

---

# Optimal betting for a multi-armed bandit

---

Kurt Ehlert  
kehlert@math.wisc.edu

## Abstract

If we are betting on a slot machine, then the natural question to ask is “how much should we bet”? If the probability of winning is  $p$ , it turns out that betting a fraction of our wealth equal to  $2p - 1$  is an optimal strategy in many senses. However, if we are faced with many slot machines with unknown probabilities of winning, then how should we bet, and which slot machines should we bet on? We will present the Kelly-UCB algorithm, which addresses this more complicated problem. We will also provide analytic and numerical results regarding the performance of the Kelly-UCB algorithm.

## 1 Introduction

A slot machine is also known as a one-armed bandit, because early slot machines were operated by pulling a lever (arm) attached to their side [Glimne, 2015]. Suppose that we are pulling an arm with a probability  $p$  of winning, where  $1/2 < p < 1$ . Also, we start with a finite amount of money, and we can bet any fraction of our wealth each time we pull the arm. We win an amount equal to our bet with probability  $p$ , and conversely we lose our bet with probability  $1 - p$ . If  $p > 1/2$ , then we are likely to gain money by playing the game, but we need to be careful with our betting strategy. If we bet a fraction of our wealth equal to  $2p - 1$ , then this strategy is known as the “Kelly bet” [Kelly, 1956, Thorp, 2006]. In some senses it maximizes our wealth, and it also avoids bankruptcy. In Section 2, we briefly discuss the Kelly bet for a slightly generalized situation and prove that it is optimal in a few basic senses.

Now consider a situation where we can choose between many arms. If we pull arm  $i$ , then we win with probability  $p_i$  and lose with probability  $1 - p_i$ . If we are only allowed to bet a fixed amount and we do not know the value of  $p_i$  for any  $i$ , then this is called the stochastic multi-armed bandit problem. There is a trade-off between exploitation and exploration: we want to exploit an arm with a seemingly high probability of winning, but we also want to explore and find the best arm. Here we consider a new version of the stochastic bandit problem. Instead of fixed-size bets, we can bet any fraction of our wealth. The KL-UCB algorithm in Cappé et al. [2013] does not address how much to bet, so we plan to extend the KL-UCB algorithm and its analysis to the variable-bet situation. We will also allow for more complicated payoffs than just a win or loss of our entire bet.

In Section 2, we review the Kelly bet for the one-armed bandit with known probabilities. Then we extend the Kelly bet to the situation where we do not know the probabilities of the possible outcomes. Section 3 outlines our Kelly-UCB algorithm for playing a variable-bet multi-armed bandit. Then we will present analytic and numerical results on the algorithm’s performance.

## 2 Optimal betting for a single-armed bandit

### 2.1 Kelly bet for a bandit with a known probability of winning

Suppose we are playing a slot machine that allows us to bet any proportion of our wealth. Let  $X_n$  be a random variable that takes on only finitely many values in  $[-1, \infty)$ . When we pull the arm, we

obtain of profit of  $X_n$  per unit bet on the  $n^{\text{th}}$  pull. For example, if we bet \$100 and  $X_1 = -1$ , then we lose our entire bet of \$100. If instead  $X_1 = 2$ , then we win \$200 (and get back our original bet of \$100). Let  $p_j$ ,  $1 \leq j \leq m$ , denote the probability that  $X_n$  is equal to  $x_j$ . Furthermore, assume that  $p_j$  is known for every  $j$ , the  $X_n$  are i.i.d., and that  $\mathbb{E}[X_1] > 0$ . Also assume that  $P(X_1 = -1) > 0$ .

If we want to maximize our expected gain, then we should bet our entire fortune on every pull, but we would go bankrupt with probability 1. If we bet a fixed amount, then there is still a positive probability of going bankrupt. Instead we will take a different approach and bet a fixed proportion of our wealth on each pull. Later we will show that this is actually optimal in two senses. Ethier [2010] analyzes the Kelly bet in much more detail and shows that it is optimal in many more ways than we consider here. Ethier [2010] also notes that some of the assumptions on  $X_1$  can be relaxed, namely it does not need to be discrete-valued and we can have  $P(X_1 = -1) = 0$ .

