

# Intro ML

## ML. 6. Support vector machines

DataLab CSIC

# Objectives and schedule

Introduce key concepts about SVMs. A case of non-probabilistic ML approach (now made probabilistic). Kernel methods

## Contents

- Maximal margin classifier
- Support Vector Classifier
- Support Vector Machines
- Extensions: More than two classes, SVM for reg, SVM vs LR,...

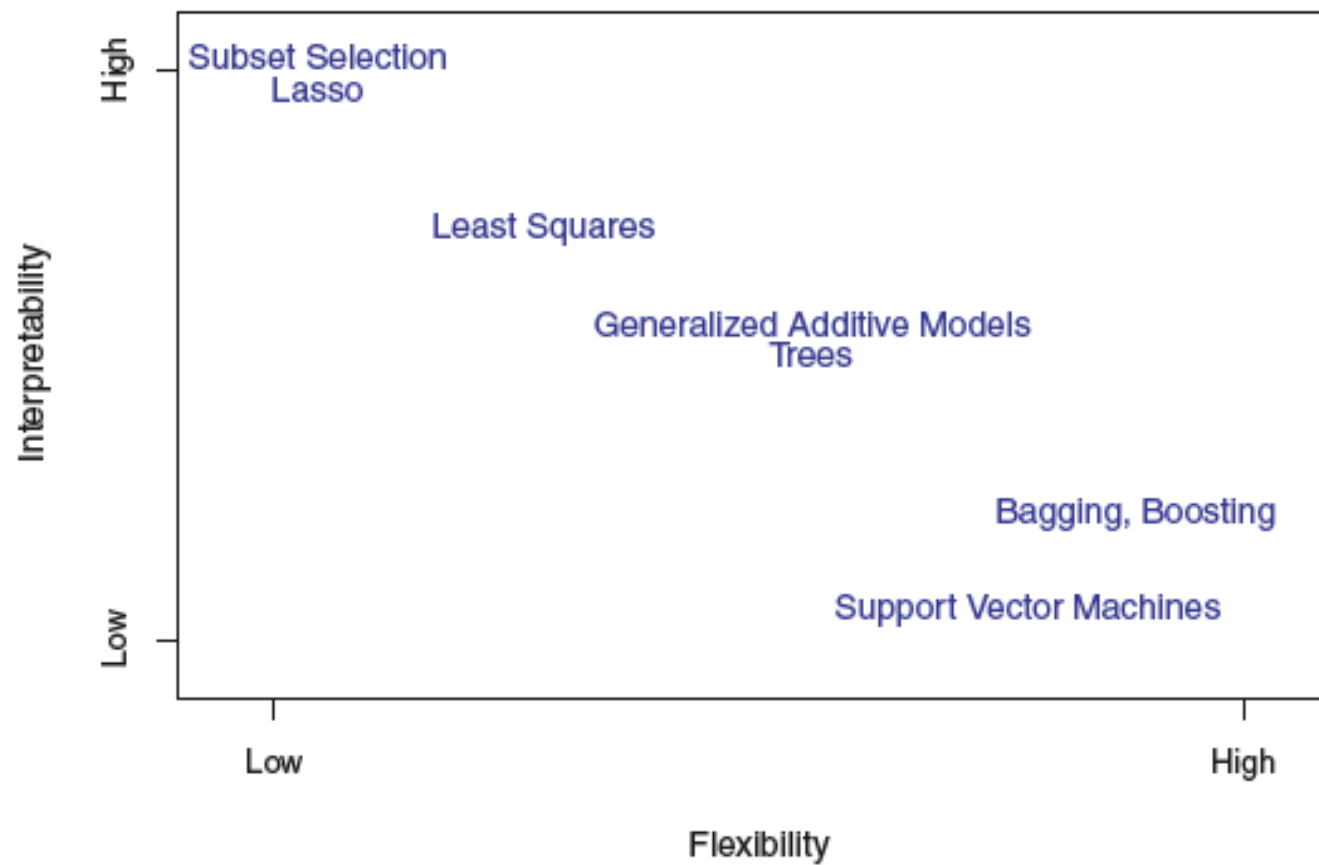
Case by Carlos Oscar Sorzano (CNB) on image processing

ISLR 9, CASI 19, Moguerza, Muñoz (2006), ESL 12, Bishop 7.1

# Labs

- Support Vector Classifier
- Support Vector Machines
- SVMs with multiple classes
- Full SVM project
- Comparison of Linear regression, Ridge regression, Lasso, Elastic Net, PLS, SV Regression

# A good old friend



# SVMs. Motivation

# Motivation

Perceptron (Rosenblatt, 1958). Minsky and Papert (1969) 'blew' them

Vapnik (Neural nets group at ATT) 1990's

SVMs lead computer vision competitions until (re)emergence of convolutional neural networks (mid 2000's)

