

Reinforcement Learning Notes

1. k -armed Bandits

1.1. The k -armed Bandit Problem

Setup

- K different actions to choose from at each time step
- $p(\text{reward}|\text{action})$ is a stationary (not changing over time) probability distribution
- Goal of the agent is to maximize the cumulative reward (over some time span)

What is difficult about solving this problem?

- Do not know what conditional distributions (over reward) are, which means that we need to
- Estimate the expected values of underlying distributions using observed data
- Given that our estimate is never perfect, an important question arises, that is, under what circumstances should we exploit / explore?

1.2. Action Value Methods

Perfect Information (Full Knowledge of Underlying Reward Distributions)

Notation:

- A_t : action (out of k actions) selected at time step t
- R_t : reward received at time step t (after an action has been taken in that time step) (a random variable)

Since we are first dealing with stationary reward distributions, values of actions do not depend on time step. Therefore, the t subscripts merely serve to emphasize that A and R happens in the same time step. In other words, A caused the agent to receive R .

Expected reward given an arbitrary action a :

$$q_*(a) = \mathbb{E}[R_t | A_t = a]$$

Obviously, the action that should be chosen at each and every time step is $\text{argmax}_a q_*(a)$.

Imperfect Information (Can Only Sample From Underlying Reward Distributions)

Since we do not know the true expected rewards (for all actions) but can only sample from their corresponding distributions, we need a way to estimate the expected rewards using the finite samples.

This is intuitively done using the **sample-average method**:

$$Q_t(a) = \frac{\text{sum of rewards when } a \text{ was taken prior to } t}{\text{number of times } a \text{ was taken prior to } t} \quad (1)$$

$$= \frac{\sum_{i=1}^{t-1} R_i \mathbb{I}_{A_i=a}}{\sum_{i=1}^{t-1} \mathbb{I}_{A_i=a}} \quad (2)$$

Different algorithms use the sample-average rewards in different ways:

- Greedy method
 - Initialization: randomly choose between actions
 - Action selection: $A_t = \operatorname{argmax}_a Q_t(a)$ (never explores)
 - Disadvantage: It is very likely that the true expected rewards of other non-greedy actions are higher, but this algorithm would never know because it never explore.
- ϵ -greedy method
 - Initialization: randomly choose between actions
 - Action selection:
 - * Behave greedily most of the time: $A_t = \operatorname{argmax}_a Q_t(a)$
 - * Once in a while, with small probability ϵ , select randomly from all other action
 - Advantage: In the limit as the number of steps increases, every action will be sampled an infinite number of times, thus ensuring that $Q_t(a)$ converges to $q_*(a)$.
 - Disadvantage: Such asymptotic guarantees say little about its practical effectiveness.
- Optimistic initial values method
 - Initialization: set the action values to very high (and non-realistic) values, randomly choose between actions
 - Action selection:
 - * $A_t = \operatorname{argmax}_a Q_t(a)$, but explores much more than greedy method without optimistic initial values and here's why:
 - * For example, let's say that $q_*(a)$'s in a multi-armed bandit problem are selected from a normal distribution with mean 0 and variance 1. In this case, an initial estimate of +5 is very optimistic for all $q_*(a)$'s.
 - * Whichever action is initially selected, the reward is less than the starting estimate and the estimate is decreased; the learner therefore switches to other actions.
 - * All actions are explored several times before the value estimates $q_*(a)$ converge.
- Upper confidence bound (UCB) method
 - Initialization: randomly choose between actions
 - Action selection:
 - * $A_t = \operatorname{argmax}_a [Q_t(a) + c\sqrt{\frac{\ln t}{N_t(a)}}]$
 - * Motivation:
 - ϵ -greedy forces non-greedy actions to be tried, but indiscriminately.
 - It would be better to select among non-greedy actions
 1. according to their potential for being actually optimal ($Q_t(a)$)
 2. take into account uncertainties in their value estimates ($\sqrt{\frac{\ln t}{N_t(a)}}$)
 - Disadvantage: hard to generalize to other reinforcement learning tasks (we do not know why yet)

1.3. The 10-armed Testbed

The 10-armed testbed consists of 1000 tasks similar to the one shown in Fig. 1.

