

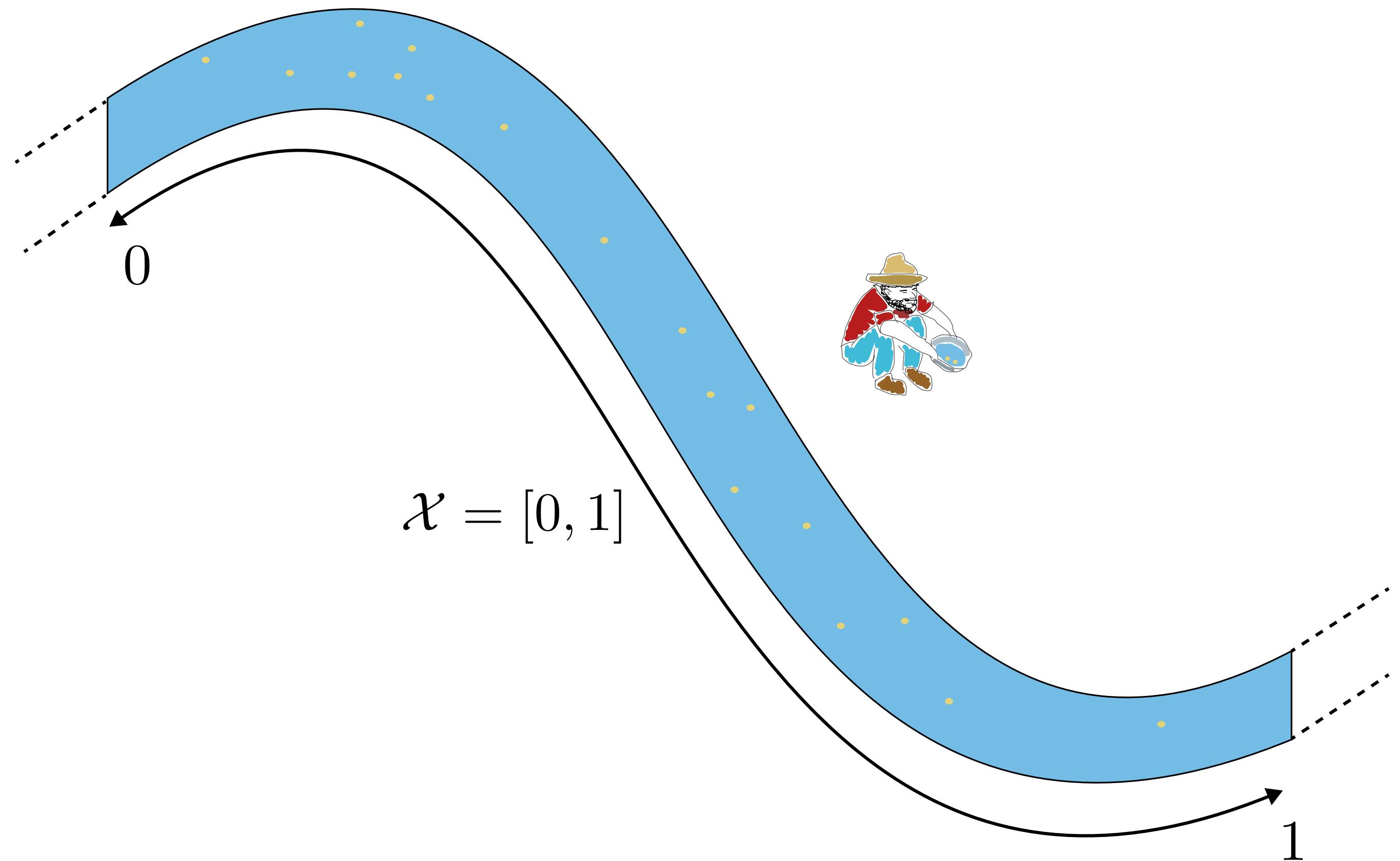
On some adaptivity questions in multi-armed bandits

