

ONLINE LEARNING:

2 Formulations

→ Mistake Bounded Model

→ Regret Minimization.

MISTAKE BOUNDED MODEL:

Cannot make more than a fixed mistake.

$$f_i = f_*(x_i)$$

↳ ASSUMED CONSTANT!

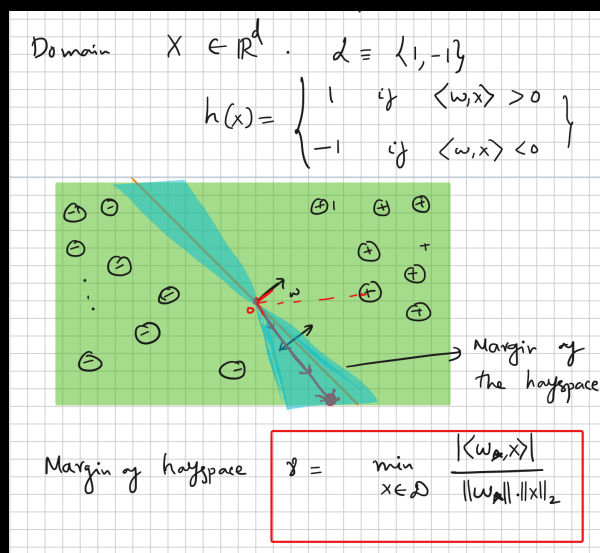
→ Let us see if this can be applied to a very simple problem.

ONLINE LEARNING OF HALF SPACES:

$$f_*(x) = \text{sign}(\langle w_*, x \rangle)$$

$$\|w_*\| = 1$$

$$\|x\| = 1$$



Given: The data is not only separated but also have a margin, γ .

$$\gamma = \min_{x \in \infty} \frac{|\langle w_*, x \rangle|}{\|w_*\| \cdot \|x\|_2} = \min_{x \in \infty} |\langle w_*, x \rangle|$$

MARGIN $\geq \gamma$: $|\langle w_*, x_i \rangle| \geq \gamma \quad \forall i$

Perceptron can learn such a model for the half spaces problem.

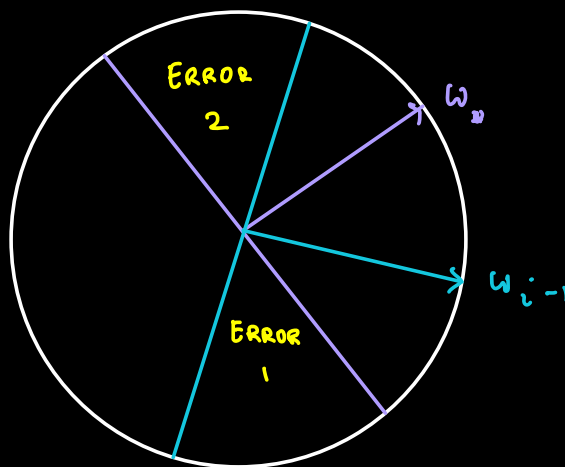
PERCEPTRON:

→ $w_0 \equiv$ random vector

→ On Day i

$$w_i = \begin{cases} w_{i-1} & \text{if no mistake} \\ w_{i-1} + y_i x_i & \text{if mistake} \end{cases}$$

Perceptron makes at most $1/\gamma^2 + 1$ mistakes
($\gamma \equiv$ margin of w_*)



$$y_i = \text{sign}(\langle w_*, x_i \rangle)$$

Proof:

we know

$$\|w_0\|^2 = \|w_1\|^2 = 1, \quad \|x_i\|^2 = 1$$

$\|w_i\|^2$ need not be 1.

