#### Lecture 13: Multiplicative Weights Update

Rafael Oliveira

University of Waterloo Cheriton School of Computer Science

rafael.oliveira.teaching@gmail.com

June 22, 2021

イロン 不得 とうほう イロン 二日

1/83

#### Overview

- Multiplicative Weights Update
- Solving Linear Programs
- Conclusion
- Acknowledgements

- Setup: investing your co-op money on stock markets (or gambling).
- Objective: to get rich, but we don't know much about stock markets
- Have access to *n* experts (news programs, newspapers, social media)

- Setup: investing your co-op money on stock markets (or gambling).
- Objective: to get rich, but we don't know much about stock markets
- Have access to *n* experts (news programs, newspapers, social media)
  - Each morning, before market opens, experts predict whether the price of a stock will go up or down
  - By the time market closes, we can check the outcome, who was right or wrong.

- Setup: investing your co-op money on stock markets (or gambling).
- Objective: to get rich, but we don't know much about stock markets
- Have access to *n* experts (news programs, newspapers, social media)
  - Each morning, before market opens, experts predict whether the price of a stock will go up or down
  - By the time market closes, we can check the outcome, who was right or wrong.
    - Experts who were right earn one dollar
    - Experts who were wrong lose one dollar

- Setup: investing your co-op money on stock markets (or gambling).
- Objective: to get rich, but we don't know much about stock markets
- Have access to *n* experts (news programs, newspapers, social media)
  - Each morning, before market opens, experts predict whether the price of a stock will go up or down
  - By the time market closes, we can check the outcome, who was right or wrong.
    - Experts who were right earn one dollar
    - Experts who were wrong lose one dollar
- Some expert did really well, and if we followed their advice we would have made a lot of money...

- Setup: investing your co-op money on stock markets (or gambling).
- Objective: to get rich, but we don't know much about stock markets
- Have access to *n* experts (news programs, newspapers, social media)
  - Each morning, before market opens, experts predict whether the price of a stock will go up or down
  - By the time market closes, we can check the outcome, who was right or wrong.
    - Experts who were right earn one dollar
    - Experts who were wrong lose one dollar
- Some expert did really well, and if we followed their advice we would have made a lot of money...
- Hindsight is 20/20 though. To make money, need to make a decision on what & how to trade every day.

- Setup: investing your co-op money on stock markets (or gambling).
- Objective: to get rich, but we don't know much about stock markets
- Have access to *n* experts (news programs, newspapers, social media)
  - Each morning, before market opens, experts predict whether the price of a stock will go up or down
  - By the time market closes, we can check the outcome, who was right or wrong.
    - Experts who were right earn one dollar
    - Experts who were wrong lose one dollar
- Some expert did really well, and if we followed their advice we would have made a lot of money...
- Hindsight is 20/20 though. To make money, need to make a decision on what & how to trade every day.
- Can we hope to do as well as the best expert in hindsight?

• Online Learning

#### • Online Learning

Experts are weak classifiers, want to choose hypothesis based on these experts

Boosting (in learning theory)

- Online Learning
  - Experts are weak classifiers, want to choose hypothesis based on these experts

Boosting (in learning theory)

• Solving linear programs! (today)

- Online Learning
  - Experts are weak classifiers, want to choose hypothesis based on these experts

Boosting (in learning theory)

- Solving linear programs! (today)
- Game Theory
- many more

Say we are trading for T days.

Why not want our algorithm to make as much money as a function of T as we can?

Say we are trading for T days.

Why not want our algorithm to make as much money as a function of T as we can?

• With the benchmark above, guessing correctly (in expectation) *T*/2 times is *trivial* (pick your trades at random)

Say we are trading for T days.

Why not want our algorithm to make as much money as a function of T as we can?

- With the benchmark above, guessing correctly (in expectation) *T*/2 times is *trivial* (pick your trades at random)
- It turns out, T/2 correct guesses (in expectation) is also optimal

Say we are trading for T days.

Why not want our algorithm to make as much money as a function of T as we can?

- With the benchmark above, guessing correctly (in expectation) *T*/2 times is *trivial* (pick your trades at random)
- It turns out, T/2 correct guesses (in expectation) is also optimal
  - Worst-case analysis.

Say we knew that there was *one expert* which will be *right every time*. What should we do?

• At each trading day, take *majority vote* of the opinions of the experts.

- At each trading day, take *majority vote* of the opinions of the experts.
- If we made the right trade, do nothing.

- At each trading day, take *majority vote* of the opinions of the experts.
- If we made the right trade, do nothing.
- If our trade was bad, at the end of the trading day *discard* all experts that *made a mistake that day*.

