

# “Machine Learning estadístico para interfaces Cerebro-Computadora”

MODULO III - parte II

**Dra. Victoria Peterson**

[vpeterson@santafe-conicet.gov.ar](mailto:vpeterson@santafe-conicet.gov.ar)



NiCALab

Nov 2023



- 1 Análisis Discriminante Lineal
  - Método de transformación lineal
  - Regularización
  - Regresión
  
- 2 Máquinas de soporte vectorial
  - Datos linealmente separables
  - Datos no-linealmente separables
  - Kernel trick

## 1 Análisis Discriminante Lineal

- Método de transformación lineal
- Regularización
- Regresión

## 2 Máquinas de soporte vectorial

- Datos linealmente separables
- Datos no-linealmente separables
- Kernel trick

- 1 **Análisis Discriminante Lineal**
  - Método de transformación lineal
  - Regularización
  - Regresión
  
- 2 Máquinas de soporte vectorial
  - Datos linealmente separables
  - Datos no-linealmente separables
  - Kernel trick

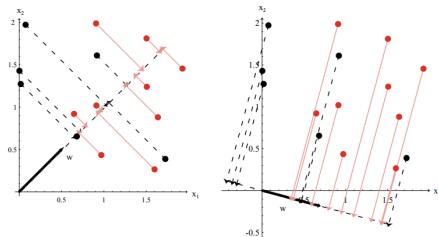
# Criterio de Fisher (FDA)

Formulación máxima/mínima varianza

Cap. 4[Webb, 2003]

## Problema

Hallar  $\mathbf{w}$  tal que  $J_F = \frac{\mathbf{w}^T \Sigma_b \mathbf{w}}{\mathbf{w}^T \hat{\Sigma}_w \mathbf{w}}$  se maximice.



Fuente: [Duda et al., 2012]

## 1 Análisis Discriminante Lineal

- Método de transformación lineal
- Regularización
- Regresión

## 2 Máquinas de soporte vectorial

- Datos linealmente separables
- Datos no-linealmente separables
- Kernel trick

# LDA regularizado



LDA se basa en estimaciones de las matrices de covarianza muestrales  $\Rightarrow$  cuando hay pocos datos  $\Sigma$  es *mal-condicionada*, y el cálculo de su inversa es *inestable*.

# LDA regularizado



LDA se basa en estimaciones de las matrices de covarianza muestrales  $\Rightarrow$  cuando hay pocos datos  $\Sigma$  es *mal-condicionada*, y el cálculo de su inversa es *inestable*.

## Regularización, ¿qué es?

Las técnicas de regularización buscan otorgar *estabilidad* a un proceso intrínsecamente inestable.



# LDA regularizado



LDA se basa en estimaciones de las matrices de covarianza muestrales  $\Rightarrow$  cuando hay pocos datos  $\Sigma$  es *mal-condicionada*, y el cálculo de su inversa es *inestable*.

## Regularización, ¿qué es?

Las técnicas de regularización buscan otorgar *estabilidad* a un proceso intrínsecamente inestable.

Permite inducir en la solución ciertas características (deseables) que, por ejemplo, mejorarán el proceso de aprendizaje, evitarán el sobre-entrenamiento y/o generarán soluciones más robustas a outliers.

# Shrinkage LDA

Formulación

[Blankertz et al., 2011]

## Shrinkage LDA

Cuando la cantidad de variables es mucho mayor que la cantidad de observaciones ( $p \gg N_t$ ) la estimación de  $\Sigma_c = \frac{1}{N_t} \sum_{i \in I_c} (\mathbf{x}_i - \boldsymbol{\mu}_c)(\mathbf{x}_i - \boldsymbol{\mu}_c)^T$  es pobre.

# Shrinkage LDA

Formulación

[Blankertz et al., 2011]

## Shrinkage LDA

Cuando la cantidad de variables es mucho mayor que la cantidad de observaciones ( $p \gg N_t$ ) la estimación de  $\Sigma_c = \frac{1}{N_t} \sum_{i \in I_c} (\mathbf{x}_i - \boldsymbol{\mu}_c)(\mathbf{x}_i - \boldsymbol{\mu}_c)^T$  es pobre.

