

---

# Generic Exploration and $K$ -armed Voting Bandits

---

Tanguy Urvoy  
Fabrice Clerot  
Raphael Féraud  
Sami Naamane

TANGUY.URVOY@ORANGE.COM  
FABRICE.CLEROT@ORANGE.COM  
RAPHAEL.FERAUD@ORANGE.COM  
SAMI.NAAMANE@ORANGE.COM

Orange-labs, 2 avenue Pierre Marzin, 22307 Lannion, FRANCE

## Abstract

We study a stochastic online learning scheme with partial feedback where the utility of decisions is only observable through an estimation of the environment parameters. We propose a generic pure-exploration algorithm, able to cope with various utility functions from multi-armed bandits settings to dueling bandits. The primary application of this setting is to offer a natural generalization of dueling bandits for situations where the environment parameters reflect the idiosyncratic preferences of a mixed crowd.

## 1. Introduction

The stochastic multi-armed bandits became popular as a stripped-down model of *exploration versus exploitation* balance in sequential decision problems. In its simplest formulation, we are facing a slot machine with several arms. The rewards of these arms are modeled by unknown but bounded and independent random variables. To maximize our long term reward, we would like to play an arm with maximal expected value but we need to explore efficiently all the arms in order to find it.

The cost of ignorance is traditionally expressed in term of *expected regret* : the expected difference of reward between a playing policy established with perfect knowledge of the environment parameters and a given "unaware" policy. Following Lai & Robbins (1985), several regret analysis have been proposed (see for instance Auer et al., 2002; Audibert et al., 2008; Auer & Ortner, 2011).

Another way to evaluate bandit algorithms is to con-

sider *pure-exploration* or PAC *sample complexity*: how to find a nearly-optimal arm with high confidence in a minimum of trials? For multi-armed bandits, several algorithms were already proposed and studied from this perspective in (Even-Dar et al., 2002; 2006). We can also reverse the PAC question to control the prediction accuracy after a fixed number of samples as in (Audibert et al., 2010; Bubeck et al., 2011).

When the number of trials is bounded by a known horizon, we can adopt an *explore then exploit* strategy to control the final regret with an efficient exploration algorithm (Yue et al., 2012). This is the perspective we adopt here.

### 1.1. Rigged bandits

As an introduction to our *generic exploration* setting we consider the *rigged bandits* problem which is a simplified model for click-fraud in online advertising.

In this variant of multi-armed bandits we know from a reliable source (the barmaid of the casino) that the  $m$  best arms of the slot machine have been rigged and only deliver counterfeit money. To maximize our gain, we want to design a sequence of experiments in order to determine and play the  $(m + 1)^{th}$  best arm as early as possible while avoiding the rigged arms.

The main characteristic of this problem lies in the absence of a direct utility feedback: to estimate our real income we need to know with enough confidence which arms were rigged. This problem also requires what we call a *generic exploration policy*: we do not want to design a new exploration algorithm for each possible fraud-detection criterion we have at hand (see section 5.1 for a formalized example).

### 1.2. Dueling bandits

The *dueling bandit* problem, introduced by Yue & Joachims (2009) to formalize online learning from pref-

erence feedback, shares the *indirect* or *parametric feedback* property with rigged bandits. The initial motivation to depart from the absolute-reward model came from information retrieval evaluation where the implicit feedback by means of click logs is strongly biased by the ranking itself. A solution was proposed by Joachims (2003) to circumvent this problem: by interleaving two ranking models, and checking where the user clicked, one obtains an unbiased – but pairwise – preference feedback. Further experiments were performed in (Chapelle et al., 2012).

The original definition of the dueling bandits problem (Yue et al., 2012; Yue & Joachims, 2011) was built upon strong assumptions about the preference matrix: existence of a *strict linear ordering*, *stochastic transitivity* and *stochastic triangular inequality* (see Yue & Joachims, 2011). An extension of this setting with restricted pairing was proposed by Di Castro et al. (2011), but this extension also assumes the preference matrix to be the byproduct of an inherent value for each arm.

In a situation where the preferences reflects the expression of a mixed crowd, there can be several inconsistencies or *voting paradoxes* which contradict these assumptions. The definition we propose here is more relaxed: we do not assume the existence of a perfect linear order, neither do we assume the existence of an inherent value of arms. We simply try to sample efficiently the preference matrix in order to propose a "best element" similar to the one we would choose with perfect knowledge of the crowd preferences.

Electing a "best element" or a "best linear ordering" from such a preference matrix is a tough but old and well-studied problem (see Charon & Hudry, 2010, for a survey), but the works about online and noisy declinations of this problem are scarce (see however Ravikumar et al., 1987; Feige et al., 1994, for related problems). If we change the election criterion according to the vast social choice theory (see Chevaleyre et al., 2007, for a survey), we can decline the dueling bandits in several unexplored flavors: for instance *Borda bandits*, *Copeland bandits*, *Slater bandits*, or *Kemeny bandits*.

### 1.3. Toward generic exploration algorithms

The traditional approach to deal with exotic sequential decision problems is to design tailor-made algorithms which handle simultaneously the exploration of the environment and the exploitation of its knowledge. One purpose of this article is to explore the possibility of a generic algorithm which automatically generates an efficient exploration policy for any given decision crite-

tion. We propose a generic algorithm with theoretical guarantees for the case of *parametric decision problems* and we evaluate its performances on relatively simple declinations of rigged and dueling bandits.

## 2. Main Problem Statement

Consider a stationary environment modeled by a vector of  $N$  unknown parameters  $\boldsymbol{\mu} = (\mu_1, \dots, \mu_N) \in [0, 1]^N$  (By convention, we write the vectors with bold faces). We have a noisy perception of  $\boldsymbol{\mu}$  modeled by a vector  $\mathbf{X}$  of  $N$  independent random variables  $X_i \in [0, 1]$  verifying  $\mathbb{E}[X_i] = \mu_i$  for each index  $i = 1, \dots, N$ . The only hint we may have about  $\boldsymbol{\mu}$  is a set of feasible environment configurations  $\mathcal{F} \subseteq [0, 1]^N$ .

