

TEMA 5: SOBREAJUSTE Y REGULARIZACION.

SOBREAJUSTE

- + Nos sobreajustamos a una muestra $\leftrightarrow \downarrow \downarrow E_{in} \rightarrow \uparrow \uparrow E_{out}$.
- * ¿Por qué sucede? ¿Tenemos suficientes ejemplos?

→ Según dimension-VC: $E_{out} = E_{in} + \Omega(N, \mathcal{H}, S)$
 ¿Hemos elegido \mathcal{H} muy compleja?

→ Según sesgo-varianza:

$$E_{\sigma}[E_{out}(g^{(D)})] = \sigma^2 + \text{bias} + \text{varianza}$$

* σ^2 = "Error estocástico": error presente en datos entrenamiento.

* Bias = "Error determinista": error que comete la mejor función de nuestro \mathcal{H} "g" estimando f (objetivo).

→ Para mejorarlo → Aumentar \mathcal{H} → \uparrow varianza (intercambio).

* Si tenemos una $\mathcal{H} > D$ fijas \Rightarrow ruido determinista \hookrightarrow indistinguibles.

* Por ejemplo, consideremos 2 casos; con $N=15$:

1. f = Polinomio grado 10 + ruido estocástico (σ^2).

→ Si intentamos ajustar polinomios:

* Con grado 2 $\Rightarrow \uparrow E_{in}$ (debido a bias), pero $\downarrow E_{out}$ (\downarrow varianza).

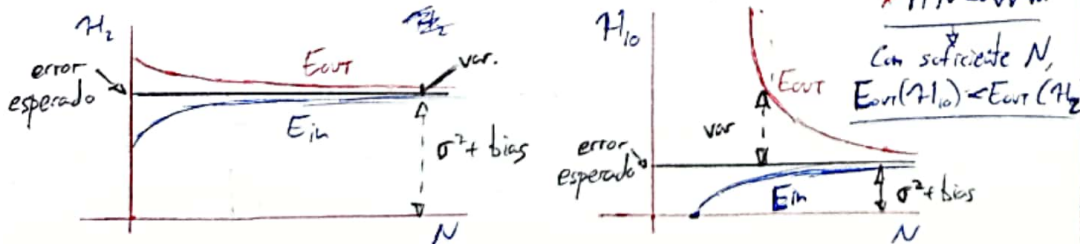
* Con grado 10 $\Rightarrow \downarrow E_{in} (\approx 0)$, pero $\uparrow \uparrow E_{out}$ ($\uparrow \uparrow$ sobreajuste, debido al error estocástico, "aprendimos el error").

2. f = Polinomio de grado 50. (sin ruido).

* Con grado 2 $\Rightarrow \uparrow E_{in}$ (bias), $\downarrow E_{out}$ (varianza baja).

* Con grado 10 $\Rightarrow \downarrow \downarrow E_{in}$, pero $\uparrow \uparrow E_{out}$ ($\uparrow \uparrow$ sobreajuste, ya que \uparrow variabilidad $>$ hay error determinista).

* En ambos casos, estas son las curvas de aprendizaje:



PROTECCIONES CONTRA SOBREAJUSTE

+ ¿Cómo decidimos la complejidad de \mathcal{H} ?

* Hay ruido en la muestra.

* La selección de $h \in \mathcal{H}$ se realiza con ERM o SRM.

+ Solución 1: sesgo inductivo.

* Limitamos \mathcal{H} para que no aprenda el ruido.

* Problema: Restringimos también la capacidad de \mathcal{H} para estimar f .

+ Solución 2: regularización.

* Añadimos a la función de error otra componente \approx complejidad.

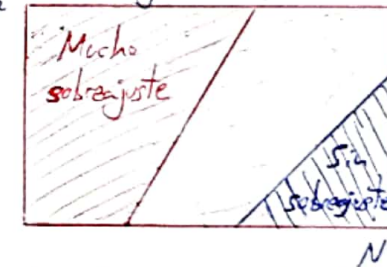
* Así \rightarrow Compromiso entre ajuste y complejidad/variabilidad.

INFLUENCIA DEL RUIDO EN SOBREAJUSTE

+ Ruido estocástico (σ^2)

\mathcal{H}_0 : Sobreajuste (comparando \mathcal{H}_0 y \mathcal{H}_2).

* $\uparrow \uparrow \sigma^2$
 \downarrow
 $\uparrow \uparrow$ sobreajuste
 (aprendemos ruido).

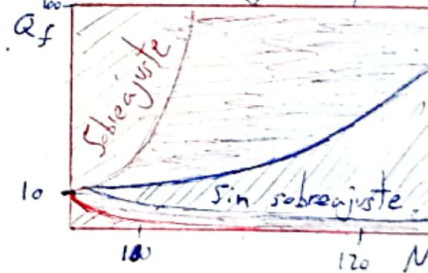


* $\uparrow \uparrow N \Rightarrow \downarrow$ sobreajuste
 (si hay suficientes ej., entonces el ruido se contrarresta, pues tiene media = 0 \rightarrow aprendemos).

