

APRENDIZAJE AUTOMÁTICO

CLASIFICACIÓN CON K-VECINOS MÁS PRÓXIMOS (K-NN)

ENTRADA: datos entrenamiento $D_{\text{train}}: \{(\underline{x}_n, c_n), n=1, \dots, N\}$

ejemplos de test $D_{\text{test}} \circ \underline{x}_{\text{test}}$

número de vecinos K

métrica o función distancia $d(\cdot, \cdot)$

SALIDA: clase predicha por K-NN para el ejemplo de test $\underline{x}_{\text{test}}$

CÓDIGO:

1. Encontrar los K -vecinos más próximos a $\underline{x}_{\text{test}}$ en D_{train} .
2. Encontrar la clase mayoritaria entre esos K -vecinos más próximos a $\underline{x}_{\text{test}}$.
3. Asignar a $\underline{x}_{\text{test}}$ esa clase mayoritaria

OBSERVACIONES:

- No funciona bien si la dimensión D es muy grande. Una forma de intentar resolver esto sería seleccionar las variables más significativas.
- Si hay ruido no funciona bien, debido a que el ruido "contamina" la distancia.

TEOREMA DE BAYES EN K-NN

Denominemos: $R(\underline{x}_{\text{test}})$ = elemento de volumen centrado en $\underline{x}_{\text{test}}$.

$N(R(\underline{x}_{\text{test}}))$ = ejemplos contenidos en $R(\underline{x}_{\text{test}})$

$N_c(R(\underline{x}_{\text{test}}))$ = ejemplos de la clase c en $R(\underline{x}_{\text{test}})$.

Entonces $P(c|\underline{x}_{\text{test}}) \approx \frac{N_c(R(\underline{x}_{\text{test}}))}{N(R(\underline{x}_{\text{test}}))}$

MAP: $C^*(\underline{x}_{\text{test}}) = \max_c P(c|\underline{x}_{\text{test}}) \approx \max_c N_c(R(\underline{x}_{\text{test}}))$

tanto, K-NN puede ser visto como un predictor MAP que utiliza estimaciones locales de las densidades de probabilidad evantes.

¿CUANTOS VECINOS? ¿VALOR DE K ÓPTIMO?

El n° de vecinos actúa como un "smoothing parameter"

K pequeño: muchas regiones y más pequeñas

K grande: menos regiones y más grandes.

El K óptimo puede ser determinado mediante algoritmos de validación cruzada:

- LOOCV: Leave One Out Cross Validation
- K-FCV: K-fold Cross Validation

CLASIFICACIÓN CON ÁRBOLES DE DECISIÓN

Buscamos la primera pregunta que mejor determine la clasificación lo mejor posible.

La mejor pregunta es la que permite obtener más información sobre la clase.

Necesitamos cuantificar la ganancia de información

$$\underline{P} = \{P_k\}_{k=1}^K \quad H(\underline{P}) = \sum_{k=1}^K P_k \cdot \underbrace{\log_2(P_k)}_{\substack{\text{unidades BITS} \\ \rightarrow \text{INATS}}} \rightarrow \text{ENTROPIA}$$

La ganancia de información aportada por la variable aleatoria X sobre la variable G se define como:

$$GI(X, G) = H(G) - H(G|X)$$

Donde: $H(G)$ = mide el desconocimiento que tenemos de la variable G antes de saber X .

$H(G|X)$ = mide el desconocimiento sobre G después de saber el valor de X .

