

# Part I : Tabular Methods

→

s	✓
1	3.75
2	1.4
:	:

Ch. 2 : Bandits

+ states

Ch. 3 : MDPs

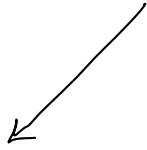
Ch. 4 : Dynamic Programming  
(Model - based)

Ch. 5 : Monte Carlo

Ch. 6 : Temporal Differences

Ch. 7 :  
N-step bootstrapping  
(Model - Free)

Ch. 9 :  
Learning and Planning

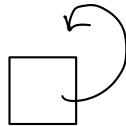


## 2.0 : Intro

- RL : differentiating factor is it uses training information that evaluates the actions taken, rather than instructs by giving correct actions.  $\hookleftarrow$  one state
- we study the nonassociative, evaluative feedback problem; k-armed bandit problem.

### Definitions

- Evaluative feedback : indicates how good the action taken was.
- Instructive feedback : gives the correct answer.
- Associative : the optimal action depends on the current state or situation,
- Non-associative : only 1 state/context to consider.  
Just action and rewards.  
This is the bandit problem!



1-state

## 2.1 : A k-armed bandit problem

- Learning problem where the agent chooses among  $k$  different options/actions (like arms in a slot machine),
- probability distribution of rewards depending on action,
- Objective : maximize expected total reward over some time period, say, 1000 time steps.

Formalization:

At time  $t$ :

Action  $A_t \in \{ \text{Arm 1, ..., } k \}$  // pick an arm

Reward  $R_t \in \text{Rewards of Bandit}_t$  // Reward distribution of the selected arm

Value of arbitrary action  $a = q_*(a)$  = the ground truth  
expected reward.

$$q_*(a) = E[R_t | A_t = a]$$

If we knew the reward of each action, then we would just choose the action with the highest value.

But we don't. We only have estimates.

Denote :

$Q_t(a)$  = estimated value of action  $a$  at time  $t$ .

would like  $Q_t(a)$  to be close to  $q_*(a)$ .

## 2.2 : Action - value methods

AVMs are methods for calculating  $Q_t(a)$  estimates. One way to do it is to average the rewards actually received:

Sample average :

$$Q_t(a) \doteq \frac{\text{sum of rewards when } a \text{ was taken at time } < t}{\# \text{ of times } a \text{ was taken at time } < t}$$

$$= \frac{\sum_{i=1}^{t-1} r_i \cdot 1_{A_i=a}}{\sum_{i=1}^{t-1} 1_{A_i=a}}$$

// sum of rewards at times when  $A_t = a$   
// count of times when  $A_t = a$ ,

By the law of large numbers  $Q_t(a)$  converges to  $q_*(a)$ ,

Now that we have an estimate for the rewards, how do we use it to select actions?

Greedy Action Selection (Definition) :

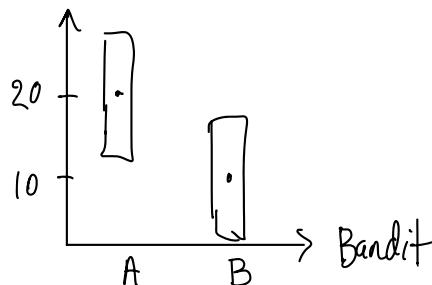
$$A_t = \underset{a}{\operatorname{argmax}} Q_t(a).$$

// pick the action  $a$  based on our current estimates of all the actions, such that action  $a$  is the one with the highest estimated reward.

$\epsilon$ -Greedy Action Selection : Behave greedily except for a small proportion  $\epsilon$ , during which the action is selected at random with equal probability (including the greedy action).

$$\epsilon \in [0, 1].$$

Let's say :

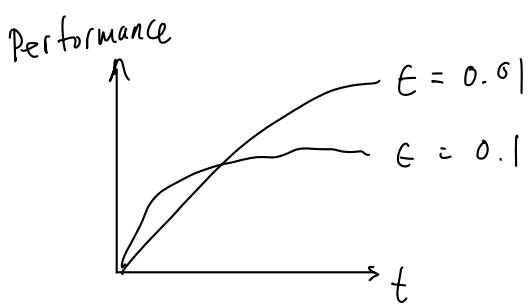


Init :  $A = B = 0$  // both are equal.

