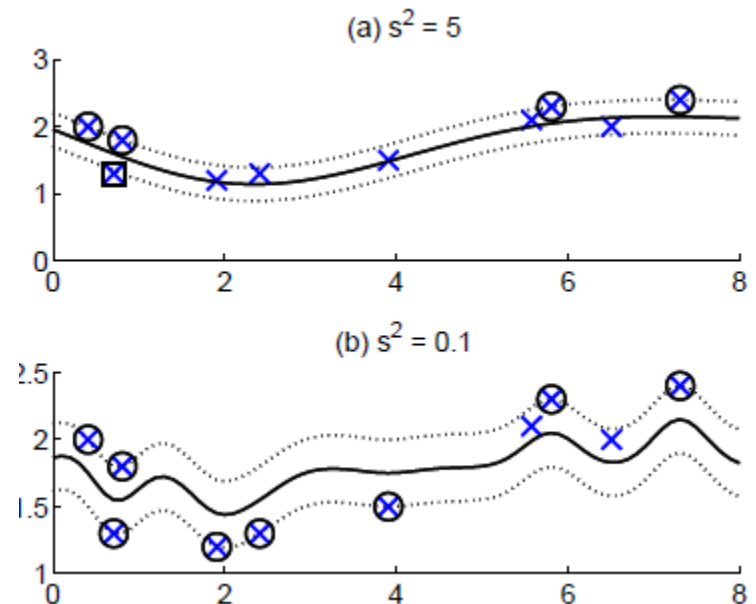


support vectors

○ inside of tube

□ outside of tube

□ Gaussian kernel



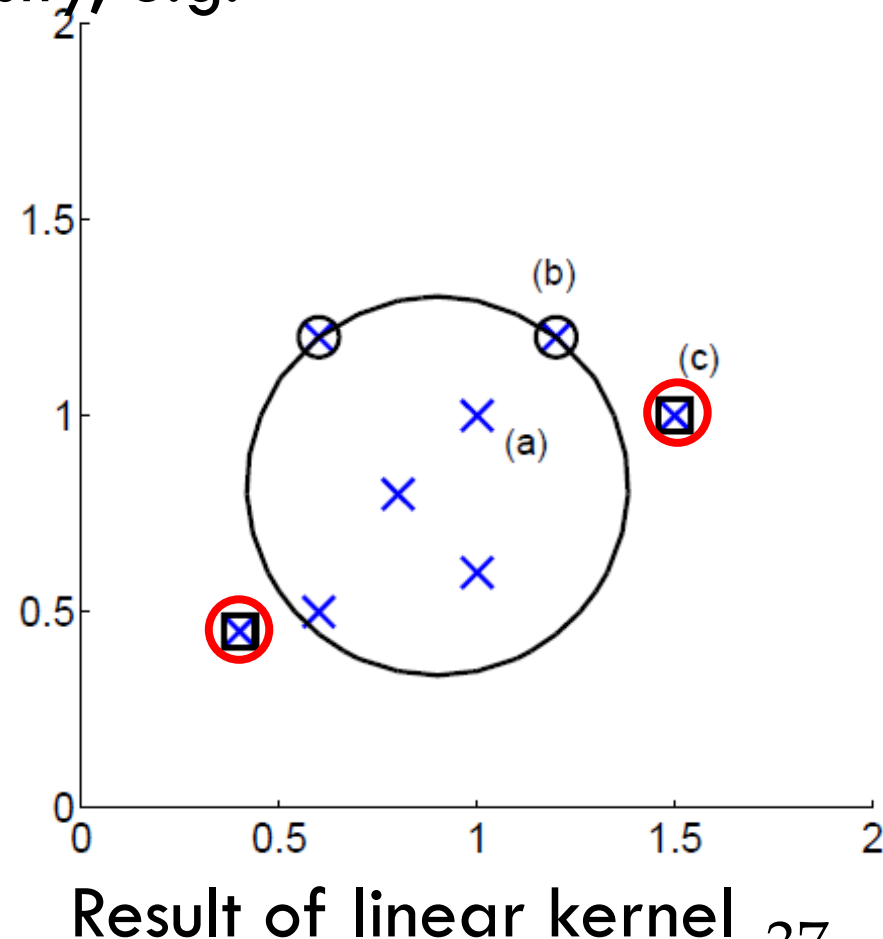
One-Class Kernel Machines

Unsupervised learning:

- Estimate regions of high density, e.g. to define outliers
- Smoothest boundary, enclosing as many points as possible
- Consider a sphere with center \mathbf{a} and radius R

