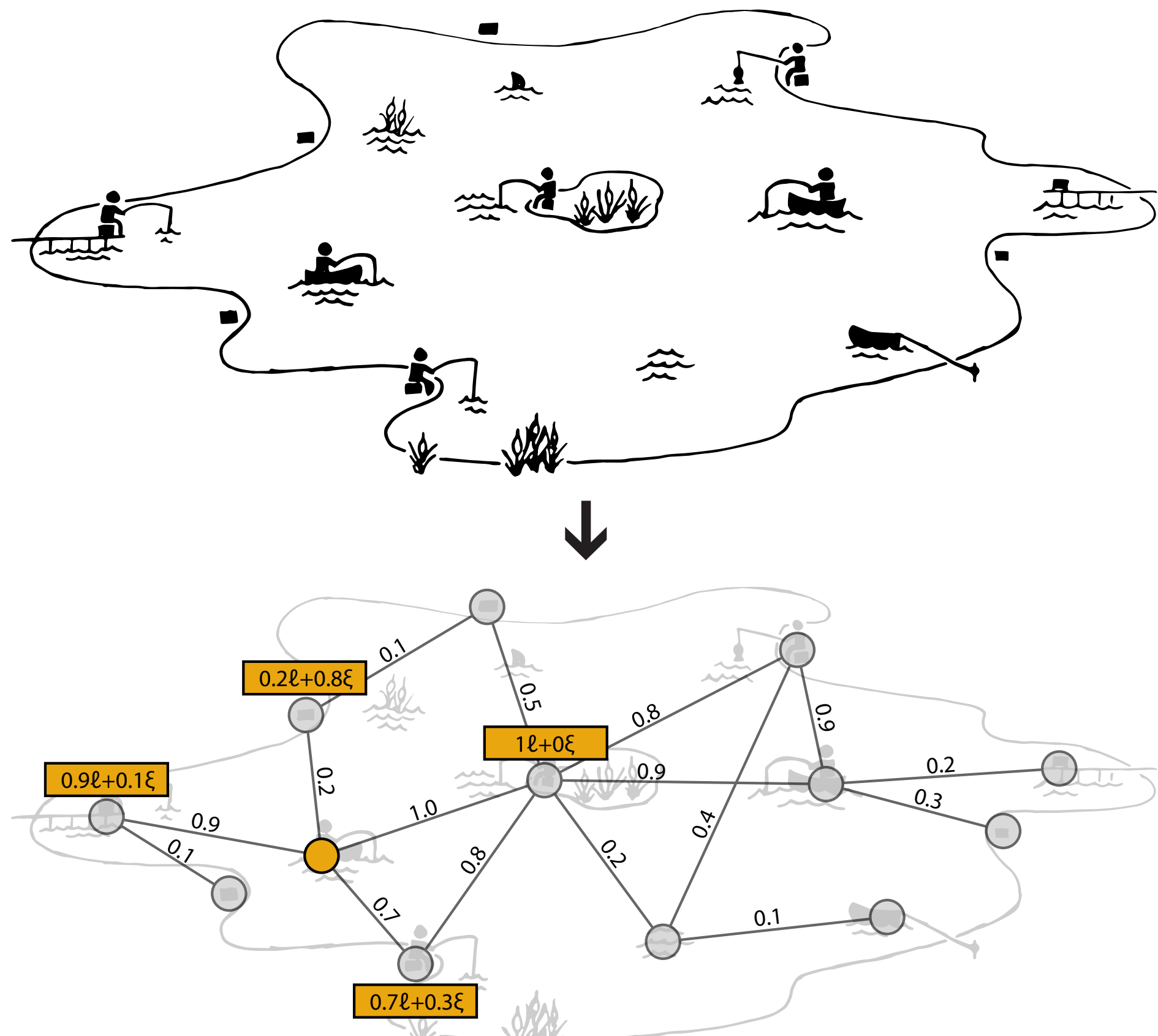


ONLINE LEARNING WITH NOISY SIDE OBSERVATIONS

TOMAS.KOČAK@INRIA.FR, NEU.GERGELY@GMAIL.COM, AND MICHAL.VALKO@INRIA.FR



MOTIVATION - FISHING SPOTS



Setup

- Pick a fishing spot (at the beginning of the day)
- Obtain some amount of fish (at the end of the day)
- Observe other fishermen (noisy side observations)
 - Does not represent what you would get exactly
 - Different fishing lure, motivation, skill...
- Goal: catch as many fish as possible over the time