Let  $W_0$  be our initial wealth, and let  $W_n$  is our wealth after pull  $n$ . Define  $f$  as the proportion of our wealth that we bet. Then

$$W_n = W_{n-1} + fW_{n-1}X_n \quad (2.1)$$

$$= W_{n-1}(1 + fX_n) \quad (2.2)$$

Using recursion, we can see that our wealth after pull  $n$  is

$$W_n = W_0 \prod_{i=1}^n (1 + fX_i) \quad (2.3)$$

Let  $r_n(f) = n^{-1} \log(W_n/W_0)$ , then we can trivially rewrite our wealth after pull  $n$  as

$$W_n = W_0 e^{r_n(f)n} \quad (2.4)$$

We can interpret  $r_n$  as the average geometric growth rate of our wealth after pull  $n$ . Intuitively, we want to maximize the growth rate of our wealth. That intuition leads to the following lemma, which is based on Lemma 10.1.1 in Ethier [2010]

**Lemma 2.1.** *Let  $\mu(f) = \mathbb{E}[\log(1 + fX_1)]$ , which is defined for  $f \in [0, 1)$ . Then  $\lim_{n \rightarrow \infty} r_n(f) = \mu(f)$ . Let  $f^* = \arg \max_{0 \leq f \leq 1} \mu(f)$ , then  $\mu(f^*) > 0$ .*

*Proof.* The strong law of large numbers implies

$$\lim_{n \rightarrow \infty} r_n(f) = \lim_{n \rightarrow \infty} \frac{1}{n} \sum_{i=1}^n \log(1 + fX_i) = \mathbb{E}[\log(1 + fX_1)] \text{ a.s.} \quad (2.5)$$

$\mu(f)$  is strictly concave, because

$$\mu''(f) = -\mathbb{E} \left[ \frac{X^2}{(1 + fX)^2} \right] < 0 \quad (2.6)$$

Therefore we can find the global maximum of  $\mu(f)$  by just setting its derivative equal to 0. If  $P(X_n = 1) = 1 - P(X_n = -1)$ , then  $\mu'(f)$  is straightforward to compute and we find that  $\arg \max_{0 \leq f \leq 1} \mu(f) = 2p - 1$ . Furthermore,  $\mu'(0) = \mathbb{E}[X] > 0$ , and  $P(X_1 = -1 > 0)$  implies that  $\mu'(1^-) = -\infty$ . Therefore  $\mu'(f) = 0$  for some  $f$ . Since  $\mu(f)$  is strictly concave, there is a unique  $f$  that satisfies  $\mu'(f) = 0$  and maximizes  $\mu(f)$ . Denote that  $f$  as  $f^*$ . Furthermore, since  $\mu(0) = 0$  and  $\mu'(0) > 0$ , we must have  $\mu(f) > 0$  for some small enough  $f$  near 0. Consequently,  $\mu(f^*) > 0$ .  $\square$

This fixed-proportion betting system is known as the *Kelly bet* or *Kelly criterion* [Kelly, 1956, Thorp, 2006]. Above we showed that the Kelly bet is optimal in an asymptotic sense, so the natural question to ask is if it is also optimal for finite times. If we wish to maximize the logarithm of our wealth after  $n$  pulls, then once again the Kelly bet is optimal.

**Lemma 2.2.**  *$f^*$  maximizes  $\mathbb{E}[\log W_n(f)]$ .*

*Proof.* From the definition of  $W_n$  and the linearity of expectations

$$\mathbb{E}[\log W_n(f)] = \log W_0 + \sum_{i=1}^n \mathbb{E}[\log(1 + fX_i)] \quad (2.7)$$

$$= \log W_0 + n\mu(f) \quad (2.8)$$

Therefore maximizing the log-wealth is equivalent to maximizing  $\mu(f)$ . In lemma 2.1 we showed that  $f = 2p - 1$  maximizes  $\mu(f)$ .  $\square$

The Kelly bet also approximately maximizes the median wealth. Also, in the long-run it will do better than any other “essentially different” strategy, which includes strategies that change betting proportions between pulls. Ethier [2010] provides the details.