Se define entonces

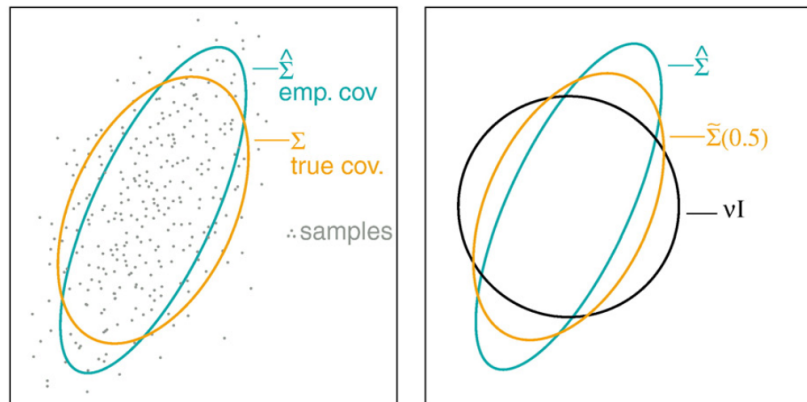
$$\tilde{\Sigma}_c = (1 - \gamma)\Sigma_c + \gamma\mathbf{I},$$

donde  $\gamma \in [0, 1]$  es el denominado parámetro de regularización

# Shrinkage LDA

## Formulación

[Blankertz et al., 2011]



Fuente: [Blankertz et al., 2011]

## 1 Análisis Discriminante Lineal

- Método de transformación lineal
- Regularización
- Regresión

## 2 Máquinas de soporte vectorial

- Datos linealmente separables
- Datos no-linealmente separables
- Kernel trick

# LDA como modelo de regresión

## Formulación

### LDA como modelo de regresión

Sea  $\mathbf{t} \in \mathbb{R}^{N_t}$  tal que  $t_i = N_2/n \quad \forall i/y_i = 1$ ;  $t_i = N_1/n \quad \forall i/y_i = 2$ .

$$\boldsymbol{\alpha}^* = \operatorname{argmin}_{\boldsymbol{\alpha} \in \mathbb{R}^p} \|\mathbf{t} - X\boldsymbol{\alpha}\|^2$$

# LDA como modelo de regresión

## Formulación

### LDA como modelo de regresión

Sea  $\mathbf{t} \in \mathbb{R}^{N_t}$  tal que  $t_i = N_2/n \quad \forall i/y_i = 1$ ;  $t_i = N_1/n \quad \forall i/y_i = 2$ .

$$\boldsymbol{\alpha}^* = \operatorname{argmin}_{\boldsymbol{\alpha} \in \mathbb{R}^p} \|\mathbf{t} - X\boldsymbol{\alpha}\|^2$$

Solución:

$$\boldsymbol{\alpha}^* = \frac{N_2 N_1}{N_t} (\boldsymbol{\Sigma}_t)^{-1} (\boldsymbol{\mu}_1 - \boldsymbol{\mu}_2)$$

# LDA como modelo de regresión

## Formulación

### LDA como modelo de regresión

Sea  $\mathbf{t} \in \mathbb{R}^{N_t}$  tal que  $t_i = N_2/n \quad \forall i/y_i = 1$ ;  $t_i = N_1/n \quad \forall i/y_i = 2$ .

$$\boldsymbol{\alpha}^* = \underset{\boldsymbol{\alpha} \in \mathbb{R}^p}{\operatorname{argmin}} \|\mathbf{t} - X\boldsymbol{\alpha}\|^2$$

Solución:

$$\boldsymbol{\alpha}^* = \frac{N_2 N_1}{N_t} (\boldsymbol{\Sigma}_t)^{-1} (\boldsymbol{\mu}_1 - \boldsymbol{\mu}_2)$$

$$\boldsymbol{\alpha}^* \propto \mathbf{w}^*$$



# LDA como modelo de regresión

## Formulación

### LDA como modelo de regresión

Sea  $\mathbf{t} \in \mathbb{R}^{N_t}$  tal que  $t_i = N_2/n \quad \forall i/y_i = 1$ ;  $t_i = N_1/n \quad \forall i/y_i = 2$ .