- At each trading day, take *majority vote* of the opinions of the experts.
- If we made the right trade, do nothing.
- If our trade was bad, at the end of the trading day *discard* all experts that *made a mistake that day*.
- **1** Every time we made a bad trade, we discard half of the experts.
  - We took majority vote, so at least half the experts also made bad trades

- At each trading day, take *majority vote* of the opinions of the experts.
- If we made the right trade, do nothing.
- If our trade was bad, at the end of the trading day *discard* all experts that *made a mistake that day*.
- Every time we made a bad trade, we discard half of the experts.
  - We took majority vote, so at least half the experts also made bad trades
- After log n bad trades, only the expert who is always right will remain! From then on, we will always be right!

Say we knew that there was *one expert* which will be *right every time*. What should we do?

- At each trading day, take *majority vote* of the opinions of the experts.
- If we made the right trade, do nothing.
- If our trade was bad, at the end of the trading day *discard* all experts that *made a mistake that day*.

Every time we made a bad trade, we discard half of the experts.

- We took majority vote, so at least half the experts also made bad trades
- After log n bad trades, only the expert who is always right will remain! From then on, we will always be right!

Total money we made:  $\geq T - \log n$ 

Total money best expert made: T

・ロト ・ 同 ト ・ 三 ト ・ 三 ・ つへの

• In general do not have a single expert who is right all the time.

- In general do not have a single expert who is right all the time.
- Algorithm from previous slide *not robust* to having "almost perfect" experts (say experts that make one mistake)

- In general do not have a single expert who is right all the time.
- Algorithm from previous slide *not robust* to having "almost perfect" experts (say experts that make one mistake)
- There is a way of making previous idea robust:

Whenever an expert makes a mistake, "consider their opinions with less importance."

- In general do not have a single expert who is right all the time.
- Algorithm from previous slide *not robust* to having "almost perfect" experts (say experts that make one mistake)
- There is a way of making previous idea robust:
   Whenever an expert makes a mistake, "consider their opinions with

less importance."

- Let  $w_t : [n] \to \mathbb{R}_+$  be a function from each expert to the non-negative reals, and  $0 < \varepsilon < 1$ 
  - $w_t(i)$  is the *weight* of expert *i* at time *t*

- In general do not have a single expert who is right all the time.
- Algorithm from previous slide *not robust* to having "almost perfect" experts (say experts that make one mistake)
- There is a way of making previous idea robust:
   Whenever an expert makes a mistake, "consider their opinions with less importance."
- Let  $w_t : [n] \to \mathbb{R}_+$  be a function from each expert to the non-negative reals, and  $0 < \varepsilon < 1$ 
  - $w_t(i)$  is the *weight* of expert *i* at time *t*
- 2 In beginning every expert has weight  $w_1(i) = 1$

- In general do not have a single expert who is right all the time.
- Algorithm from previous slide *not robust* to having "almost perfect" experts (say experts that make one mistake)
- There is a way of making previous idea robust:
   Whenever an expert makes a mistake, "consider their opinions with

less importance."

- Let  $w_t : [n] \to \mathbb{R}_+$  be a function from each expert to the non-negative reals, and  $0 < \varepsilon < 1$ 
  - $w_t(i)$  is the *weight* of expert *i* at time *t*
- 2 In beginning every expert has weight  $w_1(i) = 1$
- If an expert makes a mistake at day t, make  $w_{t+1} \mathbf{P}(i) = w_t(i) \cdot (1 \varepsilon)$

- In general do not have a single expert who is right all the time.
- Algorithm from previous slide not robust to having "almost perfect" experts (say experts that make one mistake)
- There is a way of making previous idea robust:

Whenever an expert makes a mistake, "consider their opinions with less importance."

- Let  $w_t : [n] \to \mathbb{R}_+$  be a function from each expert to the non-negative reals, and  $0 < \varepsilon < 1$ 
  - $w_t(i)$  is the *weight* of expert *i* at time *t*
- 2 In beginning every expert has weight  $w_1(i) = 1$
- 3 If an expert makes a mistake at day t, make  $w_{t+1}$ ,  $w_t(i) = w_t(i) \cdot (1 \varepsilon)$
- Each trading day, choose to trade based on *weighted majority* of the decisions of the experts

# Multiplicative Weights Update Algorithm Algorithm:

Algorithm:

Setup: we have a binary decision to make (i.e.,  $\{-1, +1\}$ ) and we have access to *n* experts, indexed by the set [n].

Algorithm:

- Setup: we have a binary decision to make (i.e.,  $\{-1, +1\}$ ) and we have access to *n* experts, indexed by the set [n].
  - At each time step t, expert takes decision  $d_t(i) \in \{-1, +1\}$
  - Parameter  $0 < \varepsilon < 1/2$

Algorithm:

- Setup: we have a binary decision to make (i.e.,  $\{-1, +1\}$ ) and we have access to *n* experts, indexed by the set [n].
  - At each time step t, expert takes decision  $d_t(i) \in \{-1,+1\}$
  - Parameter  $0 < \varepsilon < 1/2$
- 2 Let  $w_t : [n] \to \mathbb{R}_+$  weight function
  - $w_t(i)$  is the *weight* of expert *i* at time *t*
  - In beginning every expert has weight  $w_1(i) = 1$

