

TEMA 1

CONCEPTOS DE PROBABILIDAD

DOMINIO: conjunto de todos los posibles eventos de un experimento aleatorio.

VARIABLE ALEATORIA: variable que define el experimento

SUCESO/EVENTO: valor que toma la variable al realizar el experimento.

PRINCIPIO INCLUSIÓN-EXCLUSIÓN

$$P(A \cup B) = P(A) + P(B) - P(A \cap B)$$

REGLA DE LA SUMA

$$P(A) + P(\bar{A}) = 1$$

REGLA DEL PRODUCTO

$$P(A, B) = P(A|B) \cdot P(B)$$

$$P(A|B) = \frac{P(A, B)}{P(B)}$$

REGLA DE BAYES

$$P(A|B) = \frac{P(B|A) \cdot P(A)}{P(B)}$$

porque $P(A, B) = P(A|B) \cdot P(B) = P(B|A) \cdot P(A)$

$$P(H|D, I) = \frac{P(D|H, I) \cdot P(H|I)}{P(D|I)}$$

$P(H|I)$: probabilidad a priori antes de ver los datos

$P(D|H, I)$: verosimilitud. Suponiendo H cierta como de verosímil es D .

$P(H|D, I)$: probabilidad a posteriori. La probabilidad de H "a la luz" de D

$P(D|I)$: evidencia. Probabilidad de los datos.

TOMA DE DECISIONES

Queremos escoger H^* entre H_1, H_2, \dots

① Maximizar la prob. a posteriori (MAP)

$$\begin{aligned} H^* &= \underset{H_i}{\operatorname{argmax}} P(H_i|D) = \underset{H_i}{\operatorname{argmax}} \frac{P(D|H_i) \cdot P(H_i)}{P(D)} = \\ &= \underset{H_i}{\operatorname{argmax}} P(D|H_i) P(H_i) \end{aligned}$$

② Máxima verosimilitud (MV)

$$H^* = \underset{H_i}{\operatorname{argmax}} P(D|H_i)$$

obs: La ② es más fácil de calcular pero la ① da los mejores resultados.

REGLA DE BAYES SIMÉTRICA

El mejor clasificador posible es el que clasifica cada ejemplo con la clase con mayor prob. a posteriori (MAP)
 → reduce prob. error

REGLA ASIMÉTRICA DE BAYES

- Suponemos un coste R si pasamos una t por j . ← truca
- Suponemos un coste R' si pasamos una j por t . ← just

Si observamos $D = c$, ← cara ¿qué coste tienen las decisiones...?

$$\begin{cases} \textcircled{1} h(c) = j \Rightarrow P(t|c) \cdot R \\ \textcircled{2} h(c) = t \Rightarrow P(j|c) \cdot R' \end{cases}$$

Si $D = c$:

$$h(c) = \begin{cases} j & \text{si } P(t|c)R < P(j|c) \cdot R' \\ t & \text{si } P(t|c)R \geq P(j|c) \cdot R' \end{cases}$$

$$P(t|c) + P(j|c) = 1$$

$$P(j|c) > \frac{R}{R'} (1 - P(j|c)) \Rightarrow P(j|c) > \frac{R}{R+R'}$$

En conclusión:

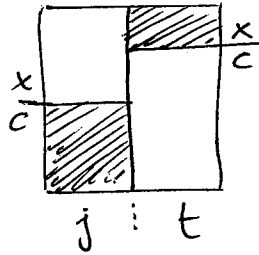
Si tenemos dos hipótesis H_1 y H_2 , y al elegir H_1 por error tenemos un coste R y si al elegir H_2 por error tenemos un coste $R' \Rightarrow$ hay que elegir:

$$\begin{cases} H_1 & \text{si } P(H_1|D) > \frac{R}{R+R'} \\ H_2 & \text{si } P(H_2|D) > \frac{R'}{R+R'} \end{cases}$$

Continuamos con el ejemplo de las monedas: $R = 3R'$

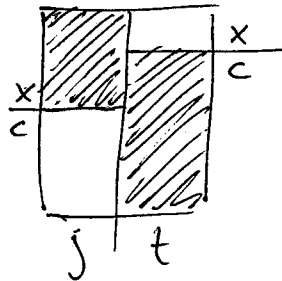
¿Modelo con menor coste? ¿ h_1, h_2, h_3 ó h_4 ?

$$\boxed{h_1} : \begin{cases} h_1(c) = t \\ h_1(x) = j \end{cases}$$



$$\text{coste}(h_1) = P(j, c) \cdot R' + P(t, x) \cdot R = \underline{0'625 \cdot R'}$$

$$\boxed{h_2} : \begin{cases} h_2(c) = j \\ h_2(x) = t \end{cases}$$



$$\begin{aligned} \text{coste}(h_2) &= P(j, c) \cdot 0 + P(j, x) \cdot R' + P(t, c) \cdot R + P(t, x) \cdot 0 = \\ &= P(t, c)R + P(j, x)R' = 0'75 \cdot 0'5 \cdot 3R' + 0'5 \cdot 0'5 R' = \\ &= \frac{9}{8} R' + \frac{1}{4} R' = \underline{1'375 R'} \end{aligned}$$

$$\boxed{h_3} : \text{coste}(h_3) = \underline{1'5 R'}$$

$$\boxed{h_4} : \text{coste}(h_4) = \underline{0'5 R'}$$

Ejercicio: problema multidimensional (lentillas)

H = no, duras, blandas

D = edad, lesión, astigmatismo, prod. lágrimas
(E) (L) (A) (P)

) Llega un paciente con $L=m, P=n$. ¿H=n, H=d, H=b?

$$P(H=n | L=m, P=n) = \frac{P(L=m, P=n | H=n) \cdot P(H=n)}{P(L=m, P=n)} = (*)$$

sacamos $P(\cdot)$

a partir de los datos

$P(H=n) = \frac{15}{24}$	$P(L=m, P=n H=n) = \frac{1}{15}$
$P(H=d) = \frac{4}{24}$	$P(L=m, P=n H=d) = \frac{3}{4}$
$P(H=b) = \frac{5}{24}$	$P(L=m, P=n H=b) = \frac{2}{5}$


$$(*) = \frac{\frac{1}{15} \cdot \frac{15}{24}}{P(L=m, P=n)} = \frac{1/24}{P(\dots)}$$

$$P(H=d | L=m, P=n) = \frac{\frac{3}{4} \cdot \frac{4}{24}}{P(L=m, P=n)} = \frac{3/24}{P(\dots)} \quad \leftarrow \text{MAYOR PROB. A POSTERIORI}$$

$$P(H=b | L=m, P=n) = \frac{\frac{2}{5} \cdot \frac{5}{24}}{P(L=m, P=n)} = \frac{2/24}{P(\dots)}$$