Thèse de doctorat (virtually) at Orsay

Hédi Hadiji supervised by **Pascal Massart** and **Gilles Stoltz**

4th December 2020

Gold panning in a river



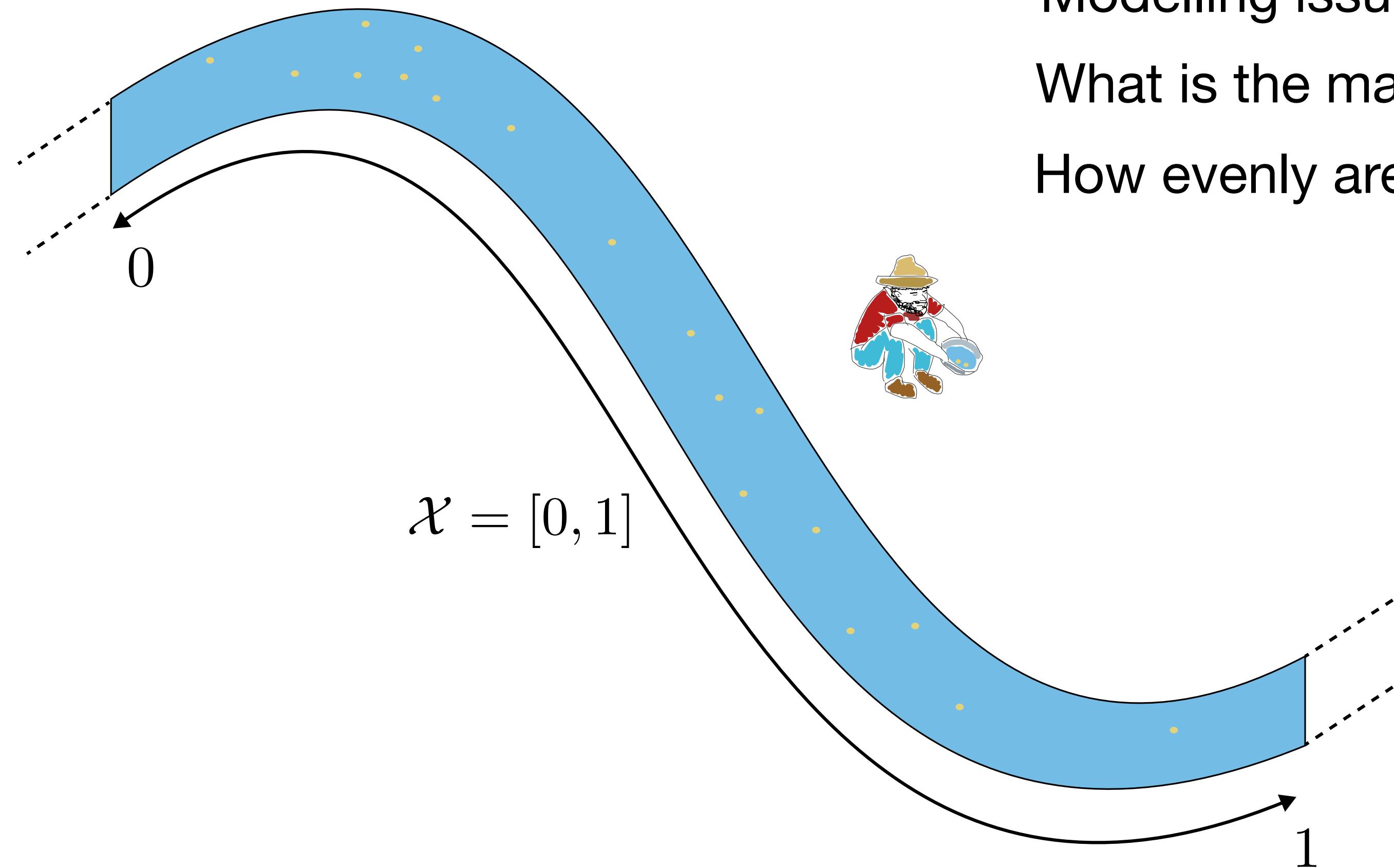
Every day $t = 1, \dots,$

- Pick a spot $X_t \in [0, 1]$ on the river
- Get Y_t grams of gold

What strategy to
get as much gold as possible?

Making as few assumptions
as necessary

Gold panning in a river



Modelling issues:

What is the maximum reward one can get in a single round?

How evenly are the rewards distributed across space?

Should we just guess?

This will crucially affect our decisions!

what is the maximum reward I can get in a single round?



how evenly are the rewards distributed across space?



How does prior information about the rewards affect optimal strategies?

Outline

- I/ Multi-armed bandits
- II/ Adapting to the unknown range of the rewards
- III/ Continuous bandits and smoothness

Multi-armed bandits

[Thompson 1933] [Robbins 1952]

K probability distributions (ν_1, \dots, ν_K) unknown to the player

$$\mu_i = \mathbb{E}(\nu_i)$$

At every time-step t the player:

- chooses an action A_t
- receives and observes reward $Y_t \sim \nu_{A_t}$ given A_t

Minimize **regret**

$$R_T = T\mu^* - \mathbb{E}\left[\sum_{t=1}^T Y_t\right] = T\mu^* - \mathbb{E}\left[\sum_{t=1}^T \mu_{A_t}\right]$$

Prior knowledge = assumption on ν_1, \dots, ν_K

Upper-Confidence Bounds

[Auer, Cesa-Bianchi, Fischer, '02]

ν_1, \dots, ν_K are supported in $[0, 1] \Leftrightarrow$ all rewards are in $[0, 1]$

Draw every arm once

Then draw arm A_t maximizing the index

$$U_a(t) = \hat{\mu}_a(t) + \sqrt{\frac{2 \log t}{N_a(t)}}$$

Empirical mean of rewards from arm a

Number of times a was picked

Distribution-free

$$R_T \leq c\sqrt{KT \ln T} + K$$

Distribution-dependent

$$R_T \leq \sum_{a=1}^K \frac{8}{\mu^* - \mu_a} \log T + \sum_{a=1}^K (\mu^* - \mu_a)$$

Improving UCB: two types of optimality

UCB explores a bit too much: we can reduce the confidence bounds and get better guarantees

MOSS

[Audibert, Bubeck '09]
[Degenne, Perchet '16]

$$R_T \leq c\sqrt{KT \ln T} + K$$

KL-UCB [Cappé, Garivier, Maillard, Munos, Stoltz '13]

$$R_T \leq \sum_{a=1}^K \frac{\mu^* - \mu_a}{\mathcal{K}_{\inf}(\nu_a, \mu^*)} \log T + \text{cst}$$

$$\mathcal{K}_{\inf}(\nu_a, \mu^*) = \inf \left\{ \text{KL}(\nu_a, \nu') \mid \mathbb{E}(\nu') > \mu^* \text{ and } \nu'([0, 1]) = 1 \right\}$$

None are further improvable [Auer, Cesa-Bianchi, Freund, Schapire '02], [Lai & Robbins 1985]

[Garivier, Ménard '17] obtain joint optimality in a parametric setting

Theorem: KL-UCB-switch [Garivier, H, Ménard, Stoltz '18]

KL-UCB-switch is both distribution-dependent and distribution-free optimal

What if the range is unknown?

rewards are in $[m, M]$ instead of $[0, 1]$, and m and M are unknown to the player?

Unknown range: initial remarks

1. $\frac{Y_t - m}{M - m} \in [0, 1]$

playing with rescaled rewards is equivalent
to standard game with regret scaled by $(M-m)$
... but we cannot rescale
2. (KL-)UCB, MOSS, etc. are scale-dependent
$$\hat{\mu}_a(t) + (M - m) \sqrt{\frac{\log t}{N_a(t)}}$$
3. Not knowing the range is harder than knowing the range
$$\sup_{\text{Problems in } [m, M]} R_T \geq c(M - m) \sqrt{KT}$$

Can we match this lower bound without knowing m and M ?

Unknown range: distribution-free adaptation

AdaHedge is from de [Rooij, van Erven, Grünwald, Koolen '13]
[Cesa-Bianchi, Mansour, Stoltz '07]

Theorem: Minimax range adaptation (H, Stoltz 2020)

AdaHedge for bandits with extra-exploration, guarantees

$$\text{for all } m < M, \quad \sup_{\substack{\text{Prob in } [m, M]}} R_T \leq 7(M-m)\sqrt{TK \ln K} + 10(M-m)K \ln K$$

Same guarantees as if we had known the range in advance!

But can we get distribution-dependent $\log T$ bounds?

Obstacle: extra-exploration forces at least $R_T \geq (M - m)\sqrt{KT}$

Adaptive rates

Distribution-free rate

For all T , for all $M > 0$,

$$\sup_{\text{Problems in } [0, M]} R_T \leq M \Phi_{\text{free}}(T)$$

e.g.

Bandit AdaHedge

enjoys $\Phi_{\text{free}}(T) = \sqrt{TK \ln K}$

Distribution-dependent rate

For all $M > 0$, for all problems in $[0, M]$,

$$\limsup_{T \rightarrow \infty} \frac{R_T}{\Phi_{\text{dep}}(T)} < +\infty$$

e.g.

$$\tilde{U}_a(t) = \hat{\mu}_a(t) + (\log t) \sqrt{\frac{\log t}{N_a(t)}}$$

enjoys $\Phi_{\text{dep}}(T) = (\log T)^2$ [Lattimore '17]

Can we get $\log T$ and \sqrt{KT} simultaneously?

Lower bound for adaptation to the range

Theorem: Lower bound for range adaptation (H, Stoltz 2020)

For any algorithm enjoying rates $\Phi_{\text{free}}(T)$ and $\Phi_{\text{dep}}(T)$

$$\Phi_{\text{dep}}(T) \Phi_{\text{free}}(T) \gtrsim T$$

in particular, $\Phi_{\text{free}}(T) \leq \mathcal{O}(\sqrt{T}) \Rightarrow \Phi_{\text{dep}}(T) \geq \Omega(\sqrt{T})$

The cost of distribution-free adaptation is hidden in distribution-dependent rates!