Algorithm:

- Setup: we have a binary decision to make (i.e.,  $\{-1, +1\}$ ) and we have access to *n* experts, indexed by the set [n].
  - At each time step t, expert takes decision  $d_t(i) \in \{-1,+1\}$
  - Parameter  $0 < \varepsilon < 1/2$
- 2 Let  $w_t : [n] \to \mathbb{R}_+$  weight function
  - $w_t(i)$  is the *weight* of expert *i* at time *t*
  - In beginning every expert has weight  $w_1(i) = 1$

#### 3 At each time step (i.e. for t = 1, ..., T):

Algorithm:

- Setup: we have a binary decision to make (i.e.,  $\{-1, +1\}$ ) and we have access to *n* experts, indexed by the set [n].
  - At each time step t, expert takes decision  $d_t(i) \in \{-1,+1\}$
  - Parameter  $0 < \varepsilon < 1/2$
- 2 Let  $w_t : [n] \to \mathbb{R}_+$  weight function
  - $w_t(i)$  is the *weight* of expert *i* at time *t*
  - In beginning every expert has weight  $w_1(i) = 1$
- At each time step (i.e. for t = 1, ..., T):
  - Make your decision based on weighted majority:

$$\begin{cases} +1, \text{ if } \sum_{i=1}^{n} w_t(i) \cdot d_t(i) \ge 0 \\ -1, \text{ otherwise} \end{cases}$$

$$\sum_{i=1}^{n} w_t(i) \cdot d_t(i) = \sum_{\substack{i \le n \\ j \le$$

## Multiplicative Weights Update Algorithm

Algorithm:

- Setup: we have a binary decision to make (i.e.,  $\{-1, +1\}$ ) and we have access to *n* experts, indexed by the set [n].
  - At each time step t, expert takes decision  $d_t(i) \in \{-1, +1\}$
  - Parameter  $0 < \varepsilon < 1/2$
- 2 Let  $w_t : [n] \to \mathbb{R}_+$  weight function
  - $w_t(i)$  is the *weight* of expert *i* at time *t*
  - In beginning every expert has weight  $w_1(i) = 1$
- At each time step (i.e. for t = 1, ..., T):
  - Make your decision based on weighted majority:

$$\begin{cases} +1, \text{ if } \sum_{i=1}^{n} w_t(i) \cdot d_t(i) \ge 0\\ -1, \text{ otherwise} \end{cases}$$

If an expert makes a mistake at time t, make

$$w_{t+1}(i) = w_t(i) \cdot (1 - \varepsilon)$$

#### Theorem (Multiplicative Weights Update)

Let  $M_t$  be the number of mistakes that our algorithm makes until time t, and let  $M_t(i)$  be the number of mistakes that expert i made until time t. Then, for any expert  $i \in [n]$ , we have:

$$M_t \leq 2 \cdot (1 + \varepsilon) M_t(i) + \frac{2 \log n}{\varepsilon}$$

#### Theorem (Multiplicative Weights Update)

Let  $M_t$  be the number of mistakes that our algorithm makes until time t, and let  $M_t(i)$  be the number of mistakes that expert i made until time t. Then, for any expert  $i \in [n]$ , we have:

$$M_t \leq 2 \cdot (1 + \varepsilon) M_t(i) + \frac{2 \log n}{\varepsilon}$$

Potential function

$$\Phi_t = \sum_{i=1}^n w_t(i)$$
polyntial
function at
time t

#### Theorem (Multiplicative Weights Update)

Let  $M_t$  be the number of mistakes that our algorithm makes until time t, and let  $M_t(i)$  be the number of mistakes that expert i made until time t. Then, for any expert  $i \in [n]$ , we have:

$$M_t \leq 2 \cdot (1+\varepsilon)M_t(i) + \frac{2\log n}{\varepsilon}$$

Potential function

$$\Phi_t = \sum_{i=1}^n w_t(i)$$

Intuition:

#### Theorem (Multiplicative Weights Update)

Let  $M_t$  be the number of mistakes that our algorithm makes until time t, and let  $M_t(i)$  be the number of mistakes that expert i made until time t. Then, for any expert  $i \in [n]$ , we have:

$$M_t \leq 2 \cdot (1+\varepsilon)M_t(i) + \frac{2\log n}{\varepsilon}$$

Potential function

$$\Phi_t = \sum_{i=1}^n w_t(i)$$

- Intuition:
  - If we make mistake,  $\Phi_{t+1}$  decreases by a multiplicative factor w.r.t.  $\Phi_t$

#### Theorem (Multiplicative Weights Update)

Let  $M_t$  be the number of mistakes that our algorithm makes until time t, and let  $M_t(i)$  be the number of mistakes that expert i made until time t. Then, for any expert  $i \in [n]$ , we have:

$$M_t \leq 2 \cdot (1+\varepsilon)M_t(i) + \frac{2\log n}{\varepsilon}$$

Potential function

$$\Phi_t = \sum_{i=1}^n w_t(i)$$

#### Intuition:

- If we make mistake,  $\Phi_{t+1}$  decreases by a multiplicative factor w.r.t.  $\Phi_t$
- $\Phi_t$  is monotonically decreasing (so if we get it right, potential does not increase either)

