

BEST ARM IDENTIFICATION: A UNIFIED APPROACH TO FIXED BUDGET & FIXED CONFIDENCE

VICTOR GABILLON, MOHAMMAD GHAVAMZADEH, ALESSANDRO LAZARIC

ABSTRACT

The problem of identifying the best arm(s) in the stochastic multi-armed bandit setting has been studied in the literature from two different perspectives: *fixed budget* and *fixed confidence*.

We propose a unifying approach that leads to a meta-algorithm with **1)** a common structure and **2)** similar theoretical analysis for these two settings.

PROBLEM FORMULATION

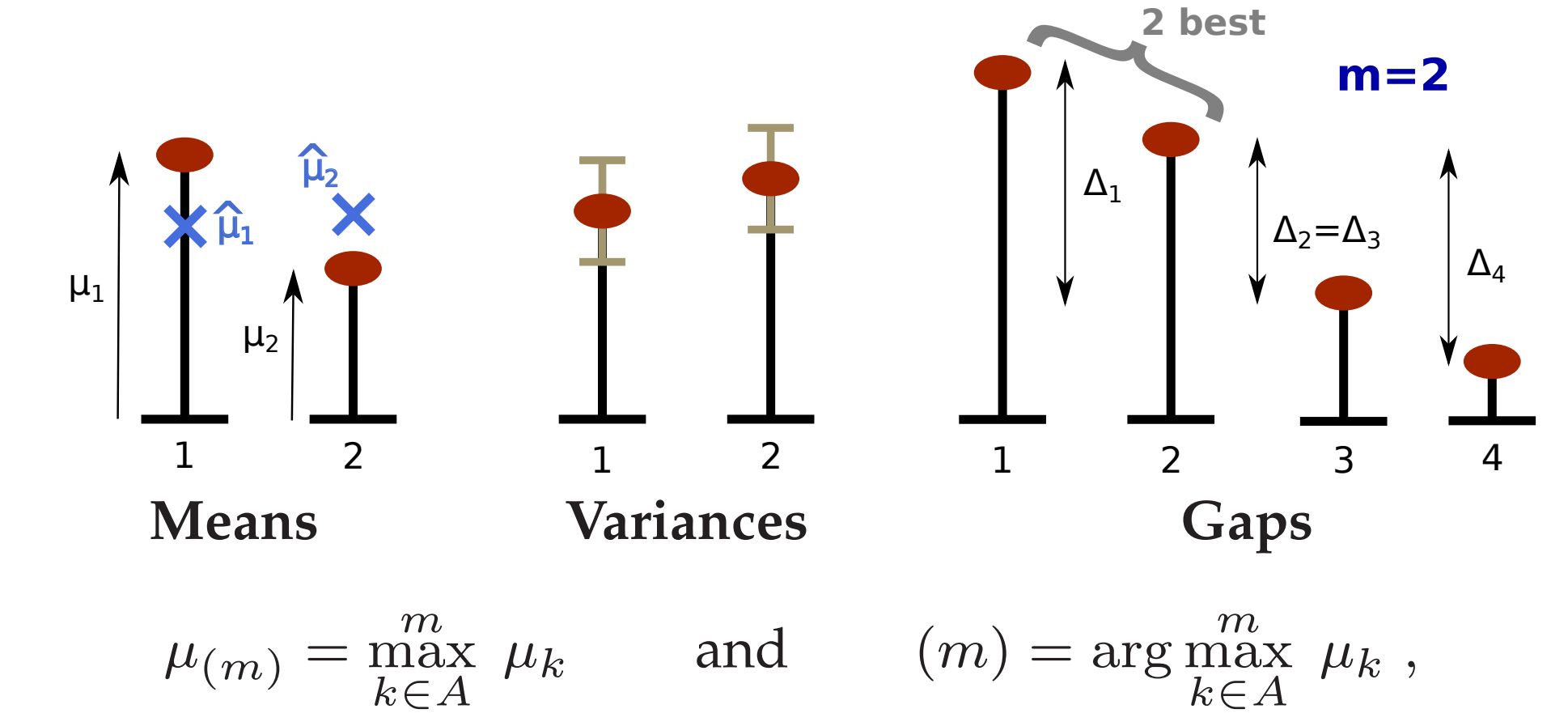


K arms: $\nu_1, \nu_2, \dots, \nu_K$ distributions (arms) in $[0, b]$.
At each time step t the player chooses an arm $I(t)$ to pull and receives a sample from $\nu_{I(t)}$.
In arm k :
Number of pulls: $T_k(t)$.
Mean-variance: μ_k and σ_k^2 .
Gap: $\Delta_k = |\max_{j \neq k} \mu_j - \mu_k|$.
 $X_{k,s}$: s^{th} sample from the arm.

