

Clasificación I

Modelos lineales de clasificación

Jonnatan Arias Garcia
jonnatan.arias@utp.edu.co
jariasg@uniquindio.edu.co

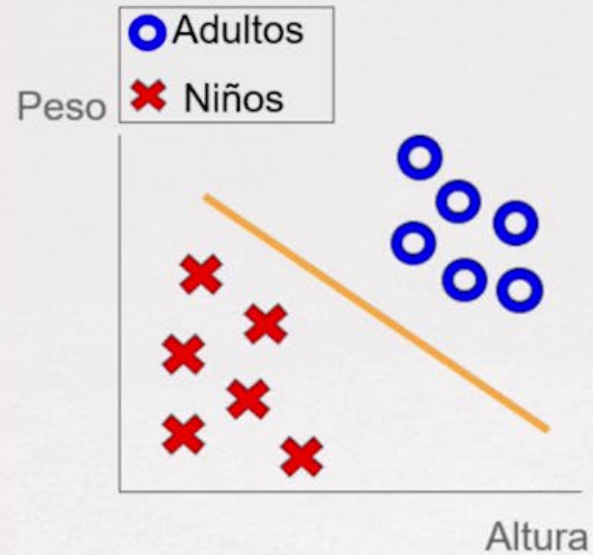
David Cardenas peña - dcardenasp@utp.edu.co
Hernán Felipe Garcia - hernanf.garcia@udea.edu.co

Contenido

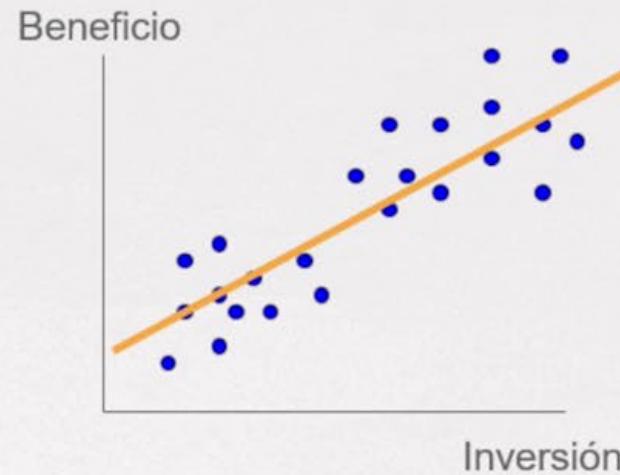
- Introducción
- Función discriminante
- Estimación de parámetros
- Modelos generativos probabilísticos

Técnicas de Machine Learning

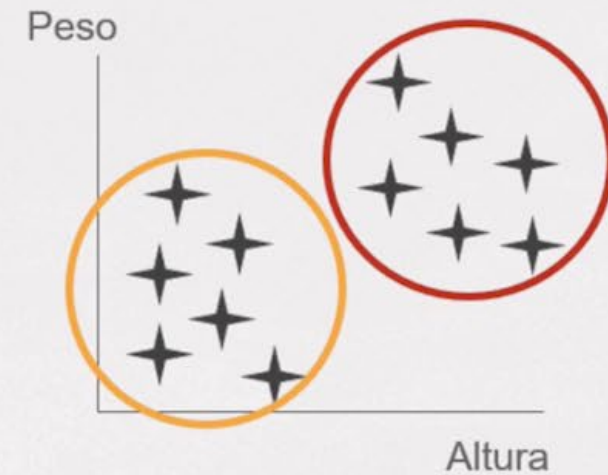
Clasificación



Regresión



Agrupación (*clustering*)

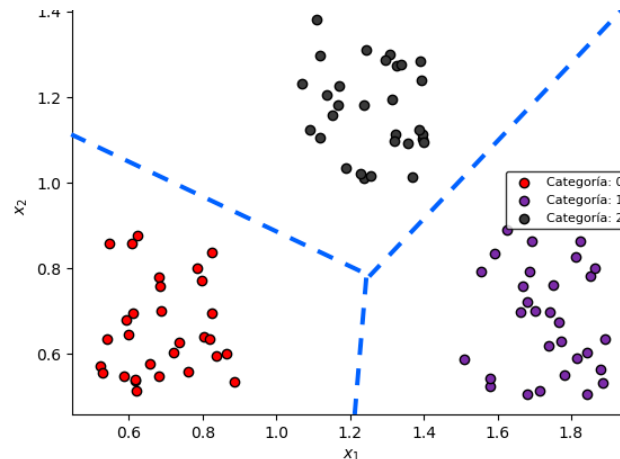


APRENDIZAJE SUPERVISADO

APRENDIZAJE NO SUPERVISADO

Definiciones (i)

- ❑ **Objetivo.** Tomar un vector de entrada, \mathbf{x} , y asignarlo a una de K clases C_k , para $k = 1, \dots, K$.
- ❑ **Espacio de entrada.** Se divide en regiones de decisión \mathcal{R}_k .
- ❑ **Líneas o superficies de decisión.** Separan las regiones de decisión.



Definiciones (ii)

- ❑ **Modelo lineal.** Las superficies de decisión son funciones lineales de \mathbf{w} .
- ❑ Las superficies están definidas por hiperplanos de dimensión $D - 1$, en un espacio de D dimensiones.
- ❑ Si los datos se pueden separar linealmente por una superficie de decisión lineal, se dice que los datos son *linealmente separables*.

Definiciones (iii)

- En regresión, t_n representaba un número real asociado a la entrada \mathbf{x}_n .

Etiqueta de
clase

- En clasificación, \mathbf{t}_n representa un vector, en codificación 1 de K .

One Hot encoding ó 1 de K: [0,0,0,1] -> 4 clases 1 es la pertenencia

- Tres enfoques al problema de clasificación
 - Funciones discriminantes.
 - Modelos generativos, $p(\mathbf{x}, \mathcal{C}_k)$.
 - Modelos discriminativos, $p(\mathcal{C}_k | \mathbf{x})$.

Función Discriminante

F. Discriminante (i)

- ❑ Función $y(\mathbf{x}): \mathbf{x} \rightarrow k, k \in \{1, \dots, K\}$, tales que la decisión de clasificación se toma como:

$$C = \operatorname{argmax}_{C_k} y_k(\mathbf{x})$$

- ❑ La clase asignada será aquella cuya función discriminante sea la mayor.

Caso Binario o biclase

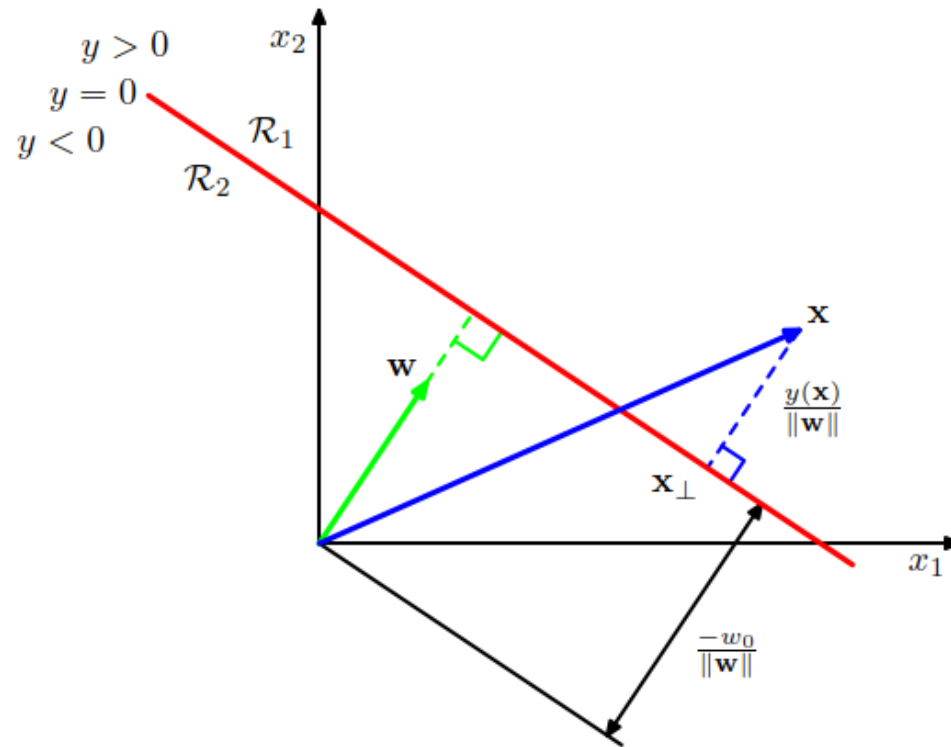
- ❑ El modelo lineal esta dado por $y(\mathbf{x}) = \mathbf{w}^T \mathbf{x} + w_0$
 - ❑ \mathbf{w} es el vector de pesos
 - ❑ w_0 es el intercepto (Bias) o tendencia
- ❑ $\mathbf{x} \in C_1$, si $y(\mathbf{x}) > 0$. De lo contrario, $\mathbf{x} \in C_2$.
- ❑ La línea de decisión o superficie de decisión es $y(\mathbf{x}) = 0$.