We get T^α and $T^{1-\alpha}$ (for $1/2 < \alpha < 1$) by tuning the extra-exploration in Bandit AdaHedge

Adapting to the smoothness

Back to the river: continuous bandits

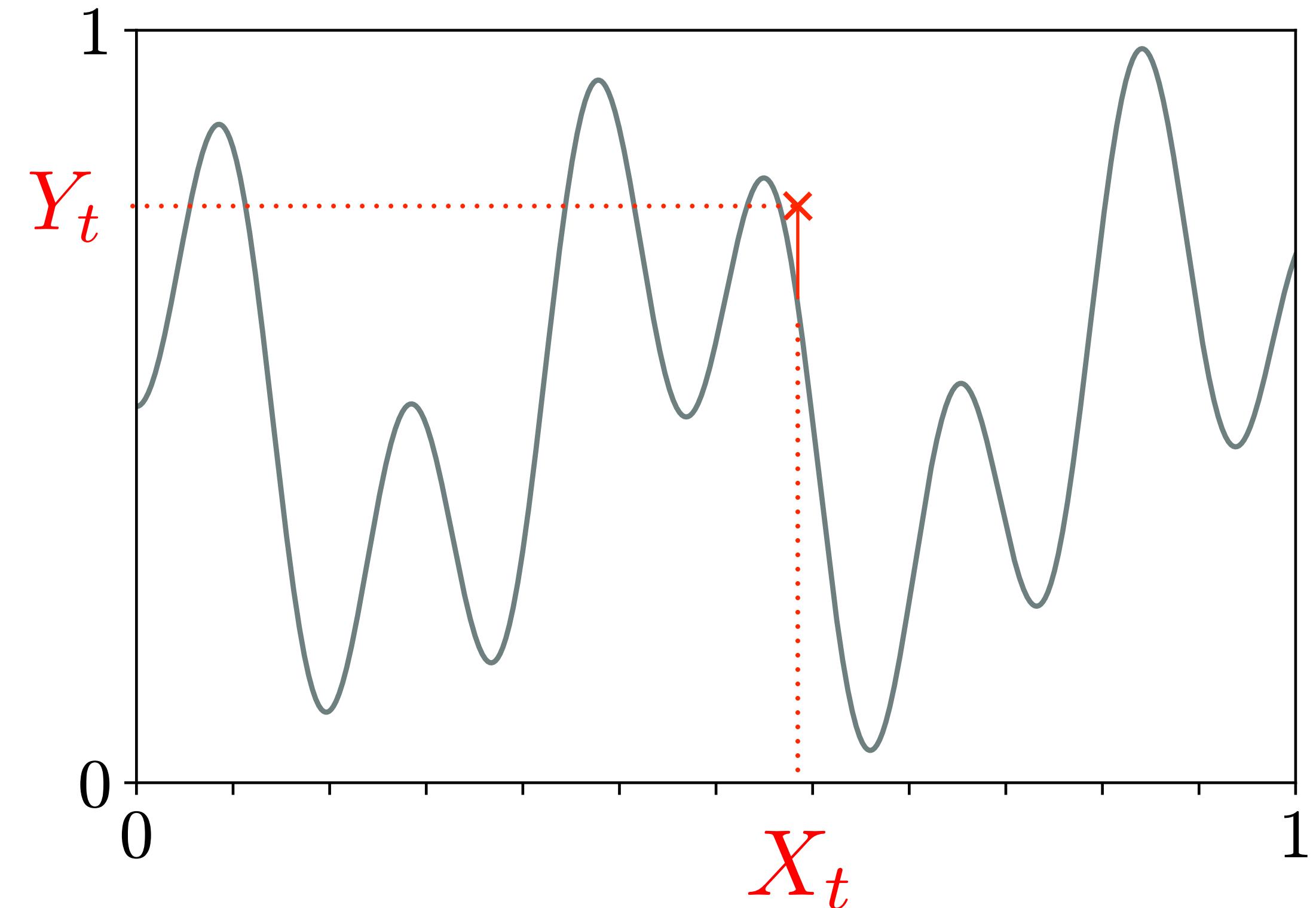
$\mathcal{X} = [0, 1]$ Unknown payoff distributions $(\nu_x, x \in \mathcal{X})$ Mean-payoff function $f : x \mapsto \mathbb{E}(\nu_x)$

For $t = 1, \dots, T, \dots$:

- pick $X_t \in \mathcal{X}$
- receive and observe $Y_t \sim \nu_{X_t}$ given X_t

Goal

$$R_T = T \max_{x \in \mathcal{X}} f(x) - \mathbb{E} \left[\sum_{t=1}^T f(X_t) \right]$$



We need assumptions

Assumptions

1. Rewards are bounded in $[0, 1]$
2. Mean-payoff function is “ (L, α) -Hölder around its max”

$$f(x^*) - f(x) \leq L |x^* - x|^\alpha$$

Lots of similar and/or more refined assumptions in the literature

[Agrawal 1995], [Kleinberg '05], [Auer, Ortner, Szepesvári, '07], [Bubeck, Munos, Stoltz, Szepesvári '11]
[Bubeck, Stoltz, Yu '11], [Kleinberg, Slivkins, Upfal, '13], [Bull '15]

Related problems: bandit optimization/simple regret, maximum estimation, maximum location

find \hat{X}
s.t. $f(\hat{X}) \approx \max f$

[Valko, Carpentier, Munos '13]
[Bartlett, Gabilon, Valko '19] [Shang, Kauffman, Valko '19]

