

Lecture 4: Bandit Algorithms

Ziyu Shao

School of Information Science and Technology
ShanghaiTech University

March 30 & April 01, 2020

Outline

- 1 Introduction
- 2 Overview of Bandit Problems
- 3 Stochastic Bandits
- 4 Greedy Algorithms
- 5 UCB Algorithms
- 6 Thompson Sampling Algorithms
- 7 Gradient Bandit Algorithms
- 8 Adversarial Bandits
- 9 Summary
- 10 References

Outline

- 1 Introduction
- 2 Overview of Bandit Problems
- 3 Stochastic Bandits
- 4 Greedy Algorithms
- 5 UCB Algorithms
- 6 Thompson Sampling Algorithms
- 7 Gradient Bandit Algorithms
- 8 Adversarial Bandits
- 9 Summary
- 10 References

One-Armed Bandit

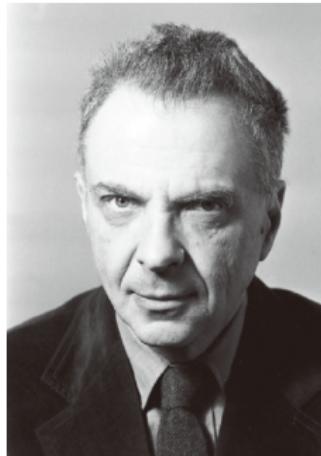


Multi-Armed Bandit



History

- 1933: William R. Thompson proposed the first algorithm (Thompson sampling method)
- 1952: Herbert Robbins (1915-2001) proposed the mathematical formulation of multi-armed bandit problem and sequential sampling algorithm for 2-armed bandit problem.



SOME ASPECTS OF THE SEQUENTIAL DESIGN
OF EXPERIMENTS

HERBERT ROBBINS

1. Introduction. Until recently, statistical theory has been restricted to the design and analysis of sampling experiments in which the size and composition of the sample are completely determined before the experimentation begins. The reasons for this are partly historical, dating back to the time when the statistician was consulted, if at all, only after the experiment was over, and partly intrinsic in the mathematical difficulty of working with anything but a finite number of independent variables. A new era appears now to be in the making with the creation of a theory of the sequential design of experiments, in which the size and composition of the samples are not fixed in advance but are functions of the observations themselves.

The first important departure from fixed sample size came in the field of industrial quality control, with the double sampling inspection method of Dodge and Romig [1]. Here there is only one population to be sampled, and the question at issue is whether the proportion of defectives in a lot is acceptably small. A preliminary sample of n_1 objects is drawn from the lot and the number x of defectives noted. If x is less than a fixed value b ($b < n_1$) the lot is accepted without further sampling, if x is greater than a fixed value a ($a < b$) the lot is rejected without further sampling, but if $a \leq x \leq b$ then a second sample of $n_2 - n_1 + x$ objects is drawn from the lot, and a decision as to the lot is made on the basis of the number of defectives in the total sample of $n_1 + x$ objects. The total sample size n is thus a random variable with two values, n_1 and $n_1 + n_2$, and the value of n is stochastically dependent on the value of x . A general extension of the idea of double sampling came during World War II, with the development, chiefly by Wald, of sequential analysis [2], in which the observations are made one by one and the decision to terminate sampling and to accept or reject the lot (or, more generally, to accept or reject whatever statistical "null hypothesis" is being tested) can come at any stage. The sample size is not fixed in advance, but is capable in principle of assuming infinitely many values, although in practice a finite upper limit on n is usually set. The advantage of sequential

An address delivered before the Akron, Alabama, meeting of the Society, November 23, 1951, by invitation of the Committee to Select Distinguished Speakers for Southeastern Sectional Meetings; revised by the editor December 16, 1951.

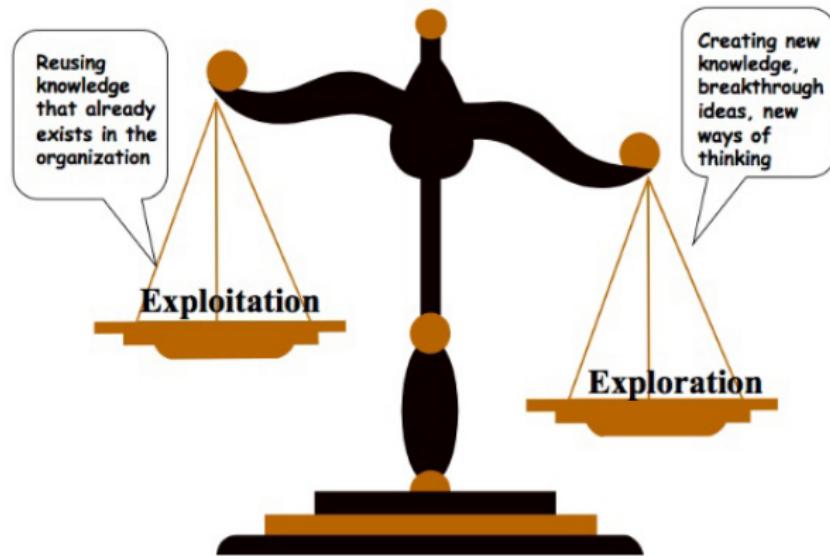
Why Bandit?

- A perfect model for online decision making under uncertainty
- Isolate an important component of reinforcement learning:
exploration-vs-exploitation
- Rich and beautiful mathematically
- Widely applications

Exploration vs. Exploitation Dilemma

- Online decision-making involves a fundamental choice
 - ▶ **Exploitation:** staying with the option that gave the highest payoffs in the past
 - ▶ **Exploration:** exploring new options that might give higher payoffs in the future
- The best long-term strategy may involve short-term sacrifices
- Gather enough information to make the best overall decisions

Exploration vs. Exploitation Dilemma



Examples in Real Life

- Restaurant Selection
 - ▶ **Exploitation:** Go to your favorite restaurant
 - ▶ **Exploration:** Try a new restaurant
- Online Banner Advertisements
 - ▶ **Exploitation:** Show the most successful advert
 - ▶ **Exploration:** Show a different advert
- Oil Drilling
 - ▶ **Exploitation:** Drill at the best known location
 - ▶ **Exploration:** Drill at a new location

Example of A Two-Armed Bandit



Time	1	2	3	4	5	6	7	8	9	10	11	12
Left arm	\$1	\$0			\$1	\$1	\$0					
Right arm			\$1	\$0								

average $\frac{3}{5} = 0.6$. exploitation left.

$\frac{1}{2} = 0.5$. exploration right.

Widely Applications

- Several companies including Adobe, Amazon, Facebook, Google, LinkedIn, Microsoft, Netflix, and Twitter
- Clinical trials/dose discovery
- Recommendation systems (movies/news/etc)
- Advertisement placement
- A/B testing
- Monte Carlo tree search algorithm (game playing including AlphaGo)
- Resource allocation
- Network routing
- Dynamic pricing (e.g.,for Amazon products)
- Ranking (e.g.,for search)

Remark



Outline

- 1 Introduction
- 2 Overview of Bandit Problems
- 3 Stochastic Bandits
- 4 Greedy Algorithms
- 5 UCB Algorithms
- 6 Thompson Sampling Algorithms
- 7 Gradient Bandit Algorithms
- 8 Adversarial Bandits
- 9 Summary
- 10 References

Bandit: Special Case of Reinforcement Learning

- What is Reinforcement Learning?
- **Wikipedia:** reinforcement learning is an area of machine learning inspired by behavioral psychology, concerned with how software agents ought to take actions in an environment so as to maximize some notion of cumulative reward.

Reward: Important Concept

- A **reward** R_t is a scalar feedback signal
- Indicates how well agent is doing at step t
- The agent's job is to maximize cumulative reward

Reinforcement learning is based on the reward hypothesis:

Definition

All goals can be described by the maximization of expected cumulative reward

Why Expectation of Rewards?

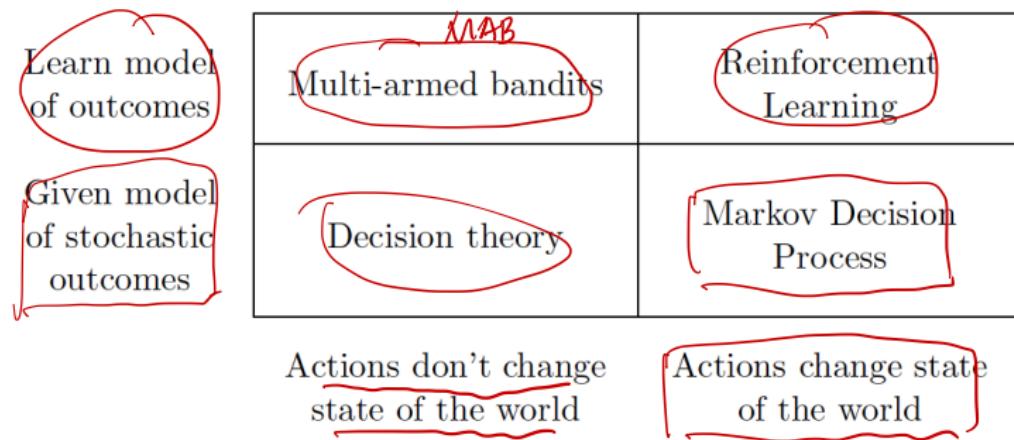
- Decision makers who base their decisions on expected values are called "risk-neutral"
- Von Neumann-Morgenstern expected utility theorem: with axioms of rational behavior under uncertainty
- A rational decision maker must choose amongst alternatives by computing the expected utility of the outcomes.

Reinforcement Learning Problems



Goal: Learn to choose actions that maximize rewards

Decision Making Under Uncertainty



Three Running Examples of Bandits

- **News website.** When a new user arrives, a website picks an article header to show, observes whether the user clicks on this header. The site's goal is maximize the total number of clicks.
- **Dynamic pricing.** A store is selling a digital good, e.g., an app or a song. When a new customer arrives, the store chooses a price offered to this customer. The customer buys (or not) and leaves forever. The store's goal is to maximize the total profit.
- **Investment.** Each morning, you choose one stock to invest into, and invest \$1. In the end of the day, you observe the change in value for each stock. The goal is to maximize the total wealth.

Examples: Action & Reward

Example	Action	Reward
News website	an article to display	1 if clicked, 0 otherwise
Dynamic pricing	a price to offer	p if sale, 0 otherwise
Investment	a stock to invest into	change in value during the day

Problem Space

- Auxiliary Feedback
- Rewards Model
- Contexts
- Bayesian Priors
- Structured Rewards
- Global Constraints
- Structured Actions

Problem Space: Auxiliary Feedback

Example	Auxiliary feedback	Rewards for any other arms?
News website	N/A	no (<i>bandit feedback</i>). yes, for some arms, but not for all (<i>partial feedback</i>).
Dynamic pricing	sale \Rightarrow sale at any lower price, no sale \Rightarrow no sale at any higher price	
Investment	change in value for all other stocks	yes, for all arms (<i>full feedback</i>).

Three Types of Feedback

- **Bandit Feedback:** when the algorithm observes the reward for the chosen arm, and no other feedback
- **Full Feedback:** when the algorithm observes the rewards for all arms that could have been chosen
- **Partial Feedback:** when some information is revealed, in addition to the reward of the chosen arm, but it does not always amount to full feedback.

Problem Space: Rewards Model

- **IID Rewards:** the reward for each arm is drawn independently from a fixed distribution that depends on the arm but not on the round t .
- **Adversarial Rewards:** rewards can be arbitrary, as if they are chosen by an “adversary” that tries to fool the algorithm.
- **Constrained Adversary:** rewards are chosen by an adversary that is subject to some constraints, e.g., reward of each arm cannot change much from one round to another, or the reward of each arm can change at most a few times, or the total change in rewards is upper-bounded.
- **Stochastic Rewards (beyond IID):** rewards evolves over time as a random process, e.g., a random walk.

Problem Space: Contexts

- In each round, an algorithm may observe some context before choosing an action.
- Such context often comprises the known properties of the current user
- Allows for personalized actions

Example	Context
News website	<u>user location and demographics</u>
Dynamic pricing	<u>customer's device (e.g., cell or laptop)</u> , location, demographics
Investment	<u>current state of the economy.</u>

Problem Space: Contexts

- The algorithm has a different high-level objective.
- It is no longer interested in learning one good arm, since any one arm may be great for some contexts, and terrible for some others.
- Instead, it strives to learn the best policy which maps contexts to arms
- While not spending too much time exploring.

Problem Space: Bayesian Priors

- Each problem can be studied under a Bayesian approach, whereby the problem instance comes from a known distribution (called Bayesian prior).
- One is typically interested in provable guarantees in expectation over this distribution.

Problem Space: Structured Rewards

- Rewards may have a known structure, e.g., arms correspond to points in \mathbb{R}^d ,
- And in each round the reward is a linear (resp., concave or Lipschitz) function of the chosen arm.

Problem Space: Global Constraints

- The algorithm can be subject to global constraints that bind across arms and across rounds.
- For example, in dynamic pricing there may be a limited inventory of items for sale.

Problem Space: Structured Actions

- An algorithm may need to make several decisions at once,
- e.g., a news website may need to pick a slate of articles,
- and a seller may need to choose prices for the entire slate of offerings.

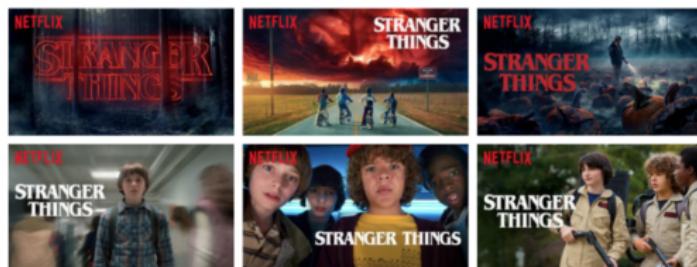
Application Domains

Application domain	Action	Reward
medical trials	which drug to prescribe	health outcome.
web design	e.g., font color or page layout	#clicks.
content optimization	which items/articles to emphasize	#clicks.
web search	search results for a given query	1 if the user is satisfied.
advertisement	which ad to display	revenue from ads.
recommender systems	e.g., which movie to watch	1 if follows recommendation.
sales optimization	which products to offer at which prices	revenue.
procurement	which items to buy at which prices	#items procured
auction/market design	e.g., which reserve price to use	revenue
crowdsourcing	which tasks to give to which workers, and at which prices	1 if task completed at sufficient quality.
datacenter design	e.g., which server to route the job to	job completion time.
Internet	e.g., which TCP settings to use?	connection quality.
radio networks	which radio frequency to use?	1 if successful transmission.
robot control	a “strategy” for a given task	job completion time.