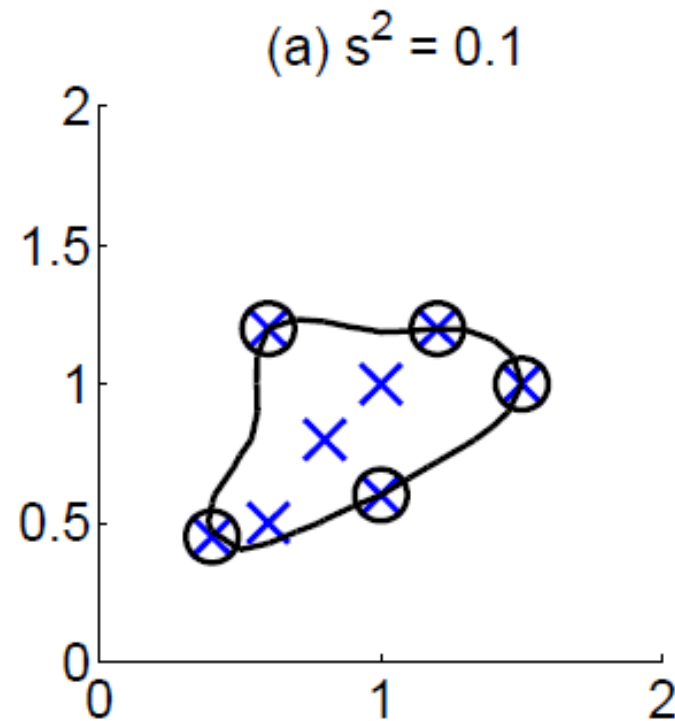
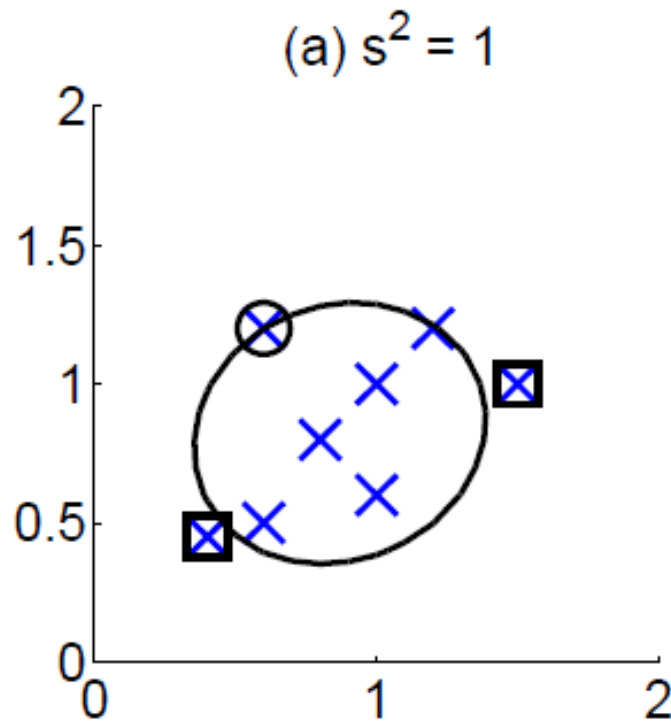
$$\min R^2 + C \sum_i \xi_i$$