On mistake: $w_i = w_{i-1} + y_i x_i \quad \text{--- (1)}$

$$1. \quad \|w_i\|^2 = \|w_{i-1} + y_i x_i\|^2$$

$$= \|w_{i-1}\|^2 + \|y_i x_i\|^2 + 2 \langle w_{i-1}, y_i x_i \rangle$$

$$\downarrow \qquad \qquad \qquad \downarrow$$

$$\underbrace{y_i^2}_{=1} \underbrace{\|x_i\|^2}_{=1} \quad 2 \cdot y_i \langle w_{i-1}, x_i \rangle$$

$$2 y_i \underbrace{\langle w_{i-1}, x_i \rangle}_{\text{---}}$$

$$y_{\text{pred}} = \text{sign}(\langle w_{i-1}, x_i \rangle)$$

mistake so $y_i \neq y_{\text{pred}}$

This is ≤ 0

$$\text{So } \|w_i\|^2 \leq \|w_{i-1}\|^2 + 1$$

$$\|w_0\|^2 = 1$$

So at any point

$$\|w_n\|^2 \leq 1 + M_n$$

$$\|w_n\| \leq \sqrt{1 + M_n}$$

2. Dot product with w^*

$$\langle w_i, w_* \rangle = \langle w_{i-1}, w_* \rangle + \underbrace{\langle y_i x_i, w_* \rangle}_{y_i \langle x_i, w_* \rangle}$$

$$y_i = \text{sign}(\langle x_i, w_* \rangle)$$

So it's always positive

$$\begin{aligned} |\langle x_i, w_* \rangle| \\ \geq \delta \quad [\text{Margin}] \end{aligned}$$

$$\langle w_i, w_* \rangle \geq \langle w_{i-1}, w_* \rangle + \delta$$

So if M_n mistakes

$$\langle w_n, w_* \rangle \geq \langle w_0, w_* \rangle + \delta M_n$$

$$\underbrace{\|w_n\| \cdot \|w_*\|}_1 \geq \langle w_0, w_* \rangle + \delta M_n \quad [\text{Cauchy-Schwarz}]$$

$$\sqrt{1 + M_n} \geq \underbrace{\langle w_0, w_* \rangle}_{\geq -1} + \delta M_n$$

$$\sqrt{1 + M_n} \geq \delta M_n - 1$$

$$1 + M_n \geq \gamma^2 M_n^2 + 1 - 2\gamma M_n$$

$$2\gamma M_n + M_n \geq \gamma^2 M_n^2$$

$$2\gamma + 1 \geq \gamma^2 M_n$$

$$M_n \leq \frac{2}{\gamma} + \frac{1}{\gamma^2} \leq \frac{1}{\gamma^2} + 2$$

$$\gamma \leq |\langle w_x, x_i \rangle| \leq 1$$

$$M_n \leq 1/\gamma^2 + 1$$

PERCEPTRON AS SGD WITH HINGE LOSS!

REGRET MINIMIZATION:

LEARNING WITH EXPERTS:

d : number of experts.

$$\text{Regret}(n) = \sum_{t=1}^n L(t) - \min_{i=1, \dots, d} \sum_{t=1}^n L(i, t).$$

* FOLLOW THE MAJORITY:

ASSUMPTION: One infallible expert.

$$\text{Regret}(T) \leq \log_2 d$$

PROOF: $w(t) = \# \text{ eligible experts}$

$$\bullet \quad w(t) \geq 1$$

$$\begin{aligned} \bullet \quad w(t) &\leq w(0) \cdot \left(\frac{1}{2}\right)^{L(t)} \\ &\leq d \cdot \left(\frac{1}{2}\right)^{L(t)} \end{aligned}$$

$$d \cdot \left(\frac{1}{2}\right)^{L(t)} \geq 1$$

$$L(t) \leq \log_2 d \quad //.$$

* WEIGHTED MAJORITY ALGORITHM:

- $w(0, i) = 1$
- Predict: Weighted majority
- If expert i made mistake:
 $w(t, i) = w(t-1, i) / 2$

$$L(T) \leq 2.4 (L_x(T) + \log_2 d)$$

↙ Total Loss till T

PROOF:

$$w(t) = \sum_{i=1}^d w(t, i) \quad [\text{Total weight of all experts}]$$

If we make a mistake on day t

$$w(t) \leq \left(\frac{3}{4}\right) w(t-1)$$

If $L(t)$ = our loss after t rounds, then

$$\begin{aligned} w(t) &\leq w(0) \cdot \left(\frac{3}{4}\right)^{L(t)} \\ &\leq d \cdot \left(\frac{3}{4}\right)^{L(t)} \end{aligned}$$