Real-World Example: Netflix

For a particular movie, we want to decide what image to show (to all the NETFLIX users)

- **Actions:** uploading one of the K movie posters to a user's home screen
- **Unknown true mean rewards:** the % of NETFLIX users that will click on the title and watch the movie
- **Estimated mean rewards:** the average click rate observed



Outline

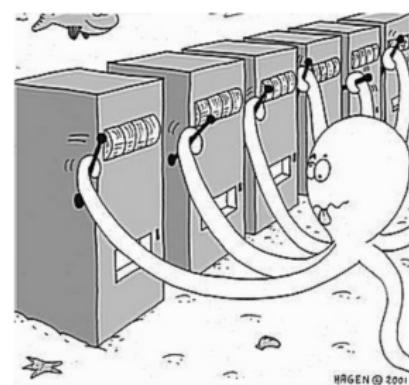
- 1 Introduction
- 2 Overview of Bandit Problems
- 3 Stochastic Bandits
- 4 Greedy Algorithms
- 5 UCB Algorithms
- 6 Thompson Sampling Algorithms
- 7 Gradient Bandit Algorithms
- 8 Adversarial Bandits
- 9 Summary
- 10 References

Stochastic Bandits

- Bandit feedback: observes only the reward for the selected (arm) action, and nothing else.
 - IID reward for each action: every time one action a is chosen, the reward is sampled independently from the reward distribution associated with the action a .
 - Per-round rewards are bounded.

Stochastic Bandits

- A multi-armed bandit is a tuple $\langle \mathcal{A}, \mathcal{R} \rangle$
- \mathcal{A} is a known set of m actions (or "arms")
- $\mathcal{R}^a(r) = \mathbb{P}[r|a]$ is an unknown probability distribution over rewards
- At each step t the agent selects an action $a_t \in \mathcal{A}$
- The environment generates a reward $r_t \sim \mathcal{R}^{a_t}$
- The goal is to maximize expectation of cumulative reward within t rounds $E(\sum_{\tau=1}^t r_\tau)$



Regret: Equivalent Goal

- The value of an arbitrary action a (action-value) is the mean reward for action a (unknown):

$$q_*(a) = \underline{Q(a)} = E[r_t | a_t = a], \forall a \in \{1, 2, \dots, K\}$$

a real number

- The optimal value V^* is

$$V^* = \underline{Q(a^*)} = \max_{a \in \mathcal{A}} Q(a)$$

a real number

- The regret is the opportunity loss for one round τ

$$L_\tau = \mathbb{E}[V^* - \underline{Q(a_\tau)}]$$

where

$$Q(a_\tau) = \underline{E[r_\tau | a_\tau]}.$$

Random Variable

Regret: Equivalent Goal

$$\begin{aligned} \text{1. } L_t &= \sum_{\tau=1}^t V^* - \sum_{\tau=1}^t E[\mathcal{L}(a_\tau)] \\ &= tV^* - \sum_{\tau=1}^t E[E[r_\tau | a_\tau]] = tV^* - \sum_{\tau=1}^t E[r_\tau] \end{aligned}$$

- The total regret is the total opportunity loss after the end of round t

$$= tV^* - \underbrace{E\left[\sum_{\tau=1}^t r_\tau\right]}$$

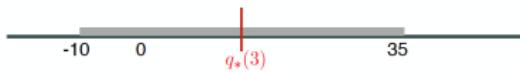
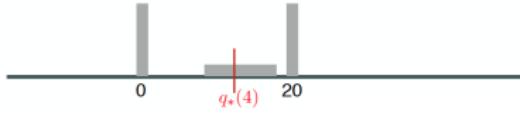
$$\underline{L_t} = \mathbb{E} \left[\sum_{\tau=1}^t \{V^* - Q(a_\tau)\} \right] \Rightarrow \text{minimize } L_t$$

\equiv maximize

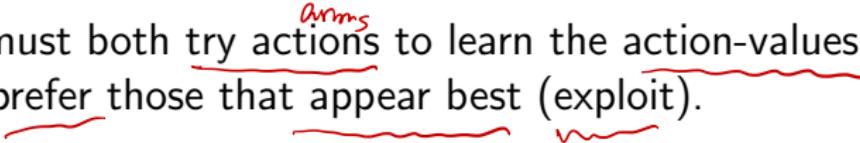
$$\underline{E\left[\sum_{\tau=1}^t r_\tau\right]}$$

- Maximize the Expectation of Cumulative Reward \equiv
Minimize Total Regret

Example on Action-Value

- Action 1 — Reward is always 8
 - value of action 1 is $q_*(1) = 8$
- Action 2 — 88% chance of 0, 12% chance of 100!
 - value of action 2 is $q_*(2) = .88 \times 0 + .12 \times 100 = 12$
- Action 3 — Randomly between -10 and 35, equiprobable
$$q_*(3) = \frac{45}{2} = 22.5$$
- Action 4 — a third 0, a third 20, and a third from {8,9,...,18}
$$q_*(4) = \frac{1}{3} \cdot 0 + \frac{1}{3} \cdot 20 + \frac{1}{3} \left(\frac{8}{11} + \frac{9}{11} + \dots + \frac{18}{11} \right) = 11$$

Example on Action-Value

- In general, those values are unknown.
- We must both try actions to learn the action-values (explore),
and prefer those that appear best (exploit).


The Exploration/Exploitation Dilemma

1^o. $Q_t(a) = \hat{Q}_t(a)$ Unknown

$$\hat{Q}_t(a) \approx Q_t(a)$$

estimated value:

2^o. Define the greedy action at time t as

$$a_t^* = \underset{a \in A}{\operatorname{argmax}} \hat{Q}_t(a)$$

If $a_t = a_t^*$, you are making exploitation

$a_t \neq a_t^*$, you are making exploration

Action-Value Methods

$$\hat{Q}_t(a) = \frac{\sum_{\tau=1}^{t-1} r_\tau I_{a_\tau=a}}{\sum_{\tau=1}^{t-1} I_{a_\tau=a}} \xrightarrow{SLN} Q(a)$$

- Methods that learn action-value estimates and nothing else
- For example, estimate action values as sample average
- The sample-average estimates converge to the true values if the action is taken an infinite number of times

Counting Regret

- The count $N_t(a)$ is the number of selections for action a after the end of round t : $N_t(a) = \sum_{\tau=1}^t \mathbb{I}_{a_\tau=a}$
- The gap Δ_a is the difference in value between action a and optimal action a^* , $\Delta_a = V^* - Q(a)$ $\Delta_a \geq 0$
- Regret is a function of gaps and the counts

$$\begin{aligned} L_t &= \mathbb{E} \left[\sum_{\tau=1}^t (V^* - Q(a_\tau)) \right] \\ &= \sum_{a \in \mathcal{A}} \mathbb{E}[N_t(a)](V^* - Q(a)) \\ &= \sum_{a \in \mathcal{A}} \mathbb{E}[N_t(a)]\Delta_a \end{aligned}$$

- A good algorithm ensures small counts for large gaps
- Problem: gaps are not known!

Regret Decomposition Lemma

Almost all bandit proofs are based on this lemma.

Lemma

$$L_t = \sum_{a \in \mathcal{A}} \mathbb{E}[N_t(a)] \Delta_a$$

Proof

1^o. $\Omega(a_2) = E[r_2 | a_2]$ is a random variable

$$E[\Omega(a_2)] = E[E[r_2 | a_2]] = E[r_2]$$

2^o. Let $S_t \triangleq \sum_{z=1}^t \Omega(a_z)$ then we have $E(S_t) = \sum_{z=1}^t E[\Omega(a_z)] = \sum_{z=1}^t E[r_z] = E\left(\sum_{z=1}^t r_z\right)$

Since $\sum_{a \in A} \mathbb{1}_{\{a_z=a\}} = 1$ for any fixed z .

$$\begin{aligned} \Rightarrow E(S_t) &= E\left(\sum_{z=1}^t r_z\right) = E\left[\sum_{z=1}^t \sum_{a \in A} r_z \mathbb{1}_{\{a_z=a\}}\right] \\ &= E\left[\sum_{a \in A} \sum_{z=1}^t r_z \mathbb{1}_{\{a_z=a\}}\right] = \sum_{a \in A} \sum_{z=1}^t E[r_z \mathbb{1}_{\{a_z=a\}}] \quad (*1) \end{aligned}$$

3^o. On the other hand, $\sum_{z=1}^t \left(\sum_{a \in A} \mathbb{1}_{\{a_z=a\}} \right) = \sum_{z=1}^t 1 = t$

$$\Rightarrow E\left[\sum_{z=1}^t \sum_{a \in A} \mathbb{1}_{\{a_z=a\}}\right] = t$$

$$\Rightarrow \sum_{z=1}^t \sum_{a \in A} E[\mathbb{1}_{\{a_z=a\}}] = t \quad \Rightarrow \sum_{a \in A} \sum_{z=1}^t E[\mathbb{1}_{\{a_z=a\}}] = t \quad (*2)$$

Proof

4°. Now the total regret

$$L_t = E \left[\sum_{z=1}^t (v^* - Q(a_z)) \right] = t v^* - E \left(\sum_{z=1}^t Q(a_z) \right)$$

$$= t v^* - E(S_t)$$

$$\stackrel{(*)1)}{=} t v^* - \sum_{a \in A} \sum_{z=1}^t E[r_z \mathbb{1}_{\{a_z=a\}}]$$

$$\stackrel{(*)2)}{=} \sum_{a \in A} \sum_{z=1}^t E[\mathbb{1}_{\{a_z=a\}} \cdot v^*] - \sum_{a \in A} \sum_{z=1}^t E[r_z \mathbb{1}_{\{a_z=a\}}]$$

$$= \sum_{a \in A} \sum_{z=1}^t E[(v^* - r_z) \mathbb{1}_{\{a_z=a\}}] \quad (*)3)$$

v^* is a real number.

$$5°. E[(v^* - r_z) \mathbb{1}_{\{a_z=a\}} | a_z] = \mathbb{1}_{\{a_z=a\}} E[v^* - r_z | a_z]$$

$$= \mathbb{1}_{\{a_z=a\}} (v^* - E(r_z | a_z)) = \mathbb{1}_{\{a_z=a\}} (v^* - Q(a_z))$$

$$= \mathbb{1}_{\{a_z=a\}} \cdot (v^* - Q(a)) = \mathbb{1}_{\{a_z=a\}} \cdot \Delta_a$$

Proof

thus we have $E[(v^* - r_2) \perp_{\{a_2=a\}} | a_2] = \perp_{\{a_2=a\}} \cdot \Delta_a$

$$\Rightarrow E[(v^* - r_2) \perp_{\{a_2=a\}}] = E[\perp_{\{a_2=a\}} \cdot \Delta_a]$$

6°. Then we have

$$\begin{aligned} L_t &\stackrel{(*)}{=} \sum_{a \in A} \sum_{z=1}^t E[(v^* - r_z) \perp_{\{a_2=a\}}] \\ &= \sum_{a \in A} \sum_{z=1}^t E[\perp_{\{a_2=a\}} \cdot \Delta_a] \\ &= \sum_{a \in A} E\left[\sum_{z=1}^t \perp_{\{a_2=a\}}\right] \cdot \Delta_a \\ &= \sum_{a \in A} E(H_t(a)) \cdot \Delta_a \end{aligned}$$

D.E.D.

Landau Notation

- Given functions $f, g : \mathbb{N} \rightarrow [0, \infty)$, we define

$$f(n) = O(g(n)) \Leftrightarrow \limsup_{n \rightarrow \infty} \frac{f(n)}{g(n)} < \infty.$$

$$f(n) = o(g(n)) \Leftrightarrow \lim_{n \rightarrow \infty} \frac{f(n)}{g(n)} = 0.$$

$$f(n) = \Omega(g(n)) \Leftrightarrow \liminf_{n \rightarrow \infty} \frac{f(n)}{g(n)} > 0.$$

$$f(n) = \omega(g(n)) \Leftrightarrow \liminf_{n \rightarrow \infty} \frac{f(n)}{g(n)} = \infty.$$

$$f(n) = \Theta(g(n)) \Leftrightarrow f(n) = O(g(n)) \text{ and } f(n) = \Omega(g(n)).$$

Good Learners

- Let L_n denote the regret over n rounds
- A good learner achieves sublinear regret: $L_n = o(n)$ or equivalently $\lim_{n \rightarrow \infty} \frac{L_n}{n} = 0$

Linear or Sublinear Regret

- If an algorithm **forever** explores it will have linear total regret
- If an algorithm **never** explores it will have linear total regret
- Is it possible to achieve sublinear total regret?

