

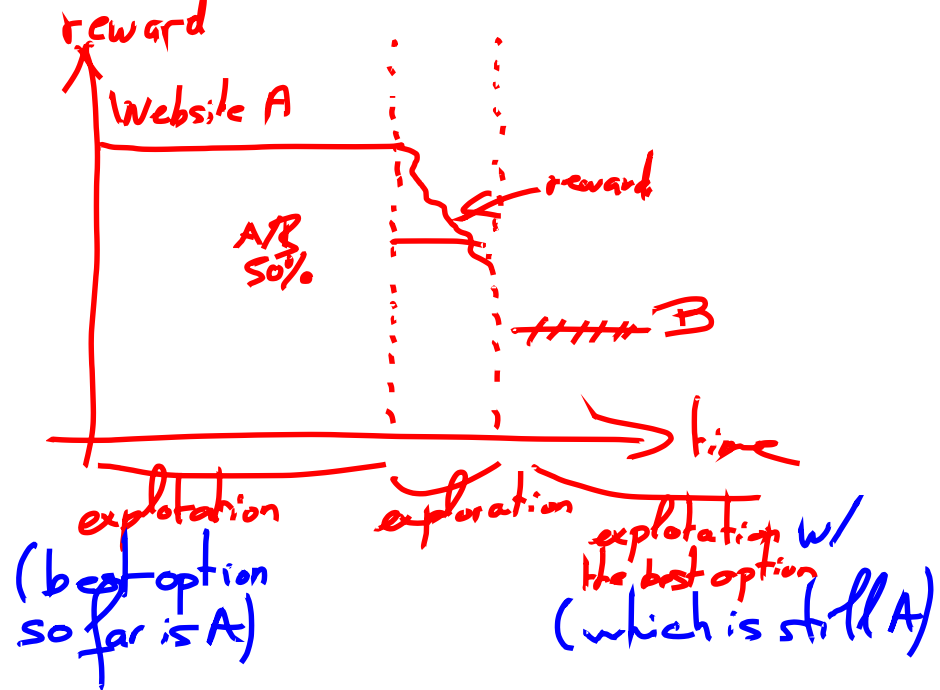
# Multi-Arm Bandit

~~N - Arm Bandit~~  
~~K - Arm Bandit~~  
Cary Goltermann

Galvanize

2016

Bandit: group of ideas to try  
Arm: one idea to try  
- pull/play/trial: one chance to try an idea  
- reward: unit of success



## Multi-Arm Bandit

- Motivation
- Exploration vs. Exploitation
- Formalization
- Regret

## Strategies

- Epsilon-Greedy
- UCB1
- Softmax
- Bayesian Bandit

## Multi-Arm Bandit

- **Motivation**
- Exploration vs. Exploitation
- Formalization
- Regret

## Strategies

- Epsilon-Greedy
- UCB1
- Softmax
- Bayesian Bandit

Consider the following scenario. (You may have seen something like this before.)

- The company you work for is testing out a new version of it's website.
- After running an A/B test for an afternoon the new version of the site appears to performing slightly better than the old version.
- After running the test over night, the new version of the site is performing better than the old version with a statistically significant p-value of 0.04.

Do you stop the test, or do you keep running it?

## Multi-Arm Bandit

- Motivation
- **Exploration vs. Exploitation**
- Formalization
- Regret

## Strategies

- Epsilon-Greedy
- UCB1
- Softmax
- Bayesian Bandit

# Exploration vs. Exploitation

When we are trying to determine which option, from a potential pool, is the best option what we are faced with is an information gathering challenge.

- **Exploration:** Testing out the different options (of the website), to determine how good each one is. Acquiring more knowledge about the reward associated with the options.
- **Exploitation:** Leveraging your current knowledge about the options to get the highest expected reward at that time.

# Traditional A/B Testing

A/B testing in terms of exploration vs. exploitation:

- Start with **pure exploration**, in which the same number of users are assigned to see each of the options. This is the testing phase.
- Once the test is complete and some conclusion has been reached you switch to **pure exploitation**. In this phase all of the users see the version that chosen as a result of the test.



