

Multi-Arm Bandits - Book

RL is different from other methods as it evaluates actions taken rather than instructs by giving correct actions. → this creates a need for active exploration, trial and error search for good behaviour.

Evaluative feedback is the basis of methods for function optimisation.

[2.1] n-Armed Bandit Problem

You are faced repeatedly with a choice among n different options.

For each choice you receive a numerical reward chosen from a stationary probability distribution acc. to your action.

- Value: Each action has an expected or mean reward for that action to be selected.
- We don't know the values with certainty.
- If you maintain estimates then at any time step there is at least one action whose estimated value is greatest.

Selecting this action is greedy approach also called Exploitation
Best to maximise reward at one step

Selecting other actions is non-greedy also called exploration.
May pursue a greater total reward.

[2.2] Action Value Methods

$a \rightarrow \text{Action } a$

$q(a) \rightarrow \text{true (actual) value of action } a$

$Q_t(a) \rightarrow \text{estimated value of action } a \text{ at } t^{\text{th}} \text{ step}$

One way to estimate is by averaging the rewards actually received.

If ~~is~~ by t^{th} time step action a has been chosen $N_t(a)$ times prior to t , yielding rewards $R_1, R_2, \dots, R_{N_t(a)}$ then

$$Q_t(a) = \frac{R_1 + R_2 + \dots + R_{N_t(a)}}{N_t(a)}$$

If $N_t(a) = 0$, then we define $Q_t(a)$ some default value such as $Q_t(a) = 0$, as $N_t(a) \rightarrow \infty$ by law of large numbers, $Q_t(a)$ converges to $q(a)$.

sample-average method

The simplest action selection rule is to select the action with highest estimated action value i.e.

A_t^* for which

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

But this again turns to exploitation. A simple solution is to select non-greedily once in a while, say with small probability ϵ .

Instead of randomly selecting amongst all actions with equal probability \rightarrow independent of the action value estimates. We call methods using this near greedy action selection rule ϵ -greedy methods.

Advantage: as the no. of plays increases, every action will be sampled an infinite number of times, guaranteeing $N_t(a) \rightarrow \infty \forall a \in A$. thus ensuring $Q_t(a)$ converges to $q^*(a) \forall a \in A$.

Probability of selecting optimal action converges to $1 - \epsilon$.

[2.3] Incremental Implementation

In the previous we need to store all the rewards received

for recalculation which is inefficient and redundant.

We say $Q_k \rightarrow$ avg. of first $k-1$ rewards

$R_k \rightarrow$ reward for k^{th} step.

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

$$= \frac{1}{k} \left(R_k + \sum_{i=1}^{k-1} R_i \right)$$

$$= \frac{1}{k} \left(R_k + (k-1)Q_k + Q_k - Q_k \right)$$

$$= \frac{1}{k} \left(R_k + kQ_k - Q_k \right)$$

$$Q_{k+1} = Q_k + \frac{1}{k} [R_k - Q_k]$$

New Estimate \leftarrow Old estimate + Step size [Target - Old estimate]

error in the
estimate
can be denoted
by α \rightarrow it is a parameter

for sample - average method

$$\alpha = \frac{1}{K}$$

[2.4] Tracking Non-stationary Problem

The averaging methods discussed until now work for stationary environment.

In this case we must weigh the recent rewards more heavily. One way \rightarrow using a constant step-size parameter.

$$Q_{k+1} = Q_k + \alpha [R_k - Q_k]$$

where $\alpha \in (0, 1)$

This makes -

$$Q_K = \alpha R_K + (1-\alpha) Q_K$$

$$= \alpha R_K + (1-\alpha) [\alpha R_{K-1} + (1-\alpha) Q_{K-1}]$$

$$= \alpha R_K + \alpha (1-\alpha) R_{K-1} + (1-\alpha)^2 Q_{K-1}$$

$$= \alpha R_K + (1-\alpha) R_{K-1} \alpha + (1-\alpha)^2 R_{K-2} \alpha$$

$$+ \dots + (1-\alpha)^{K-1} \alpha R_{K-1}$$

$$+ (1-\alpha)^K Q_1$$

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

sometimes called
exponential, heavy-weighted average

Sometimes it's convenient to vary the step-size parameter from step to step.

$\alpha_{k(a)}$ → step-size parameter to process the reward received after k^{th} selection of a .

If $\alpha_{k(a)} = \frac{1}{k}$ → this becomes sample avg. method.

→ guaranteed to converge

conditions for assured convergence

$$\sum_{k=1}^{\infty} \alpha_{k(a)} = \infty \quad \text{and} \quad \sum_{k=1}^{\infty} \alpha_{k(a)}^2 < \infty$$

guarantees steps
are large enough
to eventually overcome
initial conditions or
random fluctuations.

eventually steps will
become small enough
to assure convergence.

Noting this we see a constant $\alpha_{k(a)} = \alpha$ doesn't meet the second condition and hence will not converge, thus will keep tracking the recent reward.

[2.5] optimistic initial values

All methods discussed yet are somehow dependant on initial value estimates i.e. $q_1(a)$ i.e. these methods are biased by initial estimates. For sample average method bias disappears when all actions are visited at least once but for constant $\alpha \rightarrow$ the bias is permanent though decreasing over time.

Initial action-values encourage exploration as they have to try all methods to find the best one.

Initially optimistic method performs worse, but eventually it performs better as exploration decreases with time. → This called optimistic initial values

Good for stationary problem but not well suited for non-stationary.

[2.6] Upper-Confidence-Bound Action Selection

ϵ -greedy selection forces non-greedy actions to be tried but indiscriminately with no-preference to those who are nearly-greedy actions according to their potential for being actually optimal. One way of doing that

$$A_t = \underset{a}{\operatorname{argmax}} \left[Q_t(a) + c \sqrt{\frac{\ln t}{N_t(a)}} \right]$$

$c > 0$ controls the initial uncertainty / degree of exploration

If $N_t(a) = 0$, then a is considered to be a maximising action.

The square root term is a measure of uncertainty or variance in the estimate of a 's value. The quantity being max'ed over is thus a sort of upper bound of the true value of action a .

[2.7] Gradient Bandits

Not in $H_t \leftarrow H_t(a) \rightarrow$ Numerