

# The Multiplicative Weights Update Method: A Meta-Algorithm and its Applications

Sanjeev Arora\*  
Princeton University  
Princeton NJ 08540  
arora@cs.princeton.edu

Elad Hazan  
Technion - Israel Inst. of Tech.  
Haifa, Israel  
ehazan@ie.technion.ac.il

Satyen Kale  
Yahoo! Research  
Santa Clara, CA 95054  
satyen.kale@gmail.com

## Abstract

Algorithms in varied fields use the idea of maintaining a distribution over a certain set and use the *multiplicative update rule* to iteratively change these weights. Their analyses are usually very similar and rely on an exponential potential function.

In this survey we present a simple meta-algorithm that unifies many of these disparate algorithms and derives them as simple instantiations of the meta-algorithm. We feel that since this meta-algorithm and its analysis are so simple, and its applications so broad, it should be a standard part of algorithms courses, like “divide and conquer.”

## 1 Introduction

The *Multiplicative Weights (MW) method* is a simple idea which has been repeatedly discovered in fields as diverse as Machine Learning, Optimization, and Game Theory. The setting for this algorithm is the following. A decision maker has a choice of  $n$  decisions, and needs to repeatedly make a decision and obtain an associated payoff. The decision maker’s goal, in the long run, is to achieve a total payoff which is comparable to the payoff of that fixed decision that maximizes the total payoff with the benefit of hindsight. While this best decision may not be known *a priori*, it is still possible to achieve this goal by maintaining weights on the decisions, and choosing the decisions randomly with probability proportional to the weights. In each successive round, the weights are updated by multiplying them with factors which depend on the payoff of the associated decision in that round. Intuitively, this scheme works because it tends to focus higher weight on higher payoff decisions in the long run.

This idea lies at the core of a variety of algorithms. Some examples include: the Ada Boost algorithm in machine learning [FS97]; algorithms for game playing studied in economics (see references later), the Plotkin-Shmoys-Tardos algorithm for packing and covering LPs [PST91], and its improvements in the case of flow problems by Young [You95], Garg-Könemann [GK98], Fleischer [Fle00] and others; methods for convex optimization like exponentiated gradient (mirror descent), Lagrangian multipliers, and subgradient methods,

---

\*This project was supported by David and Lucile Packard Fellowship and NSF grants MSPA-MCS 0528414 and CCR-0205594.

Impagliazzo’s proof of the Yao XOR lemma [Imp95], etc. The analysis of the running time uses a potential function argument and the final running time is proportional to  $1/\varepsilon^2$ .

It has been clear to several researchers that these results are very similar. For example Khandekar’s PhD thesis [Kha04] makes this point about the varied applications of this idea to convex optimization. The purpose of this survey is to clarify that many of these applications are instances of the same, more general algorithm (although several specialized applications, such as [LN93], require additional technical work). This meta-algorithm is very similar to the “Hedge” algorithm from learning theory [FS97]. Similar algorithms have been independently rediscovered in many other fields; see below. The advantage of deriving the above algorithms from the same meta algorithm is that this highlights their commonalities as well as their differences. To give an example, the algorithms of Garg-Könemann [GK98] were felt to be quite different from those of Plotkin-Shmoys-Tardos [PST91]. In our framework, they can be seen as a clever trick for “width reduction” for the Plotkin-Shmoys-Tardos algorithms (see Section 3.4).

We feel that this meta-algorithm and its analysis are simple and useful enough that they should be viewed as a basic tool taught to all algorithms students together with divide-and-conquer, dynamic programming, random sampling, and the like. Note that the multiplicative weights update rule may be seen as a “constructive” version of LP duality—equivalently, von Neumann’s minimax theorem in game theory—and it gives a fairly concrete method for competing players to arrive at a solution/equilibrium (see Section 3.2). This may be an appealing feature in introductory algorithms courses, since the standard algorithms for LP such as simplex, ellipsoid, or interior point lack such a game-theoretic interpretation. Furthermore, it is a convenient stepping point to many other topics that rarely get mentioned in algorithms courses, including online algorithms (see the basic scenario in Section 1.1) and machine learning. Furthermore our proofs seem easier and cleaner than the entropy-based proofs for the same results in machine learning.

The current paper is chiefly a survey. It introduces the main algorithm, gives a few variants (mostly having to do with the range in which the payoffs lie), and surveys the most important applications—often with complete proofs. Note however that this survey does not cover *all* applications of the technique, as several of these require considerable additional technical work which is beyond the scope of this paper. We have provided pointers to some such applications which use the multiplicative weights technique at their core without going into more details. There are also a few small results that appear to be new, such as the variant of the Garg-Könemann algorithm in Section 3.4 and the lower bound in Section 4.

**Related work.** An algorithm similar in flavor to the Multiplicative Weights algorithm were proposed in game theory in the early 1950’s [BvN50, Bro51, Rob51]. Following Brown [Bro51], this algorithm was called “Fictitious Play”: at each step each player observes actions taken by his opponent in previous stages, updates his beliefs about his opponents’ strategies, and chooses the pure best response against these beliefs. In the simplest case, the player simply assumes that the opponent is playing from a stationary distribution and sets his current belief of the opponent’s distribution to be the empirical frequency of the strategies played by the opponent. This simple idea (which was shown to lead to optimal solutions in the limit in various cases) led to numerous developments in economics, including Arrow-Debreu General Equilibrium theory and more recently, evolutionary game theory.

Grigoriadis and Khachiyan [GK95] showed how a randomized variant of “Fictitious Play” can solve two player zero-sum games efficiently. This algorithm is precisely the multiplicative weights algorithm. It can be viewed as a soft version of fictitious play, when the player gives higher weight to the strategies which pay off better, and chooses her strategy using these weights rather than choosing the best response strategy.

In Machine Learning, the earliest form of the multiplicative weights update rule was used by Littlestone in his well-known Winnow algorithm [Lit87, Lit89]. It is somewhat reminiscent of the older *perceptron* learning algorithm of Minsky and Papert [MP69]. The Winnow algorithm was generalized by Littlestone and Warmuth [LW94] in the form of the Weighted Majority algorithm, and later by Freund and Schapire in the form of the Hedge algorithm [FS97]. We note that most relevant papers in learning theory use an analysis that relies on entropy (or its cousin, Kullback-Leibler divergence) calculations. This analysis is closely related to our analysis, but we use *exponential* functions instead of the *logarithm*, or entropy, used in those papers. The underlying calculation is the same: whereas we repeatedly use the fact that  $e^x \approx 1 + x$  when  $|x|$  is small, they use the fact that  $\ln(1 + x) \approx x$ . We feel that our approach is cleaner (although the entropy based approach yields somewhat tighter bounds that are useful in some applications, see Section 2.2).

The multiplicative update rule (and the exponential potential function) was also discovered in Computational Geometry in the late 1980s [Cla88, CW89] and several applications in geometry are described in Chazelle [Cha00] (p. 6, and p. 124). See also our Section 3.10.3, which also mentions some more recent applications to geometric embeddings of finite metric spaces.

The weighted majority algorithm as well as more sophisticated versions have been independently discovered in operations research and statistical decision making in the context of the *On-line decision problem*; see the surveys of Cover [Cov96], Foster and Vohra [FV99], and also Blum [Blu98] who includes applications of weighted majority to machine learning. A notable algorithm, which is different from but related to our framework, was developed by Hannan in the 1950’s [Han57]. Kalai and Vempala showed how to derive efficient algorithms via methods similar to Hannan’s [KV03]. We show how Hannan’s algorithm with the appropriate choice of parameters yields the multiplicative update decision rule in section 3.8.

Within computer science, several researchers have previously noted the close relationships between multiplicative update algorithms used in different contexts. Young [You95] notes the connection between fast LP algorithms and Raghavan’s method of pessimistic estimators for derandomization of randomized rounding algorithms; see our Section 3.5. Klivans and Servedio [KS03] relate boosting algorithms in learning theory to proofs of Yao’s XOR Lemma; see our Section 3.6. Garg and Khandekar [GK04] describe a common framework for convex optimization problems that contains Garg-Könemann and Plotkin-Shmoys-Tardos as subcases.

To the best of our knowledge our framework is the most general and arguably, the simplest. We readily acknowledge the influence of all previous papers (especially Young [You95] and Freund-Schapire [FS99]) on the development of our framework. We emphasize again that we do not claim that *every* algorithm designed using the multiplicative update idea fits in our framework, just that most do.

**Paper organization.** We proceed to define the illustrative weighted majority algorithm in this section. In section 2 we describe the general MW meta-algorithm, followed by numerous and varied applications in section 3. In section 4 we give lower bounds, followed by the more general matrix MW algorithm in section 5.

## 1.1 The weighted majority algorithm

Now we briefly illustrate the weighted majority algorithm in a simple and concrete setting, which will naturally lead to our generalized meta-algorithm. This is known as the *Prediction from Expert Advice* problem.

Imagine the process of picking good times to invest in a stock. For simplicity, assume that there is a single stock of interest, and its daily price movement is modeled as a sequence of binary events: up/down. (Below, this will be generalized to allow non-binary events.) Each morning we try to predict whether the price will go up or down that day; if our prediction happens to be wrong we lose a dollar that day, and if it's correct, we lose nothing.

The stock movements can be *arbitrary* and even *adversarial*. To balance out this pessimistic assumption, we assume that while making our predictions, we are allowed to watch the predictions of  $n$  “experts”. These experts could be arbitrarily correlated, and they may or may not know what they are talking about. The algorithm’s goal is to limit its cumulative losses (i.e., bad predictions) to roughly the same as the *best* of these experts. At first sight this seems an impossible goal, since it is not known until the end of the sequence who the best expert was, whereas the algorithm is required to make predictions all along.

Indeed, the first algorithm one thinks of is to compute each day’s up/down prediction by going with the majority opinion among the experts that day. But, this algorithm doesn’t work because a majority of experts may be consistently wrong on every single day.

The weighted majority algorithm corrects the trivial algorithm. It maintains a *weighting* of the experts. Initially all have equal weight. As time goes on, some experts are seen as making better predictions than others, and the algorithm increases their weight proportionately. The algorithm’s prediction of up/down for each day is computed by going with the opinion of the weighted majority of the experts for that day.

### Weighted majority algorithm

**Initialization:** Fix an  $\eta \leq \frac{1}{2}$ . For each expert  $i$ , associate the weight  $w_i^{(1)} := 1$ .

**For**  $t = 1, 2, \dots, T$ :

1. Make the prediction that is the weighted majority of the experts’ predictions based on the weights  $w_1^{(t)}, \dots, w_n^{(t)}$ . That is, predict “up” or “down” depending on which prediction has a higher total weight of experts advising it (breaking ties arbitrarily).
2. For every expert  $i$  who predicts wrongly, decrease his weight for the next round by multiplying it by a factor of  $(1 - \eta)$ :

$$w_i^{(t+1)} = (1 - \eta)w_i^{(t)} \quad (\text{update rule}). \quad (1)$$

**Theorem 1** *After  $T$  steps, let  $m_i^{(T)}$  be the number of mistakes of expert  $i$  and  $M^{(T)}$  be the number of mistakes our algorithm has made. Then we have the following bound for every  $i$ :*

$$M^{(T)} \leq 2(1 + \eta)m_i^{(T)} + \frac{2 \ln n}{\eta}.$$

*In particular, this holds for  $i$  which is the best expert, i.e. having the least  $m_i^{(T)}$ .*

*Remark:* When  $m_i^{(T)} \gg \frac{2 \ln n}{\eta}$  we see that the number of mistakes made by the algorithm is bounded from above by roughly  $2(1 + \eta)m_i^{(T)}$ , i.e., approximately twice the number of mistakes made by the best expert. This is tight for any deterministic algorithm. However, the factor of 2 can be removed by substituting the above deterministic algorithm by a randomized algorithm that predicts according to the majority opinion with probability proportional to its weight. (In other words, if the total weight of the experts saying “up” is  $3/4$  then the algorithm predicts “up” with probability  $3/4$  and “down” with probability  $1/4$ .) Then the number of mistakes after  $T$  steps is a random variable and the claimed upper bound holds for its *expectation* (see Section 2 for more details).

PROOF: A simple induction shows that  $w_i^{(t+1)} = (1 - \eta)^{m_i^{(t)}}$ . Let  $\Phi^{(t)} = \sum_i w_i^{(t)}$  (“the potential function”). Thus  $\Phi^{(1)} = n$ . Each time we make a mistake, the weighted majority of experts also made a mistake, so at least half the total weight decreases by a factor  $1 - \eta$ . Thus, the potential function decreases by a factor of at least  $(1 - \eta/2)$ :

$$\Phi^{(t+1)} \leq \Phi^{(t)} \left( \frac{1}{2} + \frac{1}{2}(1 - \eta) \right) = \Phi^{(t)}(1 - \eta/2).$$

Thus simple induction gives  $\Phi^{(T+1)} \leq n(1 - \eta/2)^{M^{(T)}}$ . Finally, since  $\Phi^{(T+1)} \geq w_i^{(T+1)}$  for all  $i$ , the claimed bound follows by comparing the above two expressions and using the fact that  $-\ln(1 - \eta) \leq \eta + \eta^2$  since  $\eta < \frac{1}{2}$ .  $\square$

The beauty of this analysis is that it makes no assumption about the sequence of events: they could be arbitrarily correlated and could even depend upon our current weighting of the experts. In this sense, this algorithm delivers more than initially promised, and this lies at the root of why (after obvious generalization) it can give rise to the diverse algorithms mentioned earlier. In particular, the scenario where the events are chosen adversarially resembles a zero-sum game, which we consider later in section 3.2.

## 2 The Multiplicative Weights algorithm

In the general setting, we have a set of  $n$  *decisions* and in each round, we are required to select one decision from the set. In each round, each decision incurs a certain *cost*, determined by nature or an adversary. All the costs are revealed *after* we choose our decision, and we incur the cost of the decision we chose. For example, in the prediction from expert advice problem, each decision corresponds to a choice of an expert, and the cost of an expert is 1 if the expert makes a mistake, and 0 otherwise.

To motivate the Multiplicative Weights (MW) algorithm, consider the naïve strategy that, in each iteration, simply picks a decision at random. The expected penalty will be

that of the “average” decision. Suppose now that a few decisions are clearly better in the long run. This is easy to spot as the costs are revealed over time, and so it is sensible to reward them by increasing their probability of being picked in the next round (hence the multiplicative weight update rule).

Intuitively, being in complete ignorance about the decisions at the outset, we select them uniformly at random. This maximum entropy starting rule reflects our ignorance. As we learn which ones are the good decisions and which ones are bad, we lower the entropy to reflect our increased knowledge. The multiplicative weight update is our means of skewing the distribution.

We now set up some notation. Let  $t = 1, 2, \dots, T$  denote the current round, and let  $i$  be a generic decision. In each round  $t$ , we select a distribution  $\mathbf{p}^{(t)}$  over the set of decisions, and select a decision  $i$  randomly from it. At this point, the costs of all the decisions are revealed by nature in the form of the vector  $\mathbf{m}^{(t)}$  such that decision  $i$  incurs cost  $m_i^{(t)}$ . We assume that the costs lie in the range  $[-1, 1]$ . This is the only assumption we make on the costs; nature is completely free to choose the cost vector as long as these bounds are respected, even with full knowledge of the distribution that we choose our decision from.

The expected cost to the algorithm for sampling a decision  $i$  from the distribution  $\mathbf{p}^{(t)}$  is

$$\mathbf{E}_{i \in \mathbf{p}^{(t)}}[m_i^{(t)}] = \mathbf{m}^{(t)} \cdot \mathbf{p}^{(t)}.$$

The total expected cost over all rounds is therefore  $\sum_{t=1}^T \mathbf{m}^{(t)} \cdot \mathbf{p}^{(t)}$ . Just as before, our goal is to design an algorithm which achieves a total expected cost not too much more than the cost of the best decision in hindsight, viz.  $\min_i \sum_{t=1}^T m_i^{(t)}$ .

### Multiplicative Weights algorithm

**Initialization:** Fix an  $\eta \leq \frac{1}{2}$ . For each decision  $i$ , associate the weight  $w_i^{(t)} := 1$ .

**For**  $t = 1, 2, \dots, T$ :

1. Choose decision  $i$  with probability proportional to its weight  $w_i^{(t)}$ . I.e., use the distribution over decisions  $\mathbf{p}^{(t)} = \{w_1^{(t)}/\Phi^{(t)}, \dots, w_n^{(t)}/\Phi^{(t)}\}$  where  $\Phi^{(t)} = \sum_i w_i^{(t)}$ .
2. Observe the costs of the decisions  $\mathbf{m}^{(t)}$ .
3. Penalize the costly decisions by updating their weights as follows: for every decision  $i$ , set

$$w_i^{(t+1)} = w_i^{(t)}(1 - \eta m_i^{(t)}) \tag{2}$$

Figure 1: The Multiplicative Weights algorithm.

