

Stochastic & adversarial


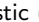



best-arm identification,

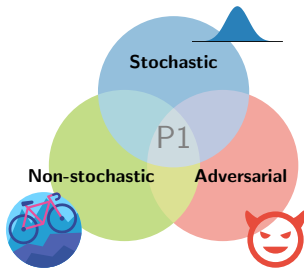
can we achieve the *best of both worlds*?

Victor Gabillon

joint work with Yasin Abbasi-Yadkori, Peter Bartlett,
Alan Malek & Michal Valko

Workshop on Quantifying Uncertainty: Stochastic, Adversarial, & Beyond,
Simons Institute, 13 September 2022

- **Setting:** A pure exploration bandit  problem
- **Question:** Can one algorithm achieve **BOB**: Perform well under data-generating regimes either stochastic () or non-stochastic (or even against an adversary ) ?
- **Contributions:**
 - a study of the  problem against 
 - an impossibility result on the **BOB** question
 - a simple algorithm P1 for **BOB** matching the lower bound.



After an **exploration phase** of T pulls, a Learner tries to **identify** the arm with highest cumulative reward out of K arms.

Bandit feedback: The learner only observes the reward/gain of the arm it chooses to explore.

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

- ▶ simultaneously, Learner picks arm $I_t \in [K]$, (K arms)
- ▶ **Adversary**  / **environment**  picks gain $g_t \in [0, g_{max}]^K$.
- ▶ Then, the Learner observes $g_{I_t, t}$.

Finally, Learner **recommends** an arm denoted $\mathbf{1}_T$, hoping $\mathbf{1}_T = \mathbf{1}_T^*$,

where $\mathbf{1}_t^* \triangleq \arg \max_{k \in [K]} G_{k, t}$ & $G_{k, t} = \frac{1}{t} \sum_{t'=1}^t g_{k, t'}$








				
G				
				
Indices	3	1	2	4


					
G					
					
Indices	3	1	2	4	
True ranks	1	2	3	4	



k is the **index** of the arm ranked k -th according to G .
 i.e. $G_{\mathbf{1}} > G_{\mathbf{2}} \geq G_{\mathbf{3}} \geq \dots \geq G_{\mathbf{k}} \geq \dots \geq G_{\mathbf{K}}$



 is the **index** of the arm ranked k -th according to G .
 i.e. $G_{\mathbf{1}} > G_{\mathbf{2}} \geq G_{\mathbf{3}} \geq \dots \geq G_{\mathbf{k}} \geq \dots \geq G_{\mathbf{K}}$

 is the **rank** of the arm (of index) k according to G .



k is the **index** of the arm ranked k -th w.r.t. to an estimate of G : \hat{G}, \tilde{G} .
 i.e. $G_{\mathbf{1}} > G_{\mathbf{2}} \geq G_{\mathbf{3}} \geq \dots \geq G_{\mathbf{k}} \geq \dots \geq G_{\mathbf{K}}$

k is the **rank** of the arm (of index) k w.r.t. to an estimate of G : \hat{G}, \tilde{G} .





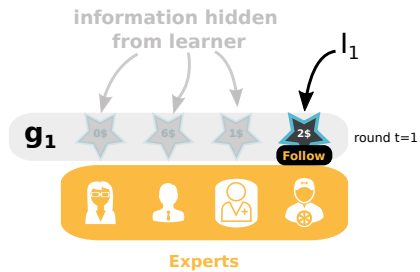
Products

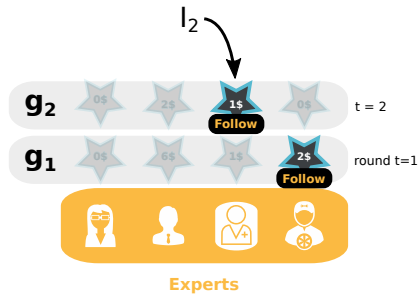


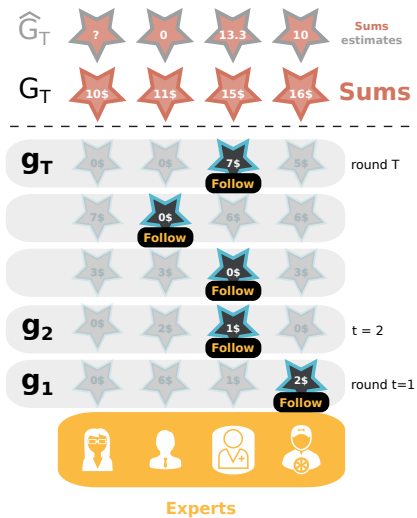
Drugs

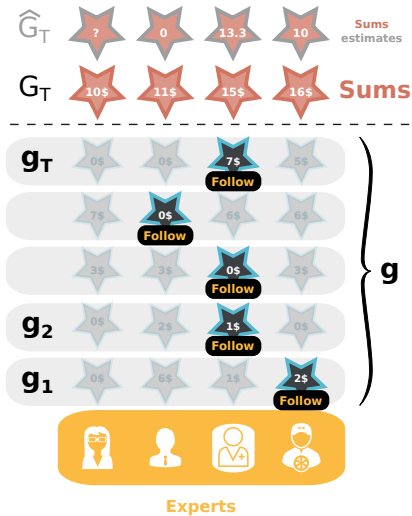


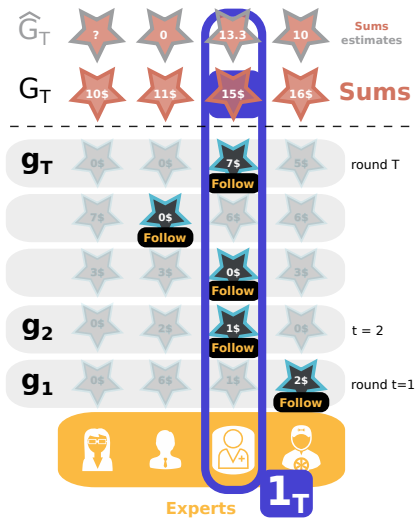
Experts

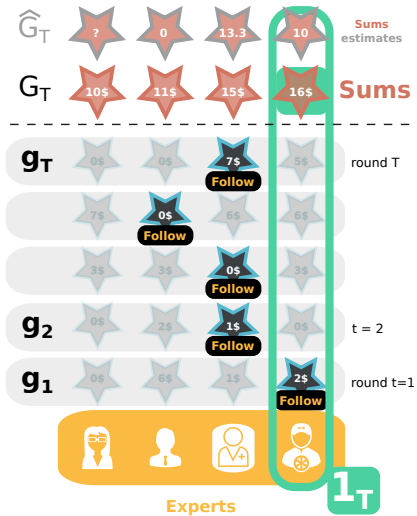












Cumulative regret

- $R(T) = \max_k (G_k) - \sum_{t=1}^T g_{I_t, t} =$




 $- \left($

 $+$

 $+ \dots +$

 $\left. \right)$
- Minimize the **cumulative** regret \Leftrightarrow Play  as often as possible
- Exploration vs Exploitation**
- Classic algorithms:** Thompson Sampling, UCB


Probability of misidentification — simple regret

- $e(T) = \mathbb{P} \left(\mathbf{1}_T \neq \hat{\mathbf{1}}_T \right)$ or $r(T) = G_{\mathbf{1}_T} - G_{\hat{\mathbf{1}}_T} =$

 $-$

- Minimize the **simple** regret \Leftrightarrow Identify 
- Pure Exploration**
- Classic algorithms:** Hoeffding Race, Successive Rejects

How are the rewards, g_1, \dots, g_T , generated?

