

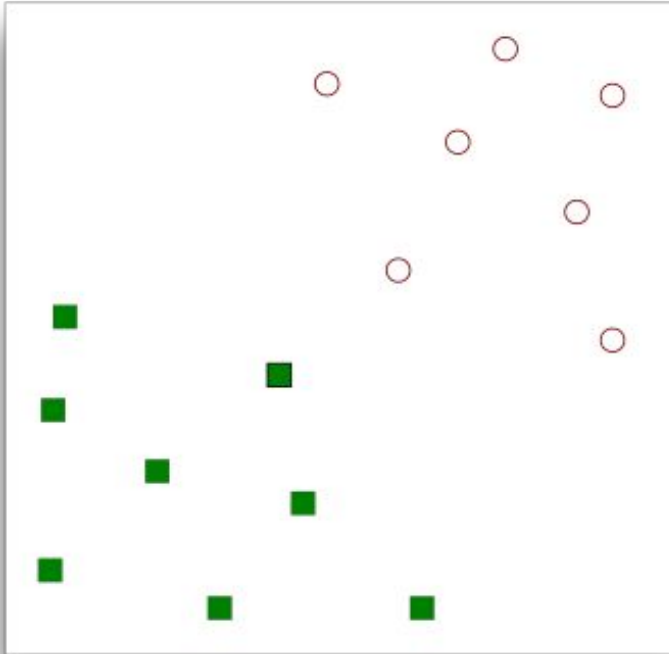
Support Vector Machines

Se formula el problema de clasificación como un problema de optimización.

Encontrar un hiperplano lineal (decision boundary).

- Representa una frontera de decisión usando un subconjunto de ejemplos de entrenamiento, conocidos como support vectors.
- Trabaja bien con datos de alta dimensionalidad y evita la "curse of dimensionality"

Support Vector Machines

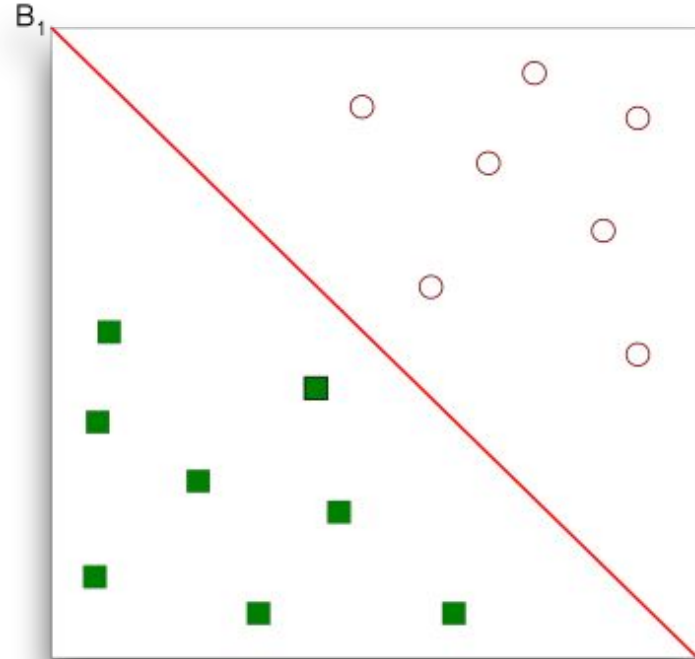
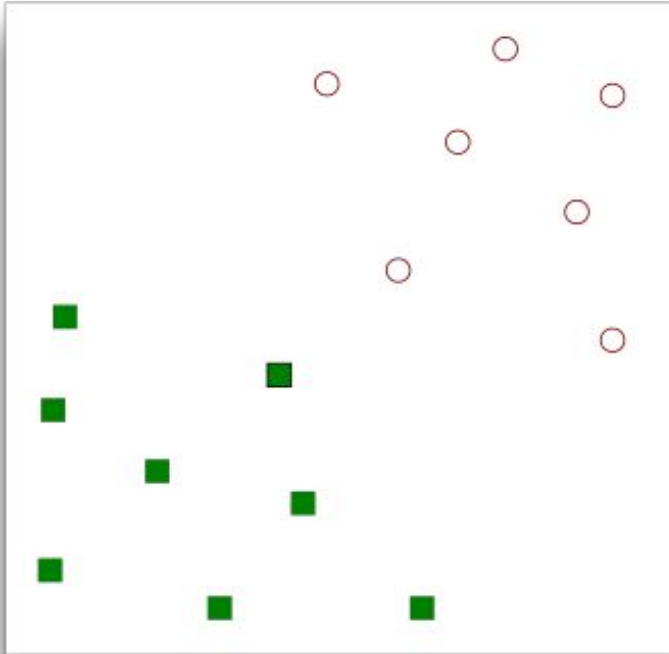


Ejemplo: Tenemos un dataset de n -dimensiones (proyectado a 2).