+ Ruido determinista (Q_f = complejidad de f = grado de f).

\mathcal{H}_0 : Sobreajuste respecto \mathcal{H}_2

* $\uparrow \uparrow Q_f \Rightarrow \uparrow$ sobreajuste.



* $\uparrow \uparrow N \Rightarrow \downarrow$ sobreajuste.

* Problema: \mathcal{H}_0 , cuando $Q_f > 10$, no es capaz de ajustarse a f . Además, como N es finito $\rightarrow \mathcal{H}_0$ trata de aprender la muestra (aunque no sea capaz de adaptarse al f real) \rightarrow Sobreajuste.
 (con ruido, tanto estocástico como determinista). (aprende patrones inexistentes)

REGULARIZACIÓN

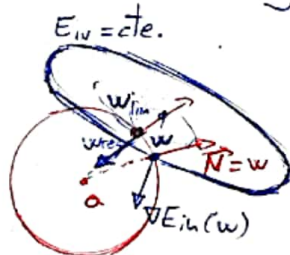
Por ello \rightarrow Heurística.

- + Consiste en restringir $H \rightarrow \downarrow \downarrow E_{out}$.
- + Tiene cierta justificación teórica:
 - * VC $\rightarrow \downarrow \downarrow dvc$, por lo que hay mejor generalización.
 - * Sesgo-Varianza \rightarrow Aumenta un poco el sesgo (limitamos H , por lo que empeora la mejor solución probablemente), pero reduce mucho la varianza.
- + También tiene una intuición práctica: evitar la influencia de ruidos:
 - * El ruido determinista suele ser irregular.
 - * El ruido estocástico tiene alta frecuencia (muchos datos están afectados).
- \rightarrow Por tanto, limitamos H para que sea simple / regular.

WEIGHT DECAY - DECAIMIENTO DE PESOS

- + Consiste en medir la complejidad de una hipótesis h por el tamaño de sus pesos.
- + Por tanto \rightarrow restringimos $H: \sum_{i=0}^n w_i^2 \leq C \rightarrow cte. (\downarrow \downarrow dvc)$.
- * Nuevo problema: $\min_w E_{in}(w)$, tal que $w^T w \leq C. \rightarrow w_{reg}$.
- * ¿Cómo calcularlo? (Aprox. gráfica)

- \rightarrow Círculo rojo = espacio de w 's que satisface $w^T w \leq C$.
- \rightarrow Elipse = puntos en el espacio con el mismo E_{in} que w .



- $\rightarrow w_{lin}$ es la solución óptima (sin restricciones).

Caso 1: si w_{lin} cumple la restricción $\rightarrow w_{reg} = w_{lin}$.

Caso 2: si w_{lin} no cumple la restricción:

- $\rightarrow w_{reg} \in$ Círculo rojo // $w_{reg}^T w_{reg} = C$. (para acercarnos lo más posible a w_{lin}).
- + Tiene que cumplirse que $-\nabla E_{in}(w_{reg})$, que es la dirección hacia donde se minimiza el error, coincide con la normal al círculo ($C = w^T w$, ya que centro $= 0$). Si no $\rightarrow -\nabla E_{in}(w)$ tiene componente en dir. tangente $\rightarrow \downarrow \downarrow E_{in}$ al mov. en 0 .

ERROR AUMENTADO

- + Dado w , definimos su error aumentado:

$$E_{aug}(w, \lambda, \Omega) = E_{in}(w) + \frac{\lambda}{N} \Omega(w)$$

\hookrightarrow La Función de reg. usada \propto Ej. weight decay: $\Omega(w) = w^T w$.
 \hookrightarrow Fuerza de la restricción de la regularización ($\uparrow \lambda \rightarrow \downarrow \downarrow dvc$).

- + Esta surge de que:

$$\min_w E_{in}(w), \text{ tal que } w^T w \leq C \xrightarrow[\text{weight decay}]{\text{equivalente}} \min_w E_{in}(w) + \lambda_c w^T w, \lambda_c > 0 \rightarrow$$

usando multiplicadores de Lagrange $\rightarrow \min_w \{E_{in}(w) + \lambda(w^T w - C)\} \rightarrow E_{aug} = E_{in}(w) + \lambda w^T w$ (sin restricción).

- + Por tanto, minimizador E_{aug} resuelve el problema con regularización.

- + E_{aug} puede interpretarse como una aproximación a E_{out} más cercana que E_{in} :

$$E_{out} \approx E_{in} + \frac{\Omega(H)}{N}; \quad E_{aug} = E_{in} + \frac{\Omega(h)}{N} \cdot \frac{\lambda}{N}$$

Estiman generalización del error

REGRESIÓN REGULARIZADA

- + Objetivo: minimizar $E_{in}(w)$ sobre la hipótesis H restringida.

