

On adaptive regret bounds for non- stochastic bandits

Gergely Neu

INRIA Lille, SequeL team

→ Universitat Pompeu Fabra, Barcelona

Outline

- Online learning and bandits
- Adaptive bounds in online learning
- Adaptive bounds for bandits
 - What we already have
 - What's new: **First-order** bounds
 - What may be possible
 - What seems **impossible***

*Opinion alert!

Online learning and non-stochastic bandits

For each round $t = 1, 2, \dots, T$

- Learner chooses **action** $I_t \in \{1, 2, \dots, N\}$
- Environment chooses **losses** $\ell_{t,i} \in [0, 1]$ for all i
- Learner suffers loss ℓ_{t, I_t}
- Learner observes **losses** $\ell_{t,i}$ for **all** i

Online learning and non-stochastic bandits

For each round $t = 1, 2, \dots, T$

- Learner chooses **action** $I_t \in \{1, 2, \dots, N\}$
- Environment chooses **losses** $\ell_{t,i} \in [0, 1]$ for all i
- Learner suffers loss ℓ_{t,I_t}
- Learner observes **losses** $\ell_{t,i}$ for **all** i

For each round $t = 1, 2, \dots, T$

- Learner chooses **action** $I_t \in \{1, 2, \dots, N\}$
- Environment chooses **losses** $\ell_{t,i} \in [0, 1]$ for all i
- Learner suffers loss ℓ_{t,I_t}
- Learner observes its own **loss** ℓ_{t,I_t}

Online learning and non-stochastic bandits

For each round $t = 1, 2, \dots, T$

- Learner chooses **action** $I_t \in \{1, 2, \dots, N\}$
- Environment chooses **losses** $\ell_{t,i} \in [0, 1]$ for all i
- Learner suffers loss ℓ_{t,I_t}
- Learner observes **losses** $\ell_{t,i}$ for **all** i

Need to explore!

For each round $t = 1, 2, \dots, T$

- Learner chooses **action** $I_t \in \{1, 2, \dots, N\}$
- Environment chooses **losses** $\ell_{t,i} \in [0, 1]$ for all i
- Learner suffers loss ℓ_{t,I_t}
- Learner observes its own **loss** ℓ_{t,I_t}

Minimax regret

- Define (expected) **regret against action i** as

$$R_{T,i} = \mathbf{E} \left[\sum_{t=1}^T \ell_{t,I_t} - \sum_{t=1}^T \ell_{t,i} \right]$$

- Goal: minimize regret against the **best action i^***

$$R_T = R_{T,i^*} = \max_i R_{T,i}$$

Minimax regret

- Define (expected) **regret against action i** as

$$R_{T,i} = \mathbf{E} \left[\sum_{t=1}^T \ell_{t,I_t} - \sum_{t=1}^T \ell_{t,i} \right]$$

- Goal: minimize regret against the **best action i^***

$$R_T = R_{T,i^*} = \max_i R_{T,i}$$

Full information

$$R_T = \Theta(\sqrt{T \log N})$$

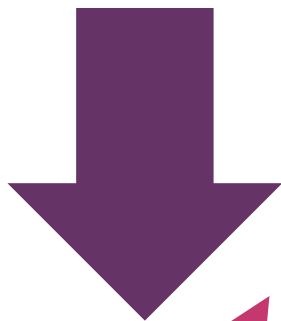
Bandit

$$R_T = \Theta(\sqrt{NT})$$

Beyond minimax: i.i.d. losses

Full information

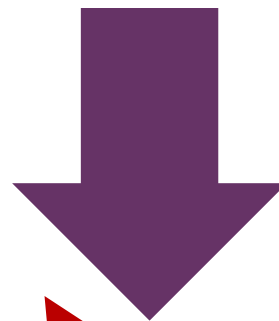
$$R_T = \Theta(\sqrt{T \log N})$$



$$\Theta(\log N)$$

Bandit

$$R_T = \Theta(\sqrt{NT})$$



$$\Theta(N \log T)$$

Beyond minimax: "Higher-order" bounds

	Full information	Bandit
minimax	$R_T = O(\sqrt{T \log N})$	$R_T = O(\sqrt{NT})$
first-order $L_{T,i} = \sum_t \ell_{t,i}$	$R_T = O(\sqrt{L_{T,i^*} \log N})$	
second-order $S_{T,i} = \sum_t \ell_{t,i}^2$	$R_T = O(\sqrt{S_{T,i^*} \log N})$ ★ Cesa-Bianchi, Mansour, Stoltz (2005)	
variance $V_{T,i} = \sum_t (\ell_{t,i} - m)^2$	$R_T = O(\sqrt{V_{T,i^*} \log N})$ ★ Hazan and Kale (2010)	

★ with a little cheating

Beyond minimax: "Higher-order" bounds

	Full information	Bandit
minimax	$R_T = O(\sqrt{T \log N})$	$R_T = O(\sqrt{NT})$
first-order $L_{T,i} = \sum_t \ell_{t,i}$	$R_T = O(\sqrt{L_{T,i^*} \log N})$	
second-order $S_{T,i} = \sum_t \ell_{t,i}^2$	$R_T = O(\sqrt{S_{T,i^*} \log N})$ ★ Cesa-Bianchi, Mansour, Stoltz (2005)	$R_T = \tilde{O}(\sqrt{\sum_i S_{t,i}})$ Auer et al. (2002) + some hacking
variance $V_{T,i} = \sum_t (\ell_{t,i} - m)^2$	$R_T = O(\sqrt{V_{T,i^*} \log N})$ ★ Hazan and Kale (2010)	$R_T = \tilde{O}(N^2 \sqrt{\sum_i V_{t,i}})$ Hazan and Kale (2011)

★ with a little cheating

Beyond minimax: “Higher-order” bounds

	Full information	Bandit
minimax	$R_T = O(\sqrt{T \log N})$	$R_T = O(\sqrt{NT})$
first-order $L_{T,i} = \sum_t \ell_{t,i}$	$R_T = O(\sqrt{L_{T,i}^* \log N})$	(it's complicated)
second-order $S_{T,i} = \sum_t \ell_{t,i}^2$	$R_T = O(\sqrt{S_{T,i}^* \log N})$ ★ Cesa-Bianchi, Mansour, Stoltz (2005)	$R_T = \tilde{O}(\sqrt{\sum_i S_{t,i}})$ Auer et al. (2002) + some hacking
variance $V_{T,i} = \sum_t (\ell_{t,i} - m)^2$	$R_T = O(\sqrt{V_{T,i}^* \log N})$ ★ Hazan and Kale (2010)	$R_T = \tilde{O}(N^2 \sqrt{\sum_i V_{t,i}})$ Hazan and Kale (2011)

★ with a little cheating

First-order bounds for bandits

(it's complicated)

- “Small-gain” bounds:
 - Consider the gain game with $g_{t,i} = 1 - \ell_{t,i}$
 - Auer, Cesa-Bianchi, Freund and Schapire (2002):

$$R_T = O(\sqrt{NG_{T,i^*} \log N})$$

$$G_{T,i} = \sum_t g_{t,i}$$

First-order bounds for bandits

(it's complicated)

- “Small-gain” bounds:
 - Consider the gain game with $g_{t,i} = 1 - \ell_{t,i}$
 - Auer, Cesa-Bianchi, Freund and Schapire (2002):

$$R_T = O(\sqrt{NG_{T,i^*} \log N})$$

$$G_{T,i} = \sum_t g_{t,i}$$

Problem:
only good if best expert is **bad!**

First-order bounds for bandits

(it's complicated)

- “Small-gain” bounds: $R_T = O(\sqrt{NG_{T,i^*} \log N})$
- A little trickier analysis gives

$$R_T = O(\sqrt{\sum_t \sum_i g_{t,i} \log N})$$

or

$$R_T = O(\sqrt{\sum_t \sum_i \ell_{t,i} \log N})$$

First-order bounds for bandits

(it's complicated)

- “Small-gain” bounds: $R_T = O(\sqrt{NG_{T,i^*} \log N})$
- A little trickier analysis gives

$$R_T = O(\sqrt{\sum_t \sum_i g_{t,i} \log N})$$

or

$$R_T = O(\sqrt{\sum_t \sum_i \ell_{t,i} \log N})$$

Problem:

one misbehaving action ruins the bound!

