

AA - TEMA 4: TEORÍA DE LA GENERALIZACIÓN

TRUQUE DE LA DISCRETIZACIÓN

- + Si $|H| = \infty$, la desigualdad de Hoeffding no nos sirve.
- + Sin embargo, H al final ajusta valores discretos (con precisión de 64 bits, por ejemplo) $\rightarrow |H| = 2^{64 \cdot d}$ N° parámetros
- $\rightarrow m_H(\epsilon, \delta) \leq \left\lceil \frac{2}{\epsilon^2} \log \frac{2|H|}{\delta} \right\rceil = \frac{128d + 2 \log(2)}{\epsilon^2}$
- + Es una cota horrible de la complejidad de la muestra.

TRAMPA DE LA DESIGUALDAD UNIFORME

- + Para llegar a la desigualdad de Hoeffding para $|H| > 1$, consideramos: que $P\left(\bigcup_{i=1:|H|} B_i\right) \leq \sum_{i=1} P(B_i)$.
- * Esto es, que no nos preocupamos por el valor real, sino que tomamos que todas las probabilidades son disjuntas ($B_i \cap B_j = \emptyset$), lo cual no es verdad.
- + Por tanto, $P(B_i: |E_{in}(g) - E_{out}(g)| > \epsilon) < 2|H|e^{-2N\epsilon^2}$ es muy mejorable.

FUNCIÓN DE CRECIMIENTO * Clasificación binaria

- + Medida del n.º de "funciones efectivas" en una clase, dadas N puntos.
- * "Función efectiva": función que crea una clasificación concreta. Si dos f crean la misma clasificación con N puntos \rightarrow Son 1 f. efectiva.
- + Función de crecimiento ($m_H(N)$) $\Rightarrow \max_{x_1, \dots, x_N} |H(x_1, \dots, x_N)| \leq 2^h$.
- * Si $m_H(N) = 2^h \rightarrow H$ "shatters" the set. $\hookrightarrow N^{\circ}$ elementos del "conjunto de las dicotomías generables por H ."
- * Es independiente de $P \rightarrow$ Peor caso.
- + A través de $m_H(N)$, podemos deducir (no trivialmente).
- * $P(B: |E_{in}(h) - E_{out}(h)| > \epsilon) \leq 4 m_H(2N) e^{-\frac{1}{2} N \epsilon^2}$.
- * Despejando ϵ : $E_{out}(h) \leq E_{in}(h) + \sqrt{\frac{8}{N} \log \frac{4 m_H(2N)}{\delta}}$
- * Problema: tenemos que calcular $m_H(N) \rightarrow$ Sabemos que $\leq 2^N$.

PUNTO DE RUPTURA

- + Si para algún valor k $m_H(k) < 2^k \rightarrow k$ es pto. ruptura de H .
- * H no puede separar una muestra de tamaño k .
- + Resultado Sauer - Shelah - Perles: si k es punto de ruptura de H : $\forall N, m_H(N) \leq \sum_{i=0}^{k-1} \binom{N}{i}$, que es un polinomio de grado $k-1$.
- * Si H no tiene break point $\rightarrow m_H(N) = 2^N$.
- + Sustituyendo:
- * Si H tiene break point $\rightarrow \delta = 4 \cdot O(N^k) \cdot e^{-\frac{1}{8} N \epsilon^2} \xrightarrow{N \rightarrow \infty} 0$.
- * Si H no tiene break point $\rightarrow \delta = 4 \cdot 2^N \cdot e^{-\frac{1}{8} N \epsilon^2} \not\rightarrow 0$.
- \hookrightarrow Aumentando N , no nos aseguramos de disminuir E_{out} .

DIMENSIÓN VC - VAPNIK & CHERNOENKIS.

- + La dimensión VC de H ($d_{VC}(H)$) es el mayor valor de N tal que N no es un punto de ruptura.
- * Si H no tiene b.p $\rightarrow d_{VC}(H) = \infty$.
- + Puede probarse que $m_H(N) \leq \sum_{i=0}^{d_{VC}} \binom{N}{i} \leq \begin{cases} N^{d_{VC}+1} \\ \left(\frac{eN}{d_{VC}}\right)^{d_{VC}} \end{cases}$
- + La dimensión VC denota el n.º de parámetros "efectivos" de H .
- * P.e: $d_{VC} = d+1$ para modelos lineales de " d " dimensiones.
- + Sustituyendo:
- $E_{out}(h) \leq E_{in}(h) + \sqrt{\frac{8}{N} \log \frac{4 (2N)^{d_{VC}+1}}{\delta}} = E_{in}(h) + O\left(\sqrt{d_{VC} \frac{\log N}{N}}\right)$
- * Si d_{VC} es finito y N muy grande \rightarrow Garantizamos generalización
- * Si $d_{VC} = \infty$, H no sirve para aplicar la regla ERM.