Still very important

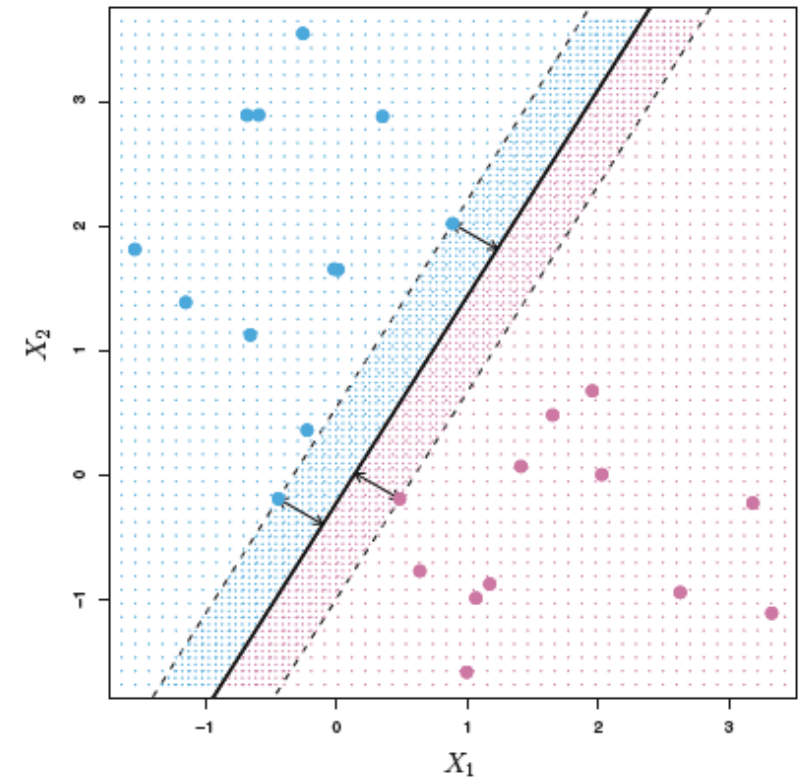
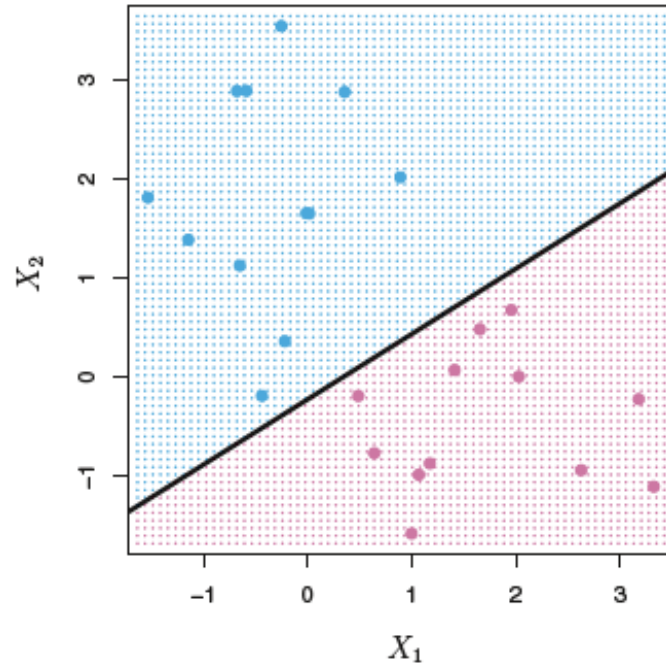
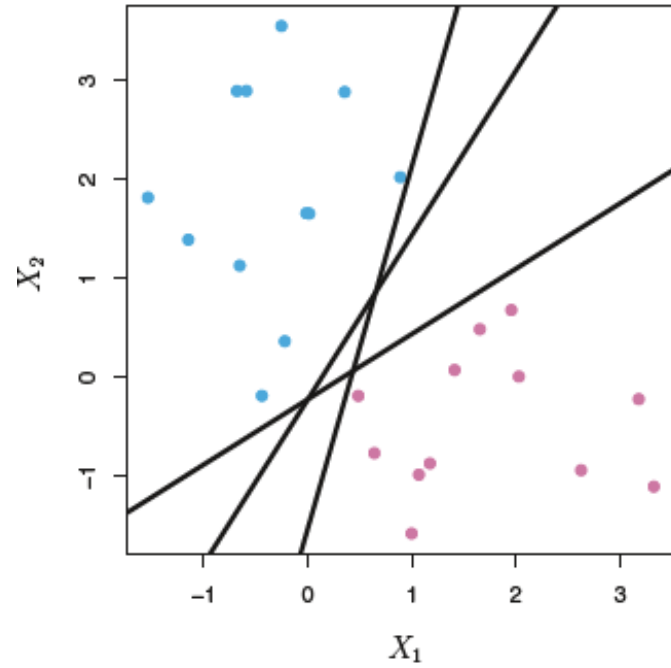
- Sparse solution. Computationally efficient

- Absence of local minima in SVM optimization

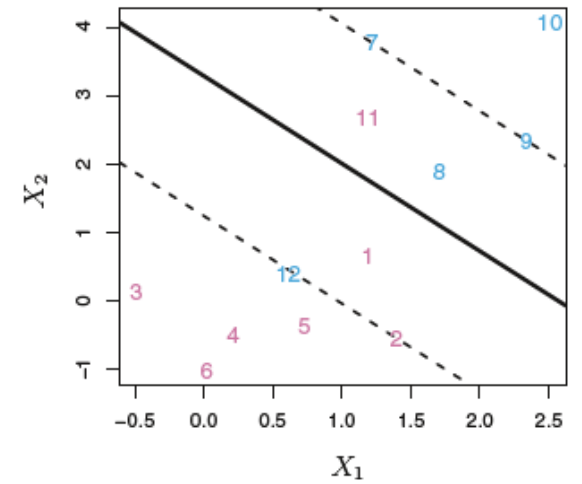
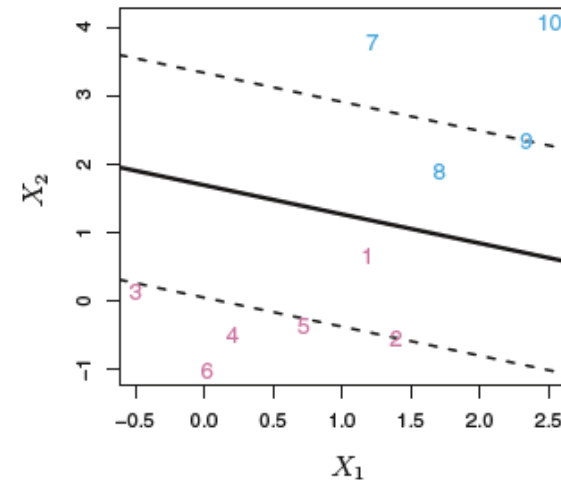
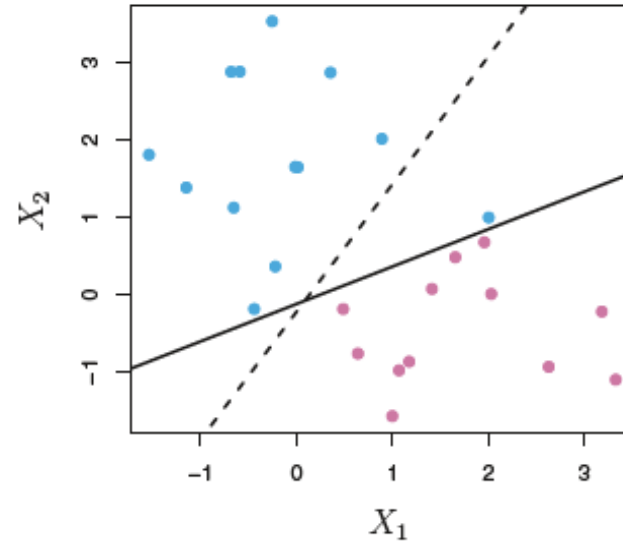
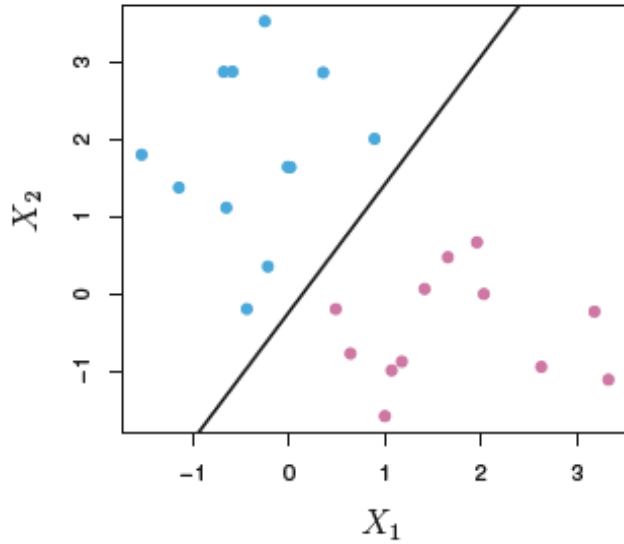
- Classification problem reduced to computation of linear decision function