Podemos decir que estos datos son linealmente separables (separables por un hiperplano)

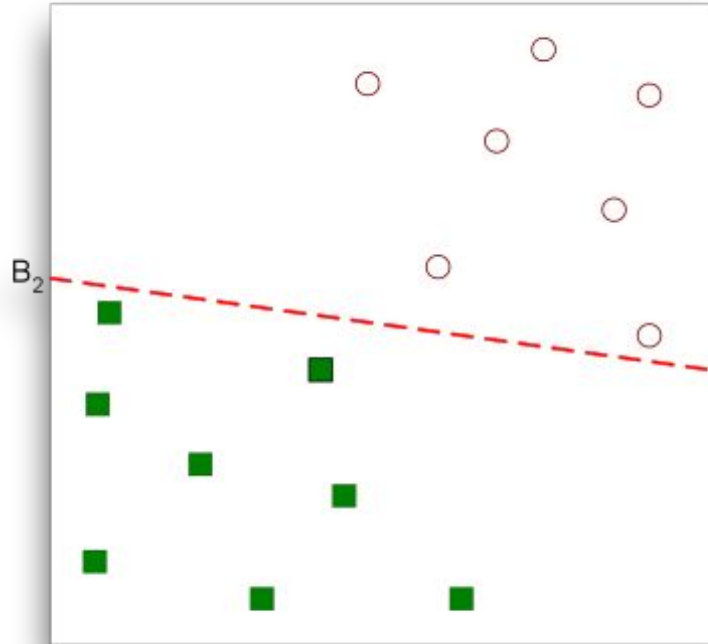
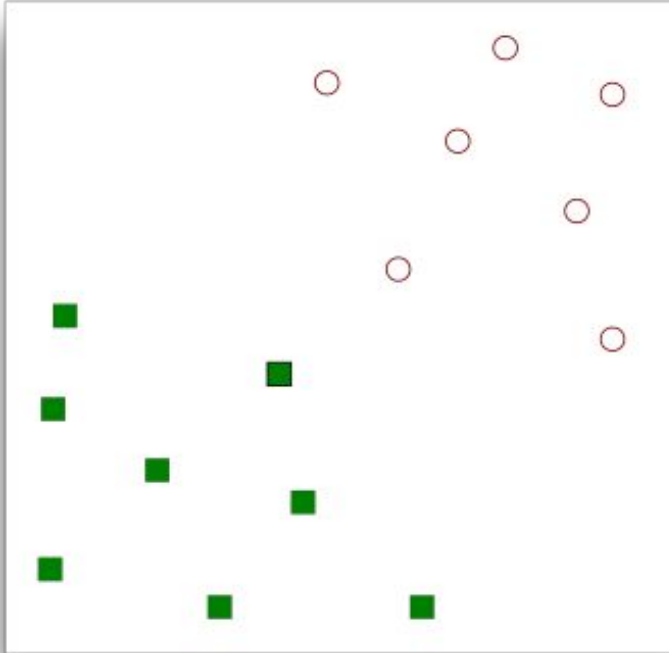
- Encontrar un hiperplano lineal (decision boundary) que separe los ejemplos positivos de los negativos.

Support Vector Machines



Una posible solución

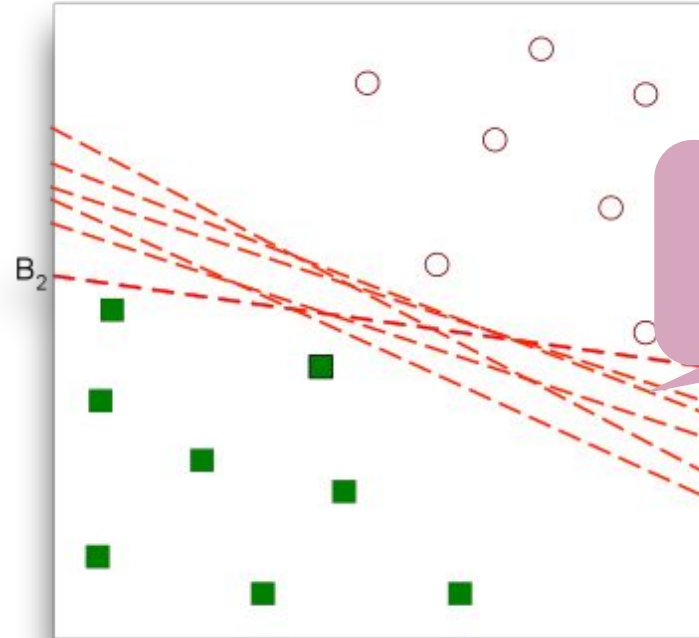
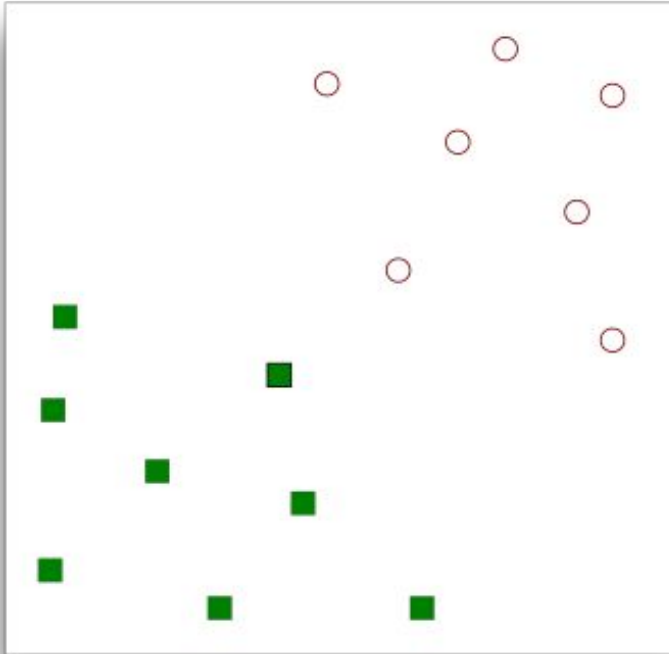
Support Vector Machines