COMPLEJIDAD DE LA MUESTRA

- + ¿N → conseguir ϵ (dif. entre dentro y fuera de la muestra). δ (prob. de que se cumpla la premisa, dado ϵ).

$$\sqrt{\frac{\delta}{N} \ln \frac{4 \ln(2N)}{\delta}} \leq \epsilon \rightarrow N \geq \frac{\delta}{\epsilon^2} \ln \left(\frac{4 \ln(2N)}{\delta} \right) \rightarrow$$

$$\rightarrow N \geq \frac{\delta}{\epsilon^2} \ln \left(\frac{4 (2N)^{dvc} + 1}{\delta} \right)$$

Ecua. implícita → resolver iterativamente.

* Ejemplo: Percepción ($dvc = 3$). (Int.)

$$N \geq \frac{\delta}{\epsilon^2} \ln \left(\frac{4 (2N)^3 + 1}{\delta} \right) \xrightarrow{N=1000} N \geq 21.100 \rightarrow N \geq 30.000 \text{ aprox.}$$

$\epsilon = \text{bil. Error}$ $\delta = \text{Prob. de encontrar "ejemplo malo"}$ punto fijo de N.

- + Según esta fórmula, las cotas crecen proporcionalmente a dvc .
- + Regla práctica: $N > 10 \cdot dvc$ (¡Qué poco!)

VC COMO RESTRICCIÓN DE LA COMPLEJIDAD DEL MODELO.

- + Normalmente → N es fijo. ¿Qué podemos esperar alcanzar con N?

$$* E_{out}(g) \leq E_{in}(g) + \sqrt{\frac{\delta}{N} \ln \frac{4 \ln(2N)}{\delta}} \leq E_{in}(g) + \sqrt{\frac{\delta}{N} \ln \frac{4 (2N)^{dvc} + 1}{\delta}}$$

- + Penalización por complejidad del modelo:

$$* \Omega(N, H, \delta) = \sqrt{\frac{\delta}{N} \ln \frac{4 (2N)^{dvc} + 1}{\delta}} = O \left(\sqrt{\frac{dvc \cdot \ln N - \ln \delta}{N}} \right)$$

- + Intercambio $\left\{ \begin{array}{l} - \uparrow \uparrow dvc \rightarrow \uparrow \uparrow \text{ probabilidad de que } f(E_{in} \approx 0). \\ - \downarrow \downarrow dvc \rightarrow \uparrow \uparrow \text{ probabilidad de generalizar mejor } (E_{in} \approx E_{out}). \end{array} \right.$

$$E_{out} \leq E_{in} + \Omega(dvc).$$

- * Análisis VC solo depende de H . → independiente de $\left\{ \begin{array}{l} - f \\ - A \end{array} \right.$
- * Principalmente aplicable a clasificación y regresión.
- * Cota superior muy amplia.

CONJUNTOS DE TEST

- + Idea: usar parte de la muestra para estimar E_{out} → Test set.
- * Estos ejemplos no se usan para entrenar.
- * Deben ser i.i.d de la misma $P(X)$ del conjunto de entrenam.
- + Error: $E_{test} \approx E_{out}$
- + ¿Por qué funciona? → Desigualdad de Hoeffding.
- * $P(|E_{test}(g) - E_{out}(g)| > \epsilon) \leq 2e^{-2N\epsilon^2}$
- Lo Para $N=1000$: $E_{test} = \pm 5\% E_{out}$ con probabilidad $\geq 98\%$.
- + Contrar: no usamos los ejemplos para entrenar → $\uparrow \uparrow E_{in}$.

DISCUSIÓN NLT

- + ¿Influyen las transformaciones no lineales en la cota VC?
- * Si fijamos H a ciegas → $dvc(H) = dvc(H)$ con probabilidad $1 - \delta$.
- ¿Qué pasa si probamos una H , fallamos, y probamos otra H ?
- * Equivale a usar Φ $\left\{ \begin{array}{l} - \text{monteizando características rectas} \\ - \text{añadiendo términos al cuadrado.} \end{array} \right.$
- * $\uparrow \uparrow dvc$.
- ¿Qué pasa si miramos los datos sin fijar antes una hipótesis?
- * Debemos incluir las hipótesis que hemos pensado → dimensión de transformación.
- * Hemos decidido que el problema es la muestra, y no la población.
- + Si necesitamos ajustar bien → NLT $\left\{ \begin{array}{l} - \uparrow \uparrow dvc \rightarrow \uparrow N \text{ para PAC.} \\ - \uparrow \uparrow X \text{ (espacio de características).} \end{array} \right.$

APRENDIZAJE NO UNIFORME

- + ¿Qué pasa si $dvc = +\infty$? → No nos sirve la regla ERM.
- + Consideramos $H = \cup_n H_n$, $dvc(H_n) < \infty \forall n$.
- * H es la unión de infinitas clases con dvc finita.
- + SRM (Structural Risk Minimization) $\left\{ \begin{array}{l} 1. \text{ Seleccionar secuencia de } H\text{'s anidadas.} \\ 2. \text{ Estima } g \text{ para cada set } H. \end{array} \right.$
- * Los tipos $\left\{ \begin{array}{l} - \text{Fija complejidad} \rightarrow \downarrow \text{ error empírico.} \\ - g^* = \arg \min (E_{in}(g) + \Omega(H)) \\ - \text{Fija error emp.} \rightarrow \downarrow VC - \text{dimensión (si } \uparrow \uparrow \text{ dim} \rightarrow \uparrow \uparrow \text{ varibilidad).} \end{array} \right.$
- * Útil cuando $\frac{N}{dvc} < 20 \rightarrow$ ERM no tiene garantías.