Normalizando:

$$P(H=n | L=m, P=n) = \frac{1}{6} \quad ; \quad P(H=d | L=m, P=n) = \frac{1}{2} \quad ; \quad P(H=b | L=m, P=n) = \frac{1}{3}$$

b) $L=h, P=n$	$\frac{1}{3}$	$\frac{1}{6}$	$\frac{1}{2}$
c) $L=m, P=r$	1	0	0
d) $L=h, P=r$	1	0	0
	$P(H=n \dots)$	$P(H=d \dots)$	$P(H=b \dots)$
MATERIAL PARCIAL 1	*		
	*		

Concepto (Maldición de la dimensionalidad): el n° de puntos necesarios para cubrir el espacio de atributos de forma uniforme crece exponencialmente con la dimensión.

Posible solución: NAÏVE-BAYES ("Bayes inocente")

Supondremos que los atributos son independientes dada la clase

Dado un vector de atributos $\vec{x} = (x_1, \dots, x_d)$, queremos saber la H^* más probable de entre las posibles. Es decir, debemos calcular $P(H_i | \vec{x}) \quad \forall i = 1, \dots, K$. $\leftarrow K \text{ hipótesis } \{H_1, \dots, H_K\}$

aplicando la regla de Bayes:

$$P(H_i | x_1, \dots, x_d) = \frac{P(x_1, \dots, x_d | H_i) \cdot P(H_i)}{P(x_1, \dots, x_d)}$$

El cálculo de la verosimilitud requiere una cantidad de datos exponencial con la dimensión.

$$\begin{aligned} \text{NAÏVE-BAYES} \quad P(H_i | x_1, \dots, x_d) &= \frac{P(x_1 | H_i) P(x_2 | H_i) \dots P(x_d | H_i) \cdot P(H_i)}{P(x_1 \dots x_d)} \\ &= \frac{\prod_{j=1}^d P(x_j | H_i) \cdot P(H_i)}{\sum_{k=1}^K \prod_{j=1}^d P(x_j | H_k) P(H_k)} \end{aligned}$$

Modelos Gráficos (Redes de Bayes)

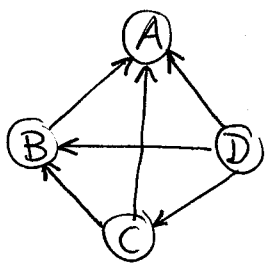
Permiten crear modelos que tengan en cuenta relaciones entre atributos

Regla de la cadena: aplicar la regla del producto múltiples veces

Podemos representar gráficamente con un grafo acíclico dirigido

Nodos: variables

Aristas: relaciones



La representación depende de como se aplique la regla de la cadena.

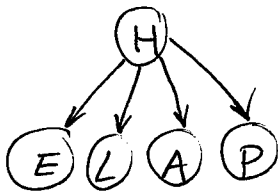
$$\text{En general, } P(x_1, \dots, x_d) = \prod_{j=1}^d P(x_j | pa_j)$$

donde pa_j = conjunto de variables padre de x_j

si x_j no tiene progenitores $P(x_j | pa_j) = P(x_j)$

Ejemplo: modelo gráfico para NB y las lentillas

$$P(E, L, A, P, H) = P(E, L, A, P | H) \cdot P(H) \stackrel{\text{NB}}{=} P(E | H) \cdot P(L | H) \cdot P(A | H) \cdot P(P | H) \cdot P(H)$$



Regla de decisión de NB:

$$H^* = \underset{H_i}{\operatorname{argmax}} \prod_{j=1}^d P(x_j | H_i) P(H_i)$$

Las $P(x_j | H_i)$ y $P(H_i)$ se estiman de los datos.

El entrenamiento consiste en crear tablas de cuentas de clases con respecto a los valores. Una tabla por atributo.

Por ejemplo, para el ejercicio de las lentillas

LESIÓN	no	d	b
miopia	7	3	2
hipermet	8	1	3

PRODUCCIÓN	no	d	b
normal	3	4	5
reducida	12	0	0

Solución para evitar ceros: CORRECCIÓN DE LAPLACE

Consiste en sumar 1 a todas las entradas de la tabla si esta contiene algún cero.

Obs: para atributos continuos las tablas pasaran a ser el cálculo de la media y varianza del atributo de cada clase

Coste computacional de Naive-Bayes

N : nº de ejemplos

D : nº de atributos

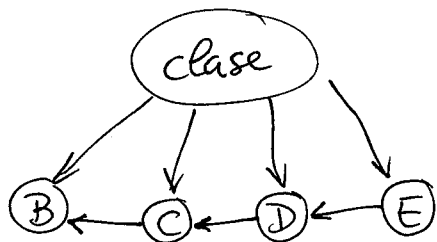
K : nº de clases

Coste de entrenamiento: $O(DN)$

Coste de clasificación: $O(DK)$

Ejemplo: Bayes gráfico

$$P(\text{clase}, B, C, D, E) = P(B|C, \text{clase}) \cdot P(C|E, \text{clase}) \cdot P(E|D) \cdot P(D|\text{clase}) \cdot P(\text{clase})$$



$$\begin{aligned} \text{Bayes} \rightarrow P(\text{clase} | x_1, \dots, x_n) &= \frac{P(x_1, \dots, x_n | \text{clase}) \cdot P(\text{clase})}{P(x_1, \dots, x_n)} = \\ &= \frac{P(\text{clase}, x_1, \dots, x_n)}{P(x_1, \dots, x_n)} \xrightarrow{\text{equiv.}} P(\text{clase}, B, C, D, E) \end{aligned}$$