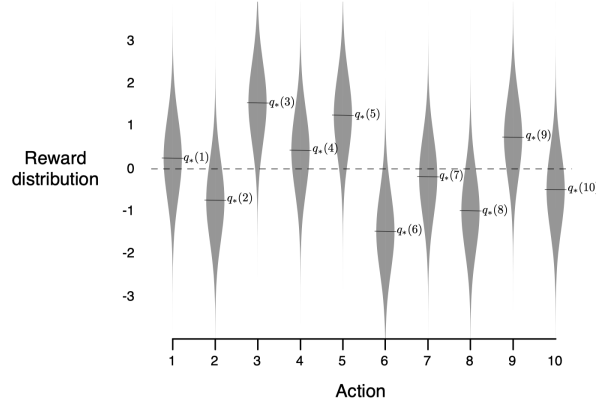


Figure 1: One task from the 10-armed testbed.

In Fig. 2, we can see that the highest average reward over first 1000 steps is achieved by UCB.

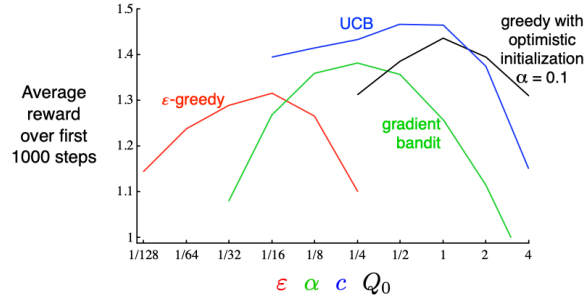


Figure 2: Performance of various algorithms on the 10-armed testbed.

1.4. Compute Sample-Averages Incrementally

If we focus on one action a_k , let R_i denote the reward received by the agent after the i th selection of this action and let Q_n denote the estimated reward:

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

Native implementation: Keep a list of R_i 's and compute the sum every time.

Incremental implementation: One can show that Q_{n+1} can be conveniently expressed as a function of Q_n :

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

where R_n is a new reward received as a consequence of a_k .

1.5. Tracking Non-stationary Problems

In this case, it makes sense to give more weight to recent rewards than to long-past rewards. One of the most popular ways of doing this is called **exponential-recency-weighted average** (still a way to estimate the value function).

The incremental update rule is modified to:

$$Q_{n+1} = Q_n + \alpha[R_n - Q_n]$$

where $\alpha \in (0, 1]$ and a common choice of α is 0.99.

To understand this update rule more intuitively, let's expand it. When expanded, this rule becomes:

$$Q_{n+1} = (1 - \alpha)^n Q_1 + \sum_{i=1}^n \alpha (1 - \alpha)^{n-i} R_i$$

We can see, as i decreases, the weight given to the i th reward decreases exponentially.

1.6. Gradient Bandit Algorithms

In previous sections we've estimated the value of each action and used those estimates to take semi-informed actions. Decision making based on value estimates can work quite well, but it is not the only selection method available. One alternative is learn a preference for each action denoted $H_t(a)$ and makes decisions based on the relative preference of one action over another. Here we will use the *soft-max distribution* to define preference as a distribution over all possible actions.

$$Pr\{A_t = a\} \doteq \frac{e^{H_t(a)}}{\sum_{b=1}^k e^{H_t(b)}} \doteq \pi_t(a) \quad (3)$$

In this case $\pi_t(a)$ is the probability of taking action a at time t . All preferences are initially set to be the same.