Side observations - prior work (Mannor and Shamir 2011)

- Playing an action reveals the losses of neighboring spots

Noisy side observations - generalization

- Side observations represented by a **weighted** graph
- Playing action reveals **noisy** the losses of neighboring spots

PROBLEM FORMALIZATION

Learning process

- N actions (nodes of a graph)
- T rounds:
 - Environment (adversary) sets losses for actions
 - Environment chooses a graph (not disclosed)
 - Learner picks an action I_t to play
 - Learner incurs the loss ℓ_{t,I_t} of the action I_t
 - Learner observes graph (second neighborhood of I_t)
 - Learner observes noisy losses $c_{t,j}$ of neighbors $j \in N(I_t)$

where

$$c_{t,j} = s_{t,(I_t,j)} \ell_{t,j} + (1 - s_{t,(I_t,j)}) \xi_{t,j}$$

- $s_{t,(I_t,j)}$: weight of the edge from node I_t to node j at time t
- $\ell_{t,j}$: loss of the action j at time t
- $\xi_{t,j}$: zero mean noise such that $|\xi_{t,j}| \leq R$

Goal of the learner: minimizing cumulative regret R_t defined as

$$R_t = \underbrace{\sum_{t=1}^T \ell_{t,I_t}}_{\text{learner}} - \underbrace{\min_{j \in [N]} \sum_{t=1}^T \ell_{t,j}}_{\text{best action}}$$

EXP3-TYPE ALGORITHM TEMPLATE

- Compute exponential weights using loss estimates $\hat{\ell}_{s,i}$

$$w_{t,i} = \exp \left(-\eta_t \sum_{s=1}^{t-1} \hat{\ell}_{s,i} \right)$$

- Create a probability distribution such that $p_{t,i} \propto w_{t,i}$
- Play action I_t such that

$$\mathbb{P}(I_t = i) = p_{t,i} = \frac{w_{t,i}}{\sum_{j=1}^N w_{t,j}}$$

- Create loss estimates (using observability graph)
 - The definition of $\hat{\ell}$ defines an EXP3-type algorithm

LOSS ESTIMATES OF ALGORITHMS

Desired property of loss estimates: $\mathbb{E}[\hat{\ell}_{t,i}] \approx \ell_{t,i}$

Typical estimates: $\hat{\ell}_{t,i} = \frac{c_{t,i} \mathbb{1}\{\text{arm } i \text{ is observed}\}}{\mathbb{P}\{\text{arm } i \text{ is observed}\}}$

Graphs without weights (known algorithms)

EXP3 (edgeless graph)

$$\hat{\ell}_{t,i} = \frac{\ell_{t,i}}{p_{t,i}} \mathbb{1}\{\text{arm } i \text{ is observed}\}$$

Hedge (full graph)

$$\hat{\ell}_{t,i} = \frac{\ell_{t,i}}{1} \mathbb{1}\{\text{arm } i \text{ is observed}\}$$