F. Discriminante (ii)

El vector \mathbf{w} es ortogonal a la superficie de decisión.

Distancia de $y(\mathbf{x})$ al origen: $-w_0/\|\mathbf{w}\|$.

Distancia de un punto \mathbf{x} a $y(\mathbf{x})$: $y(\mathbf{x})/\|\mathbf{w}\|$.

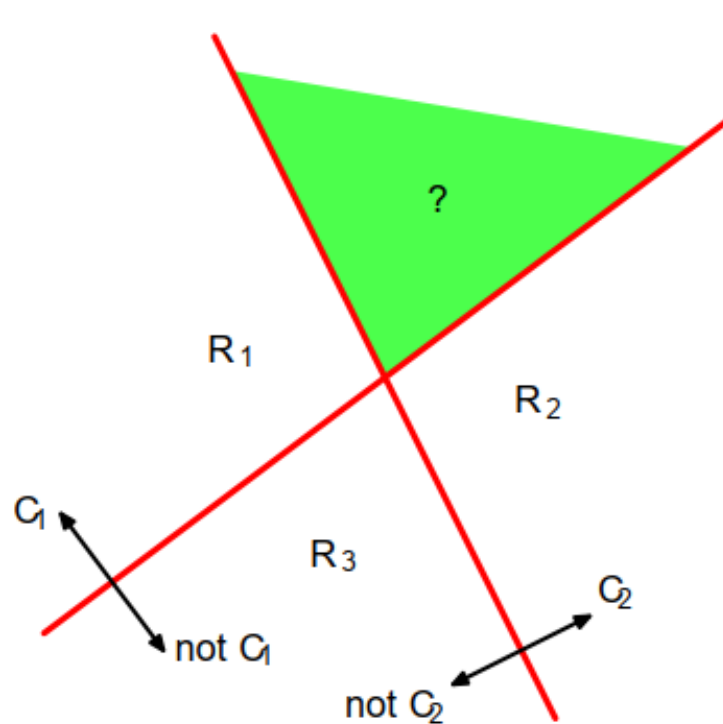


Se puede hacer $\tilde{\mathbf{x}} = (1, \mathbf{x})$, y $\tilde{\mathbf{w}} = (w_0, \mathbf{w})$, $y(\mathbf{x}) = \tilde{\mathbf{w}}^\top \tilde{\mathbf{x}}$.

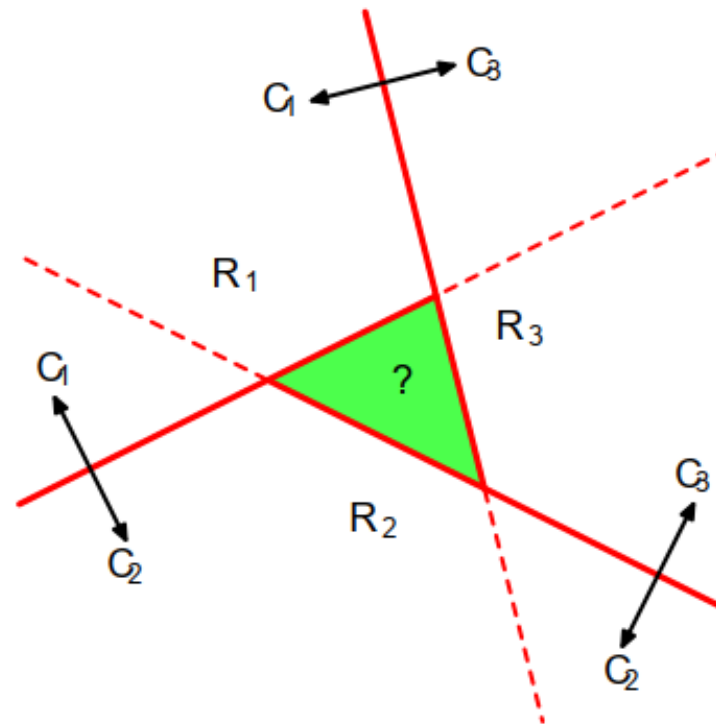
F. Discriminante (iii)

Múltiples clases (I)

Construir un clasificador de K clases a partir de clasificadores de 2 clases.



1 vs. todos



1 vs 1

Errores?

F. Discriminante (iv)

Múltiples clases (II)

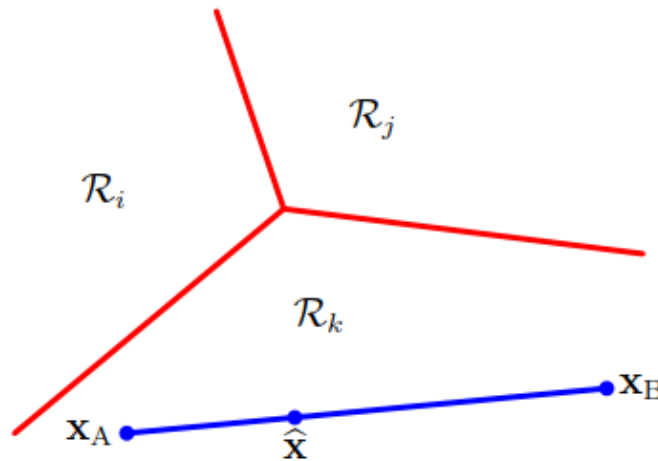
- ❑ **Solución:** discriminante de K clases con K funciones lineales

$$y_k(\mathbf{x}) = \mathbf{w}_k^\top \mathbf{x} + w_{k0}.$$

- ❑ $\mathbf{x} \in C_k$, si $y_k(\mathbf{x}) > y_j(\mathbf{x})$, $k \neq j$.
- ❑ La superficie de decisión entre las clases C_k y C_j está dada por $y_k(\mathbf{x}) = y_j(\mathbf{x})$.
- ❑ Corresponde a un plano $(D - 1)$ dimensional definido por

$$(\mathbf{w}_k - \mathbf{w}_j)^\top \mathbf{x} + (w_{k0} - w_{j0}) = 0.$$

- ❑ Regiones de decisión: conectadas de manera simple, y convexas.



Estimación de parámetros funciones discriminantes (entrenamiento)

- Por Mínimos cuadrados
- Por Fisher LDA
- Por perceptrón

Por Mínimos Cuadrado (i)

- Cada clase \mathcal{C}_k se describe por su propio modelo lineal

$$y_k(\mathbf{x}) = \mathbf{w}_k^\top \mathbf{x} + w_{k0},$$

para $k \in \{1, \dots, K\}$.

- Agrupando,

$$\mathbf{y}(\mathbf{x}) = \tilde{\mathbf{W}}^\top \tilde{\mathbf{x}},$$

donde $\tilde{\mathbf{W}}$ tiene columnas $\tilde{\mathbf{w}}_k = (w_{k0}, \mathbf{w}_k)^\top$, y $\tilde{\mathbf{x}} = (1, \mathbf{x}^\top)^\top$.

Por Mínimos Cuadrado (ii)

- Sea un conjunto de entrenamiento $\{\mathbf{x}_n, \mathbf{t}_n\}_{n=1}^N$.
- La matrix \mathbf{T} tiene vectores fila \mathbf{t}_n^\top , y la matriz $\tilde{\mathbf{X}}$ vectores fila $\tilde{\mathbf{x}}_n$.
- La función de error cuadrático está dada como

$$E_D(\tilde{\mathbf{W}}) = \frac{1}{2} \text{tr} \left\{ (\tilde{\mathbf{X}}\tilde{\mathbf{W}} - \mathbf{T})^\top (\tilde{\mathbf{X}}\tilde{\mathbf{W}} - \mathbf{T}) \right\}.$$

