



UNIVERSIDAD
NACIONAL
DE COLOMBIA

ANALÍTICA PREDICTIVA

CARLOS A. MADRIGAL

PROFESOR OCASIONAL

DEPARTAMENTO DE CIENCIAS DE LA COMPUTACIÓN Y DE LA DECISIÓN

MAESTRÍA EN INGENIERÍA - INGENIERÍA DE SISTEMAS

MAESTRÍA EN INGENIERÍA - ANALÍTICA

ESPECIALIZACIÓN EN SISTEMAS

CONTENIDO

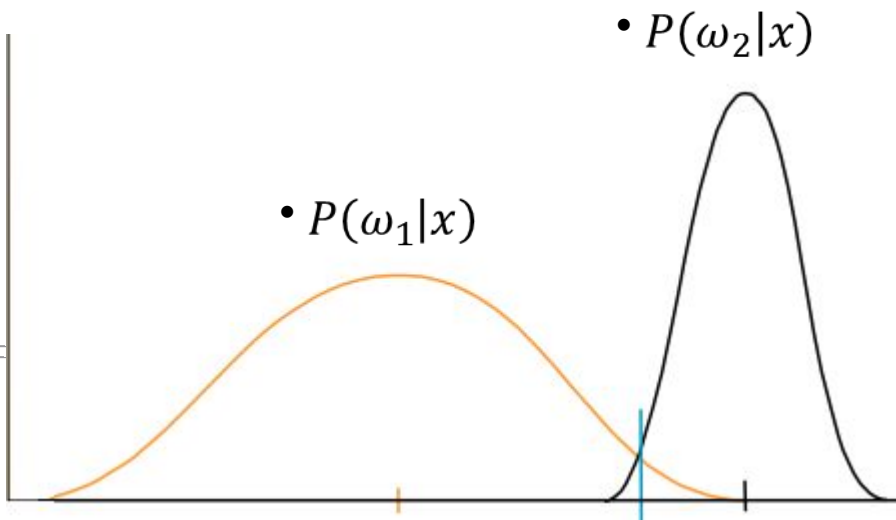
Técnicas de Clasificación y Agrupamiento

- Naive Bayes
- KNN
- K-means
- Máquina de Soporte Vectorial

NAIVE BAYES

El teorema de Bayes expresa la probabilidad a posteriori de un evento Y dado X. Usar la teoría de la probabilidad para clasificar el objeto en la clase que tenga mayor probabilidad posteriori.

$$P(\omega_i|x) = \frac{p(x|\omega_i)P(\omega_i)}{p(x)}$$



- $P(\omega_i)$ = Probabilidad de que en la población haya un objeto de clase ω_i

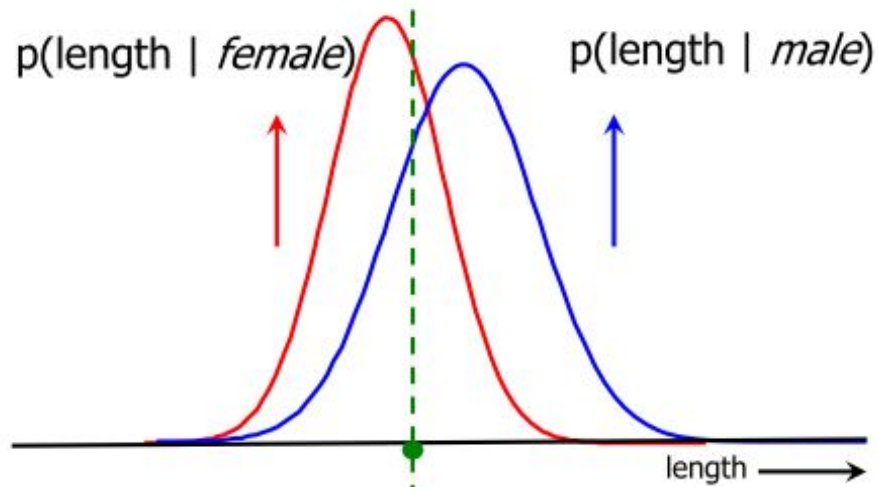
$p(x|\omega_i)$ = Probabilidad de que en la clase ω_i se de un vector de características x

$P(\omega_i|x)$ = Probabilidad de que el objeto de vector de características x pertenezca a la clase ω_i

- $g(x) = \begin{cases} 1 & \text{si } P(\omega_1|x) > P(\omega_2|x) \\ 2 & \text{en otro caso} \end{cases}$

NAIVE BAYES

- Cuál es el género de alguien con esta altura?



