

LECTURE 18

1 2

BATCH :

- TRAINING AND TEST DATA GENERATED BY SAME DISTRIBUTION
- IF MODEL CLASS NOT TOO COMPLEX AND ENOUGH EXAMPLES
MODEL THAT DOES BEST ON TRAINING DATA IS NOT TOO MUCH WORSE ON TEST DATA

ON-LINE :

- ALL IS IN FLUX
- NO STATISTICAL ASSUMPTIONS
- STILL CAN BOUND "REGRET" \equiv

TOTAL LOSS OF ON-LINE - TOTAL LOSS OF BEST OFF-LINE
CHOSEN IN HIND SIGHT

- BOUNDS HOLD FOR ARBITRARY SEQUENCES OF EXAMPLES

On-Line Learning

| | experts | | | | | | |
|---------|-----------|-----------|-----------|-----------|----------------|-----------------------------|---------------------|
| | E_1 | E_2 | E_3 | E_n | predic tion | <i>true</i> <i>label</i> | loss |
| day 1 | 1 | 1 | 0 | 0 | 0 | 1 | 1 |
| day 2 | 1 | 0 | 1 | 0 | 1 | 0 | 1 |
| day 3 | 0 | 1 | 1 | 1 | 1 | 1 | 0 |
| day t | $x_{t,1}$ | $x_{t,2}$ | $x_{t,3}$ | $x_{t,n}$ | \hat{y}_t | y_t | $ y_t - \hat{y}_t $ |

Protocol of the Master Algorithm

For $t = 1$ To T Do

Receive $x_t \in \{0, 1\}^n$

Predict $\hat{y}_t \in \{0, 1\}$

Get label $y_t \in \{0, 1\}$

Incur loss $|y_t - \hat{y}_t| \in \{0, 1\}$

CASE 1: THERE IS A CONSISTENT EXPERT

GIVEN SEQUENCE (\bar{x}_t, y_t) s.t

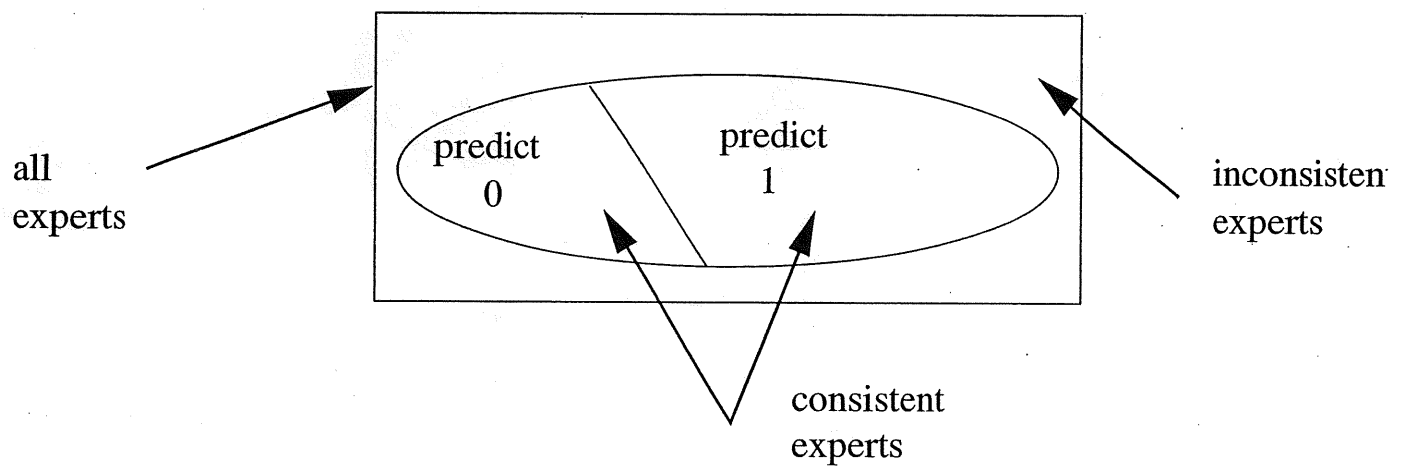
$x_{t,i} = y_t$ for all t

LOSS OF OFF-LINE COMPARATOR IS ZERO

NOISE-FREE CASE

Halving Algorithm

[BF]