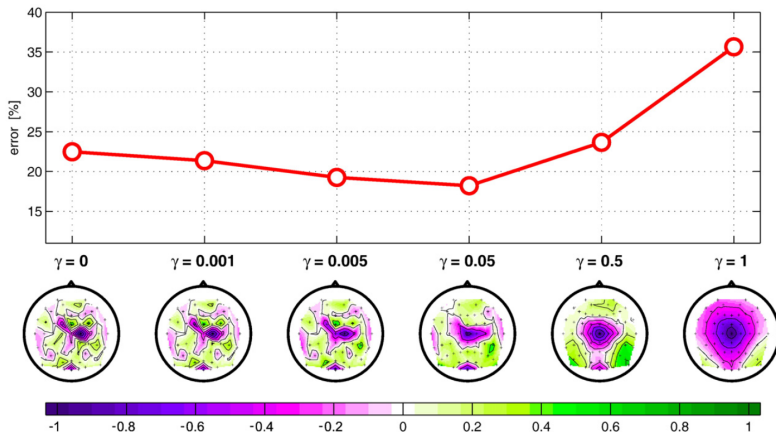
$$\boldsymbol{\alpha}^* = \operatorname{argmin}_{\boldsymbol{\alpha} \in \mathbb{R}^p} \|\mathbf{t} - X\boldsymbol{\alpha}\|^2$$

$$\boldsymbol{\alpha}^* \propto \mathbf{w}^*$$

### LDA regularizado

$$\boldsymbol{\alpha}^* = \operatorname{argmin}_{\boldsymbol{\alpha} \in \mathbb{R}^p} \|\mathbf{t} - X\boldsymbol{\alpha}\|^2 + \gamma \|\boldsymbol{\alpha}\|^2$$

# Ejemplo BCI



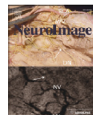
Fuente: [Blankertz et al., 2011]



Contents lists available at ScienceDirect

NeuroImage

journal homepage: [www.elsevier.com/locate/ynimg](http://www.elsevier.com/locate/ynimg)



## Single-trial analysis and classification of ERP components – A tutorial

Benjamin Blankertz<sup>a,b,\*</sup>, Steven Lemm<sup>b</sup>, Matthias Treder<sup>a</sup>, Stefan Haufe<sup>a</sup>, Klaus-Robert Müller<sup>a</sup>

<sup>a</sup> Berlin Institute of Technology, Machine Learning Laboratory, Berlin, Germany

<sup>b</sup> Fraunhofer FIRST, Intelligent Data Analysis Group, Berlin, Germany

### ARTICLE INFO

#### Article history:

Received 15 December 2009

Revised 14 June 2010

Accepted 18 June 2010

Available online 28 June 2010

#### Keywords:

EEG

ERP

BCI

Decoding

Machine learning

Shrinkage

LDA

### ABSTRACT

Analyzing brain states that correspond to event related potentials (ERPs) on a single trial basis is a hard problem due to the high trial-to-trial variability and the unfavorable ratio between signal (ERP) and noise (artifacts and neural background activity). In this tutorial, we provide a comprehensive framework for decoding ERPs, elaborating on linear concepts, namely spatio-temporal patterns and filters as well as linear ERP classification. However, the bottleneck of these techniques is that they require an accurate covariance matrix estimation in high dimensional sensor spaces which is a highly intricate problem. As a remedy, we propose to use shrinkage estimators and show that appropriate regularization of linear discriminant analysis (LDA) by shrinkage yields excellent results for single-trial ERP classification that are far superior to classical LDA classification. Furthermore, we give practical hints on the interpretation of what classifiers learned from the data and demonstrate in particular that the trade-off between goodness-of-fit and model complexity in regularized LDA relates to a morphing between a difference *pattern* of ERPs and a spatial *filter* which cancels non task-related brain activity.