Bayes: $\begin{cases} p(\text{female} \mid \text{length}) = p(\text{length} \mid \text{female}) p(\text{female}) / p(\text{length}) \\ p(\text{male} \mid \text{length}) = p(\text{length} \mid \text{male}) p(\text{male}) / p(\text{length}) \end{cases}$

NAIVE BAYES

- Regla de Clasificación de Bayes:

$$p(\textit{female} \mid \textit{length}) > p(\textit{male} \mid \textit{length}) \rightarrow \textit{female} \text{ else } \textit{male}$$

Bayes:

$$\frac{p(\textit{length} \mid \textit{female}) p(\textit{female})}{p(\textit{length})} > \frac{p(\textit{length} \mid \textit{male}) p(\textit{male})}{p(\textit{length})}$$

$$p(\textit{length} \mid \textit{female}) p(\textit{female}) > p(\textit{length} \mid \textit{male}) p(\textit{male}) \rightarrow \textit{female} \text{ else } \textit{male}$$

pdf estimated from training set

class prior probabilities
known, guessed or estimated

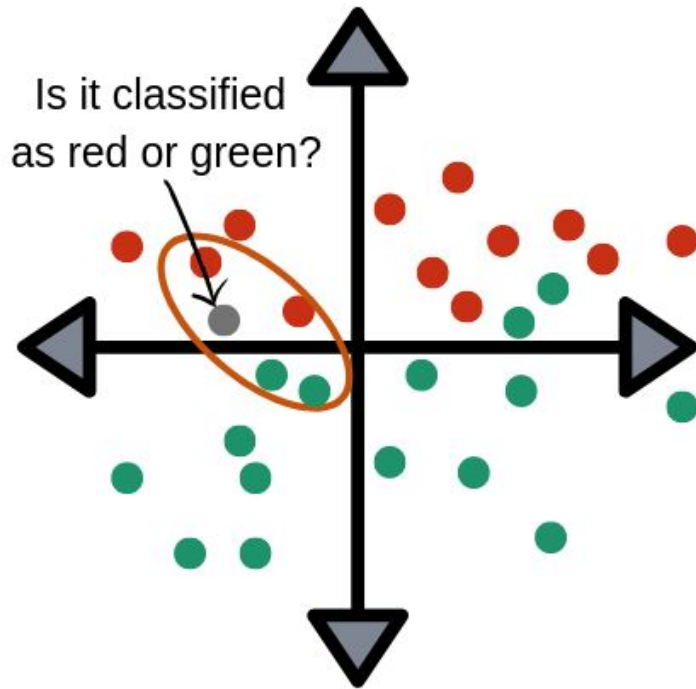
GAUSSIAN NAIVE BAYES

Cuando las características X son no categóricas, la verosimilitud se calcula como la densidad de probabilidad de una distribución normal.

$$P(x_i | y) = \frac{1}{\sqrt{2\pi\sigma_y^2}} \exp\left(-\frac{(x_i - \mu_y)^2}{2\sigma_y^2}\right)$$

CLASIFICADOR KNN

La idea básica sobre la que se fundamenta este paradigma es que un nuevo caso se va a clasificar en la clase más frecuente a la que pertenecen sus K vecinos más cercanos