Estimated mean & variance:

$$\hat{\mu}_k(t) = \frac{1}{T_k(t)} \sum_{s=1}^{T_k(t)} X_{k,s}$$

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



Given an accuracy ϵ and a number of arms m , we say that an arm k is (ϵ, m) -optimal if $\mu_k \geq \mu_{(m)} - \epsilon$.

GOAL: Identify a set S of m (ϵ, m) -optimal arms.

Regret:

The arm simple regret: $r_k = \mu_{(m)} - \mu_k$,

The simple regret of S , a set of m arms:

$$r_S = \max_{k \in S} r_k = \mu_{(m)} - \min_{k \in S} \mu_k.$$

Return Strategy: $\Omega(t)$ is the set of m arms returned after t rounds with simple regret: $r_{\Omega(t)}$

THE TWO SETTINGS

Fixed budget: The objective is to return a set of m (ϵ, m) -optimal arms with the largest possible confidence using a fixed budget of n rounds. The performance is measured by the probability $\tilde{\delta}$ of not meeting the (ϵ, m) requirement,

$$\tilde{\delta} = \mathbb{P}[r_{\Omega(n)} \geq \epsilon].$$

Fixed confidence: The goal is to stop as soon as possible and return a set of m (ϵ, m) -optimal arms with a fixed confidence. Let \tilde{n} be the time when the algorithm stops. Given a confidence level δ , the forecaster has to guarantee that

$$\mathbb{P}[r_{\Omega(\tilde{n})} \geq \epsilon] \leq \delta.$$

The performance of the forecaster is then measured by the number of rounds \tilde{n} either in expectation or high probability.

OPEN PROBLEM

Can we show an equivalence between the fixed budget setting and the fixed confidence setting when we don't assume the knowledge of the complexity in the fixed budget problem?

UNIFIED GAP-BASED EXPLORATION ALGORITHM

The unified gap-based exploration (UGapE) meta-algorithm can be implemented in the fixed-budget (UGapEb) and fixed-confidence (UGapEc) settings.

Definition of Indices:

For arm $k \in A$, $B_k(t) = \max_{i \neq k} U_i(t) - L_k(t)$

where $U_k(t) = \hat{\mu}_k(t-1) + \beta_k(t-1)$

$L_k(t) = \hat{\mu}_k(t-1) - \beta_k(t-1)$.

For a set S , $B_S(t) = \max_{i \in S} B_i(t)$

h. p. upper bound on the simple regret r_k

h. p. upper bound on the mean

h. p. lower bound on the mean

h. p. upper bound on the simple regret r_S .

$\beta_k(t-1)$ is a confidence interval. From the Chernoff-Hoeffding bound:

$$\text{UGapEb: } \beta_k(t) = b \sqrt{\frac{a}{T_k(t)}}$$

$$\text{UGapEc: } \beta_k(t) = b \sqrt{\frac{c \log \frac{4K(t-1)^3}{\delta}}{T_k(t)}}$$

Fixed-Budget: UGapEb(ϵ, n, a)

Parameters: accuracy ϵ , budget n , exploration parameter a

for $t = K + 1, \dots, n$ do

SELECT-ARM (t)

end for

Return $\Omega(n) = \arg \min_{J(t)} B_{J(t)}(t)$

Fixed-Confidence: UGapEc(ϵ, δ, c)

Parameters: accuracy ϵ , confidence level δ , exploration parameter c

while $B_{J(t)}(t) \geq \epsilon$ do

SELECT-ARM (t)