Outline

- 1 Introduction
- 2 Overview of Bandit Problems
- 3 Stochastic Bandits
- 4 Greedy Algorithms
- 5 UCB Algorithms
- 6 Thompson Sampling Algorithms
- 7 Gradient Bandit Algorithms
- 8 Adversarial Bandits
- 9 Summary
- 10 References

Greedy Algorithm

- We consider algorithms that estimate $\hat{Q}_t(a) \approx Q(a)$
- Estimate the value of each action by Monte-Carlo evaluation

$$\hat{Q}_t(a) = \frac{1}{N_{t-1}(a)} \sum_{\tau=1}^{t-1} r_\tau \mathbf{1}(a_\tau = a)$$

- The *greedy* algorithm selects action with highest value

$$a_t^* = \arg \max_{a \in \mathcal{A}} \hat{Q}_t(a)$$

- Greedy can lock onto a suboptimal action forever
- \Rightarrow Greedy has linear total regret

ϵ -Greedy Algorithm

- The ϵ -greedy algorithm continues to explore forever
 - ▶ With probability $1 - \epsilon$ select $a^* = \arg \max_{a \in \mathcal{A}} \hat{Q}(a)$
 - ▶ With probability ϵ select a random action
- The we have the lower bound of the total regret

$$L_t \geq \frac{t \cdot \epsilon}{|\mathcal{A}|} \sum_{a \in \mathcal{A}} \Delta_a$$

- $\Rightarrow \epsilon$ -greedy has linear total regret

Proof of the Lower Bound

1^o. For any fixed a , $P(a_2 = a^*) = 1 - \varepsilon + \frac{\varepsilon}{|A|}$, $P(a_2 = a) = \frac{\varepsilon}{|A|}$ $\forall a \neq a^*$.
 $\Rightarrow \forall a \in A$, $P(a_2 = a) \geq \frac{\varepsilon}{|A|}$.

$$\begin{aligned} 2^o. \quad N_t(a) &= \sum_{i=1}^t \mathbb{I}_{\{a_i=a\}} \quad \Rightarrow E(N_t(a)) = E\left(\sum_{i=1}^t \mathbb{I}_{\{a_i=a\}}\right) = \sum_{i=1}^t E(\mathbb{I}_{\{a_i=a\}}) \\ &= \sum_{i=1}^t P(a_i=a) \geq \sum_{i=1}^t \frac{\varepsilon}{|A|} = \frac{t\varepsilon}{|A|}. \end{aligned}$$

3^o. By regret decomposition lemma,

$$L_t = \sum_{a \in A} E[N_t(a)] \Delta_a \geq \underbrace{\sum_{a \in A} \frac{t\varepsilon}{|A|} \Delta_a}_{\text{Linear total regret.}} = \frac{t\varepsilon}{|A|} \sum_{a \in A} \Delta_a$$

Incremental Implementation

$$\hat{Q}_n = \frac{\overbrace{R_1 + R_2 + \dots + R_n}^n}{n}$$

incremental implementation

$$\Rightarrow \hat{Q}_{n+1} = \hat{Q}_n + \frac{1}{n} (R_n - \hat{Q}_n)$$

Target.

old estimate.

new estimate

old estimate.

Step size

Stochastic Approximation

$$\hat{Q}_{n+1} = \hat{Q}_n + \alpha_n [R_n - \hat{Q}_n]$$

With mild conditions

$$\left. \begin{aligned} \sum_{n=1}^{\infty} \alpha(n) &= \infty \\ \sum_{n=1}^{\infty} \alpha^2(n) &< \infty \end{aligned} \right\} \Rightarrow \text{Converges w.p. 1.}$$

e.g. $\alpha_n = \frac{1}{n}$

Derivation of Incremental Update

$$1^{\circ}. \quad \widehat{Q}_{n+1} = \frac{1}{n} \sum_{i=1}^n R_i = \frac{1}{n} (R_n + \sum_{i=1}^{n-1} R_i)$$

$$2^{\circ}. \quad \widehat{Q}_n = \frac{1}{n-1} \sum_{i=1}^{n-1} R_i \Rightarrow \sum_{i=1}^{n-1} R_i = (n-1) \widehat{Q}_n$$

$$\begin{aligned} 3^{\circ}. \quad \widehat{Q}_{n+1} &= \frac{1}{n} (R_n + \sum_{i=1}^{n-1} R_i) \\ &= \frac{1}{n} (R_n + (n-1) \widehat{Q}_n) \\ &= \frac{1}{n} [R_n - \widehat{Q}_n + n \widehat{Q}_n] \\ &= \frac{1}{n} [R_n - \widehat{Q}_n] + \widehat{Q}_n \\ &= \widehat{Q}_n + \frac{1}{n} (R_n - \widehat{Q}_n) \end{aligned}$$

ϵ -Greedy Bandit Algorithm

A simple bandit algorithm

Initialize, for $a = 1$ to k :

$$\hat{Q}(a) \leftarrow 0$$
$$N(a) \leftarrow 0$$

\hat{Q} : to emphasize it is an estimate value.

Repeat forever:

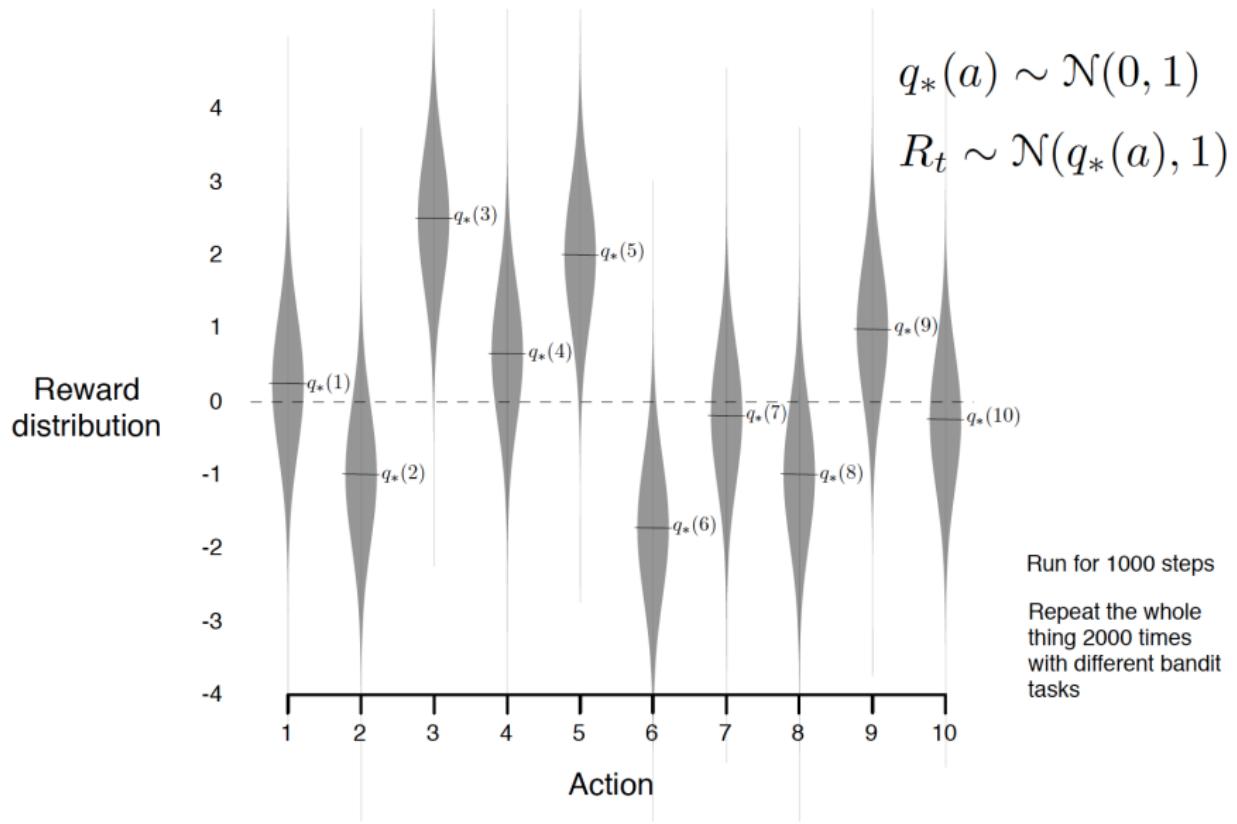
$$A \leftarrow \begin{cases} \arg \max_a \hat{Q}(a) & \text{with probability } 1 - \varepsilon \\ \text{a random action} & \text{with probability } \varepsilon \end{cases} \quad (\text{breaking ties randomly})$$

$$R \leftarrow \text{bandit}(A)$$

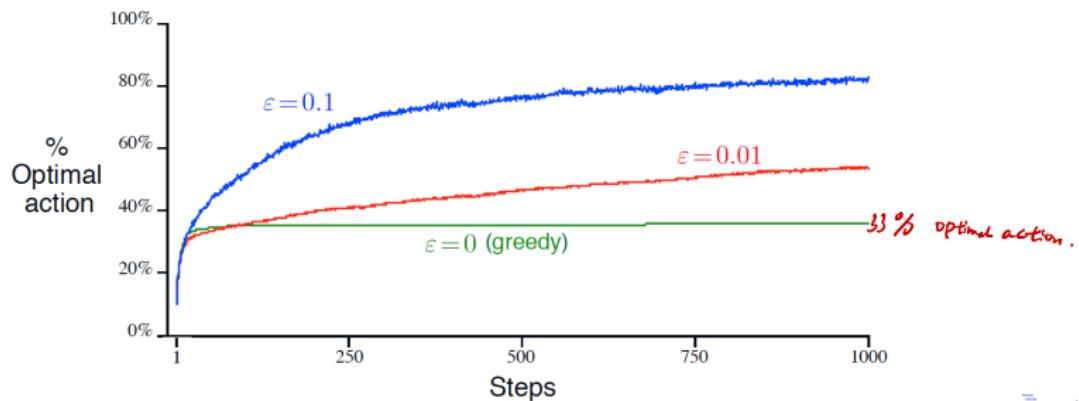
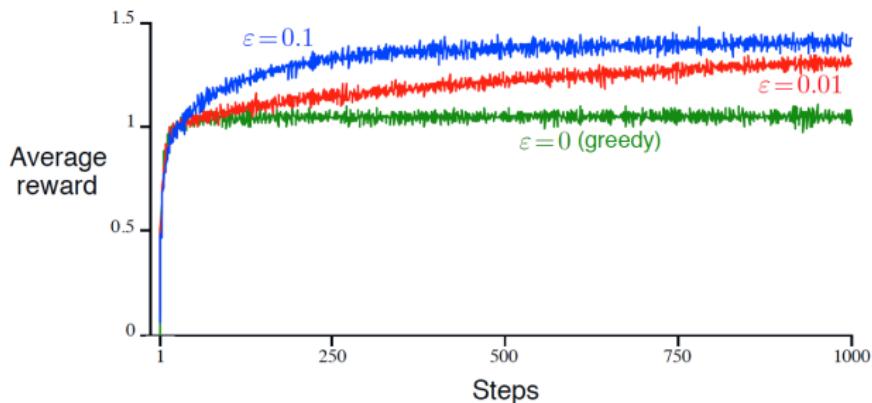
$$N(A) \leftarrow N(A) + 1$$

$$\hat{Q}(A) \leftarrow \hat{Q}(A) + \frac{1}{N(A)} [R - \hat{Q}(A)]$$

Example: 10-armed Testbed

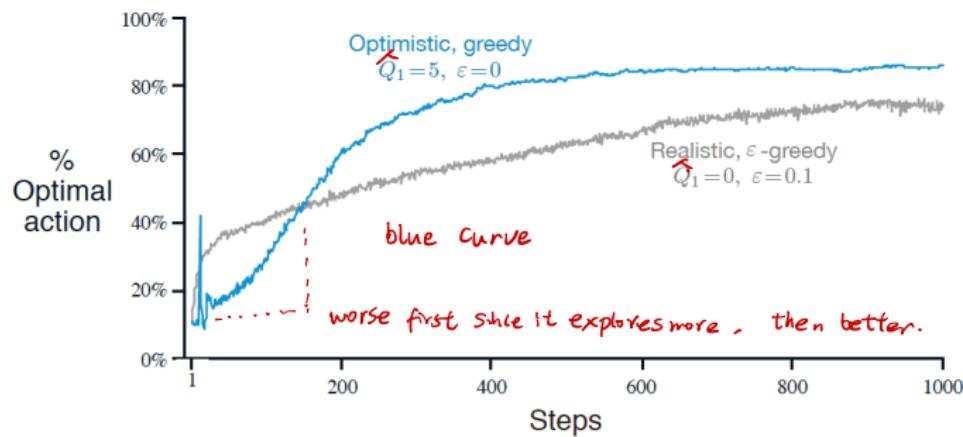


Performance of ϵ -Greedy



Optimistic Initial Values

- All methods so far depend on $\hat{Q}_1(a)$, i.e., the initial action-value estimate. So far we have used $\hat{Q}_1(a) = 0$.
- Suppose we initialize the action values optimistically ($\hat{Q}_1(a) = 5$), e.g., on the 10-armed testbed (with $\alpha = 0.1$)



Optimistic Initialization

- Simple and practical idea: initialize $\hat{Q}(a)$ to high values
- Update action value by incremental Monte-Carlo evaluation
- Starting with $N(a) > 0$

$$\hat{Q}_{t+1}(a) = \hat{Q}_t(a) + \frac{1}{N_t(a)} (r_t - \hat{Q}_t(a))$$

- Encourages systematic exploration early on
- But can still lock onto suboptimal action
- \Rightarrow greedy + optimistic initialization has linear total regret
- \Rightarrow ϵ -greedy + optimistic initialization has linear total regret

