

AA - TEMA 1: INTRODUCCIÓN

PARADIGMA

- + Aprender a partir de $\left\{ \begin{array}{l} \text{datos de problemas} \\ \text{respuestas} \end{array} \right\} \rightarrow$ responder preguntas "parecidas".
- + Para ello \rightarrow Obtener parte "general" aplicable siempre (inducción).
- + Dos fases $\left\{ \begin{array}{l} - \text{Entrenamiento: aprender, a partir de datos} \rightarrow \text{responder.} \\ - \text{Test: probar si lo aprendido resuelve casos nuevos.} \end{array} \right.$
- + Dos tipos de problemas $\left\{ \begin{array}{l} - \text{Regresión (predecir valores reales). Eligiendo 1.} \\ - \text{Clasificación (predecir categorías)} \rightarrow \text{Probabilística (\%).} \end{array} \right.$

PRINCIPALES MÉTODOS DE APRENDIZAJE

- + Machine Learning $\left\{ \begin{array}{l} + \text{Objetivo: Algoritmos} \begin{array}{l} \nearrow \text{Precios (\% Error en casos nuevos).} \\ \nearrow \text{Eficientes} \\ \searrow \text{Para problemas a gran escala.} \end{array} \\ - \text{Puede generar muchos modelos.} \end{array} \right.$
- + Aprendizaje Estadístico: aprender infiriendo con distrib. de probabilidad.
 - + Resultados buenos si cumplen hipótesis.
 - Incapaz de resolver problemas a gran escala.
- + Aprendizaje Bayesiano: aproxima estadísticamente con conoc. previo.
 - + Resultados buenos si la hipótesis se cumple.
 - Muy complejo y costoso.

APRENDIZAJE VS DISEÑO ---> información del problema?

- + Diseño: Uso de datos para fijar parámetros de un algoritmo de un problema bien especificado (tenemos mucha info y conocemos funcionamiento).
 - * Ejemplo: Calcular media/desviación de un modelo de prob. definido.
- + Aprendizaje: No conocemos (el) cómo resolver el problema, por lo que metemos los datos en un algoritmo de aprendizaje que nos (ajuste) de una hipótesis para clasificar/predecir las respuestas.

MACHINE LEARNING

+ Def. Un programa aprende de la experiencia E con respecto a una clase de tareas T y medida de desempeño P , si:

$$E \rightarrow P, \text{ también aumenta.}$$

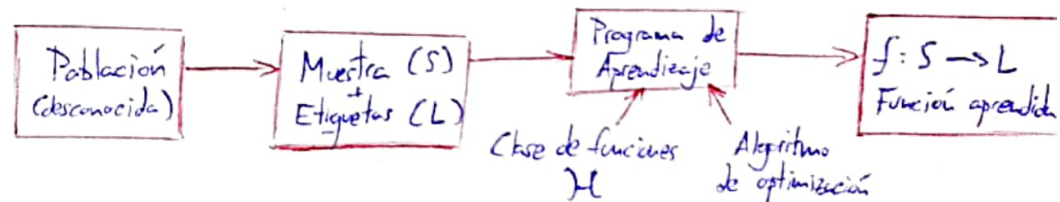
* Cuantos más datos tengamos, mayor experiencia.

PARADIGMAS DE MACHINE LEARNING

- + Aprendizaje supervisado: Aprender una función con $\left\{ \begin{array}{l} - \text{muestra de datos} \\ - \text{etiquetas. (salidas)} \end{array} \right.$
- + Aprendizaje no supervisado: Aprender propiedades de datos sin etiquetar, transformando las características observadas. (\downarrow dimensionalidad, = info).
- * Auto-supervisado: Transforma datos observados en otros (más útiles) resolviendo "tareas de pretexto" (que es supervisado).
- + Conocimiento de transferencia: cómo compartir conocimiento entre tareas.
- + Aprendizaje por refuerzo: Compuesto por datos y recompensas. Se obtiene una política de actuación, que trata de maximizar las recompensas.

APRENDIZAJE SUPERVISADO

- + Asume $\left\{ \begin{array}{l} - \text{Única información disponible} \rightarrow \text{muestra.} \\ - \text{Obtenemos, a través de un Criterio de Aprendizaje} \rightarrow f. \\ - f \text{ calcula una predicción, según su entrada.} \end{array} \right.$
- + Objetivo: elegir $\left\{ \begin{array}{l} \text{muestra} \\ \text{programa} \end{array} \right\} \rightarrow f$ que prediga correctamente toda la población.
- * Problema: ¿qué datos son relevantes para el problema? (p.e. árbol ≤ 10 m)
 - \rightarrow No lo sabemos, pero existe un patrón (ya que nosotros lo vemos).
 - Por tanto, trataremos de aprenderlo a partir de los datos.
 - \rightarrow Si elegimos buena representación, las técnicas funcionarán.



APRENDIZAJE NO SUPERVISADO

- + Objetivo: transformar muestra a un espacio de representación, donde se introduce una nueva medida de distancia entre muestras.
 → Encontrar patrones. ^{características + indep.} → facilitar posterior procesamiento.
- + Tipos
 - Clustering: usa la estructura geométrica.
 - Descubrimiento de dependencias → patrones.
 - Reducción de dimensionalidad → quitar atributos no relevantes.
 - Auto-supervisado.
- + Problema: encuentra características no relevantes → ↓↓ aprendizaje.

FORMALIZACIÓN DEL APRENDIZAJE SUPERVISADO

DATOS INICIALES

- + Obtendremos los datos de $P(D)$
 - P es la probabilidad de extraer cada elemento \rightarrow estática inmutable ($x \in D$)
 - D es la población
- * La muestra debe extraerse i.i.d (independiente, idénticamente distribuida).

+ X : características medidas en cada muestra.

FUNCIÓN OBJETIVO: $f: X \rightarrow Y$
(muestras)

+ f es la función a aproximar, que es desconocida.

MODELO

- + H : clase de funciones que usaremos para aproximar f .
- + Eligiéremos $h \in H$ que sea buena prediciendo la población.
- * Para encontrarla → Algoritmo de Aprendizaje.

ALGORITMO DE APRENDIZAJE

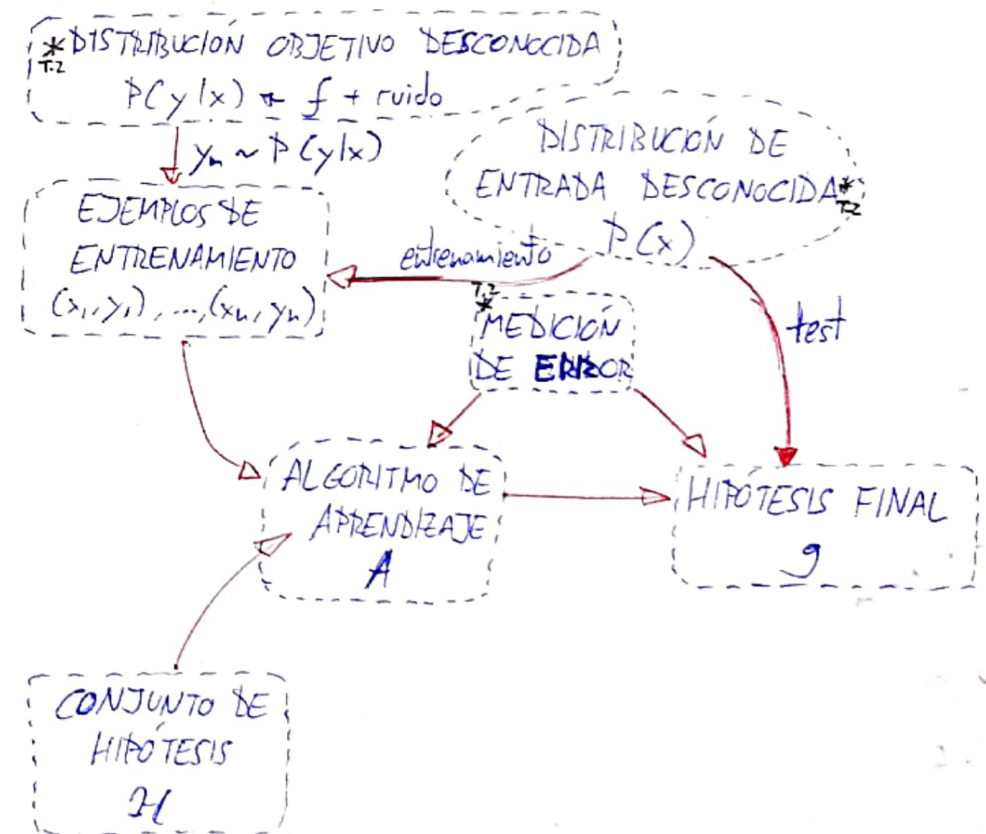
- + ¿Qué criterio optimizar para garantizar el aprendizaje? (función)
 - ERM
 - SRM
 - MBL
 - Similarity-Measure
- + ¿Qué algoritmo usar para optimizar el criterio? A
 - * Gradiente Descendente, SGD, Pseudo-inversa...

PROCESO DE APRENDIZAJE SUPERVISADO

+ Recibimos X, Y y D como punto de partida.

1. Elegimos H con los tipos de funciones que probablemente represente a f . $H =$ "conjunto de hipótesis".
2. Extraemos una muestra de D . Cada elemento de la muestra irá marcado con X e Y .
3. Elegimos, a través de un algoritmo de aprendizaje A y la muestra, una hipótesis $g \in H$.
4. Comprobamos si $g \approx f$.

CT.2) MODELO COMPLETO DE APRENDIZAJE (COMPLETO)



A A - TEMA 2 : MODELOS LINEALES

+ \mathcal{H} = clase de las f lineales de X .

$$h(x) = \sum_{i=1}^d x_i w_i + \frac{\text{bias}}{x_0 \cdot w_0} \parallel h(x) = w^T x$$

($h \in \mathcal{H}$)

$$h, f: X \rightarrow Y$$

$$X = 414 \times \mathbb{R}^d$$

$$Y = \mathbb{R}$$

+ Aplicable a todo, salvo a árboles.

+ Si un modelo lineal resuelve un problema \rightarrow mejor modelo.

* Nos da la influencia de cada característica.

+ Generalmente \rightarrow Estrategia ganadora probar primero lineal.

REGRESIÓN

+ $h(x) = w^T x$ será la predicción de la regresión.

+ Usaremos el error cuadrático para medir la pérdida de $h_w(x)$ con f .

$$* \text{Error}(x) = (h_w(x) - f(x))^2$$

$$* E_{in}(h_w) = \frac{1}{N} \sum_{i=1}^N (h_w(x_i) - y_i)^2, \text{ (error dentro de la muestra).}$$

+ Eligiéramos h_w tal que $\downarrow \downarrow E_{in}$, aunque el objetivo es $\downarrow \downarrow E_{out}$.