 **Stochastic**

$g_{k,t} \stackrel{\text{iid}}{\sim} \nu_k$, mean μ_k
 $\mathbf{1}_T = \arg \max_{k \in [K]} \mu_k$

 indifferent to $e_{\mathbf{1}_T}(T)$

 **Non stochastic**


Drop iid, non-stationary

 **Adversarial**

Arbitrary $g = \{g_{k,t}\}_{k \in [K], t \in [T]}$

$\mathbf{1}_T = \arg \max_{k \in [K]} G_k$

$G_k = \frac{1}{T} \sum_{t=1}^T g_{k,t}$

 picks g maximizing $e_{\text{devil}}(T)$

Related works:

Hoefdding race [1]

[4][5][6]

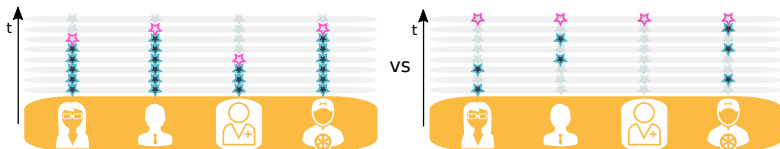
New in Best arm identification

Successive Rejects (SR) [2]

(more on the next slide!)

Similar to adversarial bandit [3]

- **Jamieson & Talwalkar, 2016** for hyperparameter optimization:
 - $\mathbf{g}_{k,t}$ are fixed by an adversary with the condition that $\mathbf{g}_{k,t}$ converge as $h_k = \lim_{t \rightarrow +\infty} \mathbf{g}_{k,t}$ exists.
 - At round t for its m -th pull of arm k , their learner observe $\mathbf{g}_{k,m}$, whereas our learner observes $\mathbf{g}_{k,t}$ (less hidden information).



- **Allesiardo, Féraud & Maillard, 2017:**
 $\mathbf{g}_{k,t}$ are sampled from a non-stationary process with the condition that the identity of the best arm so far does not change with time: $\mathbf{1}_t = \mathbf{1}_{t'}$, $\forall (t, t') \in [T]^2$.
- **Corruption/contamination, Altschuler, Brunel & Malek, 2019:**
 $\mathbf{g}_{k,t}$ are sampled i.i.d. but the learner observes $\mathbf{g}_{k,t} + \mathbf{c}_{k,t}$ where $\mathbf{c}_{k,t}$ can be an arbitrary corruption.



- **Deterministic Uniform exploration** (DETER-UNIFORM).
Pull every arm deterministically T/K times.

- **Successive Rejects** (SR) (Audibert, Bubeck & Munos, 2010)
Pull more the arms with highest estimated average reward.



- **Deterministic Uniform exploration** (DETER-UNIFORM).
Pull every arm deterministically T/K times.

- **Successive Rejects** (SR) (Audibert, Bubeck & Munos, 2010)
Pull more the arms with highest estimated average reward.

The estimated mean of arm k at time t is simply the standard average:

$$\hat{\mu}_{k,t} \triangleq \frac{\sum_{t'=1}^t \mathbf{1}\{I_{t'} = k\} g_{k,t'}}{\sum_{t'=1}^T \mathbf{1}\{I_{t'} = k\}}$$



- **Deterministic Uniform exploration** (DETER-UNIFORM).
Pull every arm deterministically T/K times.
- **Successive Rejects** (SR) (Audibert, Bubeck & Munos, 2010)
Pull more the arms with highest estimated average reward.

The estimated mean of arm k at time t is simply the standard average:

$$\widehat{\mu}_{k,t} \triangleq \frac{\sum_{t'=1}^t \mathbf{1}\{I_{t'} = k\} g_{k,t'}}{\sum_{t'=1}^t \mathbf{1}\{I_{t'} = k\}}$$

- **UCB-Exploration** (Audibert, Bubeck & Munos, 2010)

Pull $\arg \max_{k \in [K]} \widehat{\mu}_{k,t} + g_{\max} \sqrt{\frac{a}{T_k}}$, $a \in \mathbb{R}$, T_k : # of pulls of k .

$g_{\max} \sqrt{\frac{a}{T_k}}$ is the *uncertainty* on $\widehat{\mu}_{k,t}$ but requires knowledge of g_{\max} .



- **Deterministic Uniform exploration** (DETER-UNIFORM).
Pull every arm deterministically T/K times.
- **Successive Rejects** (SR) (Audibert, Bubeck & Munos, 2010)
Pull more the arms with highest estimated average reward.

The estimated mean of arm k at time t is simply the standard average:

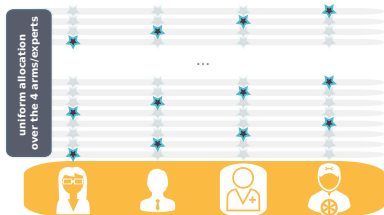
$$\hat{\mu}_{k,t} \triangleq \frac{\sum_{t'=1}^t \mathbf{1}\{I_{t'} = k\} g_{k,t'}}{\sum_{t'=1}^T \mathbf{1}\{I_{t'} = k\}}$$

	$e_{\triangle} (T)$	$e_{\ominus} (T)$
DETER-UNIFORM	✗	?
SR	✓	?

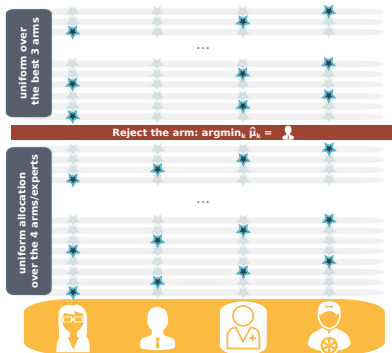
- SR is an elimination algorithm pulling uniformly over a set of remaining candidate arms.
- The arm k , ranked k -th by SR, is allocated T/k pulls



- SR is an elimination algorithm pulling uniformly over a set of remaining candidate arms.
- The arm k , ranked k -th by SR, is allocated T/k pulls



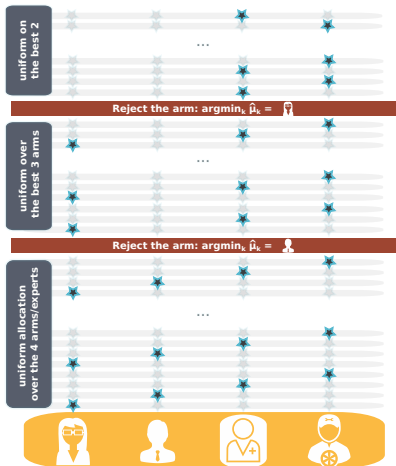
- SR is an elimination algorithm pulling uniformly over a set of remaining candidate arms.
- The arm k , ranked r -th by SR, is allocated T/r pulls



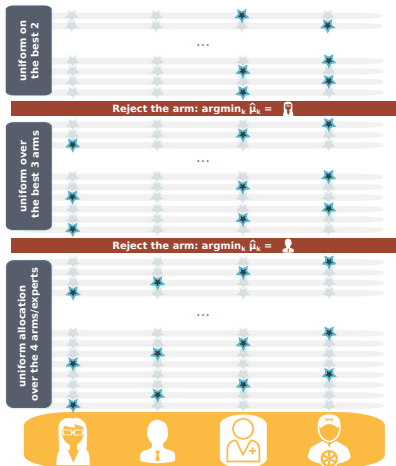
The estimated mean of arm k at time t is simply the standard average:

$$\hat{\mu}_{k,t} \triangleq \frac{\sum_{t'=1}^t \mathbf{1}\{I_{t'}=k\} g_{k,t'}}{\sum_{t'=1}^t \mathbf{1}\{I_{t'}=k\}}$$

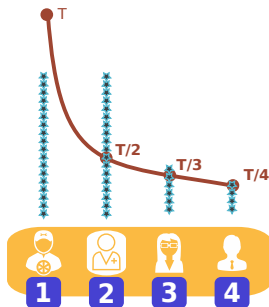
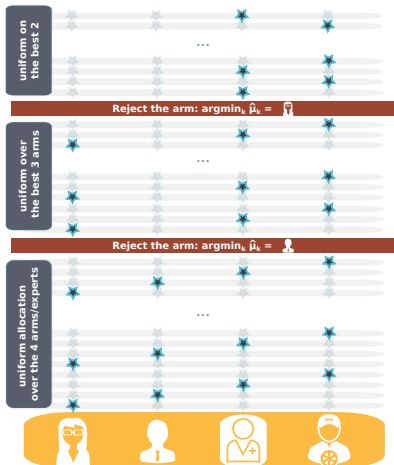
- SR is an elimination algorithm pulling uniformly over a set of remaining candidate arms.
- The arm k , ranked r -th by SR, is allocated T/r pulls



