

Lecture 20: Matrix multiplicative update

The multiplicative weight update method can be generalized to a matrix setting, and this can be used to solving SDP.

Matrix setting

There is a nice generalization of the multiplicative weights update method to a matrix setting, and we will see that it can be used to design faster primal-dual algorithms for solving SDPs.

The setting is as follows:

- There is a decision corresponding to each unit vector in  $\mathbb{R}^n$ .
- Each round, a <sup>symmetric</sup> cost matrix  $M \in \mathbb{R}^{n \times n}$  is revealed, and the cost of each unit vector  $v$  is  $v^T M v$ .

We assume that  $v^T M v \in [0, 1]$  for all  $v$ ; equivalently we assume the eigenvalues of  $M$  are in  $[0, 1]$ .

As in the original setting, this is a bound on how big a mistake we can make, and is crucial to the convergence proof. We can write  $0 \preceq M \preceq I$ .

- In each round, we can choose a probability distribution over the set of unit vectors  $p_v$  with  $\sum_v p_v = 1$ .

Then, the expected cost of this probability distribution is  $\sum_v p_v v^T M v = \sum_v M^{(t)} (p_v v v^T) = M^{(t)} \cdot (\sum_v p_v v v^T)$ .

Let  $P^{(t)} = \sum_v p_v v v^T$ . Since  $P^{(t)}$  is a convex combination of PSD matrix of trace 1,  $P^{(t)}$  is also a PSD matrix of trace 1, and such a matrix is called a density matrix. Conversely, given a density

matrix  $P$ , we can write its eigendecomposition as  $P = \sum_{i=1}^n \lambda_i v_i v_i^T$  with  $\|v_i\|=1$  and  $\sum_{i=1}^n \lambda_i = \text{tr}(P) = 1$ ,

and thus it corresponds to a probability distribution over unit vectors. Therefore, we will specify a probability distribution compactly as a density matrix.

We would like to generalize the multiplicative weight update theorem to the matrix setting: the total

expected cost of the probability distributions is not much more than the total cost of a unit vector.

The main question is what would be the multiplicative update step?

In the previous setting,  $w_i^{(t+1)} := \prod_{k=1}^t (1 - \eta m_i^{(k)}) \approx \prod_{k=1}^t (1 - \eta)^{m_i^{(k)}} = (1 - \eta)^{\sum_{k=1}^t m_i^{(k)}} = e^{-\eta \sum_{k=1}^t m_i^{(k)}}$ , where  $\eta' = -\ln(1 - \eta)$ .

It turns out that there is an appropriate exponential function for matrix.

Given a matrix  $A$ , we define  $e^A := \sum_{i \geq 0} \frac{A^i}{i!}$  as the matrix exponential of  $A$ . Some properties of  $e^A$ :

- Regardless of  $A$ ,  $\exp(A)$  is positive semidefinite, because  $\exp(A) = \exp(\frac{1}{2}A)\exp(\frac{1}{2}A)$  (exercise).
- Suppose  $A$  is real symmetric. Then  $A = VDV^T$  as the eigendecomposition. Then  $A^i = VD^iV^T$  since  $VV^T = I$ .  
Then  $e^A = Ve^DV^T = \sum_{i=1}^n e^{\lambda_i} v_i v_i^T$  where  $\lambda_i$  is the  $i$ -th eigenvalue and  $v_i$  is the corresponding eigenvector.  
So, suppose we have an inequality that holds for all  $\lambda_i$ , then it also holds for the matrix. For example, if we know that  $(1-\eta)^x \leq 1-\eta x$  for all  $x \in [0,1]$ , then we can conclude that  $(1-\eta)^A \leq (I-\eta A)$  for  $0 \leq A \leq I$  as all eigenvalues of  $A$  are in  $[0,1]$ . To see it, first rewrite  $(1-\eta)^x$  as  $e^{-\eta'x}$  where  $\eta' = -\ln(1-\eta)$  and  $(1-\eta)^A = e^{-\eta'A}$ . Then  $I-\eta A - e^{-\eta'A} = VV^T - \eta VDV^T - Ve^{-\eta'D}V^T = V(I-\eta D - e^{-\eta'D})V^T$ , and this is PSD iff  $1-\eta x - e^{-\eta'x} \geq 0$  holds for all eigenvalues of  $A$  as  $I-\eta D - e^{-\eta'D}$  is a diagonal matrix. So, the matrix exponential of a symmetric matrix behaves like the exponential of a number.