+ $E_{out}(h) = \mathbb{E}_{(x,y) \sim p} [(h(x) - y)^2]$. (esperanza / media al elegir un elem. de forma aleatoria de la población).

$$\hat{h} = \min_{h \in \mathcal{H}} E_{out}(h)$$

+ Para encontrar $\hat{h} \rightarrow$ ERM (Minimizar Riesgo Empírico).

$$E_{in}(w) = \text{media}[(w^T x_i - y_i)^2] \parallel E_{in}(w) = \frac{1}{N} \|Xw - y\|^2$$

$$w_{lin} = \min_{w \in \mathbb{R}^{d+1}} E_{in}(w) \quad \text{error cuadrático} \quad \begin{cases} X = \begin{bmatrix} x_1 \\ \vdots \\ x_n \end{bmatrix} ; y = \begin{bmatrix} y_1 \\ \vdots \\ y_n \end{bmatrix} \\ w = \begin{bmatrix} w_0 \\ \vdots \\ w_d \end{bmatrix} \end{cases}$$

+ ¿Cuándo funciona ERM?

\rightarrow Modelo: $y_i = f(x_i, \beta) + \text{ruido}$ $\xrightarrow{\text{contemplamos ruido en } y}$ (cho en X).
Estimador de los pesos w .

* Tª Gauss-Markov: El método de los mínimos cuadrados ordinarios alcanza $\hat{\beta} \rightarrow$ estimador insesgado ($E(\hat{\beta}) = \beta$) de varianza mínima. (óptimo)

* Condiciones: $\begin{cases} - E(\epsilon_i) = 0 \text{ (Media ruido=0)} \\ - E(\epsilon_i \epsilon_j) = 0 \text{ (Ruido incoherente)} \end{cases}$ $\text{Var}(\epsilon_i) = \sigma^2 < \infty$. (Var. acotada)

PSEUDO-INVERSA: ECUACIONES NORMALES DE E_{in}

$$E_{in}(w) = \frac{1}{N} \|Xw - y\|^2 \quad \nabla E_{in}(w) = 0 \text{ en puntos críticos} \quad \begin{cases} - \text{mínimo} \\ - \text{máximo} \\ - \text{p. silla} \end{cases}$$

$$\nabla E_{in}(w) = \frac{2}{N} X^T (Xw - y) = 0$$

$$X^T X w_{lin} = X^T y \rightarrow w_{lin} = \frac{(X^T X)^{-1} X^T y}{X^+ = \text{"pseudo-inversa de } X"}$$

* ¿Cómo hacer $(X^T X)^{-1}$? Descomposición en Valores Singulares.

$$X = U \Sigma V^T \rightarrow U, V^T \text{ son ortogonales } (U^{-1} = U^T)$$

Σ es diagonal ($\Sigma = \Sigma^T$).

$$X^T X = (U \Sigma V^T)^T U \Sigma V^T = V \Sigma^T U^T U \Sigma V^T = V \Sigma \Sigma^T V^T$$

$$(X^T X)^{-1} = (V \Sigma \Sigma^T V^T)^{-1} = V^{-1} \Sigma^{-1} \Sigma^{-1} V^{-1} = V \Sigma^{-2} V^T$$

$$w_{lin} = (X^T X)^{-1} X^T y = V \Sigma^{-2} V^T \cdot V \Sigma U^T y = V \Sigma^{-1} U^T y. \text{ siempre es calculable}$$

* Si $(X^T X)$ tiene rango $d \rightarrow 1$ solución

* Si $(X^T X)$ tiene rango $< d \rightarrow \infty$ soluciones (improbable).

PROPIEDADES DE REGRESIÓN

$$+ \hat{y} = X w_{lin} = X (X^T X)^{-1} X^T y = \hat{H} y$$

* \hat{H} es la "matriz de proyección". $\hat{H}: y \rightarrow \hat{y}$.

* Es relevante en el estudio de E_{in} / E_{out} . (si las respuestas son \hat{y} , crea pred. \hat{y}).

* Es idempotente ($\hat{H}^2 = \hat{H}$).

* Su traza (suma elementos diagonales) = $d+1$.

+ Generalización del error: existen fórmulas para calcular E_{out} :

$$E_{out} = E_{in}(w_{lin}) + O\left(\frac{d}{N}\right) \xrightarrow{\text{Dimensión muestras } (n^2 \text{ características})} N^2 \text{ muestras}$$

+ Dado w_0 , podemos encontrar $\hat{v}: E_{in}(w_0 + \eta \hat{v}) < E_{in}(w_0)$.

* Obtenemos \hat{v} con expansión de Taylor de grado 1. ΔE_{in} .

$$\Delta E_{in} = E_{in}(w_0 + \eta \hat{v}) - E_{in}(w_0) = \eta \nabla E_{in}(w_0)^T \hat{v} + O(\eta^2)$$

$$\eta \nabla E_{in}(w_0)^T \hat{v} = \eta \|\nabla E_{in}(w_0)\| \|\hat{v}\| \cos \theta \geq -\eta \|\nabla E_{in}(w_0)\|$$

Máximo que puede bajar

* Si tomamos $\hat{v} = -\frac{\nabla E_{in}(w_0)}{\|\nabla E_{in}(w_0)\|}$ (vector unitario en dir. contraria a ∇) $\rightarrow \Delta E_{in} = \eta$