TEMA 2

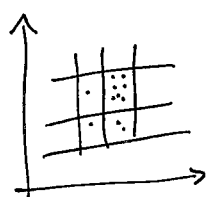
MÉTODO DE VECINOS MÁS PRÓXIMOS

Introducción & Preliminares

Para atributos no discretos

- opc. 1: discretización en intervalos/dominios
- opc. 2: ajustar distr. de prob. paramétrica

Generalización de histogramas a más dimensiones usando la PDF:



2D

$$P(\bar{x} \in \text{Región}) = \int_{\text{Región}} \text{PDF}(\bar{x}) d\bar{x}$$

Las estimaciones a la hora de discretizar son de dos formas:

① Por histogramas:

$$P(\bar{x} \in \text{Región}) \simeq \frac{K_R}{N} \quad \text{con } K_R \text{ igual al n.º de pto. de la región } R.$$

N debe ser grande para cubrir el espacio.
 \Rightarrow Maldición de dimensionalidad

② Si suponemos que la región es pequeña, podemos suponer que la PDF es constante en dicha región.

$$P(\bar{x} \in \text{Región}) = \int_{\text{Región}} \text{PDF}(\bar{x}) d\bar{x} \simeq \text{PDF}(\bar{x}_R) \int_{\text{Región}} d\bar{x} = \text{PDF}(\bar{x}_R) V_R$$

\nearrow centro región
 \nearrow volumen región

Bonus: Si ahora igualamos ① y ②:

$$\text{PDF}(\bar{x}_R) \simeq \frac{K_R}{N \cdot V_R}$$

Esta fórmula se suele usar de dos formas:

- Fijar el volumen V_R y calcular $K \Rightarrow$ MÉTODO DE NÚCLEOS
- Fijar K y V_R se define para que caigan K pto. dentro \Rightarrow
 \Rightarrow MÉTODO DE VECINOS PRÓXIMOS

Bayes: verosimilitud para la clase $K \rightarrow P(\bar{x} | C_K) \simeq \frac{K_R}{N \cdot V_R}$

Prob. a priori $\rightarrow P(C_K) = \frac{N_K}{N}$

Evidencia $\rightarrow P(\bar{x}) = \frac{K}{N \cdot V_R}$

$$P(C_K | \bar{x}) = \frac{\frac{K_R}{N \cdot V_R} \cdot \frac{N_K}{N}}{\frac{K}{N \cdot V_R}} = \frac{\frac{K_R}{N \cdot V_R}}{\frac{K}{N \cdot V_R}} = \frac{K_R}{K}$$

¿qué clase maximiza $P(C_K | \bar{x})$? Aquella que tiene más instancias en $V_R \Rightarrow$ VECINOS MÁS PRÓXIMOS (K-NN)

Algoritmo de entrenamiento K-NN \rightarrow coste = $O(D \cdot N)$ "
 \uparrow si normalizamos

• Cosas previas importantes que hay que hacer:

1. Decidir métrica/distancia (euclídea, Manhattan, ...)

2. Normalizar datos: para cada atributo se calcula media y desviación típica: $\frac{\bar{x} - \mu}{\sigma} = \bar{x}_{\text{new}}$

Se guardan los valores de la media y la desv. std.

En scikit-learn: `StandardScaler()`

Algoritmo de clasificación K-NN \rightarrow coste = $O(D \cdot N)$

• Algoritmo de K-NN: normalizamos los datos, calculamos distancias a todos los pto. y nos quedamos con los K más cercanos. Devolvemos la moda de las clases en K.

MODELOS LINEALES

[...]

Si $P(C_1) = P(C_2)$ y $\Sigma_1 = \Sigma_2 = I$ entonces:

$$\text{Si } (\bar{x} - \mu_2)^2 > (\bar{x} - \mu_1)^2 \Rightarrow C_1$$

$$\text{Si } (\bar{x} - \mu_2)^2 \leq (\bar{x} - \mu_1)^2 \Rightarrow C_2$$

Es decir, \bar{x} se clasifica como C_1 si está más cerca de la media de los ejemplos de la clase 1.

bs: ¿Qué frontera de decisión tenemos? (A partir de ahora $P(C_1)$ no tiene por qué ser igual a $P(C_2)$)

$$P(C_1 | \bar{x}) = 0.5 \Rightarrow \sigma(a) = 0.5 \Rightarrow a = 0$$

$$a = \frac{(\bar{x} - \mu_2)^2}{2} - \frac{(\bar{x} - \mu_1)^2}{2} + \ln\left(\frac{P(C_1)}{P(C_2)}\right) = 0$$

$$\underbrace{(\mu_1 - \mu_2) \bar{x}}_{w} + \underbrace{\frac{\mu_2^2 - \mu_1^2}{2} + \ln\left(\frac{P(C_1)}{P(C_2)}\right)}_{w_0 = b} \quad \text{Ecuación de un hiperplano}$$

$$w = \mu_1 - \mu_2$$

$$w_0 = b = \frac{\mu_2^2 - \mu_1^2}{2} + \ln\left(\frac{P(C_1)}{P(C_2)}\right)$$

$\uparrow \uparrow$
dependiendo
del libro

- Si $P(C_1) = P(C_2) \Rightarrow \mu_1$ y μ_2
- Si $P(C_1) > P(C_2) \Rightarrow$ la frontera se aleja de μ_1 hacia μ_2 .

Mediatriz de

Si ahora suponemos:

- 2 clases: C_1 y C_2
 - $P(\bar{x}|C_1) = N(\mu_1, \Sigma)$
 - $P(\bar{x}|C_2) = N(\mu_2, \Sigma)$
- } Es decir, $\Sigma_1 = \Sigma_2 = \Sigma$

llegaríamos a que $w = \Sigma^{-1}(\mu_1 - \mu_2)$

$$b = w_0 = \frac{1}{2}\mu_1^T \Sigma^{-1} \mu_1 + \frac{1}{2}\mu_2^T \Sigma^{-1} \mu_2 + \ln \frac{P(C_1)}{P(C_2)}$$

Para entrenar un modelo lineal podríamos estimar $\mu_1, \mu_2, \Sigma, P(C_1)$ y $P(C_2)$ y con eso sacamos w y w_0 . ¿Cuántos parámetros son?

