

# Reinforcement Learning: An Introduction

Rishabh Sambare

April 28, 2021

# Contents

<b>1</b>	<b>Introduction</b>	<b>2</b>
<b>2</b>	<b>Multi-armed Bandits</b>	<b>4</b>
2.1	A $k$ -armed Bandit Problem . . . . .	4
2.2	Action-value Methods . . . . .	4
2.3	The 10-armed Testbed . . . . .	5
2.4	Tracking a Non-stationary Problem . . . . .	6
2.5	Optimistic Initial Values . . . . .	7
2.6	Upper Confidence Bound Action Selection . . . . .	7
2.7	A Gradient Bandit Algorithm . . . . .	8

# Chapter 1

## Introduction

Reinforcement learning is learning an action to maximize numeric reward. The exploitation vs exploration problem is equivalent to asking when you should stick to your guns, or when you should take a risk.

RL agents are assumed to be capable of interaction with their environments. Some RL methods can learn with approximators which can address the curse of dimensionality on supervised learning methods.

The general elements of an RL system are as follows:

- An agent
- An environment
- A policy
- A reward signal
- A value function
- (optional) An internal environment model

A policy defines an agent's way of behaving at a certain time. Policies are functions from an environment perception to actions to be taken. Policies may be stochastic.

A reward signal defines the goal for a reinforcement learning problem. On each time step, the environment provides (the notion of the environment providing is a convention, this signal could be considered internal to the agent, but not editable by the agent) the agent with a reward. The agent attempts to maximize this reward over the long run. Reward functions may be stochastic.

The value function defines how useful a state is in the long-term. A state may net a high reward signal, but force subsequent states to yield low reward, or vice versa.

Models of the environment are the primary means by which a reinforcement learning agent can plan. Model-free methods are pure trial-and-error, while

model-based methods use model to try and infer the resultant environmental states of certain actions/policies.

I don't know why I added this next paragraph... Sutton says it's not relevant but includes it anyway. Evolutionary reinforcement learning methods apply multiple 'static' policies that interact with separate environment instances. The policies with the highest reward, and some random variations (exploration), are carried over to the 'next generation' of policies. This cycle repeats with some dynamic programming algorithm.

## Chapter 2

# Multi-armed Bandits

Reinforcement learning uses training information that evaluates actions; it does not and cannot instruct the correct action. The evaluative feedback of an action then obviously will depend on the action itself, while instructive feedback is independent of the action an agent it takes; it just provides the right one.

### 2.1 A $k$ -armed Bandit Problem

We define the  $k$ -armed bandit problem to be when an agent is faced repeatedly with a choice among  $k$  options. After the choice, the agent receives a reward from a probability distribution unique to the chosen action. The goal is to maximize the expected total reward over  $t$  time steps. We can model this expected reward as  $q_*(a)$ , where

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

We do not know  $q_*$  exactly for any action  $a$ . We can make estimates  $Q_t(a)$ , where  $t$  is the time step of that particular estimate. Taking an action  $a$  for which  $Q_t(a)$  is the maximum estimated expected reward is a 'greedy' action exploitation. Selecting a 'non-greedy' estimate is an exploration action, as now you can improve your estimate of the non-greedy action and 'explore' your action values.

### 2.2 Action-value Methods

These methods refer to the class of methods by which we can estimate the value of an action and how to select an action. The sample-average method to estimate the expected reward of an action is to take the mean of all previous rewards from the action.

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

where  $1_{A_i=a}$  is the predicate function for when the action  $A_i = a$ . We can see that as  $t \rightarrow \infty$ ,  $Q_t(a) \rightarrow q_*(a)$ .

Greedy action selection is written as

$$A_t = \arg \max_a Q_t(a).$$

We define  $\varepsilon$ -greedy action selection methods to be methods that randomly choose non-greedy actions with probability  $\varepsilon$ , then choosing between greedy actions randomly  $1 - \varepsilon$  of the time.

## 2.3 The 10-armed Testbed

An assessment of greedy and  $\varepsilon$ -greedy action selections was done on 2000 randomly generated 10-armed bandit problems. The  $q(a)$  value for  $a = 1, 2, 3 \dots 10$  was sampled from a normal distribution with variance 1 and mean 0. When an agent picked  $A_t$  at time  $t$ , the reward  $R_t$  was selected from a normal distribution with mean  $q(A_t)$  and variance 1. Basically, the greedy version sucked because it never was greedy on the right action. The 0.01-greedy action was alright initially, but did the best later as it would be non-greedy 1% of the time, so once it found the optimal action it stuck with it longer than the 0.1-greedy action, which, despite finding the optimal action sooner, didn't stick with it as much of the time.

We efficiently compute sample averages for an action's estimated reward through a recursive formula that only requires  $O(1)$  computations at each update. LET  $Q_n$  denote the  $n$ -th estimate of an action's reward. Then

$$Q_{n+1} = \frac{1}{n} \sum_{i=1}^n R_i = Q_n + (1/n)[R_n - Q_n]$$

A general computation in reinforcement learning algorithms is similar

$$newEstimate = oldEstimate + stepSize[Target - Old]$$

where the newest sample, though possibly noisy, is considered the 'target' value.

