

Experts Problem : Multiplicative Weight Updates¹

- **The Experts Problem.** Suppose you want to predict the weather. Unfortunately, you have no idea of meteorology. But you have n friends who are *experts* (that is, they have a PhD; it doesn't mean they are infallible) who are willing to tell you their opinions. Your goal is to use their opinions to make a good prediction.

Let's formalize. We denote $e_i(t) \in \{0, 1\}$ to be expert i 's prediction of whether it will rain or not on the t th day. The index i ranges from 1 to n . Before day t begins, you observe these $e_i(t)$'s and then you need to make a prediction $a(t) \in \{0, 1\}$. Formally, you have to design an algorithm \mathcal{A} which at step t takes input: $\{e_i(s) : 1 \leq s \leq t\}$ and $\{r(s), a(s) : 1 \leq s \leq t - 1\}$, and must output $a(t)$.

Then the t th day unfolds, and you get to see $r(t) \in \{0, 1\}$ whether it actually rained or not. If $a(t)$ doesn't match $r(t)$, say you pay a dollar – this is your *loss* $\ell(t)$ on day t . Thus, formally, $\ell(t) := 1 - a(t) \cdot r(t)$. We want to devise an algorithm \mathcal{A} to make the total loss incurred, $\text{loss} := \sum_{t=1}^T \ell(t)$, as small as possible. This is the experts problem.

Remark: *This is a very important problem in machine learning where the business is making predictions based on past observations. The algorithm that we will see for this problem indeed has many applications and has been rediscovered in many different guises. In particular, although we are sticking to the loss being 0 or 1, one could imagine being it real valued, and indeed, in some situations, one can imagine the experts also being a continuous set of points. Once again, a whole course can probably be taught on this. We are just looking at couple of facets this week.*

- **Benchmark : Best Expert.** Note that this is not your usual computation problem with a fixed input and output. It has a flavor of a game that one is playing with nature, and in a sense nature (who is responsible for $r(t)$) gets to “play second”. That is, our algorithm's decision has to be made under uncertainty of the outcome. Whenever this happens, it is not clear what “the best algorithm” even means. What is a good benchmark to compare one's algorithm against?

Well, since our algorithm is deriving information from the experts, perhaps we should compare our losses with those of the n experts. In particular, at time t define $\ell_i(t)$ to be 1 if expert i makes a mistake at time i , and let $\text{loss}_i := \sum_{t=1}^T \ell_i(t)$. But our experts can have differing qualities – which one should we choose our benchmark as? Well, it is clear that the “best expert” (the one with the smallest loss_i) is the most stringent among these benchmarks; let us see if we can design an algorithm which is as good as the best expert.

Remark: *Note that in some case an algorithm can make more than the best expert. Imagine an algorithm which “knew” nature, that is, knew $r(t)$ before hand. Then it would just predict $a(t) = r(t)$ and would be better than the best expert unless the best expert also knew nature.*

- **Warm-up: There is a perfect expert.** Let's start with the case there is one “perfect expert”. That is, there is some i^* with $\text{loss}_{i^*} = 0$. If we knew who this perfect expert was, then our strategy is clear – predict whatever i^* is predicting. Can we quickly recognize this expert?

¹Lecture notes by Deeparnab Chakrabarty. Last modified : 22nd April, 2021
These have not gone through scrutiny and may contain errors. If you find any, or have any other comments, please email me at deeparnab@dartmouth.edu. Highly appreciated!

The answer is yes, and the algorithm is simple. It maintains a candidate set A of active experts; any expert not in A has made some mistake in the past and thus cannot be the perfect expert. Initially $A = [m]$, and we wish to whittle this set down fast. The MAJORITY algorithm is this: always predict what the majority of the experts in A are predicting (if $|A|$ is even, break ties arbitrarily). Every time our algorithm makes a mistake, that is, everytime loss increments by $+1$, we know that at least $|A|/2$ experts are whittled out. Thus, the maximum number of mistakes MAJORITY makes is $\log_2 n$.

