Multi-arm Bandits

Sutton and Barto, Chapter 2

The simplest reinforcement learning problem



Recall: Sequential Decision Making

- At time t, agent receives an observation from set \mathcal{X} and can choose an action from set \mathcal{A} (think finite for now)
- Goal of the agent is to maximize long-term return



Simple case: One step!

- No x, take an action, observe a reward immediately
- So, a degenerate tree (not truly sequential)
- This is what we call a simple bandit problem
- No credit assignment, only exploration / exploitation
- Later: contextual bandits (there's x, feedback still immediate)
- Lots of applications in ad placement, more recently in large language models

What is a bandit?

- The simplest kind of structure: every node is a copy of every other node, and they are not connected!
- Which means there are no delayed action effects, simplifying credit assignment!
- Therefore, the main problem in bandits is exploration
- Vanilla multi-arm bandits: nodes do not have any observation
- Contextual bandits have observations (more on that later)

- Imagine you have two actions
- You play action I and get a reward of 0
- You play action 2 and get a reward of I
- Which action should you prefer?
- Which action should you try next?

- Imagine you have two actions
- You played action I three times and got rewards of 0, 1, -1
- You played action 2 three times and got a rewards of 1, 10, -10
- Which action should you prefer?
- Which action should you try next?

- Imagine you have two actions
- You played action I for 300 times and got rewards of 0 (200 times), I (50 times), -I (50 times)
- You played action 2 for 300 times and got a rewards of 1 (200 times), 10 (50 times), -10 (50 times)
- Which action should you prefer?
- Which action should you try next?

- Imagine you have two actions
- You played action 1 for 3000 times and got rewards of 0 (300 times), 1 (2000 times), -1 (600 times), +10 (100 times)
- You played action 2 for 3000 times and got a rewards of I (2000 times), 10 (500 times), -10 (500 times)
- Which action should you prefer?
- Which action should you try next?

Main Principles

- Optimize Expected Value
- Other criteria are possible, eg conditional value at risk (CVaR)
- Need to balance exploration (trying all actions) vs exploitation
- Reduce uncertainty in the mean of each action

You are the algorithm! (bandit I)

- Action I Reward is always 8
 - value of action I is $q_*(1) =$
- Action 2 88% chance of 0, 12% chance of 100!
 - value of action 2 is $q_*(2) = .88 \times 0 + .12 \times 100 =$
- Action 3 Randomly between -10 and 35, equiprobable

-10 0
$$q_*(3)$$
 35 $q_*(3) =$

• Action 4 — a third 0, a third 20, and a third from {8,9,..., 18}



 $q_*(4) =$

The k-armed Bandit Problem

- On each of an infinite sequence of time steps, t=1, 2, 3, ..., you choose an action At from k possibilities, and receive a realvalued reward Rt
- The reward depends only on the action taken; it is identically, independently distributed (i.i.d.):

$$q_*(a) \doteq \mathbb{E}[R_t | A_t = a], \quad \forall a \in \{1, \dots, k\}$$
 true values

- These true values are *unknown*. The distribution is unknown
- Nevertheless, you must maximize your total reward
- You must both try actions to learn their values (explore), and prefer those that appear best (exploit)

The Exploration/Exploitation Dilemma

• Suppose you form estimates

 $Q_t(a) \approx q_*(a), \quad \forall a$

action-value estimates

• Define the greedy action at time t as

 $A_t^* \doteq \arg\max_a Q_t(a)$

- If $A_t = A_t^*$ then you are exploiting If $A_t \neq A_t^*$ then you are exploring
- You can't do both, but you need to do both
- You can never stop exploring, but maybe you should explore less with time. Or maybe not.

Action-Value Methods

- Methods that learn action-value estimates and nothing else
- For example, estimate action values as sample averages:

 $Q_t(a) \doteq \frac{\text{sum of rewards when } a \text{ taken prior to } t}{\text{number of times } a \text{ taken prior to } t} = \frac{\sum_{i=1}^{t-1} R_i \cdot \mathbf{1}_{A_i=a}}{\sum_{i=1}^{t-1} \mathbf{1}_{A_i=a}}$

 The sample-average estimates converge to the true values If the action is taken an infinite number of times

$$\lim_{\substack{N_t(a)\to\infty\\ /}} Q_t(a) = q_*(a)$$
The number of times action *a* has been taken by time *t*

The

ε-Greedy Action Selection

- In greedy action selection, you always exploit
- In ε -greedy, you are usually greedy, but with probability ε you instead pick an action at random (possibly the greedy action again)
- This is perhaps the simplest way to balance exploration and exploitation