estimate $\max f$

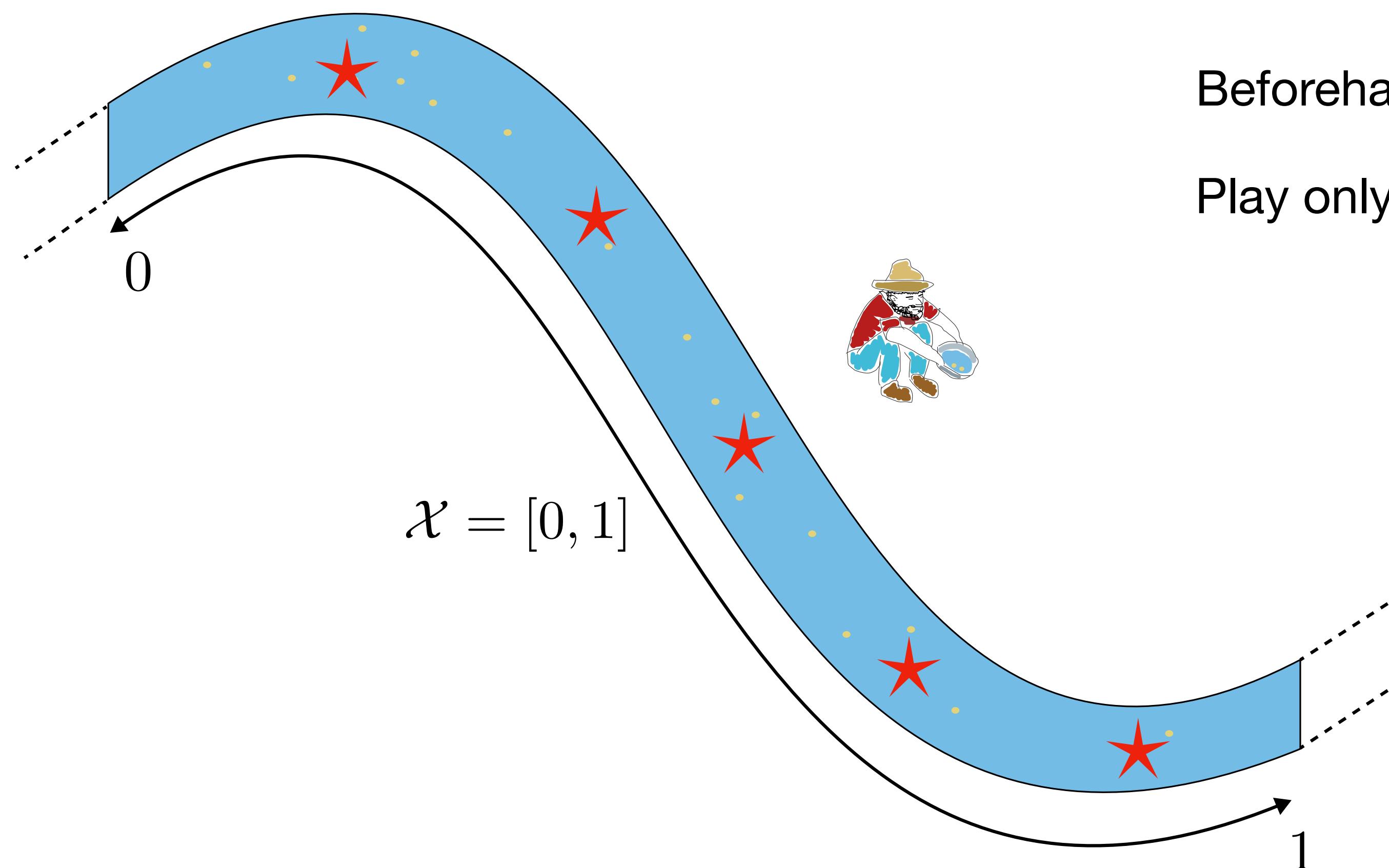
[Lepski 1994]

find \hat{X}
s.t. $\hat{X} \approx x^*$

[Muller 1989]

Discretization

[Kleinberg '05]

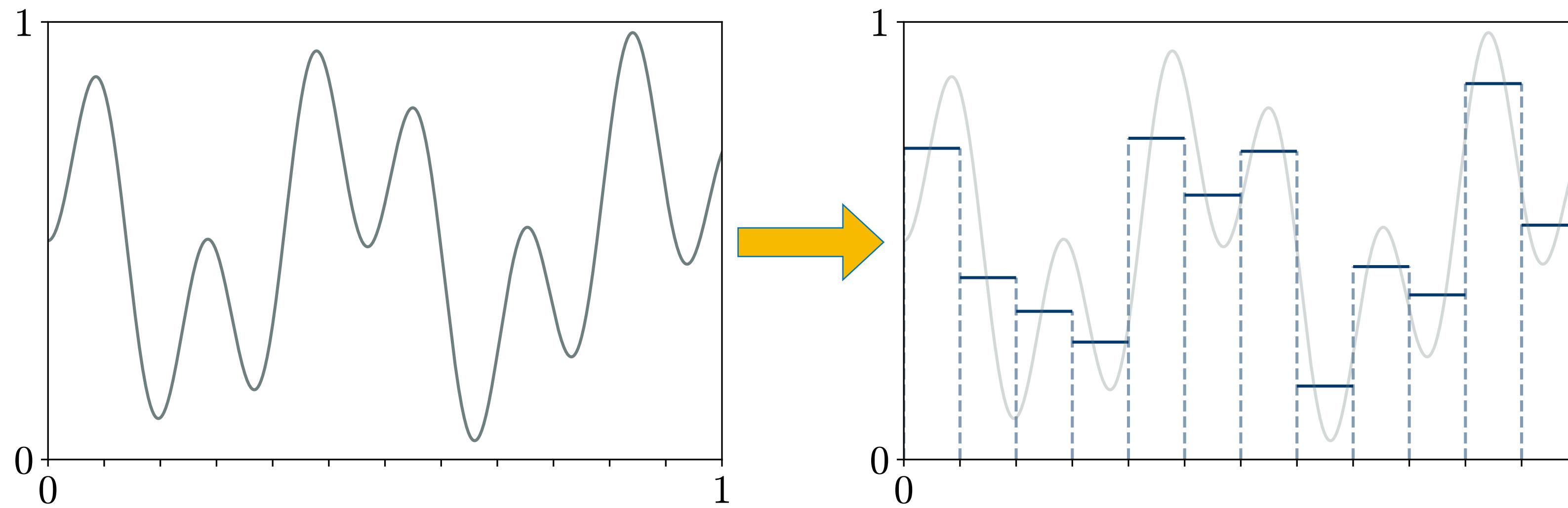


Beforehand, pick a finite number of spots

Play only at these spots, using a K-armed bandit algorithm

Back to a finite-armed bandit problem

Bounding the regret of discretization



$$R_T \leq c L^{1/(2\alpha+1)} T^{(\alpha+1)/(2\alpha+1)} \quad \text{by tuning } K = K^*(\alpha) \text{ appropriately}$$

Proof:

$$\begin{aligned} R_T &= T \left(\max f - \max_{1 \leq i \leq K} f(x_i) \right) + \max_{1 \leq i \leq K} f(x_i) - \mathbb{E} \left[\sum_{t=1}^T f(X_t) \right] \\ &\leq T \frac{L}{K^\alpha} + c\sqrt{KT} \end{aligned}$$

Impossibility of (full) adaptation

Adaptive rates

An algorithm achieves adaptive rates $\theta : \mathbb{R}_+ \rightarrow [1/2, 1]$ if

$$\sup_{\substack{f \\ \text{α-H\"older}}} R_T \leq c T^{\theta(\alpha)} \text{ for all } \alpha$$

Question: how can we obtain adaptive rates $\theta(\alpha) = \frac{\alpha + 1}{2\alpha + 1}$?