First-order bounds for bandits

(it's complicated)

- “Small-gain” bounds: $R_T = O(\sqrt{NG_{T,i^*} \log N})$
- A little trickier analysis gives $R_T = O(\sqrt{\sum_t \sum_i \ell_{t,i} \log N})$
- Some obscure **actual** first-order bounds:
 - › Stoltz (2005): $N\sqrt{L_T^*}$
 - › Allenberg, Auer, Györfi and Ottucsák (2006): $\sqrt{NL_T^*}$
 - › Rakhlin and Sridharan (2013): $N^{3/2}\sqrt{L_T^*}$

First-order bounds for bandits

(it's complicated)

- “Small-gain” bounds: $R_T = O(\sqrt{NG_{T,i^*} \log N})$
- A little trickier analysis gives $R_T = O(\sqrt{\sum_t \sum_i \ell_{t,i} \log N})$
- Some obscure **actual** first-order bounds:

›
›
›

Problem:
no real insight from analyses!

$$\sqrt{NL_T^*}$$

First-order bounds for non-stochastic bandits

A typical bandit algorithm

For every round $t = 1, 2, \dots, T$

- Choose arm $I_t = i$ with probability $p_{t,i}$
- Compute unbiased loss estimate

$$\hat{\ell}_{t,i} = \frac{\ell_{t,i}}{p_{t,i}} \mathbf{1}_{\{I_t=i\}}$$

- Use $\hat{\ell}_{t,i}$ in a **black-box** online learning algorithm to compute \mathbf{p}_{t+1}

A typical regret bound

$$R_T \leq \frac{\log N}{\eta} + \eta \mathbf{E} \left[\sum_{t=1}^T \sum_{i=1}^N p_{t,i} (\hat{\ell}_{t,i})^2 \right]$$

(η : “learning rate”)

A typical regret bound

$$R_T \leq \frac{\log N}{\eta} + \eta \mathbf{E} \left[\sum_{t=1}^T \sum_{i=1}^N p_{t,i} (\hat{\ell}_{t,i})^2 \right]$$

$$\leq \frac{\log N}{\eta} + \eta \mathbf{E} \left[\sum_{t=1}^T \sum_{i=1}^N \hat{\ell}_{t,i} \right]$$

(η : “learning rate”)

A typical regret bound

$$R_T \leq \frac{\log N}{\eta} + \eta \mathbf{E} \left[\sum_{t=1}^T \sum_{i=1}^N p_{t,i} (\hat{\ell}_{t,i})^2 \right]$$

$$\leq \frac{\log N}{\eta} + \eta \mathbf{E} \left[\sum_{t=1}^T \sum_{i=1}^N \hat{\ell}_{t,i} \right]$$

$$= \frac{\log N}{\eta} + \eta \sum_{i=1}^N L_{T,i}$$

(η : “learning rate”)

A typical regret bound

$$\begin{aligned} R_T &\leq \frac{\log N}{\eta} + \eta \mathbf{E} \left[\sum_{t=1}^T \sum_{i=1}^N p_{t,i} (\hat{\ell}_{t,i})^2 \right] \\ &\leq \frac{\log N}{\eta} + \eta \mathbf{E} \left[\sum_{t=1}^T \sum_{i=1}^N \hat{\ell}_{t,i} \right] \\ &= \frac{\log N}{\eta} + \eta \sum_{i=1}^N L_{T,i} = \tilde{O} \left(\sqrt{\sum_{i=1}^N L_{T,i}} \right) \end{aligned}$$

(η : “learning rate”)

(for appropriate η)

A typical regret bound

$$R_T = \tilde{o} \left(\sqrt{\sum_{i=1}^N L_{T,i}} \right)$$

A typical regret bound

$$R_T = \tilde{o} \left(\sqrt{\sum_{i=1}^N L_{T,i}} \right)$$

It's all because
 $\mathbf{E}[\hat{L}_{T,i}] = L_{T,i}!!!$

A typical regret bound

$$R_T = \tilde{O} \left(\sqrt{\sum_{i=1}^N L_{T,i}} \right)$$

It's all because
 $\mathbf{E}[\hat{L}_{T,i}] = L_{T,i}!!!$

Idea: try to enforce
 $\mathbf{E}[\hat{L}_{T,i}] = O(L_{T,i}^*)$

A typical regret bound

$$R_T = \sum_{i=1}^T \left(\hat{L}_{T,i} - L_{T,i} \right)$$

Need optimistic estimates!

It's all because
 $\mathbf{E}[\hat{L}_{T,i}] = L_{T,i}!!!$

Idea: try to enforce
 $\mathbf{E}[\hat{L}_{T,i}] = O(L_{T,i}^*)$

A typical algorithm – fixed!

For every round $t = 1, 2, \dots, T$

- Choose arm $I_t = i$ with probability $p_{t,i}$
- Compute unbiased loss estimate

$$\hat{\ell}_{t,i} = \frac{\ell_{t,i}}{p_{t,i}} \mathbf{1}_{\{I_t=i\}}$$

- Use $\hat{\ell}_{t,i}$ in a **black-box** online learning algorithm to compute p_{t+1}

A typical algorithm – fixed!

For every round $t = 1, 2, \dots, T$

- Choose arm $I_t = i$ with probability $p_{t,i}$
- Compute **biased** loss estimate

$$\hat{\ell}_{t,i} = \frac{\ell_{t,i}}{p_{t,i} + \gamma} \mathbf{1}_{\{I_t=i\}}$$

- Use $\hat{\ell}_{t,i}$ in a **black-box** online learning algorithm to compute p_{t+1}

“Implicit
exploration”
(Kocák, N, Valko and
Munos, 2015)

Algorithm: Follow the Perturbed Leader

(Kalai and Vempala, 2005, Poland, 2005)

For every round $t = 1, 2, \dots, T$

- Draw perturbation $Z_{t,i} \sim \text{Exp}(1)$ for all i
- Choose arm $I_t = \arg \min_i (\eta \hat{L}_{t-1,i} - Z_{t,i})$
- Compute **biased** loss estimate

$$\hat{\ell}_{t,i} = \frac{\ell_{t,i}}{p_{t,i} + \gamma} \mathbf{1}_{\{I_t=i\}}$$

- Update $\hat{L}_{t,i} = \hat{L}_{t-1,i} + \hat{\ell}_{t,i}$

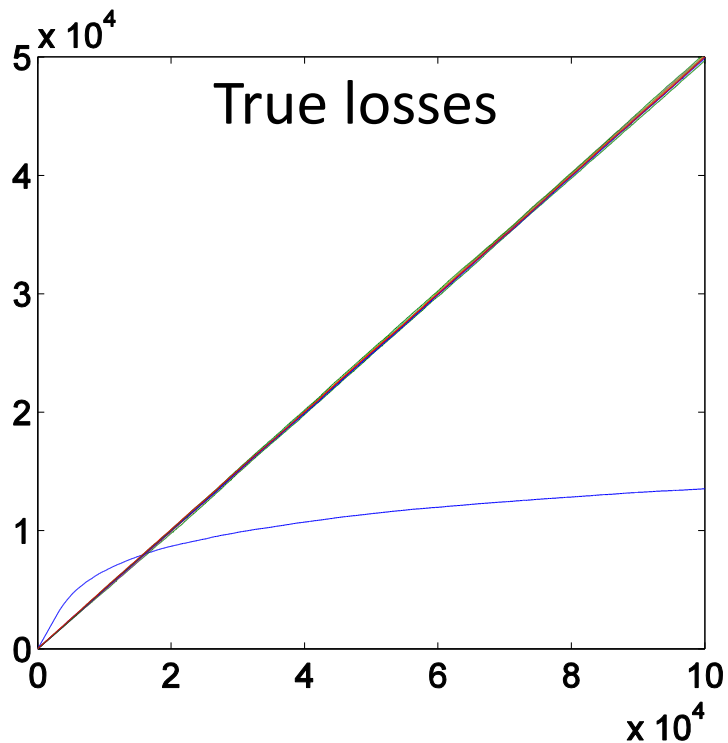
“Implicit
exploration”
(Kocák, N, Valko and
Munos, 2015)

Implicit exploration in action

$$\hat{\ell}_{t,i} = \frac{\ell_{t,i}}{p_{t,i} + \gamma} \mathbf{1}_{\{I_t=i\}}$$

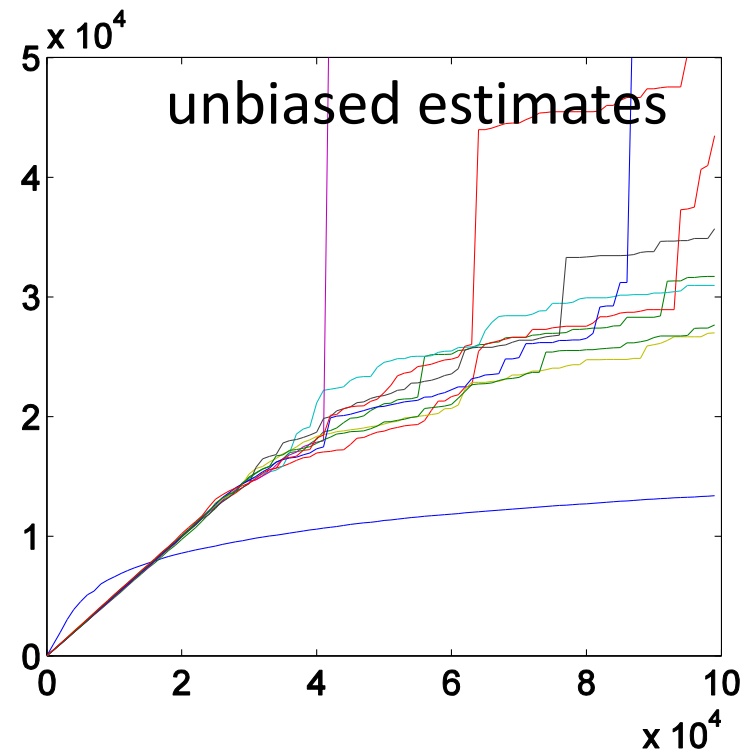
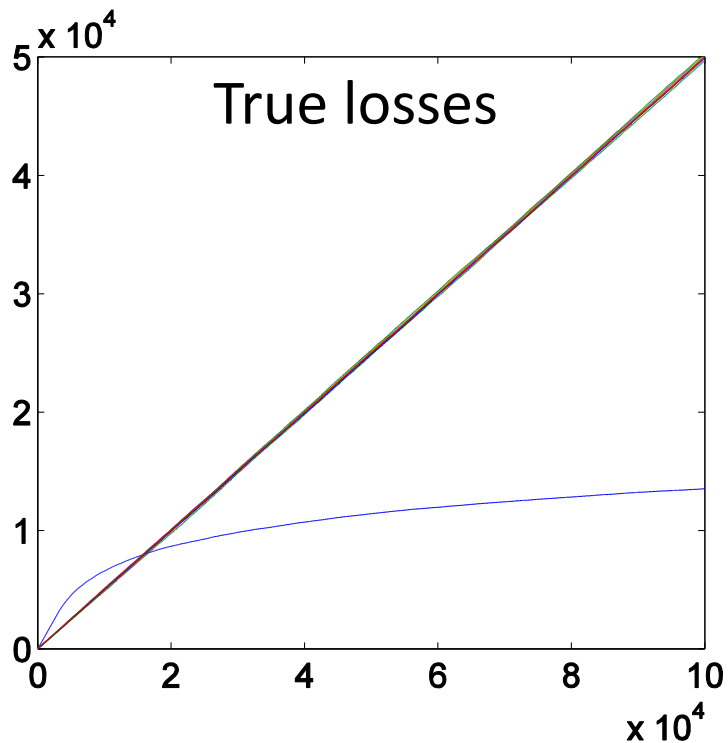
Implicit exploration in action