Now we will begin to introduce an algorithm for this setting based on stochastic gradient descent using the following pair of formulas:

$$H_{t+1}(A_t) \doteq H_t(A_t) + \alpha(R_t - \bar{R}_t)(1 - \pi_t(A_t)), \quad (4)$$

$$H_{t+1}(a) \doteq H_t(a) + \alpha(R_t - \bar{R}_t)\pi_t(a), \forall a \neq A_t \quad (5)$$

Where $\alpha > 0$ is the step size parameter and \bar{R}_t is the average reward of all rewards up to time t . Note: if \bar{R}_t is set to be a constant 0 the performance of the algorithm is reduced (e.g. any baseline works so long as it exists and isn't constant).

Now let's walk through how the gradient bandit algorithm approximates to gradient ascent.

$$H_{t+1}(a) \doteq H_t(a) + \alpha \frac{\partial \mathbb{E}[R_t]}{\partial H_t(a)}, \quad (6)$$

Where the expected reward is:

$$\mathbb{E}[R_t] = \sum_x \pi_t(x) q_*(x),$$

$$\begin{aligned} \frac{\partial \mathbb{E}[R_t]}{\partial H_t(a)} &= \frac{\partial}{\partial H_t(a)} \left[\sum_x \pi_t(x) q_*(x) \right] \\ &= \sum_x q_*(x) \frac{\partial \pi_t(x)}{\partial H_t(a)} \\ &= \sum_x (q_*(x) - B_t) \frac{\partial \pi_t(x)}{\partial H_t(a)} \end{aligned}$$

We can add the B_t during the final step above because the total change in $H_t(a)$ must be zero since all probabilities must sum to one.

$$\sum_x \frac{\partial \pi_t(x)}{\partial H_t(a)} = 0$$

Now we multiply by $\pi_t(x)/\pi_t(x)$

$$\frac{\partial \mathbb{E}[R_t]}{\partial H_t(a)} = \sum_x (q_*(x) - B_t) \frac{\partial \pi_t(x)}{\partial H_t(a)} \pi_t(x)/\pi_t(x)$$

Now we have an expectation that will allow us to sum over all possible values of x of the random variable A_t

$$\begin{aligned} &= \mathbb{E}[(q_*(A_t) - B_t) \frac{\partial \pi_t(A_t)}{\partial H_t(a)} / \pi_t(A_t)] \\ &= \mathbb{E}[(R_t - \bar{R}_t) \frac{\partial \pi_t(A_t)}{\partial H_t(a)} / \pi_t(A_t)] \end{aligned}$$

In this case we let $B_t = \bar{R}_t$ and $R_t = q_*(A_t)$, which we can do since $\mathbb{E}[R_t|A_t] = q_*(A_t)$

Briefly let $\frac{\partial \pi_t(x)}{\partial H_t(a)} = \pi_t(x)(\mathbb{1}_{a=x} - \pi_t(a))$ without proof netting us:

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

Substituting this back into (4) we get:

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

Which turns out to be equivalent to what was laid out in (2) and (3).

Now let's go back and show $\frac{\partial \pi_t(x)}{\partial H_t(a)} = \pi_t(x)(\mathbb{1}_{a=x} - \pi_t(a))$ starting with the quotient rule:

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

Thus:

$$\begin{aligned}
\frac{\partial \pi_t(x)}{\partial H_t(a)} &= \frac{\partial}{\partial H_t(a)} \pi_t(x) \\
&= \frac{\partial}{\partial H_t(a)} \left[\frac{e^{H_t(x)}}{\sum_{y=1}^k e^{H_t(y)}} \right] \\
&= \frac{\frac{\partial e^{H_t(x)}}{\partial H_t(a)} \sum_{y=1}^k e^{H_t(y)} - e^{H_t(x)} \frac{\partial \sum_{y=1}^k e^{H_t(y)}}{\partial H_t(a)}}{(\sum_{y=1}^k e^{H_t(y)})^2} \\
&= \frac{\mathbb{1}_{0=x} e^{H_t(x)} \sum_{y=1}^k e^{H_t(y)} - e^{H_t(x)} e^{H_t(a)}}{(\sum_{y=1}^k e^{H_t(y)})^2} \\
&= \frac{\mathbb{1}_{0=x} e^{H_t(x)}}{\sum_{y=1}^k e^{H_t(y)}} - \frac{e^{H_t(x)} e^{H_t(a)}}{(\sum_{y=1}^k e^{H_t(y)})^2} \\
&= \mathbb{1}_{0=x} \pi_t(x) - \pi_t(x) \pi_t(a) \\
&= \pi_t(x) (\mathbb{1}_{0=x} - \pi_t(a))
\end{aligned}$$