## 2.2 Kelly bet for a bandit with an unknown probability of winning

A simple approach to estimating the Kelly bet would be to calculate  $f^*$  based on estimated values of  $\{p_j\}_{j=1}^m$ . The obvious estimators are

$$\bar{p}_j = \frac{1}{n} \sum_{i=1}^n \mathbb{1}(X_i = x_j), \quad 1 \leq j \leq m$$

This approach is not unreasonable, because our wealth would generally grow if  $\bar{p}$  is close to the actual probability of winning  $p$ . However, we would likely go bankrupt during the first few pulls, since  $\bar{p}$  has a relatively high probability of being  $\pm 1$ . A reasonable idea is to shrink  $f_n$  toward zero. Orabona and Pál [2016] show that by using the Krichevsky and Trofimov [1981] (KT) estimate of  $p$ , then our bets converge to the Kelly bet while controlling our losses. The KT estimate of  $f_n$  essentially shrinks  $f_n$  toward 0. The KT estimator of  $p$  is

$$\hat{p}_n = \frac{1}{n+1} \left( \frac{1}{2} + \sum_{i=1}^n \mathbb{1}\{X_i = 1\} \right) \quad (2.9)$$

Let  $\hat{f}_n = 2\hat{p}_n - 1$  be our bet on pull  $n$ , then

$$\hat{f}_n = 2\hat{p}_n - 1 = \frac{1}{n+1} \sum_{i=1}^n X_i \quad (2.10)$$

Thus our wealth after pull  $n$  is

$$W_n = W_0 \prod_{i=1}^n (1 + f_i X_i) \quad (2.11)$$

Define  $f_i$  and  $p_i$  such that  $f_i = 2p_i - 1$ . Also let  $Y_i = (1 + X_i)/2$ , and define the loss  $\ell(p_i, Y_i)$  as

$$\ell(p_i, Y_i) = -x_i \log p_i - (1 - Y_i) \log(1 - p_i) \quad (2.12)$$

Combining the above definitions leads to the following

$$\log W_n = \log W_0 + \sum_{i=1}^n \log(1 + f_i X_i) \quad (2.13)$$

$$= \log W_0 + \sum_{i=1}^n \frac{1 + X_i}{2} \log(2p_i) + \frac{1 - X_i}{2} \log(2(1 - p_i)) \quad (2.14)$$

$$= \log W_0 + n \log 2 + \sum_{i=1}^n \frac{1 + X_i}{2} \log(p_i) + \frac{1 - X_i}{2} \log(1 - p_i) \quad (2.15)$$

$$= \log W_0 + n \log 2 - \sum_{i=1}^n \ell(p_i, Y_i) \quad (2.16)$$

Since the above holds for any choice of  $p_i$ , in particular it holds for  $p_i = p, \hat{p}_n$ . Let  $W_n^*$  be our wealth from using the Kelly bet, and let  $W_n^\circ$  be our wealth using  $\hat{p}_i$ , where the  $\hat{p}_i$  are  $\sigma(X_1, X_2, \dots, X_{i-1})$  measurable. Define the regret  $R_n$  as

$$R_n = \log W_n^* - \log W_n^\circ = \sum_{i=1}^n \ell(\hat{p}_i, Y_i) - \ell(p, Y_i) \quad (2.17)$$

Let  $n_1 = \sum_{i=1}^n \mathbb{1}\{X_i = 1\}$ , then according to Cesa-Bianchi and Lugosi [2006] the KT estimator satisfies

$$R_n = \log \frac{\pi p^{n_1} (1-p)^{n-n_1}}{\text{Beta}(n_1 + 1/2, n - n_1 + 1/2)} \quad (2.18)$$

The right-hand side of the above expression attains its maximum at  $n_1 = np$ . Therefore we have the following bound

$$R_n \leq -nH(p) - \log \text{Beta}(np + 1/2, n(1-p) + 1/2) + \log(\pi) \quad (2.19)$$