$$\hat{\ell}_{t,i} = \frac{\ell_{t,i}}{p_{t,i} + \gamma} \mathbf{1}_{\{I_t=i\}}$$



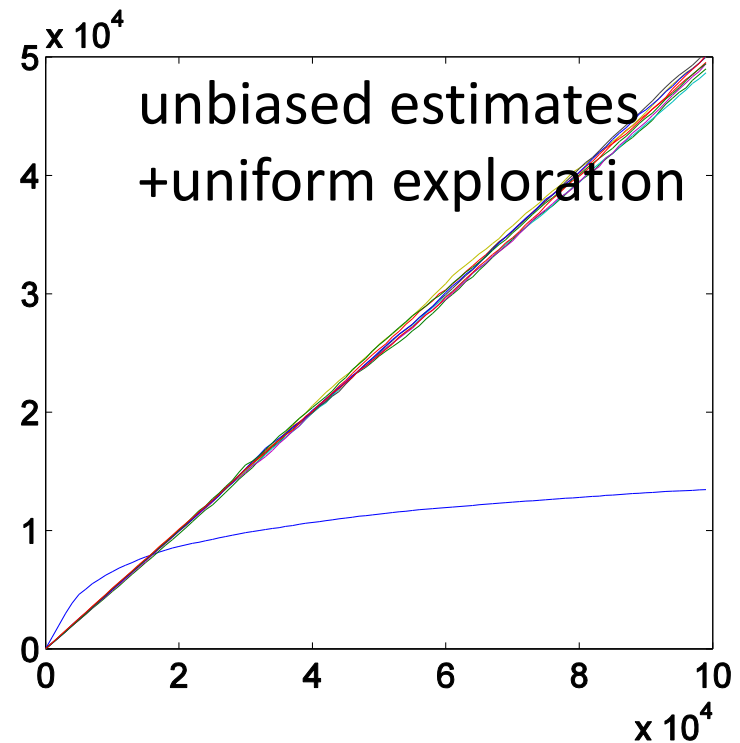
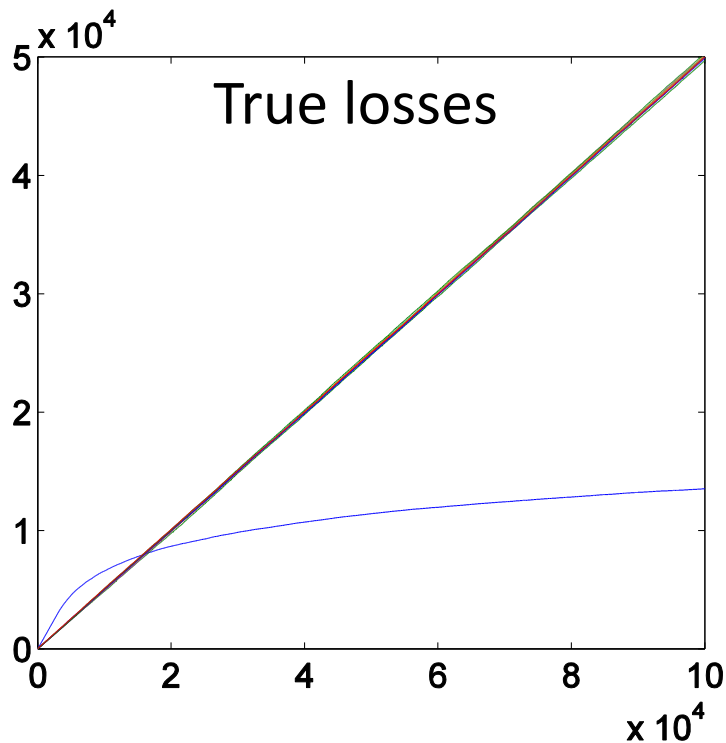
Implicit exploration in action

$$\hat{\ell}_{t,i} = \frac{\ell_{t,i}}{p_{t,i} + \gamma} \mathbf{1}_{\{I_t=i\}}$$



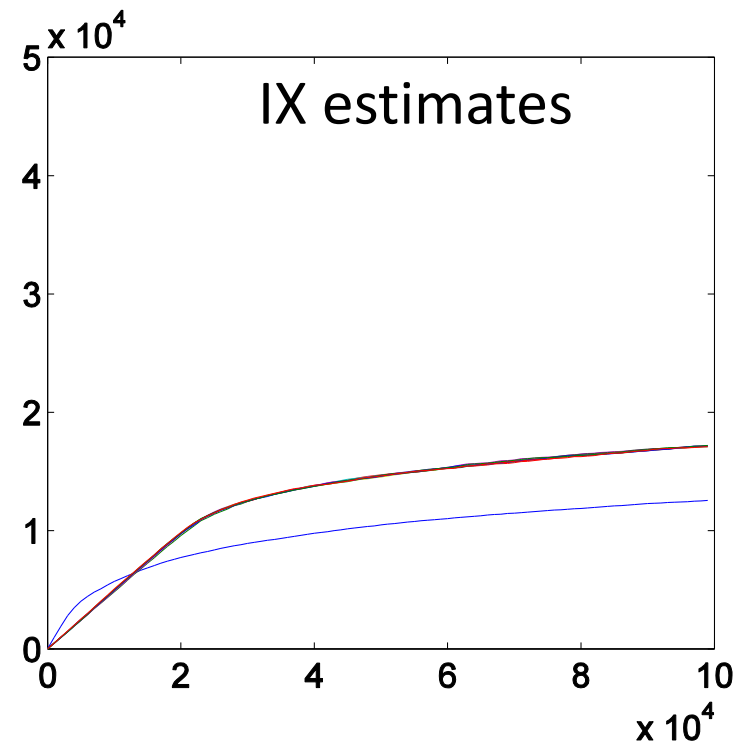
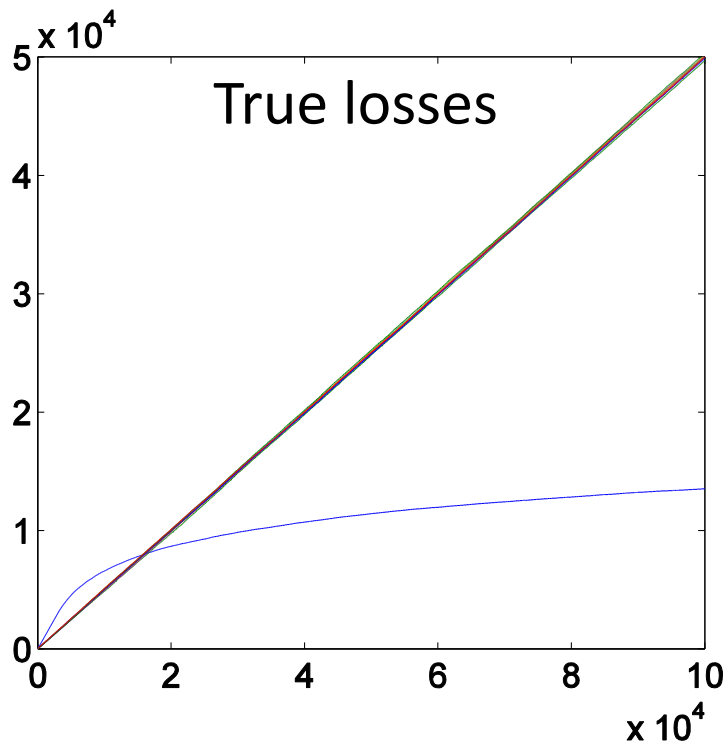
Implicit exploration in action

$$\hat{\ell}_{t,i} = \frac{\ell_{t,i}}{p_{t,i} + \gamma} \mathbf{1}_{\{I_t=i\}}$$



Implicit exploration in action

$$\hat{\ell}_{t,i} = \frac{\ell_{t,i}}{p_{t,i} + \gamma} \mathbf{1}_{\{I_t=i\}}$$



Optimistic estimates

Lemma (N, 2015a): Assume that $Z_{t,i} \leq B$ for all t and i . Then, for any i and j ,

$$\hat{L}_{T,i} \leq \hat{L}_{T,j} + \frac{(\log N + B)}{\eta} + \frac{1}{\gamma}$$

Optimistic estimates

Lemma (N, 2015a): Assume that $Z_{t,i} \leq B$ for all t and i . Then, for any i and j ,

$$\hat{L}_{T,i} \leq \hat{L}_{T,j} + \frac{(\log N + B)}{\eta} + \frac{1}{\gamma}$$

All perturbations are nicely bounded with high probability \rightarrow bad arms are suppressed!

