

Multi-Bandit Best Arm Identification

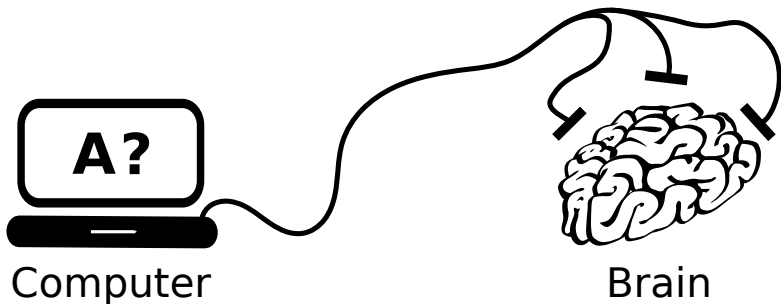
V. Gabillon, M. Ghavamzadeh, A. Lazaric & S. Bubeck



Sequel Group Meeting, 21 octobre, 2011.

An Example

A Brain-Computer Interface (BCI) problem



- A **computer** tries to guess the letter chosen by a **user**.
- It asks a **series** of **n** questions to the user.
- The answers are the recorded **noisy** brain activity signals.

Grid of letter shown to the user:

B	N	V	O
U	Z	P	M
K	H	G	A
W	F	L	T

Two types of questions:

- 1.) Is your letter in this **row**? 2.) Is your letter in this **column**?

B	N	V	O
U	Z	P	M
K	H	G	A
W	F	L	T

B	N	V	O
U	Z	P	M
K	H	G	A
W	F	L	T

- Which questions should be asked to best identify the letter?

A new bandit problem...

- Finding the right **row** is a multi-armed bandit problem.
- Finding the right **column** is a multi-armed bandit problem.

New: Solving **several** bandit problems at the same time.

The loss:

$$\begin{cases} 1 & \text{if **all** the bandits are solved,} \\ 0 & \text{otherwise.} \end{cases}$$

... based on pure exploration problems.

Known algorithms in the single bandit case.

Outline

- 1 Single-bandit best arm identification
- 2 Multi-bandit best arm identification
- 3 Experiments

Best arm identification



Identifying the best arm:

- **K** arms: $\nu_1, \nu_2, \dots, \nu_K$ distributions (arms) in $[0, b]$.
- At each time step t the player chooses an arm $I(t)$ to pull and receives a sample from $\nu_{I(t)}$
- At time \mathbf{n} the player has to recommend an arm **J(n)**.
- The player wins if he has recommended the distribution with highest mean k^* .

What is the strategy of allocation which minimizes the error?

$$\ell(n) = P(\text{error}) = P(J(n) \neq k^*)$$

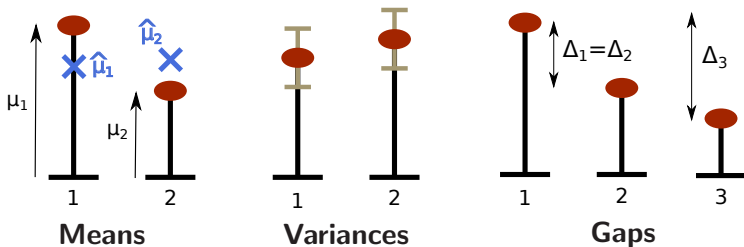
Notations

In arm k :

- **Number of pulls:** $\mathbf{T}_k(\mathbf{t})$ (at time t).
- **Mean-variance:** μ_k and σ_k^2 .
- **Estimated mean & variance:**

$$\hat{\mu}_k(t) = \frac{1}{T_k(t)} \sum_{s=1}^{T_k(t)} X_{k,s} \quad \text{and} \quad \hat{\sigma}_k^2(\mathbf{t}) = \frac{1}{T_k(t)-1} \sum_{s=1}^{T_k(t)} (X_{k,s} - \hat{\mu}_k(t))^2$$

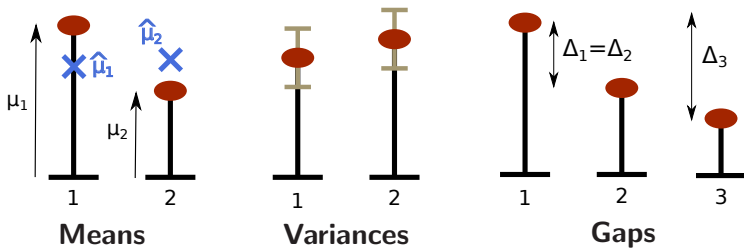
$X_{k,s}$ is the s^{th} sample from the arm.