The following theorem —completely analogous to Theorem 1— bounds the total expected cost of the Multiplicative Weights algorithm (given in Figure 1) in terms of the total cost of the best decision:

**Theorem 2** *Assume that all costs  $m_i^{(t)} \in [-1, 1]$  and  $\eta \leq 1/2$ . Then the Multiplicative*

*Weights algorithm guarantees that after  $T$  rounds, for any decision  $i$ , we have*

$$\sum_{t=1}^T \mathbf{m}^{(t)} \cdot \mathbf{p}^{(t)} \leq \sum_{t=1}^T m_i^{(t)} + \eta \sum_{t=1}^T |m_i^{(t)}| + \frac{\ln n}{\eta}.$$

PROOF: The proof is along the lines of the earlier one, using the potential function  $\Phi^{(t)} = \sum_i w_i^{(t)}$ :

$$\begin{aligned} \Phi^{(t+1)} &= \sum_i w_i^{(t+1)} \\ &= \sum_i w_i^{(t)} (1 - \eta m_i^{(t)}) \\ &= \Phi^{(t)} - \eta \Phi^{(t)} \sum_i m_i^{(t)} p_i^{(t)} \\ &= \Phi^{(t)} (1 - \eta \mathbf{m}^{(t)} \cdot \mathbf{p}^{(t)}) \\ &\leq \Phi^{(t)} \exp(-\eta \mathbf{m}^{(t)} \cdot \mathbf{p}^{(t)}). \end{aligned}$$

Here, we used the fact that  $p_i^{(t)} = w_i^{(t)} / \Phi^{(t)}$ . Thus, by induction, after  $T$  rounds, we have

$$\Phi^{(T+1)} \leq \Phi^{(1)} \exp \left( -\eta \sum_{t=1}^T \mathbf{m}^{(t)} \cdot \mathbf{p}^{(t)} \right) = n \cdot \exp \left( -\eta \sum_{t=1}^T \mathbf{m}^{(t)} \cdot \mathbf{p}^{(t)} \right). \quad (3)$$

Next we use the following facts, which follow immediately from the convexity of the exponential function:

$$\begin{aligned} (1 - \eta)^x &\leq (1 - \eta x) \quad \text{if } x \in [0, 1] \\ (1 + \eta)^{-x} &\leq (1 - \eta x) \quad \text{if } x \in [-1, 0] \end{aligned} \quad (4)$$

Since  $m_i^{(t)} \in [-1, 1]$ , we have for every decision  $i$ ,

$$\Phi^{(T+1)} \geq w_i^{(T+1)} = \prod_{t \leq T} (1 - \eta m_i^{(t)}) \geq (1 - \eta)^{\sum_{\geq 0} m_i^{(t)}} \cdot (1 + \eta)^{-\sum_{< 0} m_i^{(t)}}, \quad (5)$$

where the subscripts “ $\geq 0$ ” and “ $< 0$ ” in the summations refer to the rounds  $t$  where  $m_i^{(t)}$  is  $\geq 0$  and  $< 0$  respectively. Taking logarithms in equations (3) and (5) we get:

$$\ln n - \eta \sum_{t=1}^T \mathbf{m}^{(t)} \cdot \mathbf{p}^{(t)} \geq \sum_{\geq 0} m_i^{(t)} \ln(1 - \eta) - \sum_{< 0} m_i^{(t)} \ln(1 + \eta)$$

Negating, rearranging, and scaling by  $\frac{1}{\eta}$ :

$$\begin{aligned}
\sum_{t=1}^T \mathbf{m}^{(t)} \cdot \mathbf{p}^{(t)} &\leq \frac{\ln n}{\eta} + \frac{1}{\eta} \sum_{\geq 0} m_i^{(t)} \ln \frac{1}{1-\eta} + \frac{1}{\eta} \sum_{< 0} m_i^{(t)} \ln(1+\eta) \\
&\leq \frac{\ln n}{\eta} + \frac{1}{\eta} \sum_{\geq 0} m_i^{(t)} (\eta + \eta^2) + \frac{1}{\eta} \sum_{< 0} m_i^{(t)} (\eta - \eta^2) \\
&= \frac{\ln n}{\eta} + \sum_{t=1}^T m_i^{(t)} + \eta \sum_{\geq 0} m_i^{(t)} - \eta \sum_{< 0} m_i^{(t)} \\
&= \frac{\ln n}{\eta} + \sum_{t=1}^T m_i^{(t)} + \eta \sum_{t=1}^T |m_i^{(t)}|.
\end{aligned}$$

In the second inequality we used the facts that

$$\ln\left(\frac{1}{1-\eta}\right) \leq \eta + \eta^2 \quad \text{and} \quad \ln(1+\eta) \leq \eta - \eta^2 \quad (6)$$

for  $\eta \leq \frac{1}{2}$ .  $\square$

**Corollary 3** *The Multiplicative Weights algorithm also guarantees that after  $T$  rounds, for any distribution  $\mathbf{p}$  on the decisions,*

$$\sum_{t=1}^T \mathbf{m}^{(t)} \cdot \mathbf{p}^{(t)} \leq \sum_{t=1}^T (\mathbf{m}^{(t)} + \eta |\mathbf{m}^{(t)}|) \cdot \mathbf{p} + \frac{\ln n}{\eta},$$

where  $|\mathbf{m}^{(t)}|$  is the vector obtained by the taking the coordinate-wise absolute value of  $\mathbf{m}^{(t)}$ .

PROOF: This corollary follows immediately from Theorem 2, by taking a convex combination of the inequalities for all decisions  $i$  with the distribution  $\mathbf{p}$ .  $\square$

## 2.1 Updating with exponential factors: the Hedge algorithm

In our description of the MW algorithm, the update rule uses multiplication by a linear function of the cost (specifically,  $(1 - \eta m_i^{(t)})$  for expert  $i$ ). In several other incarnations of MW algorithm, notably the Hedge algorithm of Freund and Schapire [FS97], an exponential factor is used instead. This update rule is the following:

$$w_i^{(t+1)} = w_i^{(t)} \cdot \exp(-\eta m_i^{(t)}). \quad (7)$$

As can be seen from the analysis of the MW algorithm, Hedge is not very different. The bound we obtain for Hedge is slightly different however. While most of the applications we present in the rest of the paper can be derived using Hedge as well with a little extra calculation, some applications, such as the ones in Sections 3.3 and 3.5, explicitly need the MW algorithm rather than Hedge to obtain better bounds. Here, we state the bound



obtained for Hedge without proof – the analysis is on the same lines as before. The only difference is that instead of the inequalities (4), we use the inequality

$$\exp(-\eta x) \leq 1 - \eta x + \eta^2 x^2$$

if  $|\eta x| \leq 1$ .

**Theorem 4** *Assume that all costs  $m_i^{(t)} \in [-1, 1]$  and  $\eta \leq 1$ . Then the Hedge algorithm guarantees that after  $T$  rounds, for any decision  $i$ , we have*

$$\sum_{t=1}^T \mathbf{m}^{(t)} \cdot \mathbf{p}^{(t)} \leq \sum_{t=1}^T m_i^{(t)} + \eta \sum_{t=1}^T (\mathbf{m}^{(t)})^2 \cdot \mathbf{p}^{(t)} + \frac{\ln n}{\eta}.$$

Here,  $(\mathbf{m}^{(t)})^2$  is the vector obtained by taking coordinate-wise square of  $\mathbf{m}^{(t)}$ .

This guarantee is very similar to the one in Theorem 2, with one important difference: the term multiplying  $\eta$  is a loss which depends on the algorithm's distribution. In Theorem 2, this additional term depends on the loss of the best decision in hindsight. For some applications the latter guarantee is stronger (see Section 3.3).

## 2.2 Proof via KL-divergence

In this section, we give an alternative proof of Theorem 2 based on the Kullback-Leibler (KL) divergence, or relative entropy. While this proof is somewhat more complicated, it gives a good insight into why the MW algorithm works: the reason is that it tends to reduce the KL-divergence to the optimal solution. Another reason for giving this proof is that it yields a more nuanced form of the MW algorithm that is useful in some applications (such as the construction of hard-core sets, see Section 3.7). Readers may skip this section without loss of continuity.

For two distributions  $\mathbf{p}$  and  $\mathbf{q}$  on the decision set, the relative entropy between them is

$$\mathbf{RE}(\mathbf{p} \parallel \mathbf{q}) = \sum_i p_i \ln(p_i/q_i).$$

Consider the following twist on the basic decision-making problem from Section 2. Fix a convex subset of distributions over decisions,  $\mathcal{P}$  (note: the basic setting is recovered when  $\mathcal{P}$  is the set of *all* distributions). In each round  $t$ , the decision-maker is required to produce a distribution  $\mathbf{p}^{(t)} \in \mathcal{P}$ . At that point, the cost vector  $\mathbf{m}^{(t)}$  is revealed and the decision-maker suffers cost  $\mathbf{m}^{(t)} \cdot \mathbf{p}^{(t)}$ . Since we make the restriction that  $\mathbf{p}^{(t)} \in \mathcal{P}$ , we now want to compare the total cost of the decision-maker to the cost of the best fixed distribution in  $\mathcal{P}$ . Consider the following algorithm.

Note that in the special case when  $\mathcal{P}$  is the set of all distributions on the decisions, this algorithm is exactly the basic MW algorithm presented in Figure 1. The relative entropy projection step ensures that we always choose a distribution in  $\mathcal{P}$ . This projection is a convex program since relative entropy is convex and  $\mathcal{P}$  is a convex set, and hence can be computed using standard convex programming techniques.

We now prove a bound on the total cost of the algorithm (compare to Corollary 3). Note that in the basic setting when  $\mathcal{P}$  is the set of all distributions, the bound given below is tighter than the one in Theorem 2.

### Multiplicative Weights Update algorithm with Restricted Distributions

**Initialization:** Fix a  $\eta \leq \frac{1}{2}$ . Set  $\mathbf{p}^{(1)}$  to be an arbitrary distribution in  $\mathcal{P}$ . initialized to 1.

**For**  $t = 1, 2, \dots, T$ :

1. Choose decision  $i$  by sampling from  $\mathbf{p}^{(t)}$ .
2. Observe the costs of the decisions  $\mathbf{m}^{(t)}$ .
3. Compute the probability vector  $\hat{\mathbf{p}}^{(t+1)}$  using the usual multiplicative update rule: for every expert  $i$ ,

$$p_i^{(t+1)} = p_i^{(t)}(1 - \eta m_i^{(t)}) / \Phi^{(t)} \quad (8)$$

where  $\Phi^{(t)}$  is the normalization factor to make  $\hat{\mathbf{p}}^{(t+1)}$  a distribution.

4. Set  $\mathbf{p}^{(t+1)}$  to be the relative entropy projection of  $\hat{\mathbf{p}}^{(t)}$  on the set  $\mathcal{P}$ , i.e.

$$\mathbf{p}^{(t+1)} = \arg \min_{\mathbf{p} \in \mathcal{P}} \mathbf{RE}(\mathbf{p} \parallel \hat{\mathbf{p}}^{(t)}).$$

Figure 2: The Multiplicative Weights algorithm with Restricted Distributions

**Theorem 5** Assume that all costs  $m_i^{(t)} \in [-1, 1]$  and  $\eta \leq 1/2$ . Then the Multiplicative Weights algorithm with Restricted Distributions guarantees that after  $T$  rounds, for any  $\mathbf{p} \in \mathcal{P}$ , we have

$$\sum_{t=1}^T \mathbf{m}^{(t)} \cdot \mathbf{p}^{(t)} \leq \sum_{t=1}^T (\mathbf{m}^{(t)} + \eta |\mathbf{m}^{(t)}|) \cdot \mathbf{p} + \frac{\mathbf{RE}(\mathbf{p} \parallel \mathbf{p}^{(1)})}{\eta},$$

where  $|\mathbf{m}^{(t)}|$  is the vector obtained by taking the coordinate-wise absolute value of  $\mathbf{m}^{(t)}$ .

PROOF: We use the relative entropy between  $\mathbf{p}$  and  $\mathbf{p}^{(t)}$ ,  $\mathbf{RE}(\mathbf{p} \parallel \mathbf{p}^{(t)}) := \sum_i p_i \ln(p_i/p_i^{(t)})$  as a “potential” function. We have

$$\begin{aligned} \mathbf{RE}(\mathbf{p} \parallel \hat{\mathbf{p}}^{t+1}) - \mathbf{RE}(\mathbf{p} \parallel \mathbf{p}^{(t)}) &= \sum_i p_i \ln \frac{p_i^{(t)}}{\hat{p}_i^{(t+1)}} \\ &= \sum_i p_i \ln \frac{\Phi^{(t)}}{(1 - \eta m_i^{(t)})} \\ &\leq \ln \frac{1}{1 - \eta} \sum_{\geq 0} p_i m_i^{(t)} + \ln(1 + \eta) \sum_{< 0} p_i m_i^{(t)} + \ln \Phi^{(t)} \\ &\leq \eta (\mathbf{m}^{(t)} + \eta |\mathbf{m}^{(t)}|) \cdot \mathbf{p} + \ln \Phi^{(t)}. \end{aligned}$$

The first inequality above follows from (4), and the second follows from (6). Next, we have

$$\ln \Phi^{(t)} = \ln \left[ \sum_i p_i^{(t)} (1 - \eta m_i^{(t)}) \right] = \ln [1 - \eta \mathbf{m}^{(t)} \cdot \mathbf{p}^{(t)}] \leq -\eta \mathbf{m}^{(t)} \cdot \mathbf{p}^{(t)}$$

since  $\ln(1-x) \leq -x$  for  $x < 1$ . Thus, we get

$$\mathbf{RE}(\mathbf{p} \parallel \hat{\mathbf{p}}^{t+1}) - \mathbf{RE}(\mathbf{p} \parallel \mathbf{p}^{(t)}) \leq \eta(\mathbf{m}^{(t)} + \eta|\mathbf{m}^{(t)}|) \cdot \mathbf{p} - \eta(\mathbf{m}^{(t)} \cdot \mathbf{p}^{(t)}).$$

This inequality essentially says that if the cost of the algorithm in round  $t$ ,  $\mathbf{m}^{(t)} \cdot \mathbf{p}^{(t)}$ , is significantly larger than the cost of the comparator,  $\mathbf{m}^{(t)} \cdot \mathbf{p}$ , then  $\hat{\mathbf{p}}^{t+1}$  moves closer to  $\mathbf{p}$  (in the relative entropy distance) than  $\mathbf{p}^{(t)}$ .

Now, projection on the set  $\mathcal{P}$  using the relative entropy as a distance function is a Bregman projection, and thus it satisfies the following Generalized Pythagorean inequality (see, e.g. [HW98]), for any  $\mathbf{p} \in \mathcal{P}$ :

$$\mathbf{RE}(\mathbf{p} \parallel \mathbf{p}^{(t+1)}) + \mathbf{RE}(\mathbf{p}^{(t+1)} \parallel \hat{\mathbf{p}}^{(t+1)}) \leq \mathbf{RE}(\mathbf{p} \parallel \hat{\mathbf{p}}^{(t+1)}).$$

I.e., the projection step only brings the distribution closer to  $\mathbf{p}$ . Since relative entropy is always non-negative, we have  $\mathbf{RE}(\mathbf{p}^{(t+1)} \parallel \hat{\mathbf{p}}^{(t+1)}) \geq 0$  and so

$$\mathbf{RE}(\mathbf{p} \parallel \mathbf{p}^{(t+1)}) - \mathbf{RE}(\mathbf{p} \parallel \mathbf{p}^{(t)}) \leq \eta(\mathbf{m}^{(t)} + \eta|\mathbf{m}^{(t)}|) \cdot \mathbf{p} - \eta(\mathbf{m}^{(t)} \cdot \mathbf{p}^{(t)}).$$

Summing up from  $t = 1$  to  $T$ , dividing by  $\eta$ , and simplifying using the fact that  $\mathbf{RE}(\mathbf{p} \parallel \mathbf{p}^{(T+1)})$  is non-negative, we get the stated bound.  $\square$

### 2.3 Gains instead of losses

There are situations where it makes more sense for the vector  $\mathbf{m}^{(t)}$  to specify *gains* for each expert rather than losses. Now our goal is to get as much total expected payoff as possible in comparison to the total payoff of the best expert. We can get an algorithm for this case simply by running the Multiplicative Weights algorithm using the cost vector  $-\mathbf{m}^{(t)}$ .

The resulting algorithm is identical, and the following theorem follows directly from Theorem 2 by simply negating the quantities:

**Theorem 6** *Assume that all gains  $m_i^{(t)} \in [-1, 1]$  and  $\eta \leq 1$ . Then the Multiplicative Weights algorithm (for gains) guarantees that after  $T$  rounds, for any expert  $i$ , we have*

$$\sum_{t=1}^T \mathbf{m}^{(t)} \cdot \mathbf{p}^{(t)} \geq \sum_{t=1}^T m_i^{(t)} - \eta \sum_{t=1}^T |m_i^{(t)}| - \frac{\ln n}{\eta}.$$

We also have the following immediate corollary, corresponding to Corollary 3:

**Corollary 7** *The Multiplicative Weights algorithm also guarantees that after  $T$  rounds, for any distribution  $\mathbf{p}$  on the decisions,*

$$\sum_{t=1}^T \mathbf{m}^{(t)} \cdot \mathbf{p}^{(t)} \geq \sum_{t=1}^T (\mathbf{m}^{(t)} - \eta|\mathbf{m}^{(t)}|) \cdot \mathbf{p} - \frac{\ln n}{\eta},$$

where  $|\mathbf{m}^{(t)}|$  is the vector obtained by the taking the coordinate-wise absolute value of  $\mathbf{m}^{(t)}$ .

### 3 Applications

Typically, the Multiplicative Weights method is applied in the following manner. A prototypical example is to solve a constrained optimization problem. We then let a decision represent each constraint in the problem, with costs specified by the points in the domain of interest. For a given point, the cost of a decision is made proportional to *how well* the corresponding constraint is satisfied on the point. This might seem counterintuitive, but recall that we *reduce* a decision’s weight depending on its penalty, and if a constraint is well satisfied on points so far we would like its weight to be smaller, so that the algorithm focuses on constraints that are poorly satisfied.

In many applications (though not all) the choice of point is also under our control. Typically we will need to generate the *maximally adversarial* point, i.e. the point that maximizes the expected cost. Then the overall algorithm consists of two subprocedures: an “oracle” for generating the maximally adversarial point at each step, and the MW algorithm for updating the weights of the decisions. With this intuition, we can describe the following applications.

#### 3.1 Learning a linear classifier: the Winnow algorithm

To the best of our knowledge, the first time multiplicative weight updates were used was in the Winnow algorithm of Littlestone [Lit87]. This is an algorithmic technique used in machine learning to learn linear classifiers. Equivalently, this can also be seen as solving a linear program.