Decaying ϵ_t -Greedy Algorithm

- Pick a decay schedule for $\epsilon_1, \epsilon_2, \dots$

$$c > 0$$

$$d = \min_{a | \Delta_a > 0} \Delta_a$$

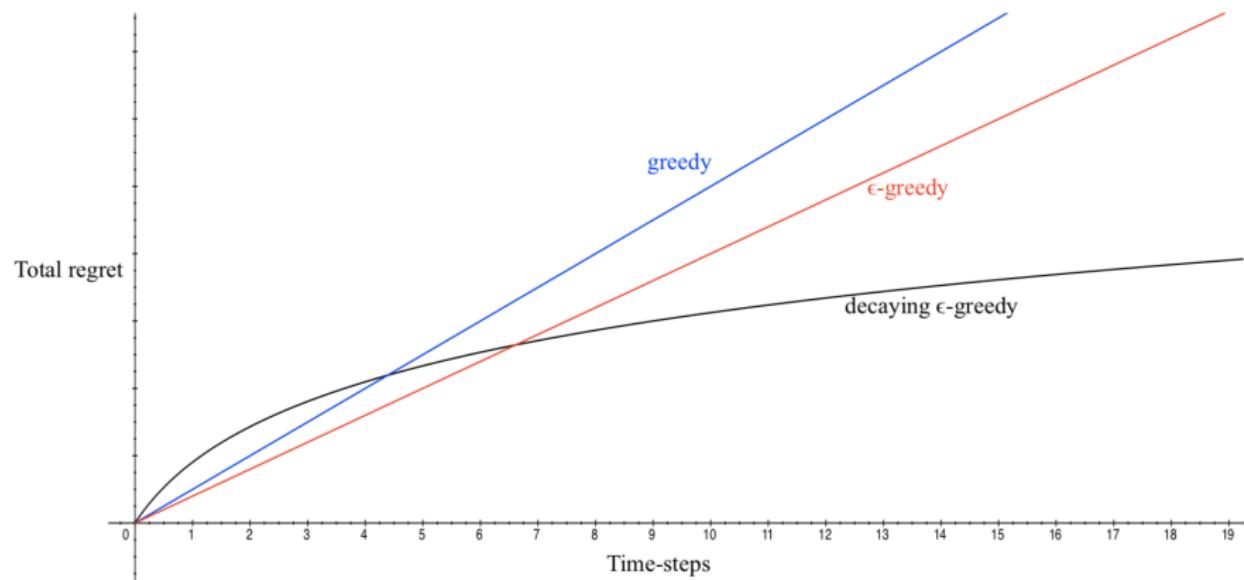
$as t \uparrow, \epsilon_t \downarrow$

$$\epsilon_t = \min \left\{ 1, \frac{c|\mathcal{A}|}{d^2 t} \right\}$$

$t \epsilon_t \uparrow$
explore first, then more and more
exploit

- Consider the following schedule
- Decaying ϵ_t -greedy has *logarithmic* asymptotic total regret!
- Unfortunately, schedule requires advance knowledge of gaps
- Goal: find an algorithm with sublinear regret for any multi-armed bandit (without knowledge of \mathcal{R})

Performance Comparison



Lower Bound

- The performance of any algorithm is determined by similarity between optimal arm and other arms
- This is described formally by the gap Δ_a and the similarity in distributions $KL(\mathcal{R}^a \parallel \mathcal{R}^{a^*})$

Theorem (Lai and Robbins)

Asymptotic total regret is at least logarithmic in number of steps

$$\lim_{t \rightarrow \infty} L_t \geq \log t \sum_{a | \Delta_a > 0} \frac{\Delta_a}{KL(\mathcal{R}^a \parallel \mathcal{R}^{a^*})}$$

Revisit Kullback-Leibler Divergence

- A finite sample space Ω
- Two probability distributions on Ω : \mathbf{p} & \mathbf{q}
- Kullback-Leibler Divergence (KL-divergence) is defined as

$$KL(\mathbf{p}, \mathbf{q}) = \sum_{x \in \Omega} p(x) \ln \frac{p(x)}{q(x)} = \mathbb{E} \left[\ln \frac{p(x)}{q(x)} \right]$$

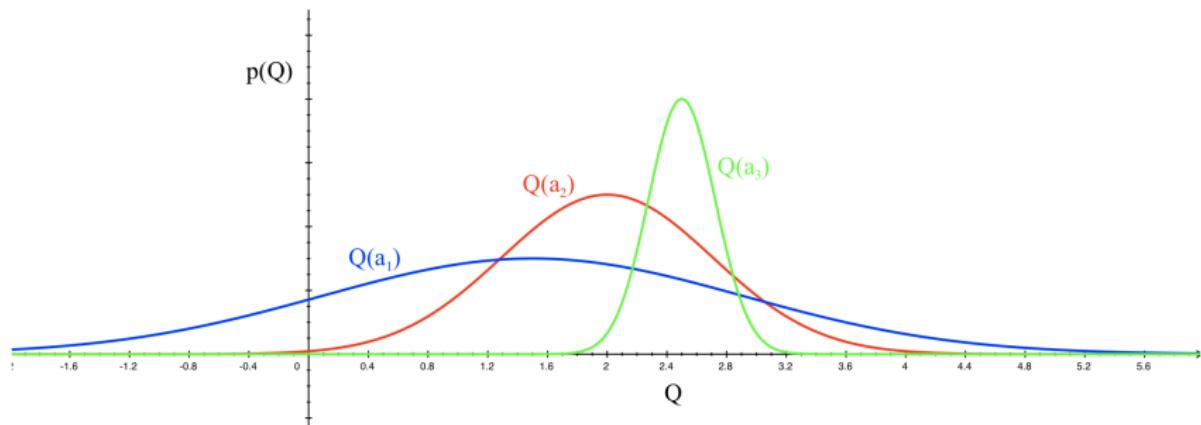
- Gibbs' Inequality: $KL(\mathbf{p}, \mathbf{q}) \geq 0$ for any two distributions \mathbf{p}, \mathbf{q} , with equality iff $\mathbf{p} = \mathbf{q}$.
- Pinsker's inequality: for any event $A \subset \Omega$, we have

$$2(p(A) - q(A))^2 \leq KL(\mathbf{p}, \mathbf{q}).$$

Outline

- 1 Introduction
- 2 Overview of Bandit Problems
- 3 Stochastic Bandits
- 4 Greedy Algorithms
- 5 UCB Algorithms
- 6 Thompson Sampling Algorithms
- 7 Gradient Bandit Algorithms
- 8 Adversarial Bandits
- 9 Summary
- 10 References

Optimism in the Face of Uncertainty



- Which action should we pick?
- The more uncertain we are about an action-value
- The more important it is to explore that action
- It could turn out to be the best action
- One should act as if the environment is as nice as plausibly possible

Upper Confidence Bounds

- Estimate an upper confidence bound $\hat{U}_t(a)$ for each action value
- Such that $Q(a) \leq \hat{Q}_t(a) + \hat{U}_t(a)$ with high probability
- This depends on the number of times $N_t(a)$ has been selected
 - ▶ Small $N_t(a) \Rightarrow$ large $\hat{U}_t(a)$ (estimated value is uncertain)
 - ▶ Large $N_t(a) \Rightarrow$ small $\hat{U}_t(a)$ (estimated value is accurate)
- Select action maximizing Upper Confidence Bound (UCB)

$$a_t = \arg \max_{a \in \mathcal{A}} \{\hat{Q}_t(a) + \hat{U}_t(a)\}.$$

Hoeffding's Inequality

Theorem (Hoeffding's Inequality)

Let X_1, \dots, X_t be i.i.d. random variables in $[0, 1]$, and let $\bar{X}_t = \frac{1}{t} \sum_{\tau=1}^t X_\tau$ be the sample mean. Then

$$\mathbb{P}[\mathbb{E}[X] > \bar{X}_t + u] \leq e^{-2tu^2}$$

W.L.O.G. renormalization of all rewards
to be within [0, 1].

- We will apply Hoeffding's Inequality to rewards of the bandit to be within [0, 1].
- Conditioned on selecting action a

$$\mathbb{P}[Q(a) > \hat{Q}_t(a) + U_t(a)] \leq e^{-2N_t(a)U_t^2(a)}$$

For general case $e^{\frac{-2N_t(a)U_t^2(a)}{c^2}}$.

Calculating Upper Confidence Bounds

- Pick a probability p that true value exceeds UCB
- Now solve for $U_t(a)$

$$e^{-2N_t(a)U_t^2(a)} = p$$

$$U_t(a) = \sqrt{\frac{-\log p}{2N_t(a)}}$$

- Reduce p as we observe more rewards, e.g. $\underline{p = t^{-4}}$
- Ensures we select optimal action as $t \rightarrow \infty$ $P(Q_a \leq Q_{t(a)} + U_{t(a)})$

$$U_t(a) = \sqrt{\frac{2 \log t}{N_t(a)}} \quad > 1-p = 1 - \frac{1}{t^4}.$$

as $t \rightarrow \infty$, $1 - \frac{1}{t^4} \rightarrow 1$

UCB1

- This leads to the UCB1 algorithm: try each action (arm) once, then select actions according to the following rule

$$a_t = \arg \max_{a \in \mathcal{A}} \left\{ \hat{Q}_t(a) + \sqrt{\frac{2 \log t}{N_t(a)}} \right\}$$

- An action a is chosen at time t for two possible reasons
- Exploitation: $\hat{Q}_t(a)$ is large, indicating a high reward w.h.p
- Exploration: $N_t(a)$ is small, indicating that this action has not been explored much. *may have a high reward (Optimism in the face of uncertainty)*
- Summing two parts up is a natural way to trade off them.

Performance of UCB1

Theorem

The UCB1 algorithm achieves logarithmic asymptotic total regret

$$\lim_{t \rightarrow \infty} L_t \leq 8 \log t \sum_{a | \Delta_a > 0} \Delta_a$$

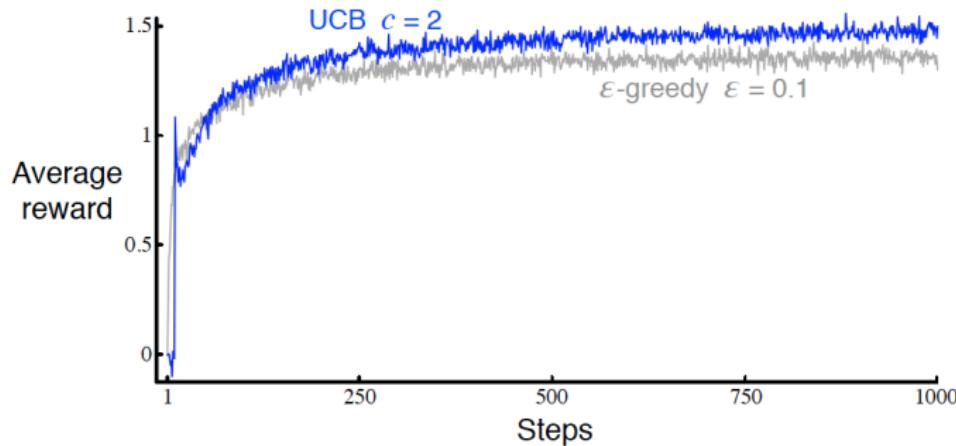
How to prove it?

Example: UCB vs. ϵ -Greedy On 10-armed Bandit

- Average performance of UCB action selection on the 10-armed testbed.
- UCB generally performs better than ϵ -Greedy
- Except in the first k steps, when it selects randomly among the as-yet-untried actions.

$$a_t = \underset{a \in A}{\operatorname{argmax}} \left[\hat{\mu}_{t(a)} + C \sqrt{\frac{\log t}{N_{t(a)}}} \right]$$

UCB-C algorithm



Example: UCB vs. ϵ -Greedy On 10-armed Bandit

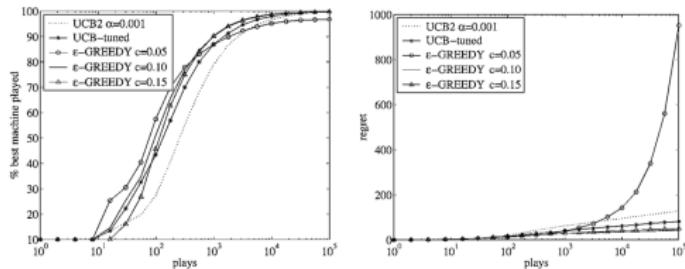


Figure 9. Comparison on distribution 11 (10 machines with parameters 0.9, 0.6, ..., 0.6).

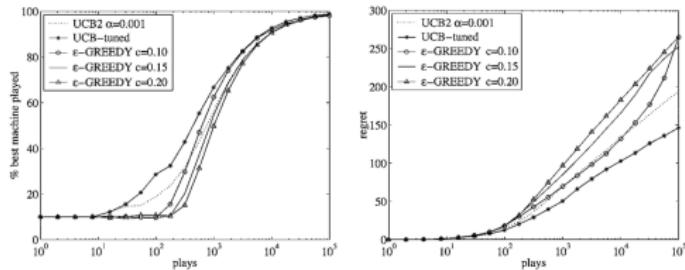


Figure 10. Comparison on distribution 12 (10 machines with parameters 0.9, 0.8, 0.8, 0.8, 0.7, 0.7, 0.7, 0.6, 0.6, 0.6).

- Peter Auer, Nicolò Cesa-Bianchi, and Paul Fischer, “Finite-time analysis of the multi-armed bandit problem”, Machine Learning, vol.47, no.2-3, pp. 235-256, 2002.

Outline

- 1 Introduction
- 2 Overview of Bandit Problems
- 3 Stochastic Bandits
- 4 Greedy Algorithms
- 5 UCB Algorithms
- 6 Thompson Sampling Algorithms
- 7 Gradient Bandit Algorithms
- 8 Adversarial Bandits
- 9 Summary
- 10 References

Bayesian Bandits

- So far we have made no assumptions about the reward distribution \mathcal{R}
 - ▶ Except bounds on rewards
- Bayesian bandits exploit prior knowledge of rewards, $p[\mathcal{R}]$
- They compute posterior distribution of rewards $p[\mathcal{R} \mid h_t]$
 - ▶ where $h_t = a_1, r_1, \dots, a_{t-1}, r_{t-1}$ is the history
- Use posterior to guide exploration
 - ▶ Upper confidence bounds (Bayesian UCB)
 - ▶ Probability matching (Thompson sampling)

Think about it.