GRADIENTE DESCENDENTE

- + Técnica de optimización iterativa general \rightarrow Mínimo local, siguiendo la dirección del vector de gradiente en cada punto.
- + Algoritmo: Comienza con w_0 . De manera iterativa:
 - * $w_j = w_j - \eta \frac{\partial E_{in}(w)}{\partial w_j}$, $\forall j \in [0, n]$, actualizando todos a la vez.
 - \downarrow Componente j del vector gradiente.
 - Tasa de aprendizaje
 - * $\frac{\partial E_{in}(w_0)}{\partial w_j} = \frac{\partial}{\partial w_j} \frac{1}{N} \sum_{n=1}^N (w^T x_n - y_n)^2 = \frac{2}{N} \sum_{n=1}^N (x_{nj} (w^T x_n - y_n))$
 - (~~* $w_{i+1} = w_i - \eta \cdot \left(\frac{2}{N} \sum_{i=1}^N (x^T (w^T x_i - y_i)^T) \right)$~~)
 - * Cada punto (x_n, y_n) contribuye a la actualización proporcionalm. a su error de predicción.
- + ¿Qué valor asignar a η ?
 - * Si es muy pequeño, tarda en converger.
 - * Si es muy grande, puede no converger y fluctuar permanentem.
 - * Ideal: empezar con valores grandes, e ir descendiendo para converger.

GRADIENTE DESCENDENTE ESTOCÁSTICO

- + En cada actualización de w , se utilizan solo un subconjunto de ejemplos (minibatch de $M \ll N$ elementos).
- + Ventajas
 - * $\uparrow \uparrow$ variabilidad en estimación del gradiente.
 - * Muy rápido computacionalmente.
 - * En funciones no convexas \rightarrow Obtiene buen mínimo local.
- + Algoritmo: Comienza con un w_0 . De manera iterativa:
 1. Reordenamos la muestra y la dividimos en minibatches.
 2. Para cada minibatch:
 - for $(j=0 \rightarrow k)$ $t = \text{minibatch}$
 - $w_j = w_j - \eta \sum_{i=0}^t x_{ij} (h(x_i) - y_i)$.

MÉTODO DE NEWTON

- + Realiza la expansión de Taylor de segundo nivel:
$$g(\Delta w) = E_{in}(w_0 + \Delta w) \approx E_{in}(w_0) + \Delta w^T \nabla E_{in}(w_0) + \frac{1}{2} \Delta w^T \nabla^2 E_{in}(w_0) \Delta w$$
$$\frac{\partial g(\Delta w)}{\partial \Delta w} = 0 \rightarrow \nabla E_{in}(w_0) + \frac{\nabla^2 E_{in}(w_0) \Delta w}{H \text{ (hessiana)}} = 0$$
$$\Delta w = -\frac{H^{-1} \nabla E_{in}(w_0)}{H \text{ (hessiana)}}$$
 - \rightarrow Define el avance en dirección del gradiente.
- Es costoso calcular H^{-1} .

ANÁLISIS DE RESULTADOS

- + Análisis de residuos (función $r = y - h$). \rightarrow detectar muestras fuera de la hipótesis de los mínimos cuadrados ordinarios.
- * Gráficos:
 - Errores correlativos \rightarrow Si los residuos muestran una tendencia.
 - \rightarrow Dependencias no controladas
 - \rightarrow Extracción de muestra no cumple i.i.d (p.e. series temporales)
 - Errores con varianza no constante:
 - \rightarrow Solución: transformar etiquetas, o dar pesos a observaciones.
 - Dependencia entre predictores
 - \rightarrow Multicolinealidad: \uparrow característica \Rightarrow combinación lineal.
 - Problema \rightarrow Inversión de matriz inestable, coef. imprecisos.
 - \rightarrow $\uparrow \uparrow$ generalización.
 - Valores atípicos / aislados
- + En dimensiones muy altas \rightarrow Escenario complejo \rightarrow Hipótesis original pierde relevancia.
- * Características categóricas \rightarrow vector "one-hot" (un elemento = 1, resto = 0)
 - \rightarrow Modelos no lineales, como árboles, manejan categorías.
- * Pérdida de valores \rightarrow Faltan en la muestra de entrenamiento.
 - \rightarrow Solución: valores sustitutos \rightarrow interpolados de histograma.
 - \rightarrow predicciones usando algún modelo.
- \rightarrow Los árboles pueden también gestionar esto.

AA - TEMA 2: MODELOS LINEALES (II)

CLASIFICACIÓN

PERCEPTRÓN - DOS CLASES

(aka. Perceptron)

+ \mathcal{H} : hiperplanos separadores, con dimensión = n° características.

* Si $h(x) \geq b \rightarrow x \in \text{Clase 1} \Rightarrow h(x) = \text{signo}((\sum_{i=1}^n x_i w_i) + b)$
 * Si $h(x) < b \rightarrow x \in \text{Clase 2}$
 $h(x) \in \{-1, 1\}$

* O, de forma equivalente: $\mathcal{H} = \{h \mid h(x) = \text{signo}(w^T x)\}$, con $w \in \mathbb{R}^n$.
 $\rightarrow w^T = (b, w_1, \dots, w_n)$
 $\rightarrow x^T = (1, x_1, \dots, x_n)$

+ Criterio de aprendizaje (ERM): $\min_{w \in \mathbb{R}^{d+1}} \frac{1}{N} \sum_{n=1}^N [\text{signo}(w^T x_n) \neq y_n]$.

* Problema: no es derivable.