The setup is as follows. We are given  $m$  labeled examples,  $(\mathbf{a}_1, \ell_1), (\mathbf{a}_2, \ell_2), \dots, (\mathbf{a}_m, \ell_m)$  where  $\mathbf{a}_j \in \mathbb{R}^n$  are *feature vectors*, and  $\ell_j \in \{-1, 1\}$  are their labels. Our goal is to find non-negative weights such that for any example, the sign of the weighted combination of the features matches its labels, i.e. find  $\mathbf{x} \in \mathbb{R}^n$  with  $x_i \geq 0$  such that for all  $j = 1, 2, \dots, m$ , we have  $\text{sgn}(\mathbf{a}_j \cdot \mathbf{x}) = \ell_j$ . Equivalently, we require that  $\ell_j \mathbf{a}_j \cdot \mathbf{x} \geq 0$  for all  $j$ . Without loss of generality, we may assume that the weights sum to 1 so that they form a distribution, i.e.  $\mathbf{1} \cdot \mathbf{x} = 1$ , where  $\mathbf{1}$  is the all 1’s vector.

Thus, for notational convenience, if we redefine  $\mathbf{a}_j$  to be  $\ell_j \mathbf{a}_j$ , then the problem reduces to finding a solution to the following LP:

$$\begin{aligned} \forall j = 1, 2, \dots, m : \quad & \mathbf{a}_j \cdot \mathbf{x} \geq 0 \\ & \mathbf{1} \cdot \mathbf{x} = 1 \\ \forall i : \quad & x_i \geq 0 \end{aligned} \tag{9}$$

Note this is a quite general form of LP, and many commonly seen LPs can be reduced to this form.

Now suppose there is a *large-margin* solution to this problem. I.e., there is an  $\varepsilon > 0$  and a distribution  $\mathbf{x}^*$  so that for all  $j$ , we have  $\mathbf{a}_j \cdot \mathbf{x}^* \geq \varepsilon$ . We now give an algorithm based on MW to solve the LP above. Define  $\rho = \max_j \|\mathbf{a}_j\|_\infty$ .

We run the MW algorithm in the gain form (see Section 2.3) with  $\eta = \frac{\varepsilon}{2\rho}$ . The decisions are given by the  $n$  features, and gains are specified by the  $m$  examples. The gain for feature  $i$  for example  $j$  is  $\frac{1}{\rho} a_{ij}$ . Note that these gains lie in the range  $[-1, 1]$  as required.

In each round  $t$ , let  $\mathbf{x}$  to be the distribution  $\mathbf{p}^{(t)}$  generated by the MW algorithm. Now, we look for a misclassified example, i.e. an example  $j$  such that  $\mathbf{a}_j \cdot \mathbf{x} < 0$ . If no such constraint exists, we are done and we can stop. Otherwise, if  $j$  is a misclassified example, then it specifies the gains for round  $t$ . Note that the gain in round  $t$  is  $\mathbf{m}^{(t)} \cdot \mathbf{p}^{(t)} = \frac{1}{\rho} \mathbf{a}_j \cdot \mathbf{x} < 0$ , whereas for the solution  $\mathbf{x}^*$ , we have  $\mathbf{m}^{(t)} \cdot \mathbf{x}^* = \frac{1}{\rho} \mathbf{a}_j \cdot \mathbf{x}^* \geq \frac{\varepsilon}{\rho}$ . We keep running the MW algorithm until we find a good solution (i.e. one that classifies all examples correctly).

To get a bound on the number of iterations until we find a good solution, we apply Corollary 7 with  $\mathbf{p} = \mathbf{x}^*$ . Using the trivial bound  $\eta |\mathbf{m}^{(t)}| \cdot \mathbf{p} \leq \eta$ , we have

$$0 > \sum_{t=1}^T \mathbf{m}^{(t)} \cdot \mathbf{p}^{(t)} \geq \sum_{t=1}^T (\mathbf{m}^{(t)} - \eta |\mathbf{m}^{(t)}|) \cdot \mathbf{p} - \frac{\ln n}{\eta} \geq \frac{\varepsilon}{2\rho} T - \frac{2\rho \ln n}{\varepsilon},$$

which implies that  $T < \frac{4\rho^2 \ln n}{\varepsilon^2}$ . Thus, in at most  $\left\lceil \frac{4\rho^2 \ln n}{\varepsilon^2} \right\rceil$  iterations, we find a good solution.

### 3.2 Solving zero-sum games approximately

We show how our general algorithm above can be used to approximately solve zero-sum games. (This is a duplication of the results of Freund and Schapire [FS99], who gave the same algorithm but a different proof of convergence that used KL-divergence. Furthermore, convergence of simple algorithms to zero-sum game equilibria were studied earlier in [Han57].)

Let  $\mathbf{A}$  be the payoff matrix of a finite 2-player zero-sum game, with  $n$  rows (the number of columns will play no role). When the row player plays strategy  $i$  and the column player plays strategy  $j$ , then the payoff to the column player is  $A(i, j) := A_{ij}$ . We assume that  $A(i, j) \in [0, 1]$ . If the row player chooses his strategy  $i$  from a distribution  $\mathbf{p}$  over the rows, then the expected payoff to the column player for choosing a strategy  $j$  is  $A(\mathbf{p}, j) := \mathbf{E}_{i \in \mathbf{p}}[A(i, j)]$ . Thus, the best response for the column player is the strategy  $j$  which maximizes this payoff. Similarly, if the column player chooses his strategy  $j$  from a distribution  $\mathbf{q}$  over the columns, then the expected payoff he gets if the row player chooses the strategy  $i$  is  $A(i, \mathbf{q}) := \mathbf{E}_{j \in \mathbf{q}}[A(i, j)]$ . Thus, the best response for the row player is the strategy  $i$  which minimizes this payoff. John von Neumann's min-max theorem says that if each of the players chooses a distribution over their strategies to optimize their worst case payoff (or payout), then the value they obtain is the same:

$$\min_{\mathbf{p}} \max_j A(\mathbf{p}, j) = \max_{\mathbf{q}} \min_i A(i, \mathbf{q}) \quad (10)$$

where  $\mathbf{p}$  (resp.,  $\mathbf{q}$ ) varies over all distributions over rows (resp., columns). Also,  $i$  (resp.,  $j$ ) varies over all rows (resp., columns). The common value of these two quantities, denoted  $\lambda^*$ , is known as the value of the game.

Let  $\varepsilon > 0$  be an error parameter. We wish to approximately solve the zero-sum game up to additive error of  $\varepsilon$ , namely, find mixed row and column strategies  $\tilde{\mathbf{p}}$  and  $\tilde{\mathbf{q}}$  such that

$$\lambda^* - \varepsilon \leq \min_i A(i, \tilde{\mathbf{q}}) \quad (11)$$

$$\max_j A(\tilde{\mathbf{p}}, j) \leq \lambda^* + \varepsilon. \quad (12)$$

The algorithmic assumption about the game is that given any distribution  $\mathbf{p}$  on decisions, we have an efficient way to pick the best response, namely, the pure column strategy  $j$  that maximizes  $A(\mathbf{p}, j)$ . This quantity is at least  $\lambda^*$  from the definition above. Call this algorithm the ORACLE.

**Theorem 8** *Given an error parameter  $\varepsilon > 0$ , there is an algorithm which solves the zero-sum game up to an additive factor of  $\varepsilon$  using  $O(\frac{\log n}{\varepsilon^2})$  calls to ORACLE, with an additional processing time of  $O(n)$  per call.*

PROOF: We map our general algorithm from Section 2 to this setting by considering (11) as specifying  $n$  linear constraints on the probability vector  $\tilde{\mathbf{q}}$ : viz., for all rows  $i$ ,  $A(i, \tilde{\mathbf{q}}) \geq \lambda^* - \varepsilon$ . Now, following the intuition given in the beginning of this section, we make our decisions to correspond to pure strategies of the row player. Thus a distribution on the decisions corresponds to a mixed row strategy. Costs of the decisions are specified by pure strategies of the column player. The cost paid by a decision  $i$  when column player chooses strategy  $j$  is  $A(i, j)$ .

In each round, given a distribution  $\mathbf{p}^{(t)}$  on the rows, we will choose the column  $j^{(t)}$  to be the best response strategy to  $\mathbf{p}^{(t)}$  for the column player, by calling ORACLE. Thus, the cost vector  $\mathbf{m}^{(t)}$  is the  $j^{(t)}$ -th column of the matrix  $\mathbf{A}$ .

Since all  $A(i, j) \in [0, 1]$ , we can apply Corollary 3 to get that after  $T$  rounds, for any distribution on the rows  $\mathbf{p}$ , we have

$$\sum_{t=1}^T A(\mathbf{p}^{(t)}, j^{(t)}) \leq (1 + \eta) \sum_{t=1}^T A(\mathbf{p}, j^{(t)}) + \frac{\ln n}{\eta}.$$

Dividing by  $T$ , and using the fact that  $A(\mathbf{p}, j^{(t)}) \leq 1$  and that for all  $t$ ,  $A(\mathbf{p}^{(t)}, j^{(t)}) \geq \lambda^*$ , we get

$$\lambda^* \leq \frac{1}{T} \sum_{t=1}^T A(\mathbf{p}^{(t)}, j^{(t)}) \leq \frac{1}{T} \sum_{t=1}^T A(\mathbf{p}, j^{(t)}) + \eta + \frac{\ln n}{\eta T}$$

Setting  $\mathbf{p} = \mathbf{p}^*$ , the optimal row strategy, we have  $A(\mathbf{p}, j) \leq \lambda^*$  for any  $j$ . By setting  $\eta = \frac{\varepsilon}{2}$  and  $T = \lceil \frac{4 \ln n}{\varepsilon^2} \rceil$ , we get that

$$\lambda^* \leq \frac{1}{T} \sum_{t=1}^T A(\mathbf{p}^{(t)}, j^{(t)}) \leq \frac{1}{T} \sum_{t=1}^T A(\mathbf{p}, j^{(t)}) + \varepsilon \leq \lambda^* + \varepsilon. \quad (13)$$

Thus,  $\frac{1}{T} \sum_{t=1}^T A(\mathbf{p}^{(t)}, j^{(t)})$  is an (additive)  $\varepsilon$ -approximation to  $\lambda^*$ .

Let  $\tilde{t}$  be the round  $t$  with the minimum value of  $A(\mathbf{p}^{(t)}, j^{(t)})$ . We have, from (13),

$$A(\mathbf{p}^{(\tilde{t})}, j^{(\tilde{t})}) \leq \frac{1}{T} \sum_{t=1}^T A(\mathbf{p}^{(t)}, j^{(t)}) \leq \lambda^* + \varepsilon.$$

Since  $j^{(\tilde{t})}$  maximizes  $A(\mathbf{p}^{(\tilde{t})}, j)$  over all  $j$ , we conclude that  $\mathbf{p}^{(\tilde{t})}$  is an approximately optimal mixed row strategy, and thus we can set  $\mathbf{p}^* := \mathbf{p}^{(\tilde{t})}$ .<sup>1</sup>

---

<sup>1</sup>Alternatively, we can set  $\mathbf{p}^* = \frac{1}{T} \sum_{t=1}^T \mathbf{p}^{(t)}$ . For let  $j^*$  be the optimal column player response to  $\mathbf{p}^*$ . Then we have  $A(\mathbf{p}^*, j^*) = \frac{1}{T} \sum_{t=1}^T A(\mathbf{p}^{(t)}, j^*) \leq \frac{1}{T} \sum_{t=1}^T A(\mathbf{p}^{(t)}, j^{(t)}) \leq \lambda^* + \varepsilon$ .

We set  $\mathbf{q}^*$  to be the distribution which assigns to column  $j$  the probability  $\frac{|\{t: j^{(t)}=j\}|}{T}$ . From (13), for any row strategy  $i$ , by setting  $\mathbf{p}$  to be concentrated on the pure strategy  $i$ , we have

$$\lambda^* - \varepsilon \leq \frac{1}{T} \sum_{t=1}^T A(i, j^{(t)}) = A(i, \mathbf{q}^*)$$

which shows that  $\mathbf{q}^*$  is an approximately optimal mixed column strategy.  $\square$

### 3.3 Plotkin, Shmoys, Tardos framework for packing/covering LPs

Plotkin, Shmoys, and Tardos [PST91] generalized some known flow algorithms to a framework for approximately solving *fractional packing and covering* problems, which are a special case of linear programming formally defined below. Their algorithm is a quantitative version of the classical *Lagrangian relaxation* idea, and applies also to general linear programs. Below, we derive the algorithm for general LPs and then mention the slight modification that yields better running time for packing-covering LPs. Also, we note that we could derive this algorithm as a special case of game solving, but for concreteness we describe it explicitly.

The basic problem is to check the feasibility of the following convex program:

$$\exists \mathbf{x} \in \mathcal{P} : \mathbf{Ax} \geq \mathbf{b} \quad (14)$$

where  $\mathbf{A} \in \mathbb{R}^{m \times n}$  is an  $m \times n$  matrix,  $\mathbf{x} \in \mathbb{R}^n$ , and  $\mathcal{P}$  is a convex set in  $\mathbb{R}^n$ . Intuitively, the set  $\mathcal{P}$  represents the “easy” constraints to satisfy, such as non-negativity, and  $\mathbf{A}$  represents the “hard” constraints to satisfy.

We wish to design an algorithm that given an error parameter  $\varepsilon > 0$ , either solves the problem to an additive error of  $\varepsilon$ , i.e., finds an  $\mathbf{x} \in \mathcal{P}$  such that for all  $i$ ,  $\mathbf{A}_i \mathbf{x} \geq b_i - \varepsilon$ , or failing that, proves that the system is infeasible. Here,  $\mathbf{A}_i$  is the  $i$ th row of  $\mathbf{A}$ .

We assume the existence of an algorithm, called ORACLE, which, given a probability vector  $\mathbf{p}$  on the  $m$  constraints, solves the following feasibility problem:

$$\exists \mathbf{x} \in \mathcal{P} : \mathbf{p}^\top \mathbf{Ax} \geq \mathbf{p}^\top \mathbf{b} \quad (15)$$

One way to implement this procedure is by maximizing  $\mathbf{p}^\top \mathbf{Ax}$  over  $\mathbf{x} \in \mathcal{P}$ . It is reasonable to expect such an optimization procedure to exist (indeed, such is the case for many applications) since we only need to check the feasibility of one constraint rather than  $m$ . If the feasibility problem (14) has a solution  $\mathbf{x}^*$ , then the same solution also satisfies (15) for *any* probability vector  $\mathbf{p}$  over the constraints. Thus, if there is a probability vector  $\mathbf{p}$  over the constraints such that no  $\mathbf{x} \in \mathcal{P}$  satisfies (15), then it is proof that the original problem is infeasible.

We assume that the ORACLE satisfies the following technical condition, which is necessary for deriving running time bounds:

**Definition 1** An  $(\ell, \rho)$ -**bounded** ORACLE, for parameters  $0 \leq \ell \leq \rho$ , is an algorithm which given a probability vector  $\mathbf{p}$  over the constraints, solves the feasibility problem (15).

Furthermore, there is a fixed subset  $I \subseteq [m]$  of constraints such that whenever the ORACLE manages to find a point  $\mathbf{x} \in \mathcal{P}$  satisfying (15), the following holds:

$$\begin{aligned} \forall i \in I: \quad & \mathbf{A}_i \mathbf{x} - b_i \in [-\ell, \rho] \\ \forall i \notin I: \quad & \mathbf{A}_i \mathbf{x} - b_i \in [-\rho, \ell] \end{aligned}$$

The value  $\rho$  is called the **width** of the problem.

In previous work, such as [PST91], only  $(\rho, \rho)$ -bounded ORACLES are considered. We separate out the upper and lower bounds in order to obtain tighter guarantees on the running time. The results of [PST91] can be recovered simply by setting  $\ell = \rho$ .

**Theorem 9** *Let  $\varepsilon > 0$  be a given error parameter. Suppose there exists an  $(\ell, \rho)$ -bounded ORACLE for the feasibility problem (15). Assume that  $\ell \geq \frac{\varepsilon}{2}$ . Then there is an algorithm which either solves the problem up to an additive error of  $\varepsilon$ , or correctly concludes that the system is infeasible, making only  $O(\frac{\ell \rho \log(m)}{\varepsilon^2})$  calls to the ORACLE, with an additional processing time of  $O(m)$  per call.*

PROOF: The condition  $\ell \geq \frac{\varepsilon}{2}$  is only technical, and if it is not met we can just redefine  $\ell$  to be  $\frac{\varepsilon}{2}$ . To map our general framework to this situation, we have a decision representing each of the  $m$  constraints. Costs are determined by points  $\mathbf{x} \in \mathcal{P}$ . The cost of constraint  $i$  for point  $\mathbf{x}$  is  $\frac{1}{\rho}[\mathbf{A}_i \mathbf{x} - b_i]$  (so that the costs lie in the range  $[-1, 1]$ ).

In each round  $t$ , given a distribution over the decisions (i.e. the constraints)  $\mathbf{p}^{(t)}$ , we run the ORACLE with  $\mathbf{p}^{(t)}$ . If the ORACLE declares that there is no  $\mathbf{x} \in \mathcal{P}$  such that  $\mathbf{p}^{(t)\top} \mathbf{A} \mathbf{x} \geq \mathbf{p}^{(t)\top} \mathbf{b}$ , then we stop, because now  $\mathbf{p}^{(t)}$  is proof that the problem (14) is infeasible.