Notations

In arm k :

- **Best arm:** $\mathbf{k}^* = \arg \max_k \mu_k$.
- **Gap:** $\Delta_k = |\max_{j \neq k} \mu_j - \mu_k|$.



Simple Regret vs. Cumulative regret

Cumulative regret

The goal is to minimize:

$$R(n) = n\mu_{k^*} - \sum_{t=1}^n X_{I(t)}$$

Dilemma: Spot the best arm (exploration) & exploit it (exploitation).

Simple regret

$$r(n) = \mu_{k^*} - \mu_{J(n)}$$

Pure exploration: No more exploration vs. exploitation trade-off.

$$\mathbb{E}r(n) \leq b \times \ell(n),$$

Budgeted Learning

The total number of pulls n is fixed and known by the player.

Here: we study how to optimize the quality of the recommendation with a fixed budget

vs.

Trying to optimize the budget in order to achieve a given confidence level (Hoeffding Races, **Even-Dar et.al. (2002)** ...).

Uniform allocation

Divide the budget equally over the arm:

Pull every arm n/K times.

Output the arm with highest estimated mean:

Output $J(n) \in \arg \max_k \hat{\mu}_k(n)$.

What is the **probability of error** of the uniform allocation?

Distinguishing two arms pulled equally:

with $\mu_1 > \mu_2$ and $\Delta = \mu_1 - \mu_2$

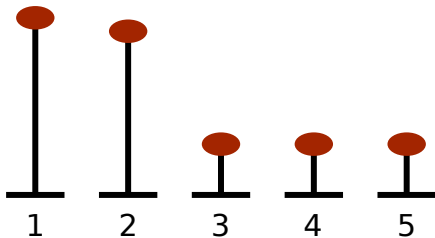
$$\ell(n) := \mathbb{P}(\text{error}) = \mathbb{P}(\hat{\mu}_1(n) < \hat{\mu}_2(n)) \leq \exp\left(-\frac{n}{2}\left(\frac{\Delta}{b}\right)^2\right)$$

For a fixed accuracy, the arms have to be pulled of order $\left(\frac{b}{\Delta}\right)^2$.

Error of Uniform allocation:

$$\begin{aligned}\ell_{Unif}(n) &:= \mathbb{P}(J(n) \neq k^*) \leq \sum_{k \neq k^*} \exp\left(-\frac{n}{K}\left(\frac{\Delta_k}{b}\right)^2\right) \\ &\leq K \exp\left(-\frac{n}{K}\left(\frac{\min_k \Delta_k}{b}\right)^2\right)\end{aligned}$$

Smarter than uniform allocation?



- We should not **waste** too much pulls on arms 3,4,5.
- We should focus on 1 and 2 to discriminate them.

UCB-Exploration (UCB-E)

Audibert, Bubeck & Munos (2010)

Parameter: exploration parameter a , range b , number of pulls n .

Index of arm k :

$$B_k(t) = \hat{\mu}_k(t-1) + b\sqrt{\frac{a}{T_k(t-1)}}$$

for each round $t = 1, 2, \dots, n$ **do**

Explore the action with the highest index:

 Draw $I(t) \in \arg \max_{k \in \{1, \dots, K\}} B_k(t)$

end for

Output the arm with highest estimated mean:

Output $J(n) \in \arg \max_k \hat{\mu}_k(n)$.

UCB-Exploration (UCB-E)

Audibert, Bubeck & Munos (2010)

Theorem: Bound on the probability of error

If $a = \frac{n}{H}$,

$$\ell_{UCB-E}(n) \leq 2nK \exp\left(-O\left(\frac{n}{H}\right)\right)$$

and $T_k(n) \approx O\left(\frac{b}{\Delta_k^2}\right) \quad \forall k$.

UCB-E pulls every arms according to its complexity.

Complexity H:

$$H = \sum_k \left(\frac{b}{\Delta_k}\right)^2$$

The global complexity is the **sum** of the complexity of the arms.

Comparing Uniform vs UCB-E

Probability of errors:

$$\begin{aligned} \ell_{Uniform}(n) & \quad \text{vs} \quad \ell_{UCB-E}(n) \\ \exp\left(-\frac{n}{K}\left(\frac{\min_k \Delta_k}{b}\right)^2\right) & \quad \exp\left(-\frac{n}{H}\right) \quad \# \text{ Comparing the bounds} \\ K \max_k \frac{1}{\Delta_k^2} & \quad \geq \quad \sum_k \frac{1}{\Delta_k^2} \end{aligned}$$

- UCB-E is expected to improve upon Uniform when the gaps are different.
- This was supported by numerical simulations.

Comparing Uniform vs UCB-E

Probability of errors:

$$\begin{aligned} \ell_{Uniform}(n) & \quad \text{vs} \quad \ell_{UCB-E}(n) \\ \exp\left(-\frac{n}{K}\left(\frac{\min_k \Delta_k}{b}\right)^2\right) & \leq \exp\left(-\frac{n}{H}\right) \quad \# \text{ Comparing the bounds} \\ K \max_k \frac{1}{\Delta_k^2} & \geq \sum_k \frac{1}{\Delta_k^2} \end{aligned}$$

- UCB-E is expected to improve upon Uniform when the gaps are different.
- This was supported by numerical simulations.

B	N	V	O
U	Z	P	M
K	H	G	A
W	F	L	T



B	N	V	O
U	Z	P	M
K	H	G	A
W	F	L	T



B	N	V	O
U	Z	P	M
K	H	G	A
W	F	L	T

New Goal: To solve **several** problems at the same time.

Outline

- ① Single-bandit best arm identification
- ② Multi-bandit best arm identification
- ③ Experiments

Multi-bandit pure exploration problem



- **M** bandits
- **K** arms in each bandit.
- At each time step t the player chooses a **bandit-arm pair** $I(t)$ to pull and receives a sample from $\nu_{I(t)}$
- At time n the player **recommends an arm for each bandit**.
- The player wins if he has recommended the distribution with highest mean **in all the bandits**.

What is the strategy of allocation which minimizes the error?

Issues

- Is UCB-E a solution for this problem?
- What type of algorithm for bandit problem : A two stage algorithm?
Allocation over bandits + allocation over arms.
- What is the complexity of this problem?

New notations

In bandit m and arm k :

- **Number of pulls:** $\mathbf{T}_{mk}(\mathbf{t})$ (at time t).
- **Mean-variance:** μ_{mk} and σ_{mk}^2 .
- **Estimated mean & variance:**

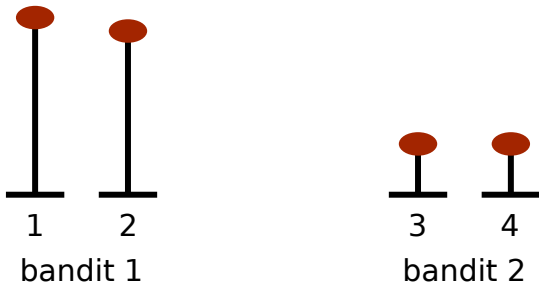
$$\hat{\mu}_{mk}(t) = \frac{1}{T_{mk}(t)} \sum_{s=1}^{T_{mk}(t)} X_{mk,s}$$

$$\hat{\sigma}_{mk}^2(\mathbf{t}) = \frac{1}{T_{mk}(t)-1} \sum_{s=1}^{T_{mk}(t)} (X_{mk,s} - \hat{\mu}_{mk}(t))^2$$

$X_{mk,s}$ is the s^{th} sample from arm.

- **Best arms:** $\mathbf{k}_m^* = \arg \max_k \mu_{mk}$ and $\hat{\mathbf{k}}_m^* = \arg \max_k \hat{\mu}_{mk}$.
- **Gap** of arm k : $\Delta_{mk} = \left| \max_{j \neq k} \mu_{mj} - \mu_{mk} \right|$.

UCB-E is not a good solution



On this problem,

- UCB-E will focus on the arms of bandits 1.
- UCB-E will **fail** to solve the two bandit problems.

Uniform strategies

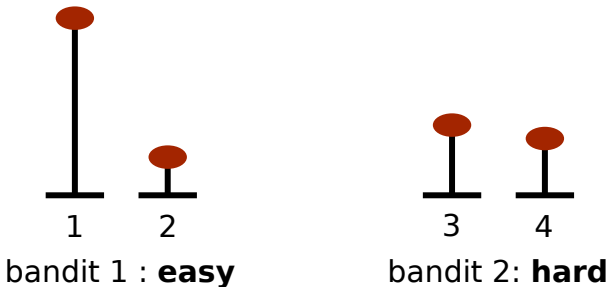
- **Uniform:** Pulls every arm n/MK .
- **Uniform + UCB-E:**
 - ① The number of pulls for each bandit is n/M .
 - ② Play UCB-E in each bandit.

Bounds on the errors: $\ell(n) := \mathbb{P}(\exists m : J_m(n) \neq k_m^*)$

$$\ell_{\text{Uniform}}(n) \leq \exp\left(-O\left(\frac{n}{MK \max_{m,k} H_{mk}}\right)\right), \quad H_{mk} = \left(\frac{b}{\Delta_{mk}}\right)^2,$$

$$\ell_{\text{Uniform+UCBE}}(n) \leq \exp\left(-O\left(\frac{n}{M \max_m H_m}\right)\right), \quad H_m = \sum_k H_{mk}.$$

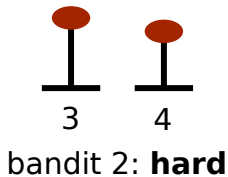
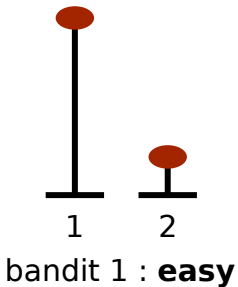
Allocation over bandits



On this problem,

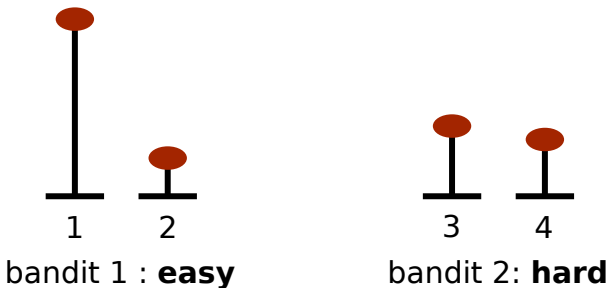
- More pulls should be allocated on bandit 2.
- Uniform allocations are not satisfactory.

Allocation other states



Allocating according to the gaps?

Allocation other states



Gaps:

- First idea: $\hat{\Delta}_{mk} = \max_j \hat{\mu}_{mj} - \hat{\mu}_{mk}$ **BUT:** $\hat{\Delta}_{mk^*_m} = 0$
- Second idea: $\hat{\Delta}_{mk} = |\max_{j \neq k} \hat{\mu}_{mj} - \hat{\mu}_{mk}|$

From UCB-E to GapE

Parameter: exploration parameter a , range b , number of pulls n .

Index of arm k :

$$B_k(t) = \hat{\mu}_k(t-1) + b\sqrt{\frac{a}{T_k(t-1)}}$$

for each round $t = 1, 2, \dots, n$ **do**

Explore the action with the highest index:

 Draw $I(t) \in \arg \max_{k \in \{1, \dots, K\}} B_k(t)$

end for

Output the arm with highest estimated mean:

Output $J(n) \in \arg \max_k \hat{\mu}_k(n)$.

From UCB-E to GapE

Parameter: exploration parameter a , range b , number of pulls n .

Index of arm k :

$$B_k(t) = \hat{\mu}_k(t-1) - \max_k \hat{\mu}_k(t-1) + b \sqrt{\frac{a}{T_k(t-1)}}$$

for each round $t = 1, 2, \dots, n$ **do**

Explore the action with the highest index:

 Draw $I(t) \in \arg \max_{k \in \{1, \dots, K\}} B_k(t)$

end for

Output the arm with highest estimated mean:

Output $J(n) \in \arg \max_k \hat{\mu}_k(n)$.

From UCB-E to GapE

Parameter: exploration parameter a , range b , number of pulls n .

Index of arm k :

$$B_k(t) = -\hat{\Delta}_k(t-1) + b\sqrt{\frac{a}{T_k(t-1)}}$$

for each round $t = 1, 2, \dots, n$ **do**

Explore the action with the highest index:

 Draw $I(t) \in \arg \max_{k \in \{1, \dots, K\}} B_k(t)$

end for

Output the arm with highest estimated mean:

Output $J(n) \in \arg \max_k \hat{\mu}_k(n)$.

From UCB-E to GapE

Parameter: exploration parameter a , range b , number of pulls n .

Index of arm k :

$$B_k(t) = -\hat{\Delta}_{mk}(t-1) + b\sqrt{\frac{a}{T_{mk}(t-1)}}$$

for each round $t = 1, 2, \dots, n$ **do**

Explore the action with the highest index:

 Draw $I(t) \in \arg \max_{k \in \{1, \dots, K\} \& m \in \{1, \dots, M\}} B_{mk}(t)$

end for

Output the arms with highest estimated mean:

Output $J_m(n) \in \arg \max_k \hat{\mu}_{mk}(n)$.

Gap-based exploration (GapE)

Theorem: Bound on the probability of error

If $a = \frac{n}{H}$,

$$\ell_{\text{GapE}}(n) \leq 2nMK \exp\left(-O\left(\frac{n}{H}\right)\right)$$

and $T_{mk}(n) \approx O\left(\frac{b}{\Delta_{mk}^2}\right) \quad \forall k, \forall m$.

GapE pulls every arms according to its complexity.

Complexity H:

$$H = \sum_m \sum_k \left(\frac{b}{\Delta_{mk}}\right)^2$$

The global complexity is the **sum** of the complexity of the arms.

Comparing Uniform allocations vs GapE

Probability of errors:

$$\begin{aligned} \ell_{Uniform}(n) & \quad \text{vs} \quad \ell_{Unif+UCB-E}(n) & \quad \text{vs} \quad \ell_{GapE}(n) \\ \exp\left(-\frac{n}{MK \max_{m,k} H_{mk}}\right) & \quad \exp\left(-\frac{n}{M \max_m H_m}\right) & \quad \exp\left(-\frac{n}{H}\right) \\ MK \max_{m,k} \frac{1}{\Delta_{mk}^2} & \quad \geq \quad M \max_m \sum_k \frac{1}{\Delta_{mk}^2} & \quad \geq \quad \sum_m \sum_k \frac{1}{\Delta_{mk}^2} \end{aligned}$$

- GapE is expected to improve upon Uniform allocations when the gaps are different.
- This will be supported by numerical simulations.