Theorem 1. In the case of an perfect expert, the MAJORITY algorithm has regret $\leq \lg n$.

- **Weighted Majority.** However, perfect experts are mythical. What if we were only promised an expert who is correct 99% of the time? How would MAJORITY perform? Well, not too well as defined since the near-perfect expert could make a mistake on the first day along with the majority and be immediately whittled out. How should we fix this?

Idea 1: Suppose we knew we were playing for T days. Keep taking MAJORITY till some expert makes $> 0.01T$ mistakes and then whittle them out. Two issues: one, it needs to know that some expert was as good as 99% (how would we know that?). More seriously, it's not a good algorithm – find an example where this algorithm fails badly.

Idea 2: When an expert makes a mistake, instead of whittling them out, just decrease their “importance”. More precisely, every expert has an importance or a weight $w_i(t)$. Initially, all the weights are 1; all experts are equal in our eyes. Each day, we go with the *weighted* majority. More precisely, for each choice we figure out the total *weight* of the experts predicting that value, and go with the larger one. Upon receiving the truth, that is, $r(t)$, we *update* the weights as follows: we penalize every expert which makes a mistake at time t by halving their weight. The full algorithm is described below.

```

1: procedure WEIGHTED MAJORITY:
2:   Maintain weights  $w_i(t)$  for each expert with  $w_i(1) = 1$  for all  $i$ .
3:   for Days  $t = 1$  to  $T$  do:
4:     Receive  $e_i(t)$  from each expert  $i$ .
5:     Let  $W_0 := \sum_{i:e_i(t)=0} w_i(t)$  and  $W_1 := \sum_{i:e_i(t)=1} w_i(t)$ .
6:     Predict 0 if  $W_0 \geq W_1$ , otherwise predict 1.
7:     Receive  $r(t)$  and incur loss  $\ell(t) = 1$  if  $a(t) \neq r(t)$ , and  $\ell(t) = 0$  otherwise.
8:     for every  $i$  with  $e_i(t) \neq r(t)$  do:
9:       Set  $w_i(t+1) = w_i(t)/2$ 

```

Theorem 2. If there exists an expert which makes at most k mistakes, then the WEIGHTED MAJORITY algorithm makes at most $2.41(k + \log_2 n)$ mistakes, where $2.41 \approx \frac{1}{\log_2(4/3)}$.

Corollary 1. If there is an expert who is correct at least 99% of the time, then if one plays for $T \gg \log n$ days, the WEIGHTED MAJORITY is correct $\geq 96\%$ of the time.

Proof. (of Theorem.) The analysis is similar in spirit to the analysis of the MAJORITY algorithm – we focus on the times t at which the algorithm makes a mistake. Consider such a time t when $a(t) \neq r(t)$.

By the algorithm's design we get to see at this t ,

$$\sum_{i:e_i(t) \neq r(t)} w_i(t) \geq \sum_{i:e_i(t) = r(t)} w_i(t) \quad (1)$$

In plain English, the total weight of the experts that agreed with nature is less than the total weight of experts which disagreed.

Now note that for all experts in the LHS, their weight $w_i(t+1) = w_i(t)/2$ while for all experts in the right, $w_i(t+1) = w_i(t)$. Therefore, the *total weight* goes down by a **multiplicative** factor.

More precisely, let's define for any t , define the total weight as

$$Z(t) := \sum_{i=1}^n w_i(t)$$

Note that $Z(1) = n$ since $w_i(1) = 1$ for all $1 \leq i \leq n$. For any time t at which WEIGHTED MAJORITY makes a mistake, (1) implies $\sum_{i:e_i(t)=r(t)} w_i(t) \leq \frac{Z(t)}{2}$. Therefore, we get

$$\begin{aligned} Z(t+1) &= \sum_{i:e_i(t) \neq r(t)} w_i(t+1) + \sum_{i:e_i(t) = r(t)} w_i(t+1) \\ &= \sum_{i:e_i(t) \neq r(t)} \frac{w_i(t)}{2} + \sum_{i:e_i(t) = r(t)} w_i(t) \\ &= \frac{Z(t)}{2} + \frac{1}{2} \sum_{i:e_i(t) = r(t)} w_i(t) \\ &\leq \frac{3Z(t)}{4} \end{aligned}$$