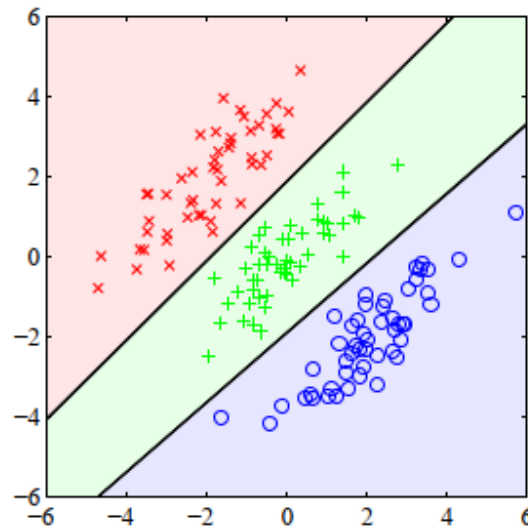
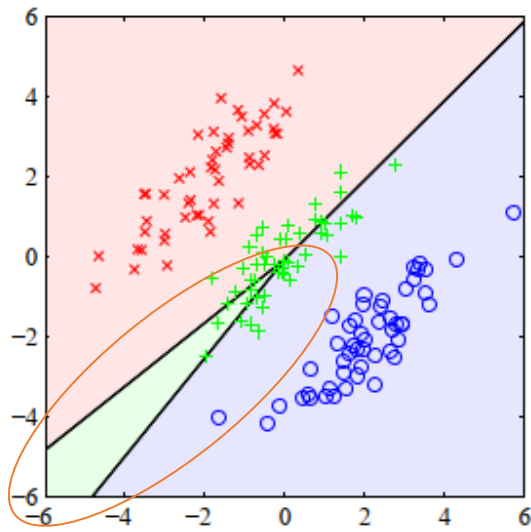
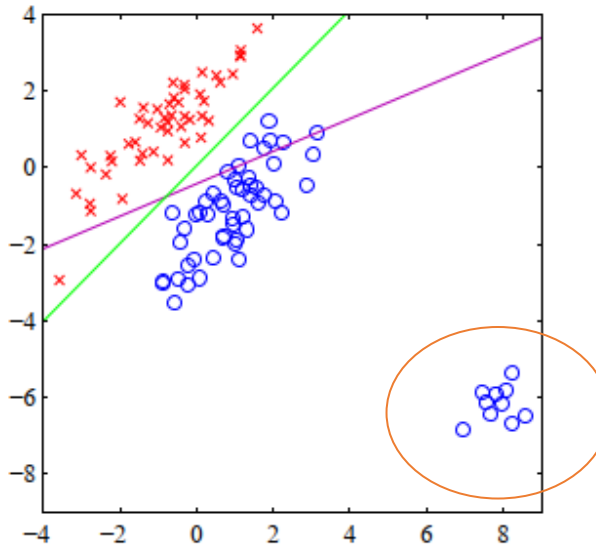
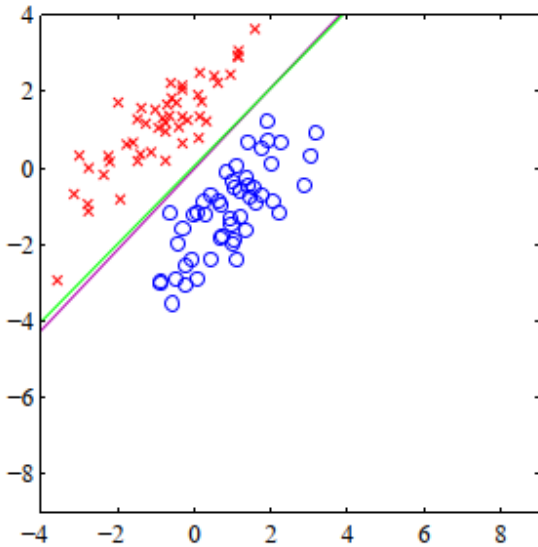
- Minimizando e igualando a cero se obtiene

$$\tilde{\mathbf{W}}_{MSE} = \left(\tilde{\mathbf{X}}^\top \tilde{\mathbf{X}} \right)^{-1} \tilde{\mathbf{X}}^\top \mathbf{T} = \tilde{\mathbf{X}}^\dagger \mathbf{T},$$

donde $\tilde{\mathbf{X}}^\dagger$ es la pseudo inversa de $\tilde{\mathbf{X}}$.

- La función discriminante está dada por $\mathbf{y}(\mathbf{x}) = \tilde{\mathbf{W}}_{MSE}^\top \tilde{\mathbf{x}} = \mathbf{T}^\top \left(\tilde{\mathbf{X}}^\dagger \right)^\top \tilde{\mathbf{x}}$.

Por Mínimos Cuadrado (iii)



- A. Sufre en la presencia de valores atípicos
- B. Se diseño para variables aleatoria gaussiana y las etiquetas no necesariamente lo son.
- C. La pérdida asociada a una clase puede ser asimétrica. Por ejemplo, los errores en una clase minoritaria pueden ser más costosos que los errores en una clase mayoritaria. Los mínimos cuadrados no capturan directamente esta asimetría en la pérdida.

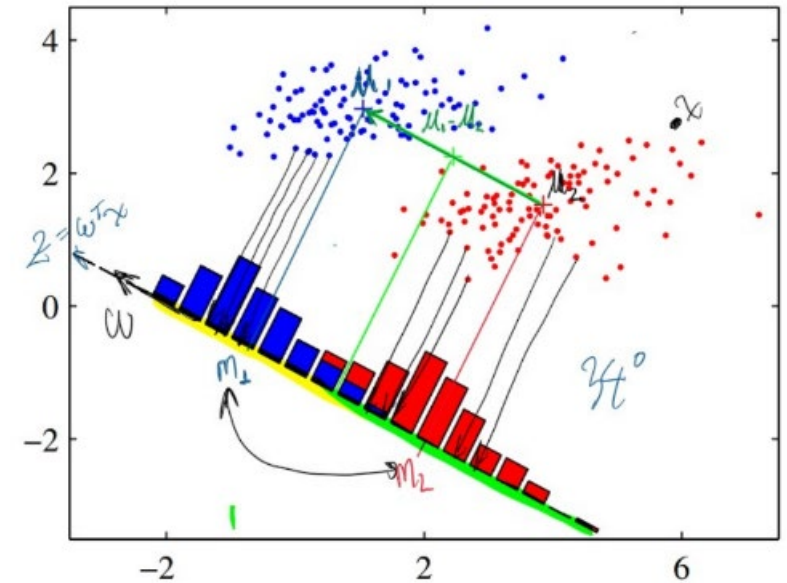
Por discriminante de Fisher (i)

- La idea es proyectar los datos a un espacio de menor dimensionalidad donde la clasificación sea más sencilla.

- Sea $\mathbf{x} \in \mathbb{R}^D$.

- Se proyecta a una dimensión usando

$$y = \mathbf{w}^T \mathbf{x}.$$



- Se establece un umbral y_0 , y se clasifica un nuevo punto como de la clase \mathcal{C}_1 si $y \geq y_0$, o de la clase \mathcal{C}_2 , si pasa lo contrario.
- La idea es escoger \mathbf{w} de manera que maximice la separabilidad de las clases.

Por discriminante de Fisher (ii)

- Sea un problema biclase, con vectores de media

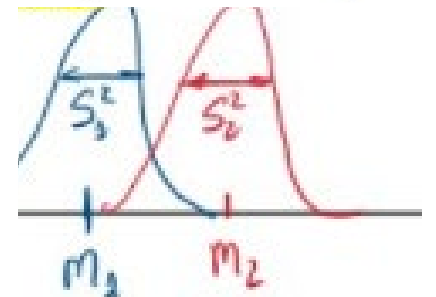
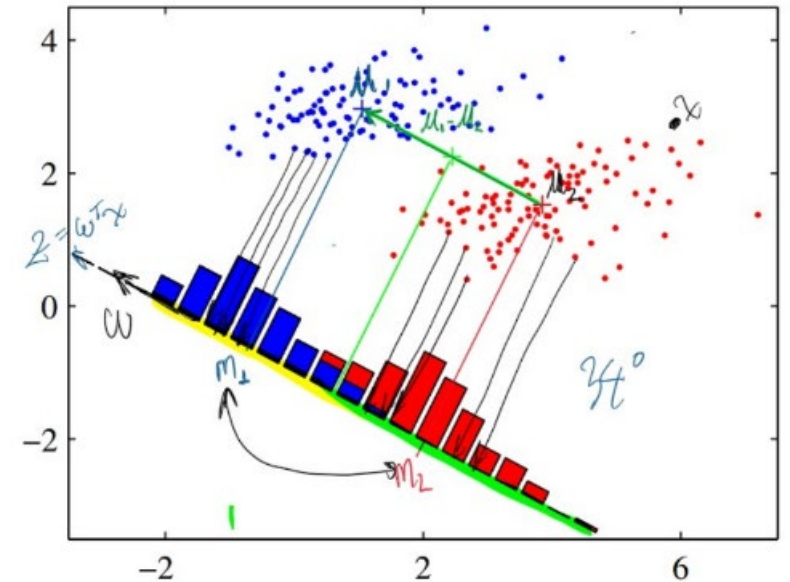
$$\mathbf{m}_1 = \frac{1}{N_1} \sum_{n \in \mathcal{C}_1} \mathbf{x}_n, \quad \mathbf{m}_2 = \frac{1}{N_2} \sum_{n \in \mathcal{C}_2} \mathbf{x}_n.$$

- Una medida de la separación entre las clases es

$$m_2 - m_1 = \mathbf{w}^\top (\mathbf{m}_2 - \mathbf{m}_1).$$

- En la expresión anterior \mathbf{w} se puede hacer muy grande, pero se limita para que tenga longitud unitaria ($\mathbf{w}^\top \mathbf{w} = 1$).