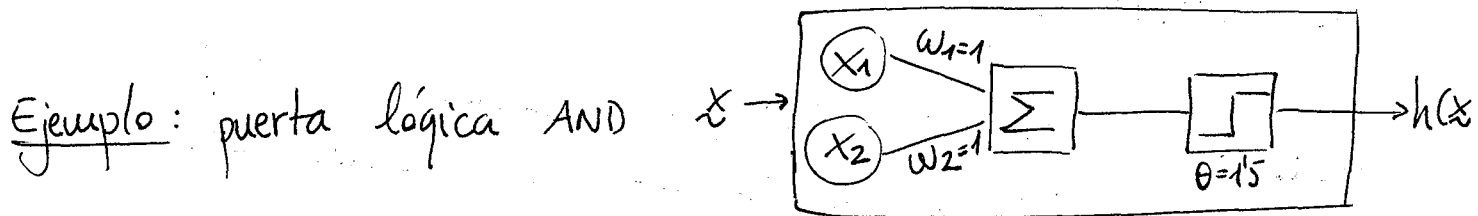
$$h(x) = w_0 + w_1 x_1 + w_2 x_2 + \dots + w_D x_D = (w_0 \ w_1 \ \dots \ w_D) \begin{pmatrix} x_0 \\ x_1 \\ \vdots \\ x_D \end{pmatrix} \text{ con } \underline{x_0 = 1}$$

↑
modelo para predecir y

$\underline{w}^T \cdot \underline{x}$

$$ECM(\underline{w}) = \frac{1}{N_{\text{train}}} \sum_{n=1}^{N_{\text{train}}} (\underline{w}^T \cdot \underline{x}_n - y_n)^2$$

$$\underline{w}_{LS}^* = \underset{\underline{w}}{\operatorname{argmin}} [ECM(\underline{w})] \quad \text{Mínimos cuadrados (Least squares)}$$



lo mismo para la puerta OR cambiando $\theta = 1.5$ por $\theta = 0.5$
Existen más soluciones.

MÉTODOS DE APRENDIZAJE

Ⓘ MÁXIMO A PRIORI : clasificación por mayoría
(maximizando priores) $h(\underline{x}) = \operatorname{argmax}_{\text{clase} \in \{c_1, \dots, c_k\}} P(\text{clase})$

Ⓙ ML $h(\underline{x}) = \operatorname{argmax}_{\text{clase} \in \{c_1, \dots, c_k\}} P(\{(\underline{x}_n, c_n)\}_{n=1}^{N_{\text{train}}} | \text{clase})$

Ⓢ MAP $h(\underline{x}) = \operatorname{argmax}_{\text{clase} \in \{c_1, \dots, c_k\}} P(\{(\underline{x}_n, c_n)\}_{n=1}^{N_{\text{train}}} | \text{clase}) \cdot P(\text{clase})$

FEED-FORWARD MULTILAYER PERCEPTRÓN

- **PARÁMETROS:** pesos sinápticos de conexiones entre neuronas

- **HIPERPARÁMETROS**

- learning rate

- tipo y peso de regularización
(TIKHONOV)

- resolución de ED's
álgebra lineal

- modelos más robustos

- **ARQUITECTURA:** grafo de la red neuronal

Observación: el ruido regulariza → el descenso del gradiente estocástico es mejor que el descenso del gradiente normal.

REGRESIÓN LOGÍSTICA (método lineal para clasificación)

SIGMOIDAL

$$\sigma(z) = \frac{1}{1 + e^{-z}} \quad \begin{array}{l} z \rightarrow -\infty \rightarrow 0 \\ z \rightarrow \infty \rightarrow 1 \end{array}$$

$$\sigma(0) = 1/2$$

LOGIT

$P(C_n = 1 | \underline{x}_n)$ Salida de la Sigmoidal \equiv
 \equiv Posterior de la clase 1

TANH

$$f(z) = \frac{e^z - e^{-z}}{e^z + e^{-z}} \quad \begin{array}{l} z \rightarrow -\infty \rightarrow -1 \\ z \rightarrow \infty \rightarrow 1 \end{array}$$