A simple bandit algorithm

 $\begin{array}{l} \mbox{Initialize, for } a=1 \mbox{ to } k: \\ Q(a) \leftarrow 0 \\ N(a) \leftarrow 0 \end{array} \\ \mbox{Repeat forever:} \\ A \leftarrow \left\{ \begin{array}{l} \arg\max_a Q(a) & \mbox{with probability } 1-\varepsilon & \mbox{(breaking ties randomly)} \\ a \mbox{ random action } & \mbox{with probability } \varepsilon \\ R \leftarrow bandit(A) \\ N(A) \leftarrow N(A) + 1 \\ Q(A) \leftarrow Q(A) + \frac{1}{N(A)} \left[R - Q(A) \right] \end{array} \right.$

One Bandit Task from

The 10-armed Testbed



ε-Greedy Methods on the 10-Armed Testbed



Averaging — learning rule

- To simplify notation, let us focus on one action
 - We consider only its rewards, and its estimate after n-1 rewards:

$$Q_n \doteq \frac{R_1 + R_2 + \dots + R_{n-1}}{n-1}$$

- How can we do this incrementally (without storing all the rewards)?
- Could store a running sum and count (and divide), or equivalently:

$$Q_{n+1} = Q_n + \frac{1}{n} \left[R_n - Q_n \right]$$

• This is a standard form for learning/update rules:

 $NewEstimate \leftarrow OldEstimate + StepSize | Target - OldEstimate |$

Derivation of incremental update

$$Q_n \doteq \frac{R_1 + R_2 + \dots + R_{n-1}}{n-1}$$

$$Q_{n+1} = \frac{1}{n} \sum_{i=1}^{n} R_i$$

= $\frac{1}{n} \left(R_n + \sum_{i=1}^{n-1} R_i \right)$
= $\frac{1}{n} \left(R_n + (n-1) \frac{1}{n-1} \sum_{i=1}^{n-1} R_i \right)$
= $\frac{1}{n} \left(R_n + (n-1)Q_n \right)$
= $\frac{1}{n} \left(R_n + nQ_n - Q_n \right)$
= $Q_n + \frac{1}{n} \left[R_n - Q_n \right],$

Averaging — learning rule

- To simplify notation, let us focus on one action
 - We consider only its rewards, and its estimate after n+1 rewards:

$$Q_n \doteq \frac{R_1 + R_2 + \dots + R_{n-1}}{n-1}$$

- How can we do this incrementally (without storing all the rewards)?
- Could store a running sum and count (and divide), or equivalently:

$$Q_{n+1} = Q_n + \frac{1}{n} \left[R_n - Q_n \right]$$

• This is a standard form for learning/update rules:

 $NewEstimate \leftarrow OldEstimate + StepSize | Target - OldEstimate |$

Tracking a Non-stationary Problem

- Suppose the true action values change slowly over time
 - then we say that the problem is *non-stationary*
- In this case, sample averages are not a good idea (Why?)
- Better is an "exponential, recency-weighted average":

$$Q_{n+1} \doteq Q_n + \alpha \left[R_n - Q_n \right]$$
$$= (1 - \alpha)^n Q_1 + \sum_{i=1}^n \alpha (1 - \alpha)^{n-i} R_i,$$

where α is a constant *step-size parameter*, $\alpha \in (0, 1]$

• There is bias due to Q_1 that becomes smaller over time

Standard stochastic approximation convergence conditions

• To assure convergence with probability I:

$$\sum_{n=1}^{\infty} \alpha_n(a) = \infty \qquad \text{and} \qquad$$

$$\sum_{n=1}^{\infty} \alpha_n^2(a) < \infty$$

• e.g.,
$$\alpha_n \doteq \frac{1}{n}$$

• not
$$\alpha_n \doteq \frac{1}{n^2}$$

if
$$\alpha_n \doteq n^{-p}$$
, $p \in (0,1)$

then convergence is at the optimal rate:

 $O(1/\sqrt{n})$

Optimistic Initial Values

- All methods so far depend on $Q_1(a)$, i.e., they are biased. So far we have used $Q_1(a) = 0$
- Suppose we initialize the action values optimistically ($Q_1(a) = 5$), e.g., on the IO-armed testbed (with $\alpha = 0.1$)



Upper Confidence Bound (UCB) action selection

- A clever way of reducing exploration over time
- Estimate an upper bound on the true action values
- Select the action with the largest (estimated) upper bound

$$A_t \doteq \underset{a}{\operatorname{arg\,max}} \left[Q_t(a) + c \sqrt{\frac{\log t}{N_t(a)}} \right]$$



• Let $H_t(a)$ be a learned preference for taking action a

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

Note that this allows us to work with unnormalized preferences and turn them into probabilities!