Otra posible solución

Support Vector Machines

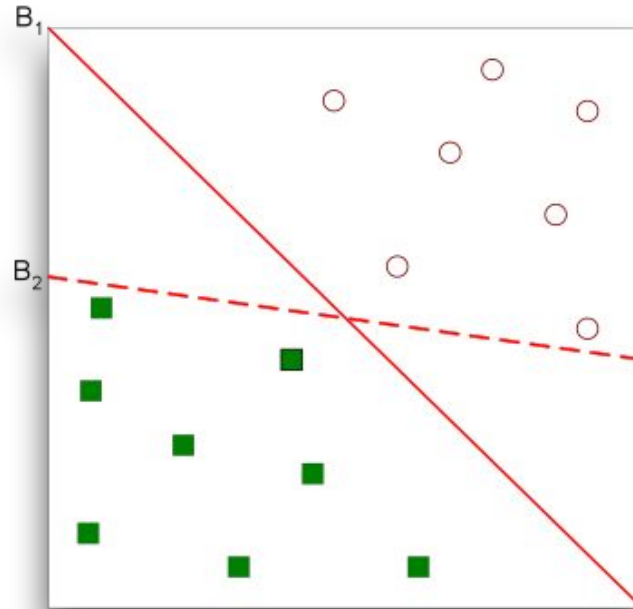
Este modelo puede extenderse a datos que no son linealmente separables.



SVM se basa en aprender a encontrar el plano marginal maximal

¡Infinitas soluciones!

Support Vector Machines



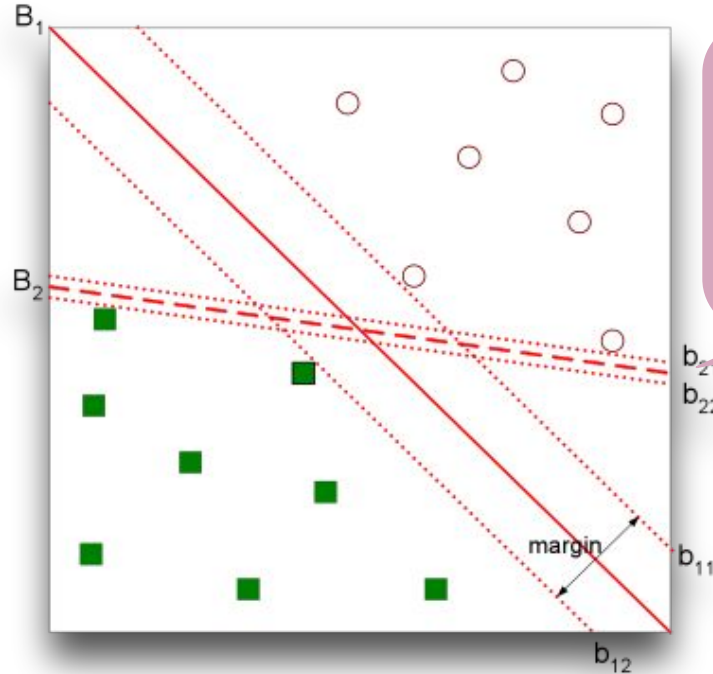
¿Cuál es mejor? ¿ B_1 o B_2 ?
¿Cómo definimos qué es mejor?

Support Vector Machines

Encontrar un hiperplano que **maximice** el margen de entrenamiento (menores errores de generalización)

- El margen del hiperplano es la distancia que hay de los puntos positivos y negativos más cercanos al hiperplano.

Entre más ancho el margen, mayor el poder de generalización.
=> B1 es mejor que B2



Menor margen hace que sea más susceptible a ruido y perturbaciones, lo que puede llevar a overfitting y mala generalización.

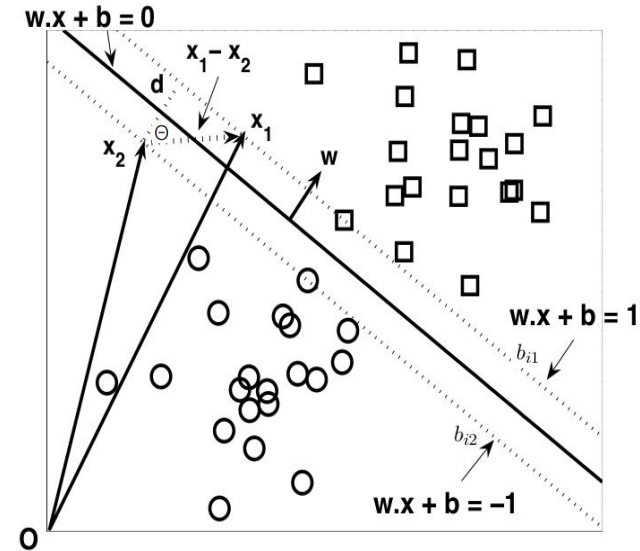
Support Vector Machines

Un SVM lineal es un clasificador que busca un hiperplano con el máximo margen.