+ Solución: usar una función parecida y derivable: u fallo

$$\text{error}(w^T x_n, y_n) = \begin{cases} -y_n w^T x_n & \text{si } \text{signo}(w^T x_n) \neq y_n \\ 0 & \text{si } \text{signo}(w^T x_n) = y_n \end{cases}$$

$$\nabla \text{error} = \frac{\partial \max(0, -y_n w^T x_n)}{\partial w} = \begin{cases} -y_n x_n & \text{si fallamos} \\ 0 & \text{si acertamos} \end{cases}$$

+ Algoritmo Perceptrón (PLA).

* Equivale a SGD con $\text{tam_minibatch} = 1$ y $m = 1$.

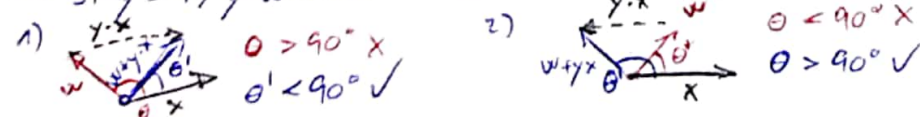
* Regla de actualización $\begin{cases} \text{Si falla: } w_{i+1} = w_i + \gamma \cdot x \\ \text{Si acierta: } w_{i+1} = w_i \end{cases}$ $\left\{ \begin{array}{l} \text{Hasta que clasifique} \\ \text{todo bien} \rightarrow \text{FIN} \end{array} \right.$

* Si el conjunto es linealmente separable \rightarrow solución que clasifica todos los ejemplos bien en t finito. Si no, no converge.

* Interpretación geométrica: $\gamma \cdot x = \|\gamma\| \|\gamma\| \cos(\theta)$

$\hat{1}$ Si $\gamma = 1$, $\gamma \cdot w \cdot x < 0 \rightarrow \theta > 90^\circ \rightarrow (w + \gamma \cdot x) \cdot x > 0$

$\hat{2}$ Si $\gamma = -1$, $\gamma \cdot w \cdot x > 0 \rightarrow \theta < 90^\circ \rightarrow (w + \gamma \cdot x) \cdot x < 0$



MODELO DE CLASIFICACIÓN LINEAL

+ $X = \{1\} \times \mathbb{R}^d$; $Y = \{-1, 1\}$, $f: X \rightarrow Y$

+ $\mathcal{H} = \{h \mid h(x) = \text{signo}(w^T x)\}$

+ De acuerdo a ERM y la teoría de aprendizaje:

$$E_{\text{out}}(g) \leq E_{\text{in}}(g) + O\left(\sqrt{\frac{\log N \cdot d}{N}}\right) \rightarrow \text{Con } N \text{ suficientemente grande, } E_{\text{out}} \approx E_{\text{in}}$$

\rightarrow Reducción de E_{in} en datos linealmente separables.

+ Resultado PLA: dada una muestra lin. sep., PLA encuentra una separación tras $(R/B)^2$ iteraciones (como máximo).

* $R = \max_i \|x_i\|$

* $B = \min_i \|w\|: \forall i \in [n], y_i w^T x_i \geq 1, w \in \mathbb{R}^d$.

Por tanto, la convergencia depende de B (que puede ser exponencialmente grande).

+ Sin embargo, no tenemos por qué llegar tan lejos (la población puede ser no separable, aunque la muestra lo sea).

\rightarrow Reducción de E_{in} en datos lin. no separables.

+ Caso 1: Errores en f lin. sep.

* Si queremos encontrar solución $\rightarrow \downarrow \downarrow E_{\text{in}}$.

$$E_{\text{in}} = \min_{w \in \mathbb{R}^d} \frac{1}{N} \sum_{i=1}^N [\text{signo}(w^T x_i) \neq y_i] \rightarrow \text{Problema NP-Duro.}$$

* Alternativa: PLA Pocket

\rightarrow En cada iteración, calcula E_{in} . Si es más bajo que el mejor E_{in} hasta ahora \rightarrow Guarda solución.

$\rightarrow \downarrow \downarrow$ eficiencia/iteración, pero acota n^2 iteraciones.

+ Caso 2: transformaciones no lineales \rightarrow Hacer que sean sep.

REGRESIÓN LINEAL PARA CLASIFICACIÓN

+ Regresión lineal aprende una $f(x) \in \mathbb{R}$.

+ Podemos calcular w : $w^T x_n \approx y_n = \pm 1$.

+ Así, $\text{sign}(w^T x_n) \rightarrow y_n = \pm 1$.

MUESTRAS CON RUIDO

- + Pueden haber muestras con $y_n \neq f(x_n)$ ^{per falta de info.} f es no determinista. (influye punto vista).
- + En vez de una función, buscaremos distribución $P(y|x)$.

+ Cada muestra $\rightarrow P(x, y) = P(x) P(y|x)$

Info. de tal población \rightarrow % de x y y ca la vez (impureza de x) \rightarrow % x conociendo y (que es la entrada). Es lo que queremos conocer.

$$\rightarrow P(x, y) = \frac{P(y)}{P(x)} P(x|y)$$

% y (etiquetas + probables) % x conociendo su clase y .