After these discussions, it is then natural to update the weight as  $w^{(t+1)} = (1-\eta)^{\sum_{\ell=1}^t M^{(\ell)}} = \exp(-\eta' \sum_{\ell=1}^t M^{(\ell)})$ .

The formal algorithm is as follows.

Fix an  $\eta < \frac{1}{2}$ . Let  $\eta' = -\ln(1-\eta)$ .

For  $t=1, 2, \dots, T$  do

1) Compute  $w^{(t)} = (1-\eta)^{\sum_{\ell=1}^{t-1} M^{(\ell)}} = \exp(-\eta' \sum_{\ell=1}^{t-1} M^{(\ell)})$ .

2) Set  $P^{(t)} = \frac{w^{(t)}}{\text{Tr}(w^{(t)})}$  and observe the event  $M^{(t)}$ .

The scaling is to make sure that  $P^{(t)}$  is a density matrix.

Theorem The matrix multiplicative weights update algorithm will produce density matrices such that

$$\sum_{t=1}^T M^{(t)} \cdot P^{(t)} \leq (1+\eta) V^T \left( \sum_{t=1}^T M^{(t)} \right) V + \frac{\ln n}{\eta} \text{ for all unit vector } v \in \mathbb{R}^n.$$

Equivalently,  $\sum_{t=1}^T M^{(t)} \cdot P^{(t)} \leq (1+\eta) \lambda_{\min} \left( \sum_{t=1}^T M^{(t)} \right) + \frac{\ln n}{\eta}$ . We assume  $0 \leq M^{(t)} \leq I$  for all  $t$ .

Proof In the previous setting, we use  $\phi^{(t)} = \sum_i w_i^{(t)}$  as a potential function.

In the matrix setting, we use  $\phi^{(t)} = \text{Tr}(W^{(t)})$  as the potential function.

$$\begin{aligned} \text{Tr}(W^{(t+1)}) &= \text{Tr} \left( \exp \left( -\eta' \sum_{\ell=1}^t M^{(\ell)} \right) \right) \\ &\leq \text{Tr} \left( \exp \left( -\eta' \sum_{\ell=1}^{t-1} M^{(\ell)} \right) \exp \left( -\eta' M^{(t)} \right) \right) \quad (\text{Golden-Thompson inequality: } \text{Tr}(\exp(A+B)) \leq \text{Tr}(\exp(A)\exp(B))) \\ &= W^{(t)} \cdot \exp \left( -\eta' M^{(t)} \right) \quad (\text{Tr}(AB) = A \cdot B \text{ for symmetric } A) \\ &\leq W^{(t)} \cdot (I - \eta M^{(t)}) \quad ((1-\eta)^A \leq (I-\eta A) \text{ if } 0 \leq A \leq I \text{ as explained above}) \end{aligned}$$

$$= W \cdot \exp(-\eta M)$$

$$\leq W^{(t)} \cdot (I - \eta M^{(t)})$$

$$= \text{Tr}(W^{(t)})(I - \eta M^{(t)} \cdot P^{(t)})$$

$$\leq \text{Tr}(W^{(t)}) e^{-\eta M^{(t)} \cdot P^{(t)}}$$

$$( \text{Tr}(AB) = A \cdot B \text{ for symmetric } A )$$

$$( (I - \eta)^A \leq (I - \eta A) \text{ if } 0 \leq A \leq I \text{ as explained above} )$$

$$( \text{use } W^{(t)} = P^{(t)} \cdot \text{Tr}(W^{(t)}) )$$

By induction, since  $\text{Tr}(W^{(1)}) = \text{Tr}(I) = n$ , this implies that  $\text{Tr}(W^{(t+1)}) \leq n \exp(-\eta \sum_{t=1}^t M^{(t)} \cdot P^{(t)})$ .