In  $\epsilon = 0$  : Let's say we pull B first, and got :

$B = 15 \rightarrow$  we are stuck w/ B since  $B > 0$  (initial value).

$\Rightarrow \epsilon = 0$  is Not a good strategy.



However,  $\epsilon$  too big is not good bc too much exploration:

### Mitigation Strategies

- Set  $\epsilon = 0$  after certain point.
- Decaying  $\epsilon$  over time.

2.3 : The 10-armed testbed :

TLDR ;  $\epsilon$ -greedy > greedy.

## 2.4: Incremental Implementation

How can the computation of the  $Q_t(a)$  estimates be done in a computationally effective manner?

For a single action, the estimate of this action value after it has been selected  $n-1$  times is:

$$Q_n \doteq \frac{R_1 + R_2 + \dots + R_{n-1}}{n-1} \quad // \text{average the rewards.}$$

We're not going to store all the values of  $R$  to recompute the sum at every time step, so we use the following update rule:

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

Reward  
at  $t=n$ 
Estimate  
of  $Q$  at  
 $t=n$

*// derivation in the book,*

The general form is:

$$\text{New Estimate} \leftarrow \text{Old Estimate} + \underbrace{\text{Step Size}}_{\text{learning rate}} [\underbrace{\text{Target} - \text{Old Estimate}}_{\text{error in the estimate.}}]$$

This applies to value function, action-value functions.

## A simple bandit algorithm

Initialize, for  $a = 1$  to  $k$ :

$$Q(a) \leftarrow 0$$

$$N(a) \leftarrow 0$$

Loop forever:

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

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

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

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

## 2.5 : Tracking a non-stationary problem

- Now the reward probabilities can change over time.
- We want to give more weight to the most recent rewards.

One easy way to do this is to use a constant step size so  $Q_{n+1}$  becomes a weighted average of past rewards and the initial estimate  $Q_n$ :

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

That is equivalent to:

$$Q_{n+1} = \underbrace{(1-\alpha)^n Q_1}_{< 1 \text{ initial estimate}} + \sum_{i=1}^n \alpha (1-\alpha)^{n-i} r_i \quad // \text{exponential recency-weighted average}$$

- sum of weights add up to 1
- weight given to  $r_i$  depends on how many rewards ago ( $n-i$ ) it was observed,
- since  $(1-\alpha) < 1$ , the weight given to  $r_i$  decreases as the number of intervening rewards increases.

criteria for  $\alpha$  (learning rate schedule) :

$$C1) \sum_{n=1}^{\infty} \alpha_n = \infty$$

$$C2) \sum_{n=1}^{\infty} \alpha_n^2 < \infty$$

Ex:  $\alpha_n = \frac{1}{n}$  : converges. // sample-average case.

C1) Harmonic series  $\rightarrow$  diverges.

C2) converges to  $\pi^2/6$ .

Ex:  $\alpha = \text{constant}$ .

$$C1) \infty \cdot C = \infty$$

$$C2) \infty \cdot C^2 = \infty$$

## 2.6: Optimistic initial values

- Setting optimistically high initial values is a way to encourage exploration, which is effective in stationary problems.
- Select initial estimates to a high number (higher than the reward we can actually get in the environment), so that unexplored states have a higher value than explored ones, and are selected more often.
- How to set the values? What is optimistic?

## 2.7: Upper - Confidence - Bound (UCB) Action Selection

- Another way to tackle the exploration problem.
- $\epsilon$ -greedy methods choose randomly the actions.
- maybe there is a better way to select those actions according to their potential of being optimal:

$$A_t \doteq \underset{\substack{\text{next} \\ \text{action} \\ \text{to choose}}}{\underset{\substack{\text{choose action} \\ \text{a that has} \\ \text{highest value}}}{\arg\max_a}} \left[ Q_t(a) + C \sqrt{\frac{\ln t}{N_t(a)}} \right]$$

↑  
Exploration bonus.  
measure of uncertainty or variance in estimate of  $a$ 's value.

