

# Máquinas de Vectores de Soporte

Alvaro J. Riascos Villegas

Junio 2018

# Contenido

- 1 **Introducción**
- 2 Hiperplano separador óptimo
- 3 Máquinas de vectores de soporte
- 4 SVM y Kernels

# Introducción

- SVM es una técnica que generaliza las técnicas de hiperplanos separadores.
- Construye hiperplanos separadores en un espacio de variables transformado por funciones no lineales. Estas funciones llevan los features a espacios potencialmente mucho más grandes.
- SVM es una de las técnicas más estudiadas actualmente por ser un problema convexo y se conoce la dimensión VC del clasificador.

# Introducción

- SVM es una técnica que generaliza las técnicas de hiperplanos separadores.
- Construye hiperplanos separadores en un espacio de variables transformado por funciones no lineales. Estas funciones llevan los features a espacios potencialmente mucho más grandes.
- SVM es una de las técnicas más estudiadas actualmente por ser un problema convexo y se conoce la dimensión VC del clasificador.

# Introducción

- SVM es una técnica que generaliza las técnicas de hiperplanos separadores.
- Construye hiperplanos separadores en un espacio de variables transformado por funciones no lineales. Estas funciones llevan los features a espacios potencialmente mucho más grandes.
- SVM es una de las técnicas más estudiadas actualmente por ser un problema convexo y se conoce la dimensión VC del clasificador.

# Contenido

- 1 Introducción
- 2 Hiperplano separador óptimo
- 3 Máquinas de vectores de soporte
- 4 SVM y Kernels

# Hiperplano separador óptimo

- Los datos de entrenamiento son  $\{(x_1, y_1), \dots, (x_N, y_N)\}$ ,  $x_i \in R^p$ ,  $y_i \in \{-1, 1\}$ .
- Un hiperplano separador es un conjunto de la forma:  $\{x : f(x) = x^T \beta + \beta_0 = 0\}$ . Donde  $\beta$  se puede normalizar a que tenga norma 1.
- La regla de clasificación que induce el hiperplano es:  $\text{sign}[x^T \beta + \beta_0]$ .
- En la figura los datos de entrenamiento son perfectamente separables por el hiperplano y se puede garantizar que para todo  $i$ ,  $f(x_i)y_i > 0$ .

# Hiperplano separador óptimo

- Los datos de entrenamiento son  $\{(x_1, y_1), \dots, (x_N, y_N)\}$ ,  $x_i \in R^p$ ,  $y_i \in \{-1, 1\}$ .
- Un hiperplano separador es un conjunto de la forma:  $\{x : f(x) = x^T \beta + \beta_0 = 0\}$ . Donde  $\beta$  se puede normalizar a que tenga norma 1.
- La regla de clasificación que induce el hiperplano es:  $\text{sign}[x^T \beta + \beta_0]$ .
- En la figura los datos de entrenamiento son perfectamente separables por el hiperplano y se puede garantizar que para todo  $i$ ,  $f(x_i)y_i > 0$ .

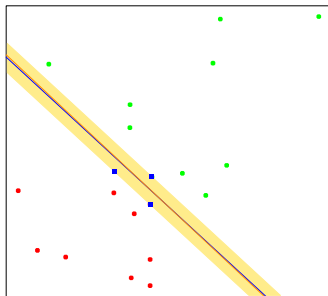
# Hiperplano separador óptimo

- Los datos de entrenamiento son  $\{(x_1, y_1), \dots, (x_N, y_N)\}$ ,  $x_i \in R^p$ ,  $y_i \in \{-1, 1\}$ .
- Un hiperplano separador es un conjunto de la forma:  $\{x : f(x) = x^T \beta + \beta_0 = 0\}$ . Donde  $\beta$  se puede normalizar a que tenga norma 1.
- La regla de clasificación que induce el hiperplano es:  $\text{sign}[x^T \beta + \beta_0]$ .
- En la figura los datos de entrenamiento son perfectamente separables por el hiperplano y se puede garantizar que para todo  $i$ ,  $f(x_i)y_i > 0$ .