- SR is an elimination algorithm pulling uniformly over a set of remaining candidate arms.
- The arm k , ranked k -th by SR, is allocated T/k pulls



- SR is an elimination algorithm pulling uniformly over a set of remaining candidate arms.
- The arm k , ranked k -th by SR, is allocated T/k pulls

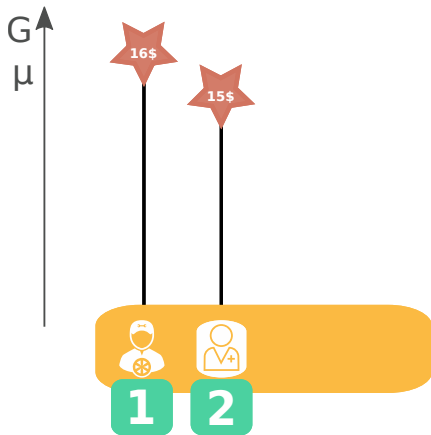


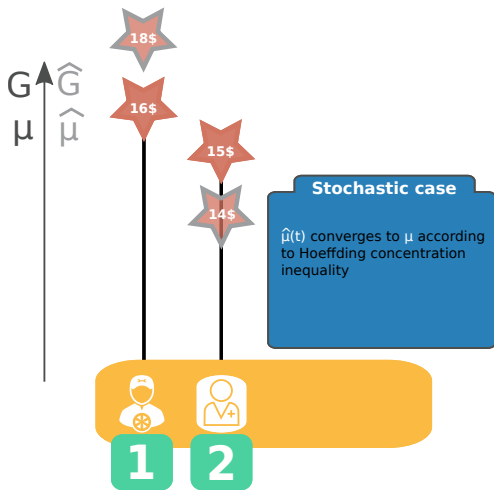


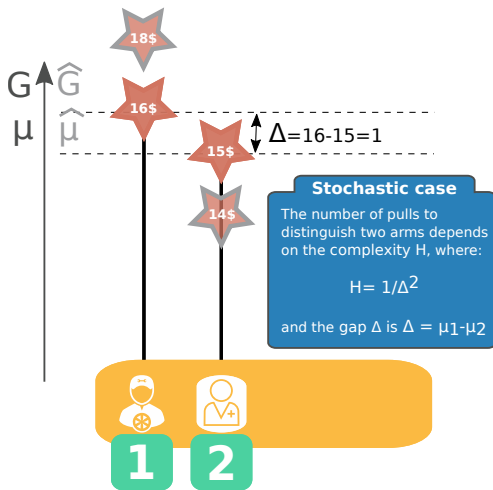
	$e_{\triangle}(T)$	$e_{\text{cat}}(T)$
DETER-UNIFORM	✗ ??	?
SR	✓ ??	?

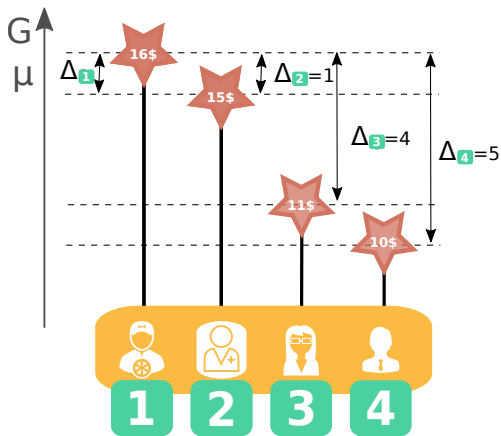
And now...

Let us precise the $e_{\triangle}(T)$ for the uniform and SR algorithms.

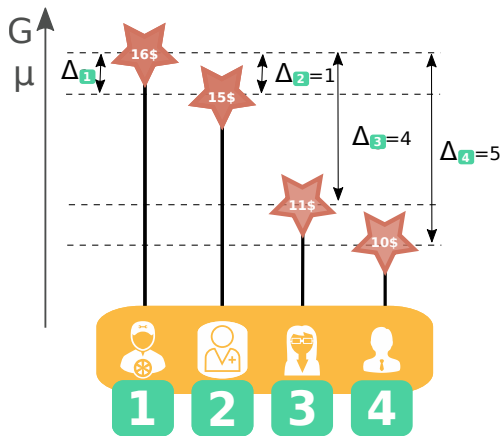




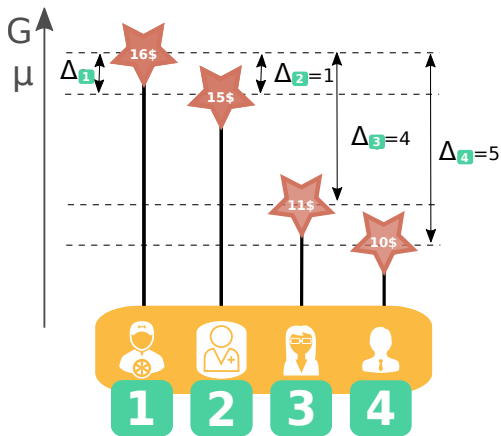




To distinguish arm k from arm 1 , the learner must have its *uncertainty* on μ_k (or G_k) smaller than Δ_k , i.e. $|\hat{\mu}_k - \mu_k| \leq \Delta_k/2$.

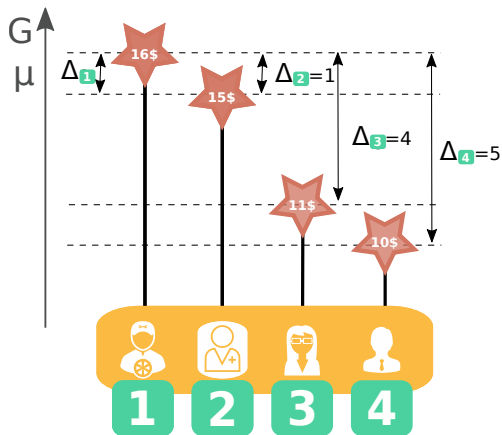


$$H_{\text{UNIF}} \triangleq \frac{K}{\Delta_1^2}$$







$$H_{\text{UNIF}} \triangleq \frac{K}{\Delta_1^2} \quad \& \quad H_{\text{SR}} \triangleq \max_{k \in [K]} \frac{k}{\Delta_k^2}$$

Stochastic case



$$H_{\text{UNIF}} \triangleq \frac{K}{\Delta_1^2} \geq H_{\text{SR}} \triangleq \max_{k \in [K]} \frac{k}{\Delta_k^2} = \tilde{O} \left(\sum_{k \in [K]} \frac{1}{\Delta_k^2} \right)$$


Stochastic case




		$e_{\triangle}(T)$	$e_{\text{smiley}}(T)$
DETER-UNIFORM		$e^{-T/H_{\text{UNIF}}}$	
SR		$e^{-T/H_{\text{SR}} \log K}$	

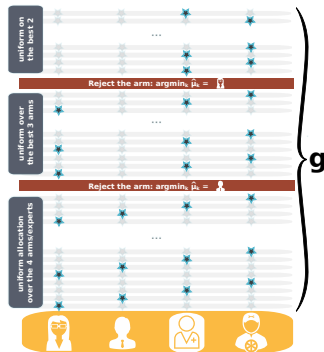
And now...

Let us discuss SR against an adversary.












SR can be tricked by an adversary  arranging g

- SR pulls the arm deterministically ( will hide rewards easily)
- SR stops pulling arms (reject) during the game ( hides rewards)
- SR uses the standard estimation of the average $\hat{\mu}_{k,t}$ (**biased** against )



- The learner needs to use internal randomization
- The learner should be careful about rejecting arm: no rejection!
- Be careful of the bias of the reward estimates.

	$e_{\triangle}(T)$	$e_{\text{cat}}(T)$
DETER-UNIFORM	 $e^{\frac{-T}{H_{\text{UNIF}}}}$	 1
SR	 $e^{\frac{-T}{H_{\text{SR}} \log K}}$	 1
   		

	$e_{\triangle}(T)$	$e_{\text{devil}}(T)$
DETER-UNIFORM	$\times e^{-T/H_{\text{UNIF}}}$	$\times 1$
SR	$\checkmark e^{-T/H_{\text{SR}} \log K}$	$\times 1$
? ? ? ?		\checkmark

And now...