Probability Matching

- Probability matching selects action a according to probability that a is the optimal action

Similar to MAP rule.

$$\pi(a \mid h_t) = \mathbb{P}[Q(a) > Q(a'), \forall a' \neq a \mid h_t]$$

- Probability matching is optimistic in the face of uncertainty
 - ▶ Uncertain actions have higher probability of being max

Thompson Sampling

- Thompson sampling (also called “posterior sampling”) implements “probability matching”

$$\begin{aligned}\pi(a \mid h_t) &= \mathbb{P}[Q(a) > Q(a'), \forall a' \neq a \mid h_t] \\ &= \mathbb{E}_{\mathcal{R} \mid h_t} \left[\mathbf{1}(a = \arg \max_{a \in \mathcal{A}} Q(a)) \right]\end{aligned}$$

- Use Bayes’ law to compute posterior distribution $p[\mathcal{R} \mid h_t]$
- Sample a reward distribution \mathcal{R} from posterior
- Compute action-value function $Q(a) = \mathbb{E}_{\substack{\mathcal{R}_a \\ \text{Sample function}}}^f[\mathcal{R}_a]$ *f: some functions of R_a .*
- Select action maximizing value on sample, $a_t = \arg \max_{a \in \mathcal{A}} Q(a)$

Thompson Sampling

- The first bandit algorithm proposed by W. R. Thompson in 1933
- For the Bernoulli case with two arms (no theory)
- Almost ignored for nearly 80 years!
- Rediscovered from 2010 & 2011 due to its superior empirical performance over UCB algorithms & others
- Theoretical analysis in 2012 showed that Thompson sampling achieves Lai-Robbins lower bound! (asymptotically optimal)
- Can be difficult to compute analytically from posterior unless the cases of conjugate prior

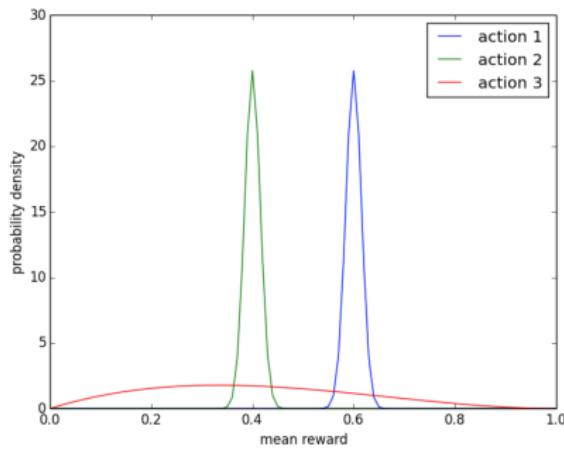
Beta-Bernoulli Bandit

- Stochastic bandits with K actions (arms)
- Action $k \in \{1, \dots, K\}$ being played produces a reward satisfying Bernoulli distribution $\text{Bern}(\theta_k)$: reward 1 w.p. θ_k and reward 0 w.p. $1 - \theta_k$.
- θ_k : an action's success probability or mean reward
- The mean rewards $\theta = (\theta_1, \dots, \theta_K)$ are unknown constants.
- Prior of each θ_k satisfies Beta distribution $\text{Beta}(\alpha_k, \beta_k)$
- x_t denotes the action selected at time t and r_t denotes the corresponding reward of action x_t
- Each action's posterior distribution is also Beta with parameters updated as follows:

$$(\alpha_k, \beta_k) \leftarrow \begin{cases} (\alpha_k, \beta_k) & \text{if } x_t \neq k \\ (\alpha_k, \beta_k) + (r_t, 1 - r_t) & \text{if } x_t = k. \end{cases}$$

Beta Distribution

- When $\alpha_k = \beta_k = 1$, the Beta distribution is $\text{Unif}(0,1)$
- (α_k, β_k) are also called “pseudo counts”
- Beta distribution $\text{Beta}(\alpha_k, \beta_k)$ has mean $\frac{\alpha_k}{\alpha_k + \beta_k}$
- Beta distribution becomes more concentrated as $\alpha_k + \beta_k$ grows
- Example: $\text{Beta}(601, 401)$, $\text{Beta}(401, 601)$ and $\text{Beta}(2, 3)$.



Greedy vs. Thompson Sampling

Algorithm 1 BernGreedy(K, α, β)

```
1: for  $t = 1, 2, \dots$  do
2:   #estimate model:
3:   for  $k = 1, \dots, K$  do
4:      $\hat{\theta}_k \leftarrow \alpha_k / (\alpha_k + \beta_k)$  E[ $\text{Beta}(\alpha_k, \beta_k)$ ]
5:   end for
6:
7:   #select and apply action:
8:    $x_t \leftarrow \text{argmax}_k \hat{\theta}_k$ 
9:   Apply  $x_t$  and observe  $r_t$ 
10:
11:  #update distribution:
12:   $(\alpha_{x_t}, \beta_{x_t}) \leftarrow (\alpha_{x_t} + r_t, \beta_{x_t} + 1 - r_t)$ 
13: end for
```

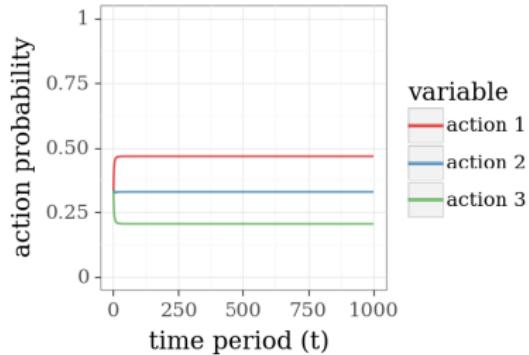
Algorithm 2 BernTS(K, α, β)

```
1: for  $t = 1, 2, \dots$  do
2:   #sample model:
3:   for  $k = 1, \dots, K$  do
4:     Sample  $\hat{\theta}_k \sim \text{beta}(\alpha_k, \beta_k)$ 
5:   end for
6:
7:   #select and apply action:
8:    $x_t \leftarrow \text{argmax}_k \hat{\theta}_k$ 
9:   Apply  $x_t$  and observe  $r_t$ 
10:
11:  #update distribution:
12:   $(\alpha_{x_t}, \beta_{x_t}) \leftarrow (\alpha_{x_t} + r_t, \beta_{x_t} + 1 - r_t)$ 
13: end for
```

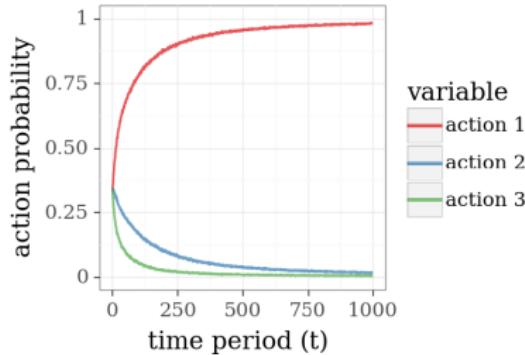
Simulation

- Three-armed beta-Bernoulli bandit with mean rewards $\theta_1 = 0.9$, $\theta_2 = 0.8$ and $\theta_3 = 0.7$.
- Prior distributions over each mean reward is uniform.

action 1 is the
optimal action



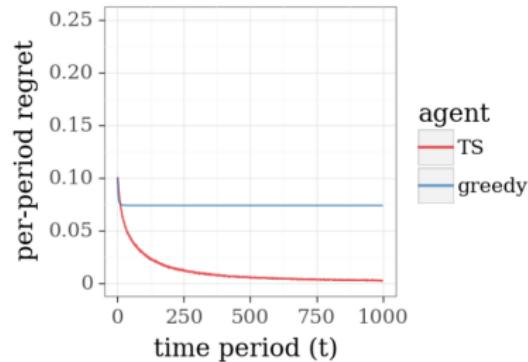
(a) greedy algorithm



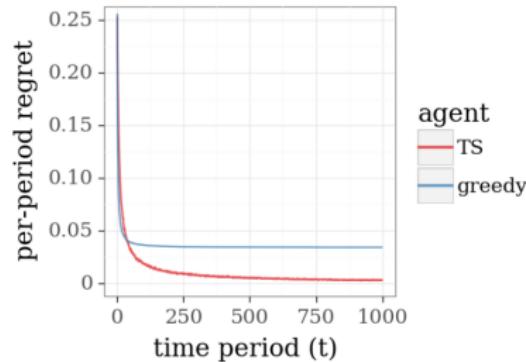
(b) Thompson sampling

Simulation

- Per-period regret: $\text{regret}_t(\theta) = \max_k \theta_k - \theta_{x_t}$



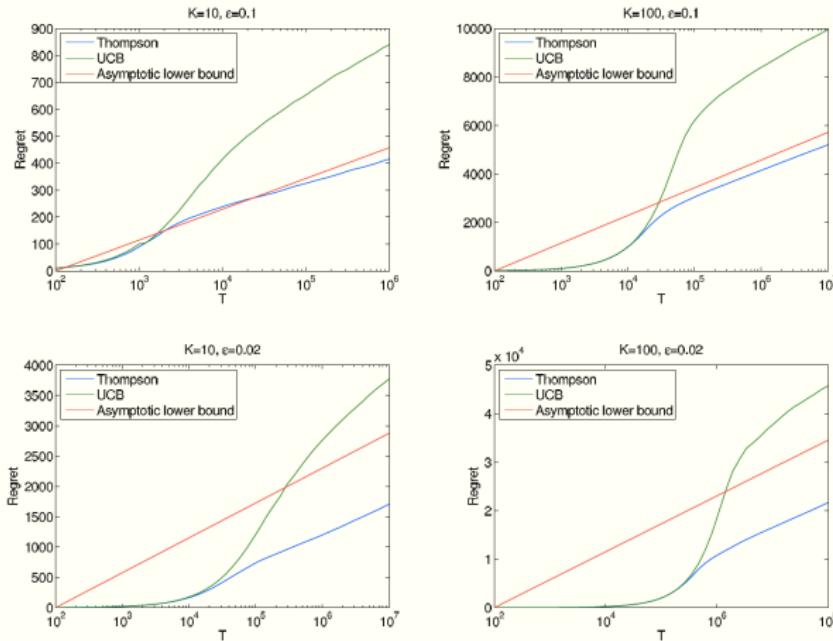
(a) $\theta = (0.9, 0.8, 0.7)$



(b) average over random θ

UCB vs. Thompson Sampling

empirically better



- O. Chapelle & L. Li, "An Empirical Evaluation of Thompson Sampling", NeuralIPS 2011.

Summary of Conjugate Priors: Discrete Case

How to reduce the computational complexity of posterior sampling?

Likelihood	Model parameters	Conjugate prior distribution	Prior hyperparameters	Posterior hyperparameters	Interpretation of hyperparameters [note 1]	Posterior predictive [note 2]
Bernoulli	p (probability)	Beta	α, β	$\alpha + \sum_{i=1}^n x_i, \beta + n - \sum_{i=1}^n x_i$	$\alpha - 1$ successes, $\beta - 1$ failures [note 1]	$p(\tilde{x} = 1) = \frac{\alpha'}{\alpha' + \beta'}$
Binomial	p (probability)	Beta	α, β	$\alpha + \sum_{i=1}^n x_i, \beta + \sum_{i=1}^n N_i - \sum_{i=1}^n x_i$	$\alpha - 1$ successes, $\beta - 1$ failures [note 1]	BetaBin($\tilde{x} \alpha', \beta'$) (Beta-binomial)
Negative binomial with known failure number, r	p (probability)	Beta	α, β	$\alpha + \sum_{i=1}^n x_i, \beta + rn$	$\alpha - 1$ total successes, $\beta - 1$ failures [note 1] (i.e., $\frac{\beta - 1}{r}$ experiments, assuming r stays fixed)	BetaNegBin($\tilde{x} \alpha', \beta'$) (Beta-negative binomial)
Poisson	λ (rate)	Gamma	k, θ	$k + \sum_{i=1}^n x_i, \frac{\theta}{\theta + 1}$	λ total occurrences in $\frac{1}{\theta}$ intervals	NB($\tilde{x} \lambda', \theta'$) (negative binomial)
			α, β [note 3]	$\alpha + \sum_{i=1}^n x_i, \beta + n$	λ total occurrences in β intervals	NB($\tilde{x} \alpha', \frac{1}{1+\beta}$) (negative binomial)
Categorical	μ (probability vector), k (number of categories; i.e., size of μ)	Dirichlet	α	$\alpha + (c_1, \dots, c_k)$, where c_i is the number of observations in category i	$\alpha_i - 1$ occurrences of category i [note 1]	$p(\tilde{x} = i) = \frac{\alpha'_i}{\sum_i \alpha'_i}$ $= \frac{\alpha_i + c_i}{\sum_i \alpha_i + n}$
Multinomial	μ (probability vector), k (number of categories; i.e., size of μ)	Dirichlet	α	$\alpha + \sum_{i=1}^n x_i$	$\alpha_i - 1$ occurrences of category i [note 1]	DirMult($\tilde{x} \alpha'$) (Dirichlet-multinomial)
Hypergeometric with known total population size, N	M (number of target members)	Beta-Ishomial ^[4]	$n = N, \alpha, \beta$	$\alpha + \sum_{i=1}^n x_i, \beta + \sum_{i=1}^n N_i - \sum_{i=1}^n x_i$	$\alpha - 1$ successes, $\beta - 1$ failures [note 1]	
Geometric	ρ_0 (probability)	Beta	α, β	$\alpha + n, \beta + \sum_{i=1}^n x_i$	$\alpha - 1$ experiments, $\beta - 1$ total failures [note 1]	

- https://en.wikipedia.org/wiki/Conjugate_prior
- A Compendium of Conjugate Priors, Daniel Fink, 1997.

