

## ① the Mistake Bounded Model

- difference from CRatio

- measure the difference not ratio

- compare not to the best sequence of adap. decisions the algo could have made

but to the best fixed "expert" in hindsight.

So fixed set of decision profiles could follow (N of them)  
which is best?

Eg. have several paging algs can use, which one should we use?

Get: Bounds of the form

$$\# \text{mistakes of algo} \leq \# \text{mistakes of best expert} + \underbrace{\hspace{10em}}_{\substack{\uparrow \\ \text{mult. loss}}} + \underbrace{\hspace{5em}}_{\substack{\uparrow \\ \text{additive loss.}}}$$

Ideally: get mult loss =  $(1 + \epsilon)$

and additive loss  $\frac{\log N}{\epsilon}$

$\Rightarrow$  almost as good as best expert "in the long run"

Start slow:

① Show can get

$$A_T \leq O(\log N) (\text{Best expert} + 1)$$

(a) What if best expert makes no errors. ?

(b) — — — —  $m^*$  errors. ?

So both additive and mult terms are  $O(\log N)$  ! 😞

② Then introduce the multiplicative weights (MW) or Weighted Majority (WM) algo.

$$A_T \leq \underbrace{2.4 m^*}_{\text{mult}} + \underbrace{O(\log N)}_{\text{add}} !!$$

So much better now!

However: it is a deterministic algorithm, and no deterministic algo can do

$$A_T \leq (2 - \epsilon) m^*$$



### ③ Randomized Weighted Majority

In the bad example, should have hedged bets.

Suppose make random choices of experts.

$$\text{then } E[\text{Algo Mistakes}] \leq (1+\epsilon) m^* + \frac{O(\lg N)}{\epsilon}$$

↑  
best expert.

This is the holy grail, — cannot do better. 😊

---

BTW: in all these arguments, show that  $\forall$  each expert  $i$

$$E[\text{Algo}] \leq \underbrace{\alpha}_{\uparrow} m_i + \underbrace{\beta}$$

↑  
expert  $i$ 's mistakes

which implies statement for best expert.

---

See the notes for proof / algo.

assume losses are in  $[0,1]$

Slight generalization / restatement of model.

Done simultaneously

• Each time expert gives a loss value  $l_i^t$   
⇒ put them in vector  $(l_1^t, l_2^t, \dots, l_N^t) = \vec{l}(t)$

• We choose probability of picking  $1, 2, \dots, N$  as

$$(p_1^t, p_2^t, \dots, p_N^t) = \vec{p}(t)$$

$$\text{Then } E[\text{loss for us}] = \langle \vec{p}(t), \vec{l}(t) \rangle$$

↑ inner product

$$\sum_i p_i^t = 1$$
$$p_i^t \geq 0$$

$$= \sum_i p_i^t l_i^t$$

So imagine:

Each time the expert panel / Nature / adversary / world chooses a loss vector  $l^t$   
- can depend on  $p^1, p^2, \dots, p^{t-1}$ , not on  $p^t$

We choose a prob. vector  $p^t$

- which can depend on the past losses  $l^1, l^2, \dots, l^{t-1}$  but not this time

- also can depend on our choices if we want.  
(but there are implied by losses)

And then our actual loss at time  $t$  is dot product  $\langle p^t, \ell^t \rangle$  <sup>or realized</sup>

$\langle p^t, \ell^t \rangle$ . ← (can think as "expected" loss)

We get that for all experts  $i$

$$\sum_t \langle p^t, \ell^t \rangle \leq \underbrace{\sum_t \ell_i^t}_{\substack{\uparrow \\ \text{loss of expert } i}} \times \underbrace{(1 + \epsilon)} + \underbrace{\frac{\log N}{\epsilon}}$$

our "expected" loss

———— x ————

This view is useful:-

Nature chooses a loss vector  $\ell^t$  in  $[0, 1]^N$ .

We choose an "action"  $p^t$  in the "probability simplex"

$$\Delta_N = \{ p \in [0, 1]^N : \sum_i p_i = 1 \}$$

Our loss is  $\langle p^t, \ell^t \rangle$

want to compare with the best coordinate

$$\min_i \ell_i^t = \min_{q \in \Delta_N} \langle \ell^t, q^* \rangle$$

make sure you believe this (its simple)

Good: why is this worth studying?

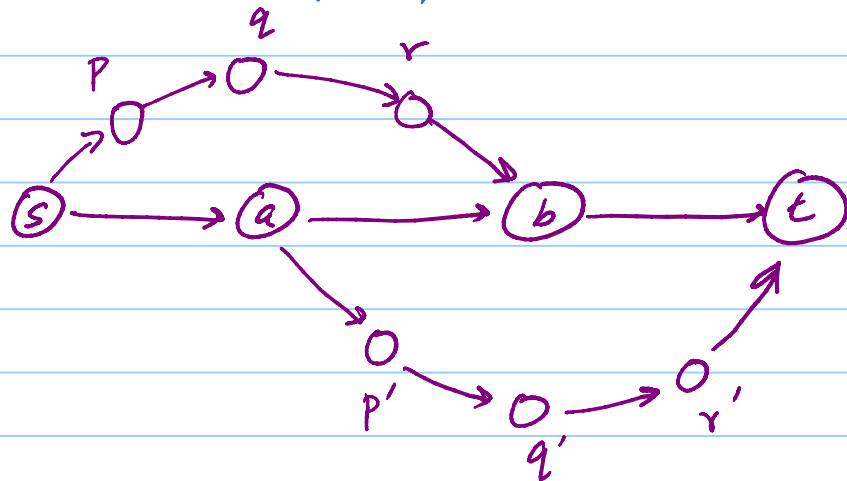
① Choose best of multiple options online

(really choosing the best of  $N$  "experts")

↑ maybe self-styled,

② Algorithmic idea of multiplicative weights is very pervasive

- e.g. here's a max-flow problem.



all edges have capacity 1 (say).

say send flow on  $s \rightarrow a \rightarrow b \rightarrow t$  (unit amount)  
then saturate those edges.

But not a max flow (since  $F^* = 2 \Rightarrow$

How to proceed. ① either use residual graph, etc.

② Or use MW.

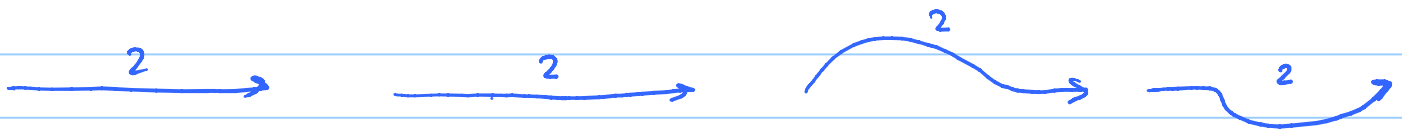
may depend on capacities,  
for now assume unit cap.

Algo: put lengths 1 on all edges.

→ push  $F^*$  flow on shortest path.  
Now set lengths to be  $(1 + \epsilon)^{\text{flow on edge}}$

repeat  $\frac{F^* \log n}{\epsilon}$  times

Take average of all these flows



Etc: each flow under-uses many edges  
and over-uses many others.

But averaging them ensures that each edge gets  
about the right amount of flow.

Another Example: Boosting in Learning. (HW?)

Suppose you have a weak classifier

Given any weights on data points, correctly classifies 51%  
of data.

Then can use a "majority of <sup>weak</sup>  $n$  classifiers" to get a  
strong classifier  
that gets 99% of data right.