$$\begin{array}{cccccc} D & + & D & + & \frac{D(D+1)}{2} & + & 1 & + & 1 \\ \uparrow & & \uparrow & & \uparrow & & \uparrow & & \uparrow \\ \mu_1 & & \mu_2 & & \Sigma & & P(C_1) & & P(C_2) \end{array}$$

Realmente solo necesitamos saber w y w_0 : $D+1$ parámetros

Vamos a redefinir los vectores para simplificar notación:

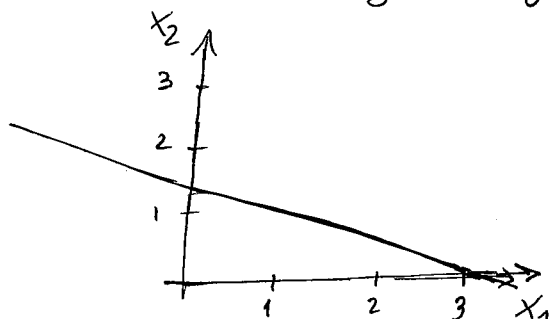
$$w = (w_0, \underbrace{w_1, \dots, w_d}_{\text{anterior } w})$$

$$\bar{x} = (1, \underbrace{x_1, \dots, x_d}_{\text{anterior } x})$$

de esta forma $P(C_1|\bar{x})$ queda $P(C_1|x) = \sigma(wx)$

Ejemplo: $w = (-1, 1/3, 2/3)$

¿frontera? $\sigma(wx) = 1/2 \Rightarrow wx = 0 \Rightarrow \langle (-1, 1/3, 2/3), (1, x_1, x_2) \rangle = 0 \Rightarrow$
 $\Rightarrow \frac{1}{3}x_1 + \frac{2}{3}x_2 - 1 = 0$



- frontera de decisión

Tenemos un conjunto de N ejemplos: $\mathcal{D} = \{ \bar{x}_j, t_j \}, j=1, \dots, N \}$ donde $t_j = 1$ si \bar{x}_j es C_1 , o $t_j = 0$ si \bar{x}_j es C_2 .

Vamos a definir el modelo de REGRESIÓN LOGÍSTICA.

TEMA 3

REGRESIÓN LOGÍSTICA

Se estima \bar{w} por MV: $P(t_1, \dots, t_N | \bar{w}, \bar{x}_1, \dots, \bar{x}_N) =$

$$= \prod_{j=1}^N P(t_j | \bar{w}, \bar{x}_j) \quad \text{donde } P(t_j | \bar{w}, \bar{x}_j) = \begin{cases} \sigma(\bar{w} \bar{x}_j) & \text{si } t_j = 1 \\ 1 - \sigma(\bar{w} \bar{x}_j) & \text{si } t_j = 0 \end{cases}$$

Podemos escribir: $P(t_j | \bar{w}, \bar{x}_j) = \sigma(\bar{w} \bar{x}_j)^{t_j} \cdot (1 - \sigma(\bar{w} \bar{x}_j))^{1-t_j}$

$$\Rightarrow P(t_1, \dots, t_N | \bar{w}, \bar{x}_1, \dots, \bar{x}_N) = \prod_{j=1}^N P(t_j | \bar{w}, \bar{x}_j) = \prod_{j=1}^N \sigma(\bar{w} \bar{x}_j)^{t_j} \cdot (1 - \sigma(\bar{w} \bar{x}_j))^{1-t_j}$$

Cómo queremos calcular \bar{w} que maximice lo anterior, derivamos:

$$\text{Antes observar: } \prod_{j=1}^N \sigma(\bar{w} \bar{x}_j)^{t_j} \cdot (1 - \sigma(\bar{w} \bar{x}_j))^{1-t_j} = \sum_{j=1}^N t_j \log(\sigma(\bar{w} \bar{x}_j)) + (1-t_j) \log(1 - \sigma(\bar{w} \bar{x}_j))$$

Ahora minimizamos $-\sum_{j=1}^N t_j \cdot \log(\sigma(\bar{w} \bar{x}_j)) + (1-t_j) \log(1 - \sigma(\bar{w} \bar{x}_j)) = E$

$$\frac{\partial E}{\partial \bar{w}} = -\sum_{j=1}^N \left[t_j \frac{1}{\sigma(\bar{w} \bar{x}_j)} - (1-t_j) \frac{1}{1 - \sigma(\bar{w} \bar{x}_j)} \right] \frac{\partial \sigma(\bar{w} \bar{x}_j)}{\partial \bar{w}} =$$

$$= -\sum_{j=1}^N \left[t_j \frac{1}{\sigma(\bar{w} \bar{x}_j)} - (1-t_j) \frac{1}{1 - \sigma(\bar{w} \bar{x}_j)} \right] \sigma(\bar{w} \bar{x}_j) (1 - \sigma(\bar{w} \bar{x}_j)) \bar{x}_j =$$

$$= -\sum_{j=1}^N \left[t_j (1 - \sigma(\bar{w} \bar{x}_j)) - (1-t_j) \sigma(\bar{w} \bar{x}_j) \right] \bar{x}_j =$$

$$= \sum_{j=1}^N (\sigma(\bar{w} \bar{x}_j) - t_j) \bar{x}_j = 0$$

Para mejorar la verosimilitud del ejemplo \bar{x}_j hay que mover la recta en sentido opuesto al gradiente $(\sigma_j - t_j) \bar{x}_j$ y proporcional a una constante de aprendizaje η : $\bar{w} = \bar{w} - \eta (\sigma_j - t_j) \bar{x}_j$

Algoritmo de regresión logística:

▷ Entrenamiento: parámetros η , épocas

- Generamos \bar{w} aleatorio con coef. en $[-0.5, 0.5]$.

