

Online Learning ; Mistake-bounded Learning

- So far we've assumed that there is a distribution over some instance space X , and the learner gets data from such a distribution.
- This is often not realistic. Data typically does not come from a stationary distribution.
- Is there a "meaningful" model of learning which doesn't require any fixed distribution.
- We will think of an online/sequential learning framework.
We have discrete time-steps
 - (a) At time t an example $x_t \in X$ is given to the learning algorithm.
 - (b) Learning algorithm predicts $\hat{y}_t \in \{0,1\}$
(using all possible historical data)
 - (c) The true label $y_t \in \{0,1\}$ is revealed..

what's the notion of performance?

$$\text{For time } t, \text{ MISTAKES}(t) = \sum_{s=1}^t \mathbb{1}(\hat{y}_s \neq y_s).$$

#MISTAKES remains finite even as $t \rightarrow \infty$.

§ Assume that there is a class of functions C from $X \rightarrow \{0,1\}$ and that $\exists c \in C$, s.t. $y_t = c(x_t) \forall t$.

§ Formally, at time t

$$(x_1, y_1, \dots, x_{t-1}, y_{t-1}, x_t) \xrightarrow{\text{Learner}} \hat{y}_t.$$

Learning CONJUNCTIONS

- start with $h_0 \equiv x_1 \wedge \bar{x}_1 \wedge x_2 \wedge \bar{x}_2 \dots \wedge x_n \wedge \bar{x}_n$

- For $t=1, 2, \dots$

$$\hat{y}_t = h_{t-1}(x_t)$$

if $(y_t = 1 \text{ and } \hat{y}_t = 0)$ {

remove all literals that were false/0 in x_t

}

Theorem: $\forall \epsilon \text{ MISTAKES}(t) \leq 2n$

Proof: Suppose c is the target conjunction.

Maintain the invariant that $c(x) = 0 \Rightarrow h_t(x) = 0 \quad \forall t$

Every time we make a mistake, we remove at least one literal from h_t .

literals $\leq 2n$. □

MISTAKES $\leq n+1$.

Exercise, show that $n+1$ is tight. i.e. $\exists c$, and a sequence x_1, x_2, \dots st. ~~any~~^{this} algorithm would make at least $n+1$ mistakes.

Def: MISTAKE-BOUNDED LEARNING: We will say that a concept class C is learnable with a mistake bound M , if there exists a learning algorithm L , that when input any sequence $(x_t, y_t)_{t=1}^n$ satisfies $\forall t, \text{MISTAKES}(L, t) \leq M$.

Efficiency:

- Most generous: At time t , running time is $\text{poly}(t, n, \text{size}(c))$, where n is the size of examples in X , & $\text{size}(c)$ is the rep. size of the target.
- Space-bounded algorithms: Learning algorithm keeps some "sketch" S_t of data up to time t , $|S_t| \leq \text{poly}(n, \text{size}(c))$.
- Want that $S_{t+1} = f(S_t, x_t, y_t)$ and f is efficiently computable.
- Want that $\hat{y}_t = g(S_t, x_t)$ and g is efficiently computable.

PAC model: $EX(C, D)$ - example oracle

MQ - (Membership Queries) $x \rightarrow$ receive $c(x)$ x is chosen by L

EQ (Equivalence Queries) L produces h , \rightarrow get success if $h \equiv c$ or counterexample x st $h(x) \neq c(x)$.

• Mistake bound of M , implies that C is learnable using at most $M+1$ equivalence queries only

• $g(S_t, x_t) = \hat{y}_t$: $g(S_t, \cdot) \equiv h$ (hypothesis)

If C is learnable using q equivalence queries only then, C is learnable with mistake bound of at most q .

- Simulate L that only uses EQs.
- at any point it has a hypothesis h .

Defn: Conservative Online Learner: We say that an online learning algorithm is conservative if it only changes its prediction rule after making a mistake.

Lemma: If A learns a concept class C with mistake-bound M , then there exists a conservative algorithm A' which also learns C with mistake bound M .

