

Reinforcement Learning Notes

Kiran Kannar

July 17, 2025

Contents

1	Introduction to Reinforcement Learning	2
1.1	Elements of RL	2
2	Multi-Armed Bandit Problems	3
2.1	Action-Value Methods	3
2.1.1	Sample-Average Method:	3
2.1.2	Non-Stationary Problems: Exponential Recency Weighting	4
2.2	Gradient Bandit	5
2.3	Contextual Bandits	5

1 Introduction to Reinforcement Learning

Computational approach to **goal-directed** learning by an agent from interactions with an environment.

- Trial-and-error search
- Delayed reward
- Exploration vs exploitation
- Model-based vs model-free

Broadly, maximize reward signal, despite uncertainty about the environment → Optimization problem.

1.1 Elements of RL

A Markov Decision Process (MDP) involving:

- Policy : Mapping from states to actions
- Reward Signal : Scalar signal that indicates the **desirability** of the state i.e. immediate value
- Value Function : Expected cumulative reward from a state i.e. far-sighted judgment of value of starting from a particular state.
- Model (of the environment): Predicts the next state and reward given the current state and action; used for planning.

Value function is computed without **explicit search** over possible sequences of future states and actions. Focus is on highest value, and not on highest reward, even though value is computed from rewards.

2 Multi-Armed Bandit Problems

Setting: A single state (aka situation) with multiple actions, one of which must be selected repeatedly. Each arm is slot machine with a stationary probability distribution of rewards, say, $\mathcal{N}(\mu_i, \Sigma_i)$. The goal is to maximize the expected reward over some time horizon.

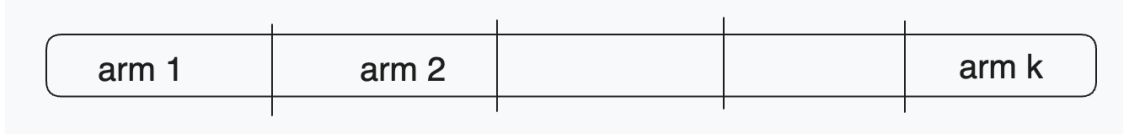


Figure 1: A single state with multiple actions.

The expected reward of taking an action is given by:

$$q_*(a) = \mathbb{E}[R_t | A_t = a] \quad (1)$$

This is the value of taking that action. We do not however know the value of each action, but we can estimate it by sampling the reward. Thus, we approximate the value of an action as $Q_t(a) \approx q_*(a)$.

2.1 Action-Value Methods

Methods to estimate the value of an action and then use that estimate to select actions.

2.1.1 Sample-Average Method:

Each estimate is the average of all the sample of rewards collected up to that point. Suppose $Q_t(a) = \frac{1}{N_t(a)} \sum_{i=1}^{N_t(a)} R_i$, where $N_t(a)$ is the number of times action a has been selected up to time t .

$$Q_t(a) = \frac{\sum_{i=1}^{t-1} R_i \cdot 1(A_i = a)}{\sum_{i=1}^{t-1} 1(A_i = a)} \quad (2)$$

Action Selection:

- Greedy: $A_t = \operatorname{argmax}_a Q_t(a)$
- ϵ -greedy: $A_t = \begin{cases} \operatorname{argmax}_a Q_t(a) & \text{with probability } 1 - \epsilon \\ \text{random action} & \text{with probability } \epsilon \end{cases}$

With ϵ -greedy, as $t \rightarrow \infty$, each action is sampled infinite times, and therefore the estimate $Q_t(a)$ converges to $q_*(a)$. The greedy selection, however, does not have strong convergence properties.

With 'k' arms in ϵ -greedy method, the probability of selecting the greedy action is,

$$\begin{aligned}
P(\text{greedy action}) &= (1 - \epsilon).P(\text{action is greedy}) + \epsilon.P(\text{action is random}) \\
&= (1 - \epsilon).1 + \epsilon.\frac{1}{k} \\
&= 1 - \epsilon + \frac{\epsilon}{k}
\end{aligned} \tag{3}$$

The expected reward using ϵ -greedy method is $(1 - \epsilon).q_*(a^*)$, because of the strong convergence properties of the sample-average method. That is, in the long run, the expected reward will be equal to the expected value of selecting the greedy action.

- If the variance of rewards increases across arms, it's always good to explore quickly. ϵ -greedy method is better than greedy method.

- If the task is non-stationary, then the ϵ -greedy method is good because we can explore previously explored actions whose rewards have potentially increased.