Comparing Uniform allocations vs GapE

Probability of errors:

$$\begin{aligned} \ell_{Uniform}(n) & \text{ vs } \ell_{Unif+UCB-E}(n) & \text{ vs } & \ell_{GapE}(n) \\ \exp\left(-\frac{n}{MK \max_{m,k} H_{mk}}\right) & \leq \exp\left(-\frac{n}{M \max_m H_m}\right) & \leq & \exp\left(-\frac{n}{H}\right) \\ MK \max_{m,k} \frac{1}{\Delta_{mk}^2} & \geq M \max_m \sum_k \frac{1}{\Delta_{mk}^2} & \geq & \sum_m \sum_k \frac{1}{\Delta_{mk}^2} \end{aligned}$$

- GapE is expected to improve upon Uniform allocations when the gaps are different.
- This will be supported by numerical simulations.

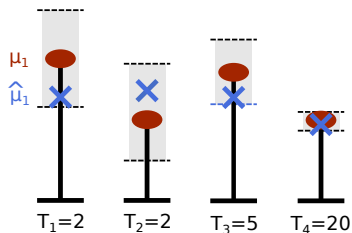
Sketch of the Proof:

Difficulty: Having upper bound on the estimated mean of the estimated best arm.

- 1 Characterize an event on which the algorithm succeeds and compute its probability.
- 2 On this event, give a condition on the number of pulls so that the algorithm succeeds.
- 3 Prove, on that event, an induction formula for each time t showing the dependence between the pulls of each arm and their gaps.
- 4 With the induction formula, compute the number of pulls at time n and show that it satisfies the requirement to succeed.

1.) The event where GapE succeeds

$$\mathcal{E} = \left\{ \forall m \in \{1, \dots, M\}, \forall k \in \{1, \dots, K\}, \forall t \in \{1, \dots, n\}, \right. \\ \left. |\hat{\mu}_{mk}(t) - \mu_{mk}| < bc \sqrt{\frac{a}{T_{mk}(t)}} \right\}.$$

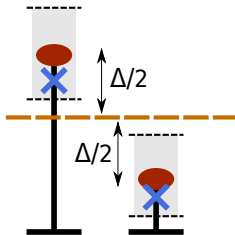


From Chernoff-Hoeffding's inequality and a union bound, we have $\mathbb{P}(\xi) \geq 1 - 2MKn \exp(-2ac^2)$.

2.) Condition on the number of pulls

Success in arm k of bandit m if,

$$\begin{aligned} J_m(n) = k_m^* &\iff \forall k, \hat{\mu}_{mk}(n) \leq \hat{\mu}_{mk_m^*}(n) \\ &\iff \forall k, bc\sqrt{a/T_{mk}(n)} \leq \Delta_{mk}/2 \\ &\iff \forall k, T_{mk}(n) \geq \frac{4ab^2c^2}{\Delta_{mk}^2} \end{aligned}$$



3.) Induction formula

For any bandits (m, q) and arms (k, j) , and for any $t \geq MK$,

$$-\Delta_{mk} + (1+d)b\sqrt{\frac{a}{\max(T_{mk}(t) - 1, 1)}} \geq -\Delta_{qj} + (1-d)b\sqrt{\frac{a}{T_{qj}(t)}},$$

It shows the dependence between the number of pulls of **all** the pair of arms from **all** the bandit.

4.) At time n

Proof by contradiction,

- Let us assume $\exists(m, k) : T_{mk}(n) \leq \frac{ab^2(1-d)^2}{\Delta_{mk}^2}$

$$\Rightarrow -\Delta_{mk} + (1-d)b\sqrt{\frac{a}{T_{mk}(n)}} > 0$$

$$\Rightarrow -\Delta_{qj} + (1+d)b\sqrt{\frac{a}{T_{qj}(n)-1}} > 0, \quad \forall(q, j) \quad \text{induction formula}$$

$$\Rightarrow T_{qj}(n) < \frac{ab^2(1+d)^2}{\Delta_{qj}^2} + 1$$

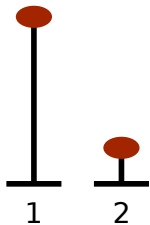
- And we have

$$\sum_{q,j} T_{qj}(n) = n \quad \Rightarrow \quad n - MK < ab^2(1+d)^2 \sum_{q,j} \frac{1}{\Delta_{qj}^2}$$

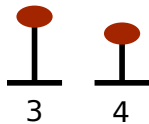
- So, if we select a such that $n - MK \geq ab^2(1+d)^2 \sum_{q,j} \frac{1}{\Delta_{qj}^2}$, we contradict the first assumption and this concludes the proof.