On the other hand, we have  $\text{Tr}(W^{(t+1)}) = \text{Tr}(\exp(-\eta' \sum_{t=1}^t M^{(t)})) \geq \exp(-\eta' \lambda_{\min}(\sum_{t=1}^t M^{(t)}))$ ,

$$\text{because } \text{Tr}(e^A) = \text{Tr}(V e^D V^T) = \text{Tr}(e^D) = \sum_{i=1}^n e^{\lambda_i(A)} \geq e^{\lambda_{\min}(A)}.$$

Therefore, we have  $\exp(-\eta' \lambda_{\min}(\sum_{t=1}^t M^{(t)})) \leq n \exp(-\eta \sum_{t=1}^t M^{(t)} \cdot P^{(t)})$

$$\Rightarrow \eta \sum_{t=1}^t M^{(t)} \cdot P^{(t)} \leq \ln\left(\frac{1}{1-\eta}\right) \lambda_{\min}\left(\sum_{t=1}^t M^{(t)}\right) + \ln n$$

$$\Rightarrow \eta \sum_{t=1}^t M^{(t)} \cdot P^{(t)} \leq (\eta + \eta^2) \lambda_{\min}\left(\sum_{t=1}^t M^{(t)}\right) + \ln n \quad (\text{using } \ln\left(\frac{1}{1-\eta}\right) \leq \eta + \eta^2 \text{ for } \eta < \frac{1}{2})$$

$$\Rightarrow \sum_{t=1}^t M^{(t)} \cdot P^{(t)} \leq (1+\eta) \lambda_{\min}\left(\sum_{t=1}^t M^{(t)}\right) + \frac{\ln n}{\eta}. \quad \square$$

## Primal-dual SDP

The matrix multiplicative update method can be used nicely to design a fast algorithm for solving SDP.

Consider the following primal dual pair.

$$\begin{array}{lll}
 \max C \cdot X & & \min b^T y \\
 \text{(primal)} \quad A_i \cdot X \leq b_i \text{ for } 1 \leq i \leq m & \Leftrightarrow & \text{(dual)} \quad \sum_{i=1}^m y_i A_i - C \leq 0 \\
 \text{Tr}(X) \leq 1 & & (\text{where } A_0 = I, b_0 = 1) \\
 X \succeq 0 & & y_i \geq 0 \text{ for } 0 \leq i \leq m \\
 & & X \succeq 0
 \end{array}$$

Focus on solving the dual program. Reduce it to a decision problem.

$$\begin{array}{ll}
 \sum_{i=1}^m y_i A_i - C \leq 0 & \left( \sum_{i=1}^m y_i A_i - C \right) \cdot X \leq 0 \text{ for all } X \succeq 0 \\
 b^T y \leq \alpha & \Leftrightarrow \quad b^T y \leq \alpha \\
 y \geq 0 & y \geq 0.
 \end{array}$$

As in the previous approach, we identify some constraints as hard and use multiplicative update to combine them, while some constraints as easy and directly implement an oracle to solve them.

In this case, the easy constraints are  $\{ b^T y \leq \alpha, y \geq 0 \}$  and the hard constraints as the PSD constraint.

To satisfy the hard constraint, we need to find  $y$  that satisfies all possible density matrices  $X$ .

Equivalently, we need to find  $y$  that satisfies  $v^T (\sum_{i=0}^m y_i A_i - C) v \geq 0$  for all unit vectors  $v \in \mathbb{R}^n$ .

There are infinitely many unit vectors to consider. It seems too complicated.

But this is exactly what the matrix multiplicative update algorithm does!