ROBIT

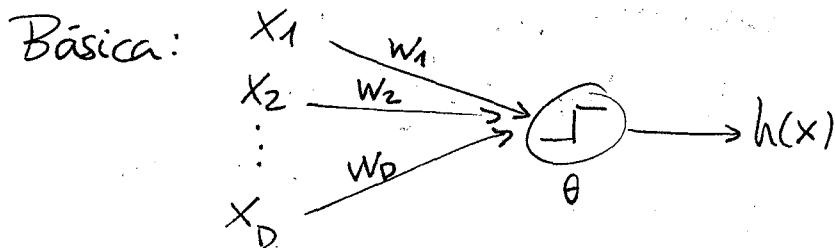
$$\text{normpdf} = \frac{e^{-z^2/2}}{\sqrt{2\pi}} \quad N(0,1)$$

$$\text{normcdf}(z) = \int_{-\infty}^z \text{normpdf}(z') dz'$$

$$P(Z \leq z)$$

$$Z \sim N(0,1)$$

REDES NEURONALES



$$h(x) = \begin{cases} 1 & \text{si } w_1x_1 + \dots + w_dx_d \geq \theta \\ 0 & \text{si } w_1x_1 + \dots + w_dx_d < \theta \end{cases}$$

umbral: θ

$h(x) :=$ hipótesis

σ función de activación

$$h(x) = \begin{cases} 1 & \text{si } \underline{w}^T \underline{x} \geq \theta \\ 0 & \text{si } \underline{w}^T \underline{x} < \theta \end{cases}$$

FUNCIONES DE ACTIVACIÓN PRINCIPALES:

- FUNCIÓN SIGMOIDE; $\sigma(z) = \frac{1}{1+e^{-z}}$ (LOGIT)

$$\sigma(-\infty) = 0 \quad \sigma(0) = \frac{1}{2}$$

$$\sigma(\infty) = 1$$

monótona creciente

$$\sigma'(z) = \sigma(z)\sigma(-z)$$

- FUNCIÓN DE ACTIVACIÓN LINEAL: $f(z) = z$
Neurona de salida para regresión

- TANGENTE HIPERBÓLICA: $\tanh(z) = \frac{e^z - e^{-z}}{e^z + e^{-z}}$

$$\tanh(-\infty) = -1$$

$$\tanh(0) = 0$$

$$\tanh(\infty) = 1$$

- FUNCIÓN DE ACTIVACIÓN PROBIT: (GAUSSIANA)

$$\text{normpdf}(z) = \frac{1}{\sqrt{2\pi\sigma^2}} \int_{-\infty}^z e^{-\frac{y^2}{2\sigma^2}} dy$$

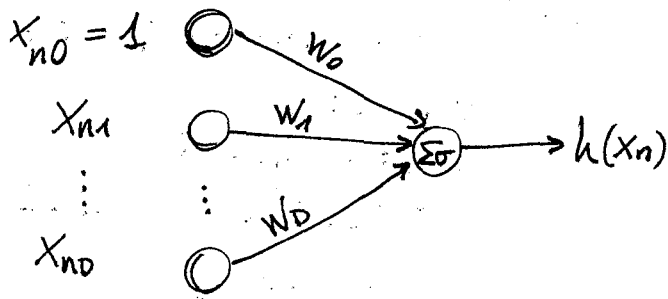
ARQUITECTURA BÁSICA

Una red neuronal es un conjunto de neuronas interconectadas.

La arquitectura de una red neuronal está determinada por la forma en que las neuronas están conectadas entre sí.

- Nodos fuente o neuronas de entrada (input layer): proporcionan la señal de entrada.
- Nodos de computación o neuronas ocultas (hidden layer): computan las señales de la capa anterior junto con sus pesos.
- Nodos de salida o neuronas de salida (output layer)

REGRESIÓN LOGÍSTICA MEDIANTE UN PERCEPTRÓN LINEAL



$$w_0 x_1 + \dots + w_d x_d = \underline{w}^T \underline{x} \Rightarrow$$

$$\Rightarrow \underline{w}^T \underline{x}_n = \sum_{d=0}^D w_d x_{nd}$$

Función de activación: $\sigma(\underline{w}^T \underline{x}_n)$

Salida: $h(x_n) = \sigma(\underline{w}^T \underline{x}_n)$

APRENDIZAJE POR MÁXIMA VEROSIMILITUD

Datos: $\{(x_n, t_n)\}_{n=1}^N$

$$\underline{w}_{ML} = \underset{\underline{w}}{\operatorname{argmax}} \left\{ \mathcal{L}(\{(x_n, t_n)\}_{n=1}^N; \underline{w}) \right\}$$
$$= \mathbb{P}(\{t_n\}_{n=1}^N | \{x_n\}_{n=1}^N, \underline{w})$$

