

TEMA 3 - FUNDAMENTOS DE AA

(de forma exacta)

- + Es imposible aprender f viendo solo una muestra (por inducción)
- + Sin embargo, podemos aprender propiedades que la estimamos, asumiendo % error.
- * Si los elementos de una muestra D son \sum independientes idénticamente distribuidos (no cambia la distrib.)
 \rightarrow Hay dependencia probabilística entre \rightarrow variable aleatoria (población) muestra
- * Teorema No-Free-Lunch
 \rightarrow Todos los algoritmos de aprendizaje son igual de buenos de media si consideramos todos los problemas de aprendizaje.
 \rightarrow Por tanto, cada aprendiz (A, \mathcal{H}) debe aplicarse a distribuciones P que pueda aprender correctamente \rightarrow Explotar car. específicos.
 $\rightarrow \forall$ algoritmo A , \exists distribución P con el que falla.

INFIRIENDO DE LA MUESTRA: SOLUCIÓN PAC

MODELO BIN

- + Tenemos una bolsa de canicas $\left\{ \begin{array}{l} - \text{verdes} \\ - \text{rojas} \end{array} \right.$
 \rightarrow ¿% canicas verdes en la bolsa?
 \rightarrow Hay demasiadas "canicas" \rightarrow Sacamos muestra y vemos %.
- + ¿Podemos garantizar que % verdes muestra \approx % verdes bolsa?
 \rightarrow No al 100%, pero es probable.
- * Desigualdad de Hoeffding:

$$P(D: |\% \text{ muestra} - \% \text{ población}| > \epsilon) \leq 2e^{-2\epsilon^2 N} \quad \forall \epsilon > 0.$$
 \rightarrow Dada una muestra i.i.d. de tamaño N , la probabilidad de que la diferencia entre la medición en la muestra y la real sea mayor que ϵ está acotada.
 $\rightarrow \epsilon$: la diferencia entre "dentro" y "fuera" (binario)
- * La cota depende solo de $\sum N$: el tamaño de la muestra.
 $\rightarrow \uparrow \uparrow \uparrow N \rightarrow \downarrow$ cota (mayor precisión).

RELACIÓN BIN - APRENDIZAJE

- + Ahora:
 \rightarrow Cada canica es un ejemplo de la población, con sus características (x) .
 \rightarrow Esta será "verde" o "roja" dependiendo de si una hipótesis h clasifica correctamente el ejemplo o no (clasificación binaria).
 \rightarrow % muestra $\Rightarrow E_{in}(h)$; % población $\Rightarrow E_{out}(h) = P_x[h(x) \neq f(x)]$.
 \rightarrow Aquí influye la distribución de la población.

- * Desigualdad de Hoeffding para una única hipótesis.

$$P(D: |E_{in}(h) - E_{out}(h)| > \epsilon) \leq 2e^{-2N\epsilon^2}$$
 para cualquier $\epsilon > 0$.
 \rightarrow Importante: esta probabilidad se cumple solo cuando $|\mathcal{H}| = 1$ (solo tomamos una hipótesis), y h es elegida sin ver los datos.
- * Este es el "resultado PAC" (Probably Approximately Correct), que cuantifica N tal que un algoritmo de aprendizaje produzca resultado PAC
 $\rightarrow E_{in}$ es "aleatorio", pero conocido.
 \rightarrow Si $E_{in}(h) \approx 0 \xrightarrow{N \gg \frac{1}{\epsilon^2}} E_{out}(h) \approx 0$ con alta probabilidad $\rightarrow f \approx h$ en X .
 \rightarrow Si $E_{in} \gg 0 \rightarrow E_{out}(h)$ también. (ii)
 \rightarrow * Por $E_{in}^{(h)} = 0.5$ (equivale a elegir aleatoriamente).
 \rightarrow * Si $E_{in}^{(h)} > 0.5 \rightarrow \bar{h}$ (función complementaria de h): $E_{in}(\bar{h}) < 0.5$.
- + Desarrollando la desigualdad:

$$P(D: |E_{out}(h) - E_{in}(h)| \leq \epsilon) > 1 - \delta \rightarrow E_{out}(h) \leq E_{in}(h) + \epsilon$$

$$\rightarrow E_{out}(h) \leq E_{in}(h) + \sqrt{\frac{1}{2N} \log \frac{2}{\delta}}$$
 $\rightarrow \uparrow \uparrow N \rightarrow \downarrow$ intervalo $\rightarrow \downarrow \delta \rightarrow \uparrow \uparrow$ intervalo
 \rightarrow Reducir $\delta \rightarrow$ aumentar precisión

- * Desigualdad Hoeffding con varias hipótesis.
 \rightarrow Esta vez, fijamos la hipótesis g a partir de los datos \rightarrow "Hemos mirado más de una hipótesis" $\rightarrow \uparrow \uparrow$ % de que $E_{in} = 0$, pero no que $E_{out} \approx E_{in}$.
- + Si $|\mathcal{H}|$ es finito, estamos fijando ese n° de modelos:

$$P(D: |E_{in}(g) - E_{out}(g)| > \epsilon) = P(\bigcup_{h \in \mathcal{H}} (D: |E_{in}(h) - E_{out}(h)| > \epsilon))$$

$$\text{Como } P(\bigcup_{i=1}^M B_i) \leq \sum_{i=1}^M P(B_i) \rightarrow P(D: |E_{in}(g) - E_{out}(g)| > \epsilon) \leq 2|\mathcal{H}|e^{-2\epsilon^2 N}, \quad \epsilon > 0.$$

RELACIÓN BIN-APRENDIZAJE (II)

$$O\left(\sqrt{\frac{\ln |\mathcal{H}|}{N}}\right)$$

* Cota de inecuación de Hoeffding.

+ Con probabilidad $1 - \delta$, $E_{out}(g) \leq E_{in}(g) + \sqrt{\frac{1}{2N} \ln \frac{2|\mathcal{H}|}{\delta}}$.

→ Si $|\mathcal{H}|$ es pequeño y $N \gg \ln |\mathcal{H}|$, entonces $E_{out}(g) \approx E_{in}(g)$.

* Solo requiere $P(x)$ para generar muestras i.i.d.

HIPÓTESIS DE REALIZABILIDAD

+ Consiste en suponer que $\exists h^* \in \mathcal{H}$ tal que $E_{out, f}(h^*) = 0$.

* Entonces, \mathcal{H} incluye al menos una función que $\forall P, f \rightarrow E_{in} = 0$.

+ Bajo esta hipótesis \rightarrow ERM en \mathcal{H} alcanza $h_s: E_{in}(h_s) = 0$.

* Por tanto $\rightarrow E_{out}(h_s) \leq \frac{1}{N} \log \frac{|\mathcal{H}|}{\delta}$ (con prob. $\geq 1 - \delta$).

DEFINICIÓN FORMAL DEL APRENDIZAJE PAC.

+ Una clase de funciones \mathcal{H} es PAC-aprendible si existen:

1. Una función $m_{\mathcal{H}}(\epsilon, \delta) \rightarrow \mathbb{N}$ (n° muestras para tener $< \epsilon$).
2. Un algoritmo de aprendizaje A .

tales que $\forall \epsilon, \delta \in (0, 1)$, $\forall P(x)$, $\forall f$: al ejecutar A con $N \geq m_{\mathcal{H}}(\epsilon, \delta)$ ejemplos i.i.d. generados por P y etiq. por f , A devuelve $h \in \mathcal{H}$: $P(ERM_{P, f}(h) \leq \epsilon) > 1 - \delta$.

+ ¿ $m_{\mathcal{H}}(\epsilon, \delta)$? \rightarrow depende de $\begin{cases} \epsilon \\ \delta \\ \mathcal{H} \end{cases}$ rango de ERM (u otra f pérdida).