Summary of Conjugate Priors: Continuous

Likelihood	Model parameters	Conjugate prior distribution	Prior hyperparameters	Posterior hyperparameters	Interpretation of hyperparameters	Posterior predictive [14]
Normal with known variance σ^2	μ (mean)	Normal	μ_0, σ_0^2	$\frac{1}{\sigma_0^2} + \frac{n}{\sigma^2}, \left(\frac{\mu_0}{\sigma_0^2} + \frac{\sum_{i=1}^n x_i}{\sigma^2} \right)^{-1}$	mean was estimated from observations with total precision (sum of all individual precisions) $1/\sigma_0^2$ and with sample mean \bar{x}_0	$N(\bar{x}_0 \mu_0^*, \sigma_0^2 + \sigma^2)$
Normal with known precision τ	μ (mean)	Normal	μ_0, η_0	$\frac{\eta_0}{\eta_0 + n\tau}, \frac{\eta_0 + \sum_{i=1}^n x_i}{\eta_0 + n\tau}$	mean was estimated from observations with total precision (sum of all individual precisions) η_0 and with sample mean \bar{x}_0	$N\left(\bar{x}_0 \mu_0^*, \frac{1}{\eta_0^* + \frac{1}{\tau}}\right)$
Normal with known mean μ	σ^2 (variance)	Inverse gamma	α_0, β_0 (rate β_0)	$\alpha_0 + n, \beta_0 + \sum_{i=1}^n (x_i - \mu)^2$	variance was estimated from $2n$ observations with sample variance δ^2 (i.e. sum of squared deviations $2S$, where deviations are from known mean μ)	$\text{Inv-G}(2\bar{x}_0^2 + \delta^2 \alpha_0^*)$
Normal with known mean μ	σ^2 (variance)	Student's t-distribution	ν, σ_0^2	$\nu + n, \frac{\nu + n}{\nu + n - 2} (\bar{x}_0 - \mu)^2$	variance was estimated from ν observations with sample variance δ^2	$t_{\nu}(\bar{x}_0 \mu, \sigma_0^2)$
Normal with known mean μ	τ (precision)	Gamma	α_0, β_0 (rate β_0)	$\alpha_0 + n, \beta_0 + \sum_{i=1}^n (x_i - \mu)^2$	precision was estimated from $2n$ observations with sample variance δ^2 (i.e. sum of squared deviations $2S$, where deviations are from known mean μ)	$\text{Inv-G}(\delta^2 \alpha_0^*, \beta_0^*)$
Normal (prior E)	μ and σ^2 Assuming exchangeability	Normal-inverse gamma	$\mu_0, \nu_0, \eta_0, \beta$	$\mu_0 + \frac{\nu_0}{\nu_0 + n} \bar{x}_0, \frac{\nu_0 + n}{\nu_0 + n - 2} \delta^2 + \frac{\nu_0}{\nu_0 + n} (\bar{x}_0 - \mu_0)^2$	mean was estimated from n observations with sample mean \bar{x}_0 ; variance was estimated from $2n$ observations with sample mean μ_0 and sum of squared deviations $2S$	$t_{\nu_0 + n}(\bar{x}_0 \mu_0^*, \frac{\delta^2(\nu_0 + 1)}{\nu_0^*})$
Normal	μ and τ Assuming exchangeability	Normal-inverse gamma	$\mu_0, \nu_0, \eta_0, \beta$	$\mu_0 + \frac{\nu_0}{\nu_0 + n} \bar{x}_0, \frac{\nu_0 + n}{\nu_0 + n - 2} \delta^2$	mean was estimated from n observations with sample mean \bar{x}_0 ; precision was estimated from $2n$ observations with sample mean μ_0 and sum of squared deviations $2S$	$t_{\nu_0 + n}(\bar{x}_0 \mu_0^*, \frac{\delta^2(\nu_0 + 1)}{\nu_0^*})$
Multivariate normal with known covariance matrix Σ	μ (mean vector)	Multivariate normal	μ_0, Σ_0	$\mu_0 + \frac{1}{n} \sum_{i=1}^n (x_i - \mu_0)^T \Sigma_0^{-1} + \frac{\Sigma_0}{n + 1} (\bar{x}_0 - \mu_0)^T \Sigma_0^{-1} \bar{x}_0$	δ is the sample mean	$A[\lambda \mu_0, \Sigma_0^{-1} + 2\Sigma]$
Multivariate normal with known precision matrix A	μ (mean vector)	Multivariate normal	μ_0, A_0	$(A_0 + nA)^{-1} (\bar{x}_0 + nA\bar{x}_0) = (A_0 + nA)^{-1} \bar{x}_0$	δ is the sample mean	$A[\lambda \mu_0, (A_0^{-1} + A)^{-1}]$
Multivariate normal with known mean μ	Σ (covariance matrix)	Inverse-Wishart	ν, Ψ	$\nu + n, \bar{x}_0 + \frac{1}{\nu + n} \sum_{i=1}^n (x_i - \mu)(x_i - \mu)^T$	covariance matrix was estimated from n observations with sum of pairwise deviation products \bar{V}^{-1}	$t_{\nu + n - 1}(\bar{V} \nu + p - 1, \Psi)$
Multivariate normal with known mean μ	A (precision matrix)	Wishart	ν, V	$\nu + n, V + \frac{1}{\nu + n} \sum_{i=1}^n (x_i - \mu)(x_i - \mu)^T$	covariance matrix was estimated from n observations with sum of pairwise deviation products V^{-1}	$t_{\nu + n - 1}(\bar{V} \nu + p - 1, V^{-1})$
Multivariate normal	μ (mean vector) and Σ (covariance matrix)	Normal-inverse-Wishart	$\mu_0, \nu_0, \eta_0, \Psi$	$\mu_0 + \frac{\nu_0}{\nu_0 + n} \bar{x}_0, \frac{\nu_0 + n}{\nu_0 + n - 2} \delta^2 + \frac{\nu_0}{\nu_0 + n} (\bar{x}_0 - \mu_0)^2$	mean was estimated from n observations with sample mean \bar{x}_0 ; covariance matrix was estimated from n observations with sample mean \bar{x}_0 and with sum of pairwise deviation products Ψ	$t_{\nu_0 + n - 1}(\bar{x}_0 \mu_0^*, \frac{\nu_0^* + 1}{\nu_0^*(\nu_0^* + p - 1)} \Psi)$
Multivariate normal	μ (mean vector) and A (precision matrix)	Normal-Wishart	μ_0, ν_0, η_0, V	$\mu_0 + \frac{\nu_0}{\nu_0 + n} \bar{x}_0, \frac{\nu_0 + n}{\nu_0 + n - 2} \delta^2 + \frac{\nu_0}{\nu_0 + n} (\bar{x}_0 - \mu_0)^2$	mean was estimated from n observations with sample mean \bar{x}_0 ; covariance matrix was estimated from n observations with sample mean \bar{x}_0 and with sum of pairwise deviation products V^{-1}	$t_{\nu_0 + n - 1}(\bar{x}_0 \mu_0^*, \frac{\nu_0^* + 1}{\nu_0^*(\nu_0^* + p - 1)} V^{-1})$
Uniform	$E(\theta, R)$	Point	x_{\min}, R	$\max(x_1, \dots, x_n), R$	θ observations with minimum value x_{\min}	
Parzen with known minimum x_{\min}	A (shape)	Gaussian	μ_0, β	$\mu_0 + \nu_0, \sum_{i=1}^n \frac{x_i - x_{\min}}{\beta}$	ν observations with sum β of the \log of each observation (i.e. the logarithm of the ratio of each observation to the minimum x_{\min})	
Weibull with known shape β	θ (scale)	Inverse gamma [14]	α_0, β_0	$\alpha_0 + n, \sum_{i=1}^n \frac{1}{x_i}$	ν observations with sum β of the β th power of each observation	
Log-normal	(Same as for the normal distribution after exponentiating the data)	Inverse gamma [14]	α_0, β_0			
Exponential	λ (rate)	Gamma	α_0, β_0 (rate β)	$\alpha_0 + n, \beta + \sum_{i=1}^n x_i$	ν observations that sum to β [14]	
Gamma with known shape α	λ (rate)	Gamma	α_0, β_0	$\alpha_0 + n\lambda, \beta_0 + \sum_{i=1}^n x_i$	ν observations with sum β_0	
Inverse Gamma with known shape α	β (inverse scale)	Gamma	α_0, β_0	$\alpha_0 + n\beta_0, \beta_0 + \sum_{i=1}^n \frac{1}{x_i}$	ν observations with sum β_0	
Gamma	α (shape)	$\sigma^{-1} \cdot \text{Exp}$	ν, k, c	$\nu \prod_{i=1}^k x_i, \delta + \alpha_k + \nu$	δ or ν observations (for estimating α_k or estimating δ with product ν)	
Gamma with known rate β						

- https://en.wikipedia.org/wiki/Conjugate_prior
- *A Compendium of Conjugate Priors, Daniel Fink, 1997.*

Approximate Posterior Sampling

- In many cases exact Bayesian inference is computational intractable
- Four approaches to approximate posterior sampling
 - ▶ Gibbs sampling
 - ▶ Langevin Monte Carlo
 - ▶ Sampling from a Laplace approximation
 - ▶ Bootstrap

Outline

- 1 Introduction
- 2 Overview of Bandit Problems
- 3 Stochastic Bandits
- 4 Greedy Algorithms
- 5 UCB Algorithms
- 6 Thompson Sampling Algorithms
- 7 Gradient Bandit Algorithms
- 8 Adversarial Bandits
- 9 Summary
- 10 References

Gradient Bandit Algorithm

Previously our bandit algorithms are based on estimate of action value.

- Let $H_t(a)$ be a learned preference for taking action a (to select action)

$$\Pr\{A_t = a\} \doteq \frac{e^{H_t(a)}}{\sum_{b=1}^k e^{H_t(b)}} \doteq \pi_t(a) \propto e^{H_t(a)}$$

Stochastic Gradient Ascent

$$H_{t+1}(a) \doteq H_t(a) + \alpha \left(R_t - \bar{R}_t \right) \left(\mathbb{1}\{A_t = a\} - \pi_t(a) \right), \quad \forall a$$

$$\bar{R}_t \doteq \frac{1}{t} \sum_{i=1}^t R_i$$

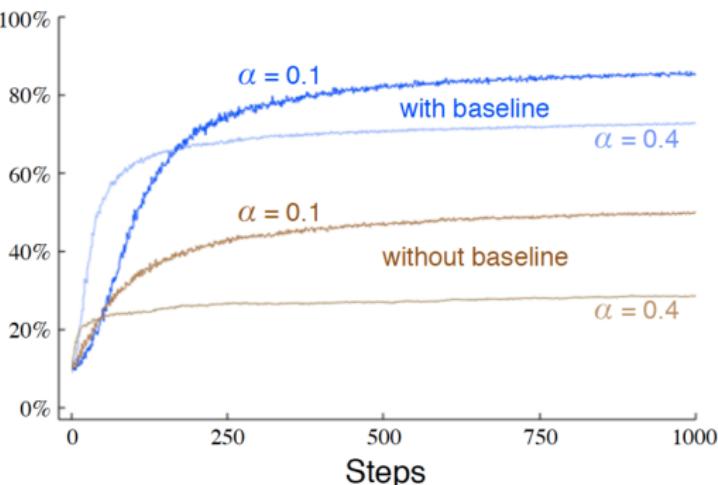
Baseline.

When Baseline = 0;

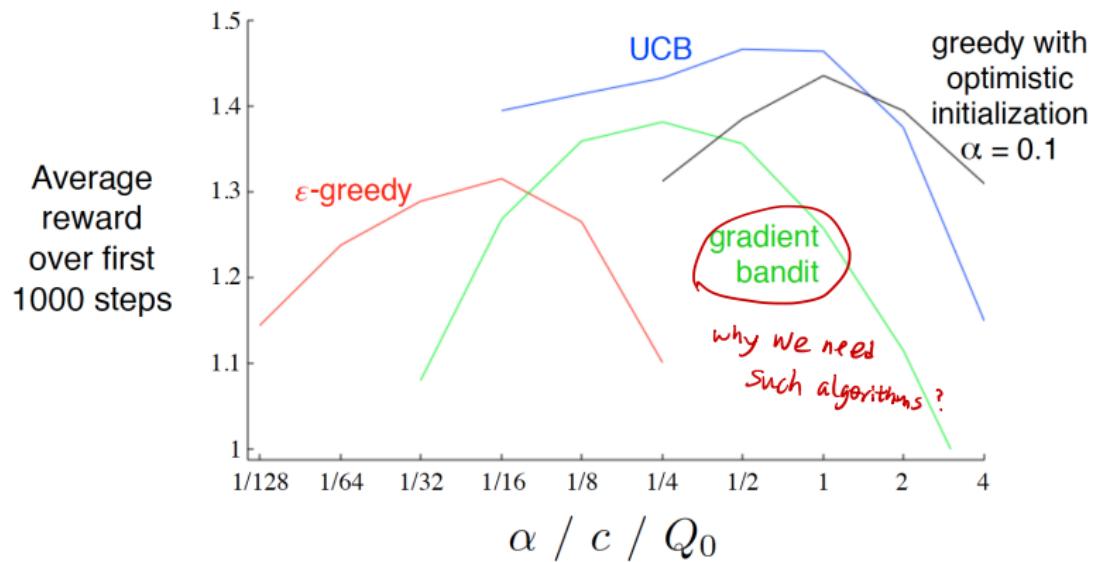
$$H_{t+1}(a) = H_t(a) + \alpha R_t (\mathbb{1}_{A_t=a} - \pi_t(a))$$

$$\alpha = \Delta t, \quad H_{t+1}(a) > H_t(a) \quad \text{prob. } T$$

$$\alpha \neq \Delta t, \quad H_{t+1}(a) < H_t(a). \quad \text{prob. } \downarrow$$



Performance Comparison with Other Bandit Algorithms



Why We Study Gradient Bandit Algorithm

- A special case of the gradient-based RL algorithms
- Precursor of policy-gradient & actor-critic algorithms in RL

How to improve gradient bandit algorithm?

$$\pi_t(a) = \frac{e^{H_t(a)}}{\sum_{b=1}^K e^{H_t(b)}}$$

parameterized version

β or $\beta(t)$ new degrees of freedom

$$\pi_t^\beta(a) = \frac{e^{\beta H_t(a)}}{\sum_{b=1}^K e^{\beta H_t(b)}}$$

β is a positive constant

$$\pi_t^{\beta(t)}(a) = \frac{e^{\beta(t) H_t(a)}}{\sum_{b=1}^K e^{\beta(t) H_t(b)}}.$$

rules to adjust $\beta(t)$.