Función de error: entropía cruzada

$$E(\underline{w}) = \sum_{n=1}^N E_n(\underline{w}) \quad \text{con} \quad E_n(\underline{w}) = t_n \log \underset{\uparrow \underline{w}}{h(x_n)} + (1-t_n) \log \underset{\uparrow \underline{w}}{(1-h(x_n))}$$

Aprendizaje ML \equiv Minimizar $E(\underline{w})$

Utilizando esto, basamos los siguientes algoritmos de aprendizaje usando redes neuronales.

APRENDIZAJE POR LOTES (Batch)

INPUT: datos de entrenamiento: $\{(x_n, t_n)\}_{n=1}^N$ donde $t_n = \begin{cases} 1 & \text{si } c_n = C_1 \\ 0 & \text{si } c_n = C_0 \end{cases}$

parámetro de aprendizaje: η

OUTPUT: parámetros del modelo (pesos de la red): \underline{w}_{ML}

PSEUDOCÓDIGO: 1. Inicializa los pesos aleatoriamente $\underline{w} \sim U[-0.5, 0.5]^{D+}$

2. Mientras no haya convergencia (o $\text{niter} < \text{max-iter}$):

2.1. Calcular salida de la red con los pesos actuales
 $h(x_n) = \sigma(\underline{w}^T x_n)$ para $n=1, 2, \dots, N$

2.2. Calcular gradiente (descenso del gradiente).

$$\frac{\partial E(\underline{w})}{\partial \underline{w}} = \sum_{n=1}^N (h(x_n) - t_n) x_n$$

2.3. Actualizar pesos:

$$\underline{w} = \underline{w} - \eta \frac{\partial E(\underline{w})}{\partial \underline{w}}$$

APRENDIZAJE EN LÍNEA (Online)

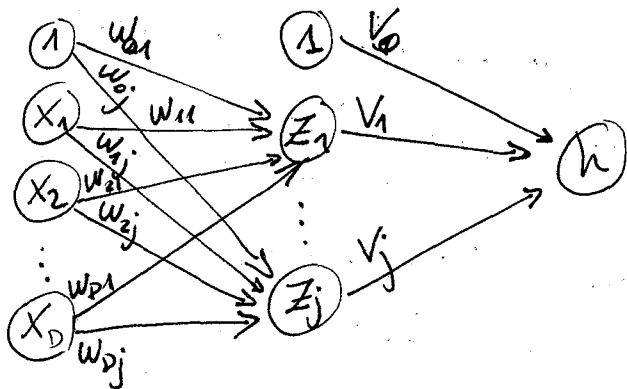
INPUT: datos de entrenamiento $\{(x_n, t_n)\}_{n=1}^N$ $t_n = \begin{cases} 1 & \text{si } C_n = C_1 \\ 0 & \text{si } C_n = C_0 \end{cases}$
parámetro de aprendizaje: η

OUTPUT: parámetros del modelo (pesos de la red): \underline{w}_{ML}

PSEUDOCÓDIGO:

1. Inicializar los pesos aleatoriamente $\underline{w} \sim U(0.5, 0.5)^{D+1}$
2. Inicializar el contador de épocas: $K = 0$
3. Mientras no haya convergencia (o $\text{niter} < \text{max_iter}$):
 - 3.1. Actualizar contador de épocas: $K += 1$
 - 3.2. Calcular la salida de la red:
$$h(x_n) = \sigma(\underline{w}^T x_n) \quad n=1, \dots, N$$
 - 3.3. Gradiente:
$$\frac{\partial E(\underline{w})}{\partial \underline{w}} = \sum_{n=1}^N (h(x_n) - t_n) x_n$$
 - 3.4. Actualizar pesos:
$$\underline{w} = \underline{w} - \eta \frac{\partial E}{\partial \underline{w}}(\underline{w})$$