### Theorem (Multiplicative Weights Update)

Let  $M_t$  be the number of mistakes that our algorithm makes until time t, and let  $M_t(i)$  be the number of mistakes that expert i made until time t. Then, for any expert  $i \in [n]$ , we have:

$$M_t \leq 2 \cdot (1 + \varepsilon) M_t(i) + \frac{2 \log n}{\varepsilon}$$

#### Potential function

$$\Phi_t = \sum_{i=1}^n w_t(i)$$

#### Intuition:

- If we make mistake,  $\Phi_{t+1}$  decreases by a multiplicative factor w.r.t.  $\Phi_t$
- $\Phi_t$  is monotonically decreasing (so if we get it right, potential does not increase either)
- Initially  $\Phi_1 = n$
- $\Phi_t \ge 0$  for all t

#### Theorem (Multiplicative Weights Update)

Let  $M_t$  be the number of mistakes that our algorithm makes until time t, and let  $M_t(i)$  be the number of mistakes that expert i made until time t. Then, for any expert  $i \in [n]$ , we have:

$$M_t \leq 2 \cdot (1 + \varepsilon) M_t(i) + \frac{2 \log n}{\varepsilon}$$

#### Theorem (Multiplicative Weights Update)

Let  $M_t$  be the number of mistakes that our algorithm makes until time t, and let  $M_t(i)$  be the number of mistakes that expert i made until time t. Then, for any expert  $i \in [n]$ , we have:

$$M_t \leq 2 \cdot (1 + \varepsilon) M_t(i) + \frac{2 \log n}{\varepsilon}$$

• Potential function  $\Phi_t = \sum_{i=1}^n w_t(i)$ 

#### Theorem (Multiplicative Weights Update)

Let  $M_t$  be the number of mistakes that our algorithm makes until time t, and let  $M_t(i)$  be the number of mistakes that expert i made until time t. Then, for any expert  $i \in [n]$ , we have:

$$M_t \leq 2 \cdot (1 + \varepsilon) M_t(i) + \frac{2 \log n}{\varepsilon}$$

#### Theorem (Multiplicative Weights Update)

Let  $M_t$  be the number of mistakes that our algorithm makes until time t, and let  $M_t(i)$  be the number of mistakes that expert i made until time t. Then, for any expert  $i \in [n]$ , we have:

$$M_t \leq 2 \cdot (1 + \varepsilon) M_t(i) + \frac{2 \log n}{\varepsilon}$$

- Potential function  $\Phi_t = \sum_{i=1}^n w_t(i)$
- If we made a mistake, at least half the weight was on the wrong answer. Thus

$$\Phi_{t+1} = \sum_{i \text{ right}} w_t(i) + (1 - \varepsilon) \cdot \sum_{j \text{ wrong}} w_t(j) \le \left(1 - \frac{\varepsilon}{2}\right) \cdot \Phi_t$$

Thus,

$$\Phi_t \leq \Phi_1 \cdot \left(1 - \frac{\varepsilon}{2}\right)^{M_t} = n \cdot \left(1 - \frac{\varepsilon}{2}\right)^{M_t}$$

イロト 不得 トイヨト イヨト 二日

We have

$$\Phi_t \leq n \cdot \left(1 - \frac{\varepsilon}{2}\right)^{M_t}$$

We have

$$\Phi_t \leq n \cdot \left(1 - \frac{\varepsilon}{2}\right)^{M_t}$$

On the other hand, have:

$$\Phi_t = \sum_{j=1}^n w_t(j) > w_t(i) = (1 - \varepsilon)^{M_t(i)}$$
  
definition by algorithm

We have

$$\Phi_t \leq n \cdot \left(1 - \frac{\varepsilon}{2}\right)^{M_t}$$

On the other hand, have:

$$\Phi_t = \sum_{j=1}^n w_t(j) > w_t(i) = (1-\varepsilon)^{M_t(i)}$$

 $\bigcirc$  Putting (1) and (2) together

$$n \cdot \left(1 - \frac{\varepsilon}{2}\right)^{M_t} > (1 - \varepsilon)^{M_t(i)} \Rightarrow \log(1 - \varepsilon/2) \cdot M_t + \log n > M_t(i) \cdot \log(1 - \varepsilon)$$