Let  $\mathcal{D}$  be a set of decisions and  $U : \mathcal{D} \times \mathcal{F} \rightarrow \mathbb{R}^+$  be a given utility function. From this utility function we can derive a decision function  $f : \mathcal{F} \rightarrow \mathcal{D}$  which computes an optimal option  $f(\mathbf{x}) \in \arg \max_d U(d, \mathbf{x})$  for each feasible realization  $\mathbf{x}$  of the random vector  $\mathbf{X}$ .

The *parametric decision problem* consists in finding the best decision  $d^* := f(\boldsymbol{\mu})$  with a high probability after a minimum amount of samples of the environment parameters.

If we have a metric  $L : \mathcal{D} \times \mathcal{D} \rightarrow \mathbb{R}_+$  to compare decisions, we can also search an  $\varepsilon$ -approximation of  $d^*$  (i.e. a decision  $d \in \mathcal{D}$  such that  $L(d^*, d) \leq \varepsilon$ ).

We use a "budgeted" version of PAC learning:

**Definition 1.** *An algorithm is an  $(\varepsilon, \delta)$ -PAC algorithm with horizon  $T$  for the parametric decision problem if it outputs an  $\varepsilon$ -approximation of  $d^*$  with probability at least  $1 - \delta$  when it terminates with strictly less than  $T$  samples. We call exploration time the number of parameter samples required for termination.*

This definition extends PAC learning to finite horizons: to avoid confusion we use the term *exploration time* instead of *sample complexity* when  $T$  is finite. The exploration time at horizon  $T$  with  $\delta = 1/T$  provides an upper bound for the expected cumulative regret (see Yue et al., 2012, section 4).

The decision function may be the result of quite a complex algorithm, but in the problems we consider here,  $\mathcal{D}$  is finite and the decision function is partitioning the input space into single-decision areas. For these problems we will assume that the environment state  $\boldsymbol{\mu}$  falls outside of the decision frontiers. In other words, we will assume that there exists a neighborhood of  $\boldsymbol{\mu}$  where  $f$  is constant. In order to analyze the performance of our algorithm in the next section, we will need a finer description of this neighborhood. For instance, the binary decision function defined in Figure 1

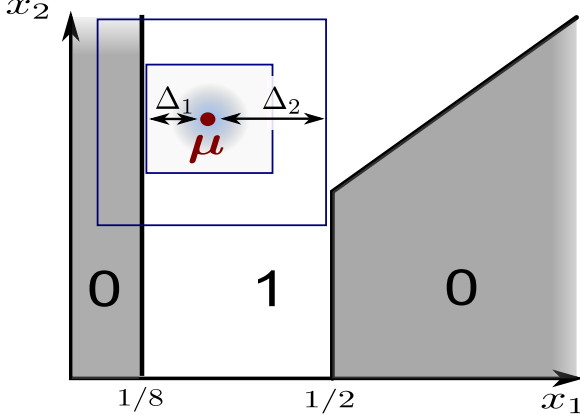


Figure 1. A simple example of binary decision function defined by  $f(x_1, x_2) = [x_1 > 1/8 \wedge (x_1 < 1/2 \vee x_2 > x_1)]$  (We use the square brackets  $[\cdot]$  to denote the characteristic function of a predicate). How accurately do we need to estimate  $\mu$  in order to take the good decision? The environment state  $\mu$  is at distance  $\Delta_1$  from the nearest decision frontier but we can stop exploring  $x_2$  as soon as we know that  $x_1 < 1/2$ .

is constant on the neighborhood of  $\mu$  but it is also independent of  $x_2$  for any configuration where  $x_1 < 1/2$ .

Let us introduce some more notations: hereafter  $\|\mathbf{x}\|_\infty := \max_i |x_i|$  will denote the  $l_\infty$  norm,  $B(\mu, r) := \{\mathbf{x} \mid \|\mathbf{x} - \mu\|_\infty < r\}$  will denote the  $l_\infty$  ball (or box) of radius  $r$  around  $\mu$ , and  $\mathbf{e}_i$  will denote the standard basis vector with a 1 in the  $i^{\text{th}}$  coordinate and 0's elsewhere.

**Definition 2.** Let  $\mathcal{H}$  be a subset of  $\mathcal{F}$ . The decision function  $f$  is independent of its parameter  $i$  on  $\mathcal{H}$  if for any  $\alpha \in [-1, +1]$  we have:

$$\mathbf{x}, \mathbf{x} + \alpha \mathbf{e}_i \in \mathcal{H} \Rightarrow f(\mathbf{x}) = f(\mathbf{x} + \alpha \mathbf{e}_i)$$

For instance on Figure 1 the decision is independent of  $x_2$  on the set  $B(\mu, \Delta_2)$ . This local *independence*, parametrized by the  $\Delta_i$  radii, captures the sensitivity of the decision to its input parameters around the environment state  $\mu$ .