So let us assume that this doesn't happen, i.e. in all rounds  $t$ , the ORACLE manages to find a solution  $\mathbf{x}^{(t)}$  such  $\mathbf{p}^{(t)\top} \mathbf{A} \mathbf{x} \geq \mathbf{p}^{(t)\top} \mathbf{b}$ . Since the cost vector to the Multiplicative Weights algorithm is specified to be  $\mathbf{m}^{(t)} := \frac{1}{\rho}[\mathbf{A} \mathbf{x}^{(t)} - \mathbf{b}]$ , we conclude that the expected cost in each round is non-negative:

$$\mathbf{m}^{(t)} \cdot \mathbf{p}^{(t)} = \frac{1}{\rho}[\mathbf{A} \mathbf{x}^{(t)} - \mathbf{b}] \cdot \mathbf{p}^{(t)} = \frac{1}{\rho}[\mathbf{p}^{(t)\top} \mathbf{A} \mathbf{x} - \mathbf{p}^{(t)\top} \mathbf{b}] \geq 0.$$

Let  $i \in I$ . Then Theorem 2 tells us that after  $T$  rounds,

$$\begin{aligned} 0 &\leq \sum_{t=1}^T \frac{1}{\rho}[\mathbf{A}_i \mathbf{x}^{(t)} - b_i] + \eta \sum_{t=1}^T \frac{1}{\rho} |\mathbf{A}_i \mathbf{x}^{(t)} - b_i| + \frac{\ln m}{\eta} \\ &= (1 + \eta) \sum_{t=1}^T \frac{1}{\rho}[\mathbf{A}_i \mathbf{x}^{(t)} - b_i] + 2\eta \sum_{<0} \frac{1}{\rho} |\mathbf{A}_i \mathbf{x}^{(t)} - b_i| + \frac{\ln m}{\eta} \\ &\leq (1 + \eta) \sum_{t=1}^T \frac{1}{\rho}[\mathbf{A}_i \mathbf{x}^{(t)} - b_i] + \frac{2\eta\ell}{\rho} T + \frac{\ln m}{\eta} \end{aligned}$$

Here, the subscript “ $< 0$ ” refers to the rounds  $t$  when  $\mathbf{A}_i \mathbf{x}^{(t)} - b_i < 0$ . The last inequality follows because if  $\mathbf{A}_i \mathbf{x}^{(t)} - b_i < 0$ , then  $|\mathbf{A}_i \mathbf{x}^{(t)} - b_i| \leq \ell$ . Dividing by  $T$ , multiplying by  $\rho$ , and letting  $\bar{\mathbf{x}} = \frac{1}{T} \sum_{t=1}^T \mathbf{x}^{(t)}$  (note that  $\bar{\mathbf{x}} \in \mathcal{P}$  since  $\mathcal{P}$  is a convex set), we get that

$$0 \leq (1 + \eta)[\mathbf{A}_i \bar{\mathbf{x}} - b_i] + 2\eta\ell + \frac{\rho \ln(m)}{\eta T}.$$



Now, if we choose  $\eta = \frac{\varepsilon}{4\ell}$  (note that  $\eta \leq \frac{1}{2}$  since  $\ell \geq \frac{\varepsilon}{2}$ ), and  $T = \lceil \frac{8\ell\rho\ln(m)}{\varepsilon^2} \rceil$ , we get that

$$0 \leq (1 + \eta)[\mathbf{A}_i \bar{\mathbf{x}} - b_i] + \varepsilon \implies \mathbf{A}_i \bar{\mathbf{x}} \geq b_i - \varepsilon.$$

Reasoning similarly for  $i \notin I$ , we get the same inequality. Putting both together, we conclude that  $\bar{\mathbf{x}}$  satisfies the feasibility problem (14) up to an additive  $\varepsilon$  factor, as desired.  $\square$

### 3.3.1 Concave constraints

The algorithm of Section 3.3 works not just for linear constraints over a convex domain, but also for concave constraints. Imagine that we have the following feasibility problem:

$$\exists \mathbf{x} \in \mathcal{P} : \quad \forall i \in [m] : f_i(\mathbf{x}) \geq 0 \quad (16)$$

where, as before,  $\mathcal{P} \subseteq \mathbb{R}^n$  is a convex domain, and for  $i \in [m]$ ,  $f_i : \mathcal{P} \rightarrow \mathbb{R}$  are concave functions. We wish to satisfy this system approximately, up to an additive error of  $\varepsilon$ . Again, we assume the existence of an ORACLE, which, when given a probability distribution  $\mathbf{p} = \langle p_1, p_2, \dots, p_m \rangle^\top$ , solves the following feasibility problem:

$$\exists \mathbf{x} \in \mathcal{P} : \quad \sum_i p_i f_i(\mathbf{x}) \geq 0 \quad (17)$$

An ORACLE would be called  $(\ell, \rho)$ -bounded there is a fixed subset of constraints  $I \subseteq [m]$  such that whenever it returns a feasible solution  $\mathbf{x}$  to (17), all constraints  $i \in I$  take values in the range  $[-\ell, \rho]$  on the point  $\mathbf{x}$ , and all the rest take values in  $[-\rho, \ell]$ .

**Theorem 10** *Let  $\varepsilon > 0$  be a given error parameter. Suppose there exists an  $(\ell, \rho)$ -bounded ORACLE for the feasibility problem (16). Assume that  $\ell \geq \frac{\varepsilon}{2}$ . Then there is an algorithm which either solves the problem up to an additive error of  $\varepsilon$ , or correctly concludes that the system is infeasible, making only  $O(\frac{\ell\rho\log(m)}{\varepsilon^2})$  calls to the ORACLE, with an additional processing time of  $O(m)$  per call.*

PROOF: Just as before, we have a decision for every constraint, and costs are specified by points  $\mathbf{x} \in \mathcal{P}$ . The cost of constraint  $i$  for point  $\mathbf{x}$  is  $\frac{1}{\rho} f_i(\mathbf{x})$ .

Now we run the Multiplicative Weights algorithm with this setup. Again, if at any point the ORACLE declares that (17) is infeasible, we immediately halt and declare the system (16) infeasible. So assume this never happens. Then as before, the expected cost in each round is  $\mathbf{m}^{(t)} \cdot \mathbf{p}^{(t)} \geq 0$ . Now, applying Theorem 2 as before, we conclude that for any  $i \in I$ , we have

$$0 \leq (1 + \eta) \sum_{t=1}^T \frac{1}{\rho} f_i(\mathbf{x}^{(t)}) + \frac{2\eta\ell}{\rho} T + \frac{\ln m}{\eta}.$$

Dividing by  $T$ , multiplying by  $\rho$ , and letting  $\bar{\mathbf{x}} = \frac{1}{T} \sum_{t=1}^T \mathbf{x}^{(t)}$  (note that  $\bar{\mathbf{x}} \in \mathcal{P}$  since  $\mathcal{P}$  is a convex set), we get that

$$0 \leq (1 + \eta) \rho f_i(\bar{\mathbf{x}}) + 2\eta\ell + \frac{\rho \ln(m)}{\eta T},$$

since  $\frac{1}{T} \sum_{t=1}^T f_i(\mathbf{x}^{(t)}) \leq f_i(\frac{1}{T} \sum_{t=1}^T \mathbf{x}^{(t)})$ , by Jensen's inequality, since all the  $f_i$  are concave.

Now, if we choose  $\eta = \frac{\varepsilon}{4\ell}$  (note that  $\eta \leq \frac{1}{2}$  since  $\ell \geq \frac{\varepsilon}{2}$ ), and  $T = \lceil \frac{8\ell\rho \ln(m)}{\varepsilon^2} \rceil$ , we get that

$$0 \leq (1 + \eta)f_i(\bar{\mathbf{x}}) + \varepsilon \implies f_i(\bar{\mathbf{x}}) \geq -\varepsilon.$$

Reasoning similarly for  $i \notin I$ , we get the same inequality. Putting both together, we conclude that  $\bar{\mathbf{x}}$  satisfies the feasibility problem (16) up to an additive  $\varepsilon$  factor, as desired.  $\square$

### 3.3.2 Approximate Oracles

The algorithm described in the previous section allows some slack for the implementation of the ORACLE. This slack is very useful in designing efficient implementations for the ORACLE.

Define a  $\varepsilon$ -approximate ORACLE for the feasibility problem (14) to be one that solves the feasibility problem (15) up to an additive error of  $\varepsilon$ . That is, given a probability vector  $\mathbf{p}$  on the constraints, either it finds an  $\mathbf{x} \in \mathcal{P}$  such that  $\mathbf{p}^\top \mathbf{A}\mathbf{x} \geq \mathbf{p}^\top \mathbf{b} - \varepsilon$ , or it declares correctly that (15) is infeasible.

**Theorem 11** *Let  $\varepsilon > 0$  be a given error parameter. Suppose there exists an  $(\ell, \rho)$ -bounded  $\frac{\varepsilon}{3}$ -approximate ORACLE for the feasibility problem (14). Assume that  $\ell \geq \frac{\varepsilon}{3}$ . Then there is an algorithm which either solves the problem up to an additive error of  $\varepsilon$ , or correctly concludes that the system is infeasible, making only  $O(\frac{\ell\rho \log(m)}{\varepsilon^2})$  calls to the ORACLE, with an additional processing time of  $O(m)$  per call.*

PROOF: We run the algorithm of the previous section with the given ORACLE, setting  $\eta = \frac{\varepsilon}{6\ell}$ . Now, in every round, the expected payoff is at least  $-\frac{\varepsilon}{3\rho}$ . Simplifying as before, we get that after  $T$  rounds, we have, the average point  $\bar{\mathbf{x}} = \frac{1}{T} \sum_{t=1}^T \mathbf{x}^{(t)}$  returned by the ORACLE satisfies

$$-\frac{\varepsilon}{3} \leq (1 + \eta)[\mathbf{A}_i \bar{\mathbf{x}} - b_i] + 2\eta\ell + \frac{\rho \ln(m)}{\eta T}.$$

Now, if  $T = \lceil \frac{18\ell\rho \ln(m)}{\varepsilon^2} \rceil$ , then we get that for all  $i$ ,  $\mathbf{A}_i \bar{\mathbf{x}} \geq b_i - \varepsilon$ , as required.  $\square$

### 3.3.3 Fractional Covering Problems

In fractional covering problems, the framework is the same as above, with the crucial difference that the coefficient matrix  $\mathbf{A}$  is such that  $\mathbf{A}\mathbf{x} \geq 0$  for all  $\mathbf{x} \in \mathcal{P}$ , and  $\mathbf{b} > 0$ . A  $\varepsilon$ -approximation solution to this system is an  $\mathbf{x} \in \mathcal{P}$  such that  $\mathbf{A}\mathbf{x} \geq (1 - \varepsilon)\mathbf{b}$ .

We assume without loss of generality (by appropriately scaling the inequalities) that  $b_i = 1$  for all rows, so that now we desire to find an  $\mathbf{x} \in \mathcal{P}$  which satisfies the system within an additive  $\varepsilon$  factor. Since for all  $\mathbf{x} \in \mathcal{P}$ , we have  $\mathbf{A}\mathbf{x} \geq 0$ , and since all  $b_i = 1$ , we conclude that for any  $i$ ,  $\mathbf{A}_i \mathbf{x} - b_i \geq -1$ . Thus, we assume that there is a  $(1, \rho)$ -bounded ORACLE for this problem. Now, applying Theorem 9, we get the following:

**Theorem 12** *Suppose there exists a  $(1, \rho)$ -bounded ORACLE for the program  $\mathbf{A}\mathbf{x} \geq \mathbf{b}$  with  $\mathbf{x} \in \mathcal{P}$ . Given an error parameter  $\varepsilon > 0$ , there is an algorithm which computes a  $\varepsilon$ -approximate solution to the program, or correctly concludes that it is infeasible, using  $O(\frac{\rho \log(m)}{\varepsilon^2})$  calls to the ORACLE, plus an additional processing time of  $O(m)$  per call.*

### 3.3.4 Fractional Packing Problems

A fractional packing problem can be written as

$$\exists \mathbf{x} \in \mathcal{P} : \quad \mathbf{Ax} \leq \mathbf{b}$$

where  $\mathcal{P}$  is a convex domain such that  $\mathbf{Ax} \geq 0$  for all  $\mathbf{x} \in \mathcal{P}$ , and  $\mathbf{b} > 0$ . A  $\varepsilon$ -approximate solution to this system is an  $\mathbf{x} \in \mathcal{P}$  such that  $\mathbf{Ax} \leq (1 + \varepsilon)\mathbf{b}$ .

Again, we assume that  $b_i = 1$  for all  $i$ , scaling the constraints if necessary. Now by rewriting this system as

$$\exists \mathbf{x} \in \mathcal{P} : \quad -\mathbf{Ax} \geq -\mathbf{b}$$

we cast it in our general framework, and a solution  $\mathbf{x} \in \mathcal{P}$  which satisfies this up to an additive  $\varepsilon$  is a  $\varepsilon$ -approximate solution to the original system. Since for all  $\mathbf{x} \in \mathcal{P}$ , we have  $\mathbf{Ax} \geq 0$ , and since all  $b_i = 1$ , we conclude that for any  $i$ ,  $-\mathbf{A}_i\mathbf{x} + b_i \leq 1$ . Thus, we assume that there is a  $(1, \rho)$ -bounded ORACLE for this problem. Now, applying Theorem 9, we get the following:

**Theorem 13** *Suppose there exists a  $(1, \rho)$ -bounded ORACLE for the program  $-\mathbf{Ax} \geq -\mathbf{b}$  with  $\mathbf{x} \in \mathcal{P}$ . Given an error parameter  $\varepsilon > 0$ , there is an algorithm which computes a  $\varepsilon$ -approximate solution to the program, or correctly concludes that it is infeasible, using  $O(\frac{\rho \log(m)}{\varepsilon^2})$  calls to the ORACLE, plus an additional processing time of  $O(m)$  per call.*

## 3.4 Approximating Multicommodity Flow Problems

Multicommodity flow problems are represented by packing/covering LPs and thus can be approximately solved using the PST framework outlined above. The resulting flow algorithm is outlined below together with a brief analysis. Unfortunately, the algorithm is not polynomial-time because its running time is bounded by a polynomial function of the edge capacities (as opposed to the logarithm of the capacities, which is the number of bits needed to represent them). Garg and Könemann [GK98] fixed this problem with a better algorithm whose running time does not depend upon the edge capacities.

Here we derive the Garg-Könemann algorithm using our general framework. This will highlight the essential new idea, namely, a reweighting of penalties to reduce the *width* parameter. Note that algorithm is not quite the same as in [GK98] (the termination condition is slightly different) and neither is the proof; the running time bound is the same however.

For illustrative purposes we focus on the *maximum multicommodity flow problem*. In this problem, we are given a graph  $G = (V, E)$  with capacities  $c_e$  on edges, and set of  $k$  source-sink pairs of nodes. Let  $P$  be the set of all paths between the source-sink pairs. The objective is to maximize the total flow between these pairs. The LP formulation is as follows:

$$\begin{aligned} & \max \sum_{p \in P} f_p \\ & \forall e \in E : \quad \sum_{p \ni e} f_p \leq c_e \\ & \forall p \in P : \quad f_p \geq 0 \end{aligned} \tag{18}$$

Here, the variable  $f_p$  represents the flow on path  $p$ .

Before presenting the Garg-Könemann idea we first present the algorithm one would obtain by applying our packing-covering framework (Section 3.3) in the obvious way.

First, note that by using binary search we can reduce the optimization problem to feasibility, by iteratively introducing a new constraint that gives a lower bound on the objective. So assume without loss of generality that we know the value  $F^{\text{opt}}$  of the total flow in the optimum solution. Then we want to solve the following feasibility problem:

$$\exists \mathbf{f} \in \mathcal{P} : \forall e \in E : \sum_{p \ni e} f_p \leq c_e,$$

where

$$\mathcal{P} := \left\{ \mathbf{f} : \forall p \in P : f_p \geq 0, \sum_{p \in P} f_p = F^{\text{opt}} \right\}.$$

In this form, the feasibility problem given above is a packing LP, thus, we can apply the Multiplicative Weights algorithm of Section 3.3.4.

As outlined in Section 3.3, the obvious algorithm would maintain at each step  $t$  a weight  $w_e^{(t)}$  for each edge  $e$ . The ORACLE can be implemented by finding the flow in  $\mathcal{P}$  which minimizes  $\sum_e w_e^{(t)} \sum_{p \ni e} f_p / c_e = \sum_p f_p \sum_{e \in p} w_e^{(t)} / c_e$ . The optimal flow is supported on a single path, namely, the path  $p^{(t)} \in P$  that has minimum length, when every edge  $e \in E$  is given length  $w_e^{(t)} / c_e$ . Thus in every round we find this path  $p^{(t)}$  and pass a flow  $F^{\text{opt}}$  on this path. Note that the final flow will be an average of the flows in each event, and hence will also have value  $F^{\text{opt}}$ . Costs for the edges are defined as in Section 3.3.

Unfortunately the width parameter is  $\rho = \max_{\mathbf{f} \in \mathcal{P}} \max_e \sum_{p \ni e} f_p / c_e = F^{\text{opt}} / c_{\min}$  where  $c_{\min}$  is the capacity of the minimum capacity edge in the graph. The algorithm requires  $T = \rho \ln(n) / \varepsilon^2$  iterations to get an  $(1 - \varepsilon)$ -approximation to the optimal flow. The overall running time is  $\tilde{O}(F^{\text{opt}} T_{\text{sp}} / c_{\min})$  where  $T_{\text{sp}} \leq O(mk)$  is the time needed to compute  $k$  shortest paths. As already mentioned, this is not polynomial-time since it depends upon  $1/c_{\min}$  rather than the logarithm of this value.

Now we describe the Garg-Könemann modification. We continue to maintain weights  $w_e^{(t)}$  for every edge  $e$ , where initially,  $w_e^{(1)} = 1$  for all  $e$ . The costs are determined by flows as before, however, we consider a larger set of flows, viz.

$$\mathcal{P}' := \{ \mathbf{f} : \forall p \in P : f_p \geq 0 \}.$$

Note that we no longer need to use the value  $F^{\text{opt}}$ . Again, in each round we choose a flow  $\mathbf{f} \in \mathcal{P}'$  that is supported on the shortest path  $p^{(t)} \in P$  with edge lengths  $w_e / c_e$ .

The main idea behind width reduction is the following: instead of routing the same flow  $F^{\text{opt}}$  at each time step, we route only as much flow as is allowed by the minimum capacity edge on the path. In other words, at time  $t$  we route a flow of value  $c^{(t)}$  on path  $p^{(t)}$ , where  $c^{(t)}$  is the minimum capacity of an edge on the path  $p^{(t)}$ . The cost incurred by edge  $e$  is  $m_e^{(t)} = c^{(t)} / c_e$ . (In other words, a cost of  $1/c_e$  per unit of flow passing through  $e$ .) *The width is therefore automatically upper bounded by 1.*

We run the MW algorithm with  $\eta = \varepsilon/2$ . The update rule in this setting consists of updating the weights of all edges in path  $p^{(t)}$  and leaving other weights unchanged at that

step:

$$\forall e \in p^{(t)} : \quad w_e^{(t+1)} = w_e^{(t)} \left( 1 + \eta \cdot \frac{c^{(t)}}{c_e} \right)$$

The *termination rule* for the algorithm is to stop when as soon as for some edge  $e$ , the *congestion*  $\frac{f_e}{c_e} \geq \frac{\ln m}{\eta^2}$ , where  $f_e$  is the total amount of flow routed by the algorithm so far on edge  $e$ .