Initially, not a statistical learning theory

# Concept. Linearly separable classes. MMC

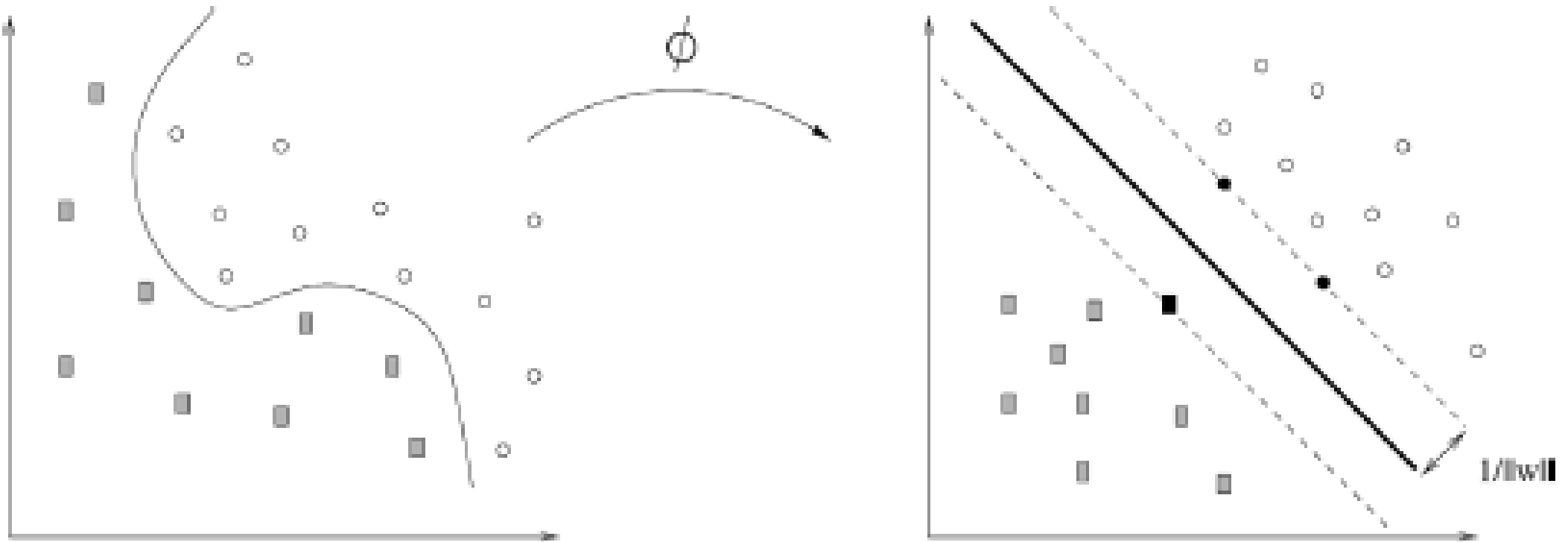


# Concept. Soft margins. SVC





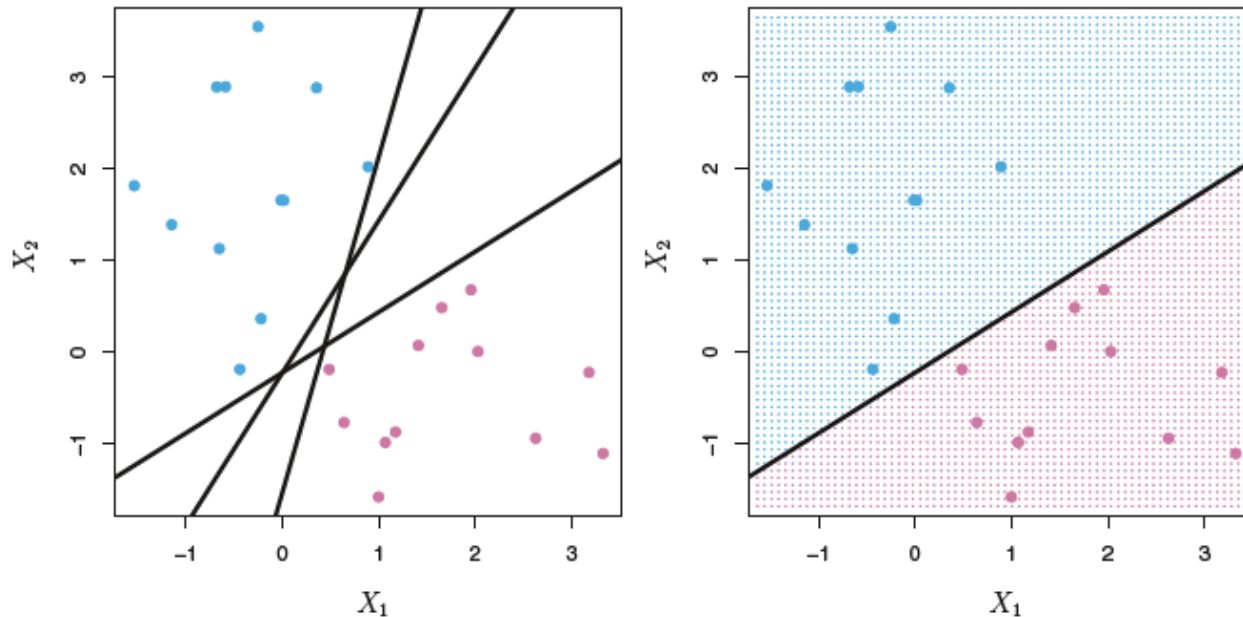
# Concept. Kernel trick. SVM



# Maximal margin classifier

# Maximal margin classifier. Basics

A hyperplane separates two regions



$$\begin{cases} \beta_0 + \beta_1 x_1 + \dots + \beta_p x_p = 0 \\ \beta_0 + \beta_1 x_1 + \dots + \beta_p x_p > 0 \\ \beta_0 + \beta_1 x_1 + \dots + \beta_p x_p < 0 \end{cases}$$

