

# Nonstochastic Bandits with Unrestricted Delays

Nicolò Cesa-Bianchi<sup>1</sup>   Yevgeny Seldin<sup>2</sup>   Tobias Thune<sup>2</sup>

<sup>1</sup>Università degli Studi di Milano

<sup>2</sup>University of Copenhagen



# The bandit problem

- $K$  actions
- Unknown **deterministic** assignment of **losses** to actions

$$\ell_t = (\ell_t(1), \dots, \ell_t(K)) \in [0, 1]^K \quad \text{for } t=1, 2, \dots$$



# The bandit problem

- $K$  actions
- Unknown **deterministic** assignment of **losses** to actions

$$\ell_t = (\ell_t(1), \dots, \ell_t(K)) \in [0, 1]^K \quad \text{for } t=1, 2, \dots$$



For  $t = 1, 2, \dots$

- 1 Player selects an action  $I_t$  (possibly using randomization) and incurs loss  $\ell_t(I_t)$



# The bandit problem

- $K$  actions
- Unknown **deterministic** assignment of **losses** to actions

$$\ell_t = (\ell_t(1), \dots, \ell_t(K)) \in [0, 1]^K \quad \text{for } t=1, 2, \dots$$



For  $t = 1, 2, \dots$

- 1 Player selects an action  $I_t$  (possibly using randomization) and incurs loss  $\ell_t(I_t)$
- 2 Player observes  $\ell_t(I_t)$  and updates selection strategy



# Regret

The average loss of the player should converge to that of the best action

$$R_T \stackrel{\text{def}}{=} \underbrace{\mathbb{E} \left[ \sum_{t=1}^T \ell_t(I_t) \right]}_{\text{player's loss}} - \min_{i=1, \dots, K} \underbrace{\sum_{t=1}^T \ell_t(i)}_{\text{best action's loss}} \stackrel{\text{want}}{=} o(T)$$



The average loss of the player should converge to that of the best action

$$R_T \stackrel{\text{def}}{=} \underbrace{\mathbb{E} \left[ \sum_{t=1}^T \ell_t(I_t) \right]}_{\text{player's loss}} - \min_{i=1, \dots, K} \underbrace{\sum_{t=1}^T \ell_t(i)}_{\text{best action's loss}} \stackrel{\text{want}}{=} o(T)$$

- Expectation is with respect to the player's internal randomization
- Random play typically results in linear regret



# The Exp3 algorithm

At time  $t$  draw action  $I_t$  with probability  $p_t(i) \stackrel{\text{def}}{=} \mathbb{P}(I_t = i \mid I_1, \dots, I_{t-1})$



# The Exp3 algorithm

At time  $t$  draw action  $I_t$  with probability  $p_t(i) \stackrel{\text{def}}{=} \mathbb{P}(I_t = i \mid I_1, \dots, I_{t-1})$

- $p_t(i) \propto \exp\left(-\eta \sum_{s=1}^{t-1} \hat{\ell}_s(i)\right) \quad i = 1, \dots, K$





# The Exp3 algorithm

At time  $t$  draw action  $I_t$  with probability  $p_t(i) \stackrel{\text{def}}{=} \mathbb{P}(I_t = i \mid I_1, \dots, I_{t-1})$

- $p_t(i) \propto \exp\left(-\eta \sum_{s=1}^{t-1} \hat{\ell}_s(i)\right) \quad i = 1, \dots, K$

- $\hat{\ell}_t(i) = \begin{cases} \frac{\ell_t(i)}{p_t(i)} & \text{if } I_t = i \\ 0 & \text{otherwise} \end{cases}$

only one non-zero component



# The Exp3 algorithm

At time  $t$  draw action  $I_t$  with probability  $p_t(i) \stackrel{\text{def}}{=} \mathbb{P}(I_t = i \mid I_1, \dots, I_{t-1})$

$$\bullet p_t(i) \propto \exp\left(-\eta \sum_{s=1}^{t-1} \hat{\ell}_s(i)\right) \quad i = 1, \dots, K$$

$$\bullet \hat{\ell}_t(i) = \begin{cases} \frac{\ell_t(i)}{p_t(i)} & \text{if } I_t = i \\ 0 & \text{otherwise} \end{cases} \quad \text{only one non-zero component}$$