Extensions

The use of variance

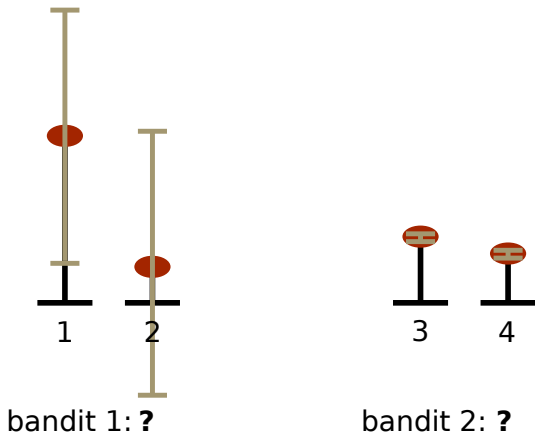


bandit 1: **easy**

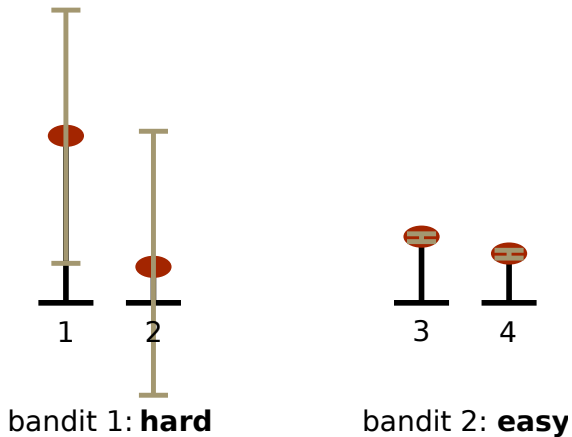


bandit 2: **hard**

The use of variance



The use of variance



GapE-Variance (GapE-V)

Index of arm k :

$$B_k(t) = -\hat{\Delta}_{mk}(t-1) + \sqrt{\frac{2a \hat{\sigma}_{mk}^2(t-1)}{T_{mk}(t-1)}} + \frac{7ab}{3(T_{mk}(t-1)-1)}$$

Theorem: Bound on the probability of error

$$\text{If } a = \frac{n}{H^\sigma}, \quad \ell_{\text{GapE-V}}(n) \leq 6nMK \exp\left(-O\left(\frac{n}{H^\sigma}\right)\right)$$

Variance-complexity H^σ :

$$H^\sigma = \sum_m \sum_k \frac{(\sigma_{mk} + \sqrt{\sigma_{mk}^2 + (16/3)b\Delta_{mk}})^2}{\Delta_{mk}^2}.$$

H^σ is smaller than H when the variance is small compared to b .

Adaptive version of GapE algorithms

Issues

- a should be tuned according to the complexities H and H^σ
- H and H^σ are rarely known in advance.

Online estimation of the complexities:

$$\text{UCB}_{\Delta_i}(t) = \hat{\Delta}_i(t-1) + \sqrt{\frac{1}{2T_i(t-1)}}$$
$$\hat{H}(t) = \sum_{m,k} \frac{b^2}{\text{UCB}_{\Delta_i}(t)^2},$$

Outline

- ① Single-bandit best arm identification
- ② Multi-bandit best arm identification
- ③ Experiments

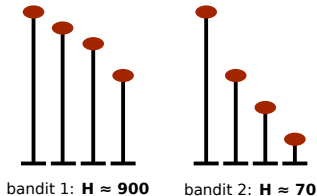
Parameters

- Theoretically, a should be proportional to $\frac{n}{H}$.
- In the experiments $a = \eta \frac{n}{H}$, η is a parameter to tune.
- η is suppose to correct the inaccuracy of the constants in the analysis
- the range of its nearly-optimal values should be constant across different problems