- + Objetivo con ruido = Objetivo determinista + ruido
 $f(x) = E(y|x)$ $y - f(x)$
- + Objetivo determinista = Objetivo con ruido: $P(y|x) = 0$, excepto en $y = f(x)$.
- + ¿ $P(y|x)$? $\left\{ \begin{array}{l} \text{- Aprenderlo directamente: si consideramos muestra i.i.d;} \\ \text{discriminativa} \end{array} \right. P(x, y) = P(y|x) P(x) \rightarrow \text{cte.}$
- $\left\{ \begin{array}{l} \text{- Aprenderlo a través de } P(x|y) \text{ con el T}^a \text{ de Bayes;} \\ \text{Generativa} \end{array} \right. P(y|x) = \frac{P(x|y) \cdot P(y)}{P(x)} \rightarrow \text{cte.}$

CLASIFICADORES PROBABILÍSTICOS

REGRESIÓN LOGÍSTICA (LGR)

+ $h(x) = \sigma(w^T x) \in [0, 1]$ // $h(x) = \tanh(w^T x) \in [-1, 1]$

\hookrightarrow Función logística (sigmoide): $\sigma(x) = \frac{e^x}{e^x + 1}$

- + Podemos considerar la salida \approx probabilidad de clasificación.
- + Nos da flexibilidad al asignar muestras a las etiquetas.

+ Caso binario $\left\{ \begin{array}{l} \text{* Probabilidad } C_1: h(x) = \sigma(w^T x) \\ \text{* Probabilidad } C_2: 1 - h(x) = \sigma(-w^T x) \end{array} \right.$

+ Ratio "odds" = $\frac{h(x)}{1-h(x)} \in (0, +\infty) \rightarrow$ Función logit = $\ln(\text{odds}) = \ln(e^{w^T x}) = w^T x \in (-\infty, \infty)$

CRITERIO DE APRENDIZAJE DE MÁXIMA VEROSIMILITUD

- + ~~Máxima~~ Trataremos de maximizar la probabilidad de la muestra:

$$L(S) = \prod_{i=1}^N P_w(y_i | x_i) \left\{ \begin{array}{l} \text{- Si clasificamos bien} \rightarrow \sim 1 \\ \text{- Si clasificamos mal} \rightarrow \sim 0 \end{array} \right.$$

+ $P_w(y|x) = \begin{cases} f(x) & \text{para } y=1 \\ 1-f(x) & \text{para } y=-1 \end{cases} = \begin{cases} h_w(x) & \text{para } y=1 \\ 1-h_w(x) & \text{para } y=-1 \end{cases} = \sigma(y w^T x)$

N_s de el sign

$$h_w(x) = \sigma(w^T x) = 1 - \sigma(-w^T x)$$

+ Por tanto, maximizaremos $L(w) = \prod_{i=1}^N P_w(y_i | x_i) = \prod_{i=1}^N \sigma(y_i w^T x_i)$

- + Equivalentemente, podemos calcular E_{in} y minimizarlo:

$$\begin{aligned} * E_{in}(w) &= -\frac{1}{N} \ln(L(w)) = -\frac{1}{N} \sum_{i=1}^N \ln(\sigma(y_i w^T x_i)) = \\ &= -\frac{1}{N} \sum_{i=1}^N \ln(1 + e^{-y_i w^T x_i}) \end{aligned}$$

② Véase "Entropía Cruzada"

$$* \nabla_w E_{in}(w) = \frac{1}{N} \sum_{i=1}^N -y_i x_i \sigma(-y_i w^T x_i)$$

CLASIFICACIÓN MULTIETIQUETA

- + Podemos generalizar regresión logística $\rightarrow k$ etiquetas.

* En vez de $P(x=1) = p$ y $P(x=0) = q$, tal que $p+q=1$.

* $P(x_1, \dots, x_N) = P(x_1) \cdot \dots \cdot P(x_N)$

* $L(x_1, \dots, x_N) = \prod_{i=1}^N \sigma(w^T x_i)^{[x_i=1]} \sigma(-w^T x_i)^{[x_i=0]}$

Iterative clases

- + Ahora, $y = 1$ -of- k vector (1 elemento es 1, resto 0).

- + La verosimilitud se define ahora usando N clases (asumimos i.i.d).

$$* P(y | w_1, \dots, w_k, x) = \prod_{k=1}^k P(y_k | w_k, x) = \prod_{k=1}^k \sigma(w_k^T x)^{y_k}$$

Estimadores de cada clase Etiqueta a 1 en la muestra Const. muestra

Para cada muestra \rightarrow Prob. en el modelo de que la muestra ancestral exista

+ $* L(y | w_1, \dots, w_k) = \prod_{n=1}^N \prod_{k=1}^k (\sigma(w_k^T x_n))^{y_{nk}} \rightarrow$ Función a maximizar (con k w 's, será k veces más divertido!)

AA - TEMA 2: MODELOS LINEALES (III)

CLASIFICACIÓN MULTITIQUETA (CONTINUACIÓN)

- + SOFTMAX: algoritmo que crea distribución de probabilidad a partir de los $w_k^T x$ resultantes de los distintos $w_k^T x$.

$$P(y_j \neq 1 | x) = \frac{e^{w_j x}}{\sum_{i=1}^K e^{w_i x}}, \text{ para cada clase } j.$$

$$\sum P(y_j | x) = 1.$$

- + Regla de Bayes: Dada una distribución P_x binaria, la mejor función de predicción será:

$$h_p(x) = \begin{cases} 1 & \text{si } P[y=1|x] \geq 0.5 \\ 0 & \text{en otro caso.} \end{cases}$$

- * Por tanto, siempre asignaremos a una muestra x la clase + probable.