$$\sum_{t=1}^{M_t} \sum_{t=1}^{M_t} \sum_{i=1}^{M_t} \sum_{j=1}^{M_t} \sum_{i=1}^{M_t} \sum_{j=1}^{M_t} \sum_{i=1}^{M_t} \sum_{j=1}^{M_t} \sum_{j=1}^{M_t} \sum_{i=1}^{M_t} \sum_{j=1}^{M_t} \sum_{i=1}^{M_t} \sum_{j=1}^{M_t} \sum_{j=1}^{M_t} \sum_{i=1}^{M_t} \sum_{j=1}^{M_t} \sum_{i=1}^{M_t} \sum_{j=1}^{M_t} \sum_{i=1}^{M_t} \sum_{j=1}^{M_t} \sum_{j=1}^{M_t} \sum_{i=1}^{M_t} \sum_{j=1}^{M_t} \sum_{i=1}^{M_t} \sum_{j=1}^{M_t} \sum_{j=1}^{M_t} \sum_{i=1}^{M_t} \sum_{j=1}^{M_t} \sum_{i=1}^{M_t} \sum_{j=1}^{M_t} \sum_{j=1}^{M_t} \sum_{j=1}^{M_t} \sum_{j=1}^{M_t} \sum_{i=1}^{M_t} \sum_{j=1}^{M_t} \sum_{j=1$$

$$\frac{1}{\log\left(\frac{1}{1-\epsilon_{1}}\right)}\left(\mathcal{M}_{t}(i)\cdot\log\left(\frac{1}{1-\epsilon_{1}}\right)+\log n\right)>\mathcal{M}_{t}$$

• Putting (1) and (2) together  $n \cdot \left(1 - \frac{\varepsilon}{2}\right)^{M_t} > (1 - \varepsilon)^{M_t(i)} \Rightarrow \log(1 - \varepsilon/2) \cdot M_t + \log n > M_t(i) \cdot \log(1 - \varepsilon)$ 

• Using inequality  $-x - x^2 < \log(1 - x) < -x$  for  $x \in (0, 1/2)$ , we get:

$$-\varepsilon/2 \cdot M_t + \log n > M_t(i) \cdot (-\varepsilon - \varepsilon^2)$$

The same algorithm and argument above, applied to the setting:

**Setup:** have access to *n* experts, indexed by the set [*n*].

- **§** Setup: have access to n experts, indexed by the set [n].
  - At each time step t, each expert will guess a value  $m_t(i) \in [-1, +1]$ .
  - Cost of  $i^{th}$  expert answer at time t is  $m_t(i)$
  - Parameter  $0 < \varepsilon < 1/2$

- **§** Setup: have access to n experts, indexed by the set [n].
  - At each time step t, each expert will guess a value  $m_t(i) \in [-1, +1]$ .
  - Cost of  $i^{th}$  expert answer at time t is  $m_t(i)$
  - Parameter  $0 < \varepsilon < 1/2$
- **2** Let  $p_t : [n] \to \mathbb{R}_+$  weight function (normalized to sum to 1)
  - $p_t(i)$  is the *weight* of expert *i* at time *t*
  - In beginning every expert has weight  $p_1(i) = 1/n$
  - Our total cost is  $\sum_t p_t \cdot m_t$

- **Setup:** have access to *n* experts, indexed by the set [*n*].
  - At each time step t, each expert will guess a value  $m_t(i) \in [-1, +1]$ .
  - Cost of  $i^{th}$  expert answer at time t is  $m_t(i)$
  - Parameter  $0 < \varepsilon < 1/2$
- 2 Let  $p_t : [n] \to \mathbb{R}_+$  weight function (normalized to sum to 1)
  - $p_t(i)$  is the *weight* of expert *i* at time *t*
  - In beginning every expert has weight  $p_1(i) = 1/n$
  - Our total cost is  $\sum_t p_t \cdot m_t$
- **③** Our goal is to minimize our total cost:  $\sum_{t=1}^{T} p_t \cdot m_t$

The same algorithm and argument above, applied to the setting:

- **()** Setup: have access to n experts, indexed by the set [n].
  - At each time step t, each expert will guess a value  $m_t(i) \in [-1, +1]$ .
  - Cost of  $i^{th}$  expert answer at time t is  $m_t(i)$
  - Parameter  $0 < \varepsilon < 1/2$

2 Let  $p_t : [n] \to \mathbb{R}_+$  weight function (normalized to sum to 1)

- $p_t(i)$  is the *weight* of expert *i* at time *t*
- In beginning every expert has weight  $p_1(i) = 1/n$
- Our total cost is  $\sum_t p_t \cdot m_t$

**③** Our goal is to minimize our total cost:  $\sum_{t=1}^{T} p_t \cdot m_t$ 

#### Theorem (Multiplicative Weights Update)

With the setup above, after t rounds, for any expert  $i \in [n]$ , we have:  $\sum_{t=1}^{T} p_t \cdot m_t \leq \sum_{t=1}^{T} m_t(i) + \varepsilon \cdot \sum_{t=1}^{T} |m_t(i)| + \frac{\ln n}{\varepsilon}$  • Multiplicative Weights Update

• Solving Linear Programs

Conclusion

Acknowledgements

Assume we are given LP in feasibility version:

 $Ax \ge b$  $x \ge 0$ 

Assume we are given LP in feasibility version:

min cTx

 $Ax \ge b$  $x \ge 0$ 

• Optimization version reduces to feasibility version by binary search.