1.7. Associative Search (Contextual Bandit)

Up to this point the formulations of the bandit problem have been non-associative tasks, meaning there is no need to tie certain actions to certain situations. The algorithms we have seen thus far have sought to find a best action that is either stationary or very slowly changing, but if we want to expand our horizons to the general reinforcement learning task we need to be able to learn a policy (a mapping from a situation a set of actions).

If for example we saw a version of the bandit problem where the value of each action changed randomly and drastically with each step the algorithms we've previously seen would be worthless. But if we have some distinct and predictable clue that queues us into how the situation has changed (but not its action values) then we may be able to begin to learn a policy associating that clue with a best action.

1.8. Summary

Despite the simplicity of these methods they can actually be used to approach very complex problems and are in some way cutting edge. In the form we've laid out so far, however, they lack the complexity needed to tackle the full reinforcement learning problem, but as we'll see later on they can be used in that context when proper considerations are taken.

Given that in our formulations of the bandit problem we are trying to discover the expected values of a set of actions its easy to formulate a Bayesian approach to the problem. Using Bayesian methods to solve the problem can be computationally complex, but allow us to solve the problem very precisely. By setting prior distributions for each action and then iterative computing the posterior probability that a given action is optimal. It is even possible to find the optimal approach by computing the possible outcome over all values, although in practice doing so is much too complex. Approximation may help reduce that problem and will be touched on later, but is an ongoing research topic at the moment.

2. Finite Markov Decision Processes

In this section we're going to go over the basics of Markov decision processes, or MPDs. When solving MPDs we estimate the value $q_*(s, a)$ of each action a in each state s . Alternatively, we can estimate the value $v_*(s)$ of each state given the optimal action is chosen.

2.1. The Agent-Environment Interface

At each step the agent sees a state $S_t \in S$ and chooses an action $A_t \in A(s)$. The agent then receives a reward R_{t+1} during the next step and enters state S_{t+1} .

For random variables $s' \in S, r \in R$ there is a chance of those values occurring at time s given the previous state and action:

$$p(s', r | s, a) \doteq \Pr\{S_t = s', R_t = r | S_{t-1} = s, A_{t-1} = a\} \quad (7)$$

It's important to note that MDPs are *Markov* processes, because the probabilities for S_t and R_t are only dependent on the immediately preceding S_{t-1} and A_{t-1} . The above equation can be slightly augmented to produce the equations below.

2.2. Goals and Rewards

At each step the agent receives a reward $R_t \in \mathbb{R}$. The goal of the agent is to maximize the total amount of reward it receives over time. That is to say, that the agent does not seek to maximize its next reward, but the total reward it will receive over the long run.

Given this simple definition of an agents reward seeking behavior, our reward choice can be quite flexible. For example, an agent learning to play a game can be given +1 for a win, 0 for a tie, and -1 for a loss, or a maze solving agent can be given a reward of -1 for each second it doesn't find the solution. In both cases the agent learns to solve the problem by maximizing the reward, but how we choose that reward determines *what* we want the agent to do, but not *how* we want it to do it.