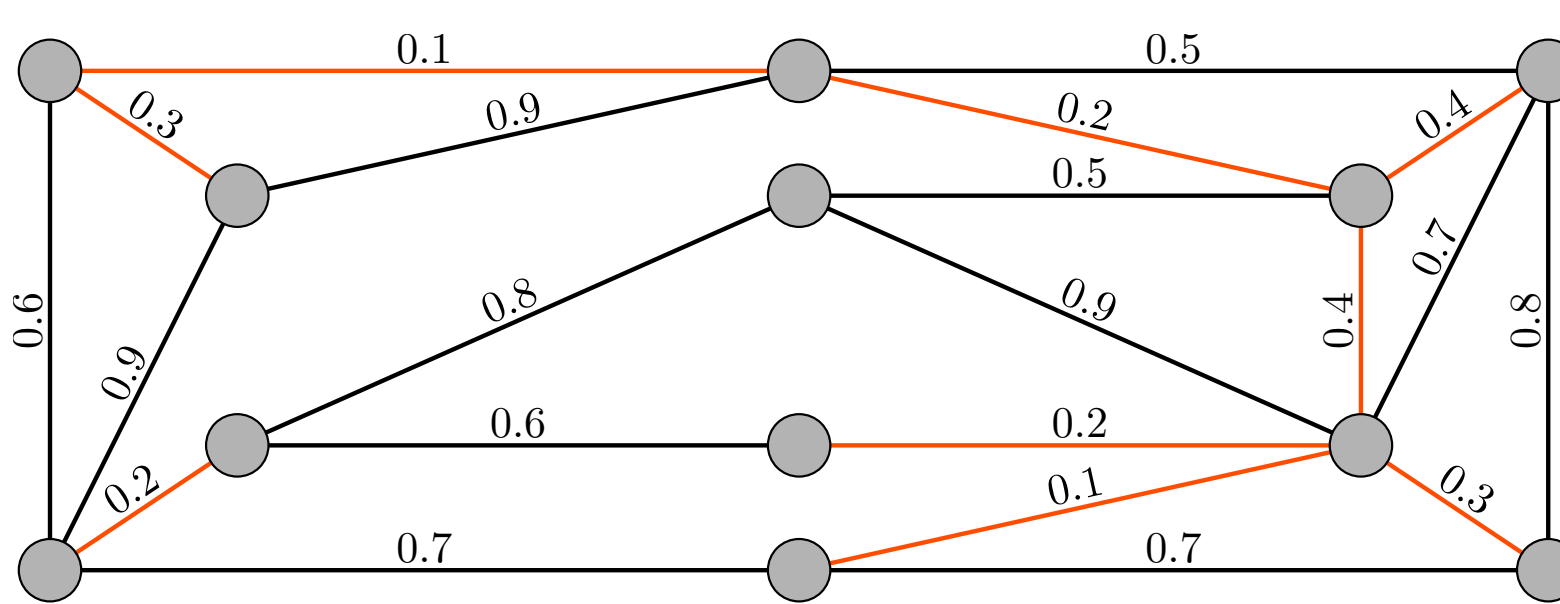
EXP3-IX (general graph)

$$\hat{\ell}_{t,i} = \frac{\ell_{t,i}}{\sum_{j \in N(j)} p_{t,j} + \gamma_t} \mathbb{1}\{\text{arm } i \text{ is observed}\}$$

Graphs with weights (new algorithms)

EXP3-IXT (general graphs with weights and thresholding)

- Use only "reliable" observations with low noise
- Delete all the edges with weights smaller than ε



$$\hat{\ell}_{t,i} = \frac{s_{t,(I_t,i)} \ell_{t,i} + (1 - s_{t,(I_t,j)}) \xi_{t,i}}{\sum_{j=1}^N s_{t,(j,i)} p_{t,j} + \gamma_t}$$

EXP3-WIX (general graph with weights)

$$\hat{\ell}_{t,i} = \frac{s_{t,(I_t,i)} [s_{t,(I_t,i)} \ell_{t,i} + (1 - s_{t,(I_t,j)}) \xi_{t,i}]}{\sum_{j=1}^N s_{t,(j,i)}^2 p_{t,j} + \gamma_t}$$

THEORETICAL GUARANTIES

EXP3-IXT

- Regret upper-bound

$$\mathbb{E}[R_t] = \tilde{O} \left(\sqrt{\frac{\alpha(\varepsilon)}{\varepsilon^2} T} \right)$$

- - Needs to know the "best" ε to optimize regret bound
- - Needs to know whole graph to find the best EXP3

EXP3-WIX (main result)

- Regret upper-bound

$$\mathbb{E}[R_t] = \tilde{O} \left(\sqrt{\min_{\varepsilon \in [0,1]} \frac{\alpha(\varepsilon)}{\varepsilon^2} T} \right)$$