Model selection? Cross-validation?
... Exploration is costly

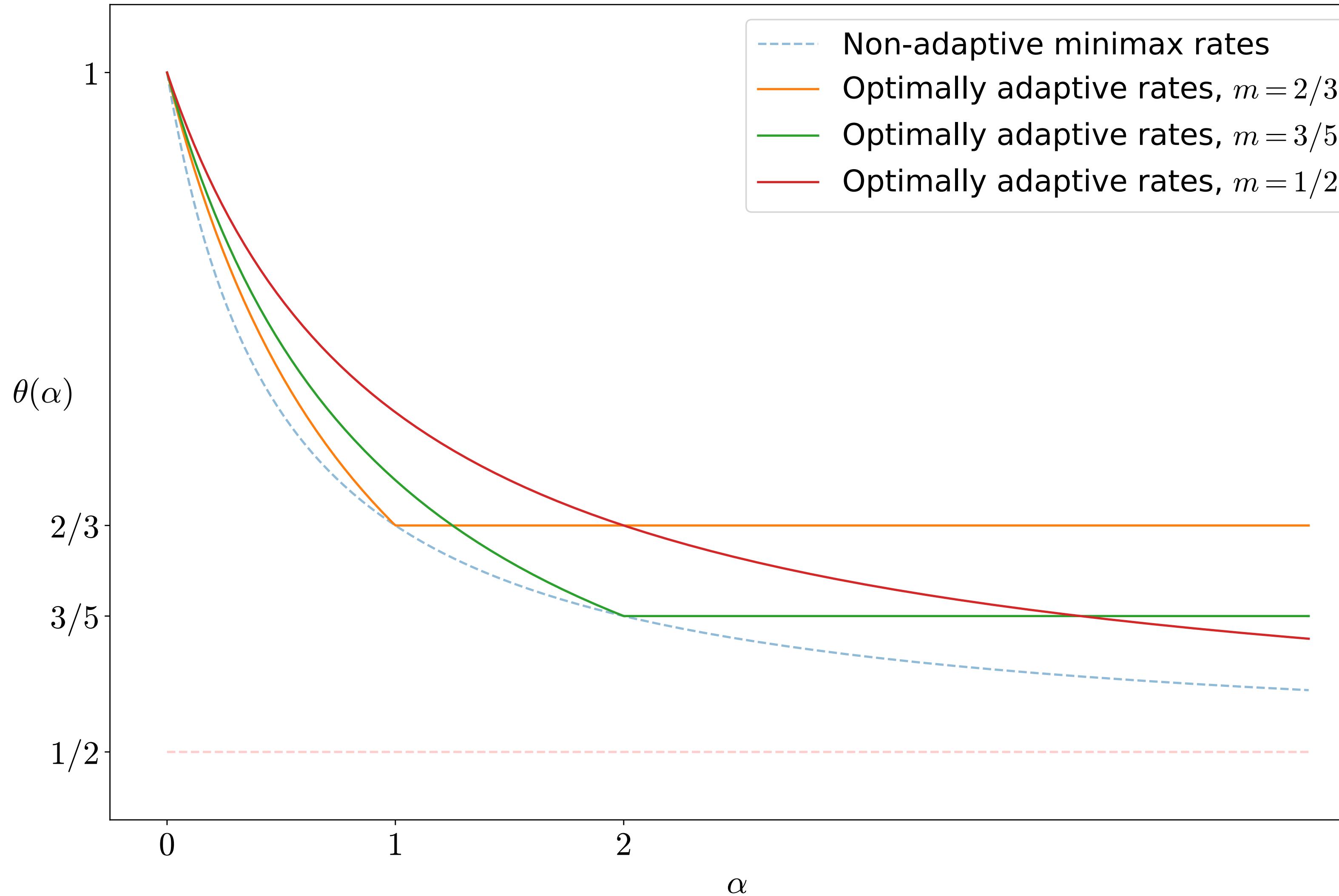
Answer: we cannot

Theorem: Consequence of [Locatelli and Carpentier '18]

If a rate function θ is achieved by some algorithm, then

$$\text{there exists } m \in [1/2, 1] \text{ s.t } \theta(\alpha) \geq \max \left(m, 1 - m \frac{\alpha}{\alpha + 1} \right)$$

Lower bound(s) on the adaptive rates



if $R_T \leq c T^{\theta(\alpha)}$
then θ is lower bounded
by one of these rates

Still, can we reach these rates?

A bird's eye view of past approaches

HOO [Bubeck, Munos, Stoltz, Szepesvári '11]