Also

$$w(t) \geq \left(\frac{1}{2}\right)^{L_x(t)}$$

Combine: $\left(\frac{1}{2}\right)^{L_*(t)} \leq d \cdot \left(\frac{3}{4}\right)^{L(t)}$

$$\left(\frac{4}{3}\right)^{L(t)} \leq 2^{L_*(t)} \cdot d$$

$$\log_2^{4/3} \cdot L(t) \leq L_*(t) + \log_2 d$$

$$L(t) \leq 2.4 \left(L_*(t) + \log_2 d \right).$$

Algorithms with $\frac{\text{Regret}(\tau)}{\tau} \rightarrow 0$ as $\tau \rightarrow \infty$ are

"NO-REURET ALGORITHM".

* MULTIPLICATIVE WEIGHTS UPDATE METHOD:

- $w(0, i) = 1$

- Predict: Select expert with probability \propto

$$\Pr[\text{expert } i] = \frac{w(t-1, i)}{\sum_{j=1}^d w(t-1, j)}$$

- If expert i made mistake:

$$w(t, i) = (1 - \epsilon) w(t-1, i)$$

[Update all experts who made mistake]

→ Total loss till T .

$$A(T) \leq (1+\epsilon) L_*(T) + \frac{\ln d}{\epsilon}$$

$(\epsilon < 1/2)$

$$\text{If } \epsilon = \sqrt{\frac{\ln d}{T}}$$

$$A(T) \leq L_*(T) + 2\sqrt{T \ln d}$$

* $L(T)$ here is loss at T unlike before.

PROOF:

$$W(t) = \sum_{i=1}^d w(t, i)$$

$$W(T, i) = (1-\epsilon) \sum_{t=1}^T L(t, i)$$

i. $L(t) = E[\text{loss incurred on day } t]$ → same as the loss of expert we chose

$$= \sum_{i=1}^d \Pr[\text{we pick expert } i] \cdot L(t, i)$$

$$= \sum_{j=1}^d \frac{w(t-1, j)}{\sum_{j=1}^d w(t-1, j)} \cdot L(t, j)$$

$\hookrightarrow W(t-1)$

$$L(t) = \frac{1}{W(t-1)} \cdot \sum_{i=1}^d w(t-1, i) \cdot L(t, i) \quad \text{--- (1)}$$

$$\text{ii. } w(t) = \sum_{i=1}^d w(t, i)$$

$$= \sum_{i=1}^d w(t-1, i) \cdot (1 - \varepsilon L(t, i))$$

$$= \sum_{i=1}^d w(t-1, i) - \varepsilon \sum_{i=1}^d w(t-1, i) \cdot L(t, i)$$

$$= w(t-1) - \varepsilon \cdot w(t-1) L(t) \quad [\text{From ①}]$$

$$w(t) = w(t-1) \cdot (1 - \varepsilon L(t))$$

$$\text{iii. } \boxed{1 - x \leq e^{-x} \quad \forall x}$$

$$w(t) = w(t-1) \cdot e^{-\varepsilon L(t)}$$

$$w(\tau) = w(0) \cdot e^{-\varepsilon (L(1) + L(2) + \dots + L(\tau))}$$

$$= w(0) \cdot e^{-\varepsilon \cdot A(\tau)}$$

↓
Total expected loss.

$$\text{iv. } w(\tau) \geq (1 - \varepsilon)^{L_{\star}(\tau)}$$

$$(1-\varepsilon)^{L_*(\tau)} \leq d \cdot e^{-\varepsilon A(\tau)}$$

$$L_*(\tau) \cdot \ln(1-\varepsilon) \leq \ln d - \varepsilon A(\tau)$$

$$\varepsilon A(\tau) \leq (-\ln(1-\varepsilon)) \cdot L_*(\tau) + \ln d$$

$$A(\tau) \leq \left(\frac{-\ln(1-\varepsilon)}{\varepsilon} \right) L_*(\tau) + \frac{\ln d}{\varepsilon}$$

v.

$$\frac{-\ln(1-x)}{x} \leq 1+x \quad \text{if } x < 1/2$$

$$\varepsilon < 1/2$$

$$A(\tau) \leq (1+\varepsilon) L_*(\tau) + \frac{\ln d}{\varepsilon}$$

BOOSTING:

ADABOOST

$\rightarrow D^{(0)} = \left(\frac{1}{d}, \frac{1}{d}, \dots, \frac{1}{d}\right)$. $w(0, i) = 1$
 for $i = 1, 2, \dots, d$.