### 3. The SAVAGE Algorithm

We propose a generic zooming algorithm to solve the  $N$  dimensional parametric decision problem with high confidence. This algorithm, called **SAVAGE** (*Sensitivity Analysis of VArables for Generic Exploration*) is described in Algorithm 1. It works by reducing progressively a box-shaped confidence set  $\mathcal{H}$  until a single decision remains in  $f(\mathcal{H})$ . The algorithm stops

---

#### Algorithm 1 SAVAGE algorithm

---

```

1: Input:  $\mathbf{X} = (X_1, \dots, X_N)$ ,  $f$ ,  $\mathcal{F}$ ,  $T$ ,  $\delta$ 
2: Initialization:
3:  $\mathcal{W} := \{1, \dots, N\}$ ,  $\mathcal{H} := \mathcal{F}$ ,  $s := 1$ 
4:  $\forall i \in \mathcal{W} : \hat{\mu}_i := 1/2$ , and  $t_i := 0$ 
5: while  $\neg \text{Accept}(f, \mathcal{H}, \mathcal{W}) \wedge s \leq T$  do
6:   Pick a variable index  $i \in \arg \min_{\mathcal{W}} \{t_1, \dots, t_N\}$ 
7:    $t_i := t_i + 1$ 
8:   Sample the  $i^{\text{th}}$  distribution  $x_i \leftarrow X_i$ 
9:    $\hat{\mu}_i := (1 - \frac{1}{t_i})\hat{\mu}_i + \frac{1}{t_i}x_i$ 
10:   $\mathcal{H} := \mathcal{H} \cap \{\mathbf{x} \mid |x_i - \hat{\mu}_i| < c(t_i)\}$ 
11:   $\mathcal{W} := \mathcal{W} \setminus \{j \mid \text{IndepTest}(f, \mathcal{H}, j)\}$ 
12:   $s := s + 1$ 
13: end while
14: return  $\hat{d} \in f(\mathcal{H})$ 
    
```

---

exploring a parameter when it knows from a sensitivity analysis subroutine **IndepTest**( $f, \mathcal{H}, i$ ) that, given our knowledge of the environment, the final decision will not change according to this parameter; in other words when  $f$  is independent of  $i$  on  $\mathcal{H}$  as formalized in Definition 2.

The boundaries of  $\mathcal{H}$  are defined by the confidence radius:

$$c(t) = \sqrt{\frac{1}{2t} \log\left(\frac{\eta(t)}{\delta}\right)}, \quad (1)$$

where the  $\eta$  function is set to  $2NT$ , when the horizon  $T$  is finite, and  $\frac{\pi^2 N t^2}{3}$  when it is infinite (PAC setting).

Termination is controlled by the predicate:

$$\text{Accept}(f, \mathcal{H}, \mathcal{W}) := " \mathcal{W} = \emptyset " \quad (2)$$

$$\text{which implies } |f(\mathcal{H})| = 1 \quad (3)$$

**Theorem 1.** If  $f$  is independent of each parameter  $i$  on  $\mathcal{F} \cap B(\mu, \Delta_i)$ , with  $\Delta_i > 0$ , then SAVAGE is a  $(0, \delta)$ -PAC algorithm with horizon  $T$  for the parametric decision problem. When  $T = \infty$ , its sample complexity is bounded by:

$$\sum_{i=1}^N \mathcal{O}\left(\frac{\log\left(\frac{N}{\delta \Delta_i}\right)}{\Delta_i^2}\right).$$

When  $T < \infty$ , its exploration time is bounded by:

$$\sum_{i=1}^N \mathcal{O}\left(\frac{\log\left(\frac{NT}{\delta}\right)}{\Delta_i^2}\right).$$

*Proof sketch.* We first establish that the unknown parameter value  $\mu$  will stay inside the confidence set  $\mathcal{H}$  during all the computation with a probability of at least  $1 - \delta$  with the union bound and Hoeffding Lemma

(Hoeffding, 1963). Then we bound the exploration time of each parameter with a projection argument. Any point  $\mathbf{x}$  in  $\mathcal{H}$  can be projected into:

$$\mathbf{x}' := \mathbf{x} + \sum_{j \notin \mathcal{W}} (\mu_j - x_j) \mathbf{e}_j \quad (4)$$

By construction we have  $f(\mathbf{x}') = f(\mathbf{x})$ . The same projection can be done for any point  $\mathbf{x} + \alpha \mathbf{e}_i \in \mathcal{H}$ . By hypotheses we have  $f(\mathbf{x}') = f(\mathbf{x}' + \alpha \mathbf{e}_i)$ , hence  $f(\mathbf{x} + \alpha \mathbf{e}_i) = f(\mathbf{x})$ . This proof is detailed step by step in the extended version.  $\square$

By definition, if  $f$  is independent of each parameter  $i$  on  $B(\boldsymbol{\mu}, \Delta_i)$ , then  $f$  is constant on the minimal box  $B(\boldsymbol{\mu}, \Lambda)$ , where  $\Lambda = \min_i \Delta_i$ . An exploration policy without elimination would reach this neighborhood after  $\mathcal{O}\left(\frac{N \log(\frac{NT}{\delta})}{\Lambda^2}\right)$  samples.

This means that SAVAGE will outperform a uniform exploration policy as soon as the  $\Delta_i$  are not equal. It is also worth noting for practical purpose, that this improvement will hold even if we replace **IndepTest**( $f, \mathcal{H}, i$ ) with a sufficient condition of independence. Such a relaxation requires however to replace (2) by (3) to ensure termination.

### 3.1. Independence predicates

The independence predicate **IndepTest**( $f, \mathcal{H}, i$ ) is a property of the decision function and its feasible set. It can thus be specialized via symbolic calculus or hand-crafted for specific problems where the properties of  $f$  and  $\mathcal{F}$  are well known.

For example in traditional multi-armed bandits settings, when  $f(x_1, \dots, x_N) \in \arg \max \{x_1, \dots, x_N\}$  and  $\mathcal{H}(\mathbf{t})$  is encoded by a product of confidence intervals  $[a_i, b_i]$ , we can use the SAVAGE algorithm with the following specialized predicate **IndepTest** $_f(\mathcal{H}(\mathbf{t}), i)$ :

$$(\exists j, b_i \leq a_j) \vee (\forall k \neq i, b_k \leq a_i) \quad (5)$$

With this predicate, we fall-back almost to the "arm elimination" of (Even-Dar et al., 2002). We however slightly depart from this algorithm by forcing inclusion of the successive confidence sets:  $[a_i, b_i] := [\max\{a_i, \hat{\mu}_i - c(t_i)\}, \min\{b_i, \hat{\mu}_i + c(t_i)\}]$ .