### 3.4.1 Analysis

We apply Theorem 6. Since we have  $m_e^{(t)} \in [0, 1]$  for all edges  $e$  and rounds  $t$ , we conclude that for any edge  $e$ , we have

$$\sum_{t=1}^T \mathbf{m}^{(t)} \cdot \mathbf{p}^{(t)} \geq (1 - \eta) \sum_{t=1}^T m_e^{(t)} - \frac{\ln(m)}{\eta}. \quad (19)$$

We now analyze both sides of this inequality. In round  $t$ , for any edge  $e$ , we have  $m_e^{(t)} = \frac{c^{(t)}}{c_e}$  if  $e \in p^{(t)}$ , and 0 if  $e \notin p$ . Thus, we have

$$\sum_{t=1}^T \mathbf{m}^{(t)} = \frac{f_e}{c_e}, \quad (20)$$

where  $f_e$  is the total amount of flow on  $e$  at the end of the algorithm, and

$$\sum_{t=1}^T \mathbf{m}^{(t)} \cdot \mathbf{p}^{(t)} = \sum_{t=1}^T \frac{\sum_{e \in p^{(t)}} \frac{c^{(t)}}{c_e} \cdot w_e^{(t)}}{\sum_e w_e^{(t)}} = \sum_{t=1}^T c^{(t)} \cdot \frac{\sum_{e \in p^{(t)}} \frac{w_e^{(t)}}{c_e}}{\sum_e w_e^{(t)}}. \quad (21)$$

Now, suppose the optimum flow assigns  $f_p^{\text{opt}}$  flow to path  $p$ , and let  $F^{\text{opt}} = \sum_p f_p^{\text{opt}}$  be the total flow. For any set of edge lengths  $w_e/c_e$ , the shortest path  $p \in P$  with these edge lengths satisfies

$$\frac{\sum_e w_e}{\sum_{e \in p} \frac{w_e}{c_e}} \geq \frac{\sum_e w_e \cdot \sum_{p' \ni e} \frac{f_{p'}^{\text{opt}}}{c_e}}{\sum_{e \in p} \frac{w_e}{c_e}} = \frac{\sum_{p'} f_{p'}^{\text{opt}} \cdot \sum_{e \in p'} \frac{w_e}{c_e}}{\sum_{e \in p} \frac{w_e}{c_e}} \geq \sum_{p'} f_{p'}^{\text{opt}} = F^{\text{opt}}.$$

The first inequality follows because for any edge  $e$ , we have  $\sum_{p' \ni e} f_{p'}^{\text{opt}} \leq c_e$ . The second inequality follows from the fact that  $p$  is the shortest path with edge lengths given by  $w_e/c_e$ . Using this bound in (21), we get that

$$\sum_{t=1}^T \mathbf{m}^{(t)} \cdot \mathbf{p}^{(t)} \leq \sum_{t=1}^T \frac{c^{(t)}}{F^{\text{opt}}} = \frac{F}{F^{\text{opt}}}, \quad (22)$$

where  $F = \sum_{t=1}^T c^{(t)}$  is the total amount of flow passed by the algorithm.

Plugging (20) and (22) into (19), we get that

$$\frac{F}{F^{\text{opt}}} \geq (1 - \eta) \max_e \left\{ \frac{f_e}{c_e} \right\} - \frac{\ln(m)}{\eta}.$$

We stop the algorithm as soon as there we have an edge  $e$  with congestion  $\frac{f_e}{c_e} \geq \frac{\ln(m)}{\eta^2}$ , so when the algorithm terminates we have

$$\frac{F}{F^{\text{opt}}} \geq (1 - 2\eta) \max_e \left\{ \frac{f_e}{c_e} \right\}.$$

Now,  $C := \max_e \left\{ \frac{f_e}{c_e} \right\}$  is the maximum congestion of the flow passed by the algorithm. So, the flow scaled down by  $C$  respects all capacities. For this scaled down flow, we have that the total flow is

$$\frac{F}{C} \geq (1 - 2\eta) F^{\text{opt}},$$

which shows that the scaled-down flow is within  $(1 - 2\eta) = (1 - \varepsilon)$  of optimal.

**Running time.** In every iteration  $t$  of the algorithm, consider the minimum capacity edge  $e$  on the chosen path  $p^{(t)}$ . It gets congested by the flow of value  $c^{(t)} = c_e$  sent in that round. Since we stop the algorithm as soon as the congestion on any edge is at least  $\frac{\ln(m)}{\eta^2}$ , any given edge can be the minimum capacity edge on the chosen path at most  $\lceil \frac{\ln(m)}{\eta^2} \rceil$  times in the entire run of the algorithm. Since there are  $m$  edges, the number of iterations is therefore at most  $m \cdot \lceil \frac{\ln(m)}{\eta^2} \rceil = O(\frac{m \log m}{\varepsilon^2})$ .

Each iteration involves  $k$  shortest path computations. Recall that  $T_{\text{sp}}$  is the time needed for this. Thus, the overall running time is  $O(\frac{m \log m}{\varepsilon^2} \cdot T_{\text{sp}})$ .

### 3.5 $O(\log n)$ -approximation for many NP-hard problems

For many NP-hard problems, typically integer versions of packing-covering problems, one can compute a  $O(\log n)$ -approximation by solving the obvious LP relaxation and then using Raghavan-Thompson [RT87] randomized rounding. This yields a randomized algorithm; to obtain a deterministic algorithm, derandomize it using Raghavan's [Rag86] method of *pessimistic estimators*.

Young [You95] has given an especially clear framework for understanding these algorithms which as a bonus also yields faster, combinatorial algorithms for approximating integer packing/covering programs. He observes that one can collapse the three ideas in the algorithm above —LP solving, randomized rounding, derandomization— into a single algorithm that uses the multiplicative update rule, and does not need to solve the LP relaxation directly. (Young's paper is titled “*Randomized rounding without solving the linear program.*”)

At the root of Young's algorithm is the observation that the analysis of randomized rounding uses the Chernoff-Hoeffding bounds. These bounds show that the sum of bounded independent random variables  $X_1, X_2, \dots, X_n$  (which should be thought of as the random variables generated in the randomized rounding algorithm) is sharply concentrated about its mean, and are proved by applying Markov's inequality to the variable  $e^{\eta(\sum_i X_i)}$  for some parameter  $\eta$ . The key observation now is that the afore-mentioned application of Markov's inequality bounds the probability of failure (of the randomized rounding algorithm) away from 0 in terms of  $\mathbf{E}[e^{\eta(\sum_i X_i)}]$ . Thus, one can treat  $\mathbf{E}[e^{\eta(\sum_i X_i)}]$  as a pessimistic estimator (up to scaling by a constant) of the failure probability, and derandomization can be achieved

by greedily (and hence, deterministically) choosing the  $X_i$  sequentially to decrease this pessimistic estimator. The resulting algorithm is essentially the MW algorithm: in each round  $t$ , the deterministic part of the pessimistic estimator, viz.  $e^{\eta \sum_{\tau < t} X_\tau}$ , plays the role of the weight.

Below, we illustrate this idea using the canonical problem in this class, SET COVER. (A similar analysis works for other problems.) Since we have developed the multiplicative weights framework already, we do not detail Young's original intuition involving Chernoff bound arguments and can proceed directly to the algorithm. In fact, the algorithm can be simplified so it becomes exactly the classical greedy algorithm, and we obtain a  $\ln n$ -approximation, which is best-possible for this problem (assuming reasonable complexity-theoretic conjectures [Fei98]).

In the SET COVER problem, we are given a universe of  $n$  elements, say  $U = \{1, 2, 3, \dots, n\}$  and a collection  $\mathcal{C}$  of subsets of  $U$  whose union equals  $U$ . We are required to pick the minimum number of sets from  $\mathcal{C}$  which cover all of  $U$ . Let this minimum number be denoted OPT. The Greedy Algorithm picks subsets iteratively, each time choosing that set which covers the maximum number of uncovered elements.

We analyze the Greedy Algorithm in our setup as follows. Each element of the universe represents a constraint that the union of sets picked by the algorithm must cover it. Following the guidelines given in the beginning of Section 3, we cast the problem in our framework by letting decisions correspond to elements in the universe, and costs determined by sets  $C_j \in \mathcal{C}$ . The cost of the constraint corresponding to element  $i$  for a given set  $C_j$  is 1 if  $i \in C_j$  and 0 otherwise.

To translate the greedy algorithm to our framework, suppose we run the Multiplicative Weights Update algorithm with this setup with  $\eta = 1$ . Since the analysis in the proof of Theorem 2 technically requires  $\eta \leq 1/2$ , in the following we repeat the same potential function analysis for the current setting for  $\eta = 1$ . For  $\eta = 1$ , the update rule  $w_i^{(t+1)} = w_i^{(t)}(1 - \eta m_i^{(t)})$  implies that elements that have been covered so far have weight 0 while all the rest have weight 1. Thus, the distribution  $\mathbf{p}^{(t)}$  is simply the uniform distribution on the uncovered elements until time  $t$ . Since cost of an element for a given set is 1 if it is in the set and 0 otherwise, the set that maximizes the expected cost under  $\mathbf{p}^{(t)}$  is the one that maximizes the number of uncovered elements. The resulting algorithm is the Greedy Set Cover algorithm.

Since OPT sets cover all the elements, for any distribution  $p_1, p_2, \dots, p_n$  on the elements, one set must cover at least  $1/\text{OPT}$  fraction of elements. This implies that if we choose set  $C$  in round  $t$  that maximizes the number of uncovered elements, we have

$$\mathbf{m}^{(t)} \cdot \mathbf{p}^{(t)} = \max_C \sum_{i \in C} p_i \geq 1/\text{OPT}.$$

Following the analysis of Theorem 2, the change in potential for each round is:

$$\Phi^{(t+1)} = \Phi^{(t)}(1 - \eta \mathbf{m}^{(t)} \cdot \mathbf{p}^{(t)}) < \Phi^{(t)} e^{-\eta/\text{OPT}} = \Phi^{(t)} e^{-1/\text{OPT}}.$$

The strict inequality holds because  $\mathbf{m}^{(t)} \cdot \mathbf{p}^{(t)} > 0$  as long as there are uncovered elements. Thus, the potential drops by a factor of  $e^{-1/\text{OPT}}$  every time.

We run this as long as some element has not yet been covered. We show that  $T = \lceil \ln n \rceil \text{OPT}$  iterations suffice, which implies that we have a  $\lceil \ln n \rceil$  approximation to OPT.

We have

$$\Phi^{(T+1)} < \Phi^{(1)} e^{-T/\text{OPT}} = n e^{-\lceil \ln n \rceil \text{OPT}/\text{OPT}} = n e^{-\lceil \ln n \rceil} \leq 1$$

Note that with  $\eta = 1$ ,  $\Phi^{(T+1)}$  is exactly the number of elements left uncovered after  $T$  iterations. So we conclude that all elements are covered.

### 3.6 Learning theory and boosting

Boosting [Sch90] —the process of combining several moderately accurate rules-of-thumb into a single highly accurate prediction rule— is a central idea in Machine Learning today. Freund and Schapire’s *AdaBoost* [FS97] uses the Multiplicative Weights Update Rule and fits in our framework. Here we explain the main idea using some simplifying assumptions.

Let  $X$  be some set (*domain*) and suppose we are trying to learn an unknown function (*concept*)  $c : \mathcal{X} \rightarrow \{0, 1\}$  chosen from a concept class  $\mathcal{C}$ . Given a sequence of training examples  $(x, c(x))$  where  $x$  is generated from a fixed but unknown distribution  $\mathcal{D}$  on the domain  $\mathcal{X}$ , the learning algorithm is required to output a *hypothesis*  $h : \mathcal{X} \rightarrow \{0, 1\}$ . The *error* of the hypothesis is defined to be  $\mathbf{E}_{x \sim \mathcal{D}}[|h(x) - c(x)|]$ .

A *strong learning algorithm* is one that, for every distribution  $\mathcal{D}$  and every  $\varepsilon, \delta > 0$ , outputs with probability  $1 - \delta$  a hypothesis whose error is at most  $\varepsilon$ . A  $\gamma$ -*weak learning algorithm* for  $\gamma > 0$  is similar, except its error is as high as  $1/2 - \gamma$ . Boosting shows that if a  $\gamma$ -weak learning algorithm exists for a concept class, then a strong learning algorithm exists. (The running time of the algorithm and the number of samples may depend on  $\gamma$ .)

We prove this result in the so-called *boosting by sampling framework*, which uses a fixed training set  $X$  of  $N$  examples drawn from the distribution  $\mathcal{D}$ . The goal is to make sure that the final hypothesis erroneously classifies at most  $\varepsilon$  fraction of this training set. Using VC-dimension theory (see [FS97]) one can then show that if the VC-dimension of the concept class  $\mathcal{C}$  is bounded, and if  $N$  is chosen large enough (in terms of the error and confidence parameters, and the VC-dimension of the concept class  $\mathcal{C}$ ), then with probability at least  $1 - \delta$  over the choice of the sample set, the error of the hypothesis over the entire domain  $\mathcal{X}$  (under distribution  $\mathcal{D}$ ) is at most  $2\varepsilon$ .

The idea in boosting is to repeatedly run the weak learning algorithm on different distributions defined on the fixed training set  $X$ . The final hypothesis has error  $\varepsilon$  under the *uniform* distribution on the  $X$ . We run the MW algorithm with  $\eta = \gamma$  for  $T = \left\lceil \frac{2}{\gamma^2} \ln\left(\frac{1}{\varepsilon}\right) \right\rceil$  rounds. The decisions correspond to samples in  $X$  and costs are specified by a hypothesis generated by the weak learning algorithm, in the following way. If hypothesis  $h$  is generated, the cost for decision point  $x$  is 1 or 0 depending on whether  $h(x) = c(x)$  or not. In other words, the cost vector  $\mathbf{m}$  (indexed by  $x$  rather than  $i$ ) is specified by  $m_x = 1 - |h(x) - c(x)|$ . Intuitively, we want the weight of an example to *increase* if the hypothesis labels it incorrectly.

In each iteration, the algorithm presents the current distribution  $\mathbf{p}^{(t)}$  on the examples to the weak learning algorithm, and in return obtains a hypothesis  $h^{(t)}$  whose error with respect to the distribution  $\mathbf{p}^{(t)}$  is not more than  $1/2 - \gamma$ , in other words, the expected cost in each iteration satisfies

$$\mathbf{m}^{(t)} \cdot \mathbf{p}^{(t)} \geq \frac{1}{2} + \gamma.$$



The algorithm is run for  $T$  rounds, where  $T$  will be specified shortly. The final hypothesis,  $h_{\text{final}}$ , labels  $x \in X$  according to the majority vote among  $h^{(1)}(x), h^{(2)}(x), \dots, h^{(T)}(x)$ .

Let  $E$  be the set of  $x \in X$  incorrectly labeled by  $h_{\text{final}}$ . The total cost of each  $x \in E$ ,

$$\sum_t \mathbf{m}_x^{(t)} = \sum_t 1 - |h^{(t)}(x) - c(x)| \leq \frac{T}{2}$$

since the majority vote gives an incorrect label for it. Instead of applying Theorem 2, we apply Theorem 5 which gives a more nuanced bound. The set  $\mathcal{P}$  in this theorem is simply the set of all possible distributions on  $X$ . Choosing  $\mathbf{p}$  to be the uniform distribution on  $E$ , we get

$$\left(\frac{1}{2} + \gamma\right) T \leq \sum_t \mathbf{m}^{(t)} \cdot \mathbf{p}^{(t)} \leq (1 + \eta) \sum_t \mathbf{m}^{(t)} \cdot \mathbf{p} + \frac{\mathbf{RE}(\mathbf{p} \parallel \mathbf{p}^{(1)})}{\eta} \leq (1 + \eta) \frac{T}{2} + \frac{\log(n/|E|)}{\eta}.$$

Since  $\eta = \gamma$  and  $T = \left\lceil \frac{2}{\gamma^2} \ln\left(\frac{1}{\varepsilon}\right) \right\rceil$ , the above inequality implies that the fraction of errors,  $\frac{|E|}{n}$ , is at most  $\varepsilon$  as desired.

### 3.7 Hard-core sets and the XOR Lemma

A boolean function  $f : X \rightarrow \{0, 1\}$ , where  $X$  is a finite domain, is  $\varepsilon$ -strongly hard, for circuits of size  $S$  if for every circuit  $C$  of size at most  $S$ ,

$$\Pr_x[C(x) = f(x)] \leq \frac{1}{2} + \gamma.$$

Here  $x \in X$  is drawn uniformly at random, and  $\gamma < 1/2$  is a parameter. For some parameter  $\varepsilon > 0$ , it is  $\varepsilon$ -weakly hard for circuits of size  $S$  if for every circuit  $C$  of size at most  $S$ , we have

$$\Pr_x[C(x) = f(x)] \leq 1 - \varepsilon.$$

Now given  $f : \{0, 1\}^n \rightarrow \{0, 1\}$ , define  $f^{\oplus k} : \{0, 1\}^{nk} \rightarrow \{0, 1\}$  to be the boolean function obtained by dividing up the input  $nk$ -bit string into  $k$  blocks of  $n$  bits each in the natural way, applying  $f$  to each block in turn, and taking the XOR of the  $k$  outputs. Yao's XOR Lemma [Yao82] shows that if  $f$  is  $\varepsilon$ -weakly hard against circuits of size  $S$  then  $f^{\oplus k}$  is  $\gamma + (1 - \varepsilon)^k$ -strongly hard for circuits of size  $S\varepsilon^2\gamma^2/8$ .