# Hiperplano separador óptimo

- Los datos de entrenamiento son  $\{(x_1, y_1), \dots, (x_N, y_N)\}$ ,  $x_i \in R^p$ ,  $y_i \in \{-1, 1\}$ .
- Un hiperplano separador es un conjunto de la forma:  $\{x : f(x) = x^T \beta + \beta_0 = 0\}$ . Donde  $\beta$  se puede normalizar a que tenga norma 1.
- La regla de clasificación que induce el hiperplano es:  $\text{sign}[x^T \beta + \beta_0]$ .
- En la figura los datos de entrenamiento son perfectamente separables por el hiperplano y se puede garantizar que para todo  $i$ ,  $f(x_i)y_i > 0$ .

# Hiperplano separador óptimo



**FIGURE 4.16.** *The same data as in Figure 4.14. The shaded region delineates the maximum margin separating the two classes. There are three support points indicated, which lie on the boundary of the margin, and the optimal separating hyperplane (blue line) bisects the slab. Included in the figure is the boundary found using logistic regression (red line), which is very close to the optimal separating hyperplane (see Section 12.3.3).*

# Hiperplano separador óptimo

- Si  $|\beta| = 1$ , la función  $f(x) = x^T \beta + \beta_0$  nos da la distancia (con signo) al hiperplano.
- Para ver esto sea  $x_0$  la proyección perpendicular al hiperplano. Entonces:  $(x - x_0) \cdot \beta = |x - x_0|$  excepto por el signo y:

$$(x - x_0) \cdot \beta = x^T \beta + \beta_0 - x_0^T \beta - \beta_0 = x^T \beta + \beta_0$$

- Luego el siguiente problema refleja la restricción de clasificar perfectamente y que el margen  $M$ , sea máximo:

$$\begin{aligned} & \underset{\beta, \beta_0, M}{\text{máx}} && M \\ & && |\beta| = 1 \\ & && y_i(x_i^T \beta + \beta_0) \geq M \end{aligned}$$

# Hiperplano separador óptimo

- Si  $|\beta| = 1$ , la función  $f(x) = x^T \beta + \beta_0$  nos da la distancia (con signo) al hiperplano.
- Para ver esto sea  $x_0$  la proyección perpendicular al hiperplano. Entonces:  $(x - x_0) \cdot \beta = |x - x_0|$  excepto por el signo y:

$$(x - x_0) \cdot \beta = x^T \beta + \beta_0 - x_0^T \beta - \beta_0 = x^T \beta + \beta_0$$

- Luego el siguiente problema refleja la restricción de clasificar perfectamente y que el margen  $M$ , sea máximo:

$$\begin{aligned} & \underset{\beta, \beta_0, M}{\text{máx}} && M \\ & && |\beta| = 1 \\ & && y_i(x_i^T \beta + \beta_0) \geq M \end{aligned}$$

# Hiperplano separador óptimo

- Si  $|\beta| = 1$ , la función  $f(x) = x^T \beta + \beta_0$  nos da la distancia (con signo) al hiperplano.
- Para ver esto sea  $x_0$  la proyección perpendicular al hiperplano. Entonces:  $(x - x_0) \cdot \beta = |x - x_0|$  excepto por el signo y:

$$(x - x_0) \cdot \beta = x^T \beta + \beta_0 - x_0^T \beta - \beta_0 = x^T \beta + \beta_0$$

- Luego el siguiente problema refleja la restricción de clasificar perfectamente y que el margen  $M$ , sea máximo:

$$\begin{aligned} & \underset{\beta, \beta_0, M}{\text{máx}} M \\ & |\beta| = 1 \\ & y_i (x_i^T \beta + \beta_0) \geq M \end{aligned}$$