COMIENZO

Entrada: $D = \{(x_1, c_1), \dots, (x_N, c_N)\}$

$x = (x_1, \dots, x_n)$ nuevo caso a clasificar

PARA todo objeto ya clasificado (x_i, c_i)

calcular $d_i = d(x_i, x)$

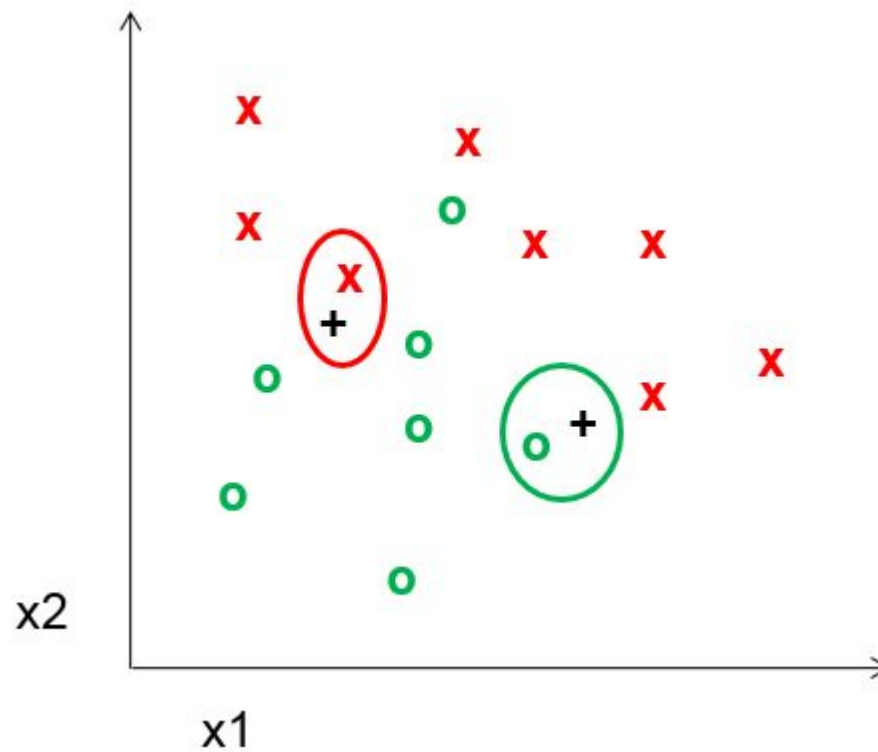
Ordenar $d_i (i = 1, \dots, N)$ en orden ascendente

Quedarnos con los K casos D_x^K ya clasificados más cercanos a x

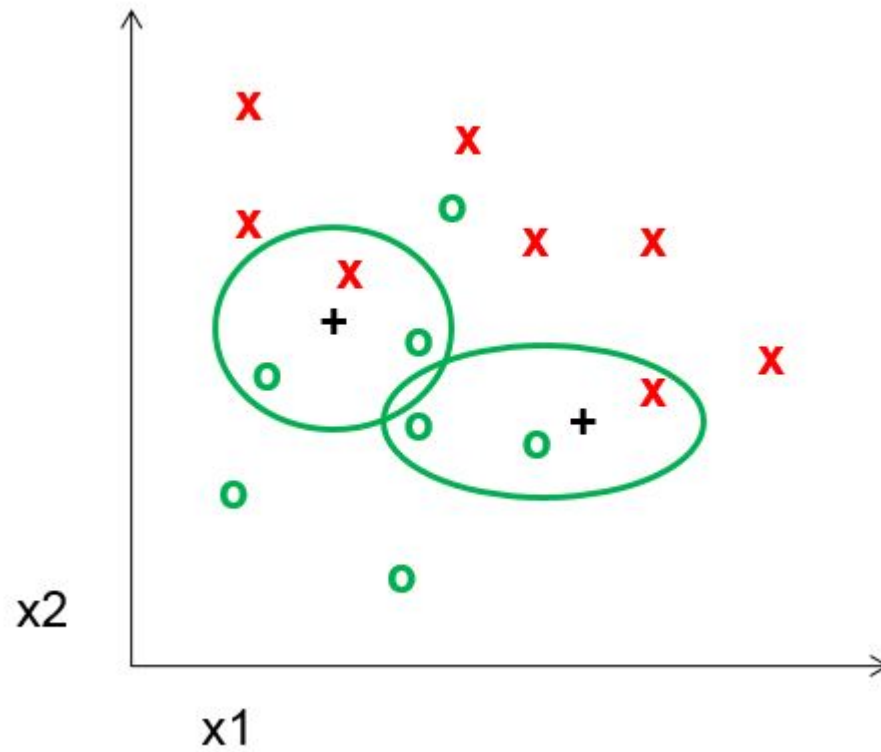
Asignar a x la clase más frecuente en D_x^K