- Usando multiplicadores de Lagrange, se puede demostrar que $\mathbf{w} \propto (\mathbf{m}_2 - \mathbf{m}_1)$.



Por discriminante de Fisher (ii-A)

Ecuación para la media_k

$$m_k = \frac{1}{N_k} \sum_{n \in \mathcal{C}_k} \mathbf{w}^\top \mathbf{x}_n = \mathbf{w}^\top \left(\frac{1}{N_k} \sum_{n \in \mathcal{C}_k} \mathbf{x}_n \right) = \mathbf{w}^\top \boldsymbol{\mu}_k$$

Distancia de las medias

$$\begin{aligned} (m_2 - m_1)^2 &= (m_1 - m_2)(m_1 - m_2) \\ (m_2 - m_1)^2 &= (\mathbf{w}^\top \boldsymbol{\mu}_1 - \mathbf{w}^\top \boldsymbol{\mu}_2)(\boldsymbol{\mu}_1^\top \mathbf{w} - \boldsymbol{\mu}_2^\top \mathbf{w}) \\ (m_2 - m_1)^2 &= \mathbf{w}^\top (\boldsymbol{\mu}_1 - \boldsymbol{\mu}_2)(\boldsymbol{\mu}_1 - \boldsymbol{\mu}_2) \mathbf{w} \\ (m_2 - m_1)^2 &= \mathbf{w}^\top \mathbf{S}_B \mathbf{w} \end{aligned}$$

Con $\mathbf{S}_B = (\boldsymbol{\mu}_1 - \boldsymbol{\mu}_2)(\boldsymbol{\mu}_1 - \boldsymbol{\mu}_2)^\top$ (Covarianza interclases)

Por discriminante de Fisher (ii-B)

La dispersión dentro de cada clase se estima como la varianza intraclases de los vectores transformados de la clase \mathcal{C}_k como:

$$s_k^2 = \sum_{n \in \mathcal{C}_k} \underbrace{(y(\mathbf{x}_n) - m_k)^2}_{\omega^\top \mathcal{X}_n} = \sum_{n \in \mathcal{C}_k} (\mathbf{w}^\top \mathbf{x}_n - m_k)^2$$

$$s_k^2 = \sum_{n \in \mathcal{C}_k} (\mathbf{w}^\top \mathbf{x}_n - m_k) (\mathbf{x}_n^\top \mathbf{w} - m_k)$$

$$s_k^2 = \sum_{n \in \mathcal{C}_k} (\mathbf{w}^\top \mathbf{x}_n - \mathbf{w}^\top \boldsymbol{\mu}_k) (\mathbf{x}_n^\top \mathbf{w} - \boldsymbol{\mu}_k^\top \mathbf{w})$$

$$s_k^2 = \sum_{n \in \mathcal{C}_k} \mathbf{w}^\top (\mathbf{x}_n - \boldsymbol{\mu}_k) (\mathbf{x}_n - \boldsymbol{\mu}_k)^\top \mathbf{w}$$

$$s_k^2 = \mathbf{w}^\top \left(\sum_{n \in \mathcal{C}_k} (\mathbf{x}_n - \boldsymbol{\mu}_k) (\mathbf{x}_n - \boldsymbol{\mu}_k)^\top \right) \mathbf{w}$$

$$s_k^2 = \mathbf{w}^\top \mathbf{S}_k \mathbf{w}$$

$\mathbf{S}_k \in \mathbb{R}^{D \times D}$ es la covarianza de la clase k y la varianza total dentro de las clases, será:

$$s_1^2 + s_2^2 = \mathbf{w}^\top \mathbf{S}_1 \mathbf{w} + \mathbf{w}^\top \mathbf{S}_2 \mathbf{w}$$

$$s_1^2 + s_2^2 = \mathbf{w}^\top (\mathbf{S}_1 + \mathbf{S}_2) \mathbf{w}$$

$$s_1^2 + s_2^2 = \mathbf{w}^\top \mathbf{S}_w \mathbf{w}$$

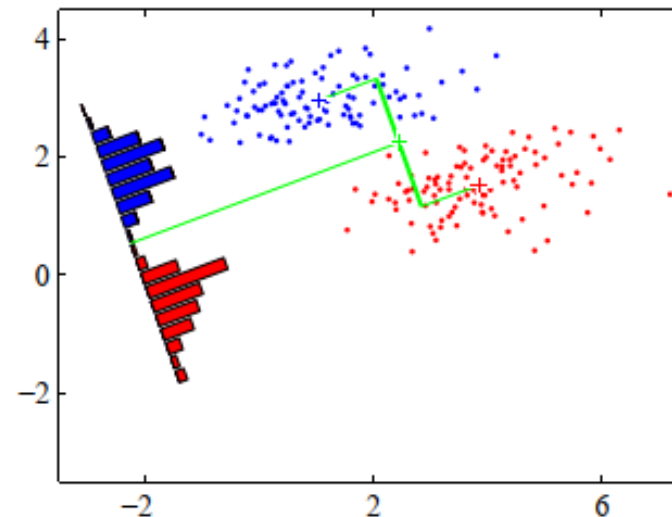
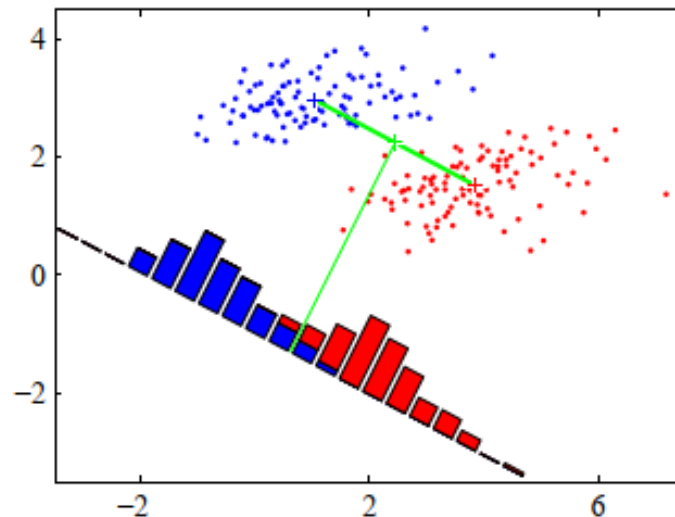
$\mathbf{S}_w \in \mathbb{R}^{D \times D}$ es la covarianza intraclase.

Por discriminante de Fisher (iii)

- ❑ No sólo se quiere maximizar la distancia entre las medias, si no también minimizar la variabilidad de las muestras en cada clase.
- ❑ La varianza intraclase se obtiene de los vectores transformados de la clase \mathcal{C}_k como

$$s_k^2 = \sum_{n \in \mathcal{C}_k} (y_n - m_k)^2,$$

donde $y_n = \mathbf{w}^\top \mathbf{x}_n$.



Por discriminante de Fisher (iv)

- El criterio de Fisher se define como el ratio de la varianza entre clases sobre la varianza intraclase

$$J(\mathbf{w}) = \frac{(m_2 - m_1)^2}{s_1^2 + s_2^2}.$$

- Haciendo los cambios apropiados se tiene

$$J(\mathbf{w}) = \frac{\mathbf{w}^\top \mathbf{S}_B \mathbf{w}}{\mathbf{w}^\top \mathbf{S}_W \mathbf{w}},$$

La idea ahora es maximizar para encontrar el \mathbf{w} optimo.

donde

$$\mathbf{S}_B = (\mathbf{m}_2 - \mathbf{m}_1)(\mathbf{m}_2 - \mathbf{m}_1)^\top,$$

$$\mathbf{S}_W = \sum_{n \in \mathcal{C}_1} (\mathbf{x}_n - \mathbf{m}_1)(\mathbf{x}_n - \mathbf{m}_1)^\top + \sum_{n \in \mathcal{C}_2} (\mathbf{x}_n - \mathbf{m}_2)(\mathbf{x}_n - \mathbf{m}_2)^\top.$$

- \mathbf{S}_B es la matriz de covarianza *entre clases*, y \mathbf{S}_W es la matriz de covarianza *intra clases*.