# Hiperplano separador óptimo

- Este problema se puede escribir como:

$$\begin{aligned} & \underset{\beta, \beta_0, M}{\text{máx}} M \\ & y_i(x_i^T \beta + \beta_0) \geq M \quad |\beta| \end{aligned}$$

- Que es equivalente a:

$$\begin{aligned} & \underset{\beta, \beta_0, M}{\text{máx}} M \\ & y_i(x_i^T \beta + \beta_0) \geq M \quad |\beta| \end{aligned}$$

# Hiperplano separador óptimo

- Este problema se puede escribir como:

$$\begin{aligned} & \underset{\beta, \beta_0, M}{\text{máx}} M \\ & y_i(x_i^T \beta + \beta_0 - M) \geq 0 \quad \forall i \end{aligned}$$

- Que es equivalente a:

$$\begin{aligned} & \underset{\beta, \beta_0, M}{\text{máx}} M \\ & y_i(x_i^T \beta + \beta_0) \geq M \quad \forall i \end{aligned}$$

- Que es equivalente a:

$$\begin{aligned} & \underset{\beta, \beta_0, M, M'}{\text{máx}} \quad M \\ & y_i(x_i^T \beta + \beta_0) \geq M' \\ & M' = M |\beta| \end{aligned}$$

# Hiperplano separador óptimo

- Que es equivalente a:

$$\begin{aligned} & \underset{\beta, \beta_0, M'}{\text{máx}} \frac{M'}{|\beta|} \\ & y_i(x_i^T \beta + \beta_0) \geq M' \end{aligned}$$

- Equivalente a:

$$\begin{aligned} & \underset{\beta, \beta_0, M'}{\text{máx}} \frac{1}{\frac{|\beta|}{M'}} \\ & y_i(x_i^T \frac{\beta}{M'} + \beta_0) \geq 1 \end{aligned}$$

- Obsérvese que podemos elegir  $M' = 1$

# Hiperplano separador óptimo

- Que es equivalente a:

$$\begin{aligned} & \underset{\beta, \beta_0, M'}{\text{máx}} \frac{M'}{|\beta|} \\ & y_i(x_i^T \beta + \beta_0) \geq M' \end{aligned}$$

- Equivalente a:

$$\begin{aligned} & \underset{\beta, \beta_0, M'}{\text{máx}} \frac{1}{\frac{|\beta|}{M'}} \\ & y_i(x_i^T \frac{\beta}{M'} + \beta_0) \geq 1 \end{aligned}$$

- Obsérvese que podemos elegir  $M' = 1$

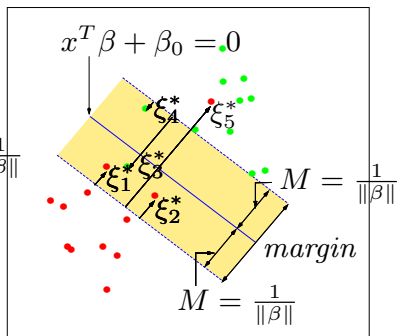
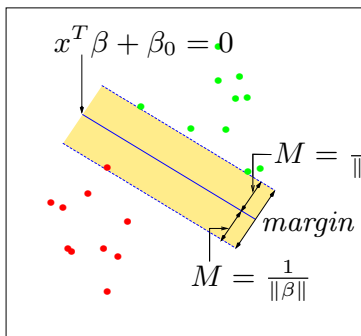
- Finalmente:

$$\begin{aligned} & \underset{\beta, \beta_0}{\text{mín}} \quad |\beta| \\ & y_i(x_i^T \beta + \beta_0) \geq 1 \end{aligned}$$

donde  $M = \frac{1}{|\beta|}$ .

- Panel izquierdo.