FIN

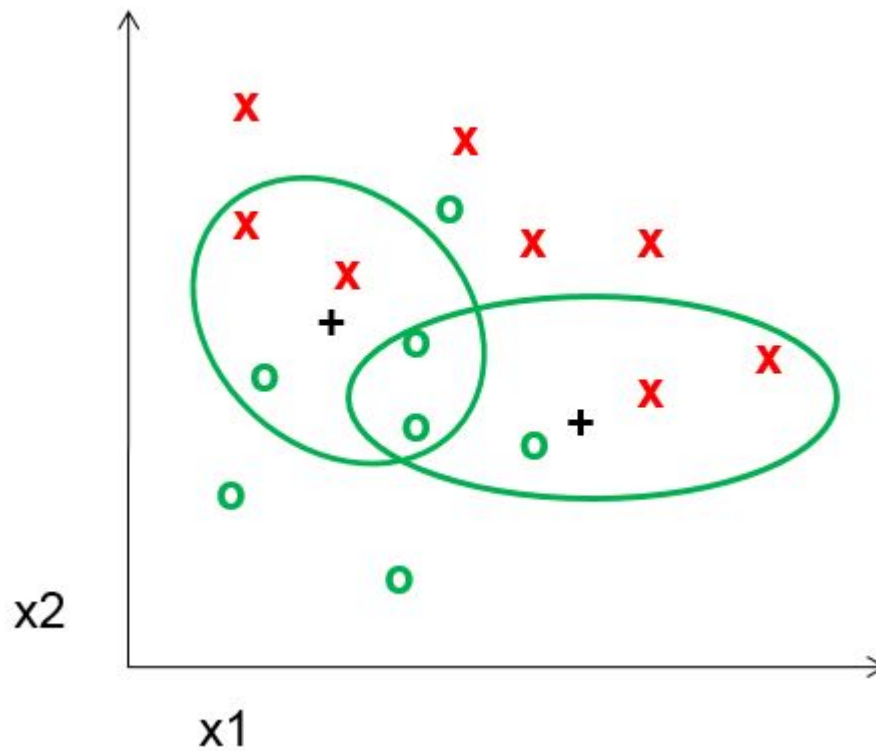
CLASIFICADOR KNN – 1 NN



CLASIFICADOR KNN – 3 NN



CLASIFICADOR KNN – 5 NN



VARIACIONES KNN

KNN- con Rechazo

Esta variación hace énfasis en las garantías que se deben dar para asignar una clase a un conjunto de características.

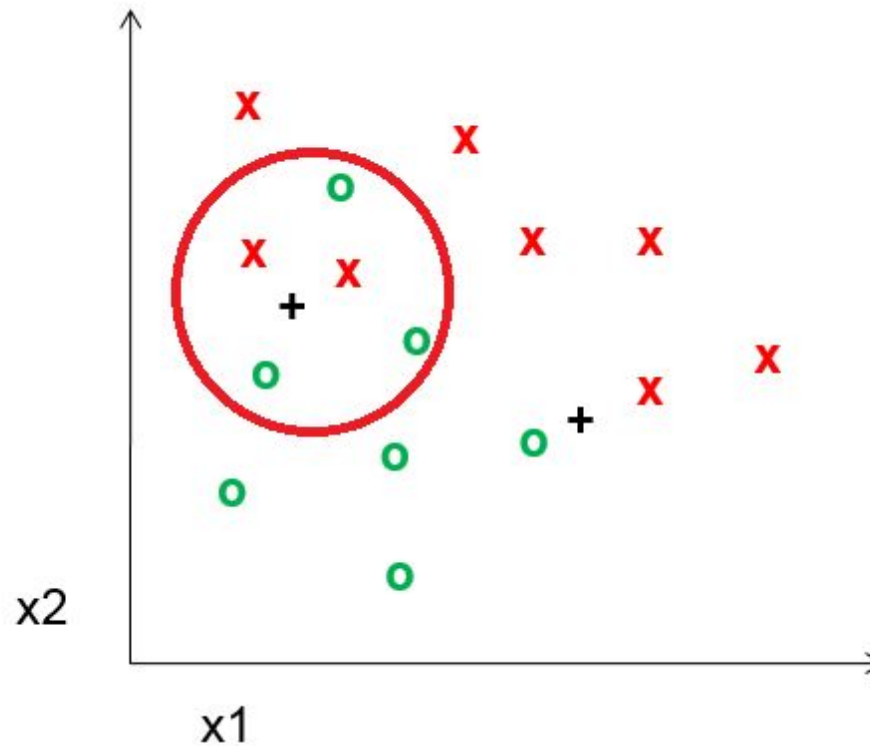
Umbral: Se refiere a que el número de votos para la clase a asignar sea superior a este valor. Ejemplo: si $K=8$, $m=2$, el umbral podría establecerse en 5 o 6.