- Para ep desde 0 a épocas:

 Para j desde 1 a N :

$\sigma = \sigma(\bar{w} \bar{x}_j)$ # prob. a posteriori de C_1

$\bar{w} = \bar{w} - \eta(\sigma - t_j) \bar{x}_j$

 return \bar{w}

VERSIÓN MAP DE REG. LOGÍSTICA

Verosimilitud:

$$P(t_1, \dots, t_N | \bar{w}, \bar{x}_1, \dots, \bar{x}_N) = \prod_{j=1}^N \sigma(\bar{w} \bar{x}_j)^{t_j} (1 - \sigma(\bar{w} \bar{x}_j))^{1-t_j}$$

Suponemos que el módulo de $|\bar{w}|$ viene dado por una gaussiana de media 0

$$P(|\bar{w}|) = \frac{1}{\sqrt{2\pi\sigma^2}} e^{-|\bar{w}|^2 / 2\sigma^2}$$

<u>OBS</u> : $\sigma^2 \equiv \text{varianza}$ $\sigma_j \equiv \sigma(\bar{w} \bar{x}_j)$

$$\text{Prob. a posteriori } P(\bar{w} | t_1, \dots, t_N, \bar{x}_1, \dots, \bar{x}_N) \simeq \underbrace{\frac{1}{\sqrt{2\pi\sigma^2}} e^{-|\bar{w}|^2 / 2\sigma^2}}_{\text{prior}} \underbrace{\prod_{j=1}^N \sigma_j^{t_j} (1 - \sigma_j)^{1-t_j}}_{\text{verosimilitud}}$$

Aplicando logaritmos:

$$E = \ln P(\bar{w} | t_1, t_2, \dots, \bar{x}_1, \dots, \bar{x}_N)$$

$$\frac{\partial E}{\partial \bar{w}} = 0 \iff \dots \iff \frac{\partial E}{\partial \bar{w}} = \sum_{j=1}^N (\sigma_j - t_j) \bar{x}_j + \frac{\bar{w}}{N\sigma^2}$$

Algoritmo reg. logística map:

reg-log-map(η , nepocas, σ^2):

- Generar \bar{w} aleatoriamente con coefs. en $[-0.5, 0.5]$ $O(D)$

- For $i: 0 \rightarrow \text{nepocas}$:

For $j: 1 \rightarrow N$:

$$\bar{w} = \bar{w} - \eta (\sigma_j - t_j) \bar{x}_j - \frac{\eta \bar{w}}{N \sigma^2}$$

$O(\text{nepocas} \cdot N \cdot D)$

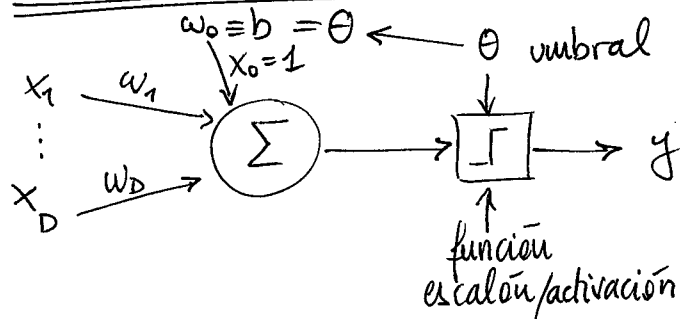
Coste de entrenamiento $\rightarrow O(\text{nepocas} \cdot N \cdot D)$

Coste de clasificación $\rightarrow O(D)$ (coste de realizar un producto escalar)

Observaciones finales:

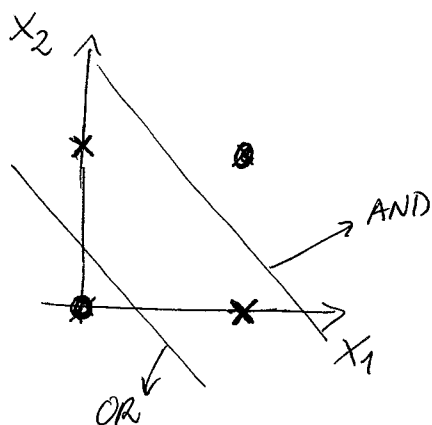
- RegLog trabaja solo con atributos numéricos
- Se recomienda normalizar atributos antes de entrenar

REDES NEURONALES

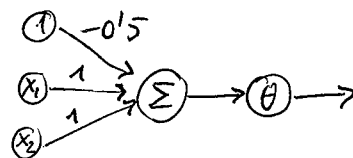


$$y(\bar{x}) = \begin{cases} 1, & \text{si } \bar{w}\bar{x} \geq 0 \\ 0, & \text{si } \bar{w}\bar{x} < 0 \end{cases}$$

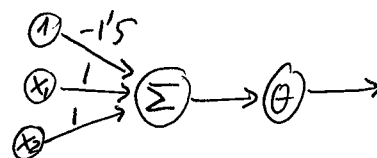
PERCEPTRÓN
DE ROSENBLATT



$$z_1 = \theta(x_1 + x_2 - 0.5)$$



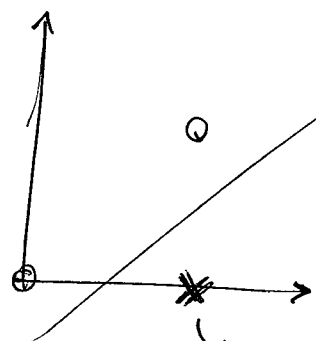
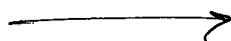
$$z_2 = \theta(x_1 + x_2 - 1.5)$$



$\Theta :=$ función de Heaviside

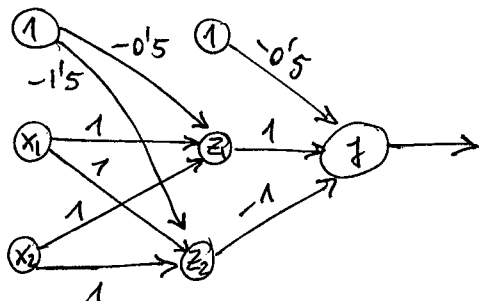
x_1	x_2	t	z_1	z_2
0	0	0	0	0
0	1	1	1	0
1	0	1	1	0
1	1	0	1	1