2.3. Returns and Episodes

As previously, mentioned we want to maximize the cumulative reward, but more generally we want to maximize the expected return. Return in this case is defined as a specific function of the reward sequence $G_t = R_{t+1} + R_{t+2} + R_{t+3} + \dots + R_T$ where T is the final time step. Tasks where we have a final step are called *episodic* tasks.

Sometimes we're able to break up a task into a series of *episodes*, but in many cases we can't. In those cases, when we have a continual process and are working on a *continuing task*, our previous return formula is no long valid given that $T = \infty$. To combat that issue we can use *discounting*, where we work with a sum of discounted rewards to maximize the *discounted return* given a chosen A_t :

$$G_t \doteq R_{t+1} + \gamma R_{t+2} + \gamma^2 R_{t+3} + \dots = \sum_{k=0}^{\infty} \gamma^k R_{t+k+1} = \quad (8)$$

where $0 \leq \gamma \leq 1$ is the discount rate.

In selecting a discount rate we choose how much less we value a future reward than an immediate reward. If γ is 0 or close to 0 than we focus near exclusively on the immediate reward, while a γ near 1 will take into account further out scenarios.

As is the case in many reinforcement learning algorithms we can define G_t in terms of successive steps:

$$\begin{aligned} G_t &\doteq R_{t+1} + \gamma R_{t+2} + \gamma^2 R_{t+3} + \dots \\ &= R_{t+1} + \gamma(R_{t+2} + \gamma R_{t+3} + \dots) \\ &= R_{t+1} + \gamma G_{t+1} \end{aligned} \quad (9)$$

We should also mention that when $\gamma < 1$ and the reward is constant

$$G_t = \sum_{k=0}^{\infty} \gamma^k = \frac{1}{1 - \gamma} \quad (10)$$

2.4. Unified Notation for Episodic and Continuing Tasks

To simplify the notation used in the section above, we define the expected return as:

$$G_t = \sum_{k=t+1}^T \gamma^{k-t-1} R_k, \quad (11)$$

with the possibility that $T = \infty$ or $\gamma = 1$ but not both.

2.5. Golf Example (Optimal Value Functions)

Imagine we can use either a putter or driver to complete a section of a golf course. At some points the driver is advantageous and at others the putter is better. We want our agent to learn how to optimally complete the hole through a combination of puts and drives.

We'll approach this problem through the Bellman optimality equation, since the value of a state under an optimal policy must equal the return of the best action.

$$\begin{aligned} v_*(s) &= \max_{a \in A(s)} q_{\pi_*}(s, a) \\ &= \max_a \mathbb{E}_{\pi_*}[G_t | S_t = s, A_t = a] \\ &= \max_a \mathbb{E}_{\pi_*}[R_{t+1} + \gamma G_{t+1} | S_t = s, A_t = a] \\ &= \max_a \mathbb{E}_{\pi_*}[R_{t+1} + \gamma v_*(S_{t+1}) | S_t = s, A_t = a] \\ &= \max_a \sum_{s', r} p(s', r | s, a) [r + \gamma v_*(s')] \end{aligned} \quad (12)$$

Then

$$\begin{aligned} q_*(s, a) &= \mathbb{E}[R_{t+1} + \gamma \max_{a'} q_*(S_{t+1}, a') | S_t = s, A_t = a] \\ &= \sum_{s', r} p(s', r | s, a) [r + \gamma \max_{a'} q_*(s', a')] \end{aligned} \quad (13)$$

With these two sets of equations we can solve the question of which action is optimal relatively easily. This can be done by either finding a v_* greedily since the optimal policy will give an optimal value. Similarly with q_* the agent can simply solve for the best action at every state in order to optimize q_* .

3. Dynamic Programming

In this section we'll discuss how dynamic programming can be used to solve the equations laid out in the previous section.