A simple bandit algorithm half-converted to Python:

```
import random
def simple_bandit(k, eps):
    Q = []
    N = []
    for a in range(k):
        Q[a] = 0
        N[a] = 0
    while True:
        greedy = findIndex(max(Q))
        rand = random.uniform(0,1)
        if rand <= eps:
            A = # random non-max member of Q
        else:
            A = Q[greedy]

        R = bandit(A) # samples a reward from action's probability distribution
        N[A] = N[A] + 1
        Q[A] = Q[A] + 1/N[A] (R-Q[A])
```

Nonstationary or stochastic environments basically eliminate the feasibility of greedy methods since the rewards are constantly changing, so enter...

## 2.4 Tracking a Non-stationary Problem

In non-stationary environments, it makes more sense to prefer newer rewards to older ones. Modifying the incremental update rules:

$$Q_{n+1} = Q_n + \alpha(R_n - Q_n)$$

where  $\alpha \in [0, 1]$  is constant.

$$Q_{n+1} = Q_n + \alpha(R_n - Q_n)$$

$$Q_{n+1} = (1 - \alpha)Q_n + \alpha R_n$$

$$Q_{n+1} = \alpha R_n + (1 - \alpha)[\alpha R_{n-1} + (1 - \alpha)Q_{n-1}]$$

This next step kinda big, if you're invested you should probably just look up the section in the book,

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

The first term is the original estimate of  $Q$ . This estimate matters less as  $n \rightarrow \infty$ . The second term with the summation sums all rewards with each time step lowering a reward's weight on the new estimate by a multiplicative factor  $1 - \alpha$ . The weights overall sum to 1, use a geometric series on the second term if you wanna check. This is called an exponential-recency weighted average.

For stationary environments, converge on true action values require your step-size parameter sequences to meet the following conditions

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

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

This is a result of Blum and is actually a kind of interesting and accessible to an extent paper. I recommend looking at it for a bit if you find the two conditions kind of arbitrary like I did. Just black box the weird measure theory parts if you don't know any, it still kinda reads well. Intuitively, the first condition guarantees that fluctuations do not affect the final expected value (if the sum of the step parameters converge on a real number, then fluctuating reward values will make up a non-trivial percentage of the final expected value approximation). The second step guarantees that the steps become small enough such that the sequence of estimates converges.

## 2.5 Optimistic Initial Values

All current methods discussed are biased on initial estimates

- For sample-average, bias disappears after one sample.
- For constant- $\alpha$  methods, the initial bias is permanent though its weight converges to 0 over repeated sampling.

Initial values can be set to initial expectations of the variable, or you can use them to explore.

Suppose you have all actions with mean 0 and variance 1. If you set the initial estimate of the rewards to be +5, and the learning agent repeatedly gets less, this lowers the overall estimate of the reward, and the learner tries another random action with estimated reward +5. This process repeats until the initial +5 estimate's weighting on the estimated reward becomes negligible. This leads to lots of initial exploration despite the agent choosing greedily. Over time, the values converge on their actual estimates (insofar as the sequence of step size parameters fulfill the aforementioned properties).

This technique is called 'optimistic initial values'. Obviously, this is ineffective for non-stationary problems since the exploration happens temporarily at the onset of the algorithm, and will not renew with the reward values.

## 2.6 Upper Confidence Bound Action Selection

$\epsilon$ -greedy actions select exploratory actions randomly from non-greedy actions. Some notion of 'potential' should exist to denote that some actions are more



probably optimal than others. Sutton introduces the following in his book,

$$A_t = \arg \max_a [Q_t(a) + c \sqrt{\frac{\ln(t)}{N_t(a)}}]$$

$c > 0$  is denoted the 'degree of exploration' as if  $c$  is larger, then the term with the least number of attempts has larger value (notice that smaller values for  $N_t(a)$  increase the overall value).

U.C.B works well for stationary environments, but is hard for non-stationary ones since each action's reward is, in a sense, imprinted permanently by its past 'viability'.

## 2.7 A Gradient Bandit Algorithm

You can consider action-selection algorithms that use some point notion of preference  $H_t(a)$  to rank each action. Use a soft-max to decide an action

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

Action preferences are decided by stochastic gradient ascent.

$$(1) \quad H_{t+1}(A_t) = H_t(A_t) + \alpha(R_t - \hat{R}_t)(1 - \pi_t(A_t))$$

$$(2) \quad H_{t+1}(a) = H_t(a) - \alpha(R_t - \hat{R}_t)(\pi_t(a)) \quad \text{for all } a \neq A_t$$

Here  $\alpha > 0$  is the step-size parameter. Notice in (1) that, if the reward is higher, exploitation is incentivized. Notice that (1) is small when  $\pi_t$  is high, detracting from exploitation if one action is clearly better than all the others. In (2), the more an action is not selected, the less it will be selected in the future, as the  $\pi_t(a)$  term denotes.