Incremental Implementation:

$$Q_{t+1}(a) = Q_t(a) + \frac{1}{N_t(a)}(R_t - Q_t(a)) \tag{4}$$

Simple bandit algorithm:

Algorithm 1 Simple Bandit Algorithm

```

 $Q(a) \leftarrow 0$ 
 $N(a) \leftarrow 0$ 
for  $t = 1$  to  $T$  do
   $A_t \leftarrow \begin{cases} \operatorname{argmax}_a Q(a) & \text{with probability } 1 - \epsilon \\ \text{random action} & \text{with probability } \epsilon \end{cases}$ 
   $R_t \leftarrow \text{reward}(A_t)$ 
   $Q(A_t) \leftarrow Q(A_t) + \frac{1}{N(A_t)}[R_t - Q(A_t)]$ 
   $N(A_t) \leftarrow N(A_t) + 1$ 
end for

```

2.1.2 Non-Stationary Problems: Exponential Recency Weighting

In the case of non-stationary problems, we would want to give more weight to recent rewards.

$$\begin{aligned}
Q_{t+1}(a) &= Q_t(a) + \alpha[R_t - Q_t(a)] \\
&= (1 - \alpha)^t Q_1(a) + \sum_{i=1}^t \alpha(1 - \alpha)^{t-i} R_i
\end{aligned} \tag{5}$$

Note that all of the above methods are biased by the initial estimate $Q_1(a)$ (aka prior knowledge). For sample average methods, the bias disappears once all actions are explored at least once. $Q_t(a) \gg R_t$ allows us to explore the actions more quickly. Wild optimism is a good strategy. For non-stationary problems, however, the bias is persistent but decreasing over time. So this drive for exploration is temporary.

2.2 Gradient Bandit

Instead of estimating the value of each action, we estimate the numerical preference for each action, $H_t(a)$. With exact gradient ascent,

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

$$\pi_t(A_t = a) = P(A_t = a) = \frac{e^{H_t(a)}}{\sum_{b=1}^k e^{H_t(b)}} \quad (7)$$

Hence, $\mathbb{E}[R_t] = \sum_x \pi_t(x) q_*(x)$; Given that we do not know $q_*(x)$, we can show that this is an instance of stochastic gradient ascent with robust convergence properties.

$$H_{t+1}(a) = \begin{cases} H_t(a) + \alpha(R_t - R_{avg_t})(1 - \pi_t(A_t)) & \text{if } a = A_t \\ H_t(a) - \alpha(R_t - R_{avg_t})\pi_t(a) & \text{for all } a \neq A_t \end{cases} \quad (8)$$

2.3 Contextual Bandits

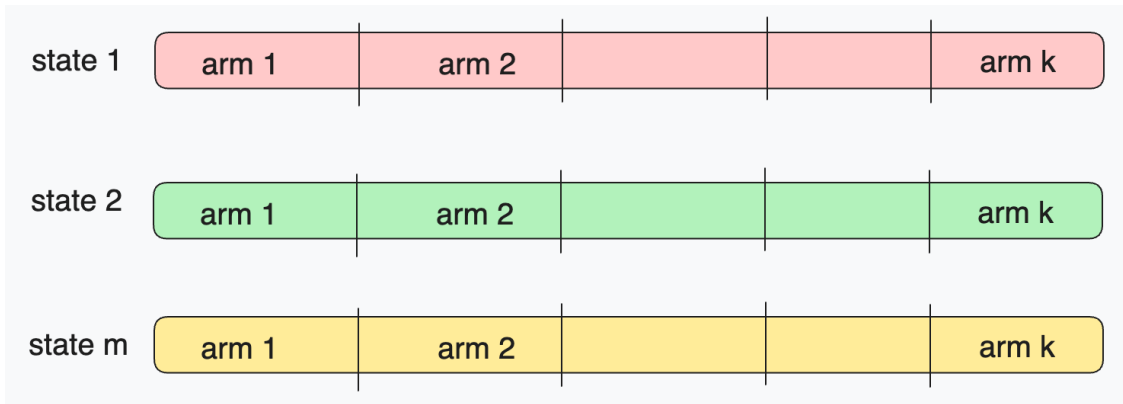


Figure 2: Contextual Bandit with many states, each with multiple actions.

Associative search is a contextual bandit problem where we do trial-and-error learning to

- search for the best action, and,
- the association of these actions with the best situation (state) to be in.

Each (s, a) pair results in a reward $R_{s,a}$. However, it does not change the state s . If the state changes, we are in a full reinforcement learning problem.