- 1 Análisis Discriminante Lineal
  - Método de transformación lineal
  - Regularización
  - Regresión
  
- 2 Máquinas de soporte vectorial
  - Datos linealmente separables
  - Datos no-linealmente separables
  - Kernel trick

## 1 Análisis Discriminante Lineal

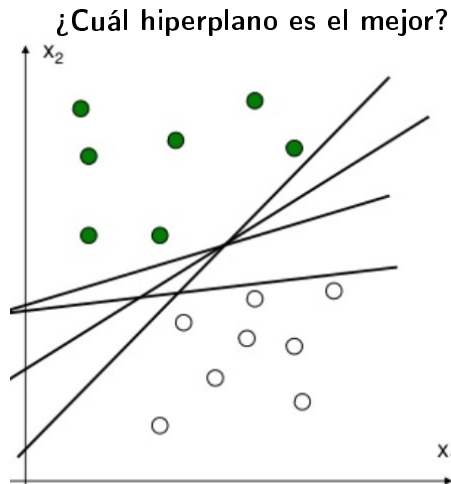
- Método de transformación lineal
- Regularización
- Regresión

## 2 Máquinas de soporte vectorial

- Datos linealmente separables
- Datos no-linealmente separables
- Kernel trick

# Máquinas de soporte vectorial

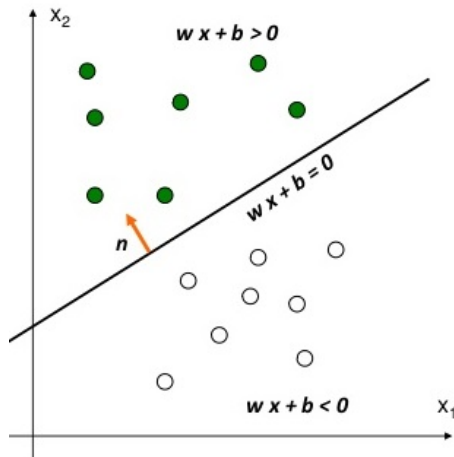
## Motivación



# Máquinas de soporte vectorial

## Motivación

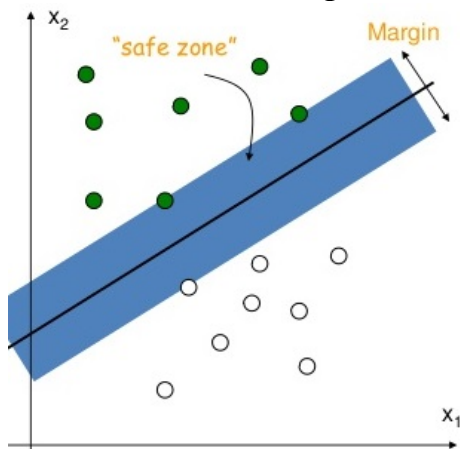
### Definición de márgenes



# Máquinas de soporte vectorial

## Motivación

### Definición de márgenes

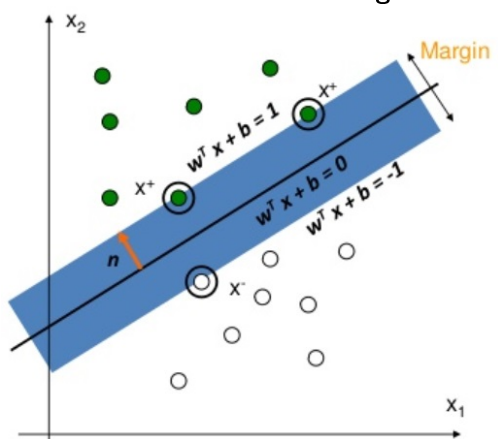




# Máquinas de soporte vectorial

## Motivación

### Maximización del margen



## 1 Análisis Discriminante Lineal

- Método de transformación lineal
- Regularización
- Regresión

## 2 Máquinas de soporte vectorial

- Datos linealmente separables
- Datos no-linealmente separables
- Kernel trick

## Datos linealmente separables

$$\mathbf{w}^T \mathbf{x} + b \begin{cases} > 0 \\ < 0 \end{cases} \Rightarrow \mathbf{x} \in \begin{cases} \mathcal{X}_1 (y_i = +1) \\ \mathcal{X}_2 (y_i = -1) \end{cases}$$