AA-TEMA 4: TEORÍA DE LA GENERALIZACIÓN (II)

INTERCAMBIO BIAS - VARIANZA

- + $E_{out}(g^{(D)}) = E_x[(g^{(D)}(x) - f(x))^2] \rightarrow$ Error cuadrático medio de $g^{(D)}$
Esperanza (media) en población $P(x)$.
- * Toma en cuenta H y A .
Hipótesis g obtenida con A a partir de muestra D .
- + $g^{(D)}(x)$ será la clasificación que la hipótesis g , creada a partir de D , dará a x . Esta varía según D de forma aleatoria.
- + $\bar{g}(x) = E_D[g^{(D)}(x)] \approx \frac{1}{K}(g^1(x) + \dots + g^K(x))$
* Es la predicción media de x .
- + $var(x) = E_D[(g^{(D)}(x) - \bar{g}(x))^2] = E_D[g^{(D)}(x)^2] - \bar{g}(x)^2$
* Es lo que varía, de media, el error cuadrático de una predicción con su valor medio. Es decir, la varianza.
- + $E_{out}^D(x) = (g^{(D)}(x) - f(x))^2$ * Error cuadrático de g^D .
- + $E_{out}(x) = E_D[E_{out}^D(x)]$ * Media de los errores de g^D respecto a f .

DESCOMPOSICIÓN BIAS - VARIANZA

$$E_{out} = E_D[E_{out}(g^D)] = E_D[E_x[(g^D(x) - f(x))^2]] \rightarrow$$

En este caso, podemos intercambiar orden.

$$\rightarrow E_D[(g^D(x) - f(x))^2] = E_D[g^D(x)^2] - 2E_D[g^D(x) \cdot f(x)] + f(x)^2 \rightarrow$$

$$* \text{Sustituimos } E_D(g^D(x)) = \bar{g}(x):$$

$$\rightarrow E_D[g^D(x)^2] - 2E_D[\bar{g}(x)] \cdot f(x) + f(x)^2 =$$

$$= \underbrace{E_D[g^D(x)^2] - \bar{g}(x)^2}_{var(x)} + \underbrace{\bar{g}(x)^2 - 2\bar{g}(x)f(x) + f(x)^2}_{bias(x)} \quad // \text{Acadimos y restamos } \bar{g}(x)^2$$

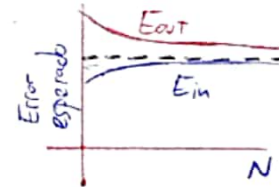
$$E_{out} = E_x[var(x) + bias(x)] = bias + variance$$

CONSECUENCIAS

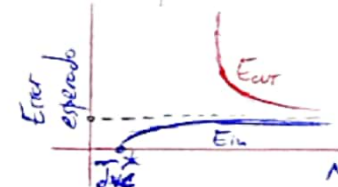
- + $E_D[E_{out}(g^D(x))] = \sigma^2 + bias + varianza$ (con señal ruidosa).
- * σ^2 es la varianza del ruido.
- * σ^2 es inevitable \rightarrow nos centramos en el resto.
- + Hay dos posibilidades:
 - \rightarrow $\downarrow \downarrow \downarrow$ varianza con + bias: Regularización.
 - \rightarrow $\downarrow \downarrow \downarrow$ bias con + varianza: Conjuncamiento previo para minimizar bias, tenemos que "conocer" f .
- + El poder del aprendizaje está vinculado a N (tamaño de la muestra).

CURVA DE APRENDIZAJE

- + Resume el comportamiento de E_{in} y E_{out} (en nuestro caso, el error medio de todas las muestras D posibles) cuando variamos N .



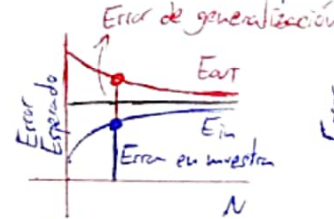
Función simple
(polinomio grado 2)



Modelo complejo
(polinomio grado 10)

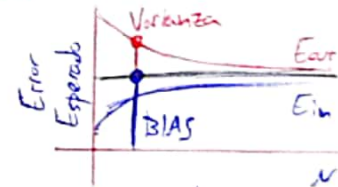
* Si $N \leq d_{vc}$, entonces la función es capaz de ajustarse totalmente $E_{in} = 0$

- + Según el enfoque, podemos ~~ver~~ interpretarlos de diferente manera:



Análisis VC

- + H_2 que generalice
acertar con datos



Análisis bias-varianza

- + $H_2, A \rightarrow$ Aproxime f
No tenga comportamiento con alta varianza.

- + En regresión lineal $\rightarrow E_{in}$ aumenta cuando el modelo absorbe info con d_{vc} para E_{out} ~~disminuye~~ hasta alcanzar el ruido residual.