$t \leftarrow t + 1$

end while

Return $\Omega(t) = J(t)$

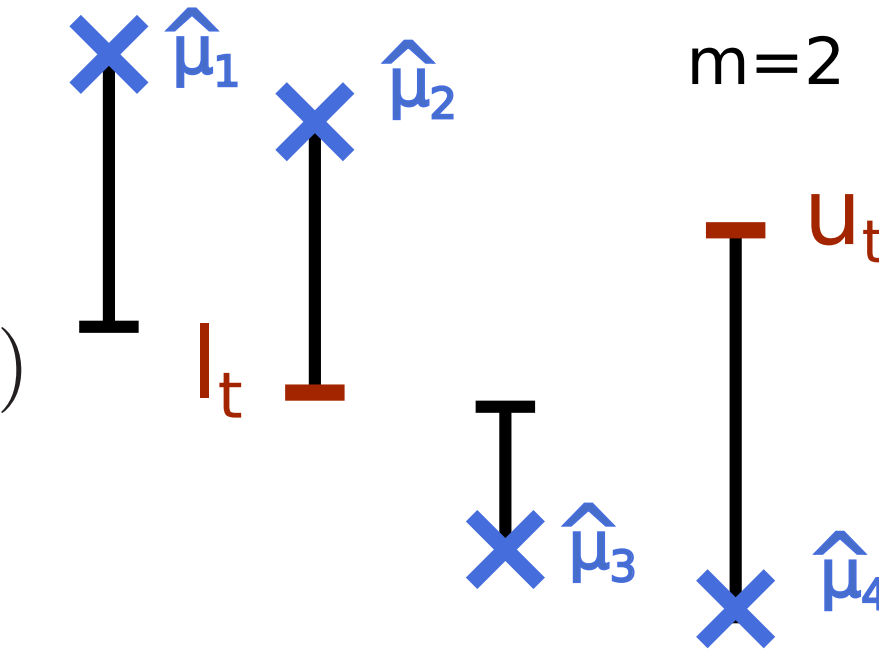
Arm selection: UGapEc & GapEb are based on the same sample routine.

$J(t)$: the set of m estimated best arms

$$J(t) = \arg \min_{k \in A}^{1..m} B_k(t).$$

$$u_t = \arg \max_{j \notin J(t)} U_j(t)$$

$$l_t = \arg \min_{i \in J(t)} L_i(t)$$



SELECT-ARM (t)

Compute $B_k(t)$ for each arm $k \in A$

Identify the set of m arms $J(t) \in \arg \min_{k \in A}^{1..m} B_k(t)$

Pull the arm $I(t) = \arg \max_{k \in \{l_t, u_t\}} \beta_k(t-1)$

Observe $X_{I(t)}(T_{I(t)}(t-1) + 1) \sim \nu_{I(t)}$

Update $\hat{\mu}_{I(t)}(t)$ and $T_{I(t)}(t)$

ANALYSIS

$B_{J(t)}(t)$ is upper bounded by the quantities of interest (the gaps) independently of the setting

Lemma 1 With high probability, if arm k is pulled at time t , $B_{J(t)}(t) \leq \min(0, -\Delta_k + 2\beta_k(t-1)) + 2\beta_k(t-1)$.

From this Lemma, we can prove an upper-bound on the simple-regret of both UGapEc & UGapEb:

Regret Bound for the Fixed-Budget Setting

Theorem 1 If we run UGapEb with parameter

$$0 < a \leq \frac{n-K}{4H_\epsilon}, \text{ its simple regret } r_{\Omega(n)} \text{ satisfies}$$

$$\tilde{\delta} = \mathbb{P}(r_{\Omega(n)} \geq \epsilon) \leq 2Kn \exp(-2a),$$

and this probability is minimized for $a = \frac{n-K}{4H_\epsilon}$.

Regret Bound for the Fixed-Confidence Setting
Theorem 2

$$\mathbb{P}(r_{\Omega(\tilde{n}+1)} \leq \epsilon \wedge \tilde{n} \leq N) \geq 1 - \delta,$$

where $N = K + O(H_\epsilon \log \frac{H_\epsilon}{\delta})$ and c has been set to its optimal value $1/2$.

H_ϵ captures the intrinsic difficulty of the (ϵ, m) -best arm(s) identification problem independently from the specific setting considered.

$$H_\epsilon = \sum_{i=1}^K \frac{b^2}{\max(\frac{\Delta_i + \epsilon}{2}, \epsilon)^2}.$$

Targeted number of pulls in each setting for arm i

$$\text{FC: } T_i(\tilde{n}) = \frac{b^2}{\Delta_i^2}, \quad \text{FB: } T_i(n) = \frac{b^2}{\Delta_i^2} * \frac{n}{H}$$