# Classifying with a separating hyperplane

A separating hyperplane has the property. Classes +1,-1

$$\begin{aligned} \beta_0 + \beta_1 x_{c1} + \beta_2 x_{c2} + \dots + \beta_p x_{cp} &> 0, \quad y_c = 1 \\ \beta_0 + \beta_1 x_{c1} + \beta_2 x_{c2} + \dots + \beta_p x_{cp} &< 0, \quad y_c = -1 \\ \Updownarrow \\ y_c (\beta_0 + \beta_1 x_{c1} + \dots + \beta_p x_{cp}) &> 0, \quad c = 1, \dots, n \end{aligned}$$

If a separating hyperplane exists, classify according to

$$\begin{aligned} x^* &\rightarrow f(x^*) = \beta_0 + \beta_1 x_1^* + \dots + \beta_p x_p^* \\ f(x^*) > 0 &\rightarrow 1 \\ f(x^*) < 0 &\rightarrow -1 \end{aligned} \quad \text{sgn}(f(x^*))$$

# Maximal margin classifier. Concept

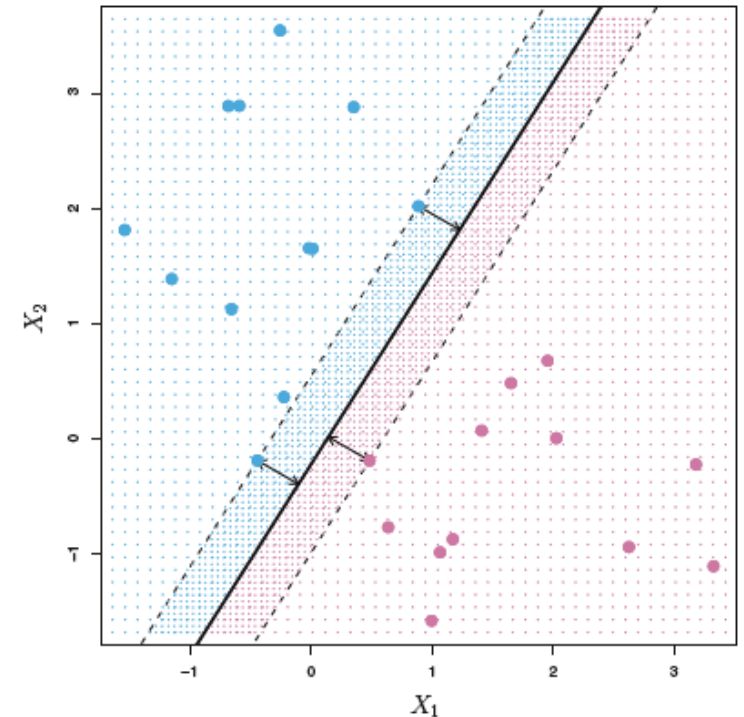
Optimal separating classifier → Maximal margin classifier

Separating hyperplane that is farthest from the training observations

Margin, minimal distance from points to hyperplane

Maximise margin

Depends only on the support vectors!!!



# Maximal margin classifier. Construction

Through the optimisation problem

$$\begin{array}{ll} \max & M \\ \beta_0, \dots, \beta_p & \\ \text{s.t.} & \sum_{j=1}^p \beta_j^2 = 1 \\ & y_i (\beta_0 + \beta_1 x_{i1} + \dots + \beta_p x_{ip}) \geq M, \quad i=1, \dots, n \end{array}$$

# Maximal margin classifier. Construction

Reformulate

$$\frac{1}{\|\beta\|} y_i (x_i^T \beta + \beta_0) \geq M \quad y_i (x_i^T \beta + \beta_0) \geq M \|\beta\|$$

$$\min_{\beta, \beta_0} \frac{1}{2} \|\beta\|^2$$

$$\text{s.t. } y_i (x_i^T \beta + \beta_0) \geq 1, i=1, \dots, n$$

Lagrangian formulation

$$\min \frac{1}{2} \|\beta\|^2 - \sum_{i=1}^n \alpha_i [y_i (x_i^T \beta + \beta_0) - 1]$$