The above was for the times when our algorithm makes a mistake. What about the times when it doesn't? Well, since the weights never *increase*, we have the trivial inequality $Z(t+1) \leq Z(t)$. Thus, if our algorithm makes ℓ mistakes (that is, $\ell = \text{loss}$), after T rounds we get

$$Z(T) \leq (3/4)^\ell \cdot \underbrace{Z(1)}_{\text{which equals } n} \leq \left(\frac{3}{4}\right)^\ell \cdot n \quad (2)$$

So if our algorithm makes a lot of mistakes, the *potential* Z falls rapidly. Why is this at all useful in comparing with the best expert's loss? This is the second insight: if there is some expert making only k mistakes, then his weight at the end is *at least* $(1/2)^k$. Even if all the other experts are duds making tons and tons of mistakes, this expert forces

$$Z(T) \geq \frac{1}{2^k} \quad (3)$$

Now rest is arithmetic. From (2) and (3) we get

$$\frac{1}{2^k} \leq \left(\frac{3}{4}\right)^\ell \cdot n \quad \Leftrightarrow \quad \ell \leq \frac{1}{\log_2(4/3)} \cdot (k + \log_2 n)$$

Taking log base 2, swapping signs, and rearranging

completing the proof. □

Remark: *It is not hard to show that no deterministic algorithm can get a better than factor 2 approximation. More precisely, for any strategy and for any T , there is a collection of expert answers and “truths” such that the best expert makes $\leq \delta$ mistakes but the algorithm needs to make 2δ mistakes. In fact, this is true even with just two experts. Can you imagine such a scenario?*

- **Getting arbitrarily close with randomization.** Reiterating the theme that we have seen throughout the course, randomness allows us to get a better algorithm because it allows us to be more robust. This time more robust to the whims of nature. The flip is that we will have our result hold in expectation, and there will be a probability that it does poorly. In this lecture, I will stick only to the analysis of the expectation.

Formally, the quantity $a(t)$ we predict is now random variable which 0 with a certain probability and 1 with the remainder. The *expected* number of mistakes made by the algorithm at time t is therefore the *probability* with which $a(t) \neq r(t)$. The total expected loss is the sum of these expectations (recall, linearity of expectation).

What would be a candidate randomized algorithm? Note that in the WEIGHTED MAJORITY if 50.001% of the weight voted for 0 and 49.999% voted for 1, the algorithm went with 0. Even with this slim lead. The “fairer” randomized algorithm that suggests itself is that we should say 0 with probability 50.001% and 1 with the remainder probability. Turns out, this algorithm leads us arbitrarily close to the best expert!

A more convenient way of looking at the RANDOMIZED WEIGHTED MAJORITY algorithm is by thinking of picking a random expert and going with their decision. That is, given the weights we pick an expert i proportional to $w_i(t)$, and then predict $a(t) = e_i(t)$. Observe that this is the same fair algorithm described above.

```

1: procedure RANDOMIZED WEIGHTED MAJORITY( $\eta$ ): $\triangleright \eta \in (0, 1/2)$ .
2:   Maintain weights  $w_i(\cdot)$  for each expert with  $w_i(1) = 1$  for all  $i$ .
3:   for Days  $t = 1$  to  $T$  do:
4:     Select expert  $i$  with probability  $\propto w_i(t)$ .
5:     Predict  $a(t) = e_i(t)$ .
6:     Receive  $r(t)$  and incur loss  $\ell(t) = 1$  if  $a(t) \neq r(t)$ , and  $\ell(t) = 0$  otherwise.
7:     for every  $i$  with  $e_i(t) \neq r(t)$  do:
8:       Set  $w_i(t + 1) = w_i(t) \cdot (1 - \eta)$  where  $\eta \in (0, 1/2)$ .