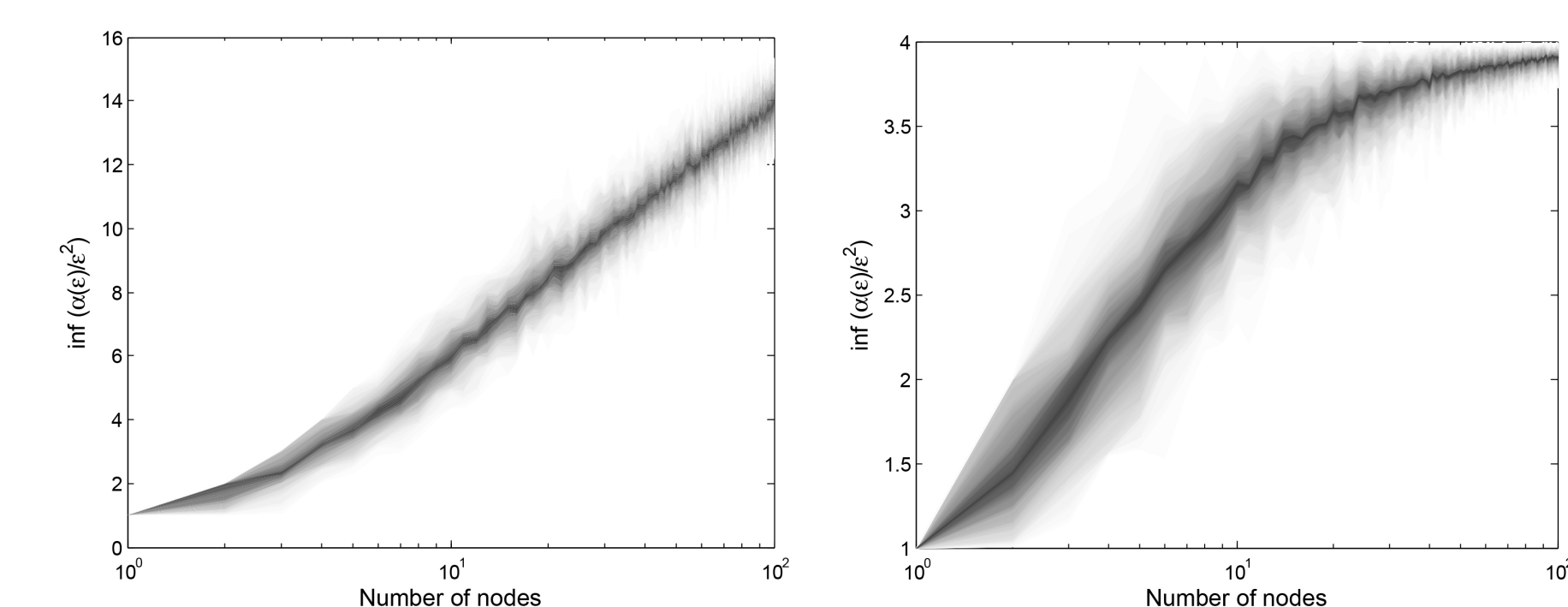
- + Does not need to set any threshold ε
- + Does not need to compute any independence numbers α
- Regret bound as good as the best bound of EXP3-IXT

EFFECTIVE INDEPENDENCE NUMBER

Effective independence number α^* is

$$\alpha^* = \min_{\varepsilon \in [0,1]} \frac{\alpha(\varepsilon)}{\varepsilon^2}$$