If we rather want to retrieve the  $(m+1)^{\text{th}}$  best arm like in rigged bandits, the independence predicate becomes:

$$\vee \left( \begin{array}{l} (\exists \mathcal{A}, |\mathcal{A}| = m+1, \forall j \in \mathcal{A}, b_i \leq a_j) \\ (\exists \mathcal{B}, |\mathcal{B}| = N-m, \forall k \in \mathcal{B}, b_k \leq a_i) \end{array} \right) \quad (6)$$

A simple formalization of the independence allows us to apply SAVAGE and Theorem 1 to several other variants of multi-armed bandits.

---

### Algorithm 2 Parameters Elimination by Sampling

---

```

1: Input:  $f, \mathcal{H}, \mathcal{W}, m, M$ 
2: Initialization:
3:  $\mathcal{S} \leftarrow \emptyset$ 
4: for  $l = 1, \dots, m$  do
5:   Sample  $\mathbf{x}$  uniformly from  $\mathcal{H}$ 
6:    $\mathbf{x}' \leftarrow \mathbf{x}$ 
7:   for  $s = 1, \dots, M$  do
8:     Pick a random parameter  $i \in \mathcal{W} \setminus \mathcal{S}$ 
9:     Re-sample  $x_i$  until  $\mathbf{x} \in \mathcal{H}$ 
10:    if  $f(\mathbf{x}) \neq f(\mathbf{x}')$  then
11:       $\mathcal{S} := \mathcal{S} \cup \{i\}$ 
12:    end if
13:     $x'_i \leftarrow x_i$ 
14:  end for
15: end for
16:  $\mathcal{W} := \mathcal{W} \cap \mathcal{S}$ 
    
```

---

When the knowledge about  $f$  or  $\mathcal{F}$  is scarce, and the dimension of the problem is not too high, another solution that we only explored empirically is to estimate the independence predicate by "introspective" simulations. We used the multi-start random-walk approximation detailed in Algorithm 2. It provides an "almost-everywhere statement" of the property with an asymmetric risk of failure which can be made arbitrary low by increasing the number of samples (parameters  $m$  and  $M$ ). This kind of method is widely used in sensitivity analysis (see Saltelli et al., 2000, for a survey).

### 3.2. Approximate decision

If the decision function  $f$  is  $\lambda$ -Lipschitz for the decision comparison metric  $L$  and the  $l_\infty$  norm with a known Lipschitz constant, we are able to relax the problem by searching only an  $\varepsilon$ -approximation of the best decision. In order to do so we replace the **Accept**( $f, \mathcal{H}, \mathcal{W}$ ) condition by:

$$\forall \mathbf{x}, \mathbf{x}' \in \mathcal{H}, \|\mathbf{x} - \mathbf{x}'\|_\infty \leq \frac{\varepsilon}{\lambda} \quad (7)$$

**Theorem 2.** *If  $f$  is  $\lambda$ -Lipschitz around the environment state  $\boldsymbol{\mu}$ , and if  $f$  is independent of each parameter  $i$  on  $\mathcal{F} \cap B(\boldsymbol{\mu}, \Delta_i)$  with radius  $\Delta_i > 0$ , then SAVAGE with (7) as acceptance condition is an  $(\varepsilon, \delta)$ -PAC algorithm with horizon  $T$  for the parametric decision problem. When  $T = \infty$ , its sample complexity is bounded by:*

$$\sum_{i: \Delta_i \geq \varepsilon/\lambda} \mathcal{O}\left(\frac{\log(\frac{N}{\delta \Delta_i})}{\Delta_i^2}\right) + \mathcal{O}\left(\frac{\lambda^2 N_{\varepsilon, \lambda}}{\varepsilon^2} \log\left(\frac{\lambda N}{\delta \varepsilon}\right)\right);$$

where  $N_{\varepsilon, \lambda} = |\{i \mid \Delta_i < \varepsilon/\lambda\}|$ .

When  $T < \infty$ , its exploration time is bounded by:

$$\sum_{i: \Delta_i > \varepsilon/\lambda} \mathcal{O}\left(\frac{\log(\frac{NT}{\delta})}{\Delta_i^2}\right) + \mathcal{O}\left(\frac{\lambda^2 N_{\varepsilon, \lambda}}{\varepsilon^2} \log\left(\frac{NT}{\delta}\right)\right).$$

See extended version for the proof.  $\square$

## 4. Application to $K$ -armed Dueling Bandits

From now on, we call  $K \times K$  preference matrix a  $K \times K$  matrix  $(x_{i,j})$  such that  $x_{i,j} + x_{j,i} = 1$  for each  $i, j \in \{1, \dots, K\}$  (we use lower-case letters to match the notations of Section 2).

The  $K$ -dueling problem, as presented in (Yue et al., 2012; Yue & Joachims, 2011), assumes the existence of an environment preference matrix  $\mu$  from which we only have a noisy perception modeled, as in (Feige et al., 1994), by a  $K \times K$ -matrix of random variables  $X_{i,j} \in [0, 1]$  verifying  $\mathbb{E}[X_{i,j}] = \mu_{i,j}$ . Our aim is to design a sequence of pairwise experiments  $(i_t, j_t)$  called *duels* for  $t = 1, \dots, T$  in order to find the best arm. They also assume the following properties for the preference matrix (WLOG for a proper indexation of the matrix):

**strict linear order:** if  $i < j$  then  $\mu_{i,j} > \frac{1}{2}$ ;

**$\gamma$ -relaxed stochastic transitivity:** if  $1 < j < k$  then  $\gamma \cdot \mu_{1,k} \geq \max\{\mu_{1,j}, \mu_{j,k}\}$ ;

**stochastic triangular inequality:** if  $1 < j < k$  then  $\mu_{1,k} \leq \mu_{1,j} + \mu_{j,k} - \frac{1}{2}$ .

These last three assumptions are realistic when the preference matrix is the result of a perturbed linear order. This is indeed the case for some generative models where the number of parameters of the environment is assumed to be  $K$ : the inherent values of arms. In a situation where the preferences may contain cycles (or *voting paradoxes*) there is no clear notion of what the best arm is, and the notion of regret is unclear.