$\rightarrow h^{(0)} = WL(D^{(0)})$.

For $t = 1, \dots, T$:

- Define $w(t, i) = \begin{cases} (1-\epsilon) w(t-1, i) & \text{if "Correct"} \\ \epsilon w(t-1, i) & \text{if wrong.} \end{cases}$
- $D^{(t)}$ = distribution proportional to weights
- $h^t = WL(D^{(t)})$

Output $h \equiv \text{MAJ}(h^0, h^1, h^2, \dots, h^T)$.

Correct means
 $h^{t-1}(x^i) = y^i$

Adaboost achieves accuracy $1 - \delta$ on the dataset if

$$T \geq \frac{2 \ln(1/\delta)}{(1/2 - \epsilon)^2}$$

Weak Learner: $\Pr[h(x^i) \neq y^i] \leq \epsilon$
 $(x^i, y^i) \sim D$

$$\Rightarrow T \propto \ln(1/\delta)$$

WL : 60% We need 90%

$$\epsilon = 0.4$$

$$\delta = 0.1$$

$$T \geq \frac{2 \ln(10)}{(0.1)^2} \approx 600$$

Generate 600 Weak Learners using
 ADABOOST!

$$w(t) \leq w(t-1) \cdot (1-\epsilon)^{L(t)}$$

$$L(i, t) = \begin{cases} 1 & \text{if } h^t(x^i) \neq y^i \\ 0 & \text{else} \end{cases}$$

$$L(t) = \mathbb{E}[\text{loss we incur at time } t]$$

↳ loss of the weak learner at that time

$$\geq 1 - \delta$$

$$w(t) \leq w(t-1) \cdot (1-\epsilon)^{1-\delta}$$

$$\leq w(t-1) \cdot e^{-\epsilon(1-\delta)}$$

$$\boxed{w(T) \leq w(0) \cdot e^{-\epsilon T(1-\delta)}}$$

$$\text{Basically } A(T) = (1-\delta)T$$

For every bad example / expert, that was wrong

for more than $T/2$ times. So weight became

$(1-\epsilon)w_{t-1}$ only less than $T/2$ times.

$$\text{BAD index : } w(i, T) \geq (1-\epsilon)^{T/2}$$

$$w(T) \geq |\text{BAD}| \cdot (1-\epsilon)^{T/2}$$

$$|B_{AD}| \cdot (1-\epsilon)^{\tau/2} \leq d \cdot e^{-\epsilon\tau(1-\delta)}$$

$$\left(\frac{|B_{AD}|}{d} \right) \leq e^{-\epsilon\tau(1-\delta)} \cdot (1-\epsilon)^{-\tau/2}$$

$$\leq e^{-\epsilon\tau(1-\delta)} \cdot e^{\epsilon(1+\epsilon)\tau/2}$$



$$\left. \begin{aligned} \frac{-\ln(1-x)}{x} &\leq 1+x \quad \text{if } x < 1/2 \\ \ln\left(\frac{1}{1-x}\right) &\leq x(1+x) \\ \frac{1}{1-x} &\leq e^{x(1+x)} \\ (1-x)^{-1} &\leq e^{x(1+x)} \end{aligned} \right\}$$

$$\frac{|B_{AD}|}{d} \leq e^{-\epsilon\tau \left((1/2 - \delta) - \epsilon/2 \right)}$$

$$\text{Set } \epsilon = 1/2 - \delta$$

$$\frac{|B_{AD}|}{d} \leq e^{-\tau \left(1/2 - \delta \right) \left(1/2 - \delta \right) / 2}$$

$$\leq e^{-\tau \left(1/2 - \delta \right)^2 / 2}$$

$$\text{So if } \tau \geq \frac{2 \ln(1/\delta)}{(\frac{1}{2} - \epsilon)^2}$$

$$\frac{|\mathcal{B}_{\text{BAD}}|}{d} \leq \delta$$

SUMMARY:

1. Mistake Bounded : Perceptron

* Halfspaces with margin:

$$|\langle w_*, x_i \rangle| \geq \gamma \quad \forall i$$

$$w_i = \begin{cases} w_{i-1} & \text{if no mistake} \\ w_{i-1} + y_i x_i & \text{if mistake} \end{cases}$$

$$y_i = \text{sign}(\langle w_{i-1}, x_i \rangle)$$

* Perceptron makes at most $\frac{1}{\gamma^2} + 1$ mistakes

Proof:

$$\|w_0\|^2 = \|w_*\|^2 = 1, \quad \|x_i\|^2 = 1$$

$$\rightarrow \|w_i\|_2^2 = \|w_{i-1} + y_i x_i\|_2^2$$

$$\|w_n\|^2 \leq 1 + M_n$$

$$\rightarrow \langle w_i, w_* \rangle = \langle w_{i-1}, w_* \rangle + \langle y_i x_i, w_* \rangle$$

$\rightarrow \gamma$

$$\langle w_n, w_* \rangle = \langle w_0, w_* \rangle + \gamma M_n$$

* Perceptron \Leftrightarrow SGD + Hinge Loss

2. Regret Minimization

i. 1 Infallible + FTM:

$$\text{Regret}(T) \leq \log_2 d$$

ii. WMA:

$$L(T) \leq 2.4 (L_*(T) + \log_2 d)$$

iii) No-Regret: MWM:

Predict, Pr $\propto w(t-1, i)$

on Mistake: $w(t, i) = (1 - \epsilon L(t, i)) w(t-1, i)$

$$A(T) \leq (1 + \epsilon) L_*(T) + \ln d / \epsilon$$

$$\epsilon = \sqrt{\frac{\ln d}{T}}$$

$$A(T) \leq L_*(T) + 2 \sqrt{T \ln d}$$

PROOF:

$$\rightarrow w(t) \leq w(t-1) \cdot (1 - \epsilon)^{L(t)}$$

$$w(T) \leq w(0) \cdot (1 - \epsilon)^{A(T)} \leq d \cdot e^{-\epsilon(A(T))}$$

$$\rightarrow w(T) \geq (1 - \epsilon)^{L_*(T)}$$

$$(1 - \epsilon)^{L_*(T)} \leq d \cdot e^{-\epsilon(A(T))}$$

$$L_n(T) \ln(1-\epsilon) \leq \ln d - \epsilon A(T)$$

$$A(T) \leq \left(\frac{-\ln(1-\epsilon)}{\epsilon} \right) L_n(T) + \frac{\ln d}{\epsilon}$$

3. Boosting : ADABOOST

Weak Learner : $\Pr[h(x^i) \neq y^i] \leq \delta$

Achieves accuracy $1-\delta$ if

$$\text{Number of WLS, } T \geq \frac{2 \ln(1/\delta)}{(1/2 - \delta)^2}$$

PROOF :

$$\rightarrow A(T) \geq T \times (1-\delta)$$

$$w(T) \leq w(0) \cdot (1-\epsilon)^{A(T)} \leq d \cdot e^{-\epsilon T(1-\delta)}$$

$$\rightarrow w(T) \geq |BAD| \cdot (1-\epsilon)^{T/2}$$

→ Properties:

$$1-x \leq e^{-x} \quad \forall x$$

$$\frac{-\ln(1-x)}{x} \leq 1+x \quad \text{if } x < 1/2$$

QUESTIONS

$$\gamma \in (0, 1)$$

at least ' γd ' good experts \rightarrow loss $\leq L$

$$L(\tau) \geq (1-\epsilon)^L \cdot \gamma d + \text{other weights}$$

$$L(\tau) \leq d \cdot e^{-\epsilon(A(\tau))}$$

$$(1-\epsilon)^L \cdot \gamma d \leq d \cdot e^{-\epsilon(A(\tau))}$$

$$L \ln(1-\epsilon) + \ln \gamma \leq -\epsilon A(\tau)$$

$$A(\tau) \geq \frac{L}{\epsilon} \ln(1-\epsilon) + \frac{\ln \gamma}{\epsilon}$$