- Sea un problema de clasificación binaria con N ejemplos, cada ejemplo se representa como la tupa x_i, y_i ($i=1,2,\dots,N$), donde x_i es un vector de d dimensiones $(x_{i1}, x_{i2}, \dots, x_{id})^T$ e $y_i \in \{-1, 1\}$ (ejemplos negativos y positivos).
- El límite de decisión de un clasificador lineal se escribe de la siguiente manera:

$$\mathbf{w} \cdot \mathbf{x} + b = 0,$$

- Donde w y b son parámetros del modelo.



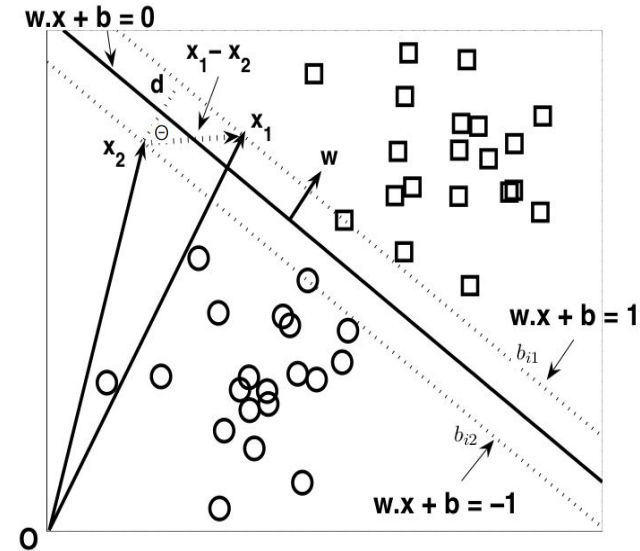
Support Vector Machines

Consiste en:

1. Formular el margen en función de los parámetros.
2. Formular un problema de optimización que permita encontrar el hiperplano de máximo margen.

Support vector machine (SVM) a detalle. Prof. Felipe Bravo.

https://www.youtube.com/watch?v=P_ArDrCQSM



Support Vector Machines

- El **hiperplano** $\mathbf{w} \cdot \mathbf{x} + b = 0$ separa los ejemplos positivos (cuadrados) de los negativos (círculos)
- Cualquier ejemplo que se encuentre en el límite de decisión debe satisfacer la ecuación del hiperplano.
- Por ejemplo, sean \mathbf{x}_a y \mathbf{x}_b dos puntos localizados en el límite de decisión:

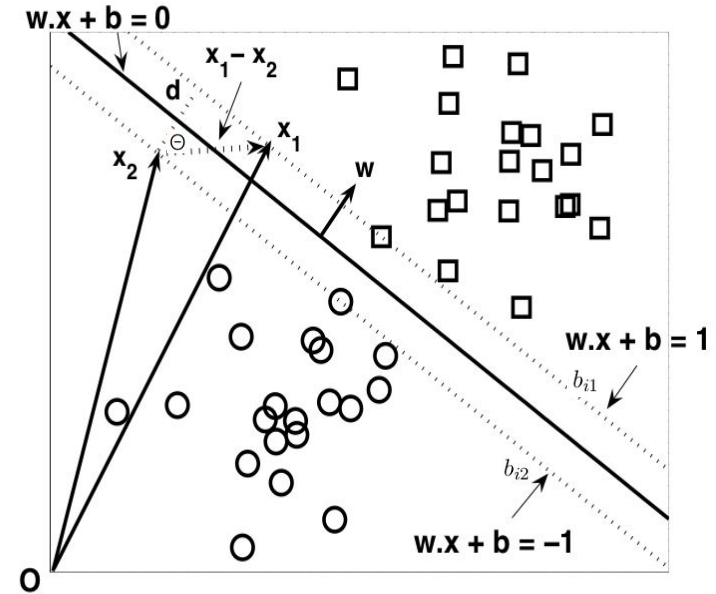
$$\mathbf{w} \cdot \mathbf{x}_a + b = 0,$$

$$\mathbf{w} \cdot \mathbf{x}_b + b = 0.$$

- Si restamos las dos ecuaciones obtenemos:

$$\mathbf{w} \cdot (\mathbf{x}_b - \mathbf{x}_a) = 0,$$

- Cómo $\mathbf{x}_b - \mathbf{x}_a$ es paralelo al límite de decisión, entonces \mathbf{w} es perpendicular al hiperplano (el producto punto de dos vectores ortogonales es cero).



Support Vector Machines

- Para cualquier cuadrado x_s localizado **sobre el límite de decisión**, se cumple que:

$$\mathbf{w} \cdot \mathbf{x}_s + b = k,$$

donde $k > 0$.