To avoid these problems, we propose to consider a "voting" variant of  $K$ -dueling bandits where a pairwise election criterion is used to determine the best candidate from the preference matrix. Several election systems can be used, but we will focus here on a simple and well-established one: the Copeland pairwise aggregation method (see Charon & Hudry, 2010).

### 4.1. Copeland bandits

If  $\mathbf{x}$  is a  $K \times K$  preference matrix, we define the *Copeland score* of an arm  $i$  by its number of one-to-one

majority victories:

$$U_{Cop}(i, \mathbf{x}) = \sum_j [x_{i,j} > \frac{1}{2}]. \quad (8)$$

Any element of  $\arg \max_i U_{Cop}(i, \mathbf{x})$  is called a *Copeland winner* of the matrix.

#### 4.1.1. GENERAL COPELAND BANDITS

With a preference matrix of size  $K$  we have  $N = K(K-1)/2$  free parameters to estimate: we can encode  $\mathcal{H}$  as a product of intervals  $[a_{i,j}, b_{i,j}]$  and apply the SAVAGE algorithm with  $\mathbf{IndepTest}_f(\mathcal{H}, (i, j)) =$

$$(a_{i,j} > \frac{1}{2} \vee b_{i,j} < \frac{1}{2}) \vee \mathbf{Cop}(\mathcal{H}, (i, j)),$$

where  $\mathbf{Cop}(\mathcal{H}, (i, j)) := \exists i^+ \text{ s. t.,}$  (9)

$$\wedge (\min U_{cop}(i^+, \mathcal{H}) > \max U_{cop}(i, \mathcal{H}))$$

$$\wedge (\min U_{cop}(i^+, \mathcal{H}) > \max U_{cop}(j, \mathcal{H})).$$

By applying Theorem 1, we obtain an exploration time bound of order:

$$\sum_{i < j} \mathcal{O}\left(\frac{\log(KT/\delta)}{\Delta_{i,j}^2}\right) \leq \mathcal{O}\left(K^2 \frac{\log(KT/\delta)}{\Lambda^2}\right), \quad (10)$$

where  $\Delta_{i,j} = |\mu_{i,j} - \frac{1}{2}|$  for any  $i < j$ , and  $\Lambda = \min_{i < j} \Delta_{i,j}$ . This bound requires weak assumptions about the preference matrix  $\mu$  but its strong dependence on the "hard" parameters (when  $\mu_{i,j}$  is close to  $\frac{1}{2}$ ) makes it quite conservative. The behavior of the algorithm is more efficient in practice.

#### 4.1.2. CONDORCET ASSUMPTION

If there exists an arm  $f(\mathbf{x})$  preferred to all the others, it is unique and verifies  $U_{Cop}(f(\mathbf{x}), \mathbf{x}) = K - 1$ . The existence of this arm, called *Condorcet winner* of the matrix, allows us to tighten the exploration bound.

**Property 1.** *If the environment state  $\mu$  admits arm  $i^*$  as a Condorcet winner with  $\Delta = \min_{j \neq i^*} \mu_{i^*,j} - \frac{1}{2}$  and  $\Delta_{i,j} = \max\{\Delta, |\mu_{i,j} - \frac{1}{2}|\}$  then  $f$  is independent of  $x_{i,j}$  on  $B(\mu, \Delta_{i,j})$  for any  $i < j$ .*

By applying Theorem 1 when Property 1 holds, we obtain a bound which is less sensitive to the presence of tight duels than (10) without changing the algorithm.

If we know the existence of a Condorcet winner, we can also tame the SAVAGE algorithm by restricting the feasible set  $\mathcal{F}$  to the  $K \times K$ -preferences matrices admitting a *Condorcet winner*:

$$\mathcal{F}_{Cond} := \{\mathbf{x} \mid \exists i^*, U_{Cop}(i^*, \mathbf{x}) = K - 1\}. \quad (11)$$



We can obtain a formal independence test in  $\mathcal{F}_{\text{cond}}$  by replacing (9) with:

$$\vee \begin{pmatrix} \max_{\mathbf{x} \in \mathcal{H}} U_{\text{Cop}}(i, \mathbf{x}) < K - 1 \\ \max_{\mathbf{x} \in \mathcal{H}} U_{\text{Cop}}(j, \mathbf{x}) < K - 1 \end{pmatrix} \quad (12)$$

To stop exploration with an  $\varepsilon$ -approximation<sup>1</sup> of the winner, we replace **Accept**( $f, \mathcal{H}, \mathcal{W}$ ) by:

$$\forall \mathbf{x} \in \mathcal{H}, K - 1 - U_{\text{Cop}}(f(\mathbf{x}), \mathbf{x}) \leq \varepsilon. \quad (13)$$

**Theorem 3.** *If the environment state  $\boldsymbol{\mu}$  is known to admit a Condorcet winner  $i^* = f(\boldsymbol{\mu})$  then SAVAGE with  $\mathcal{F}_{\text{Cond}}$  as feasible set and (13) as acceptance condition is an  $(\varepsilon, \delta)$ -PAC algorithm with horizon  $T$  for the Copeland bandits problem. When  $T = \infty$ , its samples complexity is bounded by:*

$$\sum_{j=\varepsilon+1}^{K-1} \mathcal{O} \left( \frac{j \cdot \log \left( \frac{K}{\delta \Delta_j} \right)}{\Delta_j^2} \right).$$

Where for each  $j \neq i^*$  we have  $\Delta_j = \mu_{i^*,j} - \frac{1}{2}$  (indexed WLOG by increasing values of  $\Delta_j$ ).