Por discriminante de Fisher (v)

- Derivando $J(\mathbf{w})$ con respecto a \mathbf{w} e igualando a cero se tiene

$$\frac{d}{d\mathbf{w}} J(\mathbf{w}) = \frac{(\cancel{2\mathbf{S}_B\mathbf{w}})(\mathbf{w}^\top \mathbf{S}_W \mathbf{w}) - (\mathbf{w}^\top \mathbf{S}_B \mathbf{w})(\cancel{2\mathbf{S}_W\mathbf{w}})}{(\mathbf{w}^\top \mathbf{S}_W \mathbf{w})^2} = 0 \quad (\mathbf{w}^\top \mathbf{S}_B \mathbf{w}) \mathbf{S}_W \mathbf{w} = (\mathbf{w}^\top \mathbf{S}_W \mathbf{w}) \mathbf{S}_B \mathbf{w}.$$

- Lo que importa de \mathbf{w} es su dirección, no su magnitud.
- Por lo tanto los escalares $\mathbf{w}^\top \mathbf{S}_B \mathbf{w}$ y $\mathbf{w}^\top \mathbf{S}_W \mathbf{w}$ se pueden omitir.
- Además $\mathbf{S}_B \mathbf{w}$ está siempre en la dirección de $\mathbf{m}_2 - \mathbf{m}_1$.
- Premultiplicando por \mathbf{S}_W^{-1} se encuentra que

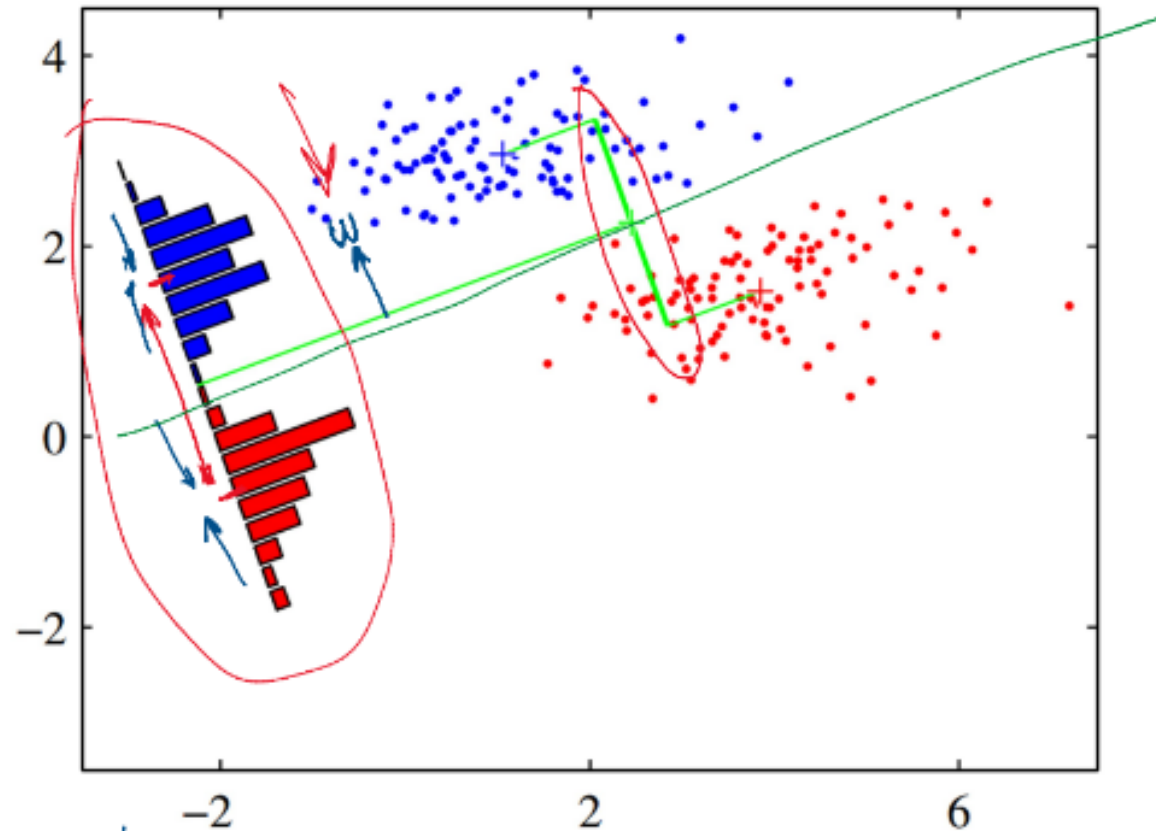
$$\mathbf{w} \propto \mathbf{S}_W^{-1} (\mathbf{m}_2 - \mathbf{m}_1).$$

- El resultado se conoce como el *discriminante lineal de Fisher*.

Por discriminante de Fisher (vi)

- El resultado se conoce como **discriminante lineal de Fisher** (Linear Discriminant Analysis - LDA).

$$\mathbf{w} \propto \mathbf{S}_W^{-1}(\mathbf{m}_2 - \mathbf{m}_1).$$



Por Algoritmo del perceptrón (i)

Es el modelo matemático de una neurona

- ❑ Dos clases. Los datos de entrada se transforman como $\phi(\mathbf{x})$.

- ❑ El modelo lineal tiene la forma

$$y(\mathbf{x}) = f(\mathbf{w}^\top \phi(\mathbf{x})),$$

donde $f(\cdot)$ es la función escalón unitario.

Se conoce como función
de activación

- ❑ En el perceptrón se asume $t = +1$, para C_1 , y $t = -1$ para C_2 .

Por Algoritmo del perceptrón (ii)

- La función a minimizar se conoce como el criterio del perceptrón.

$$E_P(\mathbf{w}) = - \sum_{n \in \mathcal{M}} \mathbf{w}^\top \phi(\mathbf{x}_n) t_n,$$

donde \mathcal{M} denota el conjunto de patrones incorrectamente clasificados.

- Aplicando el algoritmo de gradiente descendente estocástico a esta función, se tiene

$$\mathbf{w}^{(\tau+1)} = \mathbf{w}^{(\tau)} - \eta \nabla E_P(\mathbf{w}) = \mathbf{w}^{(\tau)} + \eta \phi(\mathbf{x}_n) t_n,$$

donde η se conoce como la razón de aprendizaje, y τ indexa los pasos del algoritmo.

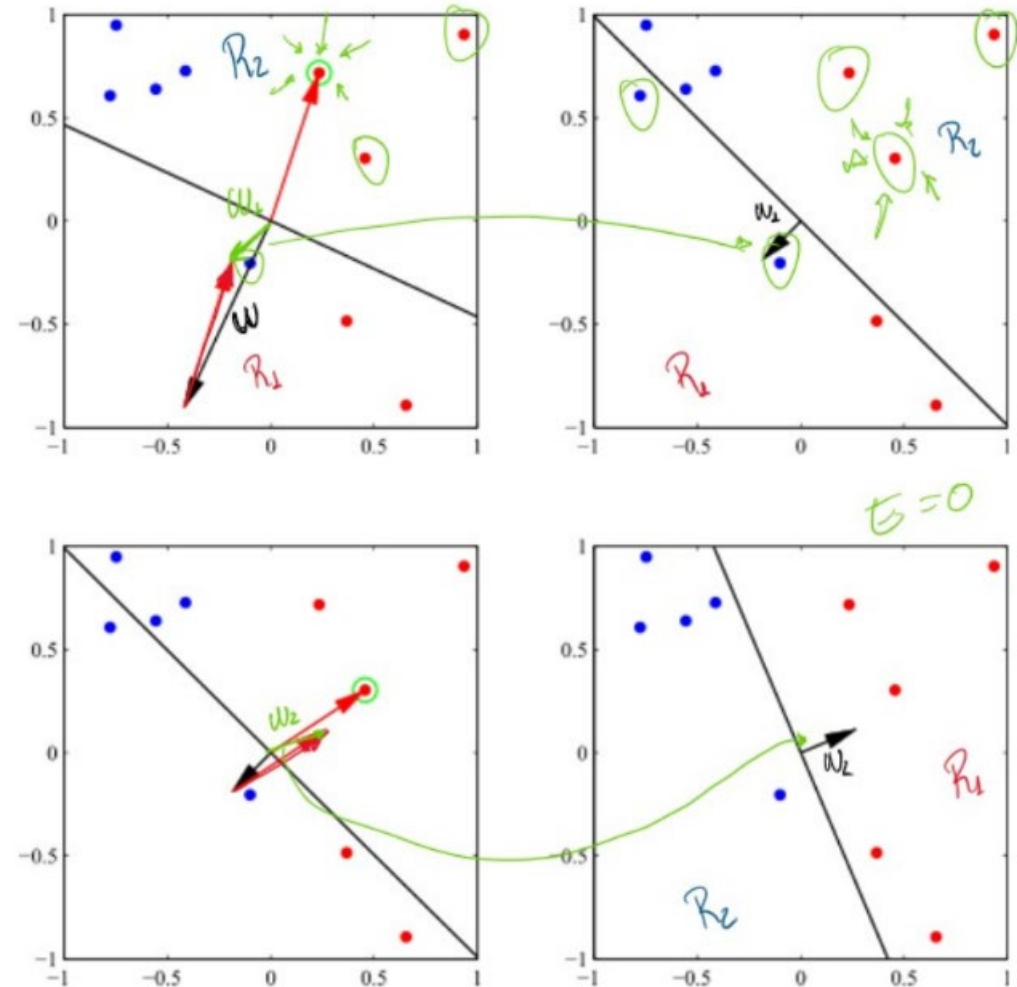
Por Algoritmo del perceptrón (iii)

Pasos:

1. Inicializar el modelo con w_0 y tasa de aprendizaje η ,
2. Realizar las predicciones con los parámetros actuales $y = f(\phi w_k)$
3. Determinar el conjunto de muestras mal etiquetadas: $\mathcal{M} = \{x_n: t_n \neq y_n\}$
4. Calcular la función de costo total para los parámetros actuales $E_p(w_k)$
5. Actualizar los parámetros para la siguiente iteración usando algunas de las muestras mal etiquetadas

$$w_{k+1} = w_k + \eta t_{\mathcal{M}}^T \phi_{\mathcal{M}}$$

1. Hasta la convergencia volver al paso 2.



Modelos generativos probabilísticos

Introducción

- En los modelos generativos se modelan las densidades de clase condicional $p(\mathbf{x}|\mathcal{C}_k)$, y las funciones de probabilidad a priori, $p(\mathcal{C}_k)$.
verosimilitud prior
- Ambas probabilidades se usan para calcular el posterior $p(\mathcal{C}_k|\mathbf{x})$.