## Datos linealmente separables

$$\mathbf{w}^T \mathbf{x} + b \begin{cases} > 0 \\ < 0 \end{cases} \Rightarrow \mathbf{x} \in \begin{cases} \mathcal{X}_1 (y_i = +1) \\ \mathcal{X}_2 (y_i = -1) \end{cases}$$
$$y_i(\mathbf{w}^T \mathbf{x} + b) > 0 \forall i$$

## Datos linealmente separables

$$\mathbf{w}^T \mathbf{x} + b \begin{cases} > 0 \\ < 0 \end{cases} \Rightarrow \mathbf{x} \in \begin{cases} \mathcal{X}_1 (y_i = +1) \\ \mathcal{X}_2 (y_i = -1) \end{cases}$$
$$y_i(\mathbf{w}^T \mathbf{x} + b) > 0 \forall i$$

## Hiperplanos canónicos

$$H_1 : \mathbf{w}^T \mathbf{x} + b = +1; \quad H_2 : \mathbf{w}^T \mathbf{x} + b = -1$$

## Datos linealmente separables

$$\mathbf{w}^T \mathbf{x} + b \begin{cases} > 0 \\ < 0 \end{cases} \Rightarrow \mathbf{x} \in \begin{cases} \mathcal{X}_1 (y_i = +1) \\ \mathcal{X}_2 (y_i = -1) \end{cases}$$
$$y_i(\mathbf{w}^T \mathbf{x} + b) > 0 \forall i$$

## Hiperplanos canónicos

$$H_1 : \mathbf{w}^T \mathbf{x} + b = +1; \quad H_2 : \mathbf{w}^T \mathbf{x} + b = -1$$

$$\mathbf{w}^T \mathbf{x} + b \geq +1; \text{ para } y_i = +1$$

$$\mathbf{w}^T \mathbf{x} + b \leq -1; \text{ para } y_i = -1$$

## Datos linealmente separables

$$\mathbf{w}^T \mathbf{x} + b \begin{cases} > 0 \\ < 0 \end{cases} \Rightarrow \mathbf{x} \in \begin{cases} \mathcal{X}_1 (y_i = +1) \\ \mathcal{X}_2 (y_i = -1) \end{cases}$$
$$y_i(\mathbf{w}^T \mathbf{x} + b) > 0 \forall i$$

## Hiperplanos canónicos

$$H_1 : \mathbf{w}^T \mathbf{x} + b = +1; \quad H_2 : \mathbf{w}^T \mathbf{x} + b = -1$$

$$\mathbf{w}^T \mathbf{x} + b \geq +1; \text{ para } y_i = +1$$

$$\mathbf{w}^T \mathbf{x} + b \leq -1; \text{ para } y_i = -1$$

La distancia entre estos dos hiperplanos y el plano de separación  $\mathbf{w}^T \mathbf{x} + b = 0$  es el vector normal  $1/|\mathbf{w}|$ , denominado **margen**.

## Datos linealmente separables

$$\mathbf{w}^T \mathbf{x} + b \begin{cases} > 0 \\ < 0 \end{cases} \Rightarrow \mathbf{x} \in \begin{cases} \mathcal{X}_1 (y_i = +1) \\ \mathcal{X}_2 (y_i = -1) \end{cases}$$
$$y_i(\mathbf{w}^T \mathbf{x} + b) > 0 \forall i$$