# Potential Inefficiencies

- Need to wait for the experiment to conclude for certain statistical guarantees to be provided.
- Only after the experiment concludes can we capitalize on a potentially better option.
- This wastes time - and money - showing users the suboptimal site.

# Multi-Arm Bandit Approach

- Show each user the site that you think is best **most** of the time. (The definition of most will be dictated by your strategy. These will be discussed shortly.)
- As the experiment runs and you send users to different sites, update your beliefs about each site.
- Run until a clear best site emerges.

Depending on your strategy you can balance **exploration** and **exploitation** instead of having to choose either one or the other.

- This problem was originally motivated from the problem that a gambler faces when deciding which slot machine (individually called one-armed bandits) to play at.
- Assuming that the gambler wants to be smart about their strategy they are faced, not only, with the choice of which machines to play, how many times they should play them and in what order.
- All they know is that each bandit will provide a reward when its lever is pulled according to its individual (and static) distribution.

The gambler's objective, then, is to **maximize** the sum of the rewards they receive through a series of lever pulls.

- Dynamic A/B Testing.
- Budget allocation amongst competing projects.
- Clinical trials.
- Adaptive routing in attempts to minimize network delays.
- Reinforcement learning.

## Multi-Arm Bandit

- Motivation
- Exploration vs. Exploitation
- **Formalization**
- Regret

## Strategies

- Epsilon-Greedy
- UCB1
- Softmax
- Bayesian Bandit

# Terminology

Arm unknown  $(\alpha, \beta)$

we don't know the distributions and use MAB to discover them

The model is given by a set of real distributions

$B = \{R_1, \dots, R_K\}$ , where each distribution is associated with a reward delivered by one of the  $K \in \mathbb{N}^+$  levers.

We will define the mean values associated with each of these distributions as  $\mu_1, \dots, \mu_K$ .

The agent plays one lever per round and observes the associated reward. The goal of the agent is to maximize the sum of the collective rewards, or alternatively, minimize the **regret**.

2 slot machines  $\mu_1 = .7$   
 $\mu_2 = .3$

machine 1 / arm 1 0.7 for 1 / 0.3 for 0  
2 0.3 for 1 / 0.7 for 0

try #1 arm 1 get 0 ) exploration  
arm 2 get 1

use 2 all the time exploitation  
best option is arm 2  
at this point

too little exploration being  
done here.

What's the trade off?

## Multi-Arm Bandit

- Motivation
- Exploration vs. Exploitation
- Formalization
- **Regret**

## Strategies

- Epsilon-Greedy
- UCB1
- Softmax
- Bayesian Bandit



# Regret

The regret,  $\rho$  that an agent experiences after  $T$  rounds is the difference between the reward expected reward sum associated with the optimal strategy and the sum of the rewards collected.

overall  $\rho = T\mu^* - \sum_{t=1}^T \hat{\rho}_t = \sum_{t=1}^T (\mu^* - \hat{\rho}_t)$

$\mu^*$  Maximal reward mean,  $\mu^* = \max_k \{\mu_k\}$ . "money left on table"

$\hat{\rho}_t$  The reward at time  $t$ . "at round  $t$ "

**Regret** can, therefore, be seen as a measure of how often you choose the suboptimal bandit. We can view this as a cost function we are trying to minimize.

# Zero-Regret Strategy

A **zero-regret strategy** is defined as one who's average regret per round,  $\rho/T$ , goes to zero in the limit where the number of rounds,  $T$  goes to infinity.

The interesting thing is that a zero-regret strategy does *not* guarantee that you will never choose a suboptimal outcome. Instead it guarantees that, as you continue to play you will tend to choose the optimal outcome.

$$\lim_{T \rightarrow +\infty} \frac{\bar{\rho}_T}{T} = 0$$

## Multi-Arm Bandit

- Motivation
- Exploration vs. Exploitation
- Formalization
- Regret

## Strategies

- Epsilon-Greedy
- UCB1
- Softmax
- Bayesian Bandit

What if we always want to explore a certain amount?