Derivation of Gradient Bandit Algorithm

$$1^{\text{st}}. \max_w E[f(w)]$$

gradient ascent algorithm
 $w \leftarrow w + \alpha \nabla_w E[f(w)]$

Stochastic Ascent Algorithm

$$w \leftarrow w + \alpha \nabla_w f(w) \quad [w(t+1) = w(t) + \alpha \cdot \nabla_w f(w_t)]$$

2nd. Now we return to gradient bandit algorithm

$\max E[R_t]$: objective function

Multi-armed Bandit (MAB)
action arm
 $x \in \{1, \dots, k\}$

$$E[R_t] = E[E[R_t | A_t]] = \sum_x E[R_t | A_t=x] \cdot P(A_t=x)$$

$$= \sum_x q_*(x) P(A_t=x) = \sum_x q_*(x) \pi_t(x)$$

where $\pi_t(x) = P(A_t=x) = \frac{e^{H_t(x)}}{\sum_{y=1}^k e^{H_t(y)}}$

reward info $q_*(x)$ is independent of preference $H_t(x)$.

$\Rightarrow E[R_t]$ is a function of $H_t(x)$

$$\Rightarrow H_{t+1}(a) = H_t(a) + \alpha \cdot \frac{\partial E[R_t]}{\partial H_t(a)}, \quad \forall a \in \{1, \dots, k\}$$

Derivation of Gradient Bandit Algorithm

3^o. we first show a result

Lemma 1 : $\frac{\partial \pi_t(x)}{\partial H_t(a)} = \pi_t(x) (\sum_{x \neq a} \pi_t(x) - \pi_t(a))$

Proof : $\frac{\partial}{\partial x} \left[\frac{f(x)}{g(x)} \right] = \frac{\frac{\partial f(x)}{\partial x} \cdot g(x) - f(x) \frac{\partial g(x)}{\partial x}}{(g(x))^2}$

$$\Rightarrow \frac{\partial \pi_t(x)}{\partial H_t(a)} = \frac{\partial}{\partial H_t(a)} \left[\frac{e^{H_t(x)}}{\sum_{y=1}^k e^{H_t(y)}} \right] = \frac{\frac{\partial e^{H_t(x)}}{\partial H_t(a)} \cdot \sum_{y=1}^k e^{H_t(y)} - e^{H_t(x)} \cdot \frac{\partial (\sum_{y=1}^k e^{H_t(y)})}{\partial H_t(a)}}{\left(\sum_{y=1}^k e^{H_t(y)} \right)^2}$$

$$= \frac{\sum_{x \neq a} e^{H_t(x)} \cdot \sum_{y=1}^k e^{H_t(y)} - e^{H_t(x)} \cdot e^{H_t(a)}}{\left(\sum_{y=1}^k e^{H_t(y)} \right)^2} = \frac{\sum_{x \neq a} e^{H_t(x)}}{\sum_{y=1}^k e^{H_t(y)}} - \frac{e^{H_t(x)}}{\sum_{y=1}^k e^{H_t(y)}} \cdot \frac{e^{H_t(a)}}{\sum_{y=1}^k e^{H_t(y)}}$$

$$= \sum_{x \neq a} \pi_t(x) - \pi_t(x) \cdot \pi_t(a) = \pi_t(x) (\sum_{x \neq a} \pi_t(x) - \pi_t(a)) \quad Q.E.D.$$

Derivation of Gradient Bandit Algorithm

$$4^{\circ} \text{ thus by Lemma 1, we know } \frac{\partial \pi_t(x)}{\partial h_t(a)} = \pi_t(x) (\mathbb{1}_{\{x=a\}} - \pi_t(a))$$

$$\Rightarrow \frac{\partial \pi_t(A_t)}{\partial h_t(a)} = \pi_t(A_t) (\mathbb{1}_{\{A_t=a\}} - \pi_t(a)) \quad (*_1) \quad (A_t: \text{action in time } t)$$

$$5^{\circ} \cdot E[R_t] = \sum_x q_{*}(x) \pi_t(x) \Rightarrow \frac{\partial E[R_t]}{\partial h_t(a)} = \sum_x q_{*}(x) \frac{\partial \pi_t(x)}{\partial h_t(a)} \quad (*_2)$$

$$\text{Since } \sum_x \pi_t(x) = 1 \Rightarrow \frac{\partial (\sum_x \pi_t(x))}{\partial h_t(a)} = 0 \Rightarrow \sum_x \frac{\partial \pi_t(x)}{\partial h_t(a)} = 0$$

thus we introduce a baseline B_t (a scalar that does not depend on x)

$$\Rightarrow \sum_x B_t \cdot \frac{\partial \pi_t(x)}{\partial h_t(a)} = 0 \quad \stackrel{(*_2)}{\Rightarrow} \frac{\partial E[R_t]}{\partial h_t(a)} = \sum_x [q_{*}(x) - B_t] \cdot \frac{\partial \pi_t(x)}{\partial h_t(a)}$$
$$= \sum_x \pi_t(x) [q_{*}(x) - B_t] \left(\frac{\partial \pi_t(x)}{\partial h_t(a)} \right) \frac{1}{\pi_t(x)}$$

$$\text{Lemma!} \quad = \sum_x \pi_t(x) [q_{*}(x) - B_t] (\mathbb{1}_{\{x=a\}} - \pi_t(a))$$

$$= E[(q_{*}(A_t) - B_t) (\mathbb{1}_{\{A_t=a\}} - \pi_t(a))]$$

Derivation of Gradient Bandit Algorithm

6^o. Now we choose $B_t = \widehat{R}_t = \frac{1}{t} \sum_{i=1}^t R_i$

$$\Rightarrow \frac{\partial E[R_t]}{\partial \pi_t(a)} = E[(q_\pi(A_t) - \widehat{R}_t)(\mathbb{I}_{\{A_t=a\}} - \pi_t(a))]$$

7^o. Next we will show $E[(q_\pi(A_t) - \widehat{R}_t)(\mathbb{I}_{\{A_t=a\}} - \pi_t(a))]$

$$= E[(R_t - \widehat{R}_t)(\mathbb{I}_{\{A_t=a\}} - \pi_t(a))]$$

$$\Leftrightarrow E[q_\pi(A_t)(\mathbb{I}_{\{A_t=a\}} - \pi_t(a))] = E[R_t(\mathbb{I}_{\{A_t=a\}} - \pi_t(a))]$$

Proof :

$$\begin{aligned} E[q_\pi(A_t)(\mathbb{I}_{\{A_t=a\}} - \pi_t(a))] &= E[E[R_t | A_t](\mathbb{I}_{\{A_t=a\}} - \pi_t(a))] \\ &= E[E[R_t(\mathbb{I}_{\{A_t=a\}} - \pi_t(a))] | A_t] \end{aligned}$$

function of A_t .

Adam's Law

$$= E[R_t(\mathbb{I}_{\{A_t=a\}} - \pi_t(a))]$$

Derivation of Gradient Bandit Algorithm

8°. thus we have

$$\begin{aligned}\frac{\partial E[R_t]}{\partial \pi_t(a)} &= E[(\hat{R}_t - \bar{R}_t) (\mathbb{1}_{A_t=a} - \pi_t(a))] \\ &= E[(R_t - \bar{R}_t) (\mathbb{1}_{A_t=a} - \pi_t(a))]\end{aligned}$$

$$\Rightarrow \text{Gradient Ascent Algorithm: } H_{t+1}(a) = H_t(a) + \alpha \cdot \frac{\partial E[R_t]}{\partial \pi_t(a)}$$
$$= H_t(a) + \alpha \cdot E[(R_t - \bar{R}_t) (\mathbb{1}_{A_t=a} - \pi_t(a))]$$

\Rightarrow Stochastic Gradient Ascent Algorithm (Gradient Bandit Algorithm)

$$H_{t+1}(a) = H_t(a) + \alpha (R_t - \bar{R}_t) (\mathbb{1}_{A_t=a} - \pi_t(a)) \quad \Theta a \in \{1, \dots, k\}$$

Derivation of Gradient Bandit Algorithm

Q⁰. Remark on policy Gradient.

$$\max_{\theta} E_X[f(x)] \quad , \quad X \text{ is a r.v. } \sim p_{\theta}^{\text{MuForPDF}} \cdot (\theta \text{ is a parameter})$$

$$\begin{aligned} D_{\theta} E_X[f(x)] &= D_{\theta} \sum_x p(x) f(x) = \sum_x (D_{\theta} p(x)) f(x) \\ &= \sum_x p(x) \cdot \frac{D_{\theta} p(x)}{p(x)} \cdot f(x) = \sum_x p(x) \cdot (D_{\theta} \log p(x)) \cdot f(x) \\ &= E_X [f(x) D_{\theta} \log p(x)] \end{aligned}$$

Monte Carlo.
 $\approx \frac{1}{n} \sum_{i=1}^n f(x_i) D_{\theta} \log p(x_i)$

$f(x)$ = reward function

D_{θ} : policy.

Gradient of expectation \rightarrow expectation of some gradient function

Derivation of Gradient Bandit Algorithm

Derivation of Gradient Bandit Algorithm

Outline

- 1 Introduction
- 2 Overview of Bandit Problems
- 3 Stochastic Bandits
- 4 Greedy Algorithms
- 5 UCB Algorithms
- 6 Thompson Sampling Algorithms
- 7 Gradient Bandit Algorithms
- 8 Adversarial Bandits
- 9 Summary
- 10 References

Critical Assumption in Stochastic Bandits So Far

- Stochastic reward
- Mean reward of individual arms does not change over time (stationary)
- Independent reward distribution among arms

Adversarial Bandit

- Abandon almost all assumptions on how rewards are generated
- The environment is often called the adversary
- The adversary has a great deal of power in this model
- including the ability to examine the code of the proposed algorithms and choose the rewards accordingly.

Simple Bandit Game

1^o. if your friend chooses $x_1=x_2$, then $R=0$
this is indeed a good friend!

2^o. If you adopt a deterministic strategy, e.g. $A=1$,
then your friend can choose $x_1=0, x_2=1$
and your regret $R=1$. (Your friend
is an adversary)

- The horizon is $n = 1$ and you have two actions.
- You tell your friend your strategy for choosing an action
- Your friend secretly chooses rewards $x_1 \in \{0, 1\}$ and $x_2 \in \{0, 1\}$.
- You implement your strategy to select $A \in \{1, 2\}$ and receives reward x_A .
- The regret is $R = \max\{x_1, x_2\} - x_A$.

$$R \leq 1$$

3^o. To improve your regret against an adversary,
you need to randomize your strategy.
 $PLA=1) = \frac{1}{2}, PLA=2) = \frac{1}{2}$.

then the best your friend can do.
choose $x_1=1, x_2=0$ or reversed.

Expectation of $E[R] = \frac{1}{2}$.
your regret

Adversarial Bandit Setting

Source of randomness: picking actions.

Reward itself is fixed, not stochastic.

- $K > 1$: the number of arms
- A k-armed adversarial bandit is an arbitrary sequence of reward vectors $\nu = (x_t)_{t=1}^n$ where $x_t \in [0, 1]^k$.
- Reward vectors are not random.
- In each round the learner chooses an action $A_t \in [k]$ and observes reward $X_t = x_{tA_t}$.
- The learner will usually randomize their decisions so that A_t and X_t are random variables.
- Performance of π is measured by the expected regret, which is the expected loss in revenue of the policy relative to the best fixed action in hindsight (in contrast to foresight)

$$R_n(\pi, \nu) = \max_{i \in [k]} \sum_{t=1}^n x_{ti} - E\left[\sum_{t=1}^n x_{tA_t}\right]$$

Importance-Weighted Estimators

- A key ingredient of all adversarial bandit algorithms is a mechanism for estimating the reward of unplayed arms
- P_t is the conditional distribution of the action played in round t .
- P_{ti} is the conditional probability H_t : history before t .
 $P_{ti} = P(A_t = i | X_1, \dots, X_{t-1}, A_1, \dots, A_{t-1})$.
- The importance-weighted estimator (unbiased) of x_{ti} is

$$P_{ti} = P(A_t = i | H_t)$$

$$E_t(\cdot) = E_t(\cdot | H_t)$$

$$\Rightarrow E_t[\hat{X}_{ti}] = x_{ti}$$

unbiased. why?

what you mean

$$\hat{X}_{ti} = \frac{I\{A_t = i\} X_t}{P_{ti}}$$

1°.

2°.

Let $A_{ti} = I\{A_t = i\}$, then $E_t[A_{ti}] = P_{ti}$. Since $X_t = X_t A_t$

$$\Rightarrow X_t A_{ti} = X_t A_t \cdot A_{ti} = X_{ti} A_{ti}$$
$$\Rightarrow E_t[\hat{X}_{ti}] = E_t\left[\frac{A_{ti}}{P_{ti}} X_t\right] = E_t\left[\frac{X_{ti} A_{ti}}{P_{ti}}\right] = X_{ti} \cdot \frac{E_t[A_{ti}]}{P_{ti}} = X_{ti}$$

The Loss-based Importance-Weighted Estimators

- The loss-based importance-weighted estimator (unbiased) of x_{ti} is

$$\hat{X}_{ti} = 1 - \frac{I\{A_t = i\}}{P_{ti}}(1 - X_t)$$

- Let $y_{ti} = 1 - x_{ti}$, $Y_t = 1 - X_t$ and $\hat{Y}_{ti} = 1 - \hat{X}_{ti}$

$$\hat{Y}_{ti} = \frac{I\{A_t = i\}}{P_{ti}} Y_t$$

Exp3: Exponential-weight algorithm for Exploration and Exploitation

- The importance-weighted estimator provides us with the means to estimate the reward
- The next step is to choose the distribution over actions
 $P_t = (P_{ti})_i$.
- The simplest algorithm for adversarial bandits is called Exp3
- \hat{S}_{ti} : total estimated reward by the end of round t .
- Exponential weighting with learning rate η

$$P_{ti} = \frac{\exp(\eta \hat{S}_{t-1,i})}{\sum_j \exp(\eta \hat{S}_{t-1,j})}$$

The Exp3 Algorithm

- 1: **Input:** n, K, η
- 2: Set $\hat{S}_{0i} = 0$ for all i
- 3: **for** $t = 1, \dots, n$ **do**
- 4: Calculate the sampling distribution P_t :

$$P_{ti} = \frac{\exp(\eta \hat{S}_{t-1,i})}{\sum_{j=1}^K \exp(\eta \hat{S}_{t-1,j})}$$