## Hiperplanos canónicos

$$H_1 : \mathbf{w}^T \mathbf{x} + b = +1; \quad H_2 : \mathbf{w}^T \mathbf{x} + b = -1$$

$$\mathbf{w}^T \mathbf{x} + b \geq +1; \text{ para } y_i = +1$$

$$\mathbf{w}^T \mathbf{x} + b \leq -1; \text{ para } y_i = -1$$

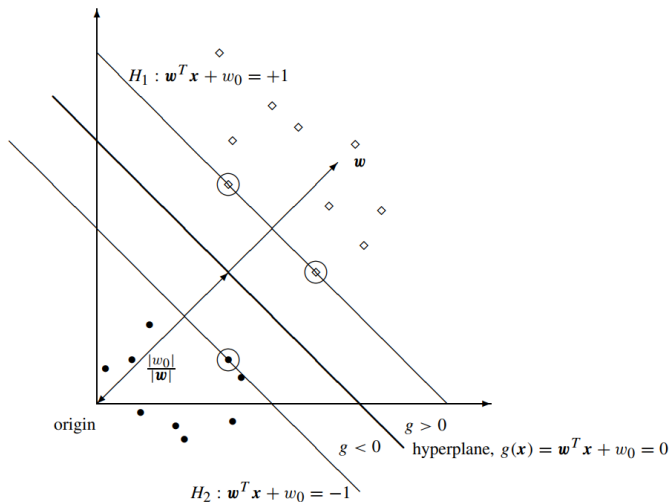
La distancia entre estos dos hiperplanos y el plano de separación  $\mathbf{w}^T \mathbf{x} + b = 0$  es el vector normal  $1/|\mathbf{w}|$ , denominado **margen**. Los puntos que caen sobre los hiperplanos canónicos se denominan **vectores soporte**.



# Máquinas de soporte vectorial

Formulación

Cap. 4[Webb, 2003]



Fuente: [Webb, 2003]

## Problema

$$\text{máx } \|\mathbf{w}\|_2^2 \quad \text{s.a} \quad y_i(\mathbf{w}^T \mathbf{x}_i + b) \geq 1 \quad \forall i$$

## Problema

$$\text{máx } \|\mathbf{w}\|_2^2 \quad \text{s.a.} \quad y_i(\mathbf{w}^T \mathbf{x}_i + b) \geq 1 \quad \forall i$$

$$\mathcal{L}_p = \frac{1}{2} \mathbf{w}^T \mathbf{w} - \sum_{i=1}^n \lambda_i (y_i(\mathbf{w}^T \mathbf{x}_i + b) - 1)$$

## Problema

$$\text{máx } \|\mathbf{w}\|_2^2 \quad \text{s.a.} \quad y_i(\mathbf{w}^T \mathbf{x}_i + b) \geq 1 \quad \forall i$$

$$\mathcal{L}_p = \frac{1}{2} \mathbf{w}^T \mathbf{w} - \sum_{i=1}^n \lambda_i (y_i(\mathbf{w}^T \mathbf{x}_i + b) - 1)$$

**Problema de optimización cuadrático con restricciones lineales**

## 1 Análisis Discriminante Lineal

- Método de transformación lineal
- Regularización
- Regresión

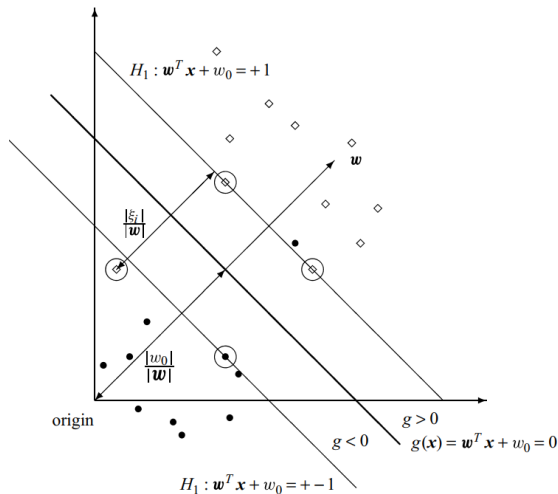
## 2 Máquinas de soporte vectorial

- Datos linealmente separables
- Datos no-linealmente separables
- Kernel trick

# Máquinas de soporte vectorial

Datos no-linealmente separables

Cap. 4[Webb, 2003]