# Hiperplano separador óptimo



- El Lagrangiano del problema es:

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

- Las CPO:

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

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

$$\alpha_i(x_i^T \beta + \beta_0 - 1) = 0$$

- El Lagrangiano del problema es:

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

- Las CPO:

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

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

$$\alpha_i (x_i^T \beta + \beta_0 - 1) = 0$$

# Hiperplano separador óptimo: CPO

- Obsérvese que:  $\alpha_i > 0$  si y solo si  $(x_i, y_i)$  está en el borde de la región de separación.
- $\beta$  está definido únicamente por los puntos que están en el borde de la región. Estos puntos se llaman vectores de soporte.
- Panel izquierdo figura.

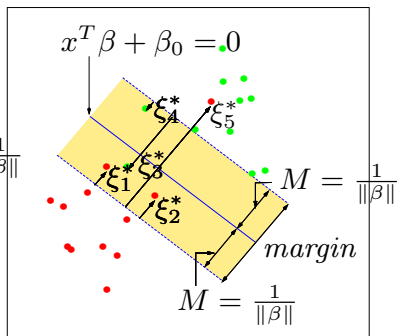
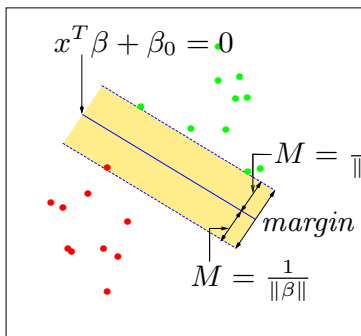
# Hiperplano separador óptimo: CPO

- Obsérvese que:  $\alpha_i > 0$  si y solo si  $(x_i, y_i)$  está en el borde de la región de separación.
- $\beta$  está definido únicamente por los puntos que están en el borde de la región. Estos puntos se llaman vectores de soporte.
- Panel izquierdo figura.

# Hiperplano separador óptimo: CPO

- Obsérvese que:  $\alpha_i > 0$  si y solo si  $(x_i, y_i)$  está en el borde de la región de separación.
- $\beta$  está definido únicamente por los puntos que están en el borde de la región. Estos puntos se llaman vectores de soporte.
- Panel izquierdo figura.

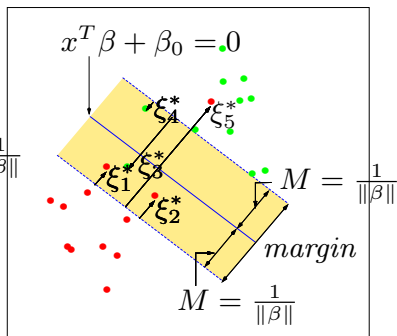
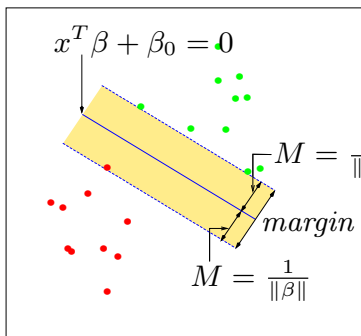
# Hiperplano separador óptimo



# Hiperplano separador óptimo: clasificación imperfecta

- Aun cuando se pueda separar perfectamente puede no ser deseable. El hiperplano puede ser muy sensible a nuevos datos (i.e., varianza alta).
- Panel derecho figura.

# Hiperplano separador óptimo



# Contenido

- 1 Introducción
- 2 Hiperplano separador óptimo
- 3 Máquinas de vectores de soporte**
- 4 SVM y Kernels

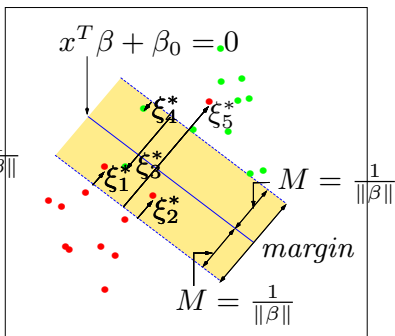
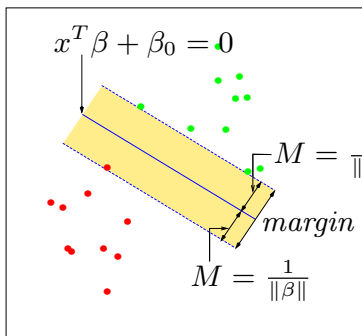
# Hiperplano separador óptimo: clasificación imperfecta