Hemos definido
nuevos atributos



$$y = \theta(z_1 - z_2 - 0.5)$$

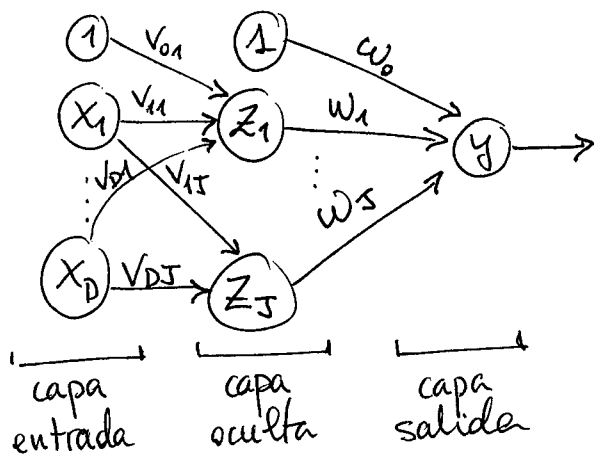
z_1	z_2	y
0	0	0
1	0	1
1	0	1
1	1	0



Regla delta

Sustituir $\Theta \rightarrow \sigma$
f. escalón f. sigmoideal

Red general de 3 capas



no. parámetros: $(D+1)J + J + 1$

v_{DJ} := peso entre la neurona D (entrada) y la J (oculta)

w_j := peso entre neurona j (oculta) y la capa de salida

$$z_j = \sigma\left(\sum_{d=0}^D x_d v_{dj}\right)$$

$$y(\bar{x}) = \sigma\left(\sum_{j=0}^J w_j z_j\right) = \sigma\left(\sum_{j=0}^J w_j \sum_{d=0}^D x_d v_{dj}\right)$$

$$\mathcal{D} = \{(\bar{x}_n, t_n)\}_{n=1}^N \text{ (conjunto de datos)}$$

$$P(t_1, \dots, t_N | v_{00} \dots v_{DJ}, w_0, \dots, w_J, \bar{x}_1, \dots, \bar{x}_N) =$$

$$= \prod_{n=1}^N P(c_n | \bar{x}_n, v, \dots, w, \dots)^{t_n} (1 - P(c_n | \bar{x}_n, v, \dots, w, \dots))^{1-t_n} =$$

$$= \prod_{n=1}^N y(\bar{x}_n)^{t_n} (1 - y(\bar{x}_n))^{1-t_n}$$

Aplicando logaritmos:

$$E(v_0, \dots, w_0, \dots) = -\log P(t_1 \dots t_N | v \dots w \dots \bar{x} \dots \bar{x}_N) =$$

$$= \sum_{n=1}^N E_n(v \dots w) \text{ con } E_n = -t_n \log y(\bar{x}_n) - (1-t_n) \log(1-y(\bar{x}_n))$$

Vamos a derivar lo anterior.

Derivada comodín: $\frac{\partial E_n}{\partial \beta} = -\frac{t_n}{y(\bar{x}_n)} \cdot \frac{\partial y(\bar{x}_n)}{\partial \beta} + \frac{1-t_n}{1-y(\bar{x}_n)} \cdot \frac{\partial y(\bar{x}_n)}{\partial \beta} = (\star)$

$$\frac{\partial y(\bar{x}_n)}{\partial \beta} = \frac{\partial}{\partial \beta} \cdot \frac{1}{1+e^{-wz}} = y(\bar{x}_n) \cdot (1-y(\bar{x}_n)) \cdot \frac{\partial wz}{\partial \beta}$$

$$\Rightarrow (\star) = -\frac{t_n}{y(\bar{x}_n)} \cdot \frac{\partial y(\bar{x}_n)}{\partial \beta} + \frac{1-t_n}{1-y(\bar{x}_n)} \cdot y(\bar{x}_n)(1-y(\bar{x}_n)) \frac{\partial wz}{\partial \beta} \Rightarrow$$

$$\Rightarrow (\star) = \left[-t_n(1-y(\bar{x}_n)) + (1-t_n)y(\bar{x}_n) \right] \frac{\partial wz}{\partial \beta} = (y(\bar{x}_n) - t_n) \frac{\partial wz}{\partial \beta}$$

Derivamos con respecto a w_i : (renombramos $\beta \rightarrow w_i$)

$$\frac{\partial wz}{\partial w_i} = \frac{\partial}{\partial w_i} (w_0 + w_1 z_1 + \dots + w_I z_I) = z_i$$

$$\Rightarrow \boxed{\frac{\partial E_n}{\partial w_i} = (y(\bar{x}_n) - t_n) z_i}$$

Derivamos con respecto a V_{pq} : (renombramos $\beta \rightarrow V_{pq}$):

$$\frac{\partial wz}{\partial V_{pq}} = \frac{\partial}{\partial V_{pq}} \sum_{j=0}^J w_j z_j = \sum_{j=0}^J w_j \frac{\partial}{\partial V_{pq}} \sigma \left(\sum_{d=0}^D x_{nd} V_{dj} \right) =$$

$$= \sum_{j=0}^J w_j z_j (1-z_j) \frac{\partial}{\partial V_{pq}} \sum_{d=0}^D x_{nd} V_{dj} = \sum_{j=0}^J w_j z_j (1-z_j) \sum_{d=0}^D x_{nd} \frac{\partial V_{dj}}{\partial V_{pq}} =$$

$$= \sum_{j=0}^J w_j z_j (1-z_j) \sum_{d=0}^D x_{nd} \overset{\substack{\text{delta de} \\ \text{Kronecker}}}{\delta_{dp} \delta_{jq}} = \sum_{j=0}^J w_j z_j (1-z_j) x_{np} \delta_{jq} =$$

$$= w_q z_q (1-z_q) x_{np}$$

$$\Rightarrow \boxed{\frac{\partial E_n}{\partial V_{pq}} = (y(\bar{x}_n) - t_n) w_q z_q (1-z_q) x_{np}}$$

Reglas de actualización :

$$w_i = w_i - \eta (y(\bar{x}_n) - t_n) z_i$$

$$V_{pq} = V_{pq} - \eta (y(\bar{x}_n) - t_n) w_q z_q (1 - z_q) x_{np}$$

Pseudocódigo

NN_model (η , nepocas) :

inicializar pesos de la red aleatoriamente $\in [-0.5, 0.5]$.