Define the cost matrix  $M^{(t)} = \sum_{i=0}^m y_i A_i - C$ . (It does not satisfy  $0 \preceq M^{(t)} \preceq I$ . We will fix it.)

Then the cost of a unit vector  $v$  is  $(\sum_{i=0}^m y_i A_i - C) \cdot (v v^T)$ . It is smaller if its constraint is violated more.

The algorithm will generate density matrices  $P^{(1)}, P^{(2)}, \dots, P^{(T)}$  such that if we can satisfy them, then we almost satisfy all unit vectors.

So, in each round, we ask the multiplicative algorithm to provide a density matrix  $P^{(t)}$  to us.

Then, we find  $y^{(t)}$  such that  $(\sum_{i=0}^m y_i^{(t)} A_i - C) \cdot P^{(t)} \geq 0$ , and return  $M^{(t)} = \sum_{i=0}^m y_i^{(t)} A_i - C$  as the cost matrix to the multiplicative update algorithm as a feedback (to tell them what unit vectors are satisfied and what unit vectors are violated).

Then, the multiplicative algorithm will cleverly update its density matrix (to put more weight on the constraints that are violated more), and send us another matrix  $P^{(t+1)}$ , and the procedure is repeated.

Now, if we can find  $y^{(1)}, y^{(2)}, \dots, y^{(T)}$  that satisfy all the density matrices  $P^{(1)}, P^{(2)}, \dots, P^{(T)}$ , by the matrix multiplicative update theorem above, they will almost satisfy all the unit vectors, and thus  $\bar{y} = \frac{1}{T} \sum_{t=1}^T y^{(t)}$  is an almost feasible dual solution.

This is the idea of this approach, very similar to the previous setting, just that different constraints are considered hard.

Now, we formally describe the algorithm. First, we specify the oracle.

ORACLE: Given  $X$ , return  $y \in \mathbb{R}^{m+1}$  such that  $(\sum_{i=0}^m y_i A_i - C) \cdot X \geq 0$  and  $b^T y \leq \alpha$  and  $y \geq 0$ .

Note that, given  $X$ , the first constraint is just a linear constraint  $\sum_{i=0}^m y_i (A_i \cdot X) \geq C \cdot X$

So, the oracle just needs to solve a linear program with two non-trivial constraints.

As it turned out, for many combinatorial problems, the oracle can be implemented by combinatorial techniques.

Clearly, if there is a dual feasible solution with value  $\leq \alpha$ , then this LP is solvable.

Indeed, if this LP is not solvable, then one can prove by LP duality that the solution  $X$  is primal feasible with value at least  $\alpha$  (exercise).

Width: We assume the oracle always returns  $y$  with  $\|\sum_{i=0}^m y_i A_i - C\| \leq w$ . Again, this is to bound the maximum violation of a unit vector, and this parameter is crucial in bounding the convergence rate.

### ALGORITHM

Set  $P^{(1)} = I/n$ .

For  $t=1, 2, \dots, T$  do

- 1) Run the oracle with  $P^{(t)}$ . If the oracle fails, return  $P^{(t)}$  as a primal feasible solution with value  $\geq \alpha$ .  
Otherwise, let  $y^{(t)}$  be the solution that the oracle returns.
- 2) Set  $M^{(t)} = \frac{1}{2w} \left( \sum_{i=0}^m y_i^{(t)} A_i - C + wI \right)$ . This guarantees that  $0 \leq M^{(t)} \leq I$  by our assumption on the width.
- 3) Use  $M^{(t)}$  as the outcome to the matrix multiplicative algorithm, and use its output as  $P^{(t+1)}$ .

All the ideas are discussed. The following details are just formal execution of the ideas.