- Predicts with majority
- If mistake then number of consistent experts is halved

A run of the Halving Algorithm

| E_1 | E_2 | E_3 | E_4 | E_5 | E_6 | E_7 | E_8 | majority | true label | loss |
|------------|-------|-------|-------|-------|-------|-------|-------|----------|------------|------|
| 1 | 1 | 0 | 0 | 1 | 1 | 0 | 0 | 1 | 0 | 1 |
| x | x | 0 | 1 | x | x | 1 | 1 | 1 | 1 | 0 |
| x | x | x | 1 | x | x | 0 | 0 | 0 | 1 | 1 |
| x | x | x | ↑ | x | x | x | x | | | |
| consistent | | | | | | | | | | |

For any sequence with a consistent expert,
HA makes $\leq \log_2 n$ mistakes

GAME AGAINST NATURE (ADVERSARY)

WHICH CHOOSES THE PREDICTION VECTOR \bar{x}_t AND LABEL y_t

IF THERE IS ONE CONSISTENT EXPERT

THEN ALG. $\leq \log_2 n$ MISTAKES

Case 2:

What if no expert is consistent?

For any sequence $S = (x_1, y_1), (x_2, y_2), \dots, (x_T, y_T)$

- $L_A(S)$ is total loss of alg. A and
- $L_i(S)$ is the total loss of expert E_i

RELATIVE LOSS

Want bounds of the form:

$$\forall S: L_A(S) \leq a \min_i L_i(S) + b \log(n)$$

where a, b are constants

Bounds loss of algorithm
relative to
loss of best expert

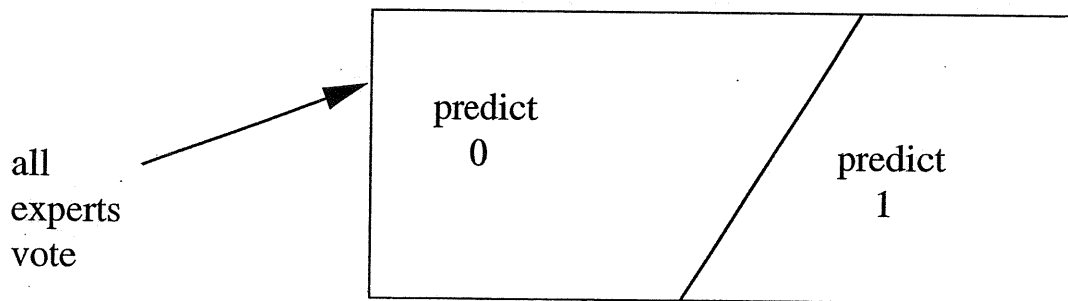
$$a = 1$$

$$L_A(S) - \min_i L_i(S) \quad \text{CALLED REGRET}$$

Can't wipe out experts!
One weight per expert

Weighted Majority Algorithm

[LW]



- Predicts with larger side
- Weights of wrong experts are multiplied by $\beta \in (0, 1]$
- β IS FITNESS FACTOR
- HA : $\beta = 0$

Number of mistakes of the WM algorithm

$M_{t-1,i}$ = # of mistakes of E_i before trial t

$w_{t-1,i}$ = $\beta^{M_{t-1,i}}$ weight of E_i at trial t , $w_{0,i} = 1$

$W_{t-1} = \sum_{i=1}^n w_{t-1,i}$ total weight at trial t

$$\text{Minority} \leq \frac{1}{2} W_{t-1}$$

$$\text{Majority} \geq \frac{1}{2} W_{t-1}$$

If no mistake then

minority multiplied by β

$$W_t \leq 1 W_{t-1}$$

If mistake then
majority multiplied by β

$$\begin{aligned}
 W_t &\leq 1 \cdot \frac{1}{2} W_{t-1} \text{ minority} + \beta \cdot \frac{1}{2} W_{t-1} \text{ majority} \\
 &= \frac{1 + \beta}{2} W_{t-1}
 \end{aligned}$$

$$\begin{aligned}
 W_T &\leq \left(\frac{1 + \beta}{2} \right)^M W_0 \\
 \text{total final weight}
 \end{aligned}$$

$$W_T = \sum_{j=1}^n w_{T,j} = \sum_{j=1}^n \beta^{M_j} \geq \beta^{M_i}$$

$$\left(\frac{1 + \beta}{2} \right)^M \underbrace{W_0}_n \geq \beta^{M_i}$$

$$M \leq \frac{-\ln \beta}{\ln \frac{2}{1+\beta}} M_i + \frac{1}{\ln \frac{2}{1+\beta}} \ln n$$

$$M_{\beta = 1/e} \leq \underbrace{2.63}_a \underbrace{\min_i M_i}_{M^*} + \underbrace{2.63}_b \ln n$$