 $\begin{array}{l} \begin{array}{c} g_{\mu\nu\nu} & c^{\mathsf{T}} \mathbf{x} \geq \alpha \\ \begin{pmatrix} A \\ c^{\mathsf{T}} \end{pmatrix} \mathbf{x} \geq \begin{pmatrix} \mathsf{b} \\ \mathbf{x} \end{pmatrix} \end{array}$ 

Assume we are given LP in feasibility version:

 $Ax \ge b$  $x \ge 0$ 

- Optimization version reduces to feasibility version by binary search.
- Think of x ≥ 0 being the easy constraints to satisfy, whereas Ax ≥ b are the hard ones

Assume we are given LP in feasibility version:

$$Ax \ge b$$
$$x \ge 0$$

- Optimization version reduces to feasibility version by binary search.
- Think of x ≥ 0 being the easy constraints to satisfy, whereas Ax ≥ b are the hard ones

#### Idea:

• Think of each inequality  $A_i x \ge b_i$  as an *expert* ( $A_i$  is  $i^{th}$  row of A)

$$A = \begin{pmatrix} -A_1 - \\ -A_2 - \\ \vdots \\ -A_n - \end{pmatrix}$$

Assume we are given LP in feasibility version:

 $Ax \ge b$  $x \ge 0$ 

- Optimization version reduces to feasibility version by binary search.
- Think of x ≥ 0 being the easy constraints to satisfy, whereas Ax ≥ b are the hard ones

#### Idea:

- Think of each inequality  $A_i x \ge b_i$  as an *expert* ( $A_i$  is  $i^{th}$  row of A)
- Each constraint would like to be the hardest constraint, i.e. the one that is violated the most by the current proposed solution x<sup>(t)</sup>

Assume we are given LP in feasibility version:

$$Ax \ge b$$
$$x \ge 0$$

- Optimization version reduces to feasibility version by binary search.
- Think of x ≥ 0 being the easy constraints to satisfy, whereas Ax ≥ b are the hard ones

#### Idea:

- **1** Think of each inequality  $A_i x \ge b_i$  as an *expert* ( $A_i$  is  $i^{th}$  row of A)
- Each constraint would like to be the hardest constraint, i.e. the one that is violated the most by the current proposed solution x<sup>(t)</sup>
- More precisely: cost of  $i^{th}$  constraint

$$A_i x - b_i$$

Assume we are given LP in feasibility version:

$$Ax \ge b$$
$$x \ge 0$$

- Optimization version reduces to feasibility version by binary search.
- Think of x ≥ 0 being the easy constraints to satisfy, whereas Ax ≥ b are the hard ones

Idea:

- Think of each inequality  $A_i x \ge b_i$  as an *expert* ( $A_i$  is  $i^{th}$  row of A)
- Each constraint would like to be the hardest constraint, i.e. the one that is violated the most by the current proposed solution x<sup>(t)</sup>
- More precisely: cost of  $i^{th}$  constraint

$$A_i x - b_i$$

We would like to propose feasible solution (i.e. lower cost of all constraints). Hard to deal with all constraints at the same time.

Would like to minimize

$$\min_{1\leq i\leq m}A_ix-b_i$$

 Multiplicative Weights Update (MWU) provides way of combining all constraints into one constraint!

Would like to minimize

$$\min_{1\leq i\leq m}A_ix-b_i$$

- Multiplicative Weights Update (MWU) provides way of combining all constraints into one constraint!
- MWU finds probability distribution over experts (normalized weights), which in our case are the inequalities.

Would like to minimize

$$\min_{1\leq i\leq m}A_ix-b_i$$

- Multiplicative Weights Update (MWU) provides way of combining all constraints into one constraint!
- MWU finds probability distribution over experts (normalized weights), which in our case are the inequalities.
- Thus, we have to deal with only the constraint  $p^{(t)}Ax \ge p^{(t)}b$ , where

$$p^{(t)} = \underbrace{\frac{1}{\sum_{i} w_{t}(i)} \cdot (w_{t}(1), \dots, w_{t}(n))}_{(m)} \underbrace{(envy_{t}(1), \dots, w_{t}(n))}_{(m)} \underbrace{(envy_{t}(1), \dots, w_{t}(n))}_{(m)}}_{(m)}$$

68 / 83

Would like to minimize

$$\min_{1\leq i\leq m}A_ix-b_i$$

- Multiplicative Weights Update (MWU) provides way of combining all constraints into one constraint!
- MWU finds probability distribution over experts (normalized weights), which in our case are the inequalities.
- Thus, we have to deal with only the constraint  $p^{(t)}Ax \ge p^{(t)}b$ , where

$$p^{(t)} = \frac{1}{\sum_i w_t(i)} \cdot (w_t(1), \ldots, w_t(n))$$

• MWU shows that over the long run:

The *total violation* of our *weighted constraints* will be close to the *total violation* of the *worst violated constraint*!