# Maximal margin classifier. Construction

Derivatives and equal 0

$$\beta = \sum_{i=1}^n \alpha_i y_i x_i \quad 0 = \sum_{i=1}^n \alpha_i y_i$$

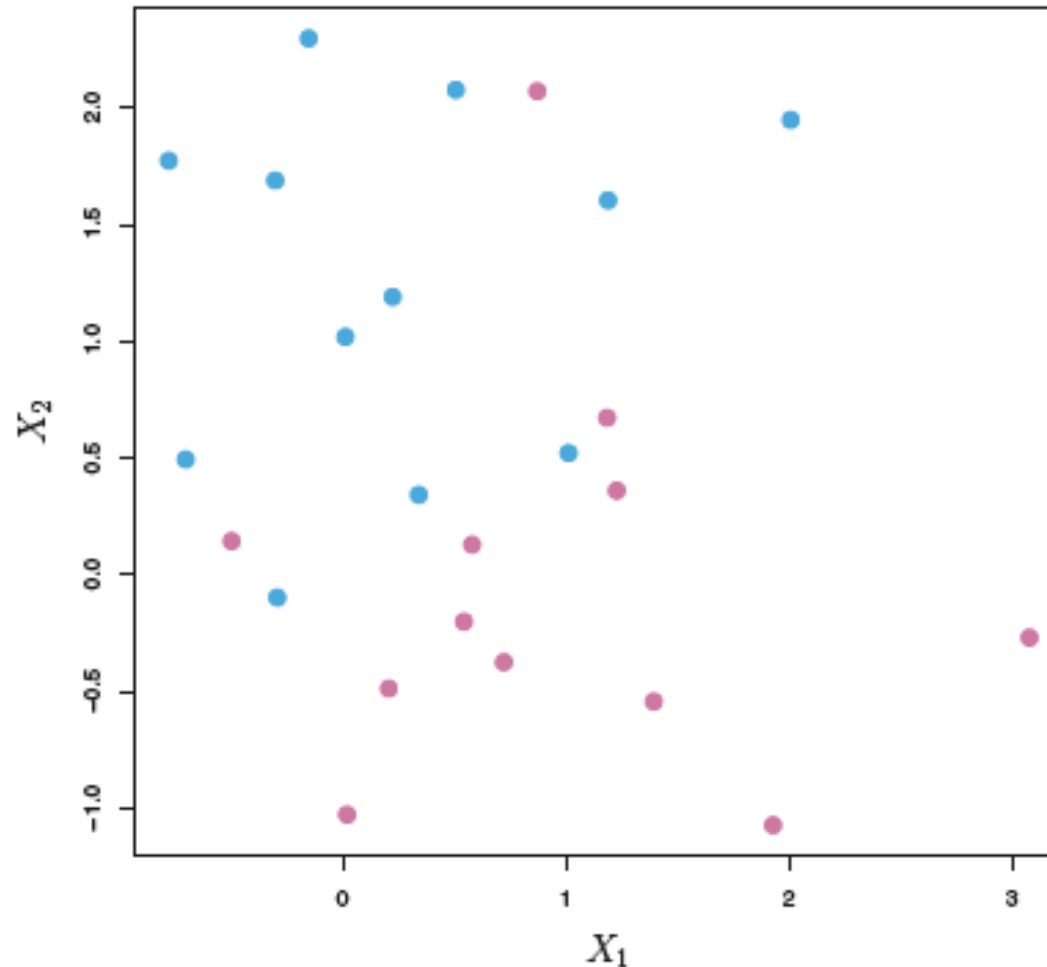
Dual

$$\begin{aligned} \max \quad & \sum \alpha_i - \frac{1}{2} \sum \sum \alpha_i \cdot \alpha_k y_i y_k x_i^T x_k \\ \text{s.t.} \quad & \alpha_i \geq 0 \quad \sum_{i=1}^n \alpha_i y_i = 0 \end{aligned}$$



# Maximal margin classifier. The non-separable case

Sometimes no solution exists....



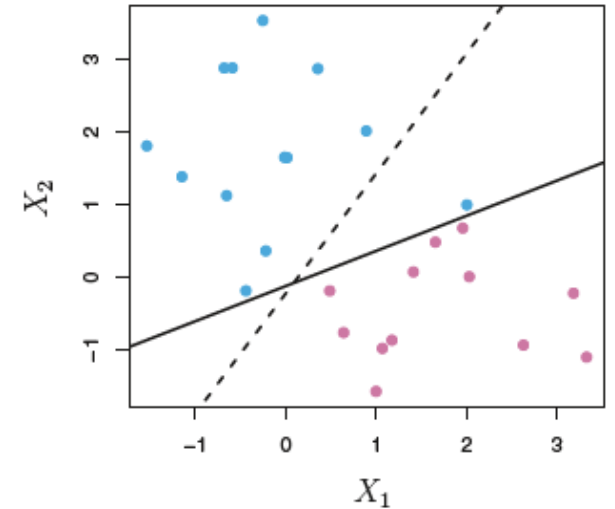
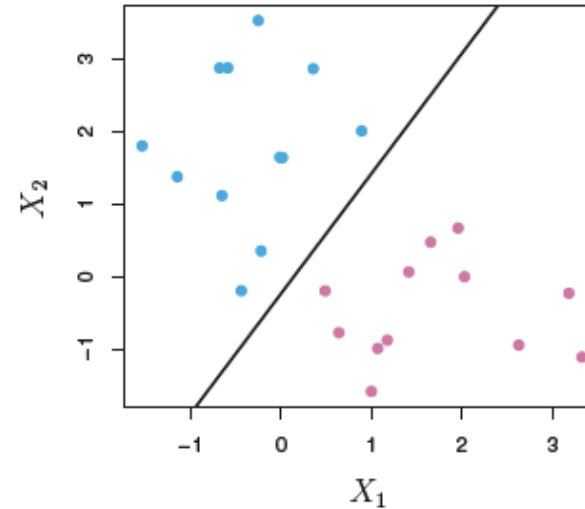
# Support vector classifiers

# Support vector classifiers. Concept I

MMC very sensitive to change in a single observation. Overfit

Willing to consider a classifier not perfectly separating to reach?

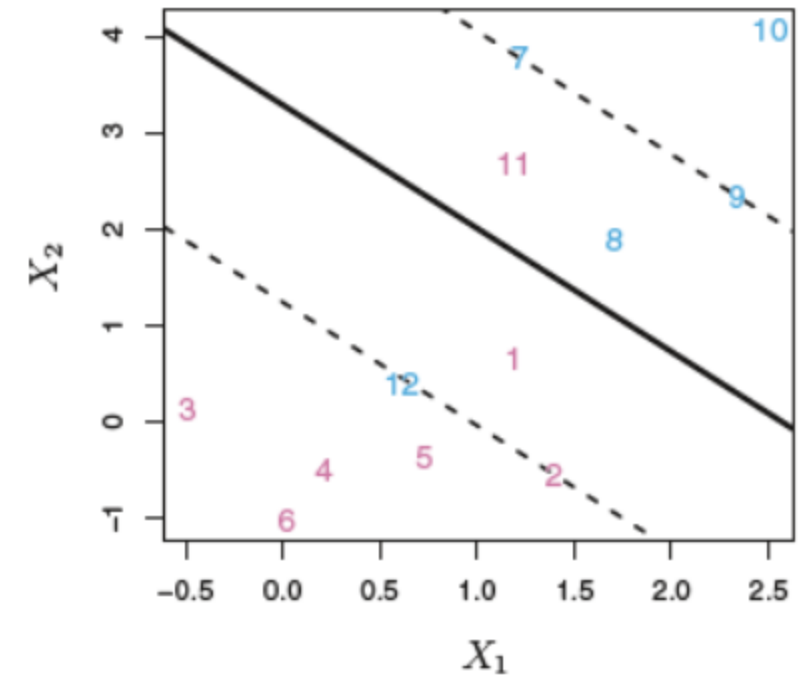
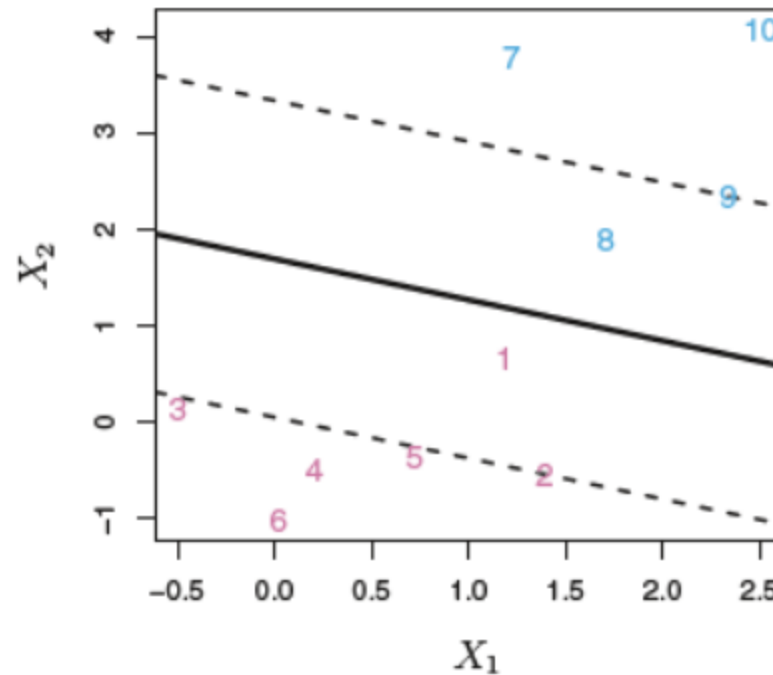
- Greater robustness
- Better classification of most training observations