If  $np$  is an integer, then the bound is sharp. We can use Stirling's formula to approximate the beta function, which gives

$$R_n \leq -nH(p) - \log \text{Beta}(np + 1/2, n(1-p) + 1/2) + \log(\pi) \quad (2.20)$$

$$\sim \frac{1}{2} \log(n+1) + \frac{1}{2} \log(\pi/2) + o(1) \quad (2.21)$$

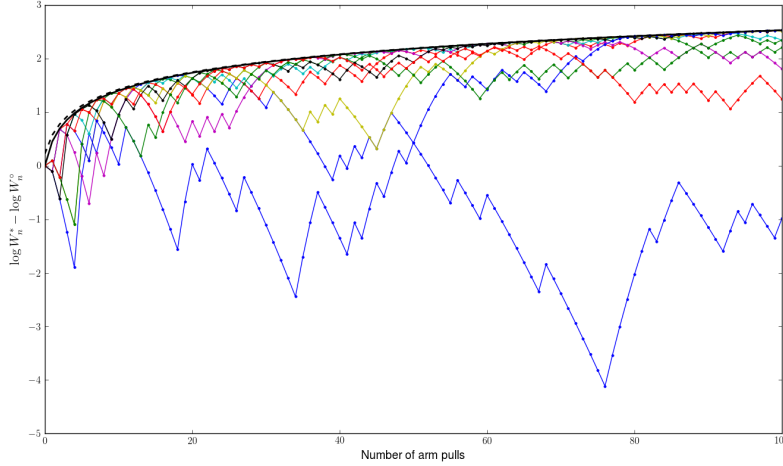


Figure 1: Comparison of regret paths to the KT regret bound for  $p = 0.55$ . The bold black curve on top is the regret bound  $-nH(p) - \log \text{Beta}(np + 1/2, n(1-p) + 1/2) + \log(\pi)$  from (2.20). The dashed black curve is the asymptotic bound (2.21). The others plots are the regrets of ten different paths.

### 3 Kelly-UCB algorithm for a multi-armed bandit

In this section, we consider a situation where we have a choice between many different arms. If we pull arm  $i$ , then we win with probability  $p_i$  and lose with probability  $1 - p_i$ . If we are only allowed to bet a fixed amount and we do not know the value of  $p_i$  for any  $i$ , then this is called the stochastic multi-armed bandit problem. There is a trade-off between exploitation and exploration: we want to exploit an arm with a seemingly high probability of winning, but we also want to explore and find the best arm.

There are many variations of the multi-armed bandit problem described in Bubeck et al. [2012]. Here we consider a new version of the stochastic bandit problem. Instead of fixed-size bets, we can bet any fraction of our wealth. The KL-UCB algorithm in Cappé et al. [2013] decides which arms to pull, but it does not address how much to bet. We also want to allow “shorting”, which means that we gain our bet if a pull is a loss (i.e. bet negative amounts). We modified the KL-UCB algorithm to allow for variable bets and shorting.

#### 3.1 Description of the Kelly-UCB algorithm

Suppose there are  $K$  levers that are indexed by  $i = 1, \dots, K$ . At each time  $t = 1, 2, \dots$ , we can choose which lever to pull. Let  $\{X_{i,t}\}$  be independent random variables indicating the payoff of lever  $i$  on pull  $t$ , where  $P(X_{i,t} = 1) = 1 - P(X_{i,t} = -1) = p_i$ . We also need to decide how much to bet. We start with a finite amount of wealth, and we stop betting if we go bankrupt.

Let  $W_0$  be our initial wealth, and let  $f_{i,t}$  be the fraction of our wealth we bet on lever  $i$  at time  $t$ . Also let  $I_t \in \{1, \dots, K\}$  be the index of the lever we pull at time  $t$ . Then

$$W_T = W_0 \prod_{t=1}^T (1 + f_{I_t,t} X_{I_t,t}) \quad (3.1)$$

Below is the Kelly-UCB algorithm. Note that our lever choice is given by a modified version of the KL-UCB algorithm in Cappé et al. [2013]. Let  $\text{kl}(p, q)$  denote the Kullback-Leibler divergence between two Bernoulli distributions with parameters  $p$  and  $q$ . We do not explicitly address the possibility of  $f = 0$  in the algorithm, but an easy fix is to just bet some extremely small minimal amount instead.