```

Theorem 3. If there is an expert making at most k mistakes, then the *expected* number of mistakes made by the RANDOMIZED WEIGHTED MAJORITY(η) algorithm is

$$\mathbf{Exp}[\text{loss}] \leq (1 + \eta) \cdot k + \frac{\ln n}{\eta}$$

Proof. The proof idea is similar to the deterministic proof at a high level. As before, let $Z(t) := \sum_{i=1}^n w_i(t)$ be the sum of weights of all the experts at any time. Note we have (3) almost as is; since

there is an expert which makes at most k mistakes, their weight at the end would be $\geq (1 - \eta)^k$. Since all weights are non-negative, we get

$$Z(T) \geq (1 - \eta)^k \geq e^{-k(\eta + \eta^2)} \quad (4)$$

where we have used our old friend $(1 - x) \geq e^{-(x+x^2)}$ for all $x \in (0, 0.5)$.

Let loss_t be the random variable indicating whether $a(t) \neq r(t)$, that is, whether we make a mistake on the t th day or not. Since $a(t) = e_i(t)$ if i is selected on day t , we get

$$\mathbf{Exp}[\text{loss}_t] = \sum_{i:r(t) \neq e_i(t)} \Pr[i \text{ selected on day } t] = \frac{1}{Z(t)} \sum_{i:r(t) \neq e_i(t)} w_i(t) \quad (5)$$

Now we relate the total weights on day $t + 1$ to those in day t as

$$\begin{aligned} Z(t+1) &= \sum_{i:e_i(t) \neq r(t)} w_i(t+1) + \sum_{i:e_i(t) = r(t)} w_i(t+1) \\ &= (1 - \eta) \sum_{i:e_i(t) \neq r(t)} w_i(t) + \sum_{i:e_i(t) = r(t)} w_i(t) \\ &= Z(t) - \eta \sum_{i:r(t) \neq e_i(t)} w_i(t) \stackrel{(5)}{=} Z(t) - \eta Z(t) \mathbf{Exp}[\text{loss}_t] \end{aligned}$$

This gives

$$Z(t+1) \leq Z(t) \cdot (1 - \eta \mathbf{Exp}[\text{loss}_t]) \leq Z(t) \cdot e^{-(\eta \mathbf{Exp}[\text{loss}_t])}$$

We have used $1 + x \leq e^x$ above. Why? Because it allows us to “telescope”. We get

$$Z(T) \leq Z(1) \cdot e^{(-\eta \sum_{t=1}^T \mathbf{Exp}[\text{loss}_t])} = n \cdot e^{(-\eta \cdot \mathbf{Exp}[\text{loss}])} \quad (6)$$

Taking natural logs on (6) and (4), gives

$$-(\eta + \eta^2)k \leq \ln n - \eta \mathbf{Exp}[\text{loss}]$$

The theorem follows by “moving stuff around”. □

- **The Notion of Regret.** The experts problem is an example of an *online learning* or *online optimization* problem, where we are designing algorithms which is competitive with unknowns in the future. In many such algorithms, one compares the performance of the algorithm and considers the *difference* with the best possible single strategy. In this case, this would be $\mathbf{Exp}[\text{loss}]$, which is the performance of our algorithm, minus $\min_i \text{loss}_i$, the best experts loss. This term is called the *regret* of the algorithm: what more it could have made if it were just simple-minded and stuck with the (unknown) best expert. $\mathbf{Reg}_T(\mathcal{A})$ is this difference after playing T rounds of this game. One desirable feature one shoots for is *sublinear* regret — $\mathbf{Reg}_T(\mathcal{A}) = o(T)$ for large T . That is, if we divide by T , then the “average difference” with the best strategy goes to 0 as we play long enough.