Exp3 regret bound

[Auer et al., 2002]

$$R_T \leq \sqrt{(\ln K) \sum_t \|\ell_t\|^2} = \mathcal{O}(\sqrt{TK \ln K})$$

# The Exp3 algorithm

At time  $t$  draw action  $I_t$  with probability  $p_t(i) \stackrel{\text{def}}{=} \mathbb{P}(I_t = i \mid I_1, \dots, I_{t-1})$

$$\bullet p_t(i) \propto \exp\left(-\eta \sum_{s=1}^{t-1} \hat{\ell}_s(i)\right) \quad i = 1, \dots, K$$

$$\bullet \hat{\ell}_t(i) = \begin{cases} \frac{\ell_t(i)}{p_t(i)} & \text{if } I_t = i \\ 0 & \text{otherwise} \end{cases} \quad \text{only one non-zero component}$$

Exp3 regret bound

[Auer et al., 2002]

$$R_T \leq \sqrt{(\ln K) \sum_t \|\ell_t\|^2} = \mathcal{O}(\sqrt{TK \ln K})$$

Tight up to log factors

# Learning with delayed losses

- In many bandit settings, like product recommendation, we observe the effects of our actions after some delay



# Learning with delayed losses

- In many bandit settings, like product recommendation, we observe the effects of our actions after some delay
- Assume the loss of action  $I_t$  is observed at time  $t + d_t$ , where  $d_1, d_2, \dots$  is an unknown and deterministic sequence of delays



# Learning with delayed losses

- In many bandit settings, like product recommendation, we observe the effects of our actions after some delay
- Assume the loss of action  $I_t$  is observed at time  $t + d_t$ , where  $d_1, d_2, \dots$  is an unknown and deterministic sequence of delays

For  $t = 1, 2, \dots$

- 1 Play  $I_t$



# Learning with delayed losses

- In many bandit settings, like product recommendation, we observe the effects of our actions after some delay
- Assume the loss of action  $I_t$  is observed at time  $t + d_t$ , where  $d_1, d_2, \dots$  is an unknown and deterministic sequence of delays

For  $t = 1, 2, \dots$

- 1 Play  $I_t$
- 2 Incur loss  $\ell_t(I_t)$



# Learning with delayed losses

- In many bandit settings, like product recommendation, we observe the effects of our actions after some delay
- Assume the loss of action  $I_t$  is observed at time  $t + d_t$ , where  $d_1, d_2, \dots$  is an unknown and deterministic sequence of delays

For  $t = 1, 2, \dots$

- 1 Play  $I_t$
- 2 Incur loss  $\ell_t(I_t)$
- 3 Observe pairs  $(s, \ell_s(I_s))$  for all  $s \leq t$  such  $s + d_s = t$





# Learning with delayed losses

- In many bandit settings, like product recommendation, we observe the effects of our actions after some delay
- Assume the loss of action  $I_t$  is observed at time  $t + d_t$ , where  $d_1, d_2, \dots$  is an unknown and deterministic sequence of delays

For  $t = 1, 2, \dots$

- 1 Play  $I_t$
- 2 Incur loss  $\ell_t(I_t)$
- 3 Observe pairs  $(s, \ell_s(I_s))$  for all  $s \leq t$  such  $s + d_s = t$

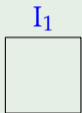
**Note:** The learner knows that  $\ell_s(I_s)$  refers to action  $I_s$  played at time  $s$



# Variable delayed feedback

Consider variable delays  $d_1, d_2, \dots$

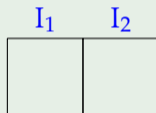
Example



# Variable delayed feedback

Consider variable delays  $d_1, d_2, \dots$

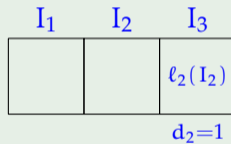
## Example



# Variable delayed feedback

Consider variable delays  $d_1, d_2, \dots$

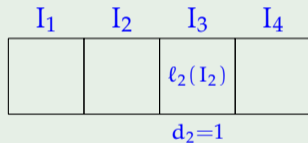
## Example



# Variable delayed feedback

Consider variable delays  $d_1, d_2, \dots$

## Example



# Variable delayed feedback

Consider variable delays  $d_1, d_2, \dots$

## Example

$I_1$	$I_2$	$I_3$	$I_4$	$I_5$
		$\ell_2(I_2)$		$\ell_3(I_3)$ $\ell_4(I_4)$
		$d_2=1$		$d_3=2$ $d_4=1$