Constant that controls how much to explore.

- $Q_t(a)$  : average observed reward for action  $a$ .
- $\ln t$  increases at each time step.
- $N_t(a)$  is the # of times action  $a$  has been selected prior to the time  $t$ , so it increases every time  $a$  is selected.
- $C > 0$ , is our confidence level that controls the degree of exploration.
  - $C = 0$  : pure exploitation
  - $C = \infty$  : pure exploration,

UCB balances exploitation (choosing actions with high average rewards) and exploration (trying actions we are uncertain about).

The quantity being max'ed over is thus a sort of upper bound on the possible true value of action  $a$ , with  $c$  determining the confidence level.

- Each time  $a$  is selected, the uncertainty is reduced :  $N_t(a)$  increments and uncertainty term decreases.
- Each time another action is selected,  $\ln t$  increases but not  $N_t(a)$ , so the uncertainty increases.

## 2.8 : Gradient Bandit Algorithms

- Another way to select actions.
  - So far we estimate action-values ( $\hat{q}$ -values) and use these to select actions.
  - Here we compute a numerical preference for each action  $a$ , which we denote  $H_t(a)$ :
- $H_t(a) \in \mathbb{R}$ .
- we select the action with a softmax.  $\pi_t(a)$ .

$$\Pr\{A_t = a\} \doteq \frac{e^{H_t(a)}}{\sum_{b=1}^k e^{H_t(b)}} \doteq \underbrace{\pi_t(a)}_{\substack{\text{probability of taking} \\ \text{action } a \text{ at time } t.}}$$

The algorithm is based on stochastic gradient ascent: at each step  $t$ , after selecting action  $A_t$  and receiving reward  $R_t$ , the action preferences are updated by:

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

and for all  $a \neq A$ :

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

•  $\alpha$  is the step size

•  $\bar{R}_t$  is the average of all rewards up through and including time  $t$ .

The  $\bar{R}_t$  term serves as a baseline with which the reward is compared. If the reward is higher than the baseline, then the probability of taking  $A_t$  in the future is increased, and if the reward is below the baseline, then the probability is decreased. The non-selected actions move in the opposite direction.

$$\frac{d E[R]}{d H(a)} = \sum_x q_*(x) \frac{d \pi(x)}{d H(a)} \quad // x = \{\text{set of possible actions}\}$$

sum over all possible actions : how good is action  $x$ ?      | how probability change?

SGD:

$$E[R|x] = q_*(x) \quad // \text{Expected reward of action } x = \text{true reward of action } x.$$

- 1) Variance of the gradient matters. Need a baseline.
- 2) What's the baseline? Why are we allowed to subtract it?

$$\sum_x \underbrace{\frac{d \pi(x)}{d H(a)}}_{\text{sum of gradients}} = 0 \quad . \quad \begin{aligned} \text{since we know that the} \\ \text{preferences are fixed sum,} \\ \text{increasing one decreases another.} \end{aligned}$$

$$\sum_x \pi(x) = 1 \quad // \text{Policy is a probability distribution.}$$

$$\sum [R - b] \frac{d\pi}{dH} \quad // \text{The influence of } b \text{ on expectation is 0.}\\ \text{Hence, there is no bias.}$$

3) Why do subtracting the baseline actually reduce the variance itself? See Below!

Ex: 2-arm bandits: Bandit 1 and 2.

Ground truth means are the same:

$$\mu(1) = \mu(2) = 9.$$

$$\text{Initialize: } H(1) = 1 \quad \pi(1) = 0.5 \quad p(a) = \frac{e^{H(a)}}{\sum_x e^{H(x)}} \\ H(2) = 1. \quad \pi(2) = 0.5$$

$$\text{Let } \alpha = 1, \beta = 9.$$

### No Baseline

$$\text{Sample } a=1, r=8.$$

$$H'(1) = H(1) + \alpha r (1 - \pi(1)) \\ = 1 + (8)(0.5) = 5$$

$$H'(2) = H(2) - \alpha r (\pi(2)) \\ = 1 - 8(0.5) = 3$$

$$\pi'(1) = e^5 / [e^5 + e^3] \\ = 0.9997$$

wild swing.

$$\pi'(2) = 1 - \pi'(1) = 0.003$$

### Baseline

$$\text{Sample } a=1, r-b = -1$$

$$H'(1) = 1 + (-1)(0.5) = 0.5$$

$$H'(2) = 1 - (-1)(0.5) = 1.5$$

$$\pi'(1) = 0.269$$

$$\pi'(2) = 0.731$$

wild swing is not as bad.