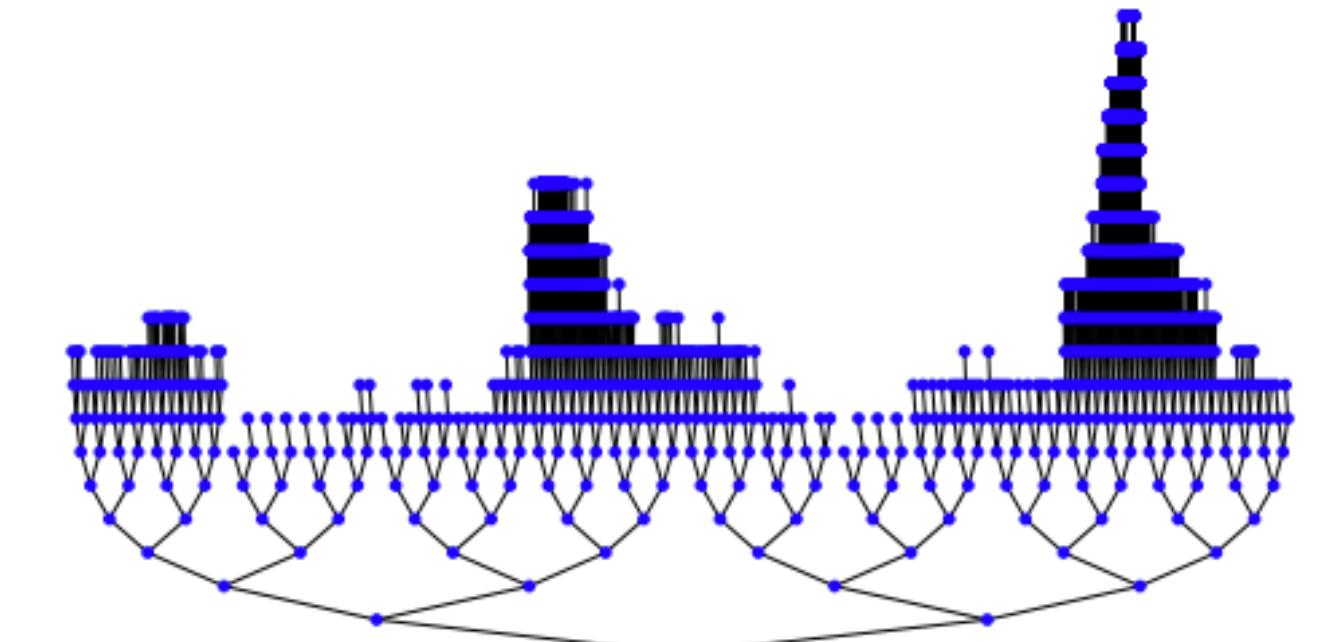
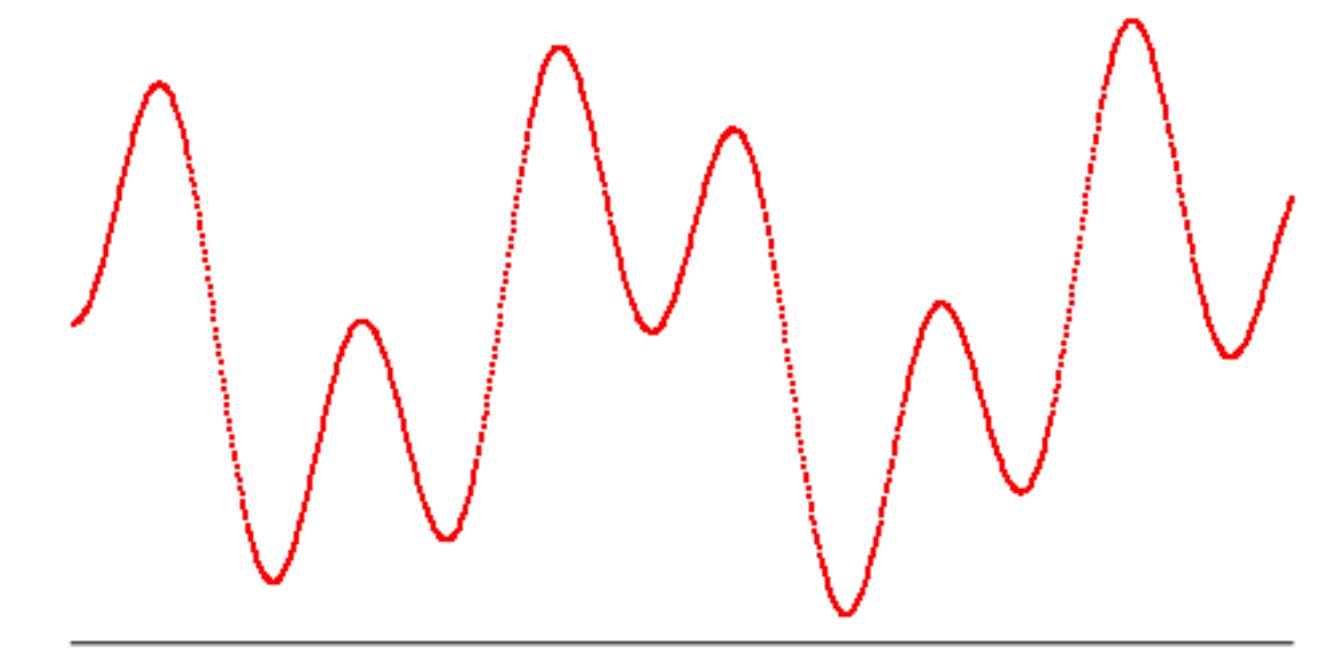
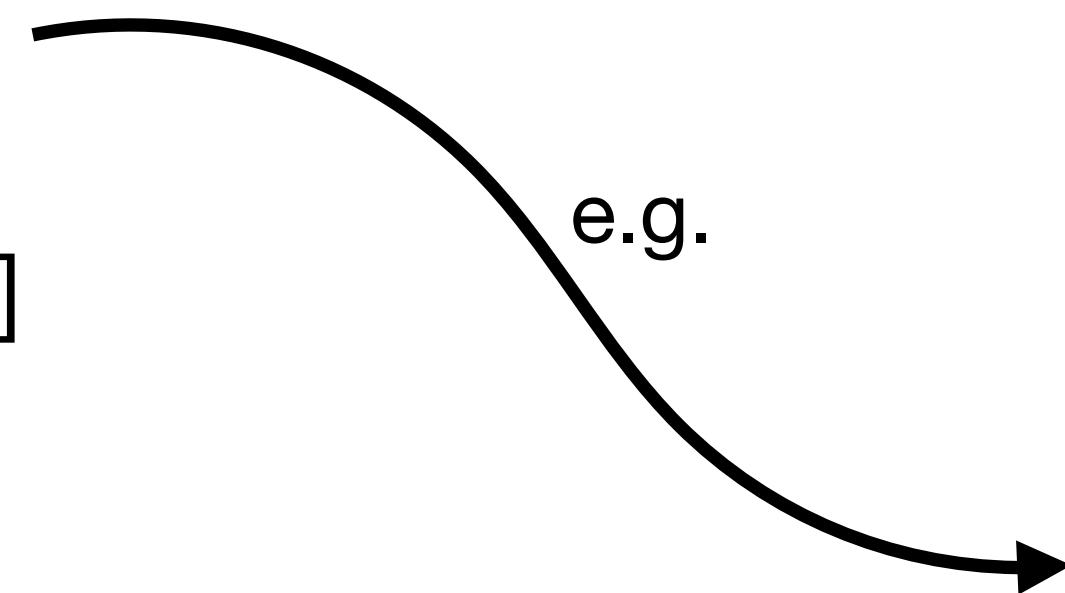
Zooming algorithm [Kleinberg, Upfal, Slivkins '13]

Adaptive-tree bandits [Bull '15]

SR [Locatelli, Carpentier '18]