# Variable delayed feedback

Consider variable delays  $d_1, d_2, \dots$

## Example

$I_1$	$I_2$	$I_3$	$I_4$	$I_5$	$I_6$
		$\ell_2(I_2)$		$\ell_3(I_3)$ $\ell_4(I_4)$	$\ell_1(I_1)$
		$d_2=1$		$d_3=2$ $d_4=1$	$d_1=5$

Note that  $\sum_t d_t$  can be of order  $T^2$  (all feedbacks received at the end)



# Learning with delayed losses

How does the regret depend on the **total delay**  $D = \sum_t d_t$ ?





# Learning with delayed losses

How does the regret depend on the **total delay**  $D = \sum_t d_t$ ?

Lower bound

$$\max \left\{ \underbrace{\sqrt{KT}}_{\text{bandit lower bound}}, \underbrace{\sqrt{(T+D) \ln K}}_{\text{delayed experts lower bound}} \right\} = \Omega \left( \sqrt{KT + D \ln K} \right)$$

Even with fixed delays  $d_t = d$  so that  $D = dT$



# Learning with delayed losses

How does the regret depend on the **total delay**  $D = \sum_t d_t$ ?

Lower bound

$$\max \left\{ \underbrace{\sqrt{KT}}_{\text{bandit lower bound}}, \underbrace{\sqrt{(T+D) \ln K}}_{\text{delayed experts lower bound}} \right\} = \Omega \left( \sqrt{KT + D \ln K} \right)$$

Even with fixed delays  $d_t = d$  so that  $D = dT$

Upper bound for **fixed delays**:  $R_T = \mathcal{O}(\sqrt{(K+d)T \ln K})$

[C-B et al., 2016]



- Online learning with delays

[Weinberger and Ordentlich, 2002], [Mesterharm, 2005], [Joulani, György, and Szepesvári, 2016], [Ghosh and Ramchandran, 2018]



- Online learning with delays

[Weinberger and Ordentlich, 2002], [Mesterharm, 2005], [Joulani, György, and Szepesvári, 2016], [Ghosh and Ramchandran, 2018]

- Bandits with delayed feedback

[Neu, Antos, György, and Szepesvári, 2010], [Julani, György, and Szepesvári, 2013], [C-B, Gentile, and Mansour, 2016], [Vernade, Cappé, and Perchet, 2017], [Li, Chen, and Giannakis, 2019]



- Online learning with delays  
[Weinberger and Ordentlich, 2002], [Mesterharm, 2005], [Joulani, György, and Szepesvári, 2016], [Ghosh and Ramchandran, 2018]
- Bandits with delayed feedback  
[Neu, Antos, György, and Szepesvári, 2010], [Julani, György, and Szepesvári, 2013], [C-B, Gentile, and Mansour, 2016], [Vernade, Cappé, and Perchet, 2017], [Li, Chen, and Giannakis, 2019]
- Contextual bandits with delayed feedback  
[Arya and Yang, 2019]



- Online learning with delays  
[Weinberger and Ordentlich, 2002], [Mesterharm, 2005], [Joulani, György, and Szepesvári, 2016], [Ghosh and Ramchandran, 2018]
- Bandits with delayed feedback  
[Neu, Antos, György, and Szepesvári, 2010], [Julani, György, and Szepesvári, 2013], [C-B, Gentile, and Mansour, 2016], [Vernade, Cappé, and Perchet, 2017], [Li, Chen, and Giannakis, 2019]
- Contextual bandits with delayed feedback  
[Arya and Yang, 2019]
- Bandits with delayed anonymous feedback  
[Pike-Burke, Agrawal, Szepesvári, and Grunewalder, 2017], [C-B, Gentile, and Mansour, 2018]