1 -  $\epsilon$  time (most of), play w/ the "best"  
the type arm so far  
 $\epsilon$  of time play any other machine  
randomly

What if we always want to explore a certain amount?

- **Explore** with some probability  $\epsilon$ , often chosen as 10%.
- All other times, **exploit**, a.k.a. choose the bandit with the best performance so far  $\rightarrow \underset{j=1,\dots,K}{argmax} \{\hat{\mu}_j(t)\}$ .

$\hat{\mu}_j(t)$  Our best guess for  $\mu_j$  at round  $t$ .

- Update our knowledge about  $\mu$ 's after each round.

Is this a zero-regret strategy? **No**

limitations of  $\epsilon$ -greedy:

Scenario #1

arm #1

pay

99%

arm #2

5%

not  
aware of  
relative  
performance

#2

arm #1

10%

arm #2

5%

## Multi-Arm Bandit

- Motivation
- Exploration vs. Exploitation
- Formalization
- Regret

## Strategies

- Epsilon-Greedy
- **UCB1**
- Softmax
- Bayesian Bandit

# UCB1 (Upper Confidence Bound)

The UCB1 algorithm greedily chooses the bandit with the highest  $\hat{\mu}(t)$  but with a clever additional factor that automatically balances exploration and exploitation.

The choice at round  $t$ , after each lever is pulled once for a baseline, is made as follows:

no.  
2 times  
8 times

$$\operatorname{argmax}_{j=1,\dots,K} \left\{ \hat{\mu}_j(t) + \sqrt{\frac{2\ln(t)}{n_j}} \right\}$$

Chernoff Hoeffding Bound

"autocontrolled exploration"

$\hat{\mu}_j(t)$  Best guess for the  $\mu_j$  at time  $t$ .

$n_j$  Number of times that bandit  $j$  has been pulled.

$t$  Number of rounds that have been played in total.

Is this a zero-regret strategy?

Yes

shrinkage goes get more information



explore options in proportion to their performance?

## Multi-Arm Bandit

- Motivation
- Exploration vs. Exploitation
- Formalization
- Regret

intuition:

## Strategies

- Epsilon-Greedy
- UCB1
- Softmax
- Bayesian Bandit

$$\begin{array}{cc} \text{arm 1} & 10\% \\ & 20\% \\ \text{arm 2} & \end{array}$$
$$P(2) = \frac{20\%}{10 + 20\%}$$

limitations:

trials 100 IB (performance)

The softmax creates a probability distribution over all the levers in proportion to how good we think each lever is.

Boltzmann distribution  
-E/kT

$$P_t(\text{choosing bandit } j) = \frac{e^{\hat{\mu}_j(t)/\tau}}{\sum_{j=1}^K e^{\hat{\mu}_j(t)/\tau}}$$

N e

E energy T temperature

same  
(independent of trials)

$\tau$ , tau

Temperature, controls the "randomness" of the distribution. Usually chosen around 0.001.

At each round the above distribution is sampled from and that bandit is chosen.

(related to SA, simulated annealing but w/ constant)

A member of probability matching algorithms, so-called for their property of creating the probability distribution.

Zero-sum strategy? Yes

## Multi-Arm Bandit

- Motivation
- Exploration vs. Exploitation
- Formalization
- Regret

## Strategies

- Epsilon-Greedy
- UCB1
- Softmax
- **Bayesian Bandit**

# Bayesian Bandit

We can use the Bayesian beta-binomial conjugate prior techniques used to model the click-through rate in the morning exercise to as our base model for each of the bandits.

This is another probability matching algorithm where we have a separate model for each of the bandits.

## Process

- Sample from the distributions for each bandit.
- Choose the bandit with the highest sampled  $\mu$
- Update the distribution of the chosen bandit with the knowledge gained from choosing it.

Zero sum strategy? Yes

# Bayesian Bandit

"True" bandit reward rates: 0.1 / 0.2 / 0.3