Mayoría Absoluta: Diferencias entre la frecuencia mayor y segunda mayor supere un valor. Ejemplo: Siendo $K=15$, $m=3$, Diferencia= 3;

CLASIFICADOR KNN

KNN – con distancia Media

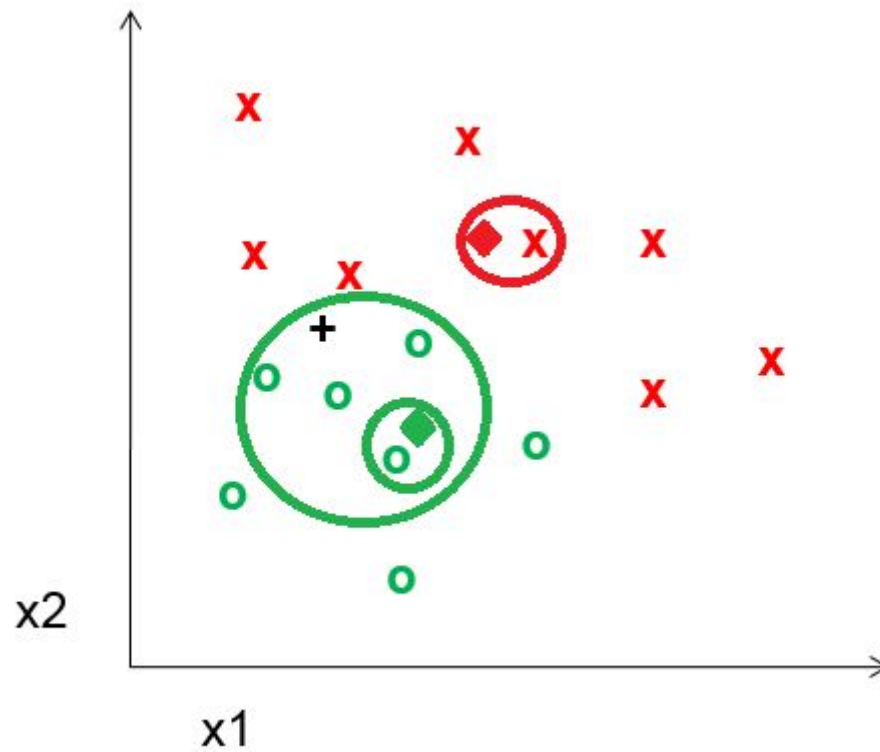
Hace referencia a la asignación de la clase con menor distancia media



CLASIFICADOR KNN

KNN – con distancia Mínima.

Se reduce el número de casos a uno por clase (baricentro) y luego se le asigna al conjunto de características la clase del baricentro más cercano.