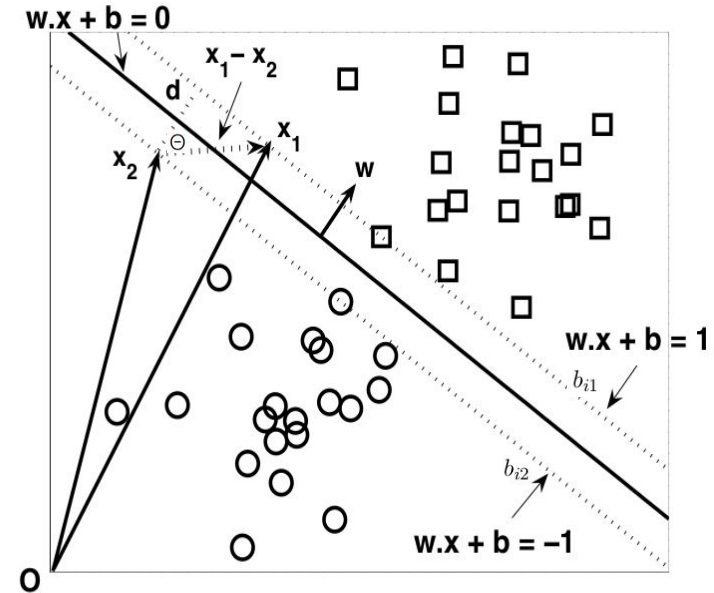
- De manera análoga, para cualquier círculo x_c localizado **bajo el límite de decisión**, se cumple que:

$$\mathbf{w} \cdot \mathbf{x}_c + b = k',$$

donde $k' < 0$.

- Si etiquetamos todos los cuadrados como la clase positiva +1 y todos los círculos como la clase negativa -1, podemos predecir la etiqueta \mathbf{y} para cualquier ejemplo de test \mathbf{z} de la siguiente manera:

$$y = \begin{cases} 1, & \text{if } \mathbf{w} \cdot \mathbf{z} + b > 0; \\ -1, & \text{if } \mathbf{w} \cdot \mathbf{z} + b < 0. \end{cases}$$



El Margen de un Clasificador Lineal

- Podemos re-escalar los parámetros w y b del límite de decisión de tal manera que los hiperplanos paralelos b_{i1} y b_{i2} puedan expresarse de la siguiente manera:

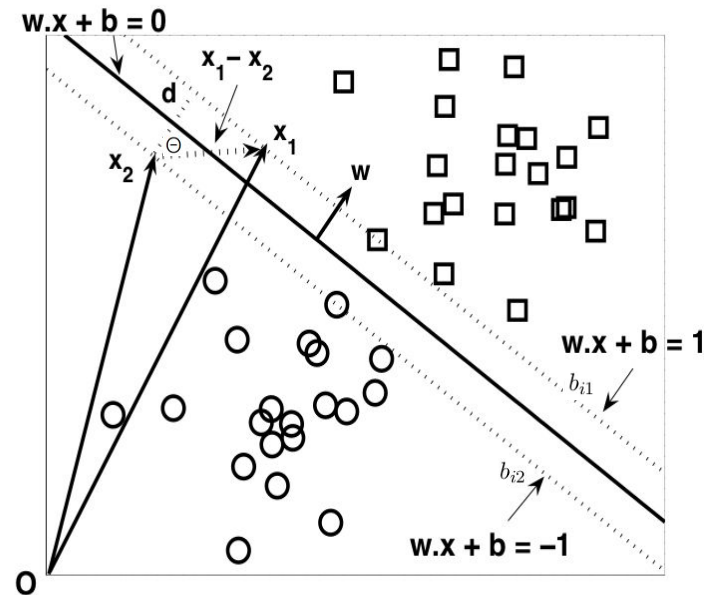
$$b_{i1} : \mathbf{w} \cdot \mathbf{x} + b = 1,$$

$$b_{i2} : \mathbf{w} \cdot \mathbf{x} + b = -1.$$

- El margen del límite de decisión se calcula como la distancia entre estos dos hiperplanos.
- Calculamos la distancia del círculo x_1 y el cuadrado x_2 . Esto se hace sustituyendo x_1 y x_2 en las ecuaciones de b_{i1} y b_{i2} respectivamente. Lo que nos da:

$$1) \mathbf{w} \cdot \mathbf{x}_1 + b = 1 \text{ y } 2) \mathbf{w} \cdot \mathbf{x}_2 + b = -1.$$

- Si sustraemos la segunda ecuación de la primera obtenemos $\mathbf{w} \cdot (\mathbf{x}_1 - \mathbf{x}_2) = 2$

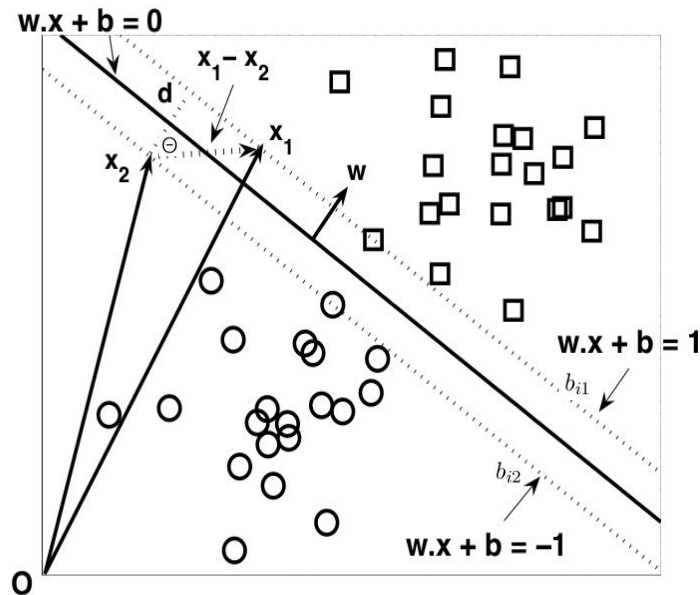


El Margen de un Clasificador Lineal

- El producto punto de dos vectores ($a \cdot b$) se puede representar como la norma $\|a\| \cdot \|b\| \cdot \cos(\Theta)$.
- Tenemos entonces que $\|w\| \cdot \|x_1 - x_2\| \cdot \cos(\Theta) = 2$
- Si miramos la figura anterior podemos ver que $\|x_1 - x_2\| \cdot \cos(\Theta) = d$.
- Entonces

$$\|w\| \times d = 2$$
$$\therefore d = \frac{2}{\|w\|}.$$

Queremos maximizar este margen



¡ Encontramos una expresión del margen que depende de w !