+ Complejidad de muestra para cierto \mathcal{H} : mínimo n° de ejemplos $m_{\mathcal{H}}(\epsilon, \delta)$ que satisface los requerimientos de PAC.

+ Si la f pérdida tiene rango $[0, 1] \rightarrow$ toda \mathcal{H} finita realizable es PAC-aprendible, con complejidad de muestra:

$$m_{\mathcal{H}}(\epsilon, \delta) \leq \left\lceil \frac{1}{\epsilon} \log \frac{|\mathcal{H}|}{\delta} \right\rceil$$

DEFINICIÓN APRENDIZAJE PAC AGNÓSTICO

+ Una clase de funciones \mathcal{H} es PAC-aprendible agnóstica si $\exists m_{\mathcal{H}}(\epsilon, \delta) \rightarrow \mathbb{N}$, $\exists A$ algoritmo de aprendizaje tal que:

* $\forall \epsilon, \delta \in (0, 1)$, $\forall P(x)$. (✓)

al ejecutar A sobre $N \geq m_{\mathcal{H}}(\delta, \epsilon)$ ejemplos \rightarrow devuelve h :

$$P[E_{in}(h) \leq \min_{h' \in \mathcal{H}} E_{in}(h') + \epsilon] > 1 - \delta.$$

* A devuelve h cercana a la mejor posible dentro de \mathcal{H} .

+ Regla ERM es un ejemplo de PAC-agnóstica.

+ En este caso $\rightarrow m_{\mathcal{H}}(\epsilon, \delta) \leq \left\lceil \frac{2}{\epsilon^2} \log \frac{2|\mathcal{H}|}{\delta} \right\rceil$

* Menor que en PAC! $O\left(\frac{1}{\epsilon^2} \log \frac{|\mathcal{H}|}{\delta}\right)$

FACTIBILIDAD DEL APRENDIZAJE

+ El aprendizaje es posible probabilísticamente \leftarrow con muestras i.i.d. \leftarrow con $P(x)$ estable.

+ Para tener éxito, debemos encontrar h : $E_{out}(h) \approx 0$.

* Por ahora, solo garantizamos la inecuación de Hoeffding (que está acotado). $(\leq 2|\mathcal{H}|e^{-\epsilon^2 N})$

¿Entonces? \rightarrow ¿Podemos asegurar que $E_{out} \approx E_{in}$ lo suficiente? \rightarrow Inecuación Hoeffding.
 \downarrow ¿Podemos reducir E_{in} lo suficiente? \rightarrow Algoritmo de aprendizaje.

+ Dilema: ¿qué \mathcal{H} elegir?

* Si \mathcal{H} es complejo $\rightarrow |\mathcal{H}| \uparrow \uparrow \rightarrow E_{in} \uparrow \uparrow \rightarrow E_{out} \uparrow \rightarrow$ Necesitamos $\uparrow \uparrow N$.

* Si \mathcal{H} es simple $\rightarrow |\mathcal{H}| \downarrow \downarrow \rightarrow E_{in} \downarrow \downarrow \rightarrow E_{out} \downarrow \rightarrow$ (no se ajusta tanto) $\rightarrow E_{in} \approx E_{out}$ con menos ejemplos.

\rightarrow NOTA: $|\mathcal{H}|$ sube cuantos más parámetros tenga que ajustar.

TEOREMA FUNDAMENTAL DEL APRENDIZAJE PAC (son equivalentes)

+ \mathcal{H} es PAC-aprendible $\Leftrightarrow \mathcal{H}$ es PAC-aprendible agnóstica \Leftrightarrow

\Leftrightarrow Cualquier ERM es un aprendiz eficaz en PAC / PAC-agnóstico \Leftrightarrow

$\Leftrightarrow \mathcal{H}$ tiene VC-dimensión finita. (T.4)