Let us discuss the adversarial e_{devil} setting



•DETER-UNIFORM•:

- ▶ Pull every arm deterministically T/K times.
- ▶ Recommend the arm with highest $\hat{\mu}_{k,t}$

Robutifying

- **Internal randomization:** pull arm k at time t with proba $p_{k,t} = \mathbb{P}(I_t = k)$
- Replace $\hat{\mu}_{k,t}$ by $\tilde{G}_{k,t}$ as $E[\tilde{G}_{k,t}] = G_{k,t}$ (**unbiased**)

$$\hat{\mu}_{k,t} \triangleq \frac{\sum_{t'=1}^t \mathbf{1}\{I_{t'} = k\} g_{k,t'}}{\sum_{t'=1}^T \mathbf{1}\{I_{t'} = k\}}$$

$$\tilde{G}_{k,t} = \frac{1}{t} \sum_{t'=1}^t \frac{g_{k,t'}}{p_{k,t'}} \mathbf{1}\{I_{t'} = k\}$$

(importance weights)

•RULE•:

- ▶ At time t , pull arm k with probability $p_{k,t} = 1/K$.
- ▶ Recommend the arm with highest $\tilde{G}_{k,t}$



Theorem (RULE vs.)

For any T and adversarial g , RULE satisfies

$$e_{\text{RULE}}(T) = \mathcal{O} \left(\exp \left(-\frac{T}{H_{\text{UNIF}}(g)} \right) \right).$$

The proof uses a Bernstein bound.

Theorem (RULE vs. )

For any T and adversarial g , RULE satisfies

$$e_{\text{RULE}}(T) = \mathcal{O} \left(\exp \left(-\frac{T}{H_{\text{UNIF}}(g)} \right) \right).$$

The proof uses a Bernstein bound.



Theorem ( Lower bound)

For any learner, a g of complexity H_{UNIF} ,

$$e_{\text{RULE}}(T) = \Omega \left(\exp \left(-\frac{T}{H_{\text{UNIF}}(g)} \right) \right).$$



RULE: optimal gap-dependent rates against .

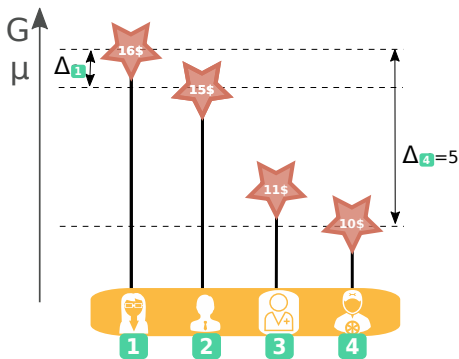


Idea: The  can force the learner to have, at $t = T/2$, an uncertainty on $\tilde{G}_{k,T/2}$ of order $\Delta_{\mathbf{1}}$, $\forall k \in [K]$ (instead of the usual $\Delta_{\mathbf{k}}$ in ).

Our proof of the lower bound uses some arguments of **Audibert & Bubeck (2010)**, **Carpentier and Locatelli (2016)** and **Auer and Chiang (2016)**





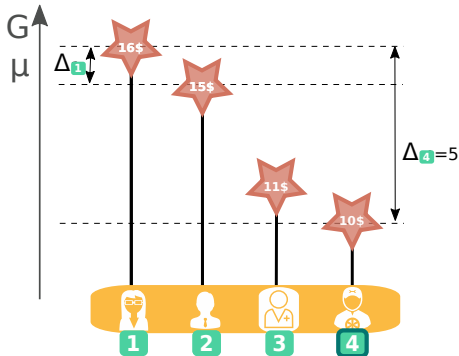
Idea: The  can force the learner to have, at $t = T/2$, an uncertainty on $\tilde{G}_{k, T/2}$ of order $\Delta_{\mathbf{1}}$, $\forall k \in [K]$ (instead of the usual $\Delta_{\mathbf{k}}$ in ).




Given a Learner and a bandit PROBLEM I defined for the first half of the game (until $t = T/2$)





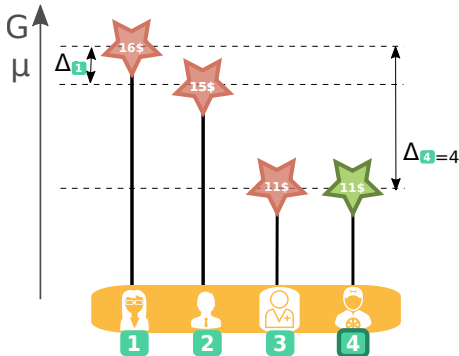
Idea: The  can force the learner to have, at $t = T/2$, an uncertainty on $\tilde{G}_{k,T/2}$ of order $\Delta_{\mathbf{1}}$, $\forall k \in [K]$ (instead of the usual $\Delta_{\mathbf{k}}$ in ).




At least one arm is pulled less than $T/(2K)$ by the Learner (not pulled enough to detect small variations of size $\Delta_{\mathbf{1}}$, of its mean \leftarrow prone to error). Here its arm 





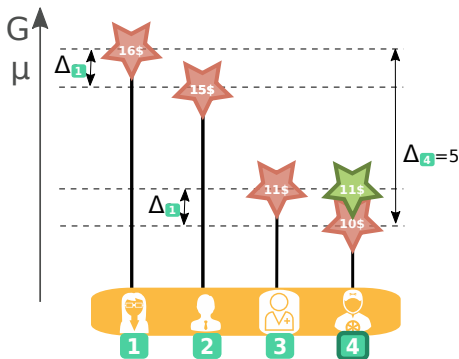
Idea: The  can force the learner to have, at $t = T/2$, an uncertainty on $\tilde{G}_{k,T/2}$ of order Δ_1 , $\forall k \in [K]$ (instead of the usual Δ_k in ).



Then, an alternative/similar PROBLEM II is created, by modifying  by Δ_1 .
 PROBLEM II is defined for $t = 1$ to $t = T/2$.



Idea: The  can force the learner to have, at $t = T/2$, an uncertainty on $\tilde{G}_{k,T/2}$ of order Δ_1 , $\forall k \in [K]$ (instead of the usual Δ_k in ).

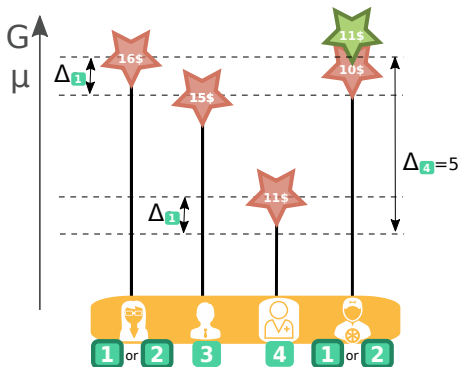


This is the superposition of PROBLEM I & II.

PROBLEM I & II are indistinguishable with proba $e^{-\frac{\tau \Delta_1^2}{K}}$



Idea: The can force the learner to have, at $t = T/2$, an uncertainty on $\tilde{G}_{k,T/2}$ of order $\Delta_{\mathbf{1}}$, $\forall k \in [K]$ (instead of the usual $\Delta_{\mathbf{k}}$ in).

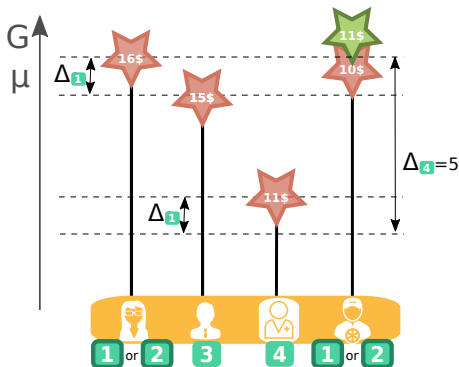


Between $t = T/2$ and $t = T$, the arm get enough reward so that:

= **1** in PROBLEM II while = **1** in PROBLEM I



Idea: The can force the learner to have, at $t = T/2$, an uncertainty on $\tilde{G}_{k,T/2}$ of order $\Delta_{\mathbf{1}}$, $\forall k \in [K]$ (instead of the usual $\Delta_{\mathbf{k}}$ in).