- Cuando los datos de entrenamiento no son perfectamente separables, o no es deseable la separación perfecta por un hiperplano, se puede formular un problema similar:

$$\begin{aligned} & \underset{\beta, \beta_0, \xi, M}{\text{máx}} && M \\ & && |\beta| = 1 \\ & y_i(x_i^T \beta + \beta_0) && \geq M - \xi_i \\ & && \xi_i \geq 0 \\ & \sum_{i=1}^n \xi_i && \leq \text{Const} \end{aligned}$$

- Panel derecho figura.

# Hiperplano separador óptimo



# Hiperplano separador óptimo: clasificación imperfecta

- O aún otra forma de hacerlo es:

$$\begin{aligned} & \underset{\beta, \beta_0, \xi, M}{\text{máx}} \quad M \\ & \quad \quad \quad |\beta| = 1 \end{aligned}$$

$$y_i(x_i^T \beta + \beta_0) \geq M(1 - \xi_i)$$

donde  $\xi_i \geq 0$  y  $\sum_i \xi_i \leq \text{Const}$ ,  $\text{Const}$  es una constante.

- El primero es más natural pero no resulta ser un problema convexo. El segundo sí lo es.
- Vamos a ver que, en la versión convexa del problema:
  - 1 Cuando  $\xi > 1$  hay un error de clasificación.
  - 2 Cuando  $0 \leq \xi \leq 1$  el punto está a una distancia del hiperplano menor o igual al margen.

# Hiperplano separador óptimo: clasificación imperfecta

- O aún otra forma de hacerlo es:

$$\begin{aligned} & \max_{\beta, \beta_0, \xi, M} M \\ & \|\beta\| = 1 \\ & y_i(x_i^T \beta + \beta_0) \geq M(1 - \xi_i) \end{aligned}$$

donde  $\xi_i \geq 0$  y  $\sum_i \xi_i \leq \text{Const}$ ,  $\text{Const}$  es una constante.

- El primero es más natural pero no resulta ser un problema convexo. El segundo sí lo es.
- Vamos a ver que, en la versión convexa del problema:
  - 1 Cuando  $\xi > 1$  hay un error de clasificación.
  - 2 Cuando  $0 \leq \xi \leq 1$  el punto está a una distancia del hiperplano menor o igual al margen.

# Hiperplano separador óptimo: clasificación imperfecta

- O aún otra forma de hacerlo es:

$$\begin{aligned} & \max_{\beta, \beta_0, \xi, M} M \\ & |\beta| = 1 \\ & y_i(x_i^T \beta + \beta_0) \geq M(1 - \xi_i) \end{aligned}$$

donde  $\xi_i \geq 0$  y  $\sum_i \xi_i \leq \text{Const}$ ,  $\text{Const}$  es una constante.

- El primero es más natural pero no resulta ser un problema convexo. El segundo sí lo es.
- Vamos a ver que, en la versión convexa del problema:
  - 1 Cuando  $\xi > 1$  hay un error de clasificación.
  - 2 Cuando  $0 \leq \xi \leq 1$  el punto está a una distancia del hiperplano menor o igual al margen.

- Un argumento similar al presentado anteriormente para el caso de separación perfecta muestra que el problema es equivalente a:

$$\begin{aligned} & \min_{\beta, \beta_0} |\beta| \\ & y_i(x_i^T \beta + \beta_0) \geq 1 - \xi \\ & \xi \geq 0 \\ & \sum_i \xi_i \leq Const \end{aligned}$$

donde  $M = \frac{1}{|\beta|}$  y  $Const$  es una constante.