Comparison: GapE vs uniform strategies

Problem 1: $n = 700$

Results



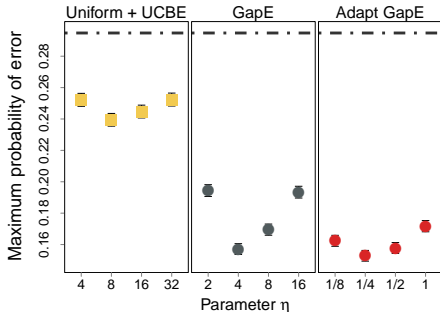
Allocation in percent in GapE:

[81]

[19]

[37, 36, 20, 7]

[36, 34, 17, 10]

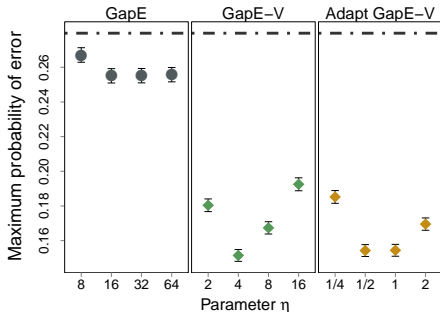
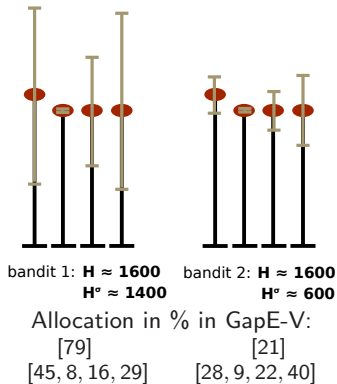


- GapE allocates according to the complexities.
- GapE improves upon Uniform and Uniform+UCB-E.

Comparison: GapE vs GapE-V

Problem 2: $n = 1000$

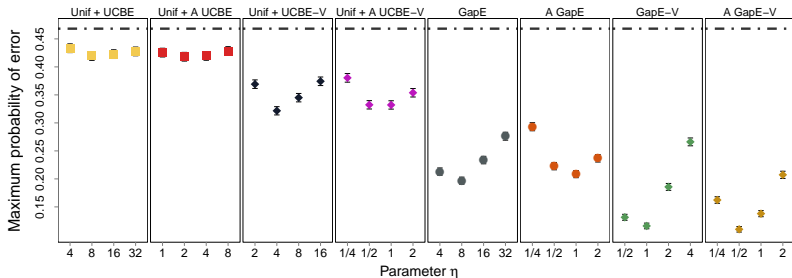
Results



- GapE-V allocates according to the variance-complexities.
- GapE-V improves upon GapE.

Final comparison

Problem 3 with strongly heterogeneous bandits & arms $M = K = 4$.



- 1st: GapE-V, 2nd: GapE, 3rd: Unif+UCBEV, 4th: Unif+UCBE
- Across the problems, in GapE-V, $\eta \in [1 - 4]$ is a near optimal choice.
- The **adaptive** versions of the algorithms have almost the same performance than their "**oracle**" counterpart.

Conclusions

- **New problem:** of best arm identification in a multi-bandit multi-armed setting.
- **New algorithm:** GapE, with an upper-bound for its probability of error showing its *potential* advantage over uniform strategies.
- **Extensions:**
 - Considering the variance
 - Adaptive version which estimates the complexity of the problem online.
- The **numerical simulations** show that GapE and GapE-V and their adaptive counterparts outperform other allocation strategies.

Perspectives

Application in RL?

In rollout allocation for classification-based policy iteration the goal is to identify the greedy action (*arm*) in each of the states (*bandit*) in a training set.

Dynamic vs Static strategy.

Gap-E is empirically better than a static strategy which would allocated according to the complexity. But they have the same theoretical guaranties. Can we improve something?

Other problems...

GapE-KL, data-dependent complexity, lower bound, pure exploration in structured set, simple regret, sum loss vs max loss, analysis of adaptive version...

Thank you!