Suppressing bad arms

$\hat{\ell}_{t,1}$

$\hat{\ell}_{t,2}$

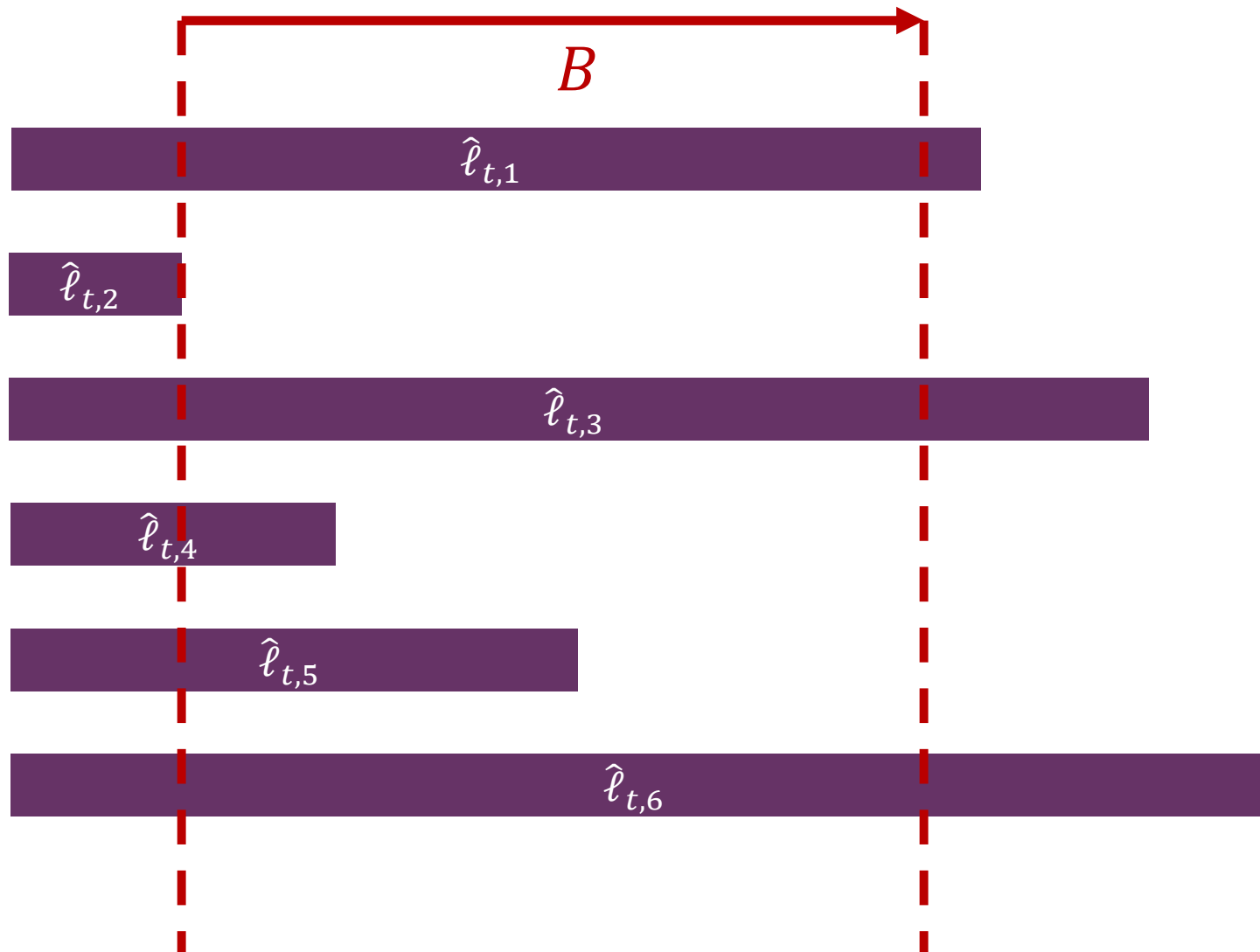
$\hat{\ell}_{t,3}$

$\hat{\ell}_{t,4}$

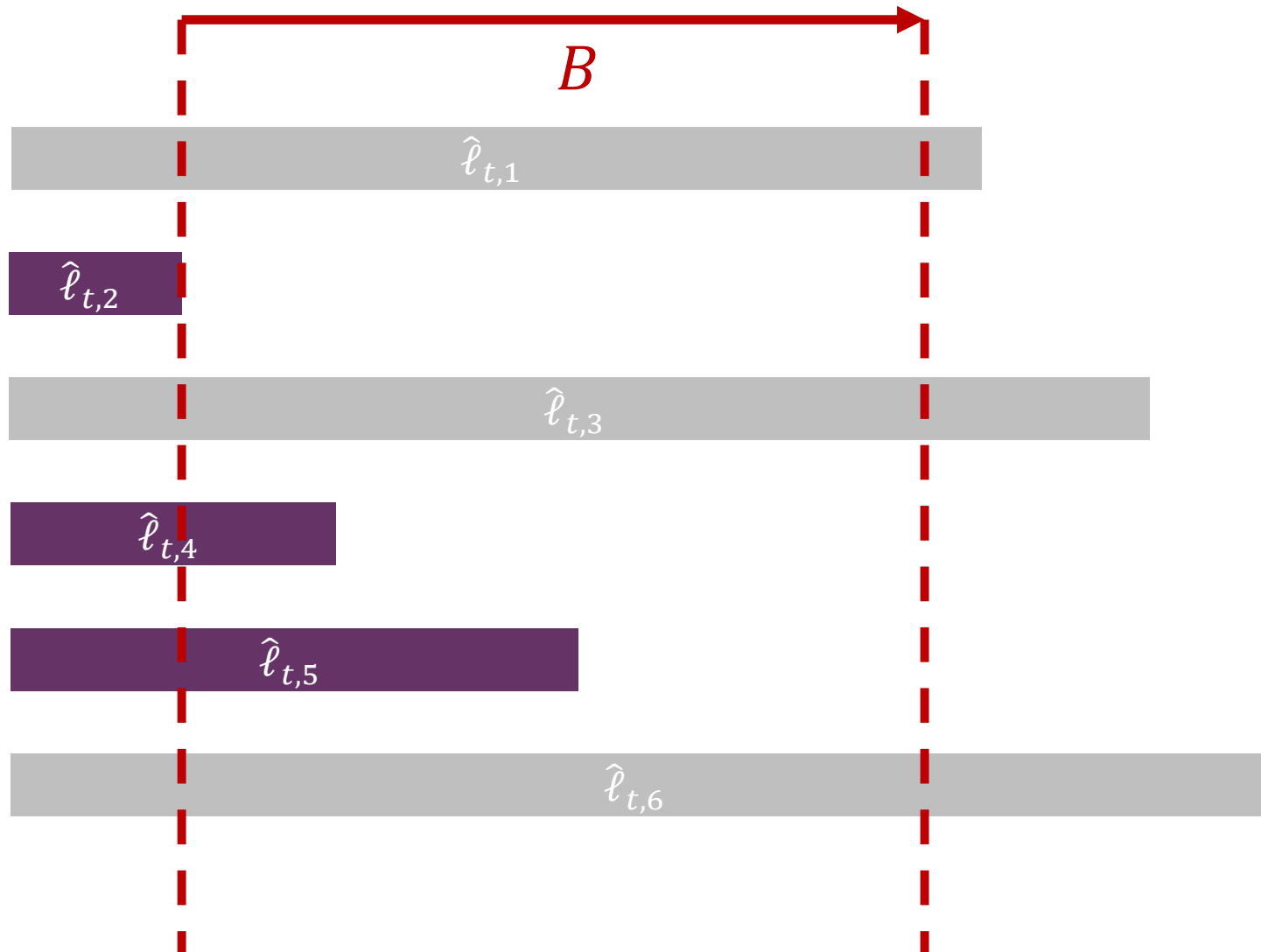
$\hat{\ell}_{t,5}$

$\hat{\ell}_{t,6}$

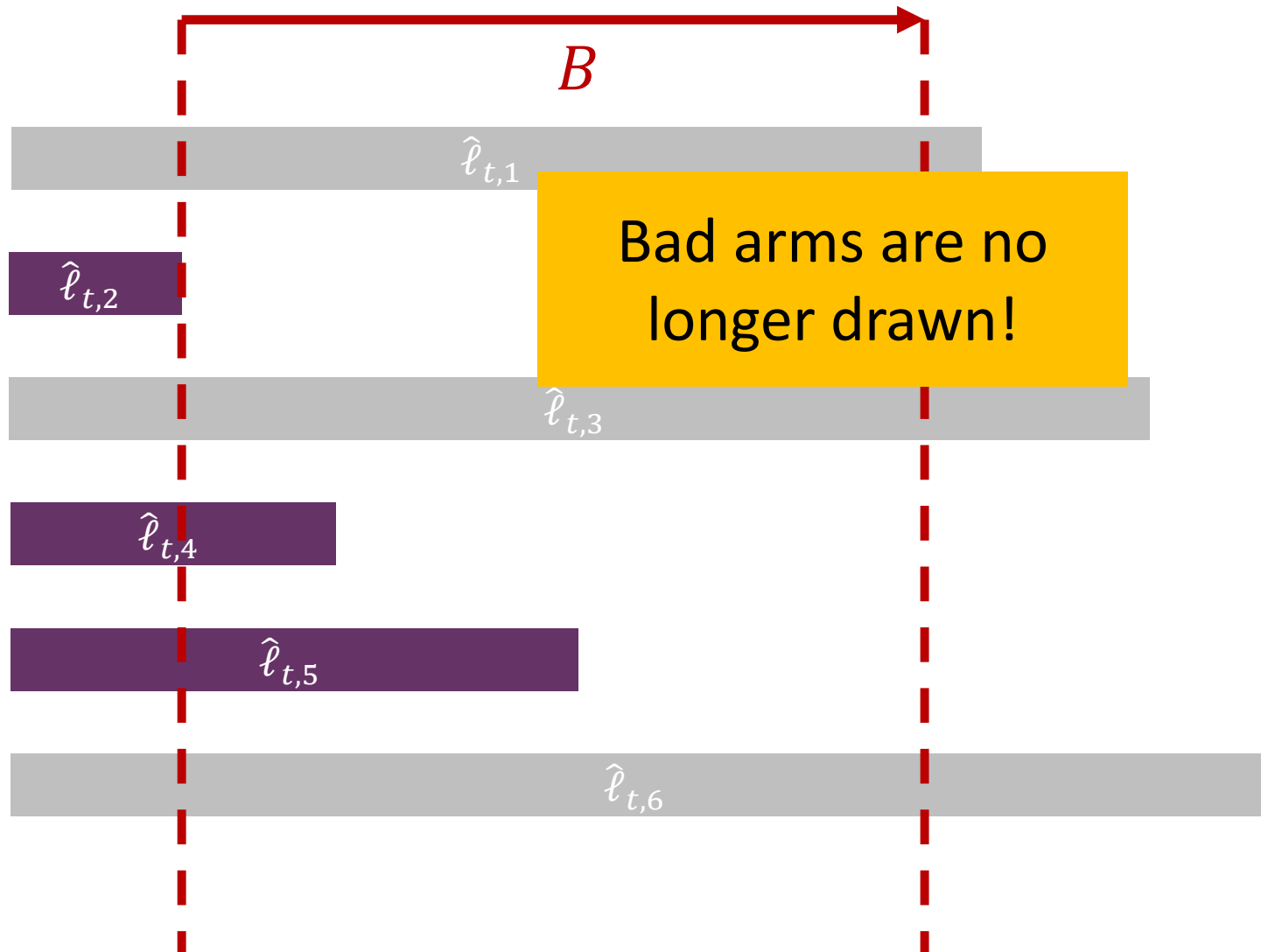
Suppressing bad arms



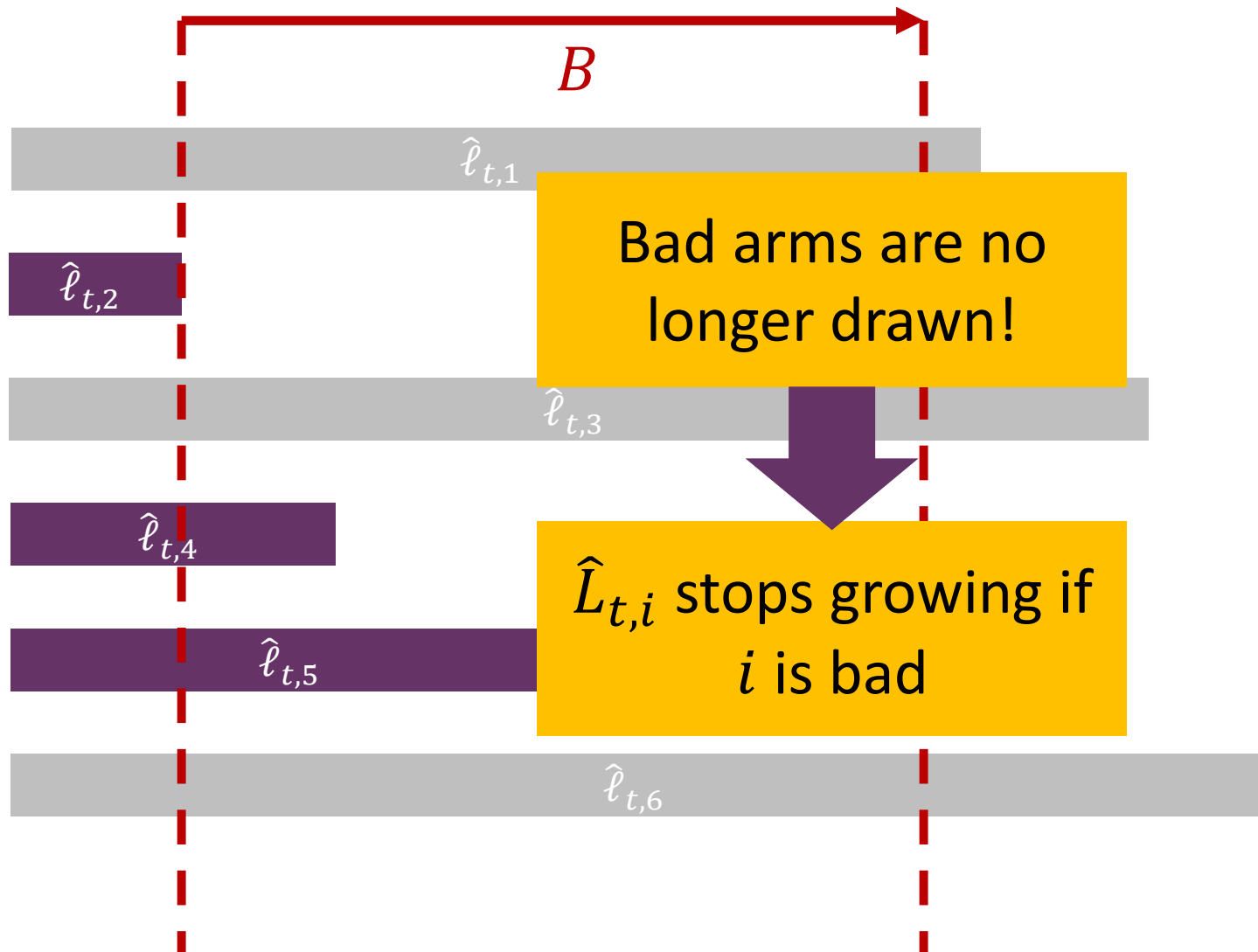
Suppressing bad arms



Suppressing bad arms



Suppressing bad arms



Finally: a first-order bound!

$$\begin{aligned} R_T &\leq \frac{\log N}{\eta} + \eta \mathbf{E} \left[\sum_{t=1}^T \sum_{i=1}^N p_{t,i} (\hat{\ell}_{t,i})^2 \right] \\ &\leq \frac{\log N}{\eta} + \eta \mathbf{E} \left[\sum_{t=1}^T \sum_{i=1}^N \hat{\ell}_{t,i} \right] \\ &= \frac{\log N}{\eta} + \eta \sum_{i=1}^N L_{T,i} = \tilde{O} \left(\sqrt{\sum_{i=1}^N L_{T,i}} \right) \end{aligned}$$

(for appropriate η)

Finally: a first-order bound!

$$R_T \leq$$

$$\frac{\log N}{\eta} + \eta \mathbf{E} \left[\sum_{t=1}^T \sum_{i=1}^N p_{t,i} (\hat{\ell}_{t,i})^2 \right]$$

$$\leq$$

$$\frac{\log N}{\eta} + \eta \mathbf{E} \left[\sum_{t=1}^T \sum_{i=1}^N \hat{\ell}_{t,i} \right]$$