The empirical value of α^* uniformly random weights

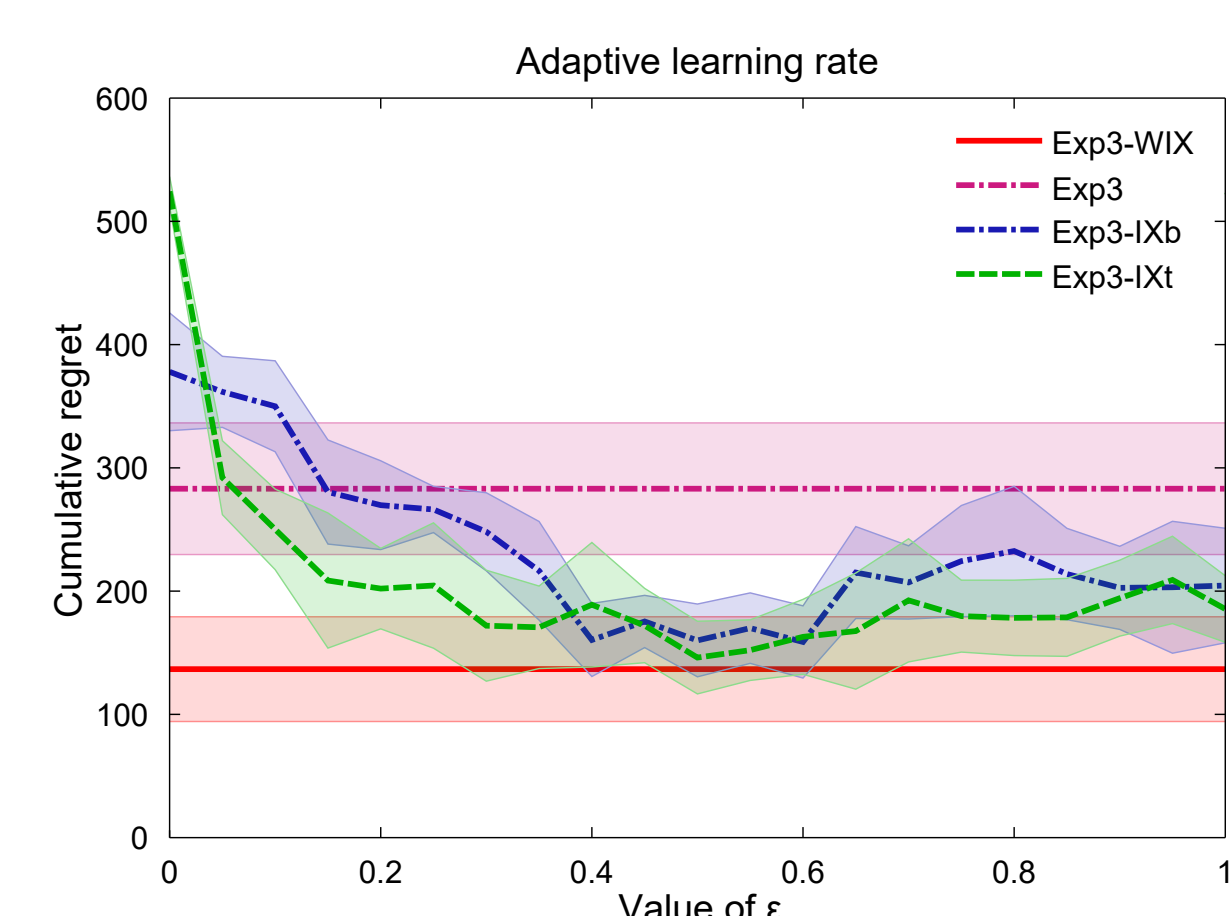


$U(0,1)$ weights

$U(\frac{1}{2},1)$ weights

EMPIRICAL RESULTS

- EXP3 - basic algorithm which ignores all side observations
- EXP3-WIX - our proposed algorithm
- EXP3-IXT - thresholded algorithm (needs to set ε)
- EXP3-IXB - algorithm ignores noise (no guaranties)



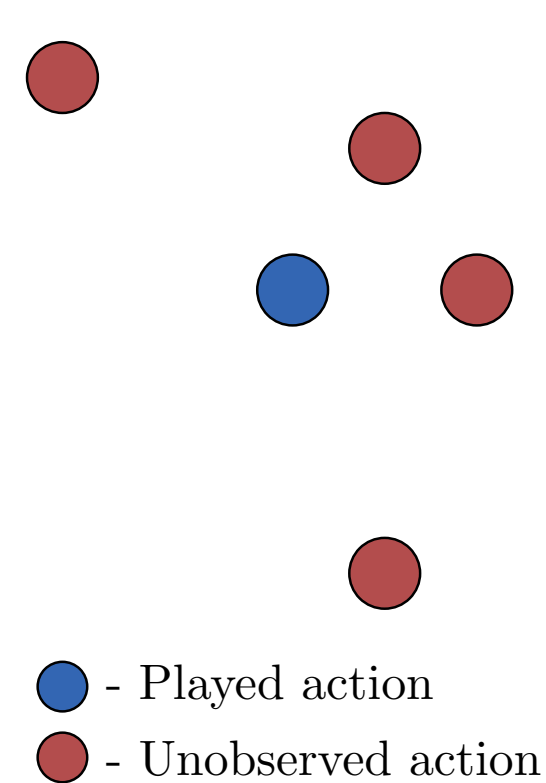
CONCLUSION

- New setting with noisy side observations
- Introduction of effective independence number α^*
- EXP3-WIX algorithm for the setting
 - Does not need to threshold
 - Does not need to know whole graph
 - Regret bound of order $\sqrt{\alpha^* T}$
- Open questions:
 - Is the effective independence number "right quantity?"
 - Is there a matching lower-bound for EXP3-WIX?
 - Upper-bound of EXP3-WIX matches lower-bound for some cases (e.g., bandits, full information, and setting of Mannor and Shamir 2011)
 - Related lower-bound (Wu et al. 2015) for a stochastic setting with Gaussian noise