PERCEPTRÓN MULTICAPA:



$$z_j = \sigma \left(\sum_{d=0}^D w_{dj} x_d \right) = \sigma(\underline{w}_j^T \underline{x})$$

$j=1, \dots, J$ donde $J := \# \text{neuronas en capa oculta}$

$$h(x) = \sigma \left(\sum_{j=0}^J v_j z_j \right) = \sigma(\underline{v}^T \underline{z})$$

PROPAGACIÓN HACIA ATRÁS (BACK PROPAGATION)

INPUT: datos de entrenamiento $\{(x_n, t_n)\}_{n=1}^N$ con $t_n = \begin{cases} 1 & \text{si } C_n = C_1 \\ 0 & \text{si } C_n = C_0 \end{cases}$
parámetro de aprendizaje: η

OUTPUT: pesos del perceptrón multicapa $\underline{W} \leftarrow$ matriz (por capa)

SEUDOCÓDIGO: 1. Inicializar pesos aleatoriamente

$$\underline{W}, \underline{V} \sim U[-0.5, 0.5]$$

2. Inicializar contador de épocas: $nEpoch = 0$

3. Mientras no haya convergencia (o $niter < max_iter$):

3.1. Actualizar contador de épocas: $nEpoch += 1$

3.2. For $n=1, \dots, N$:

$$z_{nj} = \sigma(v_j^T x_n) \quad j=1, \dots, J$$

$$h(x_n) = \sigma(\underline{W}^T \cdot \underline{Z}_n)$$

$$\delta_n = h(x_n) - t_n \quad (\text{error de predicción})$$

$$\Delta_{nj} = z_{nj}(1 - z_{nj})w_j \delta_n \quad j=1, \dots, J$$

(asigna error a ocultas)

$$\underline{W}' = \underline{W} - \eta \delta_n \underline{Z}_n \quad (\text{actualizar pesos de oculta a salida})$$

$$v_j' = v_j - \eta \Delta_{nj} x_n \quad j=1, 2, \dots, J$$

(actualizar pesos de entrada a ocultas)

APRENDIZAJE AUTOMÁTICO

$$\mathbb{E}(x^{(i)}) \equiv \text{media} \left(\{x_n^{(i)}\}_{n=1}^{N_{\text{train}}} \right) = \frac{1}{N_{\text{train}}} \sum_{n=1}^{N_{\text{train}}} x_n^{(i)}$$

$$V(\{x_n\}_{n=1}^{N_{\text{train}}}) \equiv \text{varianza}(\{x_n\}_{n=1}^{N_{\text{train}}}) = \mathbb{E} \left((x^{(i)} - \mathbb{E}(x^{(i)}))^2 \right) =$$

$$= \frac{1}{N_{\text{train}}} \sum_{n=1}^{N_{\text{train}}} (x_n^{(i)} - \mathbb{E}(x^{(i)}))^2$$

$$\text{Std}(\{x_n\}_{n=1}^{N_{\text{train}}}) = \sqrt{V(\{x_n\}_{n=1}^{N_{\text{train}}})} = \sqrt{\text{varianza}}$$

REESCALADO/NORMALIZACIÓN

$x_j^{(i)} \longrightarrow \frac{x_j^{(i)} - \text{mediana}(\{x_n^{(i)}\}_{n=1}^{N_{\text{train}}})}{\text{espacio intercuartílico}(\{x_n^{(i)}\}_{n=1}^{N_{\text{train}}})}$

estimación ROBUSTA del "centro" de la distribución
 estimación más ROBUSTA de la escala de la DIST.