Lemma

$$\leq$$

$$\frac{\log N}{\eta} + \eta N L_{T,i^*} + N(B + \log N)$$

$$= \tilde{O}(\sqrt{N L_{T,i^*}})$$

(for appropriate η, γ, B)

Finally: a first-order bound!

$$R_T = \tilde{O}\left(\sqrt{NL_{T,i^*}}\right)$$

Parameters:

- Set $\gamma = \eta/2$

Finally: a first-order bound!

$$R_T = \tilde{O}\left(\sqrt{NL_{T,i^*}}\right)$$

Parameters:

- Set $\gamma = \eta/2$
- If we know L_{T,i^*} : $\eta = \sqrt{\frac{(\log N + 1)}{NL_{T,i^*}}}$

Finally: a first-order bound!

$$R_T = \tilde{O}\left(\sqrt{NL_{T,i^*}}\right)$$

Parameters:

- Set $\gamma = \eta/2$
- If we know L_{T,i^*} : $\eta = \sqrt{\frac{(\log N + 1)}{NL_{T,i^*}}}$
- If we don't: $\eta_t = \sqrt{\frac{(\log N + 1)}{N(1 + \sum_i \hat{L}_{t-1,i})}}$

Finally: a first-order bound!

$$R_T = \tilde{O}\left(\sqrt{NL_{T,i^*}}\right)$$

Parameters:

- Set $\gamma = \eta/2$
- If we know L_{T,i^*} : $\eta = \sqrt{\frac{(\log N+1)}{NL_{T,i^*}}}$
- If we don't: $\eta_t = \sqrt{\frac{(\log N+1)}{N(1+\sum_i \hat{L}_{t-1,i})}}$

Arguments also extend to combinatorial semi-bandits!