- * Nota: emplearemos $\Omega(w) = w^T w$.

$$w_{reg} = \min_w E_{aug}(w) = \min_w \frac{\|Zw - y\|^2}{E_{in}(w)} + \lambda \frac{\|w\|^2}{N}$$

$$\nabla_w E_{aug}(w) = \nabla_w (E_{in}(w) + \lambda w^T w) = 0 \rightarrow w = w_{reg} \text{ (mínimo)}.$$

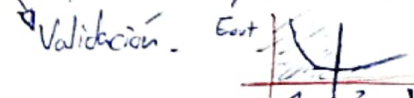
$$\rightarrow \frac{2Z^T(Zw - y)}{-\nabla E_{in}(w)} + \frac{\lambda w^T}{N} = 0 \quad (\text{coincide con intuición geométrica!})$$

$$\rightarrow w_{reg} = (Z^T Z + \lambda I)^{-1} Z^T y$$

$$\begin{aligned} \rightarrow \text{Si } \lambda = 0 &\rightarrow w_{reg} = w_{lin} \\ \rightarrow \text{Si } \lambda \rightarrow \infty &\rightarrow w_{reg} \rightarrow 0. \end{aligned}$$

- * ¿Cómo calculamos λ ?

- \rightarrow Empíricamente (ensayo y error).
 - Si es muy bajo \rightarrow Sobreajuste. (over)
 - Si es muy alto \rightarrow No se ajusta! (bias)

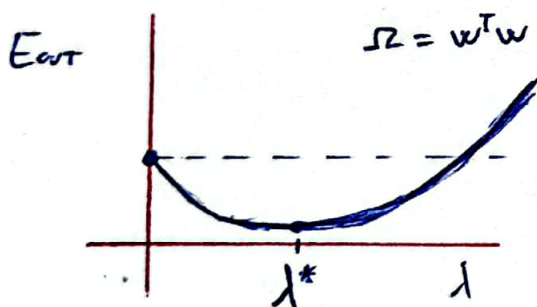


ELECCIÓN DE REGULARIZADOR

- + Regulador óptimo \rightarrow Aquel que nos lleve hacia f_i (desconocida).
- * Por tanto, nos podemos guiar por ir a hipótesis "suaves" y "simples", ya que evitan el ruido (que no es suave normal.).
- + Ejemplos de Ω
 - Weight decay
 - $\sum_{g=0}^G |w_g| \leq C \rightarrow$ Útil para hacer pesos = 0.
 - $\sum_{g=0}^G \gamma_g w_g^2 \leq C \rightarrow$ Da más importancia a ciertos g .
- + Si elegimos Ω inadecuado \rightarrow En validación, $\lambda = 0$ (sin reg.).

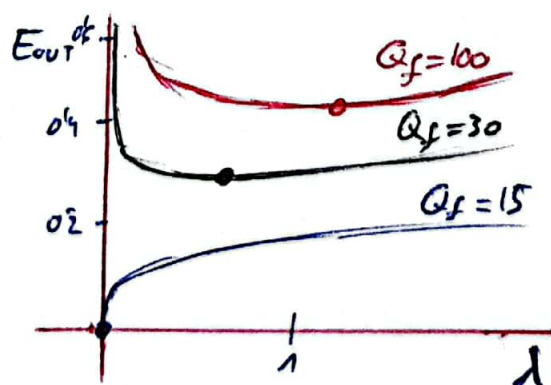
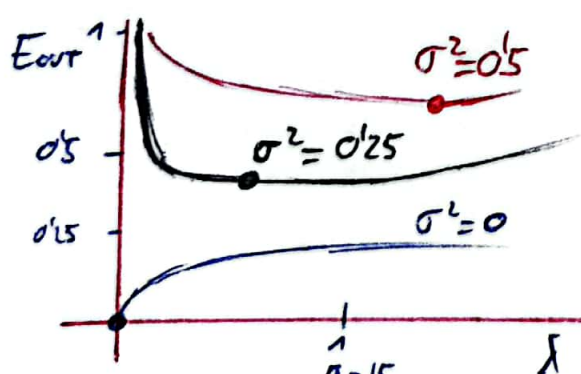
INFLUENCIA DE λ - RELACIÓN λ - RUIDO.

- + Como decíamos, debemos ajustar λ en cada problema:



- * Para $\lambda \approx 0 \rightarrow$ Sobreajuste
- * Para $\lambda \uparrow \rightarrow$ Infraajuste
- * Hay que buscar λ^* de forma experimental.

- + La cantidad de ruido influye en la elección de λ :



$$\Omega(w) = \sum_{g=0}^G w_g^2$$

- * Tanto si es error estocástico o determinista, cuanto más error haya \rightarrow Mayor λ necesitaremos para combatirlo.
- * Si no hubiera error de ningún tipo \rightarrow No necesitaremos regularización.