The original proofs were difficult but Impagliazzo [Imp95] suggested a simpler proof that as a byproduct proves an interesting fact about weakly hard functions: there is a reasonably large subset (at least  $\varepsilon$  fraction of  $X$ ) of inputs on which the function behaves like a strongly hard function, for somewhat smaller circuit size  $S'$ . This subset is called a *hard-core set*. For technical reasons, to prove such a result, it suffices to exhibit a “smooth” distribution  $\mathbf{p}$  on  $X$  (precise definition given momentarily), such for any circuit  $C$  of size at most  $S'$ , we have

$$\Pr_{x \sim \mathbf{p}}[C(x) = f(x)] \leq \frac{1}{2} + \gamma.$$

An  $\varepsilon$ -smooth distribution  $\mathbf{p}$  (the same  $\varepsilon$  from the weak hardness assumption on  $f$ ) is one that doesn't assign too much probability to any single input:  $p_x \leq \frac{1}{\varepsilon|X|}$  for any  $x \in X$ .

Such a distribution be decomposed as a convex combination of probability distributions over subsets of size at least  $\varepsilon|X|$ .

Klivans and Servedio [KS03] observed that Impagliazzo’s proof is an application of a boosting algorithm. The argument is as follows. We are given a boolean function  $f$  that is  $\varepsilon$ -weakly hard for circuits of size  $S$ . Assume for the sake of contradiction  $f$  is not strongly hard on any smooth distribution for circuits of some size  $S' < S$ . Then for any smooth distribution, we can find a small circuit of size  $S'$  that calculates  $f$  correctly with probability better than  $1/2 + \varepsilon$  when inputs are drawn from the distribution. Treat this as a weak learning algorithm, and apply boosting. Boosting combines the small circuits of size  $S'$  found by the weak learning algorithm into a larger circuit that calculates  $f$  correctly with probability at least  $1 - \varepsilon$  on the uniform distribution on  $X$ , contradicting the fact  $f$  is  $\varepsilon$ -weakly hard, if we ensure that the size of the larger circuit is smaller than  $S$ . This can be done if the  $S'$  is set to  $O(S/T)$ , where  $T$  is the number of boosting rounds.

With this insight, the problem now boils down to designing boosting algorithms that (a) are able to deal with smooth distributions on inputs and (b) have a small number of boosting rounds. The lower the number of boosting rounds, the better circuit size bound we get for showing  $\gamma$ -strong hardness.

The third author [Kal07] has shown how to construct such a boosting algorithm using the MW algorithm for restricted distributions (see Section 2.2). This boosting algorithm obtains the best known parameters in hard-core set constructions directly without having to resort to composing two different boosting algorithms as in [KS03]. This technique was extended in [BHK09] to obtain *uniform* constructions of hard-core sets with the best known parameters.

We describe the boosting algorithm of [Kal07] now. The main observation is that the set of all  $\varepsilon$ -smooth distributions is convex. Call this set  $\mathcal{P}$ . Then, exactly as in Section 3.6, the boosting algorithm simply runs the MW algorithm, with the only difference being that the distributions it generates are restricted to be in  $\mathcal{P}$  using relative entropy projections, as in the algorithm of Section 2.2. We can now apply the same analysis as in Section 3.6. Following this analysis, let  $E \subseteq X$  be the set of inputs on which the final majority hypotheses incorrectly computes  $f$ . If  $|E| \geq \varepsilon|X|$ , then since the uniform distribution on  $E$  is  $\varepsilon$ -smooth, we obtain a contradiction if  $T = \left\lceil \frac{2}{\gamma^2} \ln\left(\frac{1}{\varepsilon}\right) \right\rceil$  as before. Thus, the final majority hypothesis computes  $f$  correctly on at least a  $1 - \varepsilon$  fraction of inputs from  $X$ .

This immediately implies the following hard-core set existence result, which has the best known parameters to date:

**Theorem 14** *Given a function  $f : X \rightarrow \{0, 1\}$  that is  $\gamma$ -weakly hard for circuits of size  $S$ , there is a subset of  $X$  of size at least  $\varepsilon|X|$  on which  $f$  is  $\varepsilon$ -strongly hard for circuits of size  $S' = O\left(\frac{\gamma^2}{\log(1/\varepsilon)}S\right)$ .*

### 3.8 Hannan’s algorithm and multiplicative weights

Perhaps the earliest decision-making algorithm which attains bounds similar to the MW algorithm is Hannan’s algorithm [Han57], dubbed “follow the perturbed leader” by Kalai and Vempala [KV03]. This algorithm is given below.

**Initialization:** Pick random numbers  $r_1, r_2, \dots, r_n$ , one for each decision.

**For**  $t = 1, 2, \dots, T$ :

1. Choose the decision which minimizes the total cost including random initial cost:

$$i^{(t)} = \arg \min_i \left\{ L_i^{(t)} + r_i \right\} \quad (23)$$

Where  $L_i^{(t)} = \sum_{\tau < t} m_i^{(\tau)}$  is the total cost so far for the  $i$ th decision.

2. Observe the costs of the decisions  $\mathbf{m}^{(t)}$ .

Figure 3: Hannan's algorithm.

In this section we show that for a particular choice of random initialization numbers  $\{r_i\}$ , the algorithm above exactly reproduces the multiplicative weights algorithm, or more precisely the Hedge algorithm as in Section 2.1. This observation is due to Adam Kalai [Kal09].

**Theorem 15 ([Kal09])** *Let  $u_1, \dots, u_n$  be  $n$  independent random numbers chosen uniformly from  $[0, 1]$ , and consider the algorithm above with  $r_i = \frac{1}{\eta} \ln \ln \frac{1}{u_i}$ . Then for any decision  $j$ , we have*

$$\Pr \left[ i^{(t)} = j \right] = \frac{e^{-\eta L_j^{(t)}}}{\sum_i e^{-\eta L_i^{(t)}}}.$$

PROOF: By monotonicity of the exponential function, we have:

$$\begin{aligned} \arg \min_i \left\{ L_i^{(t)} + \frac{1}{\eta} \ln \ln \frac{1}{u_i} \right\} &= \arg \min_i \left\{ \eta L_i^{(t)} + \ln \ln \frac{1}{u_i} \right\} \\ &= \arg \max_i \left\{ \frac{e^{-\eta L_i^{(t)}}}{\ln \frac{1}{u_i}} \right\} \\ &= \arg \max_i \left\{ \frac{w_i^{(t)}}{s_i} \right\}. \end{aligned}$$

Where  $w_i^{(t)} = e^{-\eta L_i^{(t)}}$  and  $s_i = \ln \frac{1}{u_i}$  are independent exponentially distributed random variables with mean 1. The result now follows from Lemma 16 below.  $\square$

**Lemma 16** *Let  $w_1, \dots, w_n$  be arbitrary non-negative real numbers, and let  $s_1, \dots, s_n$  be independent exponential random variables with mean 1. Then*

$$\Pr \left[ \arg \max_i \frac{w_i}{s_i} = j \right] = \frac{w_j}{\sum_i w_i}.$$

PROOF: The probability density function of  $s_i$  is  $e^{-s_i}$ . Conditioning on the value of  $s_j$ , for a particular  $i \neq j$  we have

$$\Pr \left[ \frac{w_i}{s_i} \leq \frac{w_j}{s_j} \mid s_j \right] = \int_{s_i = \frac{w_i s_j}{w_j}}^{\infty} e^{-s_i} ds_i = e^{-\frac{w_i s_j}{w_j}}.$$

Since all variables  $s_1, s_2, \dots, s_n$  are independent, we have

$$\Pr \left[ \forall i \neq j : \frac{w_i}{s_i} \leq \frac{w_j}{s_j} \mid s_j \right] = e^{-\sum_{i \neq j} \frac{w_i s_j}{w_j}}.$$

Integrating to remove the conditioning on  $s_j$ , we have

$$\Pr \left[ \forall i \neq j : \frac{w_i}{s_i} \leq \frac{w_j}{s_j} \right] = \int_{s_j=0}^{\infty} e^{-\sum_{i \neq j} \frac{w_i s_j}{w_j}} e^{-s_j} ds_j = \int_{s_j=0}^{\infty} e^{-s_j \frac{\sum_i w_i}{w_j}} ds_j = \frac{w_j}{\sum_i w_i}.$$

□

### 3.9 Online convex optimization

Online convex optimization is a very general framework that can be applied to many of the problems discussed in the applications section and many more. Here “online” means that the algorithm does not know the entire input at the start, and the input is presented to it in pieces over many rounds. In this section we describe the framework and the central role of the multiplicative weights method. For a much more detailed treatment of online learning techniques see [CBL06].

In online convex optimization, we move from a discrete decision set to a continuous one. Specifically, the set of decisions is a convex, compact set  $K \subseteq \mathbb{R}^n$ . In each round  $t = 1, 2, \dots$ , the online algorithm is required to choose a decision, i.e. point  $\mathbf{p}^{(t)} \in K$ . A convex loss function  $f^{(t)}$  is presented, and the decision maker incurs a loss of  $f^{(t)}(\mathbf{p}^{(t)})$ . The goal of the online algorithm  $\mathcal{A}$  is to minimize loss compared to the best fixed offline strategy. This quantity is called *regret* in the game theory and machine learning literature.

$$\mathcal{R}_T(\mathcal{A}) = \sum_{t=1}^T f^{(t)}(\mathbf{p}^{(t)}) - \min_{\mathbf{p} \in K} \sum_{t=1}^T f^{(t)}(\mathbf{p}).$$

The basic decision-making problem described in Section 2 with  $n$  discrete decisions is recovered as a special case of online convex optimization as follows. The convex set  $K$  is the  $n$ -dimensional simplex corresponding to the set of all distributions over the  $n$  decisions, and the payoff functions  $f^{(t)}$  are defined as  $f^{(t)}(\mathbf{p}^{(t)}) = \mathbf{m}^{(t)} \cdot \mathbf{p}^{(t)}$  given the cost vector  $\mathbf{m}^{(t)}$ . It also generalizes other online learning problems such as the online portfolio selection problem and online routing (see [Haz06] for more discussion on applications). Zinkevich [Zin03] gives algorithms for the general online convex optimization problem. More efficient algorithms for online convex optimization based on strong convexity of the loss functions have also been developed [HAK07].

We now describe how to use the Multiplicative Weights algorithm to the online convex optimization problem for the special case where  $K$  is the  $n$  dimensional simplex of distributions on the coordinates. The advantage is that this algorithm has much better dependence on the dimension  $n$  than Zinkevich’s algorithm (see [Haz06] for more details). This algorithm has several applications such as in online portfolio selection [HSSW96]. Here is the algorithm. First, define

$$\rho := \max_{\mathbf{p} \in K} \max_t \|\nabla f^{(t)}(\mathbf{p})\|_{\infty},$$

where  $\nabla f^{(t)}(\mathbf{p})$  is the (sub)gradient of the function  $f^{(t)}$  at point  $\mathbf{p}$ . The parameter  $\rho$  is called the *width* of the problem. Then, run the standard Multiplicative Weights algorithm with  $\eta = \sqrt{\frac{\ln n}{T}}$  and the costs defined as

$$\mathbf{m}^{(t)} := \frac{1}{\rho} \nabla f^{(t)}(\mathbf{p}^{(t)})$$

where  $\mathbf{p}^{(t)}$  is the point played in round  $t$ . Note that for all  $t$  and all  $i$ ,  $|m_i^{(t)}| \leq 1$  as required by the MW algorithm.

**Theorem 17** *After  $T$  rounds of applying the Multiplicative Weights algorithm to the online convex optimization framework, for any  $\mathbf{p} \in K$  we have*

$$\sum_{t=1}^T f^{(t)}(\mathbf{p}^{(t)}) - \sum_{t=1}^T f^{(t)}(\mathbf{p}) \leq 2\rho\sqrt{\ln(n)T}.$$

PROOF: If  $f : K \rightarrow \mathbb{R}$  is a differentiable convex function, then for any two points  $\mathbf{p}, \mathbf{q} \in K$  we have

$$f(\mathbf{q}) \geq f(\mathbf{p}) + \nabla f(\mathbf{p}) \cdot (\mathbf{q} - \mathbf{p}),$$

where  $\nabla f(\mathbf{p})$  is the gradient of  $f$  at  $\mathbf{p}$ . Rearranging we get

$$f(\mathbf{p}) - f(\mathbf{q}) \leq \nabla f(\mathbf{p}) \cdot (\mathbf{p} - \mathbf{q}). \quad (24)$$

Applying Corollary 3, we get that for any  $\mathbf{p} \in K$ ,

$$\sum_{t=1}^T \mathbf{m}^{(t)} \cdot \mathbf{p}^{(t)} \leq \sum_{t=1}^T (\mathbf{m}^{(t)} + \eta |\mathbf{m}^{(t)}|) \cdot \mathbf{p} + \frac{\ln n}{\eta} \leq \sum_{t=1}^T \mathbf{m}^{(t)} \cdot \mathbf{p} + \eta + \frac{\ln n}{\eta}, \quad (25)$$

since  $\|\mathbf{m}^{(t)}\|_\infty \leq 1$ . Now we have

$$\begin{aligned} \sum_t [f^{(t)}(\mathbf{p}^{(t)}) - f^{(t)}(\mathbf{p})] &\leq \sum_{t=1}^T \nabla f^{(t)}(\mathbf{p}^{(t)}) \cdot (\mathbf{p}^{(t)} - \mathbf{p}) && \text{(from (24))} \\ &= \sum_{t=1}^T \rho \mathbf{m}^{(t)} \cdot (\mathbf{p}^{(t)} - \mathbf{p}) \\ &\leq \eta \rho T + \frac{\rho \ln n}{\eta}. && \text{(from (25))} \end{aligned}$$

Substituting  $\eta = \sqrt{\frac{\ln n}{T}}$  completes the proof.  $\square$

### 3.10 Other applications

#### 3.10.1 Approximately Solving certain Semidefinite Programs

*Semidefinite programming* (SDP) is a special case of convex programming. A semidefinite program is derived from a linear program by imposing the additional constraint that some

subset of  $n^2$  variables form an  $n \times n$  positive semidefinite matrix. Since the work of Goemans and Williamson [GW95], SDP has become an important tool in design of approximation algorithms for **NP**-hard optimization problems. Though this yields polynomial-time algorithms, their practicality is suspect because solving SDPs is a slow process. Therefore there is great interest in computing approximate solutions to SDPs, especially the types that arise in approximation algorithms. Since the SDPs in question are being used to design approximation algorithms anyway, it is permissible to compute approximate solutions to these SDPs.

Klein and Lu [KL96] use the PST framework (Section 3.3) to derive a more efficient 0.878-approximation algorithm for MAX-CUT than the original SDP-based method in Goemans-Williamson [GW95]. The main idea in Klein-Lu is to approximately solve the MAX-CUT SDP. However, their idea does work very well for other SDPs. The main issue is that the width  $\rho$  (see 3.3) is too high for certain SDPs of interest.

To be more precise, an SDP feasibility problem is given by:

$$\begin{aligned} \forall j = 1, 2, \dots, m : \mathbf{A}_j \bullet \mathbf{X} &\geq b_j \\ \mathbf{X} &\in \mathcal{P} \end{aligned}$$

Here, we use the notation  $\mathbf{A} \bullet \mathbf{B} = \sum_{ij} A_{ij} B_{ij}$  to denote the scalar product of two matrices thinking of them as vectors in  $\mathbb{R}^{n^2}$ . The set  $\mathcal{P} = \{\mathbf{X} \in \mathbb{R}^{n \times n} \mid \mathbf{X} \succeq 0, \text{Tr}(\mathbf{X}) \leq 1\}$  is the set of all positive semidefinite matrices with trace bounded by 1.

The Plotkin-Shmoys-Tardos framework (see Section 3.3) is suitable for approximating SDPs since all constraints are linear, and the oracle given in (15) can be implemented efficiently by an eigenvector computation. To see this, note that the oracle needs to decide, given a probability distribution  $\mathbf{p}$  on the constraints, if there exists an  $\mathbf{X} \in \mathcal{P}$  such that  $\sum_j p_j \mathbf{A}_j \bullet \mathbf{X} \geq \sum_j p_j b_j$ . This can be implemented by solving the following optimization problem:

$$\max_{\mathbf{X} \in \mathcal{P}} : \sum_j p_j \mathbf{A}_j \bullet \mathbf{X}$$

It is easily checked that an optimal solution to the above optimization problem is given by the matrix  $\mathbf{X} = \mathbf{v}\mathbf{v}^\top$  where  $\mathbf{v}$  is unit eigenvector corresponding to the largest eigenvalue of the matrix  $\sum_j p_j \mathbf{A}_j$ .

The Klein-Lu approach was of limited uses in many cases because it does not do too well when the additive error  $\varepsilon$  is required to be small. (They were interested in the MAX-CUT problem, where this problem does not arise. The reason in a nutshell is that in a graph with  $m$  edges, the maximum cut has at least  $m/2$  edges, so it suffices to compute the optimum to an additive error  $\varepsilon$  that is a fixed constant.) We have managed to extend the multiplicative weights framework to many of these settings to design efficient algorithms for SDP relaxations of many other problems. The main idea is to apply the Multiplicative Weights framework in a “nested” fashion: one can solve a constrained optimization problem by invoking the MW algorithm on an subset of constraints (the “outer” constraints) in the manner of 3.3, where the domain is now defined by the rest of the constraints (the “inner” constraints). The oracle can now be implemented by another application of the MW algorithm on the inner constraints. Alternatively, we can reduce the dependence on the width by using the observation that the Lagrangian relaxation problem on the inner

constraints can be solved by the ellipsoid method. For details of this method, refer to our paper [AHK05]. For several families of SDPs we obtain the best running time known.

More recently Arora and Kale [AK07] have designed a new approach for solving SDPs that involves a variant of the multiplicative update rule at the matrix level; see Section 5 for details.

Yet another approach for approximately solving SDPs is by reducing the SDP into a maximization problem of a single concave function over the PSD cone. The latter problem can be approximated efficiently via an iterative greedy method. The resulting algorithm is extremely similar to the MW-based algorithm of [AHK05], however its analysis is very different, see [Haz08] for more details.