When  $T < \infty$  its exploration time is bounded by:

$$\sum_{j=\varepsilon+1}^{K-1} \mathcal{O} \left( \frac{j \cdot \log \left( \frac{KT}{\delta} \right)}{\Delta_j^2} \right).$$

This is a significant improvement from (10) but it does not remove the quadratic term  $K^2$ . This leading  $K^2$  factor is the price we pay for accepting less constrained preference matrices.

## 4.2. Borda bandits

Another simple way to elect the winner of the matrix is to use Borda count. Each competitor is ranked according to its mean performance against others:

$$U_{\text{Bor}}(i, \mathbf{x}) = x_{i,\cdot} = \sum_j x_{i,j}. \quad (14)$$

The main advantage of this criterion is that it both offers stability (the utility is linear) and clearly reduces the dimension of the problem to only  $K$  parameters:  $x_{i,\cdot}$  for  $i = 1, \dots, K$ . This means that we can simply wrap a classical bandit algorithm to search for the Borda winner of the matrix. It is quite easy however to design a Condorcet preference matrix where the Borda winner is not the Condorcet winner<sup>2</sup>.

The decision criterion underlying the *Beat the Mean Bandit* algorithm proposed by (Yue & Joachims, 2011)

<sup>1</sup> $\varepsilon$  is the number of tolerated defeats.

<sup>2</sup>There exist also  $K$ -armed stochastic bandit settings which are non transitive (Gardner, 1970).

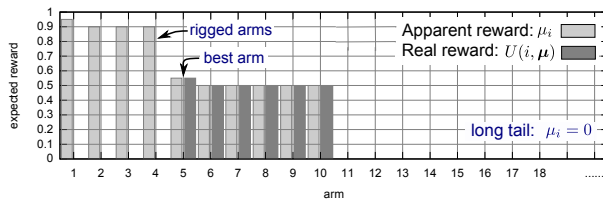


Figure 2. A 100-armed rigged bandit designed in order to game simple exploration algorithms: the four apparently best arms only deliver counterfeit money. The algorithm must sample intensively the arms 5 to 10 in order to guess that 5 is the best.

is different: the matrix rows are explored to find Borda losers which are progressively eliminated from the matrix until only one arm remains. This election procedure called *Bottom-up Borda elimination* returns the Condorcet winner if there exists one. SAVAGE being a generic algorithm, it can be applied directly to these two voting criteria.

## 5. Simulations

In order to compare algorithms of different natures, we used a pure-exploration online setting where at each time step  $t$  the algorithm choose a parameter  $i_t$  to explore, choose a decision  $d_t$  accordingly, and gets an unknown reward  $U(d_t, \boldsymbol{\mu})$ . We considered both the best-decision rate and the regret  $U(d^*, \boldsymbol{\mu}) - U(\hat{d}_t, \boldsymbol{\mu})$ . For all algorithms except *Interleave Filtering* and *Beat the Mean*, we took  $\hat{d}_t := f(\hat{\boldsymbol{\mu}}(t))$ . For PAC algorithms, we took  $\epsilon = 0$  and  $\delta = 1/T$  (explore-then-exploit setting). The sample time where the best-decision rate reaches  $1 - \delta$  gives an empirical estimation of the PAC exploration time. To avoid nasty side-effects, we shuffled the matrices/parameters at each run.

### 5.1. Bandits simulations

For bandits problems the decision space and the exploration space coincide, but we are here in a pure exploration setting where the arm we predict to be the best is not necessary the one we explore. We considered the following algorithms for our bandits simulations:

**Uniform:** baseline uniform exploration policy (each arm is explored once in a round-robin manner);

**Naive UCB:** UCB1 (as in Auer et al., 2002);

**Naive Elimination:** applies Action elimination al-

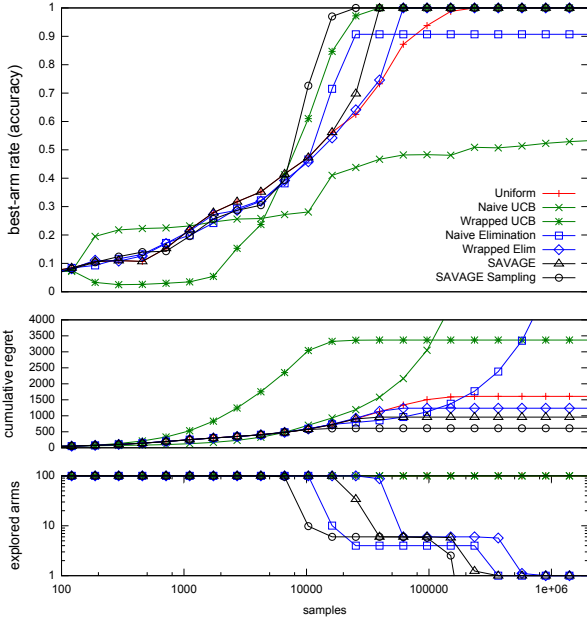


Figure 3. Behavior of the different algorithms for 1000 simulations with Figure 2 distribution. The top figure depicts the best-arm rate, the middle figure show the cumulative regret and the bottom one tracks the number of active arms for elimination algorithms. Time scale is logarithmic.

gorithm (as described in Even-Dar et al., 2002);

**Wrapped UCB:** applies UCB1 to the wrapped reward random variable  $\hat{U}_i = U(i, \hat{\boldsymbol{\mu}})$  used as a proxy for  $U(i, \boldsymbol{\mu})$ ;

**Wrapped Elimination:** applies Action elimination with the above "wrapped reward";

**SAVAGE:** applies Algorithm 1 with  $\eta(t) := 2NT$  and predicate (6) with  $m = 4$ ;

**SAVAGE Sampling:** Algorithm 1 with a sampled independence predicate and 1000 simulations by arm (see Algorithm 2).

We compared these algorithms on several Bernoulli reward distributions. We give here the simulation result for a rigged bandits problem specially designed in order to illustrate the different exploration behaviors of the algorithms. In this setting, the utility of arm  $i$  (indexed WLOG by decreasing  $\mu_i$ ) is defined by  $U(i, \boldsymbol{\mu}) = 0$  if  $i < 5$  and  $U(i, \boldsymbol{\mu}) = \mu_i$  otherwise (see Figure 2 for the reward distribution, and Figure 3 for the simulations results). As expected, the maximizing policies like *UCB* explore aggressively the head of the distribution but after around  $10^3$  samples fall into the rigged arms and neglect the second part of the distribution head. The wrapped versions are less sensitive to this trap, but the non-linearity of the utility and the