# A simple solution

## Delayed Exp3

At time  $t$ , make importance-weighted updates using the **old probabilities**

$$\hat{\ell}_t(i) = \begin{cases} \frac{\ell_s(i)}{p_s(i)} & \text{if } s + d_s = t \text{ and } I_s = i \\ 0 & \text{otherwise} \end{cases}$$



# A simple solution

## Delayed Exp3

At time  $t$ , make importance-weighted updates using the **old probabilities**

$$\hat{\ell}_t(i) = \begin{cases} \frac{\ell_s(i)}{p_s(i)} & \text{if } s + d_s = t \text{ and } I_s = i \\ 0 & \text{otherwise} \end{cases}$$

- Regret bound:

$$R_T = \mathcal{O}(\sqrt{(KT + D) \ln K})$$





# A simple solution

## Delayed Exp3

At time  $t$ , make importance-weighted updates using the **old probabilities**

$$\hat{\ell}_t(i) = \begin{cases} \frac{\ell_s(i)}{p_s(i)} & \text{if } s + d_s = t \text{ and } I_s = i \\ 0 & \text{otherwise} \end{cases}$$

- Regret bound:  $R_T = \mathcal{O}(\sqrt{(KT + D) \ln K})$
- Optimal to within  $\ln K$  factors 😊



# A simple solution

## Delayed Exp3

At time  $t$ , make importance-weighted updates using the **old probabilities**

$$\hat{\ell}_t(i) = \begin{cases} \frac{\ell_s(i)}{p_s(i)} & \text{if } s + d_s = t \text{ and } I_s = i \\ 0 & \text{otherwise} \end{cases}$$

- Regret bound:  $R_T = \mathcal{O}(\sqrt{(KT + D) \ln K})$
- Optimal to within  $\ln K$  factors ☺
- Holds only for certain sequences  $d_1, d_2, \dots$  of delays ☹



# A simple solution

## Delayed Exp3

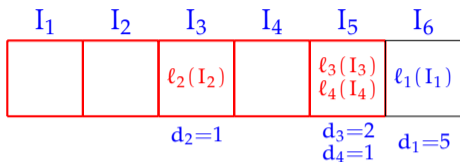
At time  $t$ , make importance-weighted updates using the **old probabilities**

$$\hat{\ell}_t(i) = \begin{cases} \frac{\ell_s(i)}{p_s(i)} & \text{if } s + d_s = t \text{ and } I_s = i \\ 0 & \text{otherwise} \end{cases}$$

- Regret bound:  $R_T = \mathcal{O}(\sqrt{(KT + D) \ln K})$
- Optimal to within  $\ln K$  factors 😊
- Holds only for certain sequences  $d_1, d_2, \dots$  of delays ☹️
- Requires tuning based on unfathomable quantities ☹️



# Why does it work?

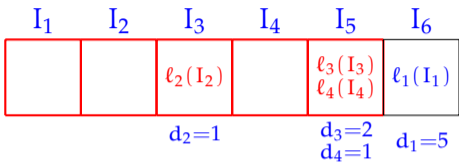


3 feedbacks received in the interval  $[1, 1 + 5)$

- The **stability span**  $N$  of a delay sequence  $d_1, d_2, \dots$  is the largest amount of feedback received in any interval  $[t, t + d_t)$



# Why does it work?

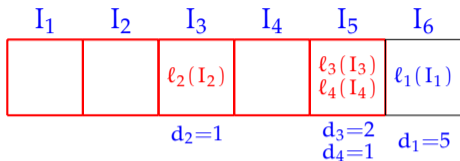


3 feedbacks received in the interval  $[1, 1 + 5)$

- The **stability span**  $N$  of a delay sequence  $d_1, d_2, \dots$  is the largest amount of feedback received in any interval  $[t, t + d_t)$
- If  $\eta \leq \frac{1}{2eN}$  then  $p_{t+d_t}(i) \leq e p_t(i)$  **stability**



# Why does it work?



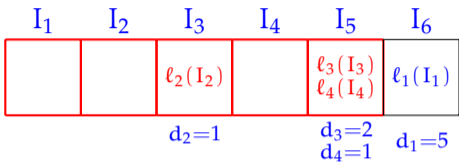
3 feedbacks received in the interval  $[1, 1 + 5)$