- Planteamos el problema como:

$$\begin{aligned} \min_{\beta, \beta_0} \quad & \frac{1}{2} (\|\beta\|)^2 + C \sum_{i=1}^n \xi_i \\ & y_i (x_i^T \beta + \beta_0) \geq 1 - \xi \\ & \xi \geq 0 \end{aligned}$$

donde  $C$  es una constante de penalización.

# Hiperplano separador óptimo: CPO

- Al hacer las condiciones de primer orden con respecto a  $\beta, \beta_0, \xi_i$  se puede demostrar que:

$$\beta_i = \sum_{i=1}^n \alpha_i y_i x_i$$

y  $\alpha$  es distinto de cero únicamente cuando:

$$y_i(x_i^T \beta + \beta_0) - (1 - \xi_i) = 0$$

es decir, cuando el ejemplo  $i$  está justo en el **márgen o en lado equivocado dependiendo de  $\xi_i$** . Estos son los vectores de soporte en este caso.

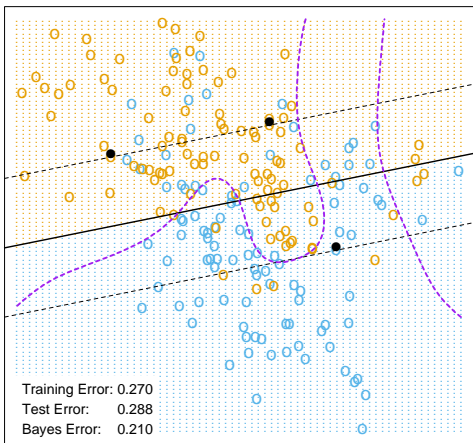
# Hiperplano separador óptimo

- La función de aprendizaje es:

$$G(x) = \text{sign}[x^T \beta + \beta_0]$$

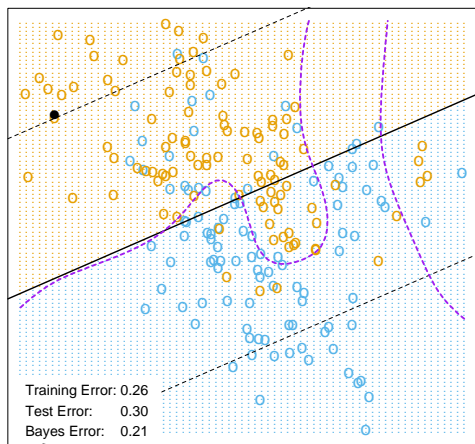
donde  $C$  es un meta parámetro.

# SVM, $C=10.000$



$C = 10000$

# SVM, $C=0.01$



$C = 0.01$

# Contenido

- 1 Introducción
- 2 Hiperplano separador óptimo
- 3 Máquinas de vectores de soporte
- 4 SVM y Kernels**

# Lagrangiano dual de Wolfe

- Una forma equivalente al problema de Lagrange introducido anteriormente (para el problema de separación imperfecta o SVM) es usar el Lagrangiano dual de Wolfe del problema (véase The Elements of Statistical Learning, página 420):

$$L_D = \sum_{i=1}^N \alpha_i - \frac{1}{2} \sum_{i,j=1}^N \alpha_i \alpha_j y_i y_j x_i^T x_j$$

- En este caso la función  $G(x)$  se puede escribir como:

$$G(x) = \text{sign}\left[x^T \sum_i^N \alpha_i y_i x_i^T + \beta_0\right]$$

- Esta representación sugiere una generalización adicional.