CLASIFICADOR KNN

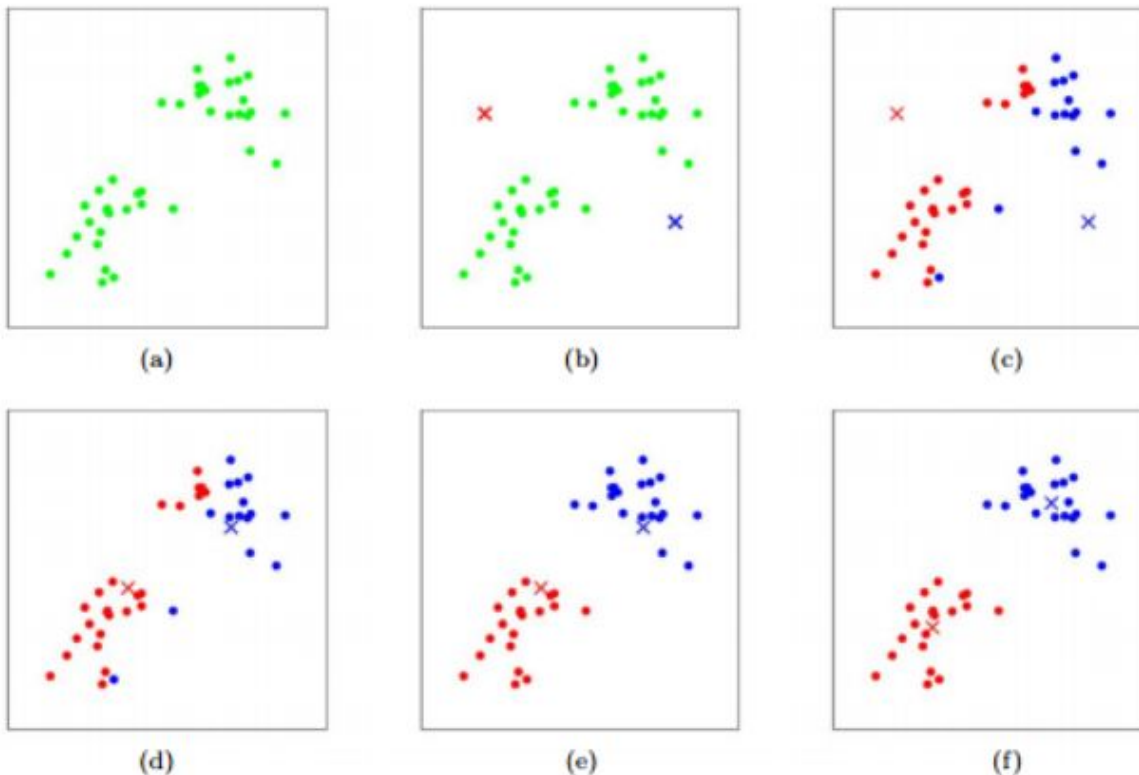
KNN – con pesado de Casos.

Hace referencia a no contabilización homogénea de los casos, sino que se genera un peso para contabilizar cada caso, por ejemplo, el inverso de la distancia entre los casos seleccionados y el nuevo caso.

$$d(\mathbf{x}, \mathbf{x}_r) = \sum_{j=1}^n w_j (x_j, x_{rj})^2$$

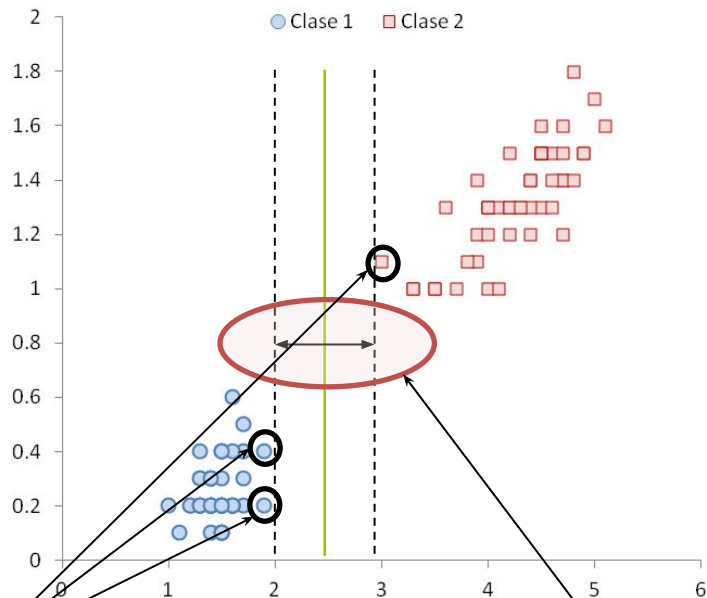
KMEANS

Permiten hacer agrupaciones entre los datos de tal manera que los casos de un cluster tengan una alta similaridad entre ellos y baja con respecto a casos de otro cluster.