- The **stability span**  $N$  of a delay sequence  $d_1, d_2, \dots$  is the largest amount of feedback received in any interval  $[t, t + d_t)$
- If  $\eta \leq \frac{1}{2eN}$  then  $p_{t+d_t}(i) \leq e p_t(i)$  stability

Regret bound when  **$N$  is known**  $R_T \leq \max \left\{ \frac{\ln K}{\eta}, N \ln K \right\} + \eta(KT + D)$



# Why does it work?



3 feedbacks received in the interval  $[1, 1 + 5]$

- The **stability span**  $N$  of a delay sequence  $d_1, d_2, \dots$  is the largest amount of feedback received in any interval  $[t, t + d_t]$
- If  $\eta \leq \frac{1}{2eN}$  then  $p_{t+d_t}(i) \leq e p_t(i)$  stability

Regret bound when  **$N$  is known**  $R_T \leq \max \left\{ \frac{\ln K}{\eta}, N \ln K \right\} + \eta(KT + D)$

If  $N = \mathcal{O} \left( \sqrt{\frac{KT + D}{\ln K}} \right)$  then **tuning based on  $N$  and  $D$**  gives  **$R_T = \mathcal{O}(\sqrt{(KT + D) \ln K})$**



# Dealing with large delays

- Large delays can be ignored by the algorithm paying only constant regret each time





# Dealing with large delays

- Large delays can be ignored by the algorithm paying only constant regret each time
- There exists a tradeoff between:
  - 1 number of skipped updates
  - 2 sum of non-skipped delays  $d_t$



# Dealing with large delays

- Large delays can be ignored by the algorithm paying only constant regret each time
- There exists a tradeoff between:
  - 1 number of skipped updates
  - 2 sum of non-skipped delays  $d_t$

## The Skipper algorithm

- Skipper wraps Delayed Exp3 ignoring feedback  $\ell_t(I_t)$  when  $d_t > \delta$  (a parameter)



# Dealing with large delays

- Large delays can be ignored by the algorithm paying only constant regret each time
- There exists a tradeoff between:
  - 1 number of skipped updates
  - 2 sum of non-skipped delays  $d_t$

## The Skipper algorithm

- Skipper wraps Delayed Exp3 ignoring feedback  $\ell_t(I_t)$  when  $d_t > \delta$  (a parameter)
- Since  $N \leq 2 \max_t d_t$ , this controls the stability span perceived by Delayed Exp3



# Dealing with large delays

- Large delays can be ignored by the algorithm paying only constant regret each time
- There exists a tradeoff between:
  - 1 number of skipped updates
  - 2 sum of non-skipped delays  $d_t$

## The Skipper algorithm

- Skipper wraps Delayed Exp3 ignoring feedback  $\ell_t(I_t)$  when  $d_t > \delta$  (a parameter)
- Since  $N \leq 2 \max_t d_t$ , this controls the stability span perceived by Delayed Exp3
- As a consequence, we can do away with the knowledge of  $N$



# Regret bound for Skipper

## Regret bound for Skipper with doubling trick

$$R_T \leq S_\delta + \max \left\{ \frac{\ln K}{\eta} + \delta \ln K + \eta(D_\delta + KT) \right\}$$

- $S_\delta$  is number of skipped delays at threshold  $\delta$
- $D_\delta$  is sum of non-skipped delays



# Regret bound for Skipper

## Regret bound for Skipper with doubling trick

$$R_T \leq S_\delta + \max \left\{ \frac{\ln K}{\eta} + \delta \ln K + \eta(D_\delta + KT) \right\}$$

- $S_\delta$  is number of skipped delays at threshold  $\delta$
- $D_\delta$  is sum of non-skipped delays

By tuning  $\delta$  as a function of  $D = \sum_t d_t$  and  $T$  only, we get  $R_T = \mathcal{O}(\sqrt{(KT + D) \ln K})$   
without any assumption on  $d_1, d_2, \dots$



# No prior knowledge?

Skipper still needs knowledge of  $\sum_t d_t$  to tune the skipping threshold  $\delta$



