

Notes 14: Weak and strong learning

1. WEAK LEARNING

Recall PAC learning definition (henceforth **strong** PAC learning):

Algorithm A PAC learns \mathcal{C} if

for any concept $c \in \mathcal{C}$ and any distribution \mathcal{D} over X

for **any** confidence parameter $\delta > 0$ and **any** accuracy parameter $\varepsilon > 0$

when A takes m samples from $\text{EX}(c, \mathcal{D})$

with prob. $\geq 1 - \delta$, A outputs hypothesis with error $\leq \varepsilon$

A needs to work for **arbitrarily small** $\delta > 0$ and $\varepsilon > 0$: stringent requirement!

What if A only is guaranteed to work for **some** $\delta > 0$ and $\varepsilon > 0$? (much weaker guarantee)

Turns out A can be boosted to a strong learning algorithm

2. BOOSTING CONFIDENCE

Suppose algorithm A , with probability $\geq 2/3$, outputs hypothesis with error $\leq \varepsilon$ (for any $\varepsilon > 0$)

A 's confidence δ bounded away from 0

Can be converted to strong PAC algorithm (with arbitrarily small δ and ε):

Strong PAC algorithm B

Repeat $t = 1, \dots, R$ times:

Run A on independent samples, with accuracy being $\varepsilon/2$, to get hypothesis h_t

Draw m' more samples S to evaluate hypotheses h_1, \dots, h_R

Output the hypothesis with least empirical error on S

$R \stackrel{\text{def}}{=} \frac{3}{2} \ln \frac{2}{\delta} = O(\ln \frac{1}{\delta})$ so that

$$\mathbb{P} \left[\text{none of } h_1, \dots, h_R \text{ has error } \leq \frac{\varepsilon}{2} \right] \leq \left(1 - \frac{2}{3} \right)^{3/2 \ln(2/\delta)} \leq e^{-\ln(2/\delta)} = \frac{\delta}{2}$$

$m' \stackrel{\text{def}}{=} O\left(\frac{1}{\varepsilon} \ln \frac{1}{\delta}\right)$ so that

Chernoff + Union Bound: with prob. $\geq 1 - \delta/2$,

all bad hypotheses among h_1, \dots, h_R have empirical error $\geq \frac{5}{6}\varepsilon$; and

some $\frac{\varepsilon}{2}$ -accurate hypothesis among h_1, \dots, h_R has empirical error $\leq \frac{4}{6}\varepsilon$

Hence any hypothesis with least empirical error must have (true) error $\leq \varepsilon$

Algorithm B succeeds with prob $\geq 1 - \delta$

A uses $m = \text{poly}\left(\frac{1}{\varepsilon}\right)$ samples $\implies B$ uses $Rm + m' = \text{poly}\left(\frac{1}{\varepsilon}, \ln \frac{1}{\delta}\right)$ samples

A runs in $T = \text{poly}\left(\frac{1}{\varepsilon}\right)$ time $\implies B$ runs in $RT + m' \text{poly}\left(\frac{1}{\varepsilon}\right) = \text{poly}\left(\frac{1}{\varepsilon}, \ln \frac{1}{\delta}\right)$ time

Summary: $O(\ln \frac{1}{\delta})$ calls to A ; $O\left(\frac{1}{\varepsilon} \ln \frac{1}{\delta}\right)$ further samples to test the hypotheses

3. BOOSTING ACCURACY

Call algorithm A weak PAC learning algorithm with **advantage** γ if

for any $c \in \mathcal{C}$, for any distribution \mathcal{D} , for any $\delta > 0$

with probability $\geq 1 - \delta$, output hypothesis h with $\text{err}_{\mathcal{D}}(h, c) \leq \frac{1}{2} - \gamma$

Getting advantage $\gamma = 0$ (i.e. $\text{err}_{\mathcal{D}}(h, c) = \frac{1}{2}$) is trivial: just output uniformly random guess

Goal: Turn any weak PAC algorithm A with advantage γ into strong PAC algorithm

with poly $\left(\frac{1}{\gamma}, \frac{1}{\varepsilon}, \frac{1}{\delta}\right)$ overhead in #samples and running time

Will show efficient boosting algorithm B with following structure

Boosting algorithm B

Draw independent training samples $S = \{(x^1, c(x^1)), \dots, (x^m, c(x^m))\}$ from $\text{EX}(c, \mathcal{D})$

Initially set $\mathcal{D}_1 =$ uniform distribution over S

Repeat $t = 1, \dots, R$ times:

 Run A on independent samples from $\text{EX}(c, \mathcal{D}_t)$ to get hypothesis h_t

 Adjust \mathcal{D}_t according to h_t to get updated distribution \mathcal{D}_{t+1} over S

Combine hypotheses h_1, \dots, h_R to get hypothesis h

Missing details:

 What are $\mathcal{D}_2, \mathcal{D}_3, \dots$?

 How to combine h_1, \dots, h_R into h ?

 Why $\text{err}_{\mathcal{D}}(h, c) \leq \varepsilon$?

History: Theory influenced practical algorithms!

 Kearns and Valiant (1989): introduced weak learning, showing weak learning may still be hard

 Freund and Schapire (1990): weak and strong learning are equivalent in distribution-free setting

 Freund and Schapire (1995): AdaBoost, now part of many machine learning libraries