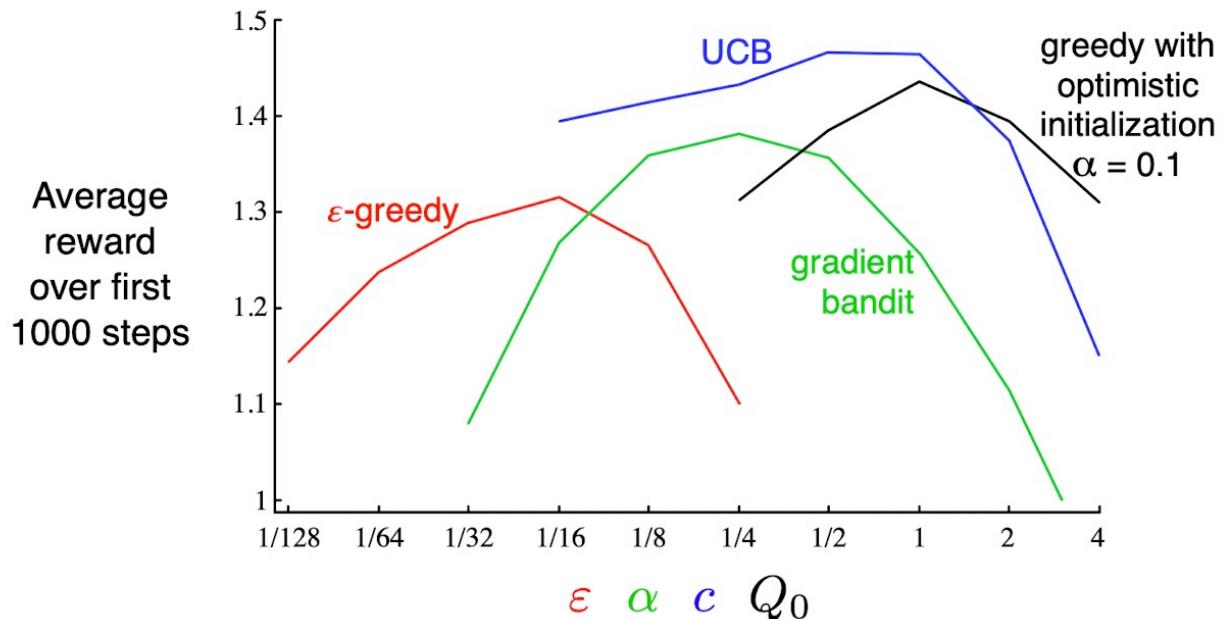


## 2.9 : Associative Search (contextual bandits)

- Nonassociative tasks : no need to associate different actions with different  $\underbrace{\text{situations}}_{\text{states}}$ .

## 2.10 : Summary

- One key topic is balancing exploitation vs exploration.
- Exploitation is straight forward : select action with the highest estimated value.
- Exploration is trickier, and we have seen several ways to deal with it:
  - $\epsilon$ -greedy choose randomly.
  - UCB favors the more uncertain actions.
  - gradient bandit estimate a preference instead of action values.



## 2.11 : Thomson Sampling (lecture)

All the explorations are costly exploration.

The exploration cost is known as regret.

Regret : expected loss after  $n$  actions due to the fact that you sometimes act suboptimally.

UCB1 : logarithmic regret,

Thompson Sampling : "

Does not extend to the full reinforcement learning problem.

Free exploration : up to a certain time  $n$ , explore the options.

- Stop pulling the bandit you know is the worse
- " the best .