Analysis: By the matrix multiplicative update theorem, we have

$$\sum_{t=1}^T M^{(t)} \cdot P^{(t)} \leq (1+\eta) \lambda_{\min} \left( \sum_{t=1}^T M^{(t)} \right) + \frac{\ln n}{\eta}, \text{ where } \eta \text{ is the parameter used in MMU algo.}$$

Assuming the oracle never fails, then  $M^{(t)} \cdot P^{(t)} = \frac{1}{2w} \left( \sum_{i=0}^m y_i^{(t)} A_i - C \right) \cdot P^{(t)} + \frac{1}{2} I \cdot P^{(t)} \geq \frac{1}{2}$ ,

since the first term  $\geq 0$  if the oracle succeeds and the second term  $= \frac{1}{2}$  as  $P^{(t)}$  is a density matrix.

Therefore, the LHS of the inequality is at least  $\frac{T}{2}$ , and thus we have

$$\frac{T}{2} - \frac{\ln n}{\eta} \leq (1+\eta) \lambda_{\min} \left( \sum_{t=1}^T \frac{1}{2w} \left( \sum_{i=0}^m y_i^{(t)} A_i - C + wI \right) \right)$$

$$\Rightarrow \lambda_{\min} \left( \sum_{t=1}^T \left( \sum_{i=0}^m y_i^{(t)} A_i - C + wI \right) \right) \geq \frac{2w}{(1+\eta)} \left( \frac{T}{2} - \frac{\ln n}{\eta} \right)$$

$$\Rightarrow \lambda_{\min} \left( \sum_{i=0}^m \left( \frac{\sum_{t=1}^T y_i^{(t)}}{T} \right) A_i - C \right) \geq \frac{-\eta w}{1+\eta} - \frac{2w \ln n}{(1+\eta)\eta T} \geq -\eta w - \frac{2w \ln n}{\eta T}$$

For an error parameter  $\delta > 0$ , set  $\eta = \frac{\delta \alpha}{2w}$  and then set  $T = \frac{8w^2 \ln n}{\alpha^2 \delta^2}$ , also let  $\bar{y} = \frac{\sum_{t=1}^T y_i^{(t)}}{T}$ ,

$$\text{we have } \lambda_{\min} \left( \sum_{i=0}^m \bar{y}_i A_i - C \right) \geq -\delta \alpha.$$

Finally, recall that  $A_0 = I$  and  $b_0 = 1$  by the trace constraint.

Set  $y^* = \bar{y} + \alpha \delta e_0$ . Then, we have  $\sum_{i=0}^m y_i^* A_i - C \preceq 0$ , and is a dual feasible solution.

The value of  $y^*$  is  $b^T y^* = b^T \bar{y} + b^T \alpha \delta e_0 = b^T \left( \frac{\sum_{t=1}^T y^{(t)}}{T} \right) + \delta \alpha \leq \alpha + \delta \alpha = (1 + \delta) \alpha$ ,

while the inequality holds because  $b^T y^{(t)} \leq \alpha$  by the definition of the oracle.

To conclude, we have the following theorem.

Theorem In  $O\left(\frac{W^2 \ln n}{\alpha^2 \delta^2}\right)$  iterations, the algorithm returns a dual feasible solution with value at most  $(1 + \delta) \alpha$ , or returns a primal feasible solution with value at least  $\alpha$ .

Remark: If  $\text{tr}(X) \leq R$  in the primal program, by the same scaling argument we've seen before, we can solve it with an  $O(R^2)$  blowup in the number of iterations.

The advantage of this method over the previous SDP solver is that there are no violations of the constraints, the violations are only in the objective function. This avoids further modification of the solution which may significantly change the objective value.

This approach is used successfully to get an  $\tilde{O}(m)$  time algorithm to solve the max-cut SDP for regular graphs. We won't discuss this result.

The matrix multiplicative update algorithm has applications in machine learning problems, and is recently used to prove an important result in quantum computing (QIP=PSPACE).

---

## References

The material is taken from the survey paper "The multiplicative weights update method: a meta-algorithm and applications" by Arora, Hazan and Kale, and the paper "A combinatorial, primal-dual approach to semidefinite programs" by Arora and Kale.

For the proof of the Golden-Thompson inequality, you are referred to the blog post of Terence Tao or the paper "Golden-Thompson from Davis" by Rivin.