# No prior knowledge?

Skipper still needs knowledge of  $\sum_t d_t$  to tune the skipping threshold  $\delta$

Doubling trick:

- 1 Run Skipper with  $\delta$  tuned using an optimistic guess  $D' \geq \sum_t d_t$





# No prior knowledge?

Skipper still needs knowledge of  $\sum_t d_t$  to tune the skipping threshold  $\delta$

## Doubling trick:

- 1 Run Skipper with  $\delta$  tuned using an optimistic guess  $D' \geq \sum_t d_t$
- 2 Restart Skipper with a doubled guess as soon as  $d_1 + \dots + d_t > D'$



# No prior knowledge?

Skipper still needs knowledge of  $\sum_t d_t$  to tune the skipping threshold  $\delta$

## Doubling trick:

- 1 Run Skipper with  $\delta$  tuned using an optimistic guess  $D' \geq \sum_t d_t$
- 2 Restart Skipper with a doubled guess as soon as  $d_1 + \dots + d_t > D'$

**Problem:** event  $d_1 + \dots + d_t > D'$  is only known at time  $t + d_t$



# Delay known at action time

**Assumption:** The delay  $d_t$  is observed before choosing  $I_t$



# Delay known at action time

**Assumption:** The delay  $d_t$  is observed before choosing  $I_t$

## Doubling trick for Skipper

For each epoch  $m = 1, 2, \dots$

- 1 Run Skipper with skipping threshold  $\delta \sim 2^m$



# Delay known at action time

**Assumption:** The delay  $d_t$  is observed before choosing  $I_t$

## Doubling trick for Skipper

For each epoch  $m = 1, 2, \dots$

- 1 Run Skipper with skipping threshold  $\delta \sim 2^m$
- 2 Restart Skipper based on a (complicated) condition involving:
  - current number of skipped rounds
  - current length of epoch
  - current sum of experienced delays



# The oracle bound

## Regret bound for Skipper with doubling trick

$$R_T \leq \min_{\delta} \left\{ S_{\delta} + \delta \ln K + \frac{KT + D_{\delta}}{\delta} \right\} + K \ln K$$

- $S_{\delta}$  is number of skipped delays at threshold  $\delta$
- $D_{\delta}$  is sum of non-skipped delays



# The oracle bound

## Regret bound for Skipper with doubling trick

$$R_T \leq \min_{\delta} \left\{ S_{\delta} + \delta \ln K + \frac{KT + D_{\delta}}{\delta} \right\} + K \ln K$$

- $S_{\delta}$  is number of skipped delays at threshold  $\delta$
- $D_{\delta}$  is sum of non-skipped delays

This bound is at most  $\sqrt{(KT + D) \ln K} + K \ln K$



# The oracle bound is better than the tuned Skipper bound

Comparing  $\min_{\delta} \left\{ \underbrace{S_{\delta} + \delta \ln K + \frac{KT + D_{\delta}}{\delta}}_{\text{Oracle}} \right\}$  with  $\underbrace{\sqrt{(KT + D) \ln K}}_{\text{Magic tuning}}$





# The oracle bound is better than the tuned Skipper bound

Comparing  $\underbrace{\min_{\delta} \left\{ S_{\delta} + \delta \ln K + \frac{KT + D_{\delta}}{\delta} \right\}}_{\text{Oracle}}$  with  $\underbrace{\sqrt{(KT + D) \ln K}}_{\text{Magic tuning}}$

- For  $\delta = o(T)$  choose delay sequence  $d_t = \begin{cases} T - t & \text{if } t < \delta \\ 0 & \text{otherwise} \end{cases}$



# The oracle bound is better than the tuned Skipper bound

Comparing  $\underbrace{\min_{\delta} \left\{ S_{\delta} + \delta \ln K + \frac{KT + D_{\delta}}{\delta} \right\}}_{\text{Oracle}}$  with  $\underbrace{\sqrt{(KT + D) \ln K}}_{\text{Magic tuning}}$

- For  $\delta = o(T)$  choose delay sequence  $d_t = \begin{cases} T - t & \text{if } t < \delta \\ 0 & \text{otherwise} \end{cases}$
- Then  $\sum_t d_t = \Theta(T\delta)$