---

**Algorithm 1** Kelly-UCB

---

**Require:**  $\epsilon > 0$  where  $\epsilon \ll 1$ , initial wealth  $W > 0$

```

1: for  $i = 1$  to  $K$  do
2:    $X \leftarrow \pm 1$  where  $P(X = 1) = 1 - P(X = -1) = p_i$ 
3:    $W \leftarrow W(1 + \epsilon X)$ 
4:    $S_i \leftarrow (X + 1)/2$ 
5:    $N_i \leftarrow 1$ 
6:    $f_i \leftarrow X/2$ 
7: end for
8: for  $t = K + 1$  to  $T$  do
9:   for  $i = 1$  to  $K$  do
10:     $q_{i,\text{long}} \leftarrow \max\{q \in [0, 1] \mid N_i \text{kl}(S_i/N_i, q) \leq \log t\}$ 
11:     $q_{i,\text{short}} \leftarrow 1 - \min\{q \in [0, 1] \mid N_i \text{kl}(S_i/N_i, q) \leq \log t\}$ 
12:     $q_i \leftarrow \max(q_{i,\text{long}}, q_{i,\text{short}})$ 
13:   end for
14:   choose  $I \in \arg \max_{i=1, \dots, K} q_i$ 
15:    $X \leftarrow \pm 1$  where  $P(X = 1) = 1 - P(X = -1) = p_I$ 
16:    $W \leftarrow W(1 + f_I X)$ 
17:    $S_I \leftarrow S_I + X + 1)/2$ 
18:    $N_I \leftarrow N_I + 1$ 
19:    $f_I \leftarrow (N_I + 1)^{-1}(N_I f_I + X)$ 
20: end for

```

---

### 3.2 Analysis of the Kelly-UCB algorithm

Similar to Cappé et al. [2013] and Bubeck et al. [2012], we will bound the regret at time  $T$ , which we denote as  $R_T$ . Let  $W_T(\Phi)$  be our wealth at time  $T$  when using strategy  $\Phi$ , where  $\Phi$  is any permissible strategy (it does not look into the future and does not bet more than we have). Let  $\hat{\Phi}$  be the strategy given in algorithm 1. Define

$$\mathbb{E}[R_T] = \max_{\Phi} E[\log W_T(\Phi) - \log W_T(\hat{\Phi})] \quad (3.2)$$

Let  $H(p)$  denote the entropy of a Bernoulli distribution with parameter  $p$ , and define  $f_i = 2p_i - 1$  ( $f_i$  is the Kelly bet for arm  $i$ ).

**Theorem 3.1.** *Let  $i^* \in \arg \max_{i=1, \dots, K} |p_i - 1/2|$ , and let  $\Phi^*$  be the strategy where we always pull lever  $i^*$  and bet  $f = 2p_{i^*} - 1$ . Then  $\Phi^* = \arg \max_{\Phi} \mathbb{E}[\log W_T(\Phi)]$  for any  $T \geq 0$ .*

*Proof.* We can define strategy  $\Phi$  as the sequence  $\{I_t, f_{I_t,t}\}$ , so (3.1) implies

$$\mathbb{E}[\log W_T(\Phi)] = \log W_0 + \mathbb{E} \left[ \sum_{i=1}^T \log(1 + f_{I_t,t} X_{I_t,t}) \right] \quad (3.3)$$

$$= \log W_0 + \sum_{i=1}^T \mathbb{E} [\mathbb{E} [\log(1 + f_{I_t,t} X_{I_t,t}) \mid I_t]] \quad (3.4)$$

$$= \log W_0 + \sum_{i=1}^T \mathbb{E} [p_{I_t} \log(1 + f_{I_t,t}) + (1 - p_{I_t}) \log(1 - f_{I_t,t})] \quad (3.5)$$