The lower bound comes from the fact that PROBLEM I & II

- have different best arms

- are indistinguishable w.p. $e^{-\frac{T\Delta_{\mathbf{1}}^2}{K}}$, i.e. $\underbrace{P_{II}(\mathbf{1} = \text{lightbulb})}_{\text{error in II}} \geq \underbrace{P_I(\mathbf{1} = \text{lightbulb})}_{\text{success in I}} e^{-\frac{T\Delta_{\mathbf{1}}^2}{K}}$

- Best arm identification against 😈 is too hard: uniform exploration (RULE) is optimal.
- However RULE is suboptimal in ▲.
- SR, optimal in ▲ fails against 😈

	$e_{\blacktriangle}(T)$	$e_{\text{😈}}(T)$
SR	✓ $e^{\frac{-T}{H_{\text{SR}} \log K}}$	✗ 1
RULE	✗ $e^{\frac{-T}{H_{\text{UNIF}}}}$	✓ $e^{\frac{-T}{H_{\text{UNIF}}}}$

The question:

Is there a learner performing optimally in both the stochastic and adversarial cases while not being aware of the nature of the rewards

	$e_{\blacktriangle}(T)$	$e_{\text{😈}}(T)$
	$e^{\frac{-T}{H_{\text{SR}} \log K}}$	$e^{\frac{-T}{H_{\text{UNIF}}}}$

- Best arm identification against 😈 is too hard: uniform exploration (RULE) is optimal.
- However RULE is suboptimal in ▲.
- SR, optimal in ▲ fails against 😈

	$e_{\blacktriangle}(T)$	$e_{\text{😈}}(T)$
SR	✓ $e^{\frac{-T}{H_{\text{SR}} \log K}}$	✗ 1
RULE	✗ $e^{\frac{-T}{H_{\text{UNIF}}}}$	✓ $e^{\frac{-T}{H_{\text{UNIF}}}}$

The BOB question:

❓ *Is there a learner performing optimally in both the stochastic and adversarial cases while not being aware of the nature of the rewards ?*

	$e_{\blacktriangle}(T)$	$e_{\text{😈}}(T)$
❓❓❓❓	✓ $e^{\frac{-T}{H_{\text{SR}} \log K}}$	✓ $e^{\frac{-T}{H_{\text{UNIF}}}}$

- Best arm identification against 😊 is too hard: uniform exploration (RULE) is optimal.
- However RULE is suboptimal in ▲.
- SR, optimal in ▲ fails against 😊

	$e_{\blacktriangle}(T)$	$e_{\😊}(T)$
SR	✓ $e^{H_{SR} \frac{-T}{\log K}}$	✗ 1
RULE	✗ $e^{H_{UNIF} \frac{-T}{}}$	✓ $e^{H_{UNIF} \frac{-T}{}}$

The **BOB** question:

❓ *Is there a learner performing optimally in **both** the *stochastic* and *adversarial* cases while not being aware of the nature of the rewards ?*

The **BOB** question was studied in the cumulative regret setting in Bubeck & Slivkins, 2012, Seldin & Slivkins, 2014, Auer & Chiang, 2016, Zimmert & Seldin, 2018...

New notion of complexity

$$H_{\text{BOB}} \triangleq \frac{1}{\Delta_{\mathbf{1}}} \max_{k \in [K]} \frac{k}{\Delta_{\mathbf{k}}}.$$

Theorem (Lower bound for the BOB challenge)

For any learner,

if for all adversarial problem g ,

$$e_{\text{smiley}}(T) \leq \frac{1}{16},$$

then there exists a stochastic problem with complexity H_{BOB} such that

$$e_{\text{sad}}(T) \geq \frac{1}{64} \exp\left(-\frac{2048T}{H_{\text{BOB}}}\right) \stackrel{\text{sometimes}}{=} \frac{1}{64} \exp\left(-\frac{2048T}{H_{\text{SR}}\sqrt{K}}\right).$$

New notion of complexity

$$H_{\text{BOB}} \triangleq \frac{1}{\Delta_{\mathbf{1}}} \max_{k \in [K]} \frac{k}{\Delta_{\mathbf{k}}}.$$

Theorem (Lower bound for the **BOB** challenge)




For any learner,

if for all adversarial problem g ,

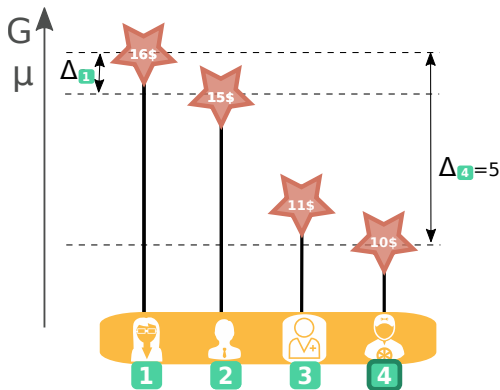
$$e_{\text{smiley}}(T) \leq \frac{1}{16},$$

then there exists a stochastic problem with complexity H_{BOB} such that


$$e_{\text{sad}}(T) \geq \frac{1}{64} \exp\left(-\frac{2048T}{H_{\text{BOB}}}\right) \stackrel{\text{sometimes}}{=} \frac{1}{64} \exp\left(-\frac{2048T}{H_{\text{SR}}\sqrt{K}}\right).$$

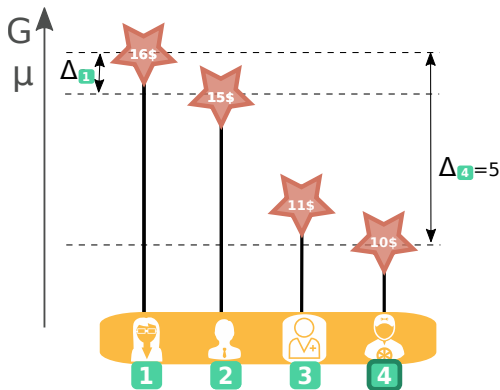
Idea: The  can force the learner to have, $\forall k \in [K]$, at $t = T \frac{\Delta_{\mathbf{1}}}{\Delta_{\mathbf{k}}}$ (instead of $t = T$ in ) , an uncertainty on  of order $\Delta_{\mathbf{k}}$.


Idea: The 😈 can force the learner to have, $\forall k \in [K]$, at $t = T \frac{\Delta_1}{\Delta_k}$ (instead of $t = T$ in \blacktriangle), an uncertainty on $\tilde{G}_{k,t}$ of order Δ_k .



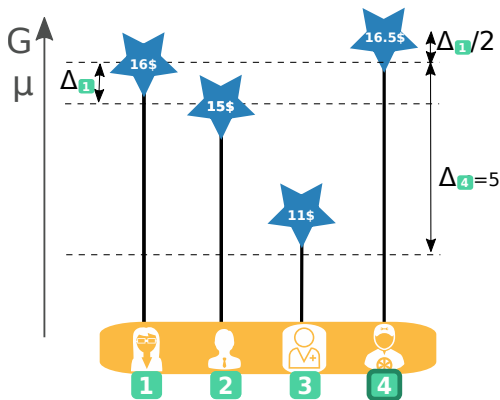
Given a Learner, a rank k , and a BASE PROBLEM defined until $t = T \frac{\Delta_1}{\Delta_k}$


Idea: The 😊 can force the learner to have, $\forall k \in [K]$, at $t = T \frac{\Delta_1}{\Delta_k}$ (instead of $t = T$ in ) , an uncertainty on $\tilde{G}_{k,t}$ of order Δ_k .