# Lagrangiano dual de Wolfe

- Una forma equivalente al problema de Lagrange introducido anteriormente (para el problema de separación imperfecta o SVM) es usar el Lagrangiano dual de Wolfe del problema (véase The Elements of Statistical Learning, página 420):

$$L_D = \sum_{i=1}^N \alpha_i - \frac{1}{2} \sum_{i,j=1}^N \alpha_i \alpha_j y_i y_j x_i^T x_j$$

- En este caso la función  $G(x)$  se puede escribir como:

$$G(x) = \text{sign}\left[x^T \sum_i^N \alpha_i y_i x_i^T + \beta_0\right]$$

- Esta representación sugiere una generalización adicional.

# Lagrangiano dual de Wolfe

- Una forma equivalente al problema de Lagrange introducido anteriormente (para el problema de separación imperfecta o SVM) es usar el Lagrangiano dual de Wolfe del problema (véase The Elements of Statistical Learning, página 420):

$$L_D = \sum_{i=1}^N \alpha_i - \frac{1}{2} \sum_{i,j=1}^N \alpha_i \alpha_j y_i y_j x_i^T x_j$$

- En este caso la función  $G(x)$  se puede escribir como:

$$G(x) = \text{sign}\left[x^T \sum_i^N \alpha_i y_i x_i^T + \beta_0\right]$$

- Esta representación sugiere una generalización adicional.

- Considere ahora el Lagrangiano

$$L_D = \sum_{i=1}^N \alpha_i - \frac{1}{2} \sum_{i,j=1}^N \alpha_i \alpha_j y_i y_j \langle h(x_i), h(x_j) \rangle$$

donde  $h$  son funciones de  $R^p$  en  $R^d$  donde  $d$  puede ser  $\infty$ .

- En este caso la función  $f(x)$  se puede escribir como:

$$G(x) = \text{sign} \left[ \sum_{i=1}^N \alpha_i y_i \langle h(x), h(x_i) \rangle + \beta_0 \right]$$

donde  $\langle , \rangle$  denotan un producto punto.

- Considere ahora el Lagrangiano

$$L_D = \sum_{i=1}^N \alpha_i - \frac{1}{2} \sum_{i,j=1}^N \alpha_i \alpha_j y_i y_j \langle h(x_i), h(x_j) \rangle$$

donde  $h$  son funciones de  $R^p$  en  $R^d$  donde  $d$  puede ser  $\infty$ .

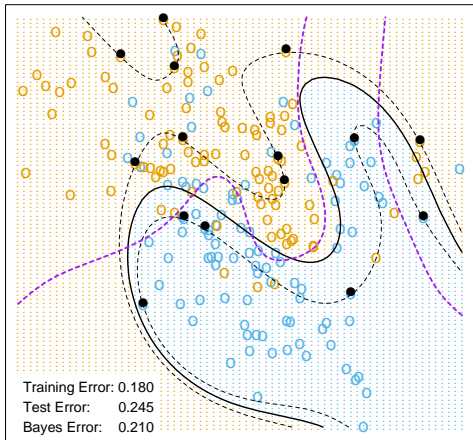
- En este caso la función  $f(x)$  se puede escribir como:

$$G(x) = \text{sign} \left[ \sum_{i=1}^N \alpha_i y_i \langle h(x), h(x_i) \rangle + \beta_0 \right]$$

donde  $\langle , \rangle$  denotan un producto punto.

- Luego, tanto el Lagrangiano como la función  $G$  solo dependen de el producto interno.
- Defina la función de Kernels:  $K(x, y) = \langle h(x), h(y) \rangle$ .
- Los Kernels más utilizados son:
  - 1 Polinomio:  $K(x, y) = (1 + \langle x, y \rangle)^d$
  - 2 Radial basis:  $K(x, y) = \exp(-\gamma(x - y)^2)$
  - 3 Red neuronal  $K(x, y) = \tanh(-\kappa_1 \langle x, y \rangle) + \kappa_2$

SVM - Degree-4 Polynomial in Feature Space



# Máquinas de vectores de soporte

LR - Radial Kernel in Feature Space