# Support vector classifiers. Concept II

SVC aka soft margin classifier

Allow some observations  
In incorrect side of margin  
or incorrect side of  
hyperplane



# Support vector classifiers. Construction I

$$\begin{aligned} & \max \quad M \\ & \beta_0, \dots, \beta_p, \epsilon_1, \dots, \epsilon_n \\ & \text{s.t.} \quad \sum \beta_i^2 \\ & \quad y_i (\beta_0 + \beta_1 x_{i1} + \dots + \beta_p x_{ip}) \geq M(1 - \epsilon_i) \\ & \quad \epsilon_i \geq 0 \quad \sum_{i=1}^n \epsilon_i \leq C \end{aligned}$$

Classify through

$$x^* \longrightarrow \text{sgn} (\beta_0 + \beta_1 x_1^* + \dots + \beta_p x_p^*)$$

# Support vector classifiers. Construction II

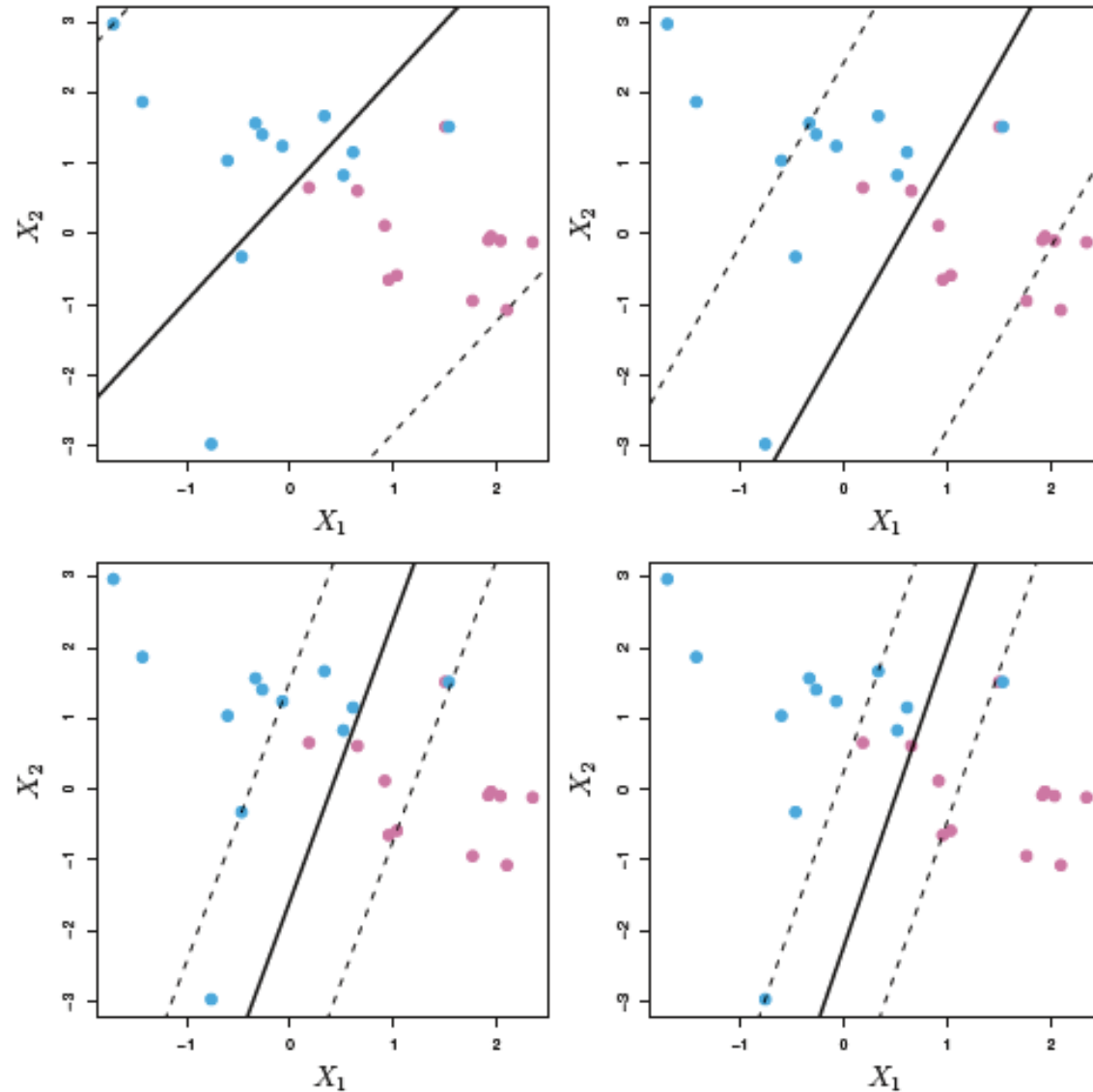
Slack  $\varepsilon$ .  $=0$ , right side of margin.  $>0$ , wrong side of margin.  $>1$ , wrong side of hyperplane