Aprendiendo una SVM Lineal

La fase de **entrenamiento de una SVM lineal** implica la estimación de los parámetros w y b del límite de decisión a partir de los datos de entrenamiento.

Los parámetros deben elegirse de manera que se cumplan las dos siguientes condiciones:

$$\begin{aligned} \mathbf{w} \cdot \mathbf{x}_i + b &\geq 1 \text{ if } y_i = 1, \\ \mathbf{w} \cdot \mathbf{x}_i + b &\leq -1 \text{ if } y_i = -1. \end{aligned}$$

- Estas condiciones imponen el requisito de que todas las instancias de entrenamiento de la clase $y = 1$ (los cuadrados) deben estar situadas en o **sobre el hiperplano $\mathbf{w} \cdot \mathbf{x} + b = 1$** .
- Mientras que las instancias de la clase $y = -1$ (los círculos) deben estar situadas **en o debajo del hiperplano $\mathbf{w} \cdot \mathbf{x} + b = -1$** .

Aprendiendo una SVM Lineal

Maximizar el margen equivale a minimizar la siguiente función objetivo:

$$f(\mathbf{w}) = \frac{\|\mathbf{w}\|^2}{2}.$$

$$\max_w \frac{2}{\|w\|} \Leftrightarrow \min_w \frac{\|w\|^2}{2}$$

La tarea de aprendizaje en SVM puede formalizarse como el siguiente problema de optimización con restricciones:

$$\begin{aligned} \min_{\mathbf{w}} \quad & \frac{\|\mathbf{w}\|^2}{2} \\ \text{subject to} \quad & y_i(\mathbf{w} \cdot \mathbf{x}_i + b) \geq 1, \quad i = 1, 2, \dots, N. \end{aligned}$$

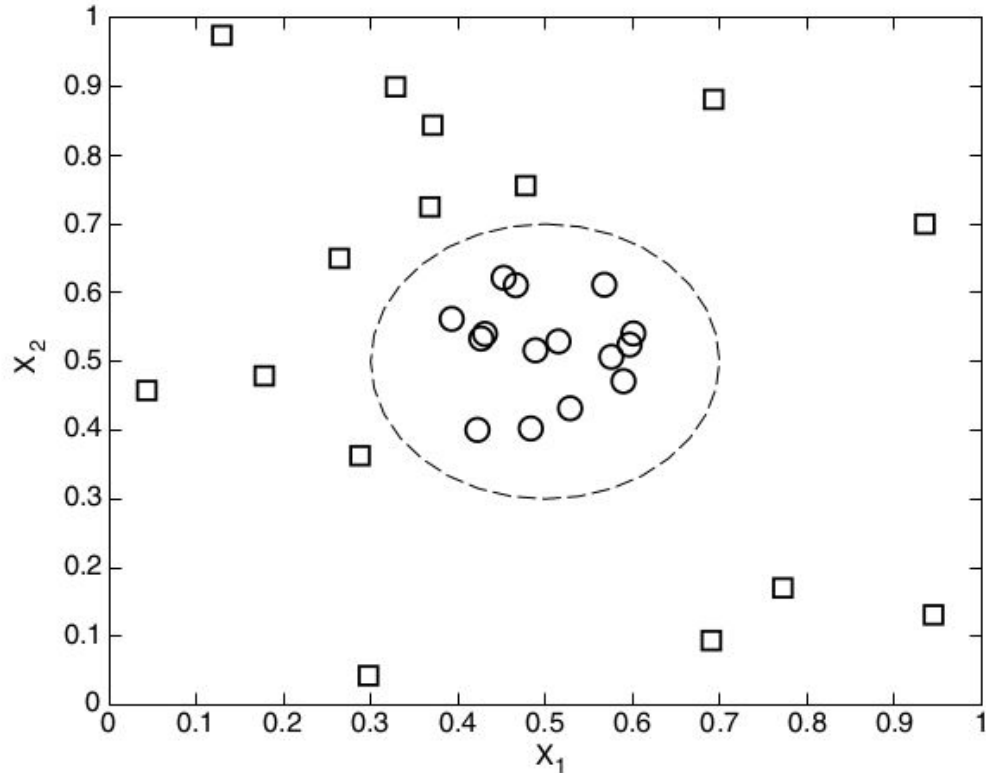
Dado que la función objetivo es cuadrática y las restricciones son lineales para los parámetros w y b , esto se conoce como un **problema de optimización convexa**.

Support Vector Machines No-Lineales

El diagrama muestra un ejemplo de un dataset bidimensional compuesto por cuadrados ($y = 1$) y círculos ($y = -1$).

Todos los círculos están agrupados cerca del centro del diagrama y todos los cuadrados se distribuyen más lejos del centro.