Zoom in on promising regions

Select the promising regions **using the knowledge of the regularity**

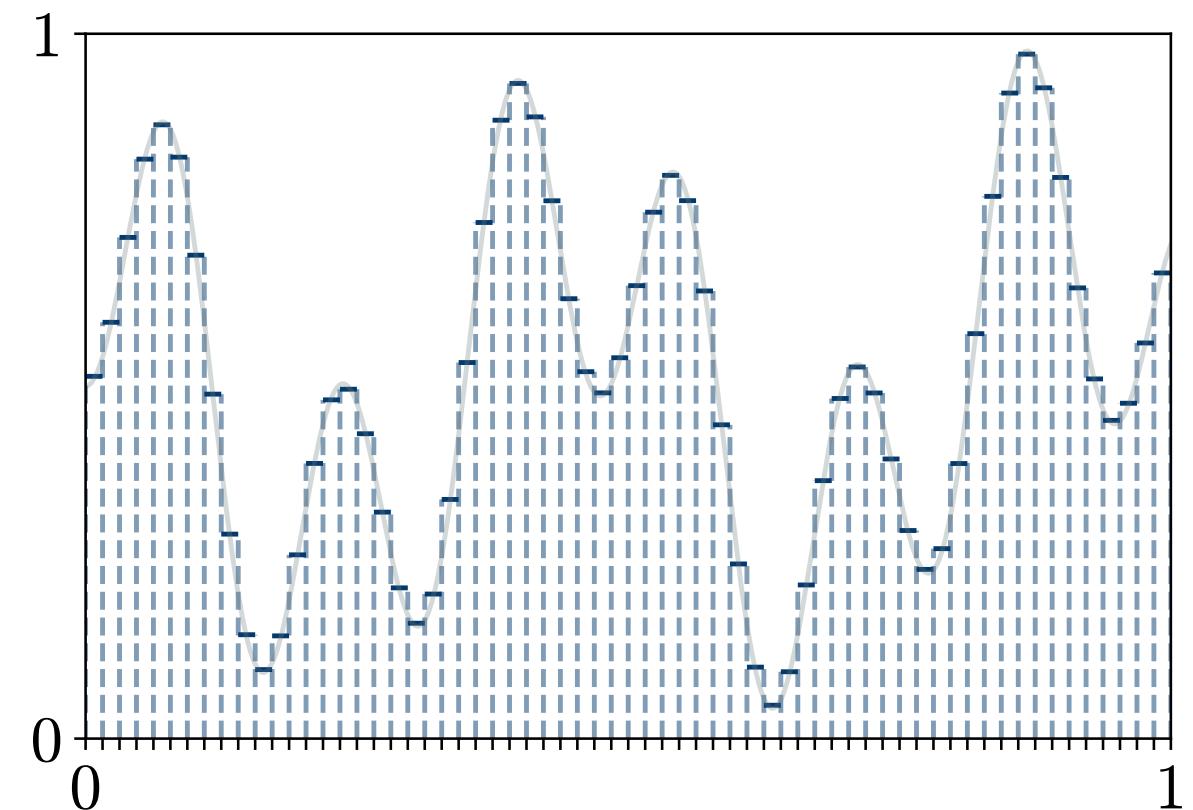


An optimally adaptive algorithm

Basic idea

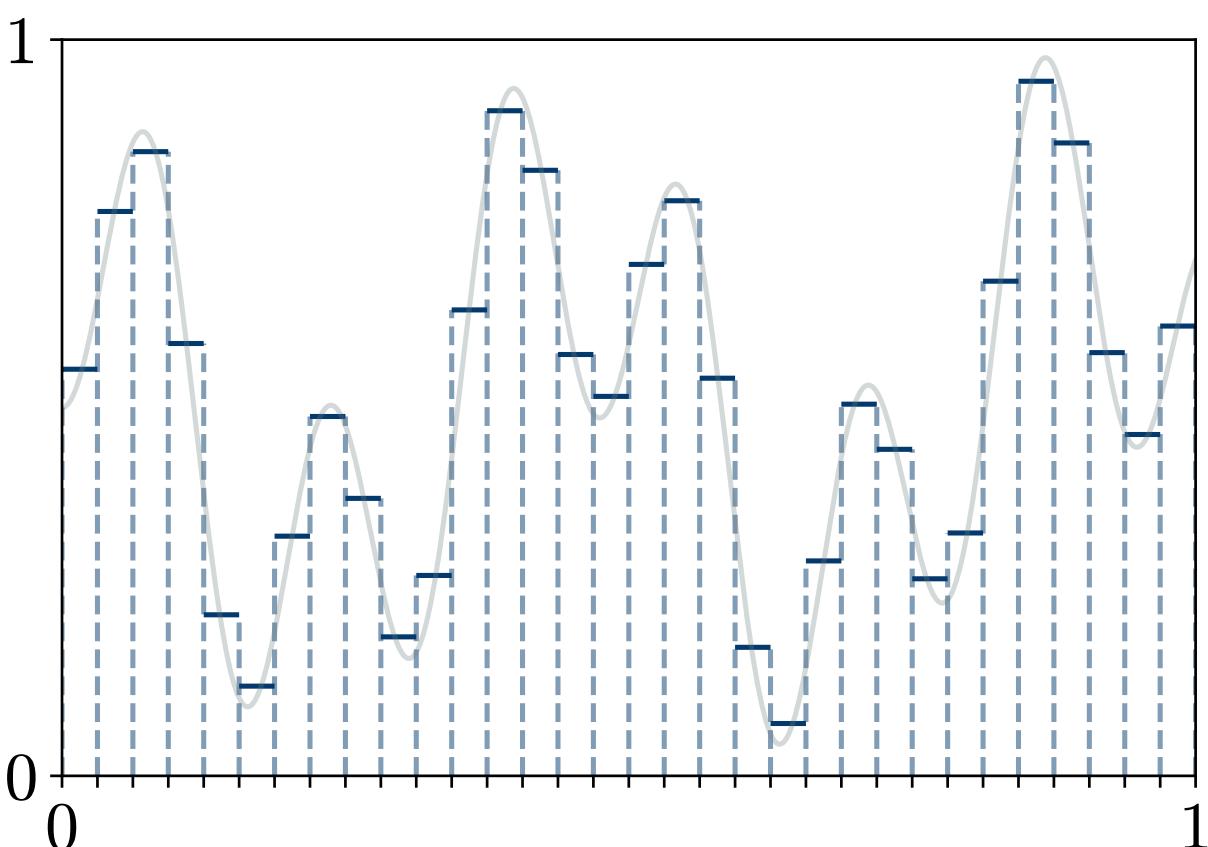
Assume the worst...

start with a very fine discretization
but not for too long



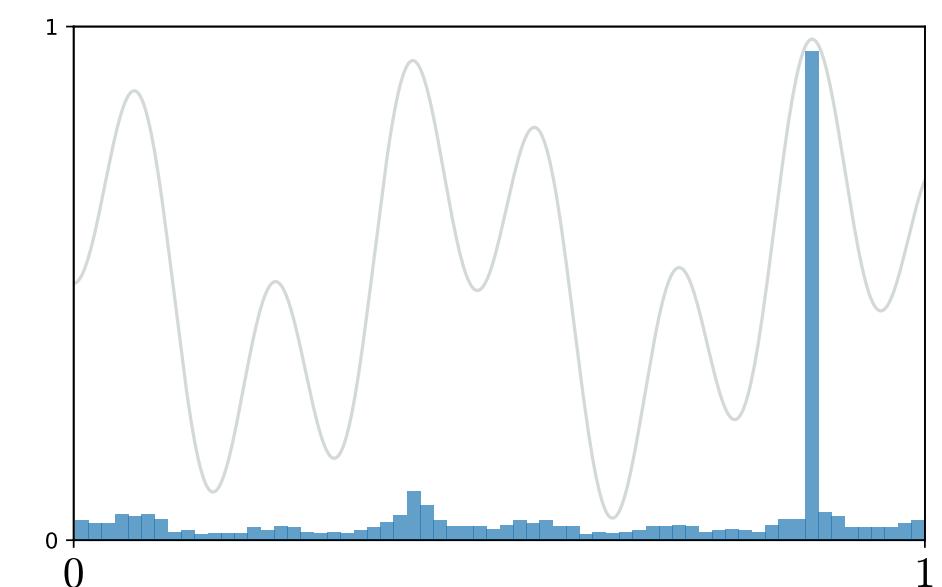
*...then **ZOOM OUT**.*

start over with a coarser discretization
but remember what you played before



Discrete algorithm choosing between: K actions from the discretization

An extra-action: playing at random
among actions selected in the past epoch



Guarantees

Algorithm : Memorize, Discretize, Zoom out

Set $K_i \approx 2^{-i}\sqrt{T}$; $D_i \approx 2^i\sqrt{T}$

For epochs 1 to $\approx \log T$

- For D_i rounds, run a K_i -discretization **with memory of previous plays**

Theorem: Medzo regret bound [H, 2019]

For any (L, α) -Hölder function, without the knowledge of L or α ,

$$R_T \leq \tilde{\mathcal{O}}(L^{1/(\alpha+1)}T^{(\alpha+2)/(2\alpha+2)})$$

with no assumption on α and L

Conclusion

Adaptation leads to interesting questions with surprising phenomena

Achieving adaptation requires new algorithmic ideas

Ultimately important in practice/getting rid of hyperparameter tuning

Bandits are fun!