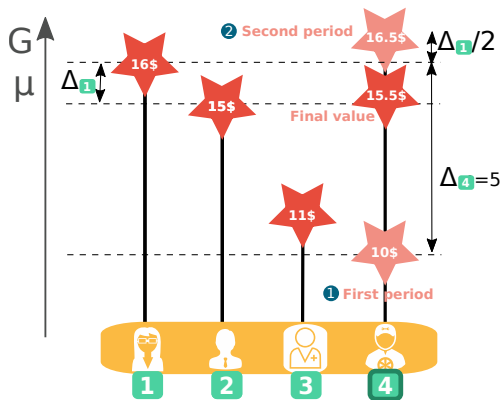
At least one arm, k' , with $\Delta_{k'} \leq \Delta_k$, is pulled less than $\frac{T}{k} \frac{\Delta_1}{\Delta_k}$ (not pulled enough). Let us illustrate $k = 4$ and assume for simplicity $k' = 4 =$ 

Idea: The 😊 can force the learner to have, $\forall k \in [K]$, at $t = T \frac{\Delta_1}{\Delta_k}$ (instead of $t = T$ in \blacktriangle), an uncertainty on $\tilde{G}_{k,t}$ of order Δ_k .



A similar **PROBLEM STO**  is created, \blacktriangle , by modifying by $\Delta_4 + \Delta_1/2$. **PROBLEM STO** is defined for $t = 1$ to $t = T$.

Idea: The 😈 can force the learner to have, $\forall k \in [K]$, at $t = T \frac{\Delta_1}{\Delta_k}$ (instead of $t = T$ in \blacktriangle), an uncertainty on $\tilde{G}_{k,t}$ of order Δ_k .

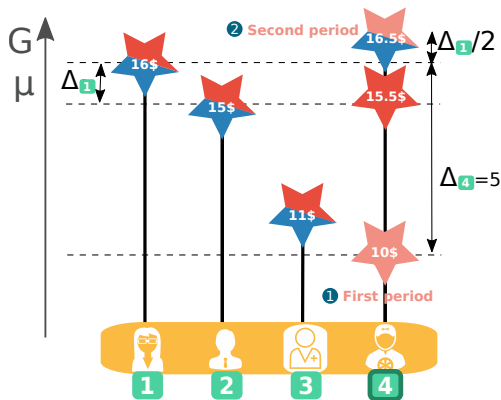


PROBLEM **ADV** 😈: ① follows **BASE** from $t = 1$ to $t = T \frac{\Delta_1}{\Delta_k}$,

② follows **STO** afterwards .

Modifying μ_k of Δ_k during $T \frac{\Delta_1}{\Delta_k}$ changes the means of Δ_1 w.r.t **STO**.

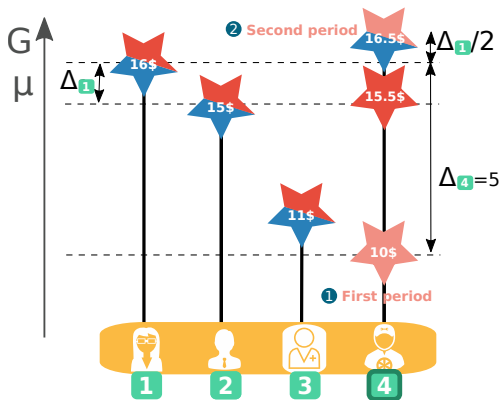
Idea: The 😈 can force the learner to have, $\forall k \in [K]$, at $t = T \frac{\Delta_1}{\Delta_k}$ (instead of $t = T$ in \blacktriangle), an uncertainty on $\tilde{G}_{k,t}$ of order Δ_k .



Superposition of **ADV** & **STO**: = **1** in **STO** and = **1** in **ADV**,

ADV & **STO** are indistinguishable w. p. $e^{-\frac{\left(T \frac{\Delta_1}{\Delta_k}\right) \Delta_k^2}{k}} = e^{-\frac{T \Delta_1 \Delta_k}{k}}$

Idea: The 😊 can force the learner to have, $\forall k \in [K]$, at $t = T \frac{\Delta_1}{\Delta_k}$ (instead of $t = T$ in \blacktriangle), an uncertainty on $\tilde{G}_{k,t}$ of order Δ_k .



The lower bound comes from, a **max** over $k \in [K]$, and the fact that **ADV** & **STO**

- have different best arms,
- are indistinguishable w.p. $e^{-\frac{T\Delta_1\Delta_k}{k}}$: $\underbrace{P_{\text{STO}}(\mathbf{1} = \text{👤})}_{\text{error in II}} \geq \underbrace{P_{\text{ADV}}(\mathbf{1} = \text{👤})}_{\text{success in I}} e^{-\frac{T\Delta_1\Delta_k}{k}}$

New notion of complexity

$$H_{\text{BOB}} \triangleq \frac{1}{\Delta_1} \max_{k \in [K]} \frac{k}{\Delta_k}.$$

Theorem (Lower bound for the BOB challenge)

For any learner,

if for all adversarial problems g ,

$$e_{\text{BOB}}(T) \leq \frac{1}{16},$$

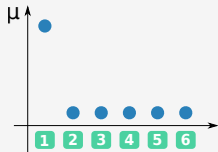
then there exists a stochastic problem with complexity H_{BOB} such that

$$e_{\text{BOB}}(T) \geq \frac{1}{64} \exp\left(-\frac{2048T}{H_{\text{BOB}}}\right) \stackrel{\text{sometimes}}{=} \frac{1}{64} \exp\left(-\frac{2048T}{H_{\text{SR}}\sqrt{K}}\right).$$

$$H_{SR} \leq H_{BOB} \leq H_{UNIF}.$$

$$\underbrace{\max_{k \in [K]} \frac{k}{\Delta_k^2}}_{H_{SR}} \leq \underbrace{\frac{1}{\Delta_1} \max_{k \in [K]} \frac{k}{\Delta_k}}_{H_{BOB}} \leq \underbrace{\frac{1}{\Delta_1^2}}_{H_{UNIF}}.$$

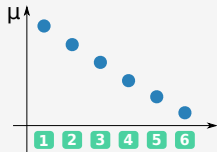
► Flat regime



$$H_{SR} = H_{BOB} = H_{UNIF}$$

BOB is achieved by Rule.

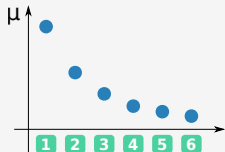
► Linear regime



$$H_{SR} = H_{BOB} = \frac{H_{UNIF}}{K}$$

**BOB can be achieved but not by Rule.
Need a new learner!**

► Square-root regime



$$H_{SR} = \frac{H_{BOB}}{\sqrt{2K}} = \frac{H_{UNIF}}{K}$$

No learner can do BOB!

There is still an open question!

The new **BOB** question:

❓ Can an algorithm achieve the following ❓

	$e_{\triangle} (T)$	$e_{\text{cat}} (T)$
❓ ❓ ❓ ❓	✓ $e^{H_{\text{BOB}} \frac{-T}{\log K}}$	✓ $e^{H_{\text{UNIF}} \frac{-T}{\log K}}$

Why is the BOB question challenging?

- ▶ Bias of estimator $\hat{\mu}_{k,t} = \frac{\sum_{t'=1}^t \mathbf{1}\{I_{t'}=k\} g_{k,t'}}{\sum_{t'=1}^t \mathbf{1}\{I_{t'}=k\}}$ (simple average)
- ▶ Variance of $\tilde{G}_{k,t} = \sum_{t'=1}^t \frac{g_{k,t'}}{p_{k,t'}} \mathbf{1}\{I_{t'}=k\}$ (importance weights)

We use :

- Pulling uniformly for too long with $p_{k,t} = \frac{1}{K}$ leads to a large variance, up to being of order K , in $\tilde{G}_{k,t}$.
- Objective: reduce the variance (**uncertainty**) of the estimators of the best arms \approx find the best arm

The new **BOB** question:

❓ Can an algorithm achieve the following ❓

	$e_{\triangle} (T)$	$e_{\text{cat}} (T)$
❓ ❓ ❓ ❓	✓ $e^{H_{\text{BOB}} \frac{-T}{\log K}}$	✓ $e^{H_{\text{UNIF}} \frac{-T}{\log K}}$

Why is the BOB question challenging?