Modelo generativo (i)

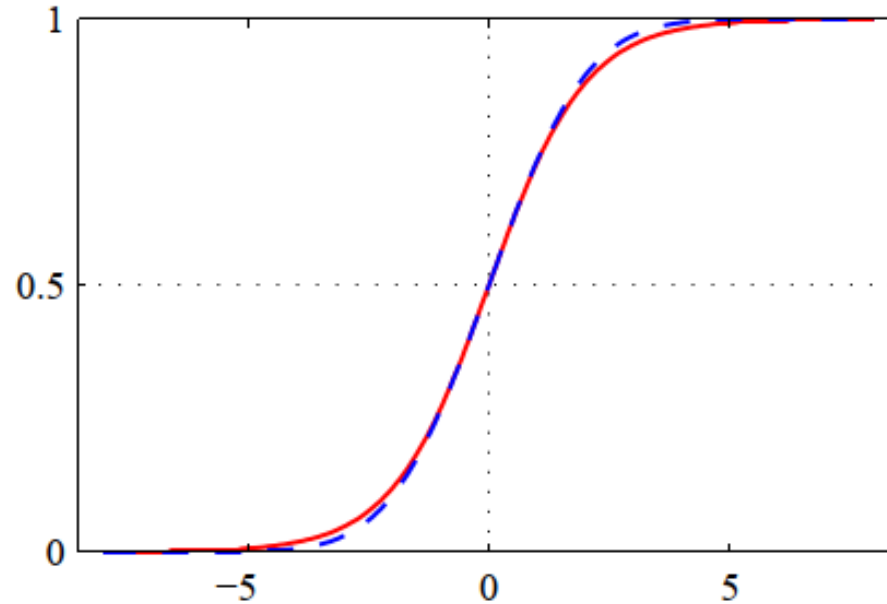
Para el caso biclase,

$$\begin{aligned} p(C_1|\mathbf{x}) &= \frac{p(\mathbf{x}|C_1)p(C_1)}{p(\mathbf{x}|C_1)p(C_1) + p(\mathbf{x}|C_2)p(C_2)} = \frac{1}{1 + \frac{p(\mathbf{x}|C_2)p(C_2)}{p(\mathbf{x}|C_1)p(C_1)}} \\ &= \frac{1}{1 + \exp \left\{ \ln \left[\frac{p(\mathbf{x}|C_2)p(C_2)}{p(\mathbf{x}|C_1)p(C_1)} \right] \right\}} = \frac{1}{1 + \exp \left\{ -\ln \left[\frac{p(\mathbf{x}|C_1)p(C_1)}{p(\mathbf{x}|C_2)p(C_2)} \right] \right\}} \\ &= \frac{1}{1 + \exp(-a)} = \sigma(a), \end{aligned}$$

donde $a = \ln \left[\frac{p(\mathbf{x}|C_1)p(C_1)}{p(\mathbf{x}|C_2)p(C_2)} \right]$, y $\sigma(a) = 1/(1 + \exp(-a))$ es la función logística sigmoïdal .

Modelo generativo (ii)

donde $\sigma(a)$ se denomina función sigmoide:

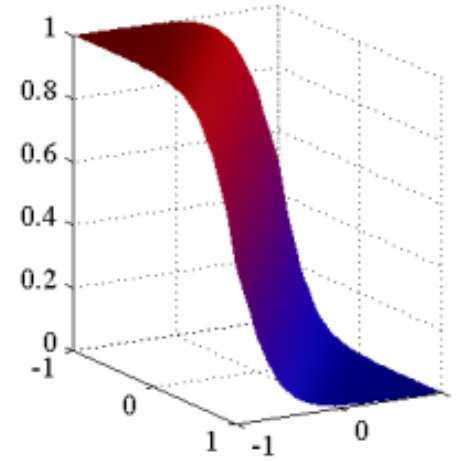
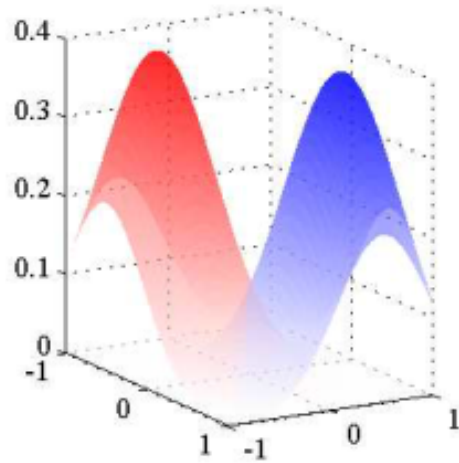


Para el caso $K > 2$ se tiene

$$p(c_k|\mathbf{x}) = \frac{p(\mathbf{x}|c_k)p(c_k)}{\sum_j p(\mathbf{x}|c_j)p(c_j)} = \frac{\exp(a_k)}{\sum_j \exp(a_j)},$$

con $a_k = \ln p(\mathbf{x}|c_k)p(c_k)$. (función soft-max).

Modelo generativo: entradas continuas (ii)



Para el caso $K > 2$ se tiene

$$a_k(\mathbf{x}) = \mathbf{w}_k^\top \mathbf{x} + w_{k0},$$

donde se ha definido $\mathbf{w}_k = \Sigma^{-1} \boldsymbol{\mu}_k$, $w_{k0} = -\frac{1}{2} \boldsymbol{\mu}_k^\top \Sigma^{-1} \boldsymbol{\mu}_k + \ln p(C_k)$.

Modelo generativo: entradas continuas

(ii) demostración (A)

CASO BINARIO

- La frontera la componen todos x los tales que

$$p(c_1|x) = p(c_2|x)$$

$$N(x) = \left(\frac{1}{2\pi}\right)^{D/2} \frac{1}{|\Sigma|^{1/2}} e^{-\frac{1}{2}Q(x)}$$

$$Q(x) = (x - \mu)^T \Sigma^{-1} (x - \mu)$$

$$p(x|c_1)p(c_1) = p(x|c_2)p(c_2) \rightarrow a=0 = \ln \left\{ \frac{p(x|c_1)p(c_1)}{p(x|c_2)p(c_2)} \right\}$$

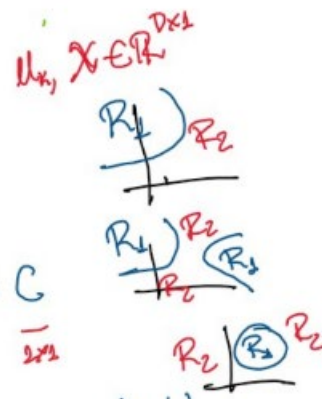
$$a = \ln N_1(x) + \ln \pi_1 - \ln N_2(x) - \ln(1 - \pi_1)$$

$$= -\frac{D}{2} \ln 2\pi - \frac{1}{2} \ln |\Sigma_1| - \frac{1}{2} Q_1(x) + \ln \pi_1 + \frac{D}{2} \ln 2\pi + \frac{1}{2} \ln |\Sigma_2| + \frac{1}{2} Q_2(x) - \ln(1 - \pi_1)$$

$$a = -\frac{1}{2} (x - \mu_1)^T \Sigma_1^{-1} (x - \mu_1) + \frac{1}{2} (x - \mu_2)^T \Sigma_2^{-1} (x - \mu_2) + C$$

$$= \left[-\frac{1}{2} x^T \Sigma_1^{-1} x + \mu_1^T \Sigma_1^{-1} x - \frac{1}{2} \mu_1^T \Sigma_1^{-1} \mu_1 \right] + \left[\frac{1}{2} x^T \Sigma_2^{-1} x - \mu_2^T \Sigma_2^{-1} x + \frac{1}{2} \mu_2^T \Sigma_2^{-1} \mu_2 \right] + C$$

$$a = \underbrace{-\frac{1}{2} x^T (\Sigma_1^{-1} - \Sigma_2^{-1}) x}_{1 \times D \quad D \times D \quad D \times 1} + \underbrace{(\mu_1^T \Sigma_1^{-1} - \mu_2^T \Sigma_2^{-1}) x}_{1 \times D \quad D \times D \quad 1 \times D \quad D \times 1} + C$$



$$a = -\frac{1}{2} x^T A x + b^T x + c = 0 : \text{Frontera es cuadrática}$$

$$A = 0 \quad a = 0 = b^T x + c : \text{Frontera lineal}$$

$$A = \Sigma_1^{-1} - \Sigma_2^{-1} = 0 \rightarrow \Sigma_1 = \Sigma_2$$

$$b^T = \mu_1^T \Sigma_1^{-1} - \mu_2^T \Sigma_2^{-1} = (\mu_1^T - \mu_2^T) \Sigma_1^{-1} \rightarrow b = \Sigma_1^{-1} (\mu_1 - \mu_2)$$