### 3.10.2 Approximating Graph Separators

A recent application of the multiplicative weights method is a combinatorial algorithm for approximating the SPARSEST CUT of graphs [AHK04]. This is a fundamental graph partitioning problem. Given a graph  $G = (V, E)$ , the expansion of a cut  $(S, \bar{S})$  where  $S \subseteq V$  and  $\bar{S} = V \setminus S$ , is defined to be

$$\frac{|E(S, \bar{S})|}{\min\{|S|, |\bar{S}|\}}.$$

Here,  $E(S, \bar{S})$  is the set of edges with one end point in  $S$  and the other in  $\bar{S}$ . The SPARSEST CUT problem is to find the cut in the input graph of minimum expansion. This problem arises as a useful subroutine in many other algorithms, such as in divide-and-conquer algorithms for optimization problems on graphs, layout problems, clustering, etc. Furthermore, the expansion of a graph is a very useful way to quantify the connectivity of a graph and has many important applications in computer science.

The work of Arora, Rao and Vazirani [ARV04] gave the first  $O(\sqrt{\log n})$  approximation algorithm to the SPARSEST CUT problem. However, their best algorithm relies on solving an SDP and runs in  $\tilde{O}(n^{4.5})$  time. They also gave an alternative algorithm based on the notion of *expander flows*, which are multicommodity flows in the graph whose demand graph has high expansion. However, their algorithm was based on the ellipsoid method, and was thus quite inefficient. In the paper [AHK04], we obtained a much more efficient algorithm for approximating the SPARSEST CUT problem to an  $O(\sqrt{\log n})$  factor in  $\tilde{O}(n^2)$  time using the expander flow idea. The algorithm casts the problem of routing an expander flow in the graph as a linear program, and then checks the feasibility of the linear program using the techniques described in Section 3.3. The oracle for this purpose is implemented using a variety of techniques: the multicommodity flow algorithm of Garg and Könemann [GK98] (and its subsequent improvement by Fleischer [Fle00]), eigenvalue computations, and graph sparsification algorithms of Benczúr and Karger [BK96] based on random sampling.

### 3.10.3 Multiplicative weight algorithms in geometry

The multiplicative weight idea has been used several times in computational geometry. Chazelle [Cha00] (p. 6) describes the main idea, which is essentially the connection between derandomization of Chernoff bound arguments and the exponential potential function noted in Section 3.5.

The geometric applications consist of derandomizing the a natural randomized algorithm by using a deterministic construction of some kind of small set cover or low-discrepancy set. Formally, the analysis is similar to our analysis of the Set Cover algorithm in Section 3.5. Clarkson used this idea to give a deterministic algorithm for Linear Programming [Cla88]. Following Clarkson, Bronnimann and Goodrich use similar methods to find Set Covers for hyper-graphs with small VC dimension [BG94].

The MW algorithm was also used in the context of geometric embeddings of finite metric spaces, specifically, embedding negative-type metrics (i.e. a set of points in Euclidean space such that the squared Euclidean distance between them also forms a metric) into  $\ell_1$ . Such embeddings are important for approximating the important non-uniform SPARSEST CUT PROBLEM.

The approximation algorithm for SPARSEST CUT in Arora et al. [ARV04] involves a “Structure Theorem”. This structure theorem was interpreted by Chawla et al. [CGR05] as saying that any  $n$ -point negative-type metric with maximum distance  $O(1)$  and average distance  $\Omega(1)$  can be embedded into  $\ell_1$  such that *average*  $\ell_1$  distance is  $\Omega(1/\sqrt{\log n})$ . Then they used the MW algorithm to construct an embedding into  $\ell_1$  in which *every* pair of points that have negative-type metric distance  $\Omega(1)$  have  $\ell_1$  distance that is off by at most an  $O(\sqrt{\log n})$  factor of the original. Using a similar idea for other distance and combining the resulting embeddings they obtained an embedding of the negative-type metric into  $\ell_1$  in which every distance distorts by at most a factor  $O(\log^{3/4} n)$ . Arora et al. [ALN05] gave a more complicated construction to improve the distortion bound to  $O(\sqrt{\log(n)} \log \log n)$ , leading to a  $O(\sqrt{\log(n)} \log \log n)$ -approximation for non-uniform sparsest cut.

### 3.11 Design of Competitive Online Algorithms

Starting with the work of Alon, Awerbuch, Azar, Buchbinder and Naor [AAA<sup>+</sup>09], a number of competitive online algorithms have been developed using an elegant primal-dual approach which involves multiplicative weight updates. While the analysis of their algorithms seems to beyond our general framework, we briefly mention this work without going into many details. We refer the readers to the survey by Buchbinder and Naor [BN09b] for an extensive discussion of the topic.

Several online problems such as the ski rental problem, caching, load balancing, ad auctions, etc. can be cast (in their fractional form) as a linear program with non-negative coefficients in the constraints and cost, where either the constraints or variables arrive online one by one. The online problem is to maintain a feasible solution at all times with a bounded competitive ratio, i.e. ensuring that the cost of the solution maintained is bounded in terms of the cost of the optimal solution in each round. The main difficulty comes from the requirement that the solution maintained is monotonic in some sense (for example, the variables are never allowed to decrease).

Buchbinder and Naor [BN09a] give a algorithm based on the primal-dual method that obtains good competitive ratios in this scenario. At the heart is a multiplicative weight update rule. Imagine we have a covering LP, i.e. all constraints are of the form  $\mathbf{a} \cdot \mathbf{x} \geq 1$ , where  $\mathbf{a}$  has non-negative coefficients, and  $\mathbf{x}$  is the vector of variables. The cost function has non-negative coefficients as well. Constraints arrive one at a time in each round and must be satisfied by the current solution. The requirement is that no variable can decrease



from round to round.

In every round, the algorithm increases primal variables in the new constraint using multiplicative updates and the corresponding dual variable additively (so essentially, each primal variable is the exponential of the dual constraint that it corresponds to, much like the Plotkin-Shmoys-Tardos algorithms of Section 3.3). This is done until the constraint gets satisfied. Clearly, we maintain a feasible solution in each round, and the variables never decrease. The analysis goes by relating the increase in the primal cost can be bounded in terms of the dual cost (this is an easy consequence of the multiplicative update, in fact, the multiplicative update can be derived using this requirement). The competitive ratio is obtained by showing that the dual solution generated simultaneously, while infeasible, is not far from being feasible, i.e. scaling down the variables by a certain factor makes it feasible. This gives us a bound on the competitive ratio via weak duality.

## 4 Lower bound

Can our analysis of the Multiplicative Weights algorithm be improved? This section shows that the answer is no, not only for the MW algorithm itself, but for any algorithm which operates in the online setting we have described.

Technically, we prove the following:

**Theorem 18** *For any online decision-making algorithm with  $n \geq 2$  decisions, there exists a sequence of cost vectors  $\mathbf{m}^{(1)}, \mathbf{m}^{(2)}, \dots, \mathbf{m}^{(T)} \in \{0, 1\}^n$  such that  $\min_i \sum_t m_i^{(t)} = \Omega(T)$ , and if the sequence of distributions over decisions produced by the algorithm are  $\mathbf{p}^{(1)}, \mathbf{p}^{(2)}, \dots, \mathbf{p}^{(T)}$ , then we have*

$$\sum_t \mathbf{m}^{(t)} \cdot \mathbf{p}^{(t)} \geq \min_i \sum_{t=1}^T m_i^{(t)} + \Omega(\sqrt{T \ln n})$$

Since in the theorem above we have  $\min_i \sum_{t=1}^T m_i^{(t)} = \Omega(T)$ , Theorem 2 implies that by choosing the optimal value of  $\eta$ , viz.  $\eta = \Theta\left(\sqrt{\frac{\log n}{T}}\right)$  we have

$$\sum_{t=1}^T \mathbf{m}^{(t)} \cdot \mathbf{p}^{(t)} \leq \min_i \sum_{t=1}^T m_i^{(t)} + O(\sqrt{T \ln n})$$

Hence our analysis is tight up to constants in the additive error term. Moreover, the above lower bound applies to *any* algorithm, efficient or not.

PROOF: The proof is via the probabilistic method. We construct a distribution over the costs so that the required bound is obtained in expectation. Interestingly, the distribution is independent of the actual algorithm used.

We now specify the costs of the decisions. The cost of decision 1 is set to 1/2 for all rounds  $t$ . For any decision  $i > 1$ , we construct its cost via the following random process: in each iteration  $t$ , choose its cost to be either one or zero uniformly at random, i.e.  $m_i^{(t)} \in \{0, 1\}$  with probability of each outcome being 1/2.

The expected cost of each decision is  $1/2$ . Hence, the expected cost of the chosen decision is also  $1/2$  irrespective of the algorithm's distribution  $\mathbf{p}^{(t)}$ , and hence:

$$\mathbf{E} \left[ \sum_{t=1}^T \mathbf{m}^{(t)} \cdot \mathbf{p}^{(t)} \right] = \frac{T}{2}. \quad (26)$$

For every decision  $i$ , define  $X_i = \sum_{t=1}^T m_i^{(t)}$ . Note that  $X_1 = T/2$ , and for  $i > 1$ , we have  $X_i \sim B(T, \frac{1}{2})$ , where  $B(T, \frac{1}{2})$  is the binomial distribution with  $T$  trials and both outcomes equally likely.

We use the following standard concentration Lemma (see Proposition 7.3.2 in [MV08]):

**Lemma 19** *Let  $X \sim B(T, \frac{1}{2})$  be a binomial random variable with  $T$  trials and probability  $\frac{1}{2}$ . Then for  $t \in [0, \frac{T}{8}]$ :*

$$\Pr \left[ X \leq \frac{T}{2} - t \right] \geq \frac{1}{15} e^{-16t^2/T}.$$

We claim that the expectation of the cost of the best decision, viz.  $\mathbf{E}[\min_i X_i]$ , in our construction is  $\frac{T}{2} - \Omega(\sqrt{T \log n})$ . Let  $t = \frac{1}{4} \sqrt{T \ln(n-1)}$ . We have

$$\Pr \left[ \min_i X_i > \frac{T}{2} - t \right] = \prod_{i=2}^n \Pr \left[ X_i > \frac{T}{2} - t \right] \leq \left( 1 - \frac{1}{15} e^{-16t^2/T} \right)^{n-1} \leq e^{-1/15} < 0.95.$$

The first equality above is by the independence of the  $X_i$ , and the first inequality is by Lemma 19. Thus, with probability at least  $1 - 0.95 = 0.05$ ,  $\min_i X_i \leq \frac{T}{2} - t$ . Since  $\min_i X_i$  is always at least  $T/2$ , we get

$$\mathbf{E}[\min_i X_i] \geq 0.95 \cdot \frac{T}{2} + 0.05 \cdot \left( \frac{T}{2} + t \right) = \frac{T}{2} + \frac{1}{80} \sqrt{T \ln(n-1)}.$$

It follows that

$$\mathbf{E} \left[ \min_i \sum_{t=1}^T m_i^{(t)} \right] \geq \frac{T}{2} - \frac{1}{80} \sqrt{T \ln(n-1)}. \quad (27)$$

From (26) and (27), we have

$$\mathbf{E} \left[ \sum_{t=1}^T \mathbf{m}^{(t)} \cdot \mathbf{p}^{(t)} - \min_i \sum_{t=1}^T m_i^{(t)} \right] \geq \frac{1}{80} \sqrt{T \ln(n-1)}.$$

Since  $\sum_{t=1}^T \mathbf{m}^{(t)} \cdot \mathbf{p}^{(t)} - \min_i \sum_{t=1}^T m_i^{(t)} \leq T$ , by Markov's inequality we conclude that

$$\Pr \left[ \sum_{t=1}^T \mathbf{m}^{(t)} \cdot \mathbf{p}^{(t)} - \min_i \sum_{t=1}^T m_i^{(t)} > \frac{1}{160} \sqrt{T \ln(n-1)} \right] \geq \frac{1}{160} \sqrt{\frac{\ln(n-1)}{T}}.$$

By the Hoeffding bound [AS92], for any decision  $i > 1$ , we have

$$\Pr \left[ \sum_{t=1}^T m_i^{(t)} < \frac{T}{4} \right] \leq \exp(-T/32).$$

By the union bound, we have

$$\Pr \left[ \exists i : \sum_{t=1}^T m_i^{(t)} < \frac{T}{4} \right] \leq n \exp(-T/32) < \frac{1}{160} \sqrt{\frac{\ln(n-1)}{T}},$$

for large enough  $T$ . This implies that there exists a sequence  $\mathbf{m}^{(1)}, \mathbf{m}^{(2)}, \dots, \mathbf{m}^{(T)}$  for which  $\min_i \sum_{t=1}^T m_i^{(t)} \geq \frac{T}{4}$  and

$$\sum_{t=1}^T \mathbf{m}^{(t)} \cdot \mathbf{p}^{(t)} - \min_i \sum_{t=1}^T m_i^{(t)} \geq \frac{1}{160} \sqrt{T \ln(n-1)}.$$

□

## 5 The Matrix Multiplicative Weights algorithm

In the preceding sections, we considered an online decision making problem. We refer to that setting as the *basic* or *scalar* case. Now we briefly consider a different online decision making problem, which is seemingly quite different from the previous one, but has enough structure that we can obtain an analogous algorithm for it. We move from cost vectors to cost matrices, and from probability vectors to density matrices. For this reason, we refer to the current setting as the *matrix* case. We call the algorithm presented in this setting the *Matrix Multiplicative Weights algorithm*. The original motivation for our interest in this matrix setting is that it leads to a constructive version of SDP duality, just as the standard multiplicative weights algorithm can be viewed as a constructive version of LP duality. In fact the standard algorithm is a special subcase of the algorithm in this section, namely, when all the matrices involved are diagonal.

Applications of the matrix multiplicative weights algorithm include solving SDPs [AK07], derandomizing constructions of expander graphs, and obtaining bounds on the sample complexity for a learning problem in quantum computing. The details of these applications are beyond the scope of this survey; please see the third author's PhD thesis [Kal06] for details. The algorithm given here is from a paper of Arora and Kale [AK07]. A very similar algorithm was discovered independently slightly earlier by Warmuth and Kuzmin [WK06], and is based on the even earlier work of Tsuda, Rätsch and Warmuth [TRW05].

We stick with our basic decision-making scenario but decisions now correspond to unit vectors  $\mathbf{v}$  in  $\mathbb{S}^{n-1}$ , the unit sphere in  $\mathbb{R}^n$ . As in the basic case, in every round, our task is to pick a decision  $\mathbf{v} \in \mathbb{S}^{n-1}$ . At this point, the costs of all decisions are revealed by nature. These costs are not arbitrary, but they are correlated in the following way. A *cost matrix*  $\mathbf{M} \in \mathbb{R}^{n \times n}$  is revealed, and the cost of a decision  $\mathbf{v}$  is then  $\mathbf{v}^\top \mathbf{M} \mathbf{v}$ . We assume that the costs of all decisions lie in  $[-1, 1]$ . Again, as in the basic case, this is the only assumption we make on the way nature chooses the costs; indeed, the costs could even be chosen adversarially. Equivalently, we assume that all the eigenvalues of the matrix  $\mathbf{M}$  are in the range  $[-1, 1]$ .

This game is repeated over a number of rounds. Let  $t = 1, 2, \dots, T$  denote the current round. In each round  $t$ , we select a distribution  $\mathcal{D}^{(t)}$  over the set of decisions  $\mathbb{S}^{n-1}$ , and select a decision  $\mathbf{v}$  randomly from it (and use his advised course of action). At this point,

the costs of all decisions are revealed by nature via the cost matrix  $\mathbf{M}^{(t)}$ . The expected cost to the algorithm for choosing the distribution  $\mathcal{D}^{(t)}$  is

$$\mathbf{E}_{\mathbf{v} \in \mathcal{D}^{(t)}}[\mathbf{v}^\top \mathbf{M}^{(t)} \mathbf{v}] = \mathbf{E}_{\mathbf{v} \in \mathcal{D}^{(t)}}[\mathbf{M}^{(t)} \bullet \mathbf{v} \mathbf{v}^\top] = \mathbf{M}^{(t)} \bullet \mathbf{E}_{\mathbf{v} \in \mathcal{D}^{(t)}}[\mathbf{v} \mathbf{v}^\top].$$

Recall that we are using the notation  $\mathbf{A} \bullet \mathbf{B} = \sum_{ij} A_{ij} B_{ij}$  to denote the scalar product of two matrices thinking of them as vectors in  $\mathbb{R}^{n^2}$ . Define the matrix  $\mathbf{P}^{(t)} := \mathbf{E}_{\mathbf{v} \in \mathcal{D}^{(t)}}[\mathbf{v} \mathbf{v}^\top]$ . Note that  $\mathbf{P}^{(t)}$  is positive semidefinite: this is because it is a convex combination of the elementary positive semidefinite matrices  $\mathbf{v} \mathbf{v}^\top$ . Let  $\text{Tr}(\mathbf{A})$  denote the trace of a matrix  $\mathbf{A}$ . Then,  $\text{Tr}(\mathbf{P}^{(t)}) = 1$ , again because for all  $\mathbf{v}$ , we have  $\text{Tr}(\mathbf{v} \mathbf{v}^\top) = \|\mathbf{v}\|^2 = 1$ . A matrix  $\mathbf{P}$  which is positive semidefinite and has trace 1 is called a *density matrix*.

We will only be interested in the expected cost to the algorithm, and all the information required for computing the expected cost for a given distribution  $\mathcal{D}$  over  $\mathbb{S}^{n-1}$  is contained in the associated density matrix.