Rango Intercuartílico $\equiv \text{percentil}(\{x_n\}_{n=1}^{N_{\text{train}}}, 0.75) - \text{percentil}(\{x_n\}_{n=1}^{N_{\text{train}}}, 0.25)$

Mediana $\equiv \text{percentil}(\{x_n\}_{n=1}^{N_{\text{train}}}, 0.5)$

$$\frac{x_j^{(i)} - \left(\frac{\max + \min}{2} \right)}{[\max(\{x_n\}_{n=1}^{N_{\text{train}}}) - \min(\{x_n\}_{n=1}^{N_{\text{train}}})]}$$

aquí suponemos $\tilde{x}_n \sim U[-1, 1]$

$$\frac{x_j^{(i)} - \mathbb{E}(\{x_n\}_{n=1}^{N_{\text{train}}})}{\text{std}(\{x_n\}_{n=1}^{N_{\text{train}}})}$$

$$\tilde{x}_n \sim \mathcal{N}(0, 1)$$

$$D_{\text{train}} = \{(X_n, C_n) \}_{n=1}^{N_{\text{train}}}$$

Inducción

PREDICTOR

$$\tilde{X}_n = \begin{pmatrix} X_{n1} \\ X_{n2} \\ \vdots \\ X_{nd} \end{pmatrix} = \begin{pmatrix} X_n^{(1)} \\ X_n^{(2)} \\ \vdots \\ X_n^{(d)} \end{pmatrix}$$

$$n = 1, \dots, N_{\text{train}}$$

$$d = 1, \dots, D$$

VECINOS PRÓXIMOS

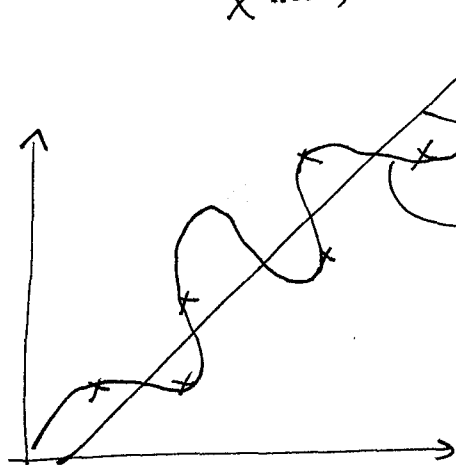
INPUT: D_{train}
distancia

K: #vecinos

$X(\text{test})$

• No funciona bien si la dimensión D es muy grande
↳ selección de variables?

• Si hay ruido no funciona bien, debido a que el ruido "contamina" la distancia.



recta (modelo lineal) Error > 0

polinomio de alto nvl Error = 0

Capacidad de generalización: estimar mediante la tasa de acierto en D_{train}

El polinomio complejo es un ej. de SOBREAJUSTE (OVERFITTING)
La recta suele SUBAJUSTAR (UNDERFITTING)

BOOTSTRAP (remuestreo):

CROSS-VALIDATION (validación-cruzada)

- loo: leave one out
- validación simple
- K-fold cross validation

Para K-NN suele ser

$$K < \sqrt{N_{\text{train}}}$$

Recordar el REESCALADO:

$$X_n^{(d)} \rightarrow \frac{X_n^{(d)} - \text{centro}}{\text{escala}}$$

centro

media

más usada

mediana

más robusta

$$\frac{\max + \min}{2}$$

valores entre 1 y n

escala

varianza $X \sim N(0, \sigma^2)$

iqf

$$\frac{\max - \min}{2} \quad X \sim U(0, 1)$$

ENTROPÍA: ¿Cuán diferente es una distribución de probabilidades de la distribución UNIFORME?

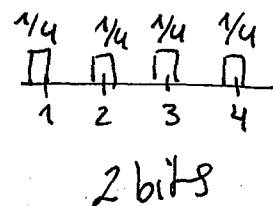
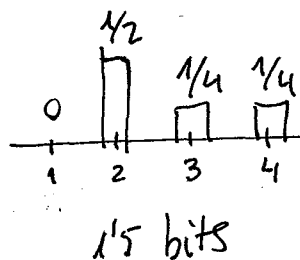
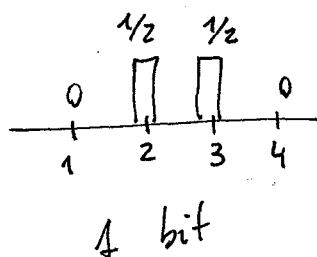
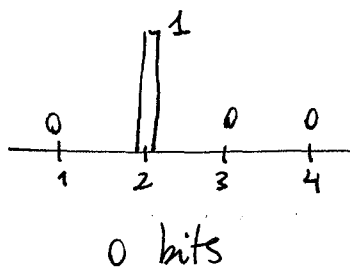
$$\underline{P} = \{P_k\}_{k=1}^K$$