Fuente: [Webb, 2003]

## Relajamos las restricciones

$$H_1 : \mathbf{w}^T \mathbf{x} + b = +1; \quad H_2 : \mathbf{w}^T \mathbf{x} + b = -1$$

$$\mathbf{w}^T \mathbf{x} + b \geq +1 - \zeta_i; \quad \text{para } y_i = +1,$$

$$\mathbf{w}^T \mathbf{x} + b \leq -1 + \zeta_i; \quad \text{para } y_i = -1,$$

$$\zeta_i \geq 0.$$

## Problema regularizado

$$\text{máx } \frac{1}{2} \|\mathbf{w}\|_2^2 + \alpha \sum_i \zeta_i \quad \text{s.a.} \quad y_i(\mathbf{w}^T \mathbf{x}_i + b) \geq 1 - \zeta_i \quad \forall i$$

## Relajamos las restricciones

$$H_1 : \mathbf{w}^T \mathbf{x} + b = +1; \quad H_2 : \mathbf{w}^T \mathbf{x} + b = -1$$

$$\mathbf{w}^T \mathbf{x} + b \geq +1 - \zeta_i; \quad \text{para } y_i = +1,$$

$$\mathbf{w}^T \mathbf{x} + b \leq -1 + \zeta_i; \quad \text{para } y_i = -1,$$

$$\zeta_i \geq 0.$$

## Problema regularizado

$$\text{máx } \frac{1}{2} \|\mathbf{w}\|_2^2 + \alpha \sum_i \zeta_i \quad \text{s.a.} \quad y_i(\mathbf{w}^T \mathbf{x}_i + b) \geq 1 - \zeta_i \quad \forall i$$

$$\mathcal{L}_p = \frac{1}{2} \mathbf{w}^T \mathbf{w} + \alpha \sum_i \zeta_i - \sum_{i=1}^n \lambda_i (y_i(\mathbf{w}^T \mathbf{x}_i + b) - 1 + \zeta_i) - \sum_{i=1}^n \delta_i \zeta_i,$$

siendo  $\lambda$  y  $\delta$  los multiplicadores de Lagrange.



## 1 Análisis Discriminante Lineal

- Método de transformación lineal
- Regularización
- Regresión

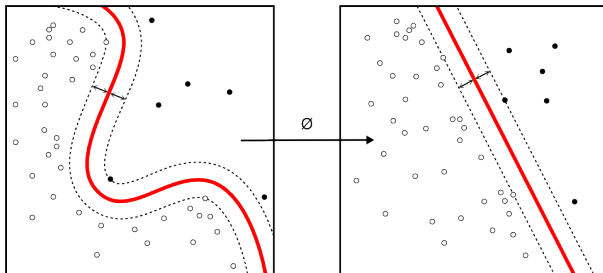
## 2 Máquinas de soporte vectorial

- Datos linealmente separables
- Datos no-linealmente separables
- Kernel trick

# Máquinas de soporte vectorial

El truco del Kernel

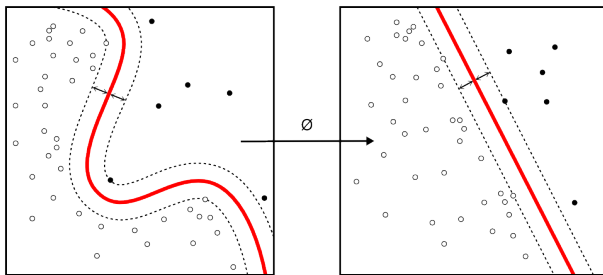
Cap. 5[Webb, 2003]



# Máquinas de soporte vectorial

El truco del Kernel

Cap. 5[Webb, 2003]