$C$  bounds sum of slacks.  $C=0$ , no slacks, MMC.  $C>0$ , at most  $C$  observations on the wrong side. Bigger  $C$ , more tolerant

$C$  tuning parameter chosen by CV

Only observations on the margin, or on the wrong side of the margin for their class matter in the optimisation. Support vectors

# Support vector classifiers. Construction III



# Support vector classifiers. Construction IV

The linear support vector classifier can be represented

$$f(x) = \beta_0 + \sum_{i=1}^n \alpha_i \langle x, x_i \rangle$$

To estimate the parameters just need

$$\langle x_i, x_i \rangle$$

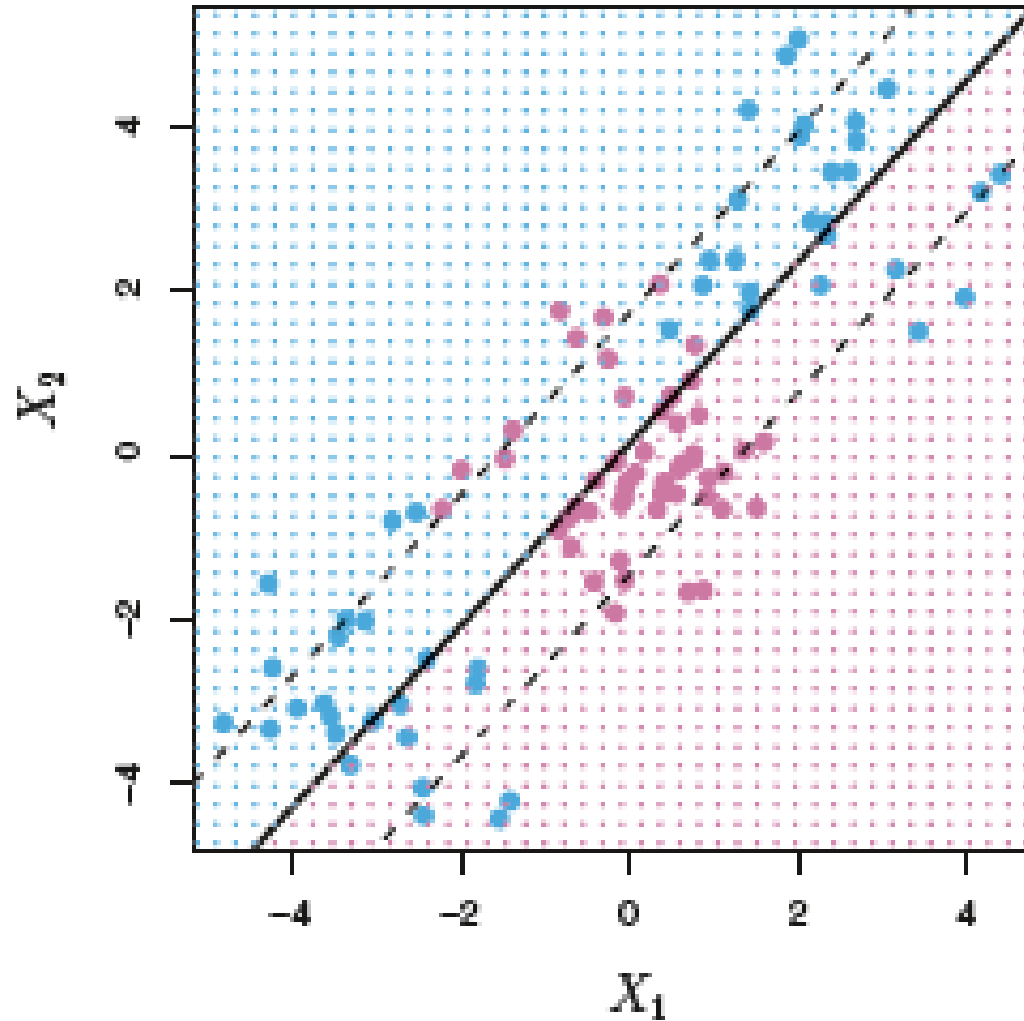
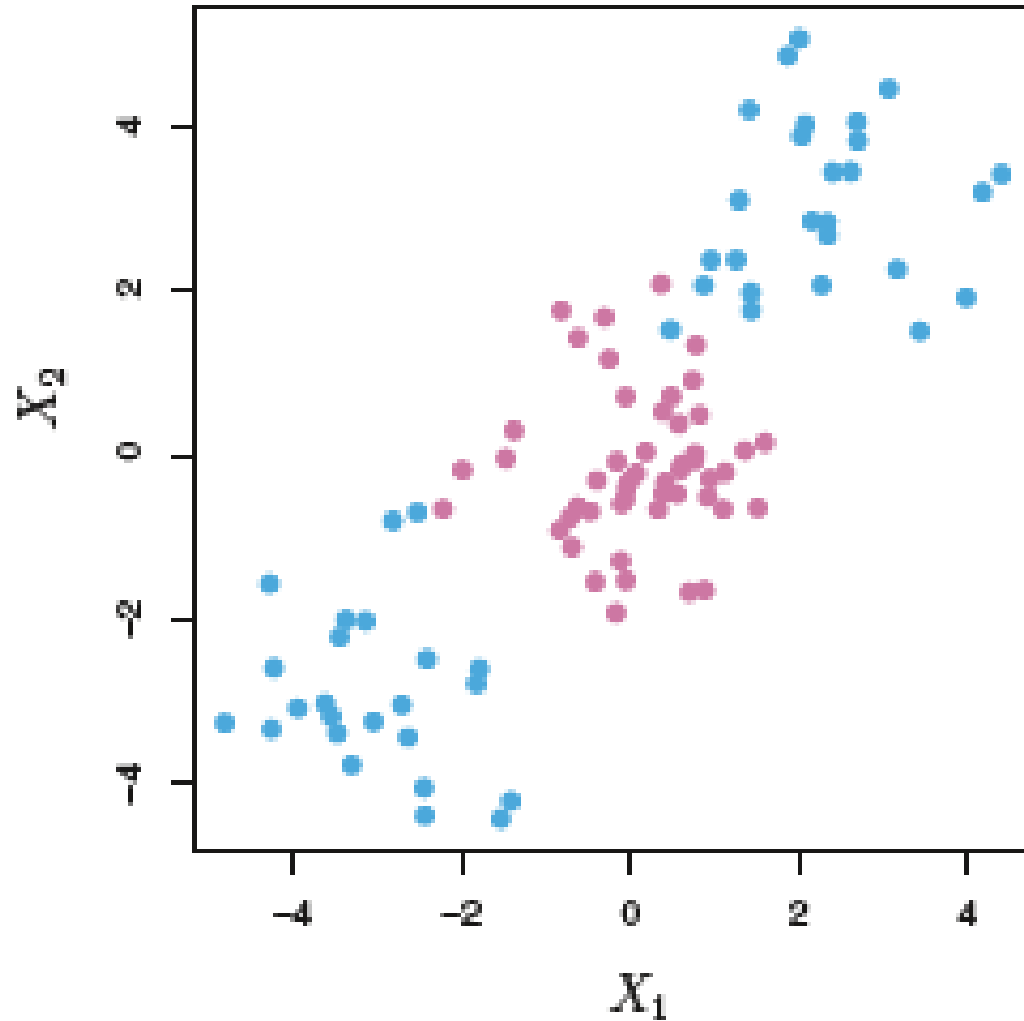
Only non-zero for support vectors