for $ic = 1 : nepocas$

for $n = 1 : N$

$$z_q = \sigma \left(\sum_{d=0}^D x_{nd} v_{dq} \right) \rightarrow O(DJ)$$

$$y = \sigma \left(\sum_{j=0}^J w_j z_j \right) \rightarrow O(J)$$

$$V_{pq} = V_{pq} - \eta (y - t_n) w_q z_q (1 - z_q) x_{np}$$

$$w_q = w_q - \eta (y - t_n) z_q \rightarrow O(J)$$

backward
prop. $\rightarrow O(DJ)$

Coste de entrenamiento : $O(\text{nepocas} \cdot N \cdot D \cdot J)$

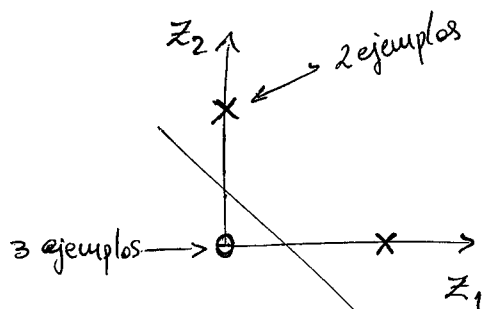
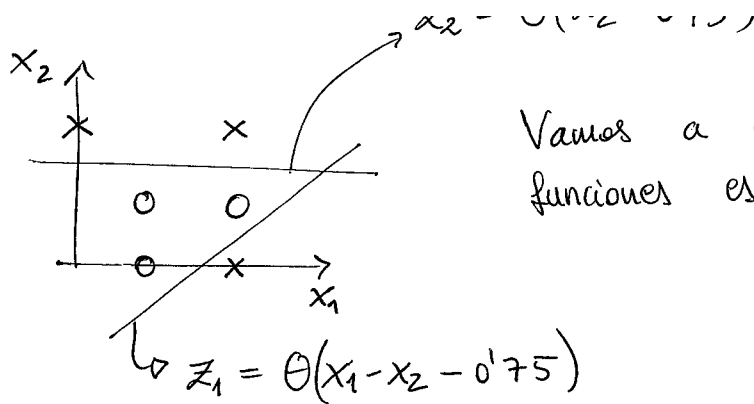
Coste de clasificación : $O(D \cdot J)$

Ejemplo 1:

x_1	x_2	t
0.5	0	0
1	0	1
1	0.5	0
0	1	1
0.5	0.5	0
1	1	1

↓

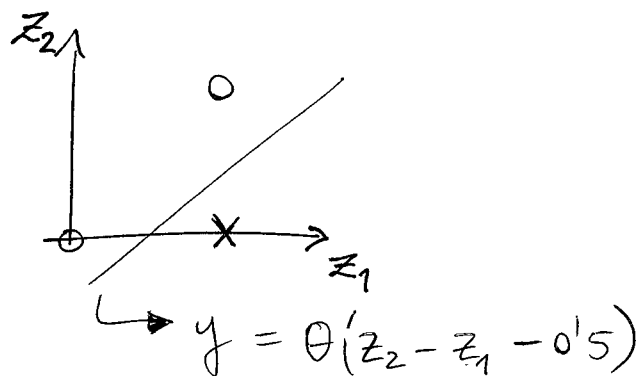
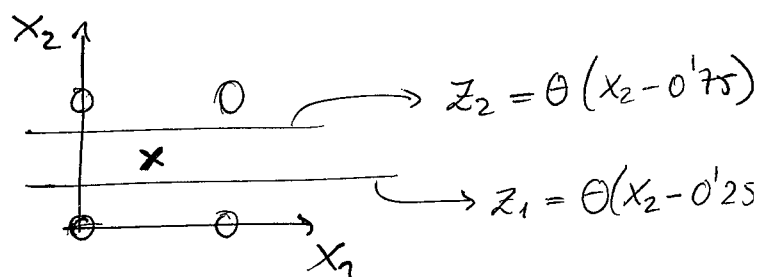
x_1	x_2	t	z_1	z_2
0.5	0	0	0	0
1	0	1	1	0
1	0.5	0	0	0
0	1	1	0	1
0.5	0.5	0	0	0
1	1	1	0	1



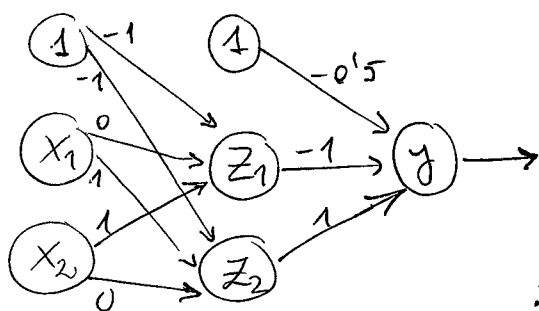
cuidado con este signo en la última capa

Ejemplo 2

x_1	x_2	t	z_1	z_2
0	0	0	0	0
0	1	0	1	1
1	0	0	0	0
1	1	0	1	1
0.5	0.5	1	1	0



Ejemplo 3



$$z_1 = \theta(x_2 - 1)$$

$$z_2 = \theta(x_1 - 1)$$

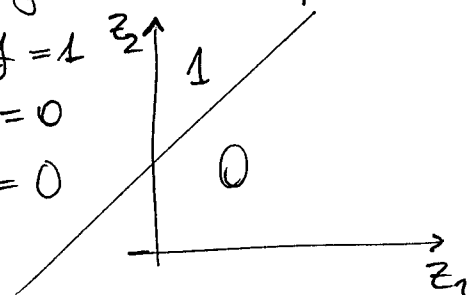
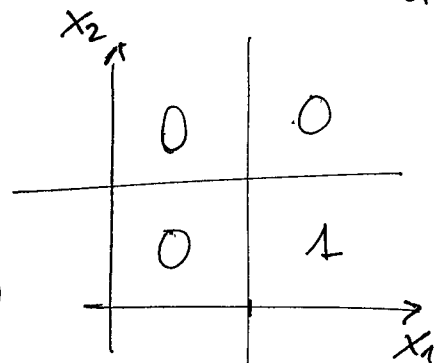
$$y = \theta(z_2 - z_1 - 0.5)$$

$$\vec{x} = (0,0) \Rightarrow \vec{z} = (0,0) \Rightarrow y = 0$$

$$\vec{x} = (2,0) \Rightarrow \vec{z} = (0,1) \Rightarrow y = 1$$

$$\vec{x} = (0,2) \Rightarrow \vec{z} = (1,0) \Rightarrow y = 0$$

$$\vec{x} = (2,2) \Rightarrow \vec{z} = (1,1) \Rightarrow y = 0$$



Pregunta: ¿Qué pasaría en estos ejemplos si la func. de activación es la sigmoideal σ ?
(examen)

TEMAS

ALGORITMOS GENÉTICOS

Los algoritmos genéticos utilizan las ideas de la evolución para resolver problemas de optimización.