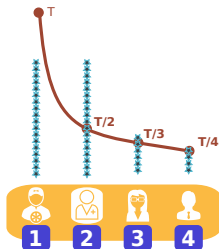
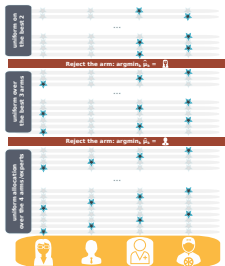
- ▶ Bias of estimator $\hat{\mu}_{k,t} = \frac{\sum_{t'=1}^t \mathbf{1}\{I_{t'}=k\} g_{k,t'}}{\sum_{t'=1}^t \mathbf{1}\{I_{t'}=k\}}$ (simple average)
- ▶ Variance of $\tilde{G}_{k,t} = \sum_{t'=1}^t \frac{g_{k,t'}}{p_{k,t'}} \mathbf{1}\{I_{t'}=k\}$ (importance weights)

We use $\tilde{G}_{k,t}$:

- Pulling uniformly for too long with $p_{k,t} = \frac{1}{K}$ leads to a large **variance**, up to being of order K , in $\tilde{G}_{k,t}$.
- Objective: reduce the variance (**uncertainty**) of the estimators of the best arms \approx find the best arm

Idea: Robustify the SR algorithm.

- We use $\tilde{G}_{k,t}$
- Cannot pull uniformly, as in SR, for almost half of the game.
- Need to pull the estimated best arms earlier.
- Need to remove the rejections
- Reuse the proportions of SR (arm k , ranked k -th by SR, is allocated T/k pulls)



At any time t , P1 pulls	• arm 1 with 'probability'	1
	• arm with 'probability'	$\frac{1}{2}$
	• arm with 'probability'	$\frac{1}{3}$
	• and so on...	
	• arm with 'probability'	$\frac{1}{k}$
	• and with 'probability'	$\frac{1}{K}$
	• (and normalize)	

$$\overline{\log K} = \sum_{k=1}^K (1/k), \text{ with } |\overline{\log K} - \log K| \leq 1$$

P1 early bets are almost **costless!** (and necessary):

- The estimated best arms are prioritized since the first pull to reduce variance.
- Up to a $\log K$ factor, all arms are pulled uniformly.
- P1 implicitly control the uncertainty of the estimates.

At any time t , P1 pulls	• arm 1 with 'probability'	1
	• arm 2 with 'probability'	$\frac{1}{2}$
	• arm with 'probability'	$\frac{1}{3}$
	• and so on...	
	• arm with 'probability'	$\frac{1}{k}$
	• and with 'probability'	$\frac{1}{K}$
	• (and normalize)	

$$\overline{\log K} = \sum_{k=1}^K (1/k), \text{ with } |\overline{\log K} - \log K| \leq 1$$

P1 early bets are almost **costless!** (and necessary):

- The estimated best arms are prioritized since the first pull to reduce variance.
- Up to a $\log K$ factor, all arms are pulled uniformly.
- P1 implicitly control the uncertainty of the estimates.

At any time t , P1 pulls	• arm 1 with 'probability'	1
	• arm 2 with 'probability'	$\frac{1}{2}$
	• arm 3 with 'probability'	$\frac{1}{3}$
	• and so on...	
	• arm with 'probability'	$\frac{1}{k}$
	• and with 'probability'	$\frac{1}{K}$
	• (and normalize)	

$$\overline{\log K} = \sum_{k=1}^K (1/k), \text{ with } |\overline{\log K} - \log K| \leq 1$$

P1 early bets are almost **costless!** (and necessary):

- The estimated best arms are prioritized since the first pull to reduce variance.
- Up to a $\log K$ factor, all arms are pulled uniformly.
- P1 implicitly control the uncertainty of the estimates.

At any time t , P1 pulls	• arm 1 with 'probability'	1
	• arm 2 with 'probability'	$\frac{1}{2}$
	• arm 3 with 'probability'	$\frac{1}{3}$
	• and so on...	
	• arm k with 'probability'	$\frac{1}{k}$
	• and K with 'probability'	$\frac{1}{K}$
	• (and normalize)	

$$\overline{\log K} = \sum_{k=1}^K (1/k), \text{ with } |\overline{\log K} - \log K| \leq 1$$

P1 early bets are almost **costless!** (and necessary):

- The estimated best arms are prioritized since the first pull to reduce variance.
- Up to a $\log K$ factor, all arms are pulled uniformly.
- P1 implicitly control the uncertainty of the estimates.

At any time t , P1 pulls	• arm 1 with 'probability'	1
	• arm 2 with 'probability'	$\frac{1}{2}$
	• arm 3 with 'probability'	$\frac{1}{3}$
	• and so on...	
	• arm k with 'probability'	$\frac{1}{k}$
	• and K with 'probability'	$\frac{1}{K}$
	• (and normalize)	

$$\overline{\log K} = \sum_{k=1}^K (1/k), \text{ with } |\overline{\log K} - \log K| \leq 1$$

P1 early bets are almost **costless!** (and necessary):

- The estimated best arms are prioritized since the first pull to reduce variance.
- Up to a $\log K$ factor, all arms are pulled uniformly.
- P1 implicitly control the uncertainty of the estimates.

- At any time t , P1 pulls
- arm **1** with probability $\frac{1}{\overline{\log K}}$
 - arm **2** with probability $\frac{1}{2 \overline{\log K}}$
 - arm **3** with probability $\frac{1}{3 \overline{\log K}}$
 - and so on...
 - arm **k** with probability $\frac{1}{k \overline{\log K}}$
 - and **K** with probability $\frac{1}{K \overline{\log K}}$
 - (and normalize)

$$\overline{\log K} = \sum_{k=1}^K (1/k), \text{ with } |\overline{\log K} - \log K| \leq 1$$

P1 early bets are almost **costless!** (and necessary):

- The estimated best arms are prioritized since the first pull to reduce variance.
- Up to a $\log K$ factor, all arms are pulled uniformly.
- P1 implicitly control the uncertainty of the estimates.

- At any time t , P1 pulls
- arm **1** with probability $\frac{1}{\overline{\log K}}$
 - arm **2** with probability $\frac{1}{2 \overline{\log K}}$
 - arm **3** with probability $\frac{1}{3 \overline{\log K}}$
 - and so on...
 - arm **k** with probability $\frac{1}{k \overline{\log K}}$
 - and **K** with probability $\frac{1}{K \overline{\log K}}$
 - (and normalize)

$$\overline{\log K} = \sum_{k=1}^K (1/k), \text{ with } |\overline{\log K} - \log K| \leq 1$$

P1 early bets are almost **costless!** (and necessary):

- The estimated best arms are prioritized since the first pull to reduce variance.
- Up to a $\log K$ factor, all arms are pulled uniformly.
- P1 implicitly control the uncertainty of the estimates.

For $t = 1, 2, \dots$

- ▶ Rank the arms according to $\tilde{G}_{k,t}$: Rank arm k as k_t .
- ▶ Select arm $I_t \in [K]$ with

$$p_{k,t} \triangleq \mathbb{P}(I_t = k) \triangleq \frac{1}{k_t \log K} \quad \text{for all } k \in [K].$$

Recommend, at any round t , $\mathbf{1}_t \triangleq \arg \max_{k \in [K]} \tilde{G}_{k,t}$.

The algorithm is anytime and parameter-free.

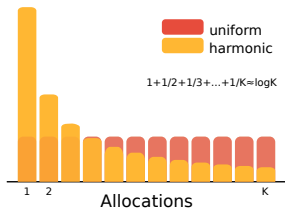
For $t = 1, 2, \dots$

- ▶ Rank the arms according to $\tilde{G}_{k,t}$: Rank arm k as k_t .
- ▶ Select arm $I_t \in [K]$ with

$$p_{k,t} \triangleq \mathbb{P}(I_t = k) \triangleq \frac{1}{k_t \overline{\log K}} \quad \text{for all } k \in [K].$$

Recommend, at any round t , $\mathbf{1}_t \triangleq \arg \max_{k \in [K]} \tilde{G}_{k,t}$.