COSTES EN EL ERROR

- + No es lo mismo fallar en unos casos u otros.
- + Dos casos:
 - Falso positivo: "Dejar a un intruso entrar en zona máxima seguridad".
 - Falso negativo: "No dejar al general entrar en cuartel".
- * Dependiendo de qué tipo de error nos sea más tolerable, ajustamos:

$$f_{\text{coste}} = c_{-1} \cdot P(\text{error}_{-1}) + c_1 \cdot P(\text{error}_1)$$

MÉTRICAS DE EVALUACIÓN

MATRIZ DE CONFUSIÓN

Muestra	Clase	
	Pos.	Neg.
Predicción	✓ True Positive	✗ False Positive
	✗ False Negative	✓ True Negative

+ Accuracy: % aciertos

$$= \frac{TP + TN}{P + N}$$

+ Sensibilidad = $\frac{TP}{P}$ (o recall)

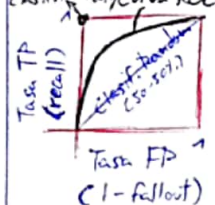
+ Especificidad = $\frac{TN}{N}$ (o fallout)

+ Precisión = $\frac{TP}{TP + FP}$

(cuanto + pérdida \rightarrow + métrica)
 * Nos interesa reducir la métrica, no la pérdida.

CURVA ROC

+ Representa si un clasificador es mejor que otro.
 Clásif. ROC / Curva ROC ("Receiver Operator Characteristic")



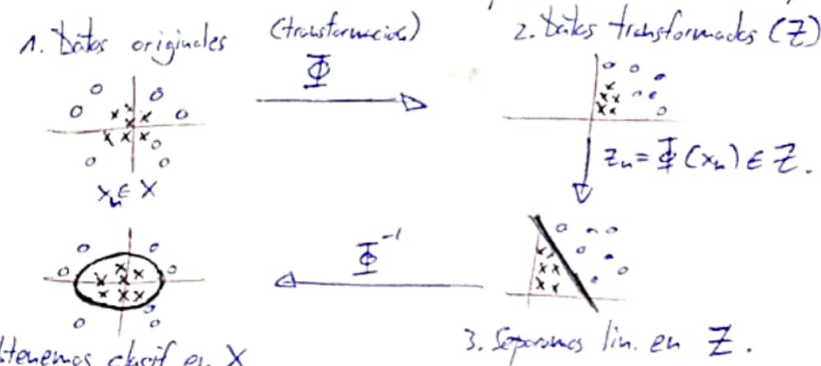
+ Cada punto de la curva equivale a un clasificador con cierto umbral.

+ Conforme ajustamos el umbral, ↑↑ tasa TP, pero también ↑↑ FP.

+ Área bajo la curva (AUC) = % de que el clasificador posicione TP sobre FP. $AUC \approx 0.5 \rightarrow$ Clásif. Random. $AUC \approx 1 \rightarrow$ Perfecto.

TRANSFORMACIONES NO LINEALES

- + Un modelo lineal solo tiene que tener w lineal.
- * resto características \rightarrow cualquier función.
- + Generalmente, añadir características no lineales mejora ajuste.
- * Recomendables cuando el residuo no es lineal.
- * Esto tiene su coste \rightarrow Necesitaremos más muestras. (\uparrow dimensionalidad $\rightarrow \uparrow N$).
- * El uso de transformaciones no lineales no cambia \mathcal{H} , sino que se transforman características a espacio donde son lin. sep.



- 4. Obtenemos clasif. en X .
 $y(x) = \text{signo}(\tilde{w}^T \Phi(x))$
 Pesos en Z Características en Z .
- Problema: debemos fijar la transformación sin ver los datos (por razones teóricas). ¿Qué Φ tomamos?
- + Función Riesgo: medida de la diferencia \approx predicciones (respecto a \mathcal{D}) valores reales.

$$L_p(h) \stackrel{\text{def}}{=} E_{z \sim p}[L(h, z)] \stackrel{\text{def}}{=} E_{\text{out}}$$

! \rightarrow Esperanza (= media)
- * Aplicable a cualquier modelo de aprendizaje.
- + Sea Φ_Q el polinomio de orden Q (y su transformación a Z).
 * $\uparrow \uparrow Q \rightarrow \uparrow \uparrow$ flexibilidad del modelo.
 * $\uparrow \uparrow Q \rightarrow \uparrow \uparrow d = O(Q^2) \rightarrow \uparrow \uparrow$ coste computacional.
 * $\uparrow \uparrow Q \rightarrow \uparrow \uparrow \uparrow N$ (N : muestras) para llegar al mismo nivel de error de generalización ($O(\frac{1}{N} \log(N))$).
 + Por tanto \rightarrow elegir Q que separe bien los ejemplos de \mathcal{D} (↓ E_{in}).
 - sea lo más pequeño posible (↓ E_{out}).

AA - TEMA 2: MODELOS LINEALES (III)

CLASIFICACIÓN MULTITIQUETA (CONTINUACIÓN)

+ SOFTMAX: algoritmo que crea distribución de probabilidad a partir de los $w_k^T x$ resultantes de los distintos $w_k^T x$.

$$P(y_j = 1 | x) = \frac{e^{w_j^T x}}{\sum_{i=1}^K e^{w_i^T x}} \text{ para cada clase } j.$$

$$* \sum P(y_j | x) = 1.$$

+ Regla de Bayes: Dada una distribución P_y binaria, la mejor función de predicción será:

$$h_p(x) = \begin{cases} 1 & \text{si } P[y=1|x] \geq 0.5 \\ 0 & \text{en otro caso.} \end{cases}$$

* Por tanto, siempre asignaremos a una muestra x la clase + probable.