- problema
- cada individuo codifica una solución
- función de ajuste (fitness)
- conjunto de soluciones
- mecanismos aleatorios de adaptación
 - los individuos con mayor fitness tienen más probabilidad de reproducirse
 - modificaciones aleatorias en las soluciones (mutaciones)

Algoritmo genético genérico

ag(fitness f , params de evolución):

Init población inicial aleatoria P

Mientras no se cumpla criterio-fin:

S = selección - progenitores (P)

S = recombinación (S)

S = mutación (S)

P = selección - supervivientes (S, P)

return best(P)

¿Por qué funcionan los AG? TEORIA DE ESQUEMAS

Los algoritmos genéticos funcionan descubriendo y recombinando buenos bloques de soluciones en paralelo.

La teoría de esquemas formaliza el concepto de bloque:

- La codificación binaria un bloque o esquema es una cadena de 0, 1 y *, donde * representa 0 o 1.

$H = 1 * * * 0 \rightarrow$ todas las soluciones de longitud 5 que empiezan por 1 y acaban por 0.

- Orden de un esquema H : $O(H) = \# \text{ bits definidos}$
- Longitud de un esquema H : $d(H) = \text{distancia entre los bits definidos más lejanos.}$
- Número de cadenas binarias posibles de long. L : 2^L
- Número de cadenas con 1, 0, * posibles de long. L : 3^L
- Cada cadena definida de longitud L es una instancia de 2^L esquemas: la cadena 101 pertenece a 2^3 esquemas.
- Una población de N individuos contiene instancias de entre 2^L esquemas (si son todos los individuos iguales) a $N2^L$ esquemas
- A pesar de estar evaluando solo N individuos, estos en realidad implícitamente ajustan más esquemas.

$$d(0 * * * 1) = 4$$

$$d(01*) = 1$$

$$d(0***) = 0$$

Notación

$n_H(t) \equiv$ número de instancias del esquema H en tiempo t .

$f_i(t) \equiv$ fitness individuo i en tiempo t .

$\bar{f}(t) \equiv$ fitness medio de la población en tiempo t .

$\bar{f}_H(t) \equiv$ fitness medio del esquema H en tiempo t .

Queremos obtener el valor esperado de instancias en $t+1$ para H : $E[n_H(t+1)]$

Paso 1: Selección proporcional a fitness

$$\mathbb{E}_s [n_i(t+1)] = \frac{f_i(t)}{\sum_i f_i(t)} \cdot N = \frac{f_i(t)}{\bar{f}(t)}$$

$$\mathbb{E}_s [n_H(t+1)] = \frac{1}{\bar{f}(t)} \sum_{i \in H} f_i(t) = n_H(t) \cdot \frac{\bar{f}_H(t)}{\bar{f}(t)}$$

Si $\bar{f}_H(t) > \bar{f}(t) \Rightarrow \mathbb{E}_s [n_H(t+1)] > n_H(t)$
en media

Paso 2: Cruce en un punto

Suponemos p_c de cruce. Tenemos que determinar si un esquema sobrevive al cruce. Decimos que sobrevive si al menos un descendiente pertenece al esquema H . Una cota inferior es:

$$S_c = 1 - p_c \cdot \frac{d(H)}{L-1}$$

Cuanto mayor sea $d(H)$ menor probabilidad de supervivencia tiene.

Paso 3: Mutación bitflip

Una instancia sobrevive si no se cambia ningún bit definido en H . $S_M = (1 - p_M)^{o(H)}$

$$\text{Finalmente: } \mathbb{E} [n_H(t+1)] = \frac{\bar{f}_H(t)}{\bar{f}(t)} n_H(t) \cdot \left(1 - p_c \frac{d(H)}{L-1}\right) \cdot (1 - p_M)^{o(H)}$$

ÁRBOLES DE DECISION

PROBLEMA DE LA SUMA (NP)

Lista de números: 3, 8, 14, 2, 25, 5, 7 (long. K)
Número objetivo: 15

Hay que seleccionar los números de la lista que sumen el número objetivo.

Queremos solucionar el problema utilizando un algoritmo genético. Necesitamos codificar las soluciones, cruce y mutación.

• Representación: lista binaria de longitud K
opc. 1 $[1, 0, 0, 0, 1, 0, 0] \equiv 3 + 25 = 28$
opc. 2 \rightarrow lista de enteros de longitud $\leq K$.
 $[3, 25] \equiv 3 + 25 = 28$

• Fitness: opc. 1 $f(x) = \frac{1}{|\text{obj} - \text{Suma}| + 1}$
opc. 2 $\rightarrow f(x) = -|\text{obj} - \text{Suma}|$

• Cruce: opc. 1 \rightarrow cruce en un punto $\begin{array}{ccc|cc} 1 & 0 & 0 & 1 & 0 & 0 \\ 0 & 1 & 1 & 0 & 1 & 0 \end{array} \rightarrow \begin{array}{ccc|cc} 1 & 1 & 0 & 1 & 1 & 0 \\ 0 & 1 & 1 & 0 & 0 & 0 \end{array}$

opc. 2 \rightarrow cruce en un punto con eliminación de repeticiones $\begin{array}{c|cc} 3 & 7 & 8 \\ 2 & 3 & 5 \end{array} \rightarrow \begin{array}{c|cc} 3 & 5 \\ 2 & 7 & 8 \end{array}$

opc. 1 \rightarrow mejorada \rightarrow cruce uniforme $\begin{array}{cccccc} C & X & X & X & C \\ 0 & 0 & 1 & 1 & 0 \\ 1 & 0 & 1 & 0 & 1 \end{array} \rightarrow \begin{array}{ccccc} 0 & 0 & 1 & 0 & 0 \\ 1 & 0 & 1 & 1 & 1 \end{array}$ $C = \text{arriba}$
 $X = \text{abajo}$

• Mutación $\xrightarrow{\text{opc. 1}}$ bit-flip
 \downarrow
 $\xrightarrow{\text{opc. 2}}$ añadir/quitar 1 número