$$\|\mathbf{x}_i - \mathbf{a}\|_2^2 \leq R^2 + \xi_i$$



One-Class Kernel Machines

Gaussian Kernel



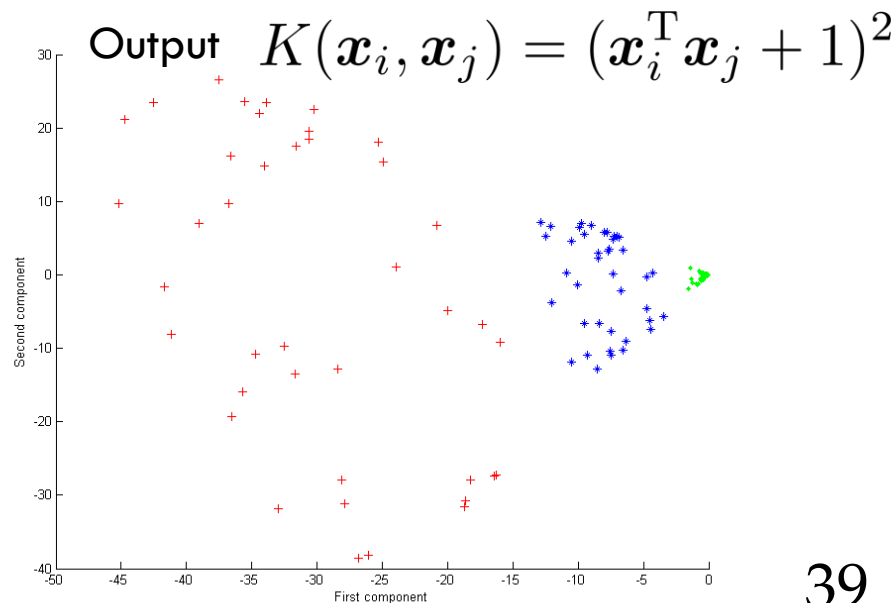
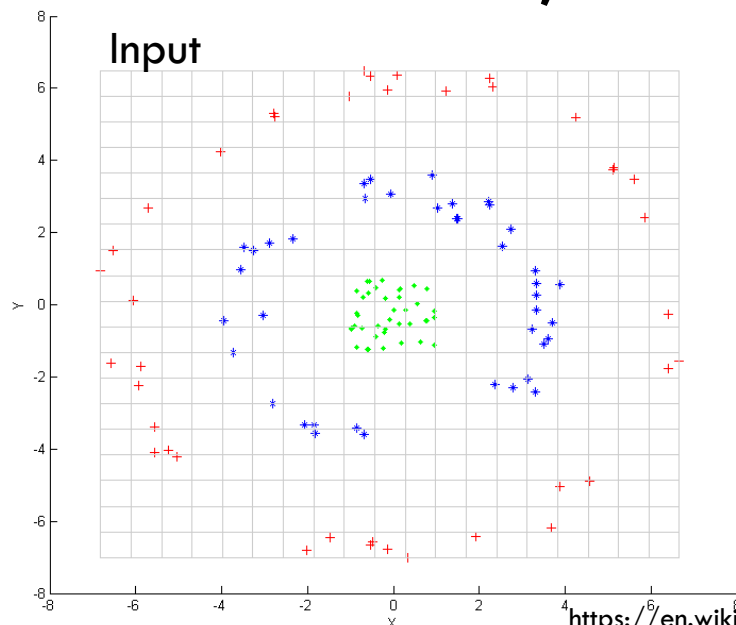
Kernel Dimensionality Reduction

□ Kernel PCA

- Previously: PCA on data covariance matrix $\hat{\Sigma} = X^T X$
- NOW PCA on kernel matrix : $N \times N$ $\hat{\Sigma} = (K(\mathbf{x}_i, \mathbf{x}_j))_{i,j=1}^N$

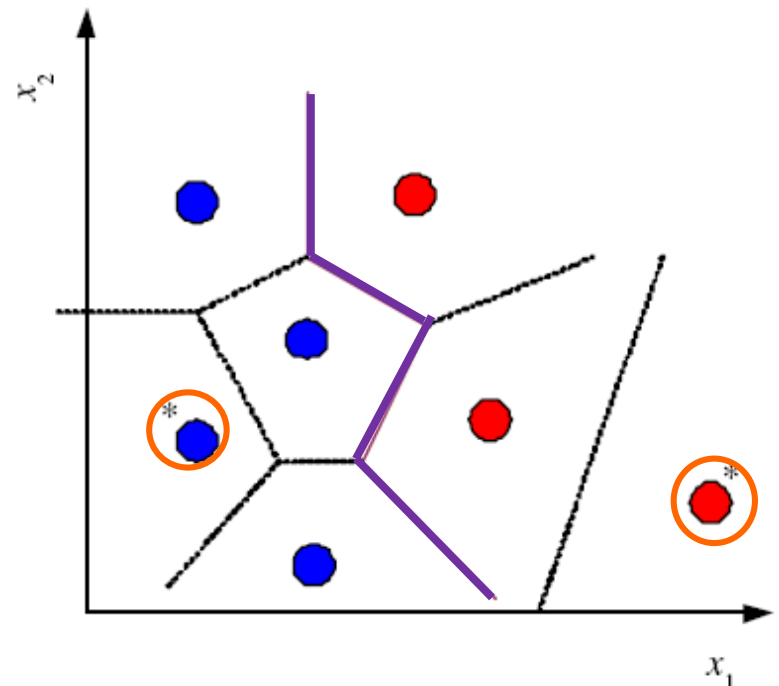
$$\hat{\Sigma} \mathbf{w}_k = \lambda_k \mathbf{w}_k, \quad \mathbf{w}_i^T \mathbf{w}_j = \begin{cases} 1 & , i = j \\ 0 & , i \neq j \end{cases}$$

□ also: Kernel LDA, CCA



Summary: Kernel Machines

- Many „kernelized“ methods
- Kernelized: if x cannot be linearly separated, use basis function $\Phi(x)$ or kernel function $K(x, x_i)$
- Various applications:
 - Classification
 - Regression
 - Dimensionality reduction
 - Outlier detection
- Few training points suffice for definition of hyperplane
- More robust to outliers



APPENDIX