Does [Theorem 3](#) imply sublinear regret? Well we see that $\mathbf{Exp}[\text{loss}] - k \leq \eta k + \frac{\ln n}{\eta}$. Note that k , the best expert, could be wrong a constant fraction of the times. Thus, ηk could be $\Omega(T)$. But if we chose $\eta = \frac{1}{\sqrt{T}}$, then we do see that both terms become $o(T)$. Indeed, choosing $\eta = \sqrt{\frac{\ln n}{T}}$ matches the two terms. So, if we knew that we were going to play for T rounds, then this indeed tells us what η to choose. This η is often called the “learning rate” or the “step size” in such algorithms.

- **Application: Minimax Theorem in Zero-Sum Games.** Consider a game between two players: the “row” player and the “column” player. Each player has n strategies and there is an $n \times n$ “pay-off” matrix A . For simplicity, let’s assume A ’s entries are in $\{0, 1\}$. If row player “plays” strategy $i \in [n]$ and column player plays strategy $j \in [n]$, then let’s say both players get a payoff of A_{ij} . The row player wants to *minimize* this payoff while the column player is trying to *maximize* it. More generally, if row player plays first and plays $i \in [n]$, col player should play $\arg \max_j A_{ij}$ and if col player plays first and play $j \in [n]$, then the row player should play $\arg \min_i A_{ij}$.

It should be clear that the player who plays second gets an advantage. For instance, let’s take the simple 2×2 identity matrix. If the row player plays first, then no matter which row she chooses, the column player can pick a column leading to a pay-off of 1. Similarly, if the column player plays first, then no matter which column he chooses, the row player can lead to a payoff of 0. And in general

$$\min_i \max_j A_{ij} \geq \max_j \min_i A_{ij}$$

The *minimax* theorem attributed to John von Neumann states that if players are allowed to *randomize*, or in other words can play *distribution over strategies*, then it doesn’t matter who plays first! More precisely, suppose \mathbf{x} is a *distribution* on $\{1, 2, \dots, n\}$; that is, $\mathbf{x} \in \mathbb{R}^n$, $x_i \geq 0$ for all $1 \leq i \leq n$, and $\sum_{i=1}^n x_i = 1$. When the row player plays \mathbf{x} , she is saying that she is going to sample a strategy i with probability x_i and play it. The column player needs to decide which column to play *only* with the knowledge of \mathbf{x} and not the actual realization. So if row player plays \mathbf{x} , then the column player should play the distribution \mathbf{y}

$$\arg \max_{\mathbf{y}} \mathbf{Exp}_{i \sim \mathbf{x}, j \sim \mathbf{y}} [A_{ij}] = \arg \max_{\mathbf{y}} \langle \mathbf{x}^\top A, \mathbf{y} \rangle$$

Exercise: Show that the maximizing \mathbf{y} actually is a “pure” strategy and in fact will be of the form \mathbf{e}_j which is 0 everywhere except some $j \in [n]$.

Similarly, if the col player play \mathbf{y} , then the row should play

$$\arg \min_{\mathbf{x}} \mathbf{Exp}_{i \sim \mathbf{x}, j \sim \mathbf{y}} [A_{ij}] = \arg \min_{\mathbf{x}} \langle \mathbf{x}, A\mathbf{y} \rangle$$

Theorem 4. For any matrix A , we have

$$\min_{\mathbf{x}} \max_{\mathbf{y}} \langle \mathbf{x}^\top A, \mathbf{y} \rangle = \max_{\mathbf{y}} \min_{\mathbf{x}} \langle \mathbf{x}, A\mathbf{y} \rangle$$

We prove this for Boolean matrices using [Theorem 3](#) above, but the theorem is true for any matrix and the proof also generalizes. Let \mathbf{x}^* be the \mathbf{x} achieving the minimum in the LHS and let \mathbf{y}^* be the \mathbf{y} achieving the maximum in the RHS. As a shorthand, let us use $f(\mathbf{x}, \mathbf{y})$ to be $\langle \mathbf{x}^\top A, \mathbf{y} \rangle = \langle \mathbf{x}, A\mathbf{y} \rangle$ and apart from brevity this will clearly show where the “bi-linearity” is used.