$$\leq \log W_0 + \sum_{i=1}^T \mathbb{E} [p_{I_t} \log(2p_{I_t}) + (1 - p_{I_t}) \log(2(1 - p_{I_t}))] \quad (3.6)$$

$$= \log W_0 + T \log 2 - \sum_{i=1}^T H(p_{I_t}) \quad (3.7)$$

$$\leq \log W_0 + T \log 2 - TH(p_{i^*}) \quad (3.8)$$

$$= \mathbb{E}[\log W_T(\Phi^*)] \quad (3.9)$$

The third equality follows from the fact that conditioned on  $I_t$ ,  $f_{I_t,t}$  and  $X_{I_t,t}$  are independent. The first inequality follows from theorem 2.1.  $\square$

**Theorem 3.2.** Let  $R_T$  be the regret defined in (3.2). Define

$$kl_i = kl(\max(p_i, 1 - p_i), \max(p_{i^*}, 1 - p_{i^*})) \quad (3.10)$$

Then for all  $T \geq 0$

$$\mathbb{E}[R_T] \leq \frac{1}{2} \log(T+1) + \sum_{i=1}^K \left[ \frac{\Delta_i \log T}{kl_i} (1 + o(1)) + \frac{1}{2} \log \left( \frac{\log T}{kl_i} + 1 \right) \right] + \frac{K}{2} \log(\pi/2) \quad (3.11)$$

*TODO proofread, and check all the little-oh terms*

*Proof.* Continuing from (3.2) and using theorem 3.1 leads to

$$\mathbb{E}[R_T] = \mathbb{E} \left[ \log W_T(\Phi^*) - \log W_T(\hat{\Phi}) \right] \quad (3.12)$$

$$= \mathbb{E} \left[ \sum_{t=1}^T \log \frac{1 + f_{i^*} X_{i^*,t}}{1 + \hat{f}_{I_t,t} X_{I_t,t}} \right] \quad (3.13)$$

$$= \underbrace{\mathbb{E} \left[ \sum_{t=1}^T \log \frac{1 + f_{i^*} X_{i^*,t}}{1 + f_{I_t,t} X_{I_t,t}} \right]}_{:=A} + \underbrace{\mathbb{E} \left[ \sum_{t=1}^T \log \frac{1 + f_{I_t,t} X_{I_t,t}}{1 + \hat{f}_{I_t,t} X_{I_t,t}} \right]}_{:=B} \quad (3.14)$$

We will simplify  $A$  and  $B$  separately. We can think of  $A$  as the regret caused by choosing the wrong arm, and  $B$  is the regret caused by using an approximation of the Kelly bet.

$$A = \sum_{t=1}^T \mathbb{E}[\log(1 + f_{i^*} X_{i^*,t})] - \sum_{t=1}^T \mathbb{E}[\log(1 + f_{I_t} X_{I_t,t})] \quad (3.15)$$

$$= \sum_{t=1}^T [\log(2) - H(p_{i^*})] - \sum_{t=1}^T \mathbb{E}[\log(2) - H(p_{I_t})] \quad (3.16)$$

$$= \sum_{t=1}^T \mathbb{E}[-H(p_{i^*}) - H(p_{I_t})] \quad (3.17)$$

$$= \sum_{i=1}^K \mathbb{E} \left[ \sum_{j=1}^{N_i(t)} -H(p_{i^*}) - H(p_i) \right] \quad (3.18)$$

$$= \sum_{i=1}^K \mathbb{E}[N_i(t)] \Delta_i \quad (3.19)$$

where  $\Delta_i = -H(p_{i^*}) + H(p_i) \geq 0$ . As for  $B$

$$B = \mathbb{E} \left[ \sum_{t=1}^T \log \frac{1 + f_{I_t} X_{I_t,t}}{1 + \hat{f}_{I_t,t} X_{I_t,t}} \right] \quad (3.20)$$

$$= \mathbb{E} \left[ \sum_{t=1}^T \frac{1 + X_{I_t,t}}{2} \log \left( \frac{1 + f_{I_t}}{1 + \hat{f}_{I_t,t}} \right) + \frac{1 - X_{I_t,t}}{2} \log \left( \frac{1 - f_{I_t}}{1 - \hat{f}_{I_t,t}} \right) \right] \quad (3.21)$$