What's next?

Beyond minimax: “Higher-order” bounds

	Full information	Bandit
minimax	$R_T = \Theta(\sqrt{T \log N})$	$R_T = O(\sqrt{NT})$
first-order $L_{T,i} = \sum_t \ell_{t,i}$	$R_T = O(\sqrt{L_{T,i^*} \log N})$	(it's complicated)
second-order $S_{T,i} = \sum_t \ell_{t,i}^2$	$R_T = O(\sqrt{S_{T,i^*} \log N})$ ★ Cesa-Bianchi, Mansour, Stoltz (2005)	$R_T = \tilde{O}(\sqrt{\sum_i S_{t,i}})$ Auer et al. (2002) + some hacking
variance $V_{T,i} = \sum_t (\ell_{t,i} - m)^2$	$R_T = O(\sqrt{V_{T,i^*} \log N})$ ★ Hazan and Kale (2010)	$R_T = \tilde{O}(N^2 \sqrt{\sum_i V_{t,i}})$ Hazan and Kale (2011)

★ with a little cheating

Beyond minimax: "Higher-order" bounds

	Full information	Bandit
minimax	$R_T = \Theta(\sqrt{T \log N})$	$R_T = O(\sqrt{NT})$
first-order $L_{T,i} = \sum_t \ell_{t,i}$	$R_T = O(\sqrt{L_{T,i^*} \log N})$	$R_T = \tilde{O}(\sqrt{NL_{T,i^*}})$
second-order $S_{T,i} = \sum_t \ell_{t,i}^2$	$R_T = O(\sqrt{S_{T,i^*} \log N})$ ★ Cesa-Bianchi, Mansour, Stoltz (2005)	$R_T = \tilde{O}(\sqrt{\sum_i S_{t,i}})$ Auer et al. (2002) + some hacking
variance $V_{T,i} = \sum_t (\ell_{t,i} - m)^2$	$R_T = O(\sqrt{V_{T,i^*} \log N})$ ★ Hazan and Kale (2010)	$R_T = \tilde{O}(N^2 \sqrt{\sum_i V_{t,i}})$ Hazan and Kale (2011)

★ with a little cheating

Beyond minimax: "Higher-order" bounds

	Full information	Bandit
minimax	$R_T = \Theta(\sqrt{T \log N})$	$R_T = O(\sqrt{NT})$
first-order $L_{T,i} = \sum_t \ell_{t,i}$	$R_T = O(\sqrt{L_{T,i^*} \log N})$	$R_T = \tilde{O}(\sqrt{NL_{T,i^*}})$
second-order $S_{T,i} = \sum_t \ell_{t,i}^2$	<div data-bbox="517 801 1251 1136" style="background-color: yellow; border: 2px solid black; padding: 10px; display: inline-block;"> <p>What about these?</p> </div>	$R_T = \tilde{O}(\sqrt{\sum_i S_{t,i}})$ Auer et al. (2002) + some hacking
variance $V_{T,i} = \sum_t (\ell_{t,i} - m)^2$	$R_T = O(\sqrt{V_{T,i^*} \log N})$ Hazan and Kale (2010)	$R_T = \tilde{O}(N^2 \sqrt{\sum_i V_{t,i}})$ Hazan and Kale (2011)

★ with a little cheating

Beyond first-order bounds?

A key tool for adaptive bounds in full-info: **PROD** (Cesa-Bianchi, Mansour and Stoltz, 2005)

$$p_{t,i} \propto \prod_{s=1}^{t-1} (1 - \eta \ell_{s,i})$$

Beyond first-order bounds?

A key tool for adaptive bounds in full-info: **PROD** (Cesa-Bianchi, Mansour and Stoltz, 2005)

$$p_{t,i} \propto \prod_{s=1}^{t-1} (1 - \eta \ell_{s,i})$$

Used for proving

- Second-order bounds (Cesa-Bianchi et al., 2005)
- Variance-dependent bounds (Hazan and Kale, 2010)
- Path-length bounds (Steinhardt and Liang, 2014)
- Quantile bounds (Koolen and Van Erven, 2015)
- Best-of-both-worlds bounds (Sani et al., 2014)
- ...

Beyond first-order bounds?

A key tool for adaptive bounds in full-info: **PROD** (Cesa-Bianchi, Mansour and Stoltz, 2005)

$$p_{t,i} \propto \prod_{s=1}^{t-1} (1 - \eta \ell_{s,i})$$

But does it work for bandits?

- ...
- ... (Cesa-Bianchi, Mansour and Stoltz, 2005)
- ... (Kale, 2010)
- ... (Kleinberg, 2014)
- ... (2015)
- ... Best of both worlds bounds (Sall et al., 2014)

Beyond first-order bounds?

A key tool for adaptive bounds in full-info: **PROD** (Cesa-Bianchi, Mansour and Stoltz, 2005)

$$p_{t,i} \propto \prod_{s=1}^{t-1} (1 - \eta \ell_{s,i})$$

But does it work for
bandits?

NO

- ...
- ... (Cesa-Bianchi, Mansour and Stoltz, 2005)
- ... (Kale, 2010)
- ... (Kleinberg, 2014)
- ... (2015)
- ... (Best of both worlds bounds (Sall et al., 2014))

Why does PROD fail?

$$p_{t,i} \propto \prod_{s=1}^{t-1} (1 - \eta \hat{\ell}_{s,i})$$

=

$$p_{t,i} \propto e^{-\eta \sum_{s=1}^{t-1} \tilde{\ell}_{s,i}}$$

with

$$\tilde{\ell}_{t,i} = -\frac{1}{\eta} \log(1 - \eta \hat{\ell}_{t,i})$$

Why does PROD fail?

EXP3 with a pessimistic estimate:

$$\mathbf{E}[\tilde{\ell}_{t,i}] \geq \ell_{t,i}$$

=

$$p_{t,i} \propto e^{-\eta \sum_{s=1}^{t-1} \tilde{\ell}_{s,i}}$$

with

$$\tilde{\ell}_{t,i} = -\frac{1}{\eta} \log(1 - \eta \hat{\ell}_{t,i})$$

Why does PROD fail?

EXP3 with a pessimistic estimate:

$$\mathbf{E}[\tilde{\ell}_{t,i}] \geq \ell_{t,i}$$

=

$$p_{t,i} \propto e^{-\eta \sum_{s=1}^{t-1} \tilde{\ell}_{s,i}}$$

with

$$\tilde{\ell}_{t,i} = -\frac{1}{\eta} \log(1 - \eta \hat{\ell}_{t,i})$$

Implicit exploration

$$\tilde{\ell}_{t,i} = \frac{\ell_{t,i}}{p_{t,i} + \gamma} \mathbf{1}_{\{I_t=i\}}$$

Why does PROD fail?

EXP3 with a pessimistic estimate:

$$\mathbf{E}[\tilde{\ell}_{t,i}] \geq \ell_{t,i}$$

=

$$p_{t,i} \propto e^{-\eta \sum_{s=1}^{t-1} \tilde{\ell}_{s,i}}$$

with

$$\tilde{\ell}_{t,i} = -\frac{1}{\eta} \log(1 - \eta \hat{\ell}_{t,i})$$

Implicit exploration

$$\tilde{\ell}_{t,i} = \frac{\ell_{t,i}}{p_{t,i} + \gamma} \mathbf{1}_{\{I_t=i\}}$$