Este problema no se puede resolver con una SVM lineal.



Support Vector Machines No-Lineales

Podríamos aplicar una **transformación no lineal** Φ para mapear los datos de su espacio de original a un nuevo espacio (de más dimensiones) donde el límite de decisión se vuelva lineal.

Supongamos que elegimos la siguiente transformación que transforma de 2 dimensiones a 5 dimensiones:

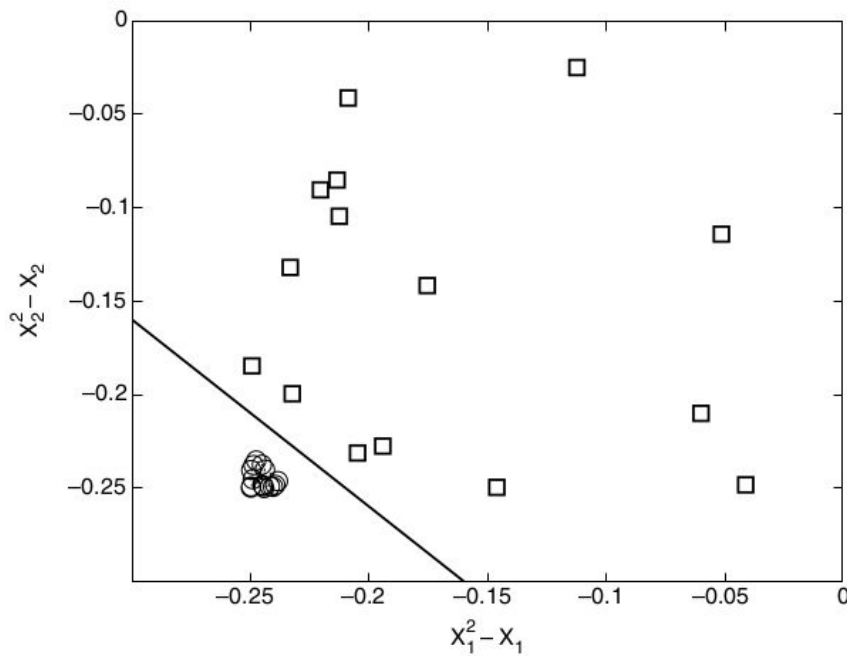
$$\Phi : (x_1, x_2) \longrightarrow (x_1^2, x_2^2, \sqrt{2}x_1, \sqrt{2}x_2, 1).$$

En este espacio transformado si es posible encontrar parámetros $\mathbf{w} = (w_0, w_1, \dots, w_4)$ que separen linealmente los datos:

$$w_4x_1^2 + w_3x_2^2 + w_2\sqrt{2}x_1 + w_1\sqrt{2}x_2 + w_0 = 0.$$

Support Vector Machines No-Lineales

- La figura muestra que en el espacio transformado se puede construir un límite de decisión lineal para separar las clases.
- Un problema potencial de este enfoque es que puede sufrir de la **maldición la dimensionalidad**: para datos de alta dimensión muchas técnicas de data mining no escalan o no funcionan bien.
- Mostraremos cómo una SVM no lineal evita este problema usando un truco llamado **Kernel Trick**.



Support Vector Machines No-Lineales

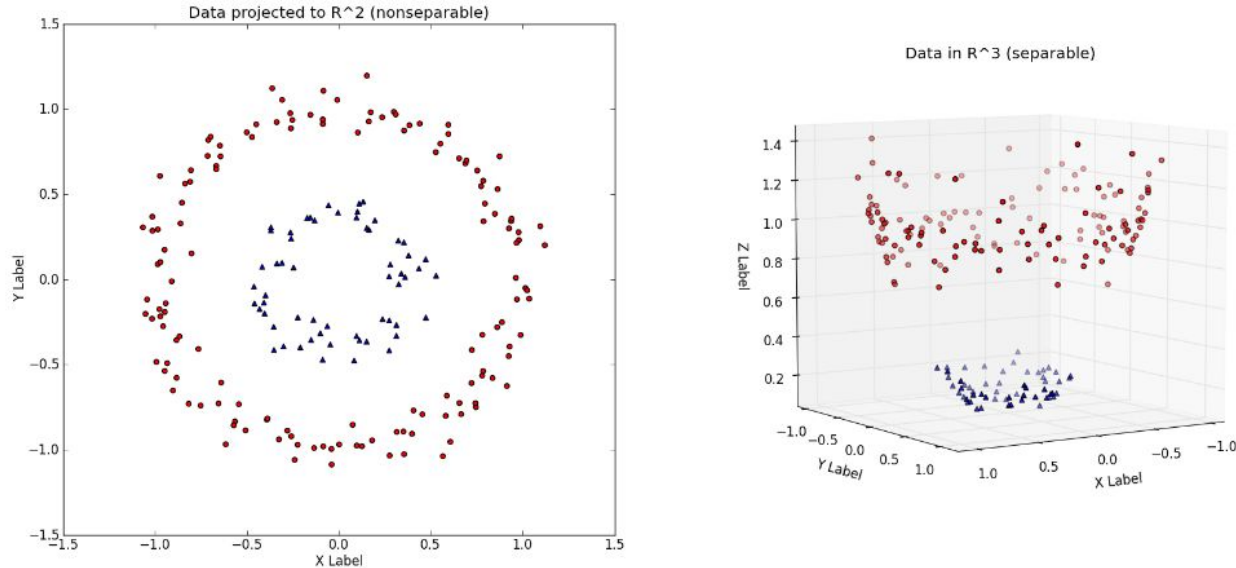


Figure 5: (Left) A dataset in \mathbb{R}^2 , not linearly separable. (Right) The same dataset transformed by the transformation:
$$[x_1, x_2] = [x_1, x_2, x_1^2 + x_2^2].$$