$$(3.22)$$

Define  $\hat{p}_{I_t,t}$  such that  $\hat{f}_{I_t,t} = 2p_{I_t,t} - 1$ , and let  $Z_i \sim \text{Bernoulli}(p_i)$ . Then continuing from above we get

$$B = \mathbb{E} \left[ \sum_{t=1}^T Z_i \log \left( \frac{p_{I_t}}{\hat{p}_{I_t,t}} \right) + (1 - Z_i) \log \left( \frac{1 - p_{I_t}}{1 - \hat{p}_{I_t,t}} \right) \right] \quad (3.23)$$

$$= \mathbb{E} \left[ \sum_{t=1}^T -\ell(p_{I_t}, Z_{I_t}) + \ell(\hat{p}_{I_t,t}, Z_{I_t}) \right] \quad (3.24)$$

$$= \sum_{i=1}^K \mathbb{E} \left[ \sum_{j=1}^{N_i(t)} -\ell(p_i, Z_i) + \ell(\hat{p}_{i,t}, Z_i) \right] \quad (3.25)$$

Note that  $\ell(p_i, y_i)$  was defined in (2.12). Combining the results for  $A$  and  $B$  leads to

$$\mathbb{E}[R_T] = \sum_{i=1}^K \mathbb{E}[N_i(t)] \Delta_i + \sum_{i=1}^K \mathbb{E} \left[ \sum_{j=1}^{N_i(t)} -\ell(p_i, Z_i) + \ell(\hat{p}_{i,t}, Z_i) \right] \quad (3.26)$$

The rightmost sum appears because we are using an estimate of the Kelly bet. If  $f_{i,t}$  is the exact Kelly bet  $f_i$ , then that sum is zero. Applying (2.20) to (3.26) leads to

$$\mathbb{E}[R_T] \leq \sum_{i=1}^K \mathbb{E}[N_i(T)] \Delta_i + \sum_{i=1}^K \mathbb{E} \left[ \frac{1}{2} \log(N_i(t) + 1) + \frac{1}{2} \log(\pi/2) + o(1) \right] \quad (3.27)$$

$$\leq \sum_{i=1}^K \left[ \mathbb{E}[N_i(T)] \Delta_i + \frac{1}{2} \log(\mathbb{E}[N_i(T)] + 1) \right] + \frac{K}{2} \log(\pi/2) + o(1) \quad (3.28)$$

The second inequality is Jensen's inequality. Cappé et al. [2013] showed that if  $f_{I_t,t} \geq 0$ , and  $i \neq i^*$

$$\mathbb{E}[N_i(t)] \leq \frac{\log T}{\text{kl}(p_i, p_{i^*})} (1 + o(1)) \quad (3.29)$$

Since we allow  $f_{I_{t,t}} < 0$ , we need to adjust the inequality. Let

$$\text{kl}_i = \text{kl}(\max(p_i, 1 - p_i), \max(p_{i^*}, 1 - p_{i^*})) \quad (3.30)$$

Lemma 3.1 shows that for  $i \neq i^*$

$$\mathbb{E}[N_i(t)] \leq \frac{\log T}{\text{kl}_i} (1 + o(1)) \quad (3.31)$$

And since  $N_{i^*}(T) \leq T$ , we have

$$\mathbb{E}[R_T] \leq \frac{1}{2} \log(T+1) + \sum_{i=1}^K \left[ \frac{\Delta_i \log T}{\text{kl}_i} (1 + o(1)) + \frac{1}{2} \log \left( \frac{\log T}{\text{kl}_i} + 1 \right) \right] + \frac{K}{2} \log \left( \frac{\pi}{2} \right) \quad (3.32)$$

□

The  $\frac{1}{2} \log(T+1)$  term is generally dominated by the other  $\log T$  term inside the sum. In other words, the error from choosing the wrong arm generally dominates the error from betting incorrectly. Numerical results for  $\vec{p} = (0.4, 0.5, 0.8)$  are shown in figure 2.