Local perspectives

KL-UCB: Unifying notion of optimality in bandits?

Range:

- Minimax and optimal distribution-dependent rates for adapting to the lower end of the range
- Getting rid of the $\sqrt{\log K}$

Continuous bandits:

- Including other types of regularity
- Is there a principled look at Medzo that could be applied elsewhere?
- What do we *need* to know to be able adapt at the usual rates [Locatelli Carpentier '18]?

General perspectives

- What about non IID payoffs? [Zimmert and Seldin '19]
- Related problem ‘Model selection in contextual multi-armed bandits’
[Foster et al. ’19], [Chatterji et al ’19], [Foster et al. COLT open problem ’20]

Observe context $C_t \in \{1, \dots, S\}$ then choose $A_t \in \{1, \dots, K\}$

A policy is mapping from context to action $\pi : \mathcal{C} \rightarrow \mathcal{A}$

A model Π is a set of policies

Easy (Exp4): $R_T(\text{Best } \pi \in \Pi) \leq c\sqrt{KT \log |\Pi|}$

Difficult: getting the same result with a sequence of nested models

-Dream general approach:

Can we combine a family of bandit algorithms and obtain one that is as good as the best?

[Agrawal et al. ’17], [Pacchiano et al. ’20], results are very specific

General perspectives

- What about non IID payoffs? [Zimmert and Seldin '19]
- Related problem ‘Model selection in contextual multi-armed bandits’
[Foster et al. ’19], [Chatterji et al ’19], [Foster et al. COLT open problem ’20]

Observe context $C_t \in \{1, \dots, S\}$ then choose $A_t \in \{1, \dots, K\}$

A policy is mapping from context to action $\pi : \mathcal{C} \rightarrow \mathcal{A}$

A model Π is a set of policies

Easy (Exp4): $R_T(\text{Best } \pi \in \Pi) \leq c\sqrt{KT \log |\Pi|}$

Difficult: getting the same result with a sequence of nested models

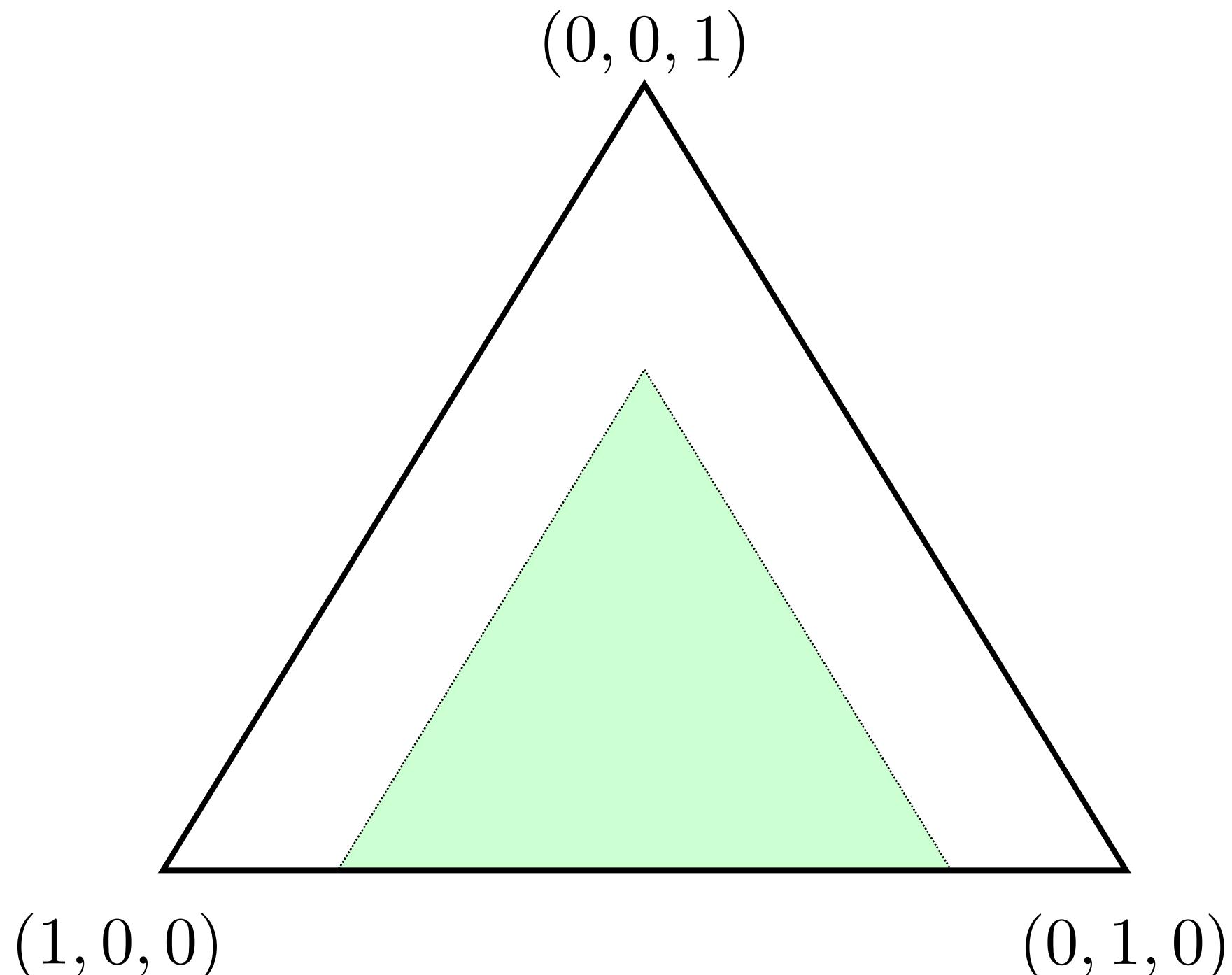
-Dream general approach:

Can we combine a family of bandit algorithms and obtain one that is as good as the best?

[Agrawal et al. ’17], [Pacchiano et al. ’20], results are very specific

Bonus: diversity-preserving bandits

with Sébastien Gerchinovitz, Jean-Michel Loubes and Gilles Stoltz



Play a probability distribution p_t over $\{1, 2, 3\}$

Observe $A_t \sim p_t$, and $Y_t \sim \nu_{A_t}$

Require that $p_{1,t}, p_{2,t} \geq 0.3$ (for example)

$$R_T = T \max_{p \text{ available}} \sum_{a=1}^K p_a \mu_a - \mathbb{E} \left[\sum_{t=1}^T Y_t \right]$$

Bounded regret is possible \Leftrightarrow the best p is in the (relative) interior of the simplex

Thank you Gilles and Pascal!

Thank you everyone!