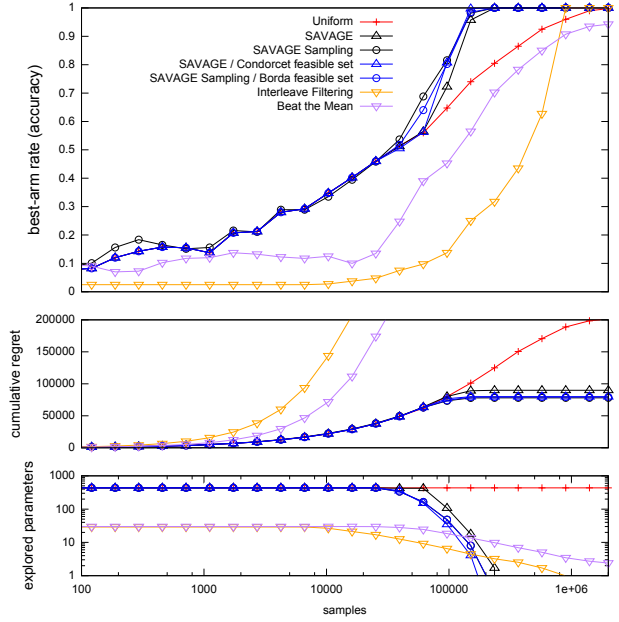


Figure 4. Average good prediction rate and regret for 500 simulations of a 30-armed Condorcet bandit instance. We used the Copeland index (8) to compute the regret.

violation of independence it induces both cripple their regret performances in the beginning of the runs. As expected, the **SAVAGE** versions perform well, more surprising is the side-effect of sampling which makes the algorithm more aggressive against weak arms.

## 5.2. Dueling bandits simulations

For the dueling bandits simulations, we considered **Uniform**, **SAVAGE**, and **SAVAGE Sampling** policies plus the following ones:

**SAVAGE/Condorcet:** Algorithm 1 with (12) for the independence test;

**SAVAGE Sampling/Borda:** Algorithm 1 with sampled oracle and a Borda relaxation of the Condorcet feasible set, i.e.  $\exists i, \sum_j x_{i,j} > K/2$ ;

**Interleave Filtering:** as in (Yue et al., 2012);

**Beat the Mean:** as in (Yue & Joachims, 2011) with  $\arg \max\{\hat{P}_b \mid b \in \mathcal{W}_i\}$  for  $\hat{d}_i$ .

### 5.2.1. CONDORCET SIMULATIONS

We first used "hard"  $K \times K$  Condorcet preference matrices  $\boldsymbol{\mu}$  defined by  $\mu_{i,j} = \frac{1}{2} + j/(2K)$  for each  $i < j$ . The matrices of this family verify all the assumptions defined in Section 4 but they also offer some difficult

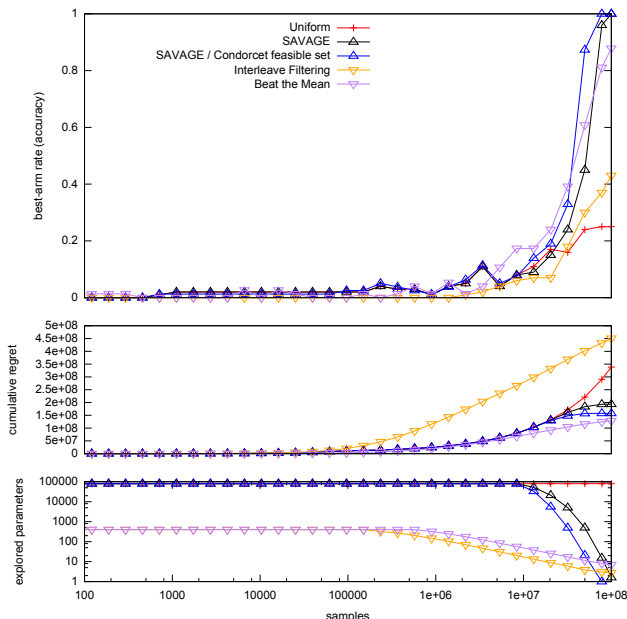


Figure 5. Same setting as in Figure 4 with  $K = 400$  and the horizon set to  $10^8$ . When we increase  $K$  the problem becomes difficult.

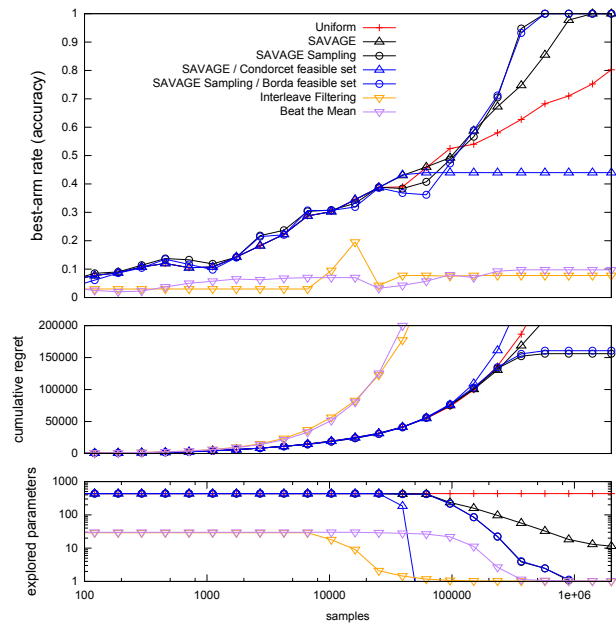


Figure 6. Result of 500 simulations with  $30 \times 30$  randomized preference matrices.

duels involving the Condorcet winner (for instance if  $K = 100$  we have  $\mu_{1,2} = 0.51$  hence  $\Delta = 0.01$ ).