Proof:

A : $\begin{matrix} \times & \times & \times \\ x_1 & x_{k+1} & x_t \\ \hat{y}_1 & \hat{y}_{k+1} & \hat{y}_t \end{matrix}$

A' : $\begin{matrix} \times & \checkmark & \checkmark & \checkmark & \times & \checkmark & \checkmark & \checkmark & \times \\ x_1 & x_2 & \dots & x_k & x_{k+1} & & & & x_t \\ \hat{y}_1 & \hat{y}_2 & \dots & \hat{y}_k & \hat{y}_{k+1} & & & & \hat{y}_t \end{matrix}$

A' copies A 's behaviour.

However it only further simulation of A if a mistake occurs

A' is clearly conservative.

#mistakes made by $A' =$ #mistakes made by A

Theorem: If C is (efficiently) online learnable with mistake bound M , where $M \leq \text{poly}(\text{size}(C), n)$, then C is (efficiently) PAC learnable.

Proof: feed examples drawn from $EX(C, D)$ to the online algorithm (conservative).

$h_1 \dots h_2 \dots h_3 \dots \dots h_m$
 if we go for $\frac{1}{\epsilon} \log(M/\delta)$ examples without making a mistake, then the h_i being used on this run has $\text{err}(h_i) \leq \epsilon$ w.p. $1-\delta$.

if $\text{err}(h_t) > \epsilon$, then $P(\text{no mistake seen in } \frac{1}{\epsilon} \log(\frac{M}{\epsilon}) \text{ examples}) \leq \frac{\epsilon}{M}$

Alg will always output one of h_1, \dots, h_m , over all failure prob $\leq \frac{M \cdot \epsilon}{M} = \epsilon$.

General Upper Bound:

• Let C be a finite concept class, then the HALVING ALGORITHM has MISTAKE bound at most $\log |C|$.

HALVING ALGORITHM:

• $C_1 = C$

for $t=1, 2, \dots$

- given x_t , computes $c(x_t) \forall c \in C_t$
and sets \hat{y}_t to be the majority of $(c(x_t) | c \in C_t)$.

- if $\hat{y}_t \neq \tilde{y}_t$, $C_{t+1} = C_t \setminus \{c \in C_t \text{ st. } c(x_t) = \tilde{y}_t\}$.

else $C_{t+1} = C_t$

- After every mistake $|C_{t+1}| \leq |C_t|/2$.

- $\therefore \# \text{ MISTAKES} \leq \log_2 |C|$.

Lower bound: $\text{VC-DIM}(C) \leq \max_t \text{MISTAKES}(t)$.

Examples: $X = \{(1, 0, \dots, 0), (0, 1, \dots, 0), \dots, (0, 0, \dots, 0, 1)\}$

$C = \text{DICTATORS}$.

$C_i(x) = x_i$.

$\text{MISTAKE-BOUND}(\text{HALVING ALG}) = 1$

Exercise: $X = \{0, 1\}^n$, then mistake bound is $\lfloor \log n \rfloor$ for DICTATORS.

Examples: Binary Search.

$X = \{1, \dots, 2^n\} \subseteq \mathbb{N}$.

C is half intervals: $C_j(x) = \begin{cases} 1 & \text{if } x \geq j \\ 0 & \text{o/w} \end{cases}$

$\text{VC}(C) = 1$.

$\text{MISTAKE-BOUND} \approx n = \log |C|$.

$$\text{err}_D(h) \geq \varepsilon$$

$$X_i \sim D \quad \mathbb{P}(h(X_i) \neq c(X_i)) \geq \varepsilon$$

$$X_i \sim D, \quad Z_i = \begin{cases} 1 & \text{if } h(X_i) \neq c(X_i) \\ 0 & \text{o/w.} \end{cases}$$

$$X_1, \dots, X_n, \quad \mathbb{P}\left(\sum_{i=1}^n Z_i = 0\right) = \prod_{i=1}^n \mathbb{P}(Z_i = 0)$$

$$1 - \varepsilon \leq e^{-\varepsilon}$$

$$\leq (1 - \varepsilon)^n \leq e^{-\varepsilon n} \leq \frac{\delta}{M}$$

$$n = \frac{1}{\varepsilon} \log\left(\frac{M}{\delta}\right)$$