- Por lo tanto, la probabilidad a posteriori de clase cuando se asumen covarianzas iguales se reduce a:

$$p(c_1|x) = \sigma(a)$$

$$a = \ln \left(\frac{p(x|c_1)p(c_1)}{p(x|c_2)p(c_2)} \right)$$

$$a = \ln(p(x|c_1)p(c_1)) - \ln(p(x|c_2)p(c_2))$$

$$a = w^T x + w_0$$

donde

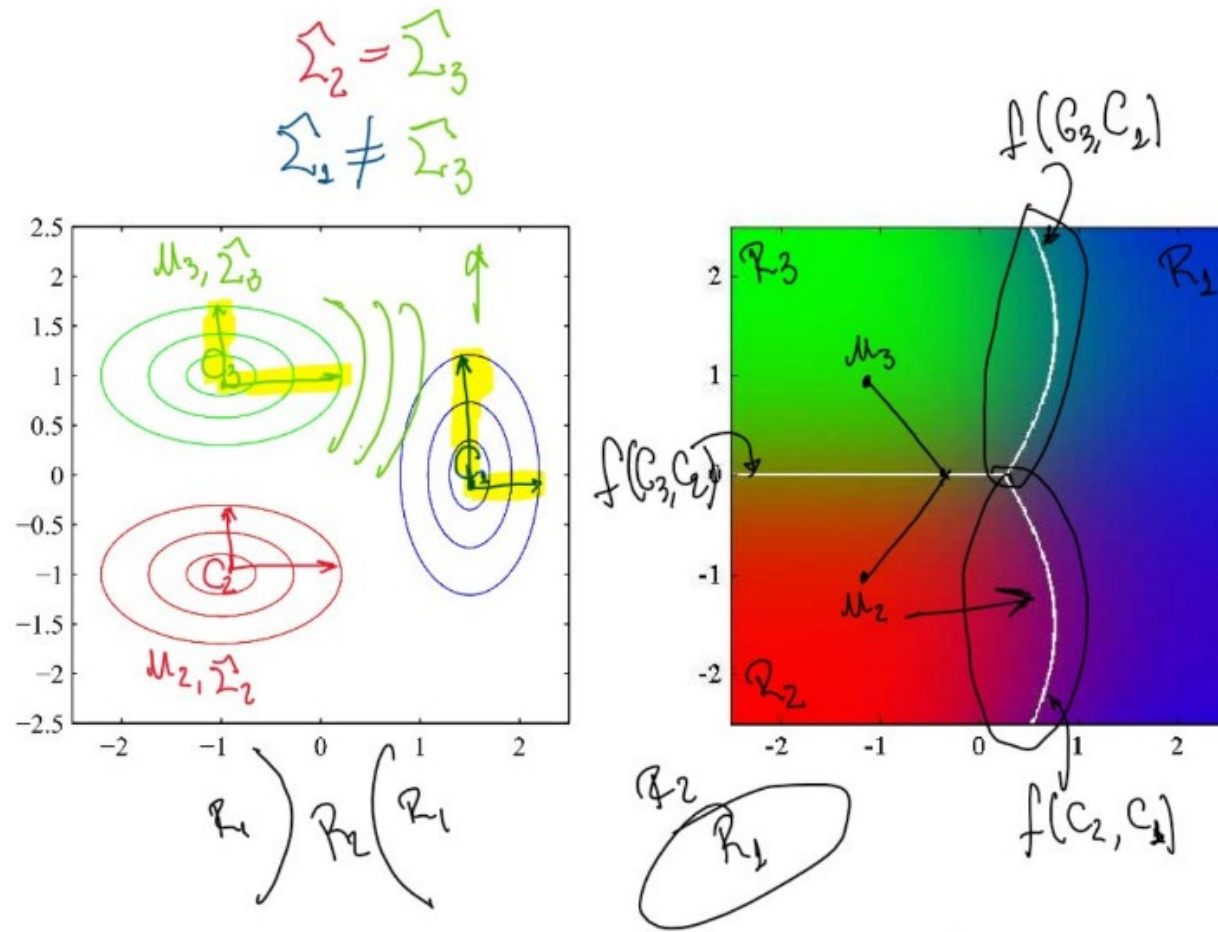
$$w = \Sigma^{-1} (\mu_1 - \mu_2) : \text{LDA}$$

$$w_0 = -\frac{1}{2} \mu_1^T \Sigma^{-1} \mu_1 + \frac{1}{2} \mu_2^T \Sigma^{-1} \mu_2 + \ln \left(\frac{\pi}{1 - \pi} \right)$$

$$\pi = p(c_1)$$

Modelo generativo: entradas continuas (iii)

Si las clases no comparten la misma matriz de covarianza, las fronteras de decisión son cuadráticas



Máxima verosimilitud (i)

- ❑ Los valores de μ_1 , μ_2 , Σ , $p(\mathcal{C}_1)$ y $p(\mathcal{C}_2)$ se determinan usando máxima verosimilitud.
- ❑ Sean dos clases y un conjunto de datos $\{\mathbf{x}_n, t_n\}_{n=1}^N$.
- ❑ $t_n = 1$ denota \mathcal{C}_1 , y $t_n = 0$ denota \mathcal{C}_2 .
- ❑ Las probabilidades a priori se escriben como $p(\mathcal{C}_1) = \pi$ y $p(\mathcal{C}_2) = 1 - \pi$.

Máxima verosimilitud (ii)

- Para un punto \mathbf{x}_n de la clase \mathcal{C}_1 se tiene $t_n = 1$ y así

$$p(\mathbf{x}_n, \mathcal{C}_1) = p(\mathbf{x}_n | \mathcal{C}_1) p(\mathcal{C}_1) = \pi \mathcal{N}(\mathbf{x}_n | \boldsymbol{\mu}_1, \boldsymbol{\Sigma}).$$

- De forma similar, para la clase \mathcal{C}_2

$$p(\mathbf{x}_n, \mathcal{C}_2) = p(\mathbf{x}_n | \mathcal{C}_2) p(\mathcal{C}_2) = (1 - \pi) \mathcal{N}(\mathbf{x}_n | \boldsymbol{\mu}_2, \boldsymbol{\Sigma}).$$

- Luego,

$$p(t_n, \mathbf{x}_n | \pi, \boldsymbol{\mu}_1, \boldsymbol{\mu}_2, \boldsymbol{\Sigma}) = \begin{cases} p(\mathbf{x}_n, \mathcal{C}_1) = \pi \mathcal{N}(\mathbf{x}_n | \boldsymbol{\mu}_1, \boldsymbol{\Sigma}), & \text{si } t_n = 1 \\ p(\mathbf{x}_n, \mathcal{C}_2) = (1 - \pi) \mathcal{N}(\mathbf{x}_n | \boldsymbol{\mu}_2, \boldsymbol{\Sigma}), & \text{si } t_n = 0. \end{cases}$$

- Lo anterior se puede escribir de forma resumida como

$$p(t_n, \mathbf{x}_n | \pi, \boldsymbol{\mu}_1, \boldsymbol{\mu}_2, \boldsymbol{\Sigma}) = [\pi \mathcal{N}(\mathbf{x}_n | \boldsymbol{\mu}_1, \boldsymbol{\Sigma})]^{t_n} [(1 - \pi) \mathcal{N}(\mathbf{x}_n | \boldsymbol{\mu}_2, \boldsymbol{\Sigma})]^{1-t_n}$$

Máxima verosimilitud (iii)

- Asumiendo que los datos son iid,

$$p(\mathbf{t}, \mathbf{X} | \pi, \mu_1, \mu_2, \Sigma) = \prod_{n=1}^N [\pi \mathcal{N}(\mathbf{x}_n | \mu_1, \Sigma)]^{t_n} [(1 - \pi) \mathcal{N}(\mathbf{x}_n | \mu_2, \Sigma)]^{1-t_n}.$$

- Se realiza la maximización de $\ln p(\mathbf{t}, \mathbf{X} | \pi, \mu_1, \mu_2, \Sigma)$, el logaritmo de la verosimilitud, con respecto a $\pi, \mu_1, \mu_2, \Sigma$,

Máxima verosimilitud: Solución para π, μ_1 y μ_2

- Se puede demostrar que

$$\pi_{ML} = \frac{1}{N} \sum_{n=1}^N t_n = \frac{N_1}{N} = \frac{N_1}{N_1 + N_2},$$

donde N_1 denota el número de puntos de la clase \mathcal{C}_1 , y N_2 el número de puntos de la clase \mathcal{C}_2 .