It’s easy to see $\text{lhs} \geq \text{rhs}$ just by definitions

$$\text{lhs} = \max_{\mathbf{y}} f(\mathbf{x}^*, \mathbf{y}) \geq f(\mathbf{x}^*, \mathbf{y}^*) \geq \min_{\mathbf{x}} f(\mathbf{x}, \mathbf{y}^*) = \text{rhs}$$

The non-trivial bit is to show $\text{lhs} \leq \text{rhs}$. This is done using the experts framework as follows.

- *Casting as experts problem.* Let us consider ourselves to be the row-player, and the n different strategies are our n experts. “Nature” will, over time, feed us columns of the matrix and this will define the loss functions. Recall, on day t , we maintain a *distribution* $\mathbf{x}^{(t)}$ over the experts (proportional to their weights). Upon seeing this, “nature” finds

$$\mathbf{y}^{(t)} := \arg \max_{\mathbf{y}} f(\mathbf{x}^{(t)}, \mathbf{y})$$

that is, it finds the *best-response* to our current mixed strategy $\mathbf{x}^{(t)}$. By the exercise above, this is going to be a single column $\mathbf{e}_{j^{(t)}}$. We define our *loss* and the experts loss to be

$$\text{loss}(t) = f(\mathbf{x}^{(t)}, \mathbf{y}^{(t)}) \quad \forall 1 \leq i \leq n, \quad \text{loss}_i(t) = f(\mathbf{e}_i, \mathbf{y}^{(t)})$$

A simple but crucial observation:

$$f(\mathbf{x}^{(t)}, \mathbf{y}^{(t)}) = \max_{\mathbf{y}} f(\mathbf{x}^{(t)}, \mathbf{y}) \geq \min_{\mathbf{x}} \max_{\mathbf{y}} f(\mathbf{x}, \mathbf{y}) = \text{lhs} \quad (\star)$$

Ok, now let’s play this above game for T days using the Multiplicative Weight Update algorithm. By [Theorem 3](#) we get that

$$\frac{1}{T} \sum_{t=1}^T f(\mathbf{x}^{(t)}, \mathbf{y}^{(t)}) \leq \min_i \frac{1}{T} \sum_{t=1}^T f(\mathbf{e}_i, \mathbf{y}^{(t)}) + \underbrace{\frac{\mathbf{Reg}(T)}{T}}_{\text{call this err}(T)} \quad (7)$$

Now we will use the fact that $f(\mathbf{x}, \mathbf{y})$ is a “bilinear function”. In particular, we observe

$$\frac{1}{T} \sum_{t=1}^T f(\mathbf{e}_i, \mathbf{y}^{(t)}) = \frac{1}{T} \sum_{t=1}^T \langle \mathbf{e}_i^\top A, \mathbf{y}^{(t)} \rangle = \langle \mathbf{e}_i^\top A, \underbrace{\frac{1}{T} \sum_{t=1}^T \mathbf{y}^{(t)}}_{\text{call this } \hat{\mathbf{y}}_T} \rangle = f(\mathbf{e}_i, \hat{\mathbf{y}}_T)$$

Crucially, note that $\hat{\mathbf{y}}_T$ is *also* a distribution over columns (it’s an average of distributions). Putting thing together, and using (\star) , we get

$$\text{lhs} \leq \frac{1}{T} \sum_{t=1}^T f(\mathbf{x}^{(t)}, \mathbf{y}^{(t)}) \leq \min_i f(\mathbf{e}_i, \hat{\mathbf{y}}_T) + \text{err}(T) \leq \max_{\mathbf{y}} \min_i f(\mathbf{e}_i, \mathbf{y}) + \text{err}(T) = \text{rhs} + \text{err}(T)$$

Just to be clear, the first inequality follows from (\star) , the second is (7) , the third is definition, and the equality follows because once the column strategy is fixed, the best strategy for row is to pick a single row...I guess this uses the bilinear form of f as well.

Now take $T \rightarrow \infty$; since $\mathbf{Reg}(T) = o(T)$, $\text{err}(T) \rightarrow 0$ proving $\text{lhs} \leq \text{rhs}$.