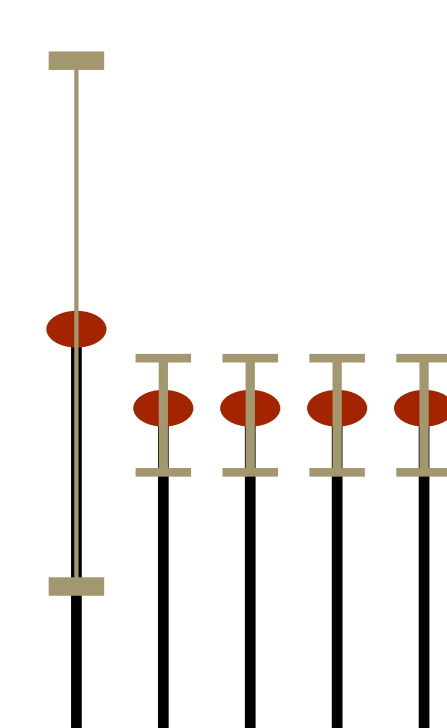
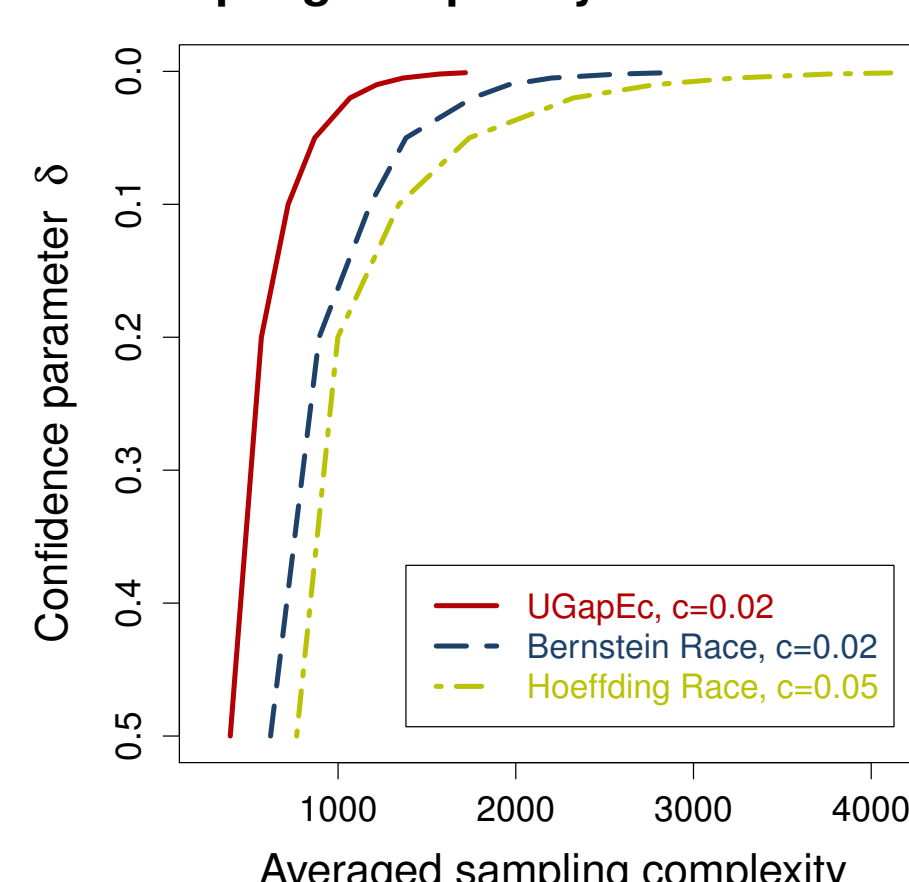
So, if we assume the knowledge of the complexity H the two problems have a very similar targeted allocation where the number of pulls of each arm does not depend on the gap of the other arms.

EXTENSION TO VARIANCE IN THE FIXED CONFIDENCE SETTING

UGapE-Variance (UGapE-V). UGapE and its analysis can be naturally and easily extended to take into account the variances of the arms in both cases. We focus here on the Fixed Confidence setting:

$$\text{UGapEc-V: } \beta_k(t) = \sqrt{\frac{2c \log \frac{Kt^3}{\delta} \hat{\sigma}_k^2(t)}{T_k(t)}} + \frac{(7/3)bc \log \frac{Kt^3}{\delta}}{T_k(t) - 1},$$

Sampling complexity w.r.t. confidence



For each algorithm, we only consider the results corresponding to the value of c for which the required confidence level δ is satisfied. We report the results for the value of c with the smallest sample complexity.

New definition of complexity,

$$H_\epsilon^\sigma = \sum_{i=1}^K \frac{(\sigma_i + \sqrt{(13/3)b\Delta_i})^2}{\max(\frac{\Delta_i + \epsilon}{2}, \epsilon)^2}.$$

Bernstein Race: $H_\epsilon^\sigma \approx (\sigma_{(1)}^2 + \sigma_i^2) / \Delta_i^2$, where $\sigma_{(1)}^2$ is the variance of the best arm. In case of $K = 2$, Bernstein Race allocate uniformly, while UGapE-V pulls the arms proportionally to their variances.