$$H[\underline{P}] = - \sum_{k=1}^K P_k \log P_k$$

base e NATS
base 2 BITS

$$0 \leq P_k \leq 1$$

$$\sum_{k=1}^K P_k = 1$$



ENTROPÍA: $H(\{P_k\}_{k=1}^K)$

$$0 \leq P_k \leq 1; \quad k=1, 2, \dots, K$$

$$\sum_{k=1}^K P_k = 1$$

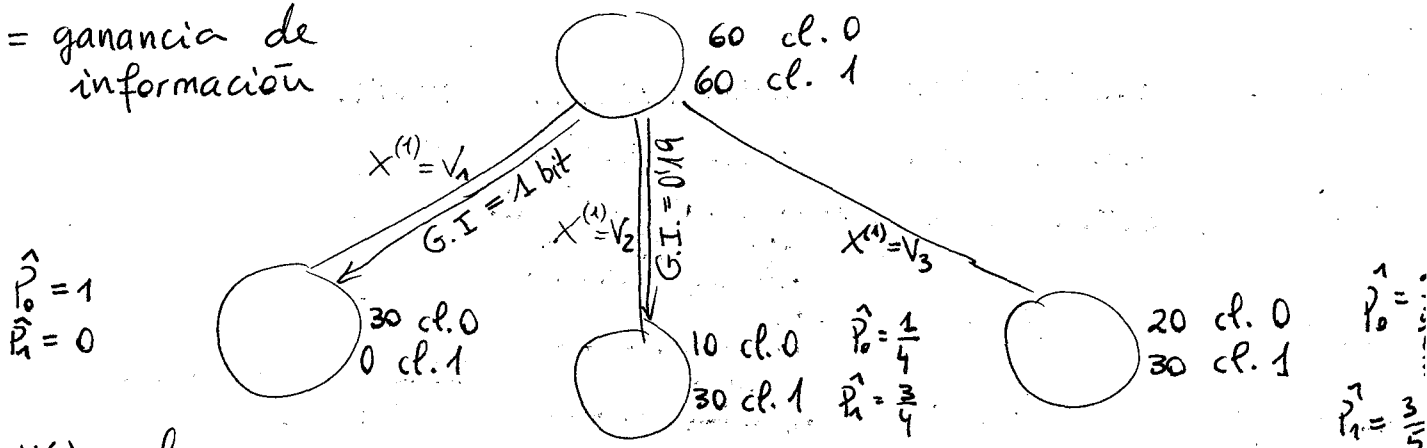
$$H(\{P_k\}_{k=1}^K) = - \sum_{k=1}^K P_k \log_2 P_k \quad (\text{bits}) \quad (\ln \rightarrow \text{NATS})$$

ENTROPÍA BINARIA: $H(P) = -P \log_2 P - (1-P) \log_2 (1-P)$

prob. clase 1 prob. clase 0

Ejemplo: $\mathcal{D}_{\text{train}} = \{(x_n, C_n)\}_{n=1}^{N_{\text{train}}}$ $C_n \in \{0, 1\}$ $N_{\text{train}} = 120$ #clase 0 = 60 #clase 1 = 60

G.I.: ganancia de información



$$H(0) = -0 \log_2 0 - 1 \log_2 1 = 0 \text{ bits}$$

$$H(\frac{1}{4}) = 0.81 \text{ bits}$$

$$H(\frac{2}{5}) = H(\frac{3}{5}) = \dots$$

Arboles de decisión (ID3)

C4.5 o C5.0

CART ← criterio de Gini (otra entropía)