• Would like to minimize

$$\min_{1 \le i \le m} \frac{A_i x - b_i}{\cot} \quad \text{each constraint}$$

Would like to minimize

$$\min_{1\leq i\leq m}A_ix-b_i$$

• MWU shows that over the long run, for any inequality  $i \in [m]$ :

$$\sum_{t=1}^{T} p^{(t)} \cdot (Ax^{(t)} - b) < \frac{\log m}{\varepsilon} + \sum_{t=1}^{T} (A_i x^{(t)} - b_i) + \varepsilon \cdot \sum_{t=1}^{T} |A_i x^{(t)} - b_i|$$

Would like to minimize

$$\min_{1 \le i \le m} \widetilde{A_i x - b_i}$$

m.(i)

• MWU shows that over the long run, for any inequality  $i \in [m]$ :

$$\sum_{t=1}^{T} p^{(t)} \cdot (A_{\mathbf{x}}^{(t)} - b) < \frac{\log m}{\varepsilon} + \sum_{t=1}^{T} (A_{i} \mathbf{x}^{(t)} - b_{i}) + \varepsilon \cdot \sum_{t=1}^{T} |A_{i} \mathbf{x}^{(t)} - b_{i}|$$

• But our theorem required  $m_t(i) \in [-1,+1]$ ... How can we fix this?

# Solving Linear Programs

Would like to minimize

$$\min_{1\leq i\leq m}A_ix-b_i$$

• MWU shows that over the long run, for any inequality  $i \in [m]$ :

$$\sum_{t=1}^{T} p^{(t)} \cdot (Ax^{(t)} - b) < \frac{\log m}{\varepsilon} + \sum_{t=1}^{T} (A_i x^{(t)} - b_i) + \varepsilon \cdot \sum_{t=1}^{T} |A_i x^{(t)} - b_i|$$

- But our theorem required  $m_t(i) \in [-1,+1]$ ... How can we fix this?
- Return solution

$$x = \frac{1}{T} \cdot \sum_{t=1}^{T} x^{(t)}$$
  
Runge of all proposed.

イロン 不得 とうほう イロン 二日

73 / 83

# Solving Linear Programs

Would like to minimize

$$\min_{1\leq i\leq m}A_ix-b_i$$

• MWU shows that over the long run, for any inequality  $i \in [m]$ :

$$\sum_{t=1}^{T} p^{(t)} \cdot (Ax^{(t)} - b) < \frac{\log m}{\varepsilon} + \sum_{t=1}^{T} (A_i x^{(t)} - b_i) + \varepsilon \cdot \sum_{t=1}^{T} |A_i x^{(t)} - b_i|$$

- But our theorem required  $m_t(i) \in [-1,+1]$ ... How can we fix this?
- Return solution

$$x = \frac{1}{T} \cdot \sum_{t=1}^{T} x^{(t)}$$

- What if there is no  $x \ge 0$  such that  $p^{(t)}Ax \ge p^{(t)}b$ ?
  - Farkas' lemma  $\Rightarrow$  the system is *infeasible*, and we are done!

Algorithm:  
initialize 
$$p^{(t)}(i) = \frac{1}{m}$$
  $\forall i \in [m]$ .  
for  $t = 1$ ,  $2$ , ...,  $T$ :  
of all any  
pt there is  
to seture  
 $p^{(t)}A \chi^{(t)} \ge 0$   $|s.t|$ .  
 $p^{(t)}A \chi^{(t)} \ge p^{(t)}b$  (only one  
integration into a i  $\in [m]$ :  
if constraint i  
integration is  
for  $i \in [m]$ :  
if constraint i  
decrease weight of  $p^{(t)}(i)$   
 $\chi(t)$  for  $\chi = \frac{1}{T} \sum_{i=1}^{T} \chi^{(t)}$  (proposed solution to  
for  $LP$ .

# Solving Linear Programs

Would like to minimize

$$\min_{1\leq i\leq m}A_ix-b_i$$

• MWU shows that over the long run, for any inequality  $i \in [m]$ :

$$\sum_{t=1}^{T} p^{(t)} \cdot (Ax^{(t)} - b) < \frac{\log m}{\varepsilon} + \sum_{t=1}^{T} (A_i x^{(t)} - b_i) + \varepsilon \cdot \sum_{t=1}^{T} |A_i x^{(t)} - b_i|$$

• But our theorem required  $m_t(i) \in [-1, +1]$ ... How can we fix this?

Return solution

$$x = \frac{1}{T} \cdot \sum_{t=1}^{T} x^{(t)}$$

- What if there is no  $x \ge 0$  such that  $p^{(t)}Ax \ge p^{(t)}b$ ?
  - Farkas' lemma  $\Rightarrow$  the system is *infeasible*, and we are done!
  - Thus, we will assume that above never happens.