3.1. Policy Evaluation

Given that equation 14 satisfies the Bellman optimality principle (which it should considering that unless something has been messed up we should achieve convergence as $k \rightarrow \infty$) we should be able to solve for v_π using *iterative policy evaluation*. That is to say, that we can repeatedly generate successive approximations until we believe convergence has been achieved. Each update is considered an *expected*

update due to our evaluation of the expectation of the function over all possible next states.

In the algorithm below we can see generally speaking how iterative policy evaluation is handled, but we still haven't seen precisely how $V(s)$ is calculated with each passing loop. In the next section we'll approach the computation behind $V(s)$.

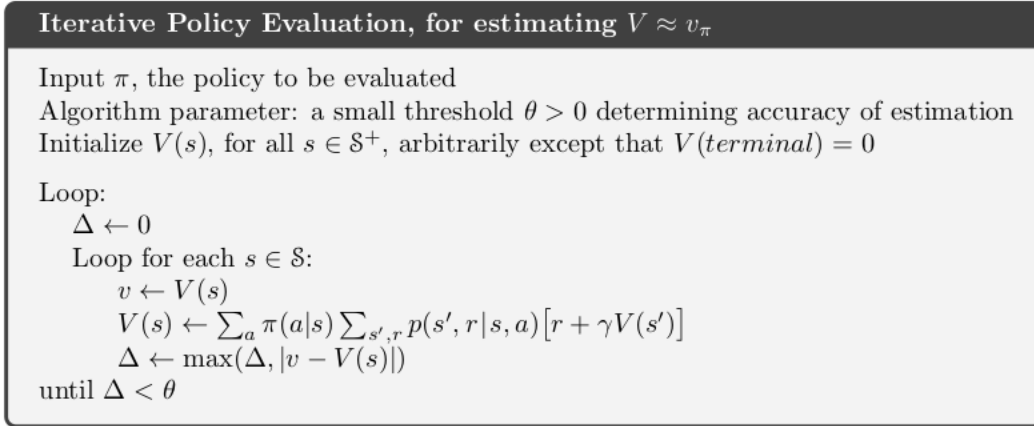


Figure 3

3.2. Policy Improvement

Using the algorithm from the previous section we can find some deterministic policy $\pi(s)$ using v_π . That policy could be very good, but what if we want to choose an action that deviates from $\pi(s)$. To do that we need to find a way to update and improve or policy.

3.3. Policy Iteration

Each policy is updated through a sequence of evaluation and improvement operations, which are guaranteed to result in a new policy that is strictly better than the previous one. Combined with the nature of MDPs (which have a finite number of policies) we can converge upon an optimal policy within a finite number of steps using solely evaluation and improvement. We call this process policy iteration.

References

Reinforcement Learning: An Introduction, 2nd Edition, Richard S. Sutton, Andrew G. Barto, 2018.

Policy Iteration (using iterative policy evaluation) for estimating $\pi \approx \pi_*$

1. Initialization
 $V(s) \in \mathbb{R}$ and $\pi(s) \in \mathcal{A}(s)$ arbitrarily for all $s \in \mathcal{S}$
2. Policy Evaluation
Loop:
 $\Delta \leftarrow 0$
 Loop for each $s \in \mathcal{S}$:
 $v \leftarrow V(s)$
 $V(s) \leftarrow \sum_{s',r} p(s', r | s, \pi(s)) [r + \gamma V(s')]$
 $\Delta \leftarrow \max(\Delta, |v - V(s)|)$
 until $\Delta < \theta$ (a small positive number determining the accuracy of estimation)
3. Policy Improvement
 policy-stable \leftarrow *true*
 For each $s \in \mathcal{S}$:
 old-action $\leftarrow \pi(s)$
 $\pi(s) \leftarrow \operatorname{argmax}_a \sum_{s',r} p(s', r | s, a) [r + \gamma V(s')]$
 If *old-action* $\neq \pi(s)$, then *policy-stable* \leftarrow *false*
 If *policy-stable*, then stop and return $V \approx v_*$ and $\pi \approx \pi_*$; else go to 2

Figure 4