Thus, in each round  $t$ , we require our online algorithm to choose a density matrix  $\mathbf{P}^{(t)}$ , rather than a distribution  $\mathcal{D}^{(t)}$  over  $\mathbb{S}^{n-1}$ . The distribution is implicit in the choice of  $\mathbf{P}^{(t)}$ . The eigendecomposition of  $\mathbf{P}^{(t)} = \sum_{i=1}^n \lambda_i \mathbf{v}_i \mathbf{v}_i^\top$ , where  $\lambda_i$  is an eigenvalue corresponding to the unit eigenvector  $\mathbf{v}_i$ , gives one such distribution (a discrete one): vector  $\mathbf{v}_i$  gets probability  $\lambda_i$ . We then observe the cost matrix  $\mathbf{M}^{(t)}$  revealed by nature, and suffer the expected cost  $\mathbf{M}^{(t)} \bullet \mathbf{P}^{(t)}$ . After  $T$  rounds, the total expected cost is  $\sum_{t=1}^T \mathbf{M}^{(t)} \bullet \mathbf{P}^{(t)}$ , while the best fixed decision in hindsight corresponds to the unit vector  $\mathbf{v}$  which minimizes  $\sum_{t=1}^T \mathbf{v}^\top \mathbf{M}^{(t)} \mathbf{v}$ . Since we minimize this quantity over all unit vectors  $\mathbf{v}$ , the variational characterization of eigenvalues implies that this minimum cost is exactly the minimum eigenvalue of  $\sum_{t=1}^T \mathbf{M}^{(t)}$ , denoted by  $\lambda_n(\sum_{t=1}^T \mathbf{M}^{(t)})$ . Our goal is to design an online algorithm whose total expected cost over the  $T$  rounds is not much more than the cost of the best decision.

Consider the following generalization of the basic MW algorithm, called the Matrix Multiplicative Weights algorithm (Matrix MW). It uses the notion of matrix exponential,  $\exp(\mathbf{A}) := \sum_{i=0}^{\infty} \frac{\mathbf{A}^i}{i!}$ . The key point here is that regardless of the matrix  $\mathbf{A}$ , its exponential  $\exp(\mathbf{A})$  is always positive definite. Note that in case all matrices involved are diagonal, the Matrix Multiplicative Weights algorithm exactly reduces to the Hedge algorithm.

#### Matrix Multiplicative Weights algorithm

**Initialization:** Fix an  $\eta \leq \frac{1}{2}$ . Initialize the weight matrix  $\mathbf{W}^{(1)} = \mathbf{I}_n$ .

**For**  $t = 1, 2, \dots, T$ :

1. Use the density matrix  $\mathbf{P}^{(t)} = \frac{\mathbf{W}^{(t)}}{\Phi^{(t)}}$ , where  $\Phi^{(t)} = \text{Tr}(\mathbf{W}^{(t)})$ .
2. Observe the cost matrix  $\mathbf{M}^{(t)}$ .
3. Update the weight matrix as follows:

$$\mathbf{W}^{(t+1)} = \exp(-\eta \sum_{\tau=1}^t \mathbf{M}^{(\tau)}).$$

Figure 4: The Matrix Multiplicative Weights algorithm.

The following theorem bounds the total expected cost of the Matrix Multiplicative Weights algorithm (given in Figure 4) in terms of the cost of the best fixed decision. This theorem is completely analogous to Theorem 4, and in fact, Theorem 4 can be obtained directly from this theorem in the case when all matrices involved are diagonal. We omit the proof of this theorem; again, it is along the lines of the proof of Theorem 4, and uses  $\Phi^{(t)}$  as a potential function. The analysis is based on the matrix analogue of the real number inequality  $\exp(-\eta x) \leq 1 - \eta x + \eta^2 x^2$  for  $|\eta x| \leq 1$ : if  $\mathbf{X}$  is a matrix such that  $\|\eta \mathbf{X}\| \leq 1$ , then

$$-\exp(-\eta \mathbf{X}) \preceq \mathbf{I} - \eta \mathbf{X} + \eta^2 \mathbf{X}^2.$$

Also, we need an additional inequality from statistical mechanics, the Golden-Thompson inequality [Gol65, Tho65], which states that for two matrices  $\mathbf{A}$  and  $\mathbf{B}$ , we have

$$\text{Tr}(\exp(\mathbf{A} + \mathbf{B})) \leq \text{Tr}(\exp(\mathbf{A}) \exp(\mathbf{B})).$$

See [Kal06, AK07] for further details.

**Theorem 20** *In the given setup, the Matrix Multiplicative Weights algorithm guarantees that after  $T$  rounds, for any decision  $\mathbf{v}$ , we have*

$$\sum_{t=1}^T \mathbf{M}^{(t)} \bullet \mathbf{P}^{(t)} \leq \sum_{t=1}^T \mathbf{v}^\top \mathbf{M}^{(t)} \mathbf{v} + \eta \sum_{t=1}^T (\mathbf{M}^{(t)})^2 \bullet \mathbf{P}^{(t)} + \frac{\ln n}{\eta}. \quad (28)$$

## 5.1 Applications in Solving SDPs

In this section, we mention the application of the Matrix MW algorithm to solving SDPs without going into any details. We refer the interested reader to [Kal06, AK07] for further details and other applications.

Consider the following quite general form of feasibility problem using SDP:

$$\begin{aligned} \forall j = 1, 2, \dots, m : \quad & \mathbf{A}_j \bullet \mathbf{X} \geq 0 \\ & \text{Tr}(\mathbf{X}) = 1 \\ & \mathbf{X} \succeq \mathbf{0} \end{aligned} \quad (29)$$

Note if all  $\mathbf{A}_j$  matrices and  $\mathbf{X}$  were diagonal, then this exactly reduces to the LP (9).

Now suppose, as in Section 3.1, that there exists a large-margin solution, i.e. a density matrix  $\mathbf{X}^*$  such that for all  $j$ , we have  $\mathbf{A}_j \bullet \mathbf{X}^* \geq \varepsilon$ . Then just as in Section 3.1, we can use the Matrix MW algorithm to find a good solution, i.e. a density matrix  $\mathbf{X}$  such that for all  $j$  we have  $\mathbf{A}_j \bullet \mathbf{X} \geq 0$ . We now need to define  $\rho = \max_j \|\mathbf{A}_j\|$ .

We run the Matrix MW algorithm (in the gain form). In each round, we set  $\mathbf{X}$  to be the current density matrix  $\mathbf{P}^{(t)}$ , and check if it is a good solution. If not, and there is a constraint  $j$  such that  $\mathbf{A}_j \bullet \mathbf{X} < 0$ , then we set  $\mathbf{M}^{(t)} = \frac{1}{\rho} \mathbf{A}_j$ . Note that in case all the matrices involved in the SDP are diagonal, then this algorithm reduces to the Winnow algorithm of Section 3.1<sup>2</sup>

---

<sup>2</sup>A minor difference being that we get the version of Winnow based on the Hedge algorithm, rather than the MW algorithm.

The analysis is exactly the same as in Section 3.1. This analysis (using the gain form of Theorem 20 corresponding to Theorem 6), implies that in at most  $\left\lceil \frac{4\rho^2 \ln n}{\varepsilon^2} \right\rceil$  iterations, we find a good solution.

This technique was used in [AK07] to obtain significantly faster algorithms than previously known for approximating several combinatorial optimization problems such as SPARSEST CUT, BALANCED SEPARATOR (both in directed and undirected graphs), MIN UNCUT, and MIN 2CNF DELETION. This also gives the first near-linear time algorithm for solving the MAX CUT SDP of [GW95] to any constant approximation factor.

## 5.2 Applications in Quantum Computing

There are several other applications of the Matrix MW algorithm in quantum computing, since the basic object in both settings is the density matrix. The celebrated result of Jain et al. [JJUW10] showing that **QIP** = **PSPACE** relied on the Matrix MW algorithm. Here **QIP** is the set of all languages which have quantum interactive proofs. This result is proven as follows. Suppose we have a language for which a quantum interactive proof exists. The main idea used there is that given a string whose membership in the language needs to be decided, the maximum acceptance probability for the verifier in the quantum interactive proof can be computed via an SDP (in an exponentially large dimension). This SDP can be generated by a **PSPACE** algorithm, and can be (approximately) solved in polynomial parallel time using the Matrix MW algorithm since the number of iterations is logarithmic in the dimension and matrix operations such as exponentiation are parallelizable. Overall, this gives a **PSPACE** algorithm for approximately computing the maximum acceptance probability, which can be used to decide membership for the string in the language.

## References

- [AAA<sup>+</sup>09] N. Alon, B. Awerbuch, Y. Azar, N. Buchbinder, and J. Naor. The online set cover problem. *SIAM Journal on Computing*, 39(2):361–370, 2009.
- [AHK04] S. Arora, E. Hazan, and S. Kale.  $O(\sqrt{\log n})$  approximation to sparsest cut in  $\tilde{O}(n^2)$  time. In *FOCS*, pages 238–247, 2004.
- [AHK05] S. Arora, E. Hazan, and S. Kale. Fast algorithms for approximate semidefinite programming using the multiplicative weights update method. In *FOCS*, pages 339–348, 2005.
- [AK07] S. Arora and S. Kale. A combinatorial, primal-dual approach to semidefinite programs. In *STOC*, 2007.
- [ALN05] S. Arora, J. Lee, and A. Naor. Euclidean distortion and the sparsest cut. In *STOC*, 2005.
- [ARV04] S. Arora, S. Rao, and U. Vazirani. Expander flows, geometric embeddings, and graph partitioning. In *STOC*, pages 222–231, 2004.
- [AS92] N. Alon and J. Spencer. *Probabilistic Method*. John Wiley, 1992.

- [BG94] H. Bronnimann and M.T. Goodrich. Almost optimal set covers in finite vc-dimension. In *SoCG*, pages 293–302, 1994.
- [BHK09] B. Barak, M. Hardt, and S. Kale. The uniform hardcore lemma via approximate bregman projections. In *SODA*, pages 1193–1200, 2009.
- [BK96] A. A. Benczúr and D. R. Karger. Approximating  $s - t$  minimum cuts in  $\tilde{O}(n^2)$  time. In *STOC*, pages 47–55, 1996.
- [Blu98] A. Blum. On-line algorithms in machine learning. In A. Fiat and G. Woeginger, editors, *Online Algorithms: The State of the Art*, pages 306–325. LNCS 1442, 1998.
- [BN09a] N. Buchbinder and J. Naor. Online primal-dual algorithms for covering and packing. *Mathematics of Operations Research*, 34(2):270–286, 2009.
- [BN09b] N. Buchbinder and J. Naor. The Design of Competitive Online Algorithms via a Primal-Dual Approach. *Foundations and Trends® in Theoretical Computer Science*, 3(2–3):93–263, 2009.
- [Bro51] G. W. Brown. *Analysis of Production and Allocation*, T. C. Koopmans, ed., chapter Iterative solution of games by fictitious play, pages 374–376. Wiley, 1951.
- [BvN50] G.W. Brown and J. von Neumann. Solutions of games by differential equations. *Annals of Mathematical Studies*, 24:73–79, 1950.
- [CBL06] N. Cesa-Bianchi and G. Lugosi. *Prediction, Learning, and Games*. Cambridge University Press, 2006.
- [CGR05] S. Chawla, A. Gupta, and H. Racke. Embeddings of negative-type metrics and an improved approximation to generalized sparsest cut. In *SODA*, 2005.
- [Cha00] B. Chazelle. *The Discrepancy Method — Randomness and Complexity*. Cambridge University Press, Cambridge, 2000.
- [Cla88] K. Clarkson. A las vegas algorithm for linear programming when the dimension is small. In *FOCS*, pages 452–456, 1988.
- [Cov96] T. M. Cover. Universal data compression and portfolio selection. In *FOCS*, pages 534–538, 1996.
- [CW89] B. Chazelle and E. Welzl. Quasi-optimal range searching in spaces of finite VC-dimension. *Discrete & Computational Geometry*, 4, 1989.
- [Fei98] U. Feige. A threshold of  $\ln n$  for approximating set cover. *Journal of the ACM*, 45(4):634–652, July 1998.
- [Fle00] L. K. Fleischer. Approximating fractional multicommodity flow independent of the number of commodities. *SIAM J. Discret. Math.*, 13(4):505–520, 2000.

- [FS97] Y. Freund and R. E. Schapire. A decision-theoretic generalization of on-line learning and an application to boosting. *Journal of Computer and System Sciences*, 55(1):119–139, August 1997.
- [FS99] Y. Freund and R. E. Schapire. Adaptive game playing using multiplicative weights. *Games and Economic Behavior*, 29:79–103, 1999.
- [FV99] D.P. Foster and R. Vohra. Regret in the on-line decision problem. *Games and Economic Behaviour*, 29:7–35, 1999.
- [GK95] M. Grigoriadis and L. Khachiyan. A sublinear-time randomized approximation algorithm for matrix games. In *Operations Research Letters*, volume 18, pages 53–58, 1995.
- [GK98] N. Garg and J. Könemann. Faster and simpler algorithms for multicommodity flow and other fractional packing problems. In *FOCS*, pages 300–309, 1998.
- [GK04] N. Garg and R. Khandekar. Fractional Covering with Upper Bounds on the Variables: Solving LPs with Negative Entries. In *ESA*, pages 371–382, 2004.
- [Gol65] S. Golden. Lower Bounds for the Helmholtz Function. *Physical Review*, 137:1127–1128, February 1965.
- [GW95] M. X. Goemans and D. P. Williamson. Improved approximation algorithms for maximum cut and satisfiability problems using semidefinite programming. *J. ACM*, 42(6):1115–1145, 1995.
- [HAK07] E. Hazan, A. Agarwal, and S. Kale. Logarithmic regret algorithms for online convex optimization. In *Machine Learning*, volume 69(2-3), pages 169–192, 2007.
- [Han57] J. Hannan. Approximation to bayes risk in repeated plays. In M. Dresher, A. Tucker, and P. Wolfe, editors, *Contributaions to the Theory of Games*, volume 3, pages 97–139. Princeton University Press, 1957.
- [Haz06] E. Hazan. *Efficient Algorithms for Online Convex Optimization and their Applications*. PhD thesis, Princeton University, 2006. Technical Report TR-766-06.
- [Haz08] E. Hazan. Sparse approximate solutions to semidefinite programs. In *LATIN’08*, pages 306–316, 2008.
- [HSSW96] D. P. Helmbold, R. E. Schapire, Y. Singer, and M. K. Warmuth. On-line portfolio selection using multiplicative updates. In *ICML*, pages 243–251, 1996.
- [HW98] M. Herbster and M. K. Warmuth. Tracking the best expert. *Machine Learning*, 32(2):151–178, 1998.
- [Imp95] R. Impagliazzo. Hard-core distributions for somewhat hard problems. In *FOCS*, pages 538–545, 1995.
- [JJUW10] R. Jain, Z. Ji, S. Upadhyay, and J. Watrous. QIP = PSPACE. In *STOC*, pages 573–582, 2010.



- [Kal06] S. Kale. *Efficient Algorithms Using The Multiplicative Weights Update Method*. PhD thesis, Princeton University, 2006. Technical Report TR-804-07.
- [Kal07] S. Kale. Boosting and hard-core set constructions: a simplified approach. *ECCC Report TR07-131*, 2007.
- [Kal09] A. Kalai. *Personal communications with M. Warmuth*, 2009.
- [Kha04] R. Khandekar. *Lagrangian Relaxation based Algorithms for Convex Programming Problems*. PhD thesis, Indian Institute of Technology, Delhi, 2004. Available at <http://www.cse.iitd.ernet.in/~rohitk>.
- [KL96] P. Klein and H.-I. Lu. Efficient approximation algorithms for semidefinite programs arising from MAX CUT and COLORING. In *STOC*, pages 338–347, 1996.
- [KS03] A. R. Klivans and R. A. Servedio. Boosting and hard-core set construction. *Machine Learning*, 51(3):217–238, 2003.
- [KV03] A. Kalai and S. Vempala. Efficient algorithms for online decision problems. In *COLT*, pages 26–40, 2003.
- [Lit87] N. Littlestone. Learning quickly when irrelevant attributes abound: A new linear-threshold algorithm. *Machine Learning*, 2(4):285–318, 1987.
- [Lit89] N. Littlestone. *Mistake bounds and logarithmic linear-threshold learning algorithms*. PhD thesis, University of California at Santa Cruz, Santa Cruz, CA, USA, 1989.
- [LN93] M. Luby and N. Nisan. A parallel approximation algorithm for positive linear programming. In *STOC*, pages 448–457, 1993.
- [LW94] N. Littlestone and M. K. Warmuth. The weighted majority algorithm. *Information and Computation*, 108(2):212–261, 1994.
- [MP69] M. Minsky and S. Papert. *Perceptrons*. MIT Press, 1969.
- [MV08] J. Matousek and J. Vondrak. The probabilistic method. In *Lecture notes at Charles University*, 2008.
- [PST91] S. A. Plotkin, D. B. Shmoys, and É. Tardos. Fast approximation algorithm for fractional packing and covering problems. In *FOCS*, pages 495–504, 1991.
- [Rag86] P. Raghavan. Probabilistic construction of deterministic algorithms: Approximating packing integer programs. In *FOCS*, pages 10–18, 1986.
- [Rob51] J. Robinson. An iterative method of solving a game. *Ann. Math.*, 54:296–301, 1951.
- [RT87] P. Raghavan and C. D. Thompson. Randomized rounding: a technique for provably good algorithms and algorithmic proofs. *Combinatorica*, 7(4):365–374, 1987.

- [Sch90] R. E. Schapire. The strength of weak learnability. *Machine Learning*, 5:197–227, 1990.
- [Tho65] C. J. Thompson. Inequality with applications in statistical mechanics. *Journal of Mathematical Physics*, 6(11):1812–1823, 1965.
- [TRW05] K. Tsuda, G. Rtsch, and M. K. Warmuth. Matrix exponentiated gradient updates for on-line learning and bregman projection. *Journal of Machine Learning Research*, 6:995–1018, 2005.
- [WK06] M. Warmuth and D. Kuzmin. Online variance minimization. In *COLT*, 2006.
- [Yao82] A. C. Yao. Theory and applications of trapdoor functions (extended abstract). In *FOCS*, pages 80–91, 1982.
- [You95] N. E. Young. Randomized rounding without solving the linear program. In *SODA*, pages 170–178, 1995.
- [Zin03] M. Zinkevich. Online convex programming and generalized infinitesimal gradient ascent. In *ICML*, pages 928–936, 2003.