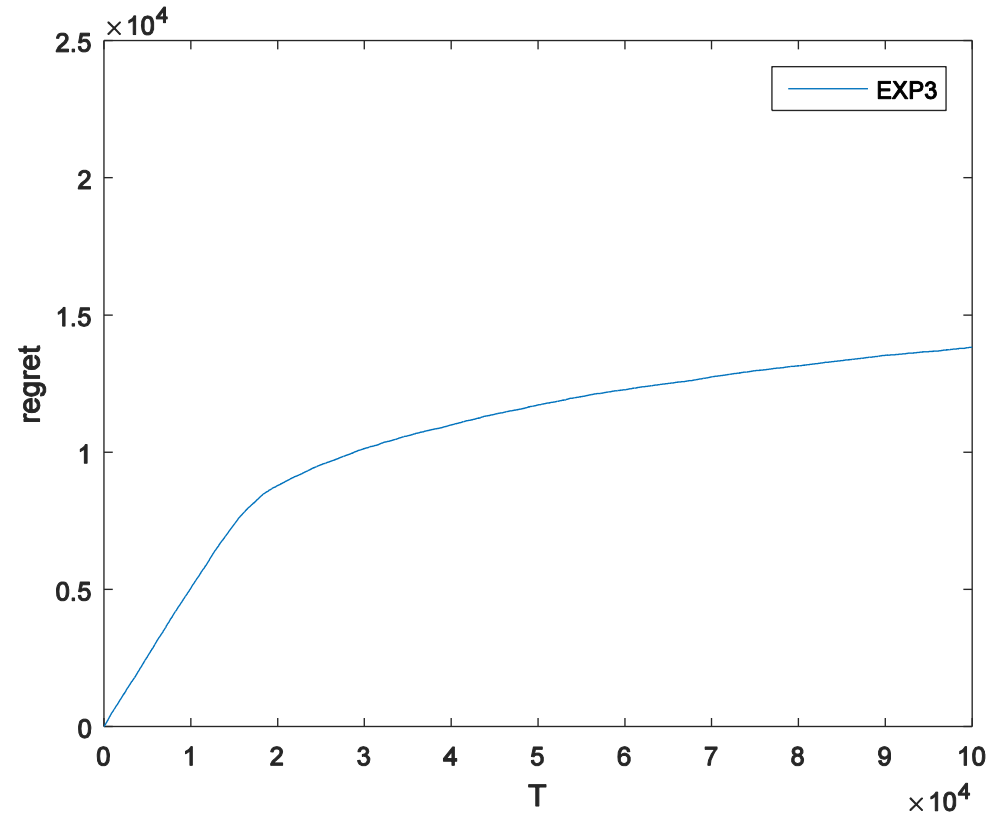
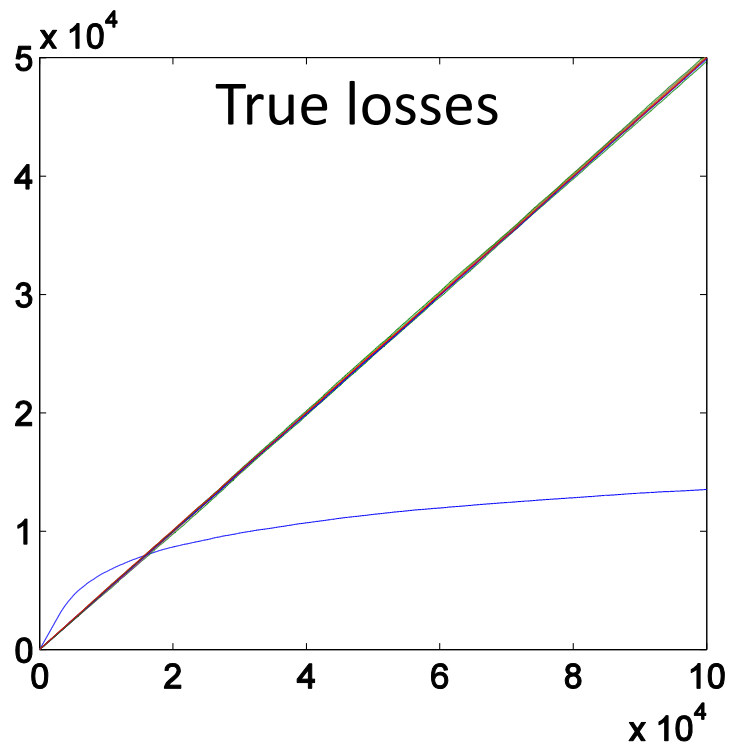
≈

$$p_{t,i} \propto e^{-\eta \sum_{s=1}^{t-1} \tilde{\ell}_{s,i}}$$

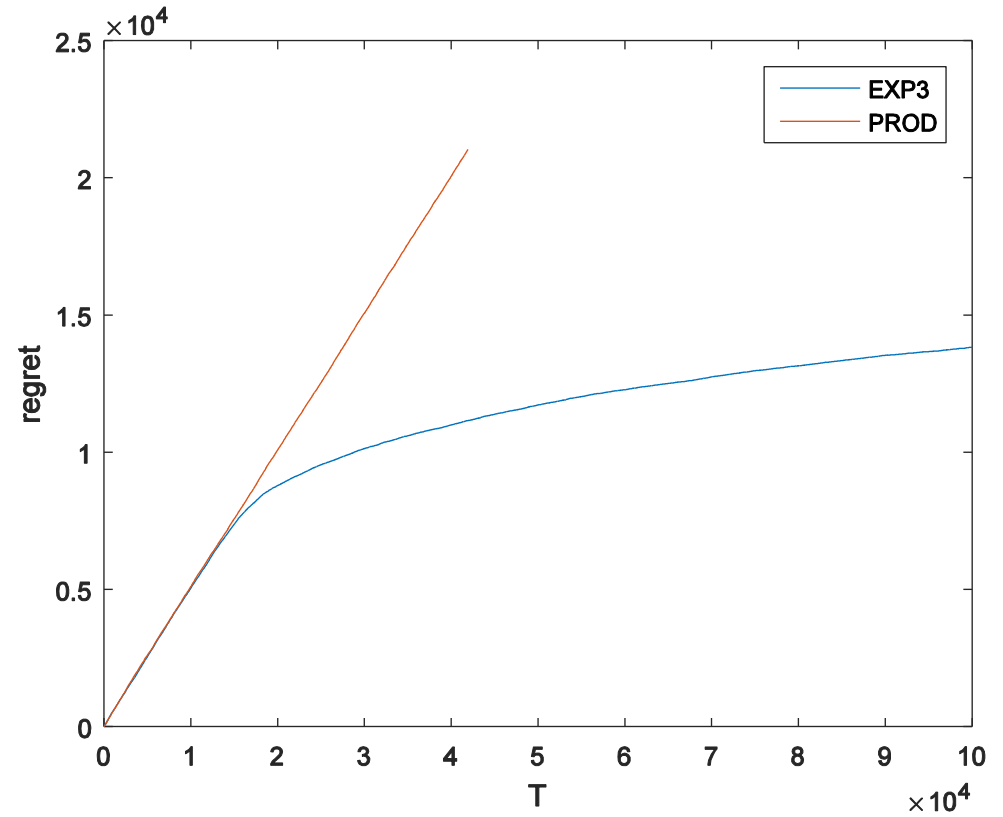
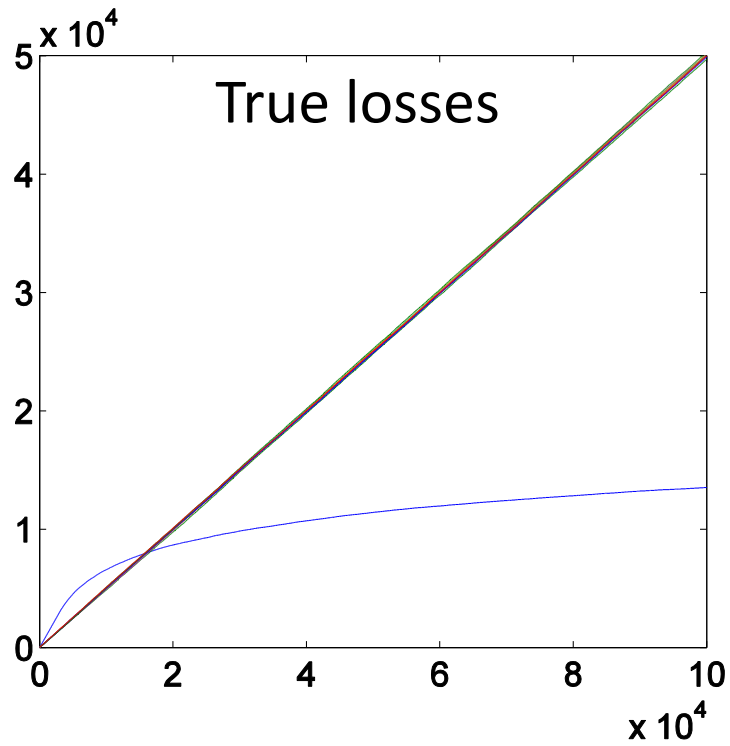
with

$$\tilde{\ell}_{t,i} = \frac{1}{\eta} \log(1 + \eta \hat{\ell}_{t,i})$$

PROD for bandits



PROD for bandits



Summary

- Key for first-order bounds: **implicit exploration**

$$\hat{\ell}_{t,i} = \frac{\ell_{t,i}}{p_{t,i} + \gamma} \mathbf{1}_{\{I_t=i\}}$$

+ also key for high-probability bounds!
(NIPS 2015)

Summary

- Key for first-order bounds: **implicit exploration**

$$\hat{\ell}_{t,i} = \frac{\ell_{t,i}}{p_{t,i} + \gamma} \mathbf{1}_{\{I_t=i\}}$$

+ also key for high-probability bounds!
(NIPS 2015)

- Further bounds seem to be difficult to prove: **smoothness** conflicts with **need to explore!**

Summary

- Key for first-order bounds: **implicit exploration**

$$\hat{\ell}_{t,i} = \frac{\ell_{t,i}}{p_{t,i} + \gamma} \mathbf{1}_{\{I_t=i\}}$$

+ also key for high-probability bounds!
(NIPS 2015)

- Further bounds seem to be difficult to prove: **smoothness** conflicts with **need to explore!**
- More depressing results by Lattimore (NIPS 2015)

Summary

- Key for first-order bounds: **implicit exploration**

$$\hat{\ell}_{t,i} = \frac{\ell_{t,i}}{p_{t,i} + \gamma} \mathbf{1}_{\{I_t=i\}}$$

+ also key for high-probability bounds!
(NIPS 2015)