**Función discriminativa**

## Clasificadores lineales

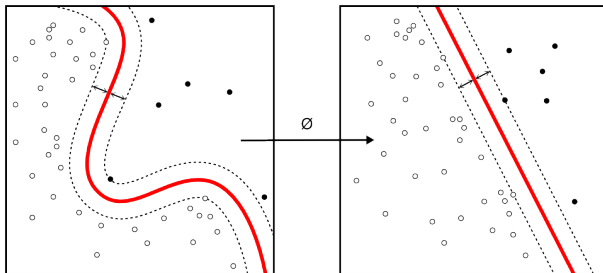
$$f(\mathbf{x}) = \sum_{i=1}^p w_i x_i + b = \mathbf{w}^T \mathbf{x} + b,$$

tal que si  $\text{sgn}(\mathbf{w}^T \mathbf{x} + b) > 0 \Rightarrow \mathbf{x} \in \mathcal{X}_1$

# Máquinas de soporte vectorial

## El truco del Kernel

Cap. 5[Webb, 2003]



## Clasificadores basados en Kernels

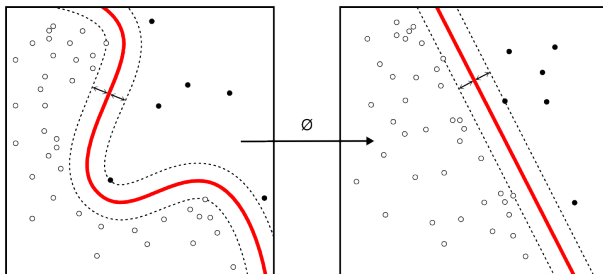
$$f(\mathbf{x}) = \sum_{i=1}^p w_i \phi(x)_i + b = \mathbf{w}^T \phi(\mathbf{x}) + b,$$

tal que si  $\text{sgn}(\mathbf{w}^T \phi(\mathbf{x}) + b) > 0 \Rightarrow \mathbf{x} \in \mathcal{W}_1$

# Máquinas de soporte vectorial

## El truco del Kernel

Cap. 5[Webb, 2003]



## Clasificadores basados en Kernels

$$f(\mathbf{x}) = \sum_{i=1}^p w_i \phi(x)_i + b = \mathbf{w}^T \phi(\mathbf{x}) + b,$$

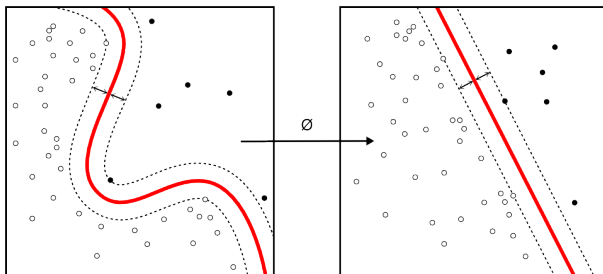
tal que si  $\text{sgn}(\mathbf{w}^T \phi(\mathbf{x}) + b) > 0 \Rightarrow \mathbf{x} \in \mathcal{W}_1$

Kernel:  $K(\mathbf{x}, \mathbf{y}) = \phi^T(\mathbf{x})\phi(\mathbf{y})$

# Máquinas de soporte vectorial

El truco del Kernel

Cap. 5[Webb, 2003]



**Table 5.2** Support vector machine kernels

Nonlinearity	Mathematical form $K(\mathbf{x}, \mathbf{y})$
Polynomial	$(1 + \mathbf{x}^T \mathbf{y})^d$
Gaussian	$\exp(- \mathbf{x} - \mathbf{y} ^2 / \sigma^2)$
Sigmoid	$\tanh(k \mathbf{x}^T \mathbf{y} - \delta)$

Fuente: [Webb, 2003]

# Bibliografía utilizada I

Blankertz, B., Lemm, S., Treder, M., Hauf, S., and Müller, K. R. (2011). Single-trial analysis and classification of ERP component- a tutorial. *Neuroimage*, 56:814–825.

Duda, R. O., Hart, P. E., and Stork, D. G. (2012). *Pattern classification*. John Wiley & Sons.

Webb, A. R. (2003). *Statistical pattern recognition*. John Wiley & Sons.