The algorithm is anytime and parameter-free.



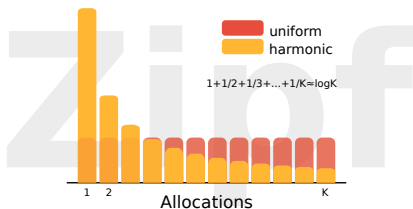
For $t = 1, 2, \dots$

- ▶ Rank the arms according to $\tilde{G}_{k,t}$: Rank arm k as k_t .
- ▶ Select arm $I_t \in [K]$ with

$$p_{k,t} \triangleq \mathbb{P}(I_t = k) \triangleq \frac{1}{k_t \overline{\log K}} \quad \text{for all } k \in [K].$$

Recommend, at any round t , $\mathbf{1}_t \triangleq \arg \max_{k \in [K]} \tilde{G}_{k,t}$.

The algorithm is anytime and parameter-free.



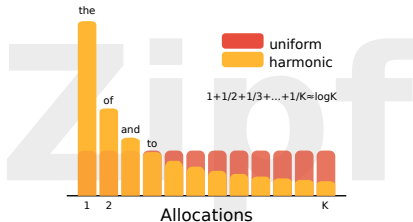
For $t = 1, 2, \dots$

- ▶ Rank the arms according to $\tilde{G}_{k,t}$: Rank arm k as k_t .
- ▶ Select arm $I_t \in [K]$ with

$$p_{k,t} \triangleq \mathbb{P}(I_t = k) \triangleq \frac{1}{k_t \overline{\log K}} \quad \text{for all } k \in [K].$$

Recommend, at any round t , $\mathbf{1}_t \triangleq \arg \max_{k \in [K]} \tilde{G}_{k,t}$.

The algorithm is anytime and parameter-free.



Theorem (Upper bounds for P1)

For any problems:


- ▶ $e_{\triangle}(T) = \mathcal{O}\left(\exp\left(-\frac{T}{H_{\text{BOB}} \log^2(K)}\right)\right)$
- ▶ $e_{\text{😊}}(T) = \mathcal{O}\left(\exp\left(-\frac{T}{H_{\text{UNIF}(g)} \log(K)}\right)\right)$

	$e_{\triangle}(T)$	$e_{\text{😊}}(T)$
SR	✓ $e^{\frac{-T}{H_{\text{SR}} \log K}}$	✗ 1
RULE	✗ $e^{\frac{-T}{H_{\text{UNIF}}}}$	✓ $e^{\frac{-T}{H_{\text{UNIF}}}}$
P1	✓ $e^{\frac{-T}{H_{\text{BOB}} \log K}}$	✓ $e^{\frac{-T}{H_{\text{UNIF}}}}$

Early bets are costless / Early bets are necessary



- $\mathbf{p}_{k,t} \geq 1/(K \overline{\log K})$ is enough to obtain the same complexity H_{UNIF} as RULE, up to a factor $\log K$.

- , $K - 1$ arbitrary 'virtual' phases that each ends at round $T_i = Ta_i$. Chosen in hindsight to minimize the upper bound (P1 is oblivious to a_i).

Intuitively, after T_i the event ξ_i happens with high probability:

$$\xi_i \triangleq \left\{ \forall t > T_i, \forall k \in [K] : \mu_{\mathbf{1}} - \mu_k < \frac{\Delta_i}{2} \implies \mathbf{k}_t < i \right\}.$$

- \implies for any such arm k , for $t > T_i$, $\mathbf{p}_{k,t} \geq 1/(i - 1)$.
- \implies smaller variance (of order $i - 1$) in their estimates $\tilde{\mathbf{g}}_{k,t}$
- \implies better estimates in the rest of the game.

The proof works iteratively over the phases.

Error = ξ_i does not hold.

Trade off in setting the length of the phases with a_i :

Trade off between event ξ_i happening fast and ξ_i happening with high probability

Short phases = not enough samples to discriminate the suboptimal arms.

Long phases = the variance of the mean estimators of good arms is increasing with the length of the early phases

$$H_{P1}(\mathbf{a}) \triangleq \max_{k \in [K]} \frac{\sum_{i=\langle k \rangle}^K (a_i - a_{i+1})i + \frac{1}{24} K a_{\langle k \rangle} \Delta_k}{a_{\langle k \rangle}^2 \Delta_k^2} \frac{1}{\log K}$$

$$H_{P1} \triangleq \min_{\mathbf{a} \in \mathcal{A}} H_{P1}(\mathbf{a}).$$

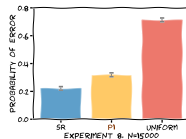
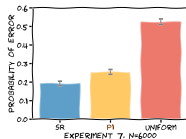
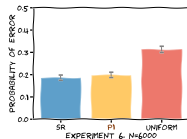
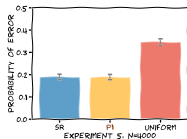
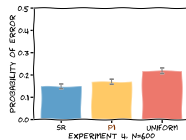
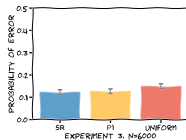
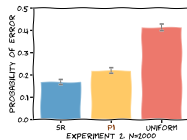
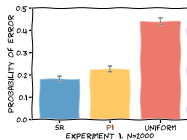
Solution: Set $T_i = T \frac{\Delta_1}{\Delta_i}$ as in the lower bound.

Corollary The complexity H_{P1} of P1 matches the complexity H_{BOB} from the lower bound of Theorem 4 of up to \log factors,

$$H_{P1} = \mathcal{O} \left(H_{\text{BOB}}^2 \log K \right).$$



Experimental setup	H_{SR}	H_{BOB}	H_{UNIF}
1. 1 group of bad arms	2000	2000	2000
2. 2 groups of bad arms	1389	2083	3125
3. Geometric progression	5540	5540	11080
4. 3 groups of bad arms	400	500	938
5. Arithmetic progression	3200	3200	24000
6. 2 good, many bad arms	5000	7692	50000
7. 3 groups of bad arms	4082	5714	12000
8. Square-root gaps	3200	22627	160000



Empirical behavior above mimics the behavior of the complexities in the table.

Thank you!

- Robin Allesiardo, Raphael Féraud, and Odalric-Ambrym Maillard. The non-stationary stochastic multi-armed bandit problem. *International Journal of Data Science and Analytics*, 2017.
- Jason Altschuler, Victor-Emmanuel Brunel, and Alan Malek. Best-arm identification for contaminated bandits. *arXiv preprint arXiv:1802.09514*, 2018.
- Jean-Yves Audibert, Sebastien Bubeck, and Rémi Munos. Best-arm identification in multi-armed bandits. In *Conference on Learning Theory*, 2010.
- Peter Auer and Chao-Kai Chiang. An algorithm with nearly optimal pseudo-regret for both stochastic and adversarial bandits. In *Conference on Learning Theory 2016*
- Peter Auer, Nicolo Cesa-Bianchi, Yoav Freund, and Robert E. Schapire. The nonstochastic multiarmed bandit problem. *SIAM Journal on Computing*, 32(1), 2002.
- Sebastien Bubeck and Aleksandrs Slivkins. The best of both worlds: stochastic and adversarial bandits. In *Conference on Learning Theory*, 2012.
- Alexandra Carpentier and Andrea Locatelli. Tight (lower) bounds for the fixed budget best-arm identification bandit problem. In *Conference on Learning Theory*, 2016.
- Kevin Jamieson and Amee Talwalkar. Non-stochastic best-arm identification and hyperparameter optimization. In *International Conference on Artificial Intelligence and Statistics*, 2016.
- Oded Maron and Andrew Moore. Hoeffding Races: Accelerating model-selection search for classification and function approximation. In *Neural Information Processing Systems*, 1993.
- Yevgeny Seldin and Aleksandrs Slivkins. One practical algorithm for both stochastic and adversarial bandits. In *International Conference on Machine Learning*, 2014.
- Julian Zimmert, Yevgeny Seldin. Tsallis-INF: An Optimal Algorithm for Stochastic and Adversarial Bandits, *JMLR* 2021.