COSTES EN EL ERROR

+ No es lo mismo fallar en unos casos u otros.

+ Dos casos:

- Falso positivo: "Dejar a un intruso entrar en zona máxima seguridad."
- Falso negativo: "No dejar al general entrar en cuartel."

* Dependiendo de qué tipo de error nos sea más tolerable, ajustamos:

$$f_{\text{COSTE}} = c_{-1} \cdot P(\text{error}_{-1}) + c_1 \cdot P(\text{error}_1)$$

MÉTRICAS DE EVALUACIÓN

MATRIZ DE CONFUSIÓN

		Muestra	
		Pos.	Neg.
Real	✓	True Positive	False Positive
	✗	False Negative	True Negative
		$P = \frac{TP+FN}{N}$	$N = \frac{TN+FP}{N}$

+ Accuracy: % aciertos

$$= \frac{TP+TN}{P+N}$$

+ Sensibilidad = $\frac{TP}{P}$
(o recall)

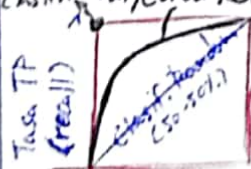
+ Especificidad = $\frac{TN}{N}$
(o fallout)

+ Precisión = $\frac{TP}{TP+FP}$

(cuanto ↓ pérdida → ↓ métrica)
* Nos interesa reducir la métrica, no la pérdida.

CURVA ROC

+ Representa si un clasificador es mejor que otro.
Chisf. Perfecto / Curva ROC ("Receiver Operator Characteristic").



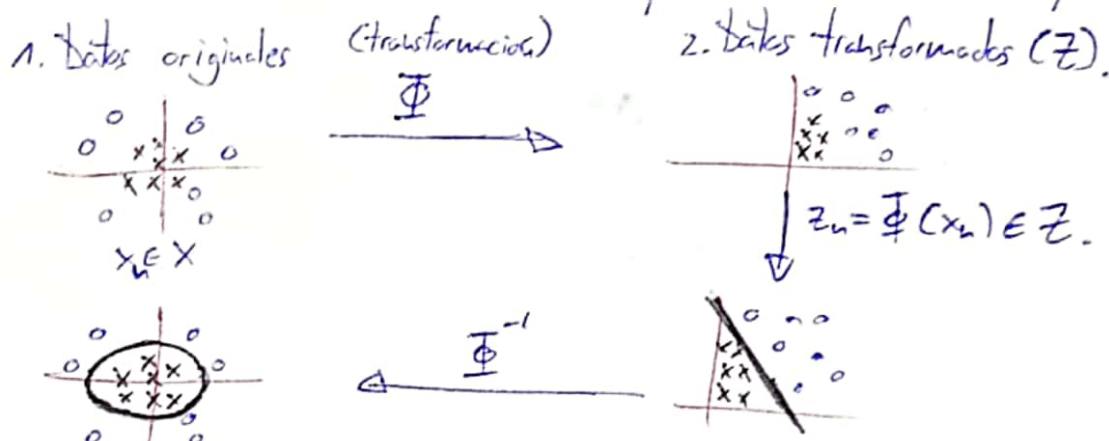
+ Cada punto de la curva equivale a un clasificador con cierto umbral.

+ Conforme ajustamos el umbral, ↑↑ tasa TP, pero también ↑↑ FP.

+ Área bajo la curva (AUC) = % de que el clasificador posicione TP sobre FP. (e[0,1]).
AUC ≈ 0.5 → Chisf. Random. AUC ≈ 1 → Perfecto.

TRANSFORMACIONES NO LINEALES

- + Un modelo lineal solo tiene que tener w lineal.
- * Resto características \rightarrow cualquier función.
- + Generalmente, añadir características no lineales mejora ajuste.
- * Recomendables cuando el residuo no es lineal.
- * Esto tiene su coste \rightarrow Necesitaremos más muestras. (\uparrow dimensionalidad $\rightarrow \uparrow N$).
- * El uso de transformaciones no lineales no cambian H , sino que se transforman características a espacio donde son lin. sep.



4. Obtenemos clasif. en X .

$$y(x) = \text{signo}(\tilde{w}^T \Phi(x))$$

Pesos en Z \leftarrow características en Z .

3. Separamos lin. en Z .

- Problema: debemos fijar la transformación sin ver los datos (por razones teóricas). ¿Qué Φ tomamos?

+ Función Riesgo: medida de la diferencia $\begin{cases} \text{predicciones} \\ \text{valores reales} \end{cases}$ respecto a D

$$L_p(h) \stackrel{\text{def}}{=} \mathbb{E}_{z \sim p} [L(h, z)] \stackrel{\text{def}}{=} E_{\text{out}}$$

\rightarrow Esperanza (= media)

* Aplicable a cualquier modelo de aprendizaje.

+ Sea Φ_Q el polinomio de orden Q (y su transformación a Z).

* $\uparrow Q \rightarrow \uparrow$ flexibilidad del modelo.

* $\uparrow Q \rightarrow \uparrow d = O(Q^2) \rightarrow \uparrow$ coste computacional.

* $\uparrow Q \rightarrow \uparrow \uparrow N$ (n = muestras) para llegar al mismo nivel de error de generalización ($O(\frac{d}{N} \log(N))$).

+ Por tanto \rightarrow elegir Q que separe bien los ejemplos de train ($\downarrow E_{\text{in}}$).
- sea lo más pequeño posible ($\downarrow E_{\text{out}}$).