**Lemma 3.1.** Define  $\text{kl}_i = \text{kl}(\max(p_i, 1 - p_i), \max(p_{i^*}, 1 - p_{i^*}))$ . Then

$$\mathbb{E}[N_i(t)] \leq \frac{\log T}{\text{kl}_i} (1 + o(1))$$

*Proof.* adapt the Cappe proof

□

We would also like to find a lower bound for the regret. Cappé et al. [2013] gives a lower bound for the number times each suboptimal arm is pulled. We do not find a lower bound here. However, we could get a lower bound for the Kelly-UCB regret by letting  $\hat{f}_{i,t} = 2p_i - 1$  for all  $t$ , and then we use the aforementioned lower bound in Cappé et al. [2013].

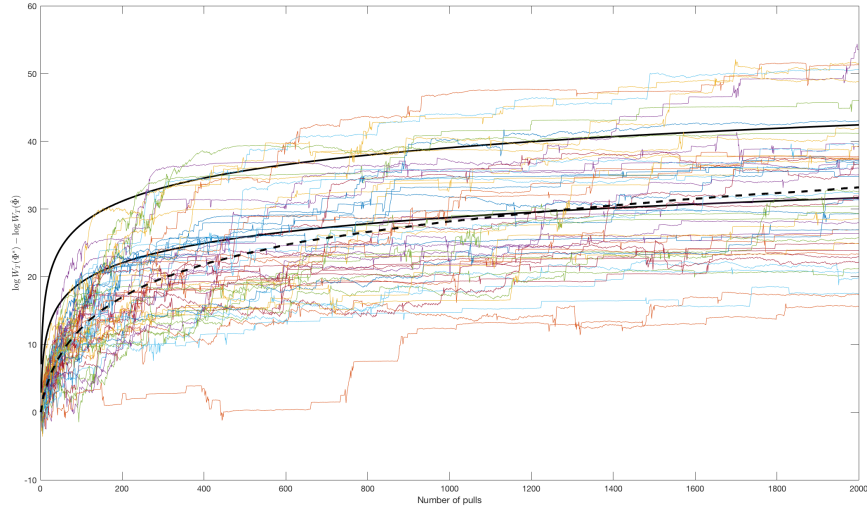


Figure 2:  $\vec{p} = (0.4, 0.5, 0.6, 0.8)$ . The top solid black curve is the expected regret bound(3.32). The bottom solid black curve is the lower bound we get if we always use the Kelly bet (i.e. the regret only accumulates from choosing the wrong arm). The dashed black line is the mean regret of 2000 paths. The other plots are 50 randomly chosen regret paths.



## References

- S. Bubeck, N. Cesa-Bianchi, et al. Regret analysis of stochastic and nonstochastic multi-armed bandit problems. *Foundations and Trends® in Machine Learning*, 5(1):1–122, 2012.
- O. Cappé, A. Garivier, O.-A. Maillard, R. Munos, G. Stoltz, et al. Kullback–leibler upper confidence bounds for optimal sequential allocation. *The Annals of Statistics*, 41(3):1516–1541, 2013.
- N. Cesa-Bianchi and G. Lugosi. *Prediction, learning, and games*. Cambridge university press, 2006.
- S. Ethier. *The doctrine of chances: probabilistic aspects of gambling*. Springer Science & Business Media, 2010.
- D. Glimne. Slot machine. <https://www.britannica.com/topic/slot-machine>, 2015. [Online; accessed 2017-10-28].
- J. L. Kelly. A new interpretation of information rate. *Bell Labs Technical Journal*, 35(4):917–926, 1956.
- R. Krichevsky and V. Trofimov. The performance of universal encoding. *IEEE Transactions on Information Theory*, 27(2):199–207, 1981.
- F. Orabona and D. Pál. Coin betting and parameter-free online learning. In *Advances in Neural Information Processing Systems*, pages 577–585, 2016.
- E. O. Thorp. The kelly criterion in blackjack, sports betting and the stock market. *Handbook of asset and liability management*, 1:385–428, 2006.