For all sequences, loss of the master algorithm is comparable to the loss of the best expert

Relative loss bounds

[F]

WITH FANCY CHOICE OF β THAT DEPENDS ON n, M^* :

$$M \leq 2 M^* + 2 \sqrt{M^* \ln(N)} + \log_2 n$$

↑
NECESSARY
FOR DETERMINISTIC
PREDICTION

STREAMLINE SETUP (NO LABELS)

FOR $t = 1$ TO T DOCHOOSE AN EXPERT i GET LOSS VECTOR $\vec{L}_t \in [0,1]^N$ INCUR LOSS $L_{t,i}$

GOAL: ACHIEVE SMALL REGRET

TOTAL LOSS OF ALG - TOTAL LOSS OF BEST

ALG I: DETERMINISTIC FOLLOW THE LEADER

- ALWAYS CHOOSE AN EXPERT OF MINIMAL LOSS

ADVERSARY:

- CHOSEN EXPERT 1 UNIT OF LOSS
 - ALL OTHERS LOSS 0
- (T IS # OF TRIALS)

| | | | | |
|-----|--|--|--|--|
| 0 | | | | |
| 1 | | | | |
| 2 | | | | |
| n | | | | |

LOSS OF ALG

 $\approx n$ LOSS OF BEST

LOSS OF ALG

 T

LOSS OF BEST

 $\lfloor T/n \rfloor$

ALG III: PERTURB LOSSES OF EXPERTS
PREDICT W. PERTURBED LEADER

LATER!

29

ALG IV: HEDGE ALGORITHM
(SIMILAR TO RANDOMIZED WEIGHTED
MAJORITY ALGORITHM)

PROBABILISTIC CHOICE OF EXPERT

\bar{w}_{t-1} PROBABILITY VECTOR USED AT TRIAL t

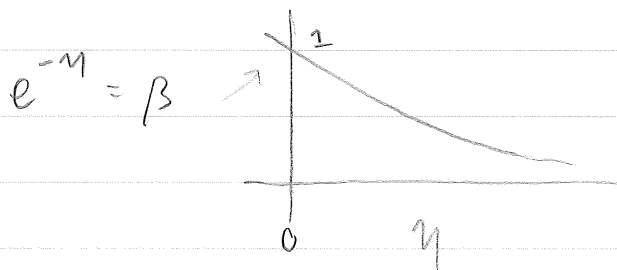
$w_{t-1,i}$ "BELIEVE" AT TRIAL t THAT i IS BEST

$$w_{0,i} = \left(\frac{1}{n}, \frac{1}{n}, \dots, \frac{1}{n} \right)$$

$$w_{t,i} = \frac{w_{t-1,i} e^{-\eta L_{t,i}}}{Z_t}$$

↑ NORMALIZATION

$\eta > 0$ LEARNING RATE



$$e^{-\infty} = 0$$

$$w_{t-1,i} = \frac{e^{-\eta L_{t-1,i}}}{Z_t}$$

AS $\eta \rightarrow \infty$, ALL WEIGHT PLACED ON BEST
& HEDGE BECOMES "FOLLOW THE LEADER"
(TIES BROKEN UNIFORMLY)

$\eta = 0$ WEIGHTS UNCHANGED

$\eta > 0$ GRADUALLY MOVE WEIGHT
TO EXPERTS W. LOW LOSS
"SOFT MIN"

$\eta < 0 \rightarrow$ HIGH LOSS
"SOFT MAX"

NEXT
CLASS:

IF η TUNED AS FUNCTION OF
 n & \hat{L} THEN

$$\sum_{t=1}^T \bar{w}_{t-1} \cdot L_t - \underbrace{\inf_i L_{t,i}}_{L^*} \leq \sqrt{2\hat{L} \ln n} + \ln n$$

IF $L^* \leq \hat{L}$

LOSS OF ALG - LOSS OF BEST

REGRET