source: http://www.eric-kim.net/eric-kim-net/posts/1/kernel_trick.html

Support Vector Machines No-Lineales

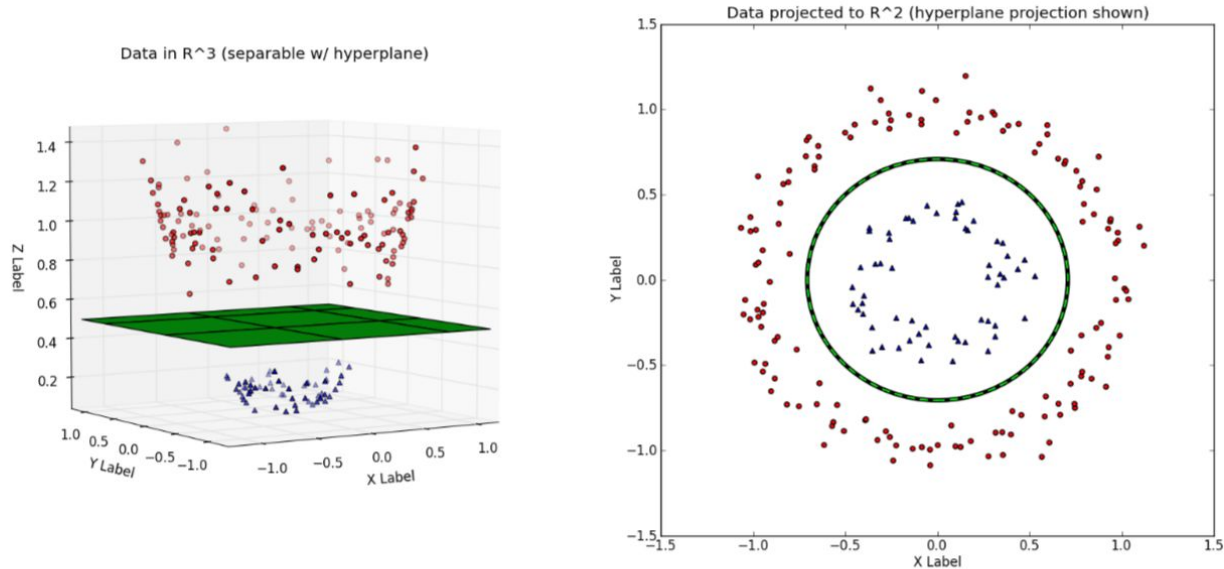
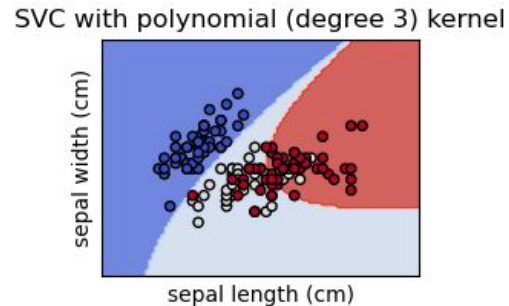
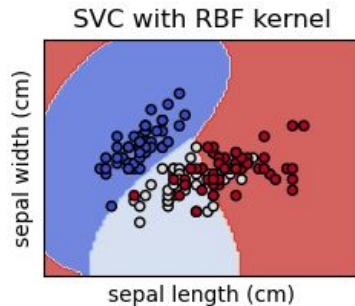


Figure 6: (Left) The decision boundary \vec{w} shown to be linear in \mathbb{R}^3 . (Right) The decision boundary \vec{w} , when transformed back to \mathbb{R}^2 , is nonlinear.

source: http://www.eric-kim.net/eric-kim-net/posts/1/kernel_trick.html

Support Vector Machines No-Lineales

- Se pueden especificar diferentes funciones de kernel para la función de decisión. Se proporcionan kernels comunes, pero también es posible especificar kernels personalizados.
 - Polinomial
 - Exponencial
 - Radial Basis Function (RBF)
 - etc.



Conclusiones

- La formulación de SVM presentada en esta clase se limita a problemas de clasificación binaria. Existen adaptaciones para trabajar con múltiples clases.
- El problema de aprendizaje de una SVM se formula como un problema de optimización convexa en donde hay algoritmos eficientes para encontrar el óptimo global. Otros métodos de clasificación como los árboles de decisión y las redes neuronales tienden a encontrar óptimos locales.
- La SVM optimiza explícitamente la capacidad de generalización al maximizar el margen del límite de decisión.
- En una SVM el usuario debe ajustar hiper-parámetros, como el tipo de función de Kernel y el costo C para las variables de holgura (esto puede ser caro).
- La gran limitación de las SVMs es que no escalan bien para datasets masivos.

Ejemplo



<https://colab.research.google.com/drive/1gIDJo0yHoTZIAboBAk4y9VGi3YyMQILU?usp=sharing>