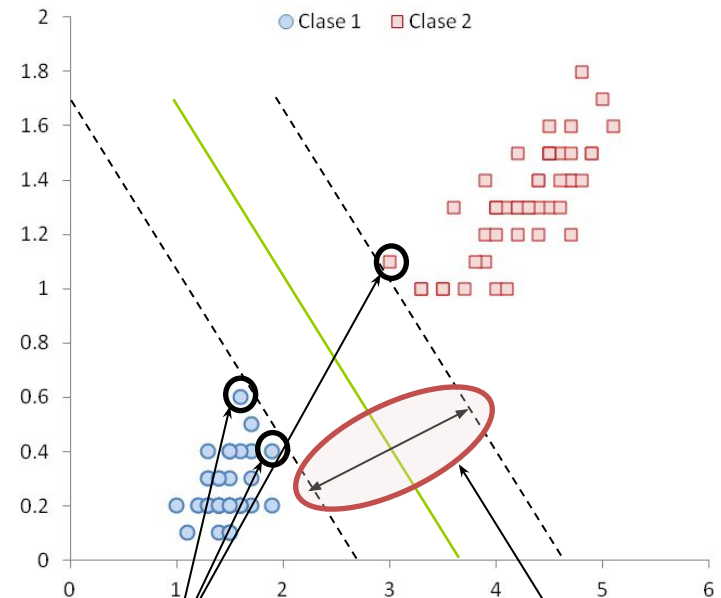
MÁQUINAS DE SOPORTE VECTORIAL

- Las **SVM** son un tipo de clasificadores de patrones basados en técnicas estadísticas de aprendizaje y están a la cabeza de los métodos de clasificación por permitir construir fronteras de decisión flexibles, y su buena capacidad de generalización.