# The oracle bound is better than the tuned Skipper bound

Comparing  $\underbrace{\min_{\delta} \left\{ S_{\delta} + \delta \ln K + \frac{KT + D_{\delta}}{\delta} \right\}}_{\text{Oracle}}$  with  $\underbrace{\sqrt{(KT + D) \ln K}}_{\text{Magic tuning}}$

- For  $\delta = o(T)$  choose delay sequence  $d_t = \begin{cases} T - t & \text{if } t < \delta \\ 0 & \text{otherwise} \end{cases}$
- Then  $\sum_t d_t = \Theta(T\delta)$
- ... but  $D_{\delta} = 0$  and  $S_{\delta} < \delta$



# The oracle bound is better than the tuned Skipper bound

Comparing  $\underbrace{\min_{\delta} \left\{ S_{\delta} + \delta \ln K + \frac{KT + D_{\delta}}{\delta} \right\}}_{\text{Oracle}}$  with  $\underbrace{\sqrt{(KT + D) \ln K}}_{\text{Magic tuning}}$

- For  $\delta = o(T)$  choose delay sequence  $d_t = \begin{cases} T - t & \text{if } t < \delta \\ 0 & \text{otherwise} \end{cases}$
- Then  $\sum_t d_t = \Theta(T\delta)$
- ... but  $D_{\delta} = 0$  and  $S_{\delta} < \delta$

Taking  $\delta = \sqrt{KT/(\ln K)}$  the corresponding regret bounds are of order

- Oracle:  $\sqrt{T}$
- Magic tuning:  $T^{3/4}$



# Experiments



# Some modeling choices

- Bounded vs. unbounded delay



# Some modeling choices

- Bounded vs. unbounded delay
- Time-stamped feedback:  $(I_s, \ell_s(I_s))$  instead of  $(s, \ell_s(I_s))$



# Some modeling choices

- Bounded vs. unbounded delay
- Time-stamped feedback:  $(I_s, \ell_s(I_s))$  instead of  $(s, \ell_s(I_s))$
- Anonymous feedback

$I_1$	$I_2$	$I_3$	$I_4$	$I_5$
		$\ell_2(I_2)$		$\ell_1(I_1)$ + $\ell_3(I_3)$





# Some modeling choices

- Bounded vs. unbounded delay
- Time-stamped feedback:  $(I_s, \ell_s(I_s))$  instead of  $(s, \ell_s(I_s))$
- Anonymous feedback
- Delay known at action vs. observation time

$I_1$	$I_2$	$I_3$	$I_4$	$I_5$
		$\ell_2(I_2)$		$\ell_1(I_1)$ + $\ell_3(I_3)$



# Summary

Setting	Regret Bound	Reference



# Summary

Setting	Regret Bound	Reference
Fixed delay	$\sqrt{(K + d)T \ln K}$	[C-B et al., 2016]



# Summary

Setting	Regret Bound	Reference
Fixed delay	$\sqrt{(K + d)T \ln K}$	[C-B et al., 2016]
Action time	$\min_{\delta} \left( S_{\delta} + \delta \ln K + \frac{KT + D_{\delta}}{\delta} \right)$	This work



# Summary

Setting	Regret Bound	Reference
Fixed delay	$\sqrt{(K + d)T \ln K}$	[C-B et al., 2016]
Action time	$\min_{\delta} \left( S_{\delta} + \delta \ln K + \frac{KT + D_{\delta}}{\delta} \right)$	This work
Observation time (known $T, D$ )	$\sqrt{(KT + D) \ln K}$	This work



# Summary

Setting	Regret Bound	Reference
Fixed delay	$\sqrt{(K + d)T \ln K}$	[C-B et al., 2016]
Action time	$\min_{\delta} \left( S_{\delta} + \delta \ln K + \frac{KT + D_{\delta}}{\delta} \right)$	This work
Observation time (known $T, D$ )	$\sqrt{(KT + D) \ln K}$	This work
Anonymous (known $d_{\max}$ )	$\sqrt{d_{\max}KT \ln K}$	[C-B et al., 2018]