The results of these experiments appear in Figure 4 and Figure 5. The SAVAGE Sampling policy slightly improves from the formal version but its heavy introspection cost makes it difficult to deploy on high-dimension problems. Low-dimension instance of Figure 4 is not favorable for *Interleave Filtering* and *Beat the Mean* which were designed to drop the  $\mathcal{O}(K^2)$  term of the bound with partial – hence risky – exploration strategies (see Yue & Joachims, 2011).

### 5.2.2. GENERAL CASE SIMULATIONS

In order to study the behavior of the algorithms with more realistic – non-Condorcet – preferences, we generated uniformly random preference matrices and performed the same experiments. As expected, when the Condorcet hypothesis is violated the performances of specialized algorithms (including SAVAGE/Condorcet) collapse well behind the baseline uniform exploration policy (see Figure 6).

## 6. Conclusion

We proposed **SAVAGE**, a flexible and generic algorithm based on sensitivity analysis of parameters for online learning with indirect feedback. We provided PAC theoretical guarantees for this algorithm when used with proper independence predicates. We also proposed a generic "introspective sampling" method to approximate these predicates.

Our simulations confirmed and reinforced the theoretical results on various parametric decision problems from classical bandits to  $K$ -armed dueling bandits.

The "voting bandits" framework we proposed naturally extends dueling bandits for realistic situations where the preferences reflects mixed and inconsistent opinions. The SAVAGE algorithm is robust and clearly outperforms state-of-the art algorithms in such situations.

The construction of a generic exploration algorithm reaching optimality for any provided decision function remains as a challenging open problem.

## Acknowledgments

We would like to thank the reviewers for their careful reading and helpful comments.



## References

- Audibert, J.Y., Munos, R., and Szepesvári, Cs. Exploration-exploitation trade-off using variance estimates in multi-armed bandits. *Theoretical Computer Science*, 2008.
- Audibert, J.Y., Bubeck, S., and Munos, R. Best arm identification in multi-armed bandits. In *COLT*, Haifa (Israel), 2010. Omnipress.
- Auer, P. and Ortner, R. UCB revisited: Improved regret bounds for the stochastic multi-armed bandit problem. *Period.Math.Hungar.*, 61(1-2):55–65, 2011.
- Auer, P., Cesa-Bianchi, N., and Fischer, P. Finite-time analysis of the multiarmed bandit problem. *Mach. Learn.*, 47(2-3):235–256, 2002.
- Bubeck, S., Munos, R., and Stoltz, G. Pure exploration in finitely-armed and continuous-armed bandits. *Theor. Comput. Sci.*, 412(19):1832–1852, 2011.
- Chapelle, O., Joachims, T., Radlinski, F., and Yue, Y. Large-scale validation and analysis of interleaved search evaluation. *ACM Trans. Inf. Syst.*, 30(1):6, 2012.
- Charon, I. and Hudry, O. An updated survey on the linear ordering problem for weighted or unweighted tournaments. *Annals OR*, 175(1):107–158, 2010.
- Chevalere, Y., Endriss, U., Lang, J., and Maudet, N. A short introduction to computational social choice. In *SOFSEM*, volume 4362 of *LNCS*, pp. 51–69. Springer-Verlag, 2007.
- Di Castro, D., Gentile, C., and Mannor, S. Bandits with an edge. *CoRR*, abs/1109.2296, 2011.
- Even-Dar, E., Mannor, S., and Mansour, Y. PAC bounds for multi-armed bandit and Markov decision processes. In *COLT*, pp. 255–270, London, UK, 2002. Springer-Verlag.
- Even-Dar, E., Mannor, S., and Mansour, Y. Action elimination and stopping conditions for the multi-armed bandit and reinforcement learning problems. *JMLR*, 7:1079–1105, 2006.
- Feige, U., Raghavan, P., Peleg, D., and Upfal, E. Computing with noisy information. *SIAM J. Comput.*, 23(5):1001–1018, 1994.
- Gardner, M. Mathematical games: The paradox of the nontransitive dice and the elusive principle of indifference. *Scientific American*, 223:110–114, dec 1970.
- Hoeffding, W. Probability inequalities for sums of bounded random variables. *J. of the American Statistical Association*, 58(301):13–30, 1963.
- Joachims, T. Evaluating retrieval performance using clickthrough data. In *Text Mining*, pp. 79–96. Physica/Springer Verlag, 2003.
- Lai, T.L. and Robbins, H. Asymptotically efficient adaptive allocation rules. *Advances in Applied Mathematics*, 6(1):4–22, 1985.
- Ravikumar, B., Ganesan, K., and Lakshmanan, K.B. On selecting the largest element in spite of erroneous information. In *STACS*, volume 247 of *LNCS*, pp. 88–99. Springer Berlin Heidelberg, 1987.
- Saltelli, A., Chan, K., and Scott, E.M. (eds.). *Sensitivity analysis*. Wiley series in probability and statistics. J. Wiley & sons, New York, Chichester, Weinheim, 2000.
- Yue, Y. and Joachims, T. Interactively optimizing information retrieval systems as a dueling bandits problem. In *ICML*, volume 382 of *ACM Proceeding Series*, pp. 1201–1208. ACM, 2009.
- Yue, Y. and Joachims, T. Beat the mean bandit. In *ICML*, pp. 241–248. Omnipress, 2011.
- Yue, Y., Broder, J., Kleinberg, R., and Joachims, T. The  $k$ -armed dueling bandits problem. *J. Comput. Syst. Sci.*, 78(5):1538–1556, 2012.