- 5: Sample $A_t \sim P_t$ and observe reward X_t
- 6: Calculate \hat{S}_{ti} :
Loss-based.
- 7: **end for**

Algorithm 8: Exp3

The Exp3-IX Algorithm

$$E_t[\hat{Y}_{ti}] = \frac{p_{ti} Y_{ti}}{p_{ti} + r} = Y_{ti} - \frac{r_{ti}}{p_{ti} + r} \varepsilon_{Y_{ti}}$$

Small losses \rightarrow Large reward

\rightarrow Smoothness of p_t \rightarrow more actions with small prob.
will be selected

- The Exp3 algorithm: a sublinear bound on the expected regret, but a large variance
- Modify Exp3 so that the regret stays small in expectation and is simultaneously well concentrated about its mean.
- EXP3-IX: Exp3 with implicit exploration why?
- Loss estimator $\hat{L}_{n,i} = \sum_{t=1}^n \hat{Y}_{ti}$
- Biased estimator with $\gamma > 0$ $r \uparrow \rightarrow$ bias ↑ variance ↓ bias-variance tradeoff.

$$\hat{Y}_{ti} = \frac{I\{A_t = i\}}{P_{ti} + \gamma} Y_t$$

The Exp3-IX Algorithm

- 1: **Input:** n, K, η, γ
- 2: Set $\hat{L}_{0i} = 0$ for all i
- 3: **for** $t = 1, \dots, n$ **do**
- 4: Calculate the sampling distribution P_t :

$$P_{ti} = \frac{\exp(-\eta \hat{L}_{t-1,i})}{\sum_{j=1}^K \exp(-\eta \hat{L}_{t-1,j})}$$

- 5: Sample $A_t \sim P_t$ and observe reward X_t
- 6: Calculate $\hat{L}_{ti} = \hat{L}_{t-1,i} + \frac{\mathbb{I}\{A_t = i\}(1 - X_t)}{P_{t-1,i} + \gamma}$
- 7: **end for**

Algorithm 9: Exp3-IX

Adversary Bandits with Expert Advice

- In this model there are M experts.
- At the beginning of each round the experts announce their predictions of which actions are the most promising
- We allow the experts to report not only a single prediction, but a probability distribution over the actions
- The Predictions of the M experts in round t is represented by a matrix $E^{(t)} \in [0, 1]^{M \times K}$ where the m th row $E_m^{(t)}$, a probability distribution of K , is the recommendation of expert m for round t .

Adversary Bandits with Expert Advice

For rounds $t = 1, 2, 3, \dots, n$:

- 1 Learner observes predictions of all experts, $E^{(t)}$.
- 2 Learner selects a distribution P_t on $[K]$ in some way.
- 3 Action A_t is sampled from P_t and the reward is $X_t = x_{tA_t}$.

The regret of the learner is with respect to the total expected reward of the best expert:

$$R_n = \mathbb{E} \left[\max_{m \in [M]} \sum_{t=1}^n E_m^{(t)} x_t - \sum_{t=1}^n X_t \right]. \quad (18.6)$$

Exp4: Exponential Weighting for Exploration and Exploitation with Experts

Maintains a probability distribution Q_t over experts.

- 1: **Input:** n, K, M, η, γ
- 2: Set $Q_1 = (1/M, \dots, 1/M) \in [0, 1]^{1 \times M}$ (a row vector)
- 3: **for** $t = 1, \dots, n$ **do**
- 4: Receive advice $E^{(t)}$
- 5: Choose the action $A_t \sim P_t$, where $P_t = Q_t E^{(t)}$
- 6: Receive the reward $X_t = x_{tA_t}$
- 7: Estimate the action rewards: $\hat{X}_{ti} = 1 - \frac{\mathbb{I}\{A_t=i\}}{P_{ti} + \gamma}(1 - X_t)$
- 8: Propagate the rewards to the experts: $\tilde{X}_t = E^{(t)} \hat{X}_t$
- 9: Update the distribution Q_t using exponential weighting:

$$Q_{t+1,i} = \frac{\exp(\eta \tilde{X}_{ti}) Q_{ti}}{\sum_j \exp(\eta \tilde{X}_{tj}) Q_{tj}} \quad \text{for all } i \in [M]$$

- 10: **end for**

Algorithm 10: Exp4 algorithm

Exp3 & Exp4 with Additional Exploration

- Uniform exploration with lower probability was proposed in 2012.
- Peter Auer, Nicolò Cesa-Bianchi, Yoav Freund, and Robert E. Schapire, “The nonstochastic multi-armed bandit Problem”, SIAM Journal on Computing, vol. 32, no. 1, pp. 48-77, 2002.
- The fact that additional exploration is not required was observed by Gilles Stoltz in 2005.

EXP3 with Additional Exploration

Algorithm Exp3

Parameters: Real $\gamma \in (0, 1]$.

Initialization: $w_i(1) = 1$ for $i = 1, \dots, K$.

For each $t = 1, 2, \dots$

1. Set

$$p_i(t) = (1 - \gamma) \frac{w_i(t)}{\sum_{j=1}^K w_j(t)} + \underbrace{\frac{\gamma}{K}}_{\text{Additional exploration.}} \quad i = 1, \dots, K.$$

2. Draw i_t randomly accordingly to the probabilities $p_1(t), \dots, p_K(t)$.
3. Receive reward $x_{i_t}(t) \in [0, 1]$.
4. For $j = 1, \dots, K$ set

$$\hat{x}_j(t) = \begin{cases} x_j(t)/p_j(t) & \text{if } j = i_t, \\ 0 & \text{otherwise,} \end{cases}$$

$$w_j(t+1) = w_j(t) \exp(\gamma \hat{x}_j(t)/K).$$

EXP4 with Additional Exploration

Algorithm Exp4

Parameters: Real $\gamma \in (0, 1]$.

Initialization: $w_i(1) = 1$ for $i = 1, \dots, N$.

For each $t = 1, 2, \dots$

1. Get advice vectors $\xi^1(t), \dots, \xi^N(t)$.
2. Set $W_t = \sum_{i=1}^N w_i(t)$ and for $j = 1, \dots, K$ set

$$p_j(t) = (1 - \gamma) \sum_{i=1}^N \frac{w_i(t) \xi_j^i(t)}{W_t} + \frac{\gamma}{K}$$

3. Draw action i_t randomly according to the probabilities $p_1(t), \dots, p_K(t)$.
4. Receive reward $x_{i_t}(t) \in [0, 1]$.
5. For $j = 1, \dots, K$ set

$$\hat{x}_j(t) = \begin{cases} x_j(t)/p_j(t) & \text{if } j = i_t, \\ 0 & \text{otherwise.} \end{cases}$$

6. For $i = 1, \dots, N$ set

$$\begin{aligned} \hat{y}_i(t) &= \xi^i(t) \cdot \hat{x}(t), \\ w_i(t+1) &= w_i(t) \exp(\gamma \hat{y}_i(t)/K) . \end{aligned}$$

Outline

- 1 Introduction
- 2 Overview of Bandit Problems
- 3 Stochastic Bandits
- 4 Greedy Algorithms
- 5 UCB Algorithms
- 6 Thompson Sampling Algorithms
- 7 Gradient Bandit Algorithms
- 8 Adversarial Bandits
- 9 Summary
- 10 References

Main Categories

- Stochastic & Adversary bandits & Markovian bandits
 - ▶ Classical multi-armed bandit
 - ▶ Contextual bandits
 - ▶ Bayesian bandits & index policy (*Gittins' index*)
 - ▶ Infinitely many-armed bandits

Regret Bounds for Main Bandit Algorithms

Algorithm	Environment	Asymptotic Bound	Finite Bound
ϵ -greedy (adaptive ϵ) [16]	K-armed, stochastic	$O(\log H)$	$\sum_t \left(\frac{\epsilon}{d^2 t} + o\left(\frac{1}{t}\right) \right)$ [16]
UCB1 [3]	K-armed, stochastic	$O(\log H)$ [16]	$8 \cdot \left[\sum_{i: \mu_i < \mu^*} \left(\frac{\log t}{\Delta_i} \right) \right] + \left(1 + \frac{\pi^2}{3} \right) \left(\sum_{i=1}^K \Delta_i \right)$ [16]
UCB2 [16]	K-armed, stochastic	$O(\log H)$	$\sum_{i: \mu_i < \mu^*} \left(\frac{(1+\alpha)(1+4\alpha)\log(2e\Delta_i^2 t)}{2\Delta_i} + \frac{c\alpha}{\Delta_i} \right)$ [16]
UCB-Tuned [16]	K-armed, stochastic	—	—
MOSS [12]	K-armed, stochastic	$O(\log H)$	$\min \left\{ \sqrt{tK}, \frac{K}{\Delta} \log \frac{t\Delta^2}{K} \right\}$ [12]
POKER [132]	K-armed, stochastic	—	—
Bayes-UCB [81]	K-armed, stochastic	$O(\log H)$	$\frac{1+\epsilon}{d(\mu_j, \mu^*)} \log(t) + o(\log(t))$ [81]
KL-UCB [58]	K-armed, stochastic	$O(\log H)$	—
Thompson sampling (TS) [124]	K-armed, stochastic	$O(\log H)$	$O \left((\sum_{a=2}^K \frac{1}{\Delta_a^2})^2 \log t \right)$ [4, 5]
Optimistic TS [40]	K-armed, stochastic	—	—
Bootstrap TS [52]	K-armed, stochastic	—	—
Simple (Optimistic) Sampler [33]	K-armed, stochastic	—	—
BESA [18]	K-armed, stochastic	$O(\log H)$ [18]	$O(\log H)$ [18]
Exp3 [17]	K-armed, adversarial	$O(\sqrt{KH})$ [17]	—
Exp4 [17, 24]	K-armed, adversarial, contextual	$\tilde{O}(\sqrt{KH \ln N})$ [24]	—
SAO [31]	K-armed, adversarial	$O(\text{polylog}(H)), O(\sqrt{H})$ [31]	—
LinUCB [91]	K-armed, contextual	$O(\sqrt{dH \ln KH \ln H})$ [91]	See Note
LinTS [32]	K-armed, contextual	—	—
RandomizedUCB [51]	K-armed, contextual	$O(\sqrt{HK \ln(N/\delta)})$ [51]	—
Banditron [78]	K-armed, contextual, nonlinear	$O(H^{2/3})$ [78]	—
NeuralBandit [7]	K-armed, contextual, nonlinear	—	—
ILOVETOCONBANDITS [2]	K-armed, contextual	$O(\sqrt{HK \ln(N/\delta)})$	—
Discounted UCB [59]	K-armed, nonstationary	$O(\sqrt{H})$	$O(\sqrt{t\Gamma \log t})$ [59]
SWUCB [59]	K-armed, nonstationary	$O(\sqrt{H})$	$O(\sqrt{t\Gamma \log t})$ [59]
Adapt-EvE [67]	K-armed, nonstationary	—	—
Exp3.R [6]	K-armed, nonstationary, possibly adversarial	$O(N \sqrt{H \log H})$	—
WLS LinTS [32]	K-armed, nonstationary, contextual	—	—
UCT [85]	∞ -armed, continuous	$O(e^{\sigma D})$ [46]	—
HOO [29]	∞ -armed, metric space	$\tilde{O}(\sqrt{H})$ [30]	—
ComBand [38]	K-armed, adversarial, multi-play	$O(m^{\frac{2}{3}} \sqrt{HK \log K})$ [129]	—
BOLOM [129]	K-armed, adversarial, multi-play	$O(mK^{\frac{2}{3}} \sqrt{H \log H})$	—
Exp3.M [129]	K-armed, adversarial, multi-play	$O(\sqrt{mHK \log(K/m)})$ [129]	—
MP-TS [75]	K-armed, multi-play	$O(\log H)$ [75]	$O(\log t)$
IMP-TS [75]	K-armed, multi-play	—	—

Perspective of Online Learning

Bandit : a research field
independent of classical
RL

Online Learning			
Statistical Learning Theory		Convex Optimization Theory	Game Theory
Online Learning with Full Feedback		Online Learning with Partial Feedback (Bandits)	
Online Supervised Learning		Stochastic Bandit	Adversarial Bandit
First-order Online Learning	Online Learning with Regularization	Stochastic Multi-armed Bandit	Adversarial Multi-armed Bandit
Second-order Online Learning	Online Learning with Kernels	Bayesian Bandit	Combinatorial Bandit
Prediction with Expert Advice	Online to Batch Conversion	Stochastic Contextual Bandit	Adversarial Contextual Bandit
Applied Online Learning		Online Active Learning	Online Semi-supervised Learning
Cost-Sensitive Online Learning	Online Collaborative Filtering	Selective Sampling	Online Manifold Regularization
Online Multi-task Learning	Online Learning to Rank	Active Learning with Expert Advice	Transductive Online Learning
Online Multi-view Learning	Distributed Online Learning	Online Unsupervised Learning (no feedback)	
Online Transfer Learning	Online Learning with Neural Networks	Online Clustering	Online Density Estimation
Online Metric Learning	Online Portfolio Selection	Online Dimension Reduction	Online Anomaly Detection

Outline

- 1 Introduction
- 2 Overview of Bandit Problems
- 3 Stochastic Bandits
- 4 Greedy Algorithms
- 5 UCB Algorithms
- 6 Thompson Sampling Algorithms
- 7 Gradient Bandit Algorithms
- 8 Adversarial Bandits
- 9 Summary
- 10 References

References

- **Reinforcement Learning: An Introduction (second edition)**, Richard Sutton & Andrew Barto, The MIT Press, 2018.
- **Introduction to Multi-Armed Bandits**, Aleksandrs Slivkins, Foundations and Trends in Machine Learning, Vol. 12, No. 1-2, pp 1-286, 2019.
- **Bandit Algorithms**, Tor Lattimore & Csaba Szepesvari, 2019.
- **A Tutorial on Thompson Sampling**, Daniel Russo, Benjamin Van Roy, Abbas Kazerouni, Ian Osband and Zheng Wen, Foundations and Trends in Machine Learning, Vol. 11, No. 1, pp. 1-96, 2018