$$R_t = O \left(\sqrt{\frac{\alpha}{\varepsilon^2} T} \right)$$

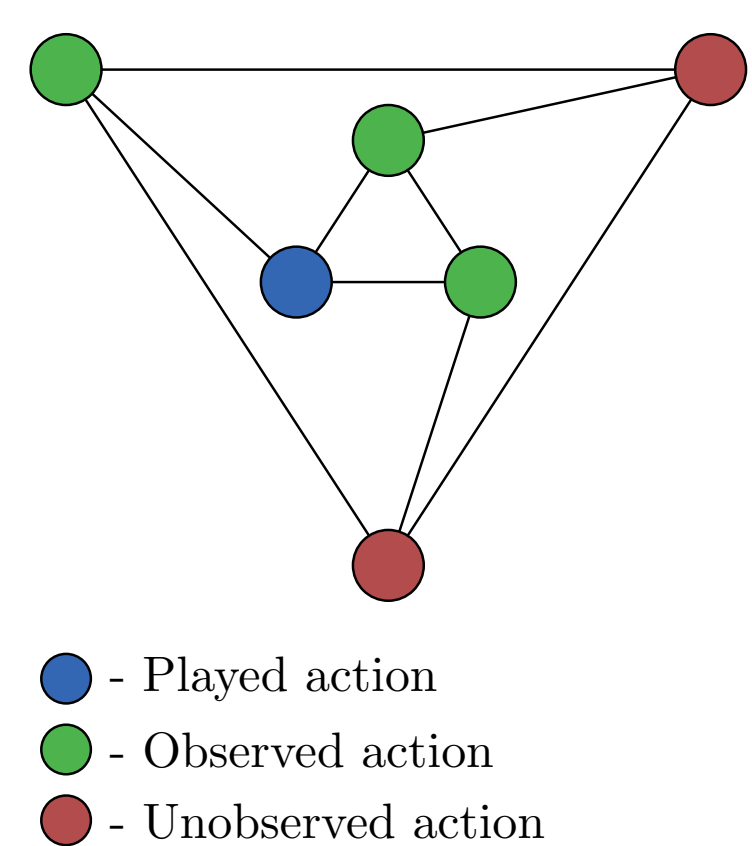
SPECIAL CASES OF THE FRAMEWORK

Edgeless graphs - Bandit setting



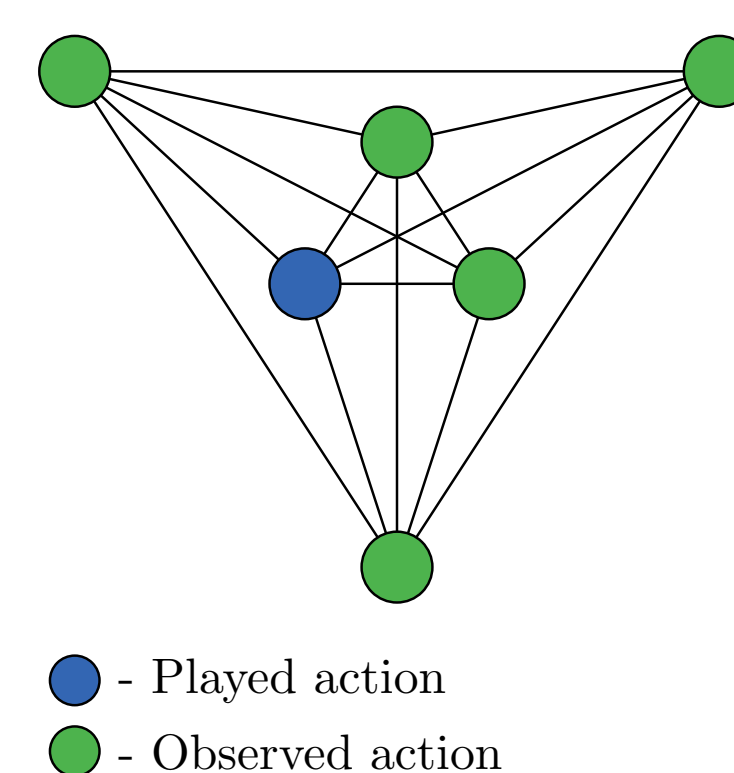
- No side observations
- Regret bound of order

General graphs - Mannor and Shamir 2011



- Some losses observed (without noise)
- Regret bound of order

Complete graphs - Full information setting



- All losses observed (without noise)
- Regret bound of order

$$\sqrt{NT}$$

\geq

$$\sqrt{\alpha T}$$

\geq

$$\sqrt{T}$$