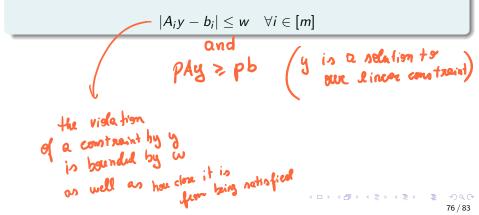
### Theorem

### Definition (Oracle)

Let  $A \in \mathbb{R}^{m \times n}$ . We say that  $\mathcal{O}$  is an oracle of *width w* for *A* if given a linear constraint

$$pAx \ge pb, x \ge 0$$

 $\mathcal{O}(p)$  will return  $y \ge 0$  such that



### Theorem

### Definition (Oracle)

Let  $A \in \mathbb{R}^{m \times n}$ . We say that  $\mathcal{O}$  is an oracle of *width w* for A if given a linear constraint

pAx > pb, x > 0

 $\mathcal{O}(p)$  will return  $y \geq 0$  such that

$$|A_iy - b_i| \le w \quad \forall i \in [m]$$

#### Theorem (Multiplicative Weights Update)

Let  $\delta > 0$  and suppose we are given an oracle with width w for A. The MWU algorithm either finds a solution  $y \ge 0$  such that · PPX mimak

$$A_i y \geq b_i - \delta \quad \forall i \in [m]$$

or concludes that the system is infeasible (and outputs a dual solution). Our algorithm makes  $O(w^2 \log(m)/\delta^2)$  oracle calls.

• As we said before, if oracle fails to find a solution, we found a separating hyperplane and we are done.

- As we said before, if oracle fails to find a solution, we found a separating hyperplane and we are done.
- Otherwise, we have that MWU algorithm with costs

- As we said before, if oracle fails to find a solution, we found a separating hyperplane and we are done.
- Otherwise, we have that MWU algorithm with costs

$$m_t(i) = rac{A_i x^{(t)} - b_i}{w}$$
 gives us that after  $T$  steps

$$O \leq \sum_{t=1}^{T} p^{(t)} \cdot \frac{A_{\boldsymbol{q}} x^{(t)} - b_{\boldsymbol{q}}}{w} < \frac{\log m}{\varepsilon} + \sum_{t=1}^{T} \frac{A_{i} x^{(t)} - b_{i}}{w} + \varepsilon \cdot \sum_{t=1}^{T} \frac{|A_{i} x^{(t)} - b_{i}|}{w}$$

Thus, we have

$$\sum_{t=1}^{T} \frac{A_i x^{(t)} - b_i}{T} \ge -\frac{w \log m}{T \cdot \varepsilon} - \varepsilon \cdot w$$

< L

- As we said before, if oracle fails to find a solution, we found a separating hyperplane and we are done.
- Otherwise, we have that MWU algorithm with costs

$$m_t(i) = \frac{A_i x^{(t)} - b_i}{w} \text{ gives us that after } T \text{ steps}$$

$$\sum_{t=1}^T p^{(t)} \cdot \frac{A_i x^{(t)} - b_i}{w} < \frac{\log m}{\varepsilon} + \sum_{t=1}^T \frac{A_i x^{(t)} - b_i}{w} + \varepsilon \cdot \sum_{t=1}^T \frac{|A_i x^{(t)} - b_i|}{w}$$

Thus, we have

$$\sum_{t=1}^{T} \frac{A_i x^{(t)} - b_i}{T} \ge -\frac{w \log m}{T \cdot \varepsilon} - \varepsilon \cdot w$$

**う**.

• Setting 
$$\varepsilon = \delta/2w$$
 and  $T = \frac{4 \cdot w^2 \cdot \log m}{\delta^2}$  we get  

$$A_i \cdot \frac{1}{T} \sum_{\substack{i=1\\ i \neq i}}^{T} \frac{\lambda^{(i)}}{1} - \sum_{\substack{i=1\\ i \neq i}}^{T} \frac{b_i}{1} = \sum_{\substack{t=1\\ t=1}}^{T} \frac{A_i x^{(t)} - b_i}{T} \ge -\delta$$

$$A_i \cdot x - b_i \cdot y - \delta$$

# Conclusion

- Online Learning
  - Experts are weak classifiers, want to choose hypothesis based on these experts
  - Ø Boosting (in learning theory)
- Solving linear programs! (today)
- Convex Optimization
- Computational Geometry
- many more

## Acknowledgement

- Lecture based largely on:
  - Lap Chi's notes
  - Yaron Singer's notes
  - Elad Hazan's survey on online optimization
- See Lap Chi's notes at https://cs.uwaterloo.ca/~lapchi/cs466/notes/L21.pdf
- See Yaron's notes https://people.seas.harvard.edu/~yaron/ AM221-S16/lecture\_notes/AM221\_lecture11.pdf
- See Elad's survey at https://arxiv.org/pdf/1909.05207.pdf
- See great survey on MWU ar https://www.cs.princeton.edu/~arora/pubs/MWsurvey.pdf