- Further bounds seem to be difficult to prove: **smoothness** conflicts with **need to explore!**
- More depressing results by Lattimore (NIPS 2015)
 - Second-order bounds (Cesa-Bianchi et al., 2005)
 - Variance-dependent bounds (Hazan and Kale, 2010)
 - Path-length bounds (Steinhardt and Liang, 2014)
 - Quantile bounds (Koolen and Van Erven, 2015)

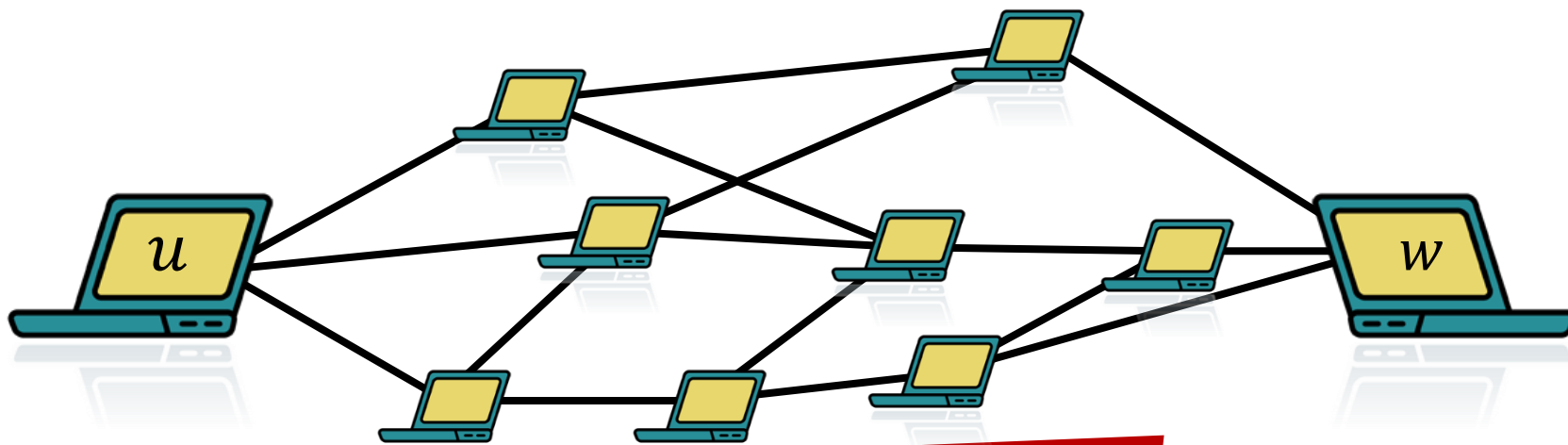


Thanks!

Appendix

First-order bounds for combinatorial semi-bandits

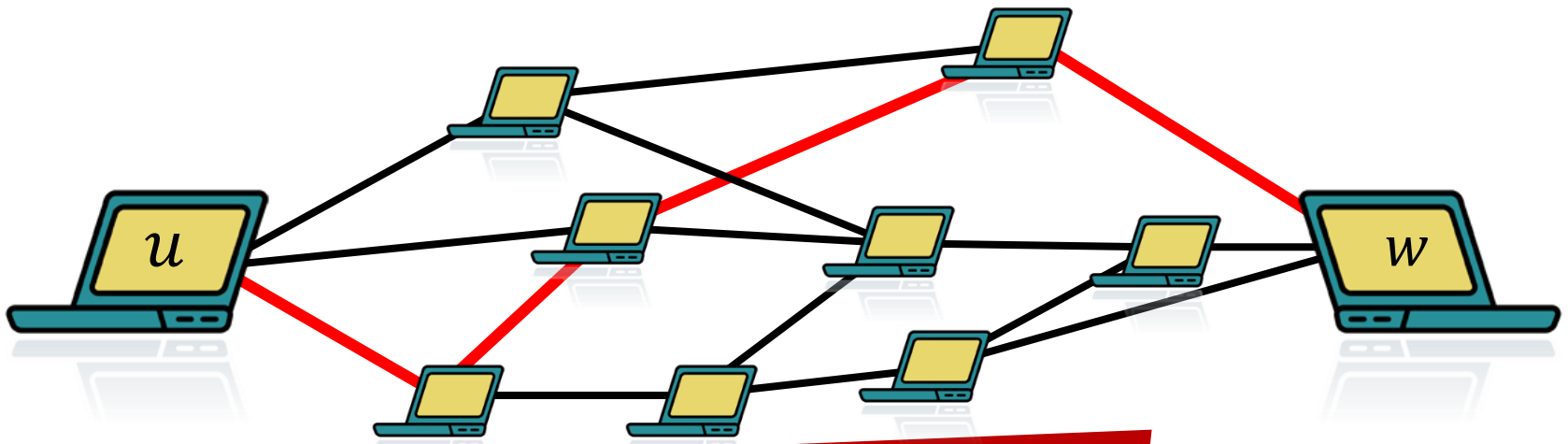
Combinatorial semi-bandits



For every round $t = 1, 2, \dots, T$

- learner picks an **action** $V_t \in S \subseteq \{0, 1\}^d$
- Environment chooses **loss vector** $\ell_t \in [0, 1]^d$
- Learner suffers loss $V_t^\top \ell_t$
- Learner observes **losses** $V_{t,i} \ell_{t,i}$

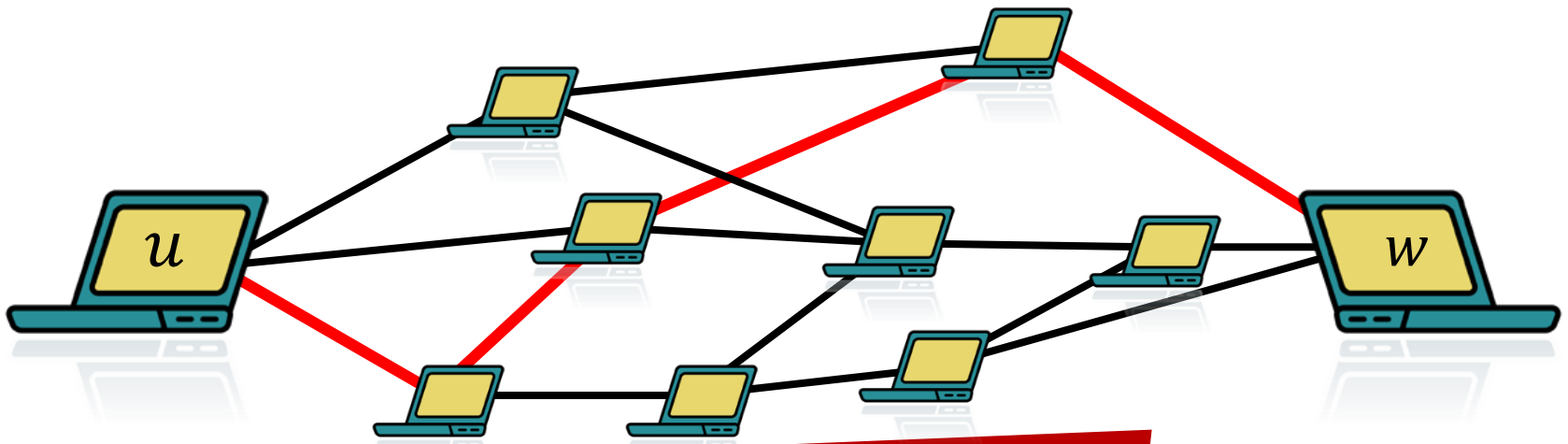
Combinatorial semi-bandits



For every round $t = 1, 2, \dots, T$

- learner picks an **action** $V_t \in S \subseteq \{0, 1\}^d$
- Environment chooses **loss vector** $\ell_t \in [0, 1]^d$
- Learner suffers loss $V_t^\top \ell_t$
- Learner observes **losses** $V_{t,i} \ell_{t,i}$

Combinatorial semi-bandits



For every round $t = 1, 2, \dots, T$

- learner picks an **action** $V_t \in S \subseteq \{0, 1\}^d$
- Environment chooses **loss vector** $\ell_t \in [0, 1]^d$
- Learner suffers loss $V_t^\top \ell_t$
- Learner observes **losses** $V_{t,i} \ell_{t,i}$

Decision set:

$$S = \{v_i\}_{i=1}^N \subseteq \{0, 1\}^d$$
$$\|v_i\|_1 \leq m$$

Combinatorial semi-bandits

- Goal: minimize (expected) **regret**

$$\hat{R}_T = \max_{v \in S} \mathbf{E} \left[\sum_{t=1}^T (V_t - v)^\top \ell_t \right]$$

- Minimax regret is

$$\hat{R}_T = \Theta(\sqrt{mdT})$$

- Best efficient algorithm (FPL) gives

$$\hat{R}_T = O(m\sqrt{dT \log(d)})$$

Combinatorial semi-bandits

- Goal: minimize (expected) **regret**

$$\hat{R}_T = \max_{v \in S} \mathbf{E} \left[\sum_{t=1}^T (V_t - v)^\top \ell_t \right]$$

- Minimax regret is

$$\hat{R}_T = \Theta(\sqrt{mdT})$$

- Best efficient algorithm (FPL) gives

$$\hat{R}_T = O(m\sqrt{dT \log(d)})$$

- Our bound:

$$\hat{R}_T = O(m\sqrt{dL_T^* \log(d)})$$