$$f(x) = \beta_0 + \sum_{i \in S} \alpha_i \langle x, x_i \rangle$$

We only need inner products!!!!



# Support vector classifiers. Failures



# Support vector machines

# Classifying with non-linear boundaries

From  $x_1, x_2, \dots, x_p$

To

Solve

$$x_1, x_1^2, x_2, x_2^2, \dots, x_p, x_p^2$$

$$\max M$$

$$\beta, \epsilon$$

$$\text{s.t. } y_i \left( \beta_0 + \sum_{j=1}^p \beta_{j1} x_{ij} + \sum_{j=1}^p \beta_{j2} x_{ij}^2 \right) \geq M(1 - \epsilon_i)$$

$$\sum_{i=1}^n \epsilon_i \leq C, \quad \epsilon_i \geq 0, \quad \sum_{j=1}^p \sum_{k=1}^2 \beta_{jk}^2 = 1$$

Linear in enlarged space, nonlinear in original space. Efficiency?

# Kernels

Map data into a higher dimensional space

**Cover's theorem.** Any data set becomes arbitrarily separable as data dimension grows

Kernels

$$\begin{aligned} K(x, y) \quad & \kappa: X \times X \rightarrow \mathbb{R} \\ & \phi: X \rightarrow Z \\ K(x, y) &= \phi^T(x) \phi(y) \end{aligned}$$

Inner product in  $Z$

$X$ , input space.  $Z$ , feature space

# Kernel. Example

$$K(x_1, x_2) = (1 + x_1^T x_2)^2$$
$$= (1 + x_{11}x_{21} + x_{12}x_{22})^2 = \phi^T(x_1) \phi(x_2)$$

$$\phi(x_i) = (1, \sqrt{2}x_{i1}, \sqrt{2}x_{i2}, x_{i1}^2, x_{i2}^2, \sqrt{2}x_{i1}x_{i2})$$

$$\phi: \mathbb{R}^2 \rightarrow \mathbb{R}^6$$

# Kernel

Positive definite function admitting an expansion

Mercer's theorem provides conditions

$$K(x, y) = \sum_{i=1}^{\infty} \lambda_i \phi_i(x) \phi_i(y) \quad \lambda_i \in \mathbb{R}^+$$

# Kernel

Linear kernel

Polynomial kernel degree  $d$

Radial kernel

$$K(x, y) = x^T y$$

$$K(x, y) = (c + x^T y)^d$$

$$K_c(x, y) = e^{-\|x - y\|^2 / c}$$

# Kernel

Given  $K$

$$f(x) = \sum \alpha_j K(x_j, x)$$

Reproducing Kernel Hilbert space (RKHS)

Adopt the form

$$f(x) = \sum_j \alpha_j K(x_j, x) + b$$

$$f(x) = \sum_S \alpha_j K(x_j, x) + b$$



# Support vector machines

SVC enlarging feature space with (nonlinear) kernel!!!!

# SVMs. Other topics

# SVMs. More than two classes

## One vs one classification

- Consider all pairs of classes

- Classify to the most frequently assigned class

## One vs all classification

- Fit SVM for each class against all others

- Assign to that maximising

# SVMs vs Logistic Regression

Rewrite SVC as

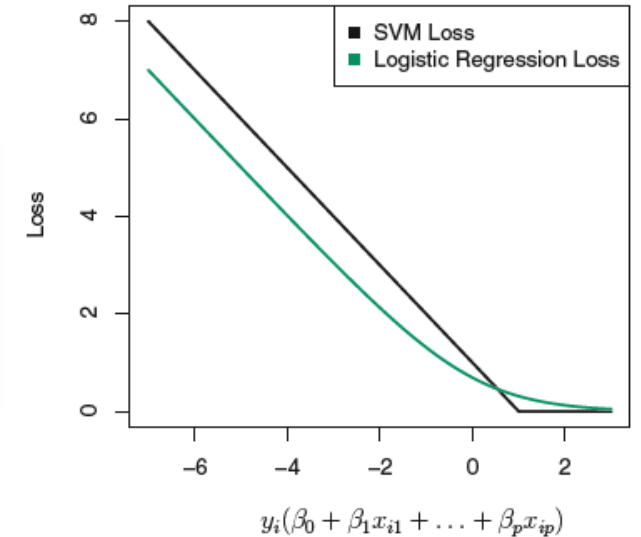
$$\min_{\beta} \underbrace{\sum_{i=1}^n \max(0, 1 - y_i f(x_i))}_{\text{LOSS}} + \underbrace{\lambda \sum_{i=1}^p \beta_i^2}_{\text{PENALTY}}$$

Loss+penalty

Hinge loss

$$L(x, y, \beta) = \sum_{i=1}^n \max[0, 1 - y_i (\beta_0 + \beta_1 x_{i1} + \dots + \beta_p x_{ip})]$$

$$\max[0, 1 - y_i f(x_i)]$$



# Further topics

Support vector regression

Computational learning theory

Platt trick to make it probabilistic

# Final comments

Perceptron

CNNs

Function fitting using kernels

Kernel smoothing and local regression