Same idea as using potentials in graphical models

• Let $H_t(a)$ be a learned preference for taking action a

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

 $H_{t+1}(A_t) \doteq H_t(A_t) + \alpha \left(R_t - \bar{R}_t \right) \left(1 - \pi_t(A_t) \right)$

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

• Let $H_t(a)$ be a learned preference for taking action a

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

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

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

• Let $H_t(a)$ be a learned preference for taking action a

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

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



Derivation of gradient-bandit algorithm

In exact gradient ascent:

$$H_{t+1}(a) \doteq H_t(a) + \alpha \frac{\partial \mathbb{E} [R_t]}{\partial H_t(a)}, \qquad (1)$$

where:

$$\mathbb{E}[R_t] \doteq \sum_b \pi_t(b) q_*(b),$$

$$\frac{\partial \mathbb{E}[R_t]}{\partial H_t(a)} = \frac{\partial}{\partial H_t(a)} \left[\sum_b \pi_t(b) q_*(b) \right]$$
$$= \sum_b q_*(b) \frac{\partial \pi_t(b)}{\partial H_t(a)}$$
$$= \sum_b \left(q_*(b) - X_t \right) \frac{\partial \pi_t(b)}{\partial H_t(a)},$$

where X_t does not depend on b, because $\sum_b \frac{\partial \pi_t(b)}{\partial H_t(a)} = 0$.

$$\begin{split} \frac{\partial \mathbb{E}[R_t]}{\partial H_t(a)} &= \sum_b \left(q_*(b) - X_t \right) \frac{\partial \pi_t(b)}{\partial H_t(a)} \\ &= \sum_b \pi_t(b) \left(q_*(b) - X_t \right) \frac{\partial \pi_t(b)}{\partial H_t(a)} / \pi_t(b) \\ &= \mathbb{E} \left[\left(q_*(A_t) - X_t \right) \frac{\partial \pi_t(A_t)}{\partial H_t(a)} / \pi_t(A_t) \right] \\ &= \mathbb{E} \left[\left(R_t - \bar{R}_t \right) \frac{\partial \pi_t(A_t)}{\partial H_t(a)} / \pi_t(A_t) \right], \end{split}$$

where here we have chosen $X_t = \overline{R}_t$ and substituted R_t for $q_*(A_t)$, which is permitted because $\mathbb{E}[R_t|A_t] = q_*(A_t)$. For now assume: $\frac{\partial \pi_t(b)}{\partial H_t(a)} = \pi_t(b)(\mathbf{1}_{a=b} - \pi_t(a))$. Then:

$$= \mathbb{E}\left[\left(R_t - \bar{R}_t\right)\pi_t(A_t)\left(\mathbf{1}_{a=A_t} - \pi_t(a)\right)/\pi_t(A_t)\right] \\= \mathbb{E}\left[\left(R_t - \bar{R}_t\right)\left(\mathbf{1}_{a=A_t} - \pi_t(a)\right)\right].$$

 $H_{t+1}(a) = H_t(a) + \alpha (R_t - \bar{R}_t) (\mathbf{1}_{a=A_t} - \pi_t(a)), \text{ (from (1), QED)}$

Thus it remains only to show that

$$\frac{\partial \pi_t(b)}{\partial H_t(a)} = \pi_t(b) \big(\mathbf{1}_{a=b} - \pi_t(a) \big).$$

Recall the standard quotient rule for derivatives:

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

Using this, we can write...

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

$$\begin{split} \frac{\partial \pi_t(b)}{\partial H_t(a)} &= \frac{\partial}{\partial H_t(a)} \pi_t(b) \\ &= \frac{\partial}{\partial H_t(a)} \left[\frac{e^{H_t(b)}}{\sum_{c=1}^k e^{H_t(c)}} \right] \\ &= \frac{\frac{\partial e^{H_t(b)}}{\partial H_t(a)} \sum_{c=1}^k e^{H_t(c)} - e^{H_t(b)} \frac{\partial \sum_{c=1}^k e^{H_t(c)}}{\partial H_t(a)}}{\left(\sum_{c=1}^k e^{H_t(c)}\right)^2} \qquad (Q.R.) \\ &= \frac{\mathbf{1}_{a=b} e^{H_t(a)} \sum_{c=1}^k e^{H_t(c)} - e^{H_t(b)} e^{H_t(a)}}{\left(\sum_{c=1}^k e^{H_t(c)}\right)^2} \qquad (\frac{\partial e^x}{\partial x} = e^x) \\ &= \frac{\mathbf{1}_{a=b} e^{H_t(b)}}{\sum_{c=1}^k e^{H_t(c)}} - \frac{e^{H_t(b)} e^{H_t(a)}}{\left(\sum_{c=1}^k e^{H_t(c)}\right)^2} \\ &= \mathbf{1}_{a=b} \pi_t(b) - \pi_t(b) \pi_t(a) \\ &= \pi_t(b) (\mathbf{1}_{a=b} - \pi_t(a)). \qquad (Q.E.D.) \end{split}$$

Softmax (Boltzmann) Exploration

• Let $H_t(a)$ be a learned preference for taking action a

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

Consider
$$H_t(a) = Q_t(a)/T$$

This is Boltzmann or softmax exploration!

If the temperature T is very large (towards infinity) - same as uniform

If temperature T goes to 0, same as greedy

Summary Comparison of Bandit Algorithms