Vectores Soporte

Margen Pequeña

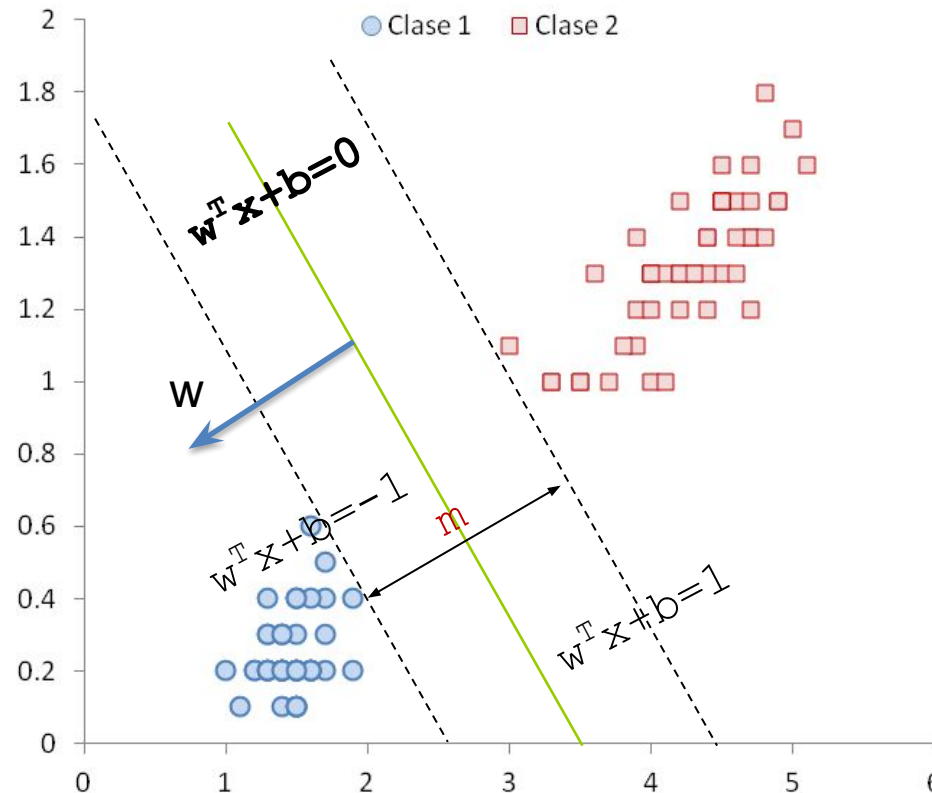


Vectores Soporte

Margen Grande

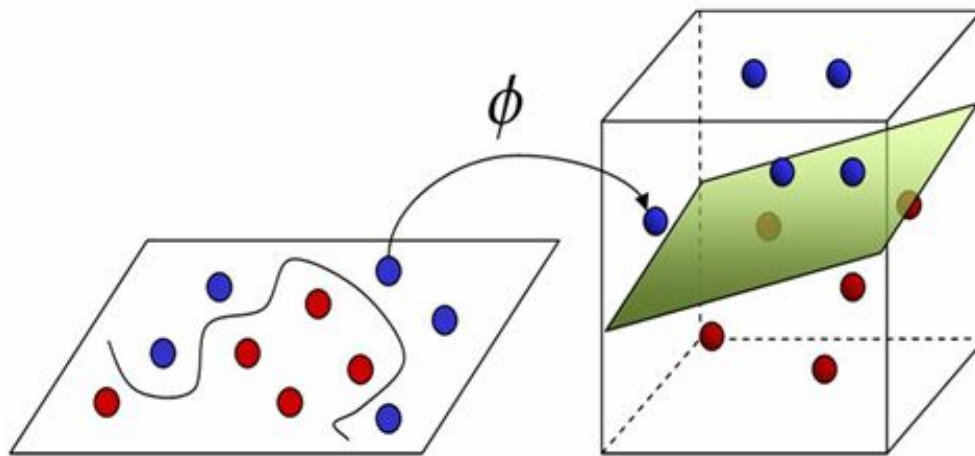
MÁQUINAS DE SOPORTE VECTORIAL

- Clasificación Lineal: Las SVM generan un hiperplano que separa el espacio en dos o más regiones, una para cada clase.



MÁQUINAS DE SOPORTE VECTORIAL

- La **Clasificación NO Lineal** con una SVM realiza una transformación del espacio de entrada a otro de dimensión más alta, en el que los datos son separables linealmente.



Lineal: $K(x_i, x_j) = x_i \cdot x_j$

Polinómico: $K(x_i, x_j) = (x_i \cdot x_j + 1)^d$

Gausiano: $K(x_i, x_j) = \exp\left(-\frac{\|x_i - x_j\|^2}{2\sigma^2}\right)$

$$\tilde{L}(\alpha) = \sum_{i=1}^n \alpha_i - \frac{1}{2} \sum_{i,j} \alpha_i \alpha_j y_i y_j \mathbf{x}_i^T \mathbf{x}_j = \sum_{i=1}^n \alpha_i - \frac{1}{2} \sum_{i,j} \alpha_i \alpha_j y_i y_j k(\mathbf{x}_i, \mathbf{x}_j)$$

- Al introducir un kernel, los parámetros α del vector w se calculan así:

Ejemplos

Naive Bayes <https://www.kaggle.com/dilip990/spam-ham-detection-using-naive-bayes-classifier>

KNN <https://www.kaggle.com/jmataya/k-nearest-neighbors-classifier>

SVM <https://www.kaggle.com/migeruj/svm-sentiment-analysis-an-lisis-de-sentimientos>

Kmeans <https://www.kaggle.com/karthickaravindan/k-means-clustering-project>
<https://www.kaggle.com/gabrielrs3/clustering-customers-k-means-algorithm/data>

PREGUNTAS





UNIVERSIDAD
NACIONAL
DE COLOMBIA