- Igualmente, se puede demostrar que

$$\mu_1 = \frac{1}{N_1} \sum_{n=1}^N t_n \mathbf{x}_n, \quad \mu_2 = \frac{1}{N_2} \sum_{n=1}^N (1 - t_n) \mathbf{x}_n.$$

Máxima verosimilitud: Solución para Σ

Finalmente, para Σ

$$\Sigma = \mathbf{S} = \frac{N_1}{N} \mathbf{S}_1 + \frac{N_2}{N} \mathbf{S}_2$$

$$\mathbf{S}_1 = \frac{1}{N_1} \sum_{n \in \mathcal{C}_1} (\mathbf{x}_n - \mu_1)(\mathbf{x}_n - \mu_1)^\top,$$

$$\mathbf{S}_2 = \frac{1}{N_2} \sum_{n \in \mathcal{C}_2} (\mathbf{x}_n - \mu_2)(\mathbf{x}_n - \mu_2)^\top.$$

Demostración solución de π, μ_1 y μ_2

$$\max_{\mu_1, \mu_2, \Sigma, \pi} \ln p(t, \mathbf{X} | \mu_1, \mu_2, \Sigma_1, \Sigma_2, \pi)$$

$\prod_n \rightarrow \sum_n$

$$L = \sum_{n=1}^N t_n \left[-\frac{D}{2} \ln 2\pi - \frac{1}{2} \ln |\Sigma_1| - \frac{1}{2} (x_n - \mu_1)^T \Sigma_1^{-1} (x_n - \mu_1) + \ln \pi_1 \right] \\ + \sum_{n=1}^N (1 - t_n) \left[-\frac{D}{2} \ln 2\pi - \frac{1}{2} \ln |\Sigma_2| - \frac{1}{2} (x - \mu_2)^T \Sigma_2^{-1} (x - \mu_2) + \ln(1 - \pi_1) \right]$$

$$L = \sum_{n=1}^N t_n \left[-\frac{D}{2} \ln 2\pi - \frac{1}{2} \ln |\Sigma_1| - \frac{1}{2} (x_n - \mu_1)^T \Sigma_1^{-1} (x_n - \mu_1) + \ln \pi_1 \right] \\ + \sum_{n=1}^N (1 - t_n) \left[-\frac{D}{2} \ln 2\pi - \frac{1}{2} \ln |\Sigma_2| - \frac{1}{2} (x - \mu_2)^T \Sigma_2^{-1} (x - \mu_2) + \ln(1 - \pi_1) \right]$$

$$\frac{\partial L}{\partial \pi_1} = \sum_n \frac{t_n}{\pi_1} + \sum_n (1 - t_n) \cdot \frac{-1}{1 - \pi_1} = 0$$

$$\frac{1}{\pi_1} \sum_n t_n - \frac{1}{1 - \pi_1} \sum_n (1 - t_n) = 0$$

$N_1 \quad N_2$

$$\frac{N_1}{\pi_1} = \frac{N_2}{1 - \pi_1} \rightarrow (1 - \pi_1) N_1 = \pi_1 N_2$$

$$\hookrightarrow N_1 = \pi_1 N_1 + \pi_1 N_2 \rightarrow \pi_1 N = \frac{N_1}{N_1 + N_2} = \frac{N_1}{N}$$

$N=6$

n	t_n	$1 - t_n$
1	0 x	1 ✓
2	1 ✓	0 ✓
3	0 x	1 ✓
4	0 x	1 ✓
5	1 ✓	0 ✓
6	0 x	1 ✓
	<u>$N_1 = 2$</u>	<u>$4 = N_2$</u>

Demostración solución de π, μ_1 y μ_2

$$\frac{dL}{d\mu_1} = \sum_n \ln\left(-\frac{1}{Z}\right) \frac{d}{d\mu_1} \left\{ \cancel{x_n^T \Sigma_1^{-1} x_n} - 2 \cancel{x_n^T \Sigma_1^{-1} \mu_1} + \underbrace{\mu_1^T \Sigma_1^{-1} \mu_1}_{\propto x^2} \right\}$$

$\frac{d}{d\mu_1} = \sum_n \ln\left(-\frac{1}{Z}\right) \left(\cancel{b^T \Sigma^{-1} a} + \cancel{b^T \Sigma^{-1} \cdot 1} \right) = 0$

$$\begin{aligned} \frac{dL}{d\mu_1} &= \sum_n \ln\left(-\frac{1}{Z}\right) \left(-2 \Sigma_1^{-1} x_n + 2 \Sigma_1^{-1} \mu_1 \right) = 0 \\ &= -\sum_n \Sigma_1^{-1} x_n \ln + \sum_n \Sigma_1^{-1} \mu_1 \ln = 0 \\ &= \cancel{\sum_n} \sum_n x_n \ln = \sum_n \sum_{n=1}^N \mu_1 \ln = \sum_n \mu_1 \sum_{n=1}^N \ln \quad N_1 \\ \sum_{n \in G_1} x_n &= N_1 \mu_1 \rightarrow \mu_1 = \frac{1}{N_1} \sum_{n \in G_1} x_n \end{aligned}$$

$$\mu_k = \frac{1}{N_k} \sum_{n \in G_k} x_n$$

$$\begin{aligned} \frac{dL}{d\Sigma_1} &= \sum_n \ln \left[-\frac{1}{Z} \frac{d}{d\Sigma_1} \left\{ \ln Z \right\} - \frac{1}{Z} \frac{d}{d\Sigma_1} \left\{ \frac{(x_n - \mu_1)^T \Sigma_1^{-1} (x_n - \mu_1)}{a} \right\} \right] \\ &\quad \text{axb} \rightarrow ab; \quad x^1 \rightarrow -x^2 \\ \frac{d \ln|x|}{dX} &= X^{-1}; \quad \frac{d}{dX} \left\{ \frac{a^T X^{-1} b}{1 \times 0 \quad 0 \times 0 \quad 0 \times 1} \right\} = \frac{-X^{-1} a b^T X^{-1}}{0 \times 0 \quad 0 \times 0 \quad 0 \times 0} \\ \frac{dL}{d\Sigma_1} &= \frac{1}{2} \sum_n \ln \left[\Sigma_1^{-1} - \Sigma_1^{-1} (x_n - \mu_1) (x_n - \mu_1)^T \Sigma_1^{-1} \right] = 0 \\ \left(\sum_n \ln \right) \Sigma_1^{-1} &= \Sigma_1^{-1} \sum_n \ln (x_n - \mu_1) (x_n - \mu_1)^T \Sigma_1^{-1} \\ N_1 \Sigma_1^{-1} &= \Sigma_1^{-1} \sum_n \ln (x_n - \mu_1) (x_n - \mu_1)^T \Sigma_1^{-1} \end{aligned}$$

$$\begin{aligned} N_1 \Sigma_1^{-1} \Sigma_1 \Sigma_1^{-1} &= \Sigma_1^{-1} \sum_n \ln (x_n - \mu_1) (x_n - \mu_1)^T \Sigma_1^{-1} \\ \Sigma_1 &= \frac{1}{N_1} \sum_{n \in G_1} (x_n - \mu_1) (x_n - \mu_1)^T \end{aligned}$$

$$\Sigma_k = \frac{1}{N_k} \sum_{n \in G_k} (x_n - \mu_k) (x_n - \mu_k)^T \quad \text{Covarianza de la clase } k$$

Demostración solución de π, μ_1 y μ_2

- Por lo tanto:

$$\pi = \frac{N_1}{N}$$

$$\mu_1 = \frac{1}{N_1} \sum_{n=1}^N t_n x_n$$

$$\mu_2 = \frac{1}{N_2} \sum_{n=1}^N (1 - t_n) x_n$$

$$\Sigma_1 = \frac{1}{N_1} \sum_{n=1}^N t_n (x_n - \mu_1) (x_n - \mu_1)^\top$$

$$\Sigma_2 = \frac{1}{N_2} \sum_{n=1}^N (1 - t_n) (x_n - \mu_2) (x_n - \mu_2)^\top$$