VALIDACIÓN PARA SELECCIÓN DE
 \swarrow ARQUITECTURA
 \searrow MODELO

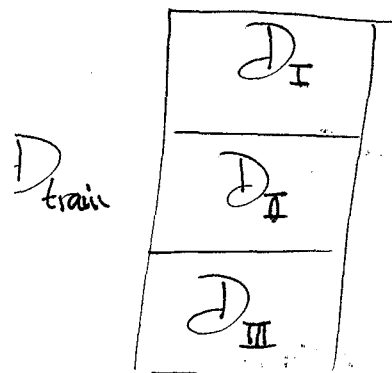
D_{test}

Entrenamiento

$D_I \cup D_{III}$

Validación

Error(D_I)



- VALIDACIÓN SIMPLE: selecciona el modelo que tenga el menor error en el conjunto de validación.

$\underset{MODELO}{\operatorname{argmin}} \operatorname{Error}_{MODELO}(D_I)$

Entrena con $D_I \cup D_{II} \cup D_{III}$ (ed. D_{train})

- VALIDACIÓN CRUZADA DE 3 BLOQUES (3-FOLD CROSS VALIDATION)

$\underset{MODELO}{\operatorname{argmin}} \operatorname{Error}_{MODELO}^{3-CV}$

$$\operatorname{Err}_{MODELO}^{3-CV} = \frac{\operatorname{Err}_{MODELO}(D_I) + \operatorname{Err}_{MODELO}(D_{II}) + \operatorname{Err}_{MODELO}(D_{III})}{3}$$

Entr	Val
$D_I \cup D_{III}$	$\operatorname{Err}(D_{II})$
$D_I \cup D_{II}$	$\operatorname{Err}(D_{III})$
$D_I \cup D_{III}$	$\operatorname{Err}(D_{II})$

MODELO	ARQUITECTURA	HIPERPARÁMETROS	PARÁM.
NN	#capas ocultas #neuronas en capa oculta	• Cte. aprendizaje • Cte. regularización • f. de activación	w
K-NN	—	• Función distancia • k • pesos para votos	—
ÁRBOLES DE DECISIÓN	ID3/C4.5/CART	• Poda • #mínimo ejemplos para división	Divisiones en nodos interiores
SVM	—	• Kernel • Anchura del Kernel	w

APRENDIZAJE POR MÁXIMA VEROSIMILITUD

$$P(C_1 | \underline{x}, \underline{w}) = h(\underline{x})$$

$$\text{Datos: } \{(x_n, t_n) ; n=1, \dots, N_{\text{train}}\} \quad t_n = \begin{cases} 1 & \text{si } C_n = C_1 \\ 0 & \text{si } C_n = C_2 \end{cases}$$

Verosimilitud:

$$\begin{aligned} \mathcal{L}(\underline{w} ; \{(x_n, t_n)\}_{n=1}^{N_{\text{train}}}) &\equiv P\left(\{t_n\}_{n=1}^{N_{\text{train}}} \mid \{x_n\}_{n=1}^{N_{\text{train}}}, \underline{w}\right) = \\ &= \prod_{n=1}^{N_{\text{train}}} \left[P(C_1 | x_n, \underline{w}) \right]^{t_n} \cdot \left[1 - P(C_1 | x_n, \underline{w}) \right]^{(1-t_n)} = \\ &= \prod_{n=1}^{N_{\text{train}}} \left[h(x_n) \right]^{t_n} \cdot \left[1 - h(x_n) \right]^{(1-t_n)} \end{aligned}$$

Función de error

$$E(\underline{w}) = -\log \mathcal{L}(\underline{w} ; \{(x_n, t_n)\}_{n=1}^{N_{\text{train}}}) = \sum_{n=1}^{N_{\text{train}}} \bar{E}_n(\underline{w})$$

$$\text{donde } \bar{E}_n(\underline{w}) = -t_n \log h(x_n) - (1-t_n) \log(1-h(x_n))$$

Aprendizaje ML = minimizar $\bar{E}(\underline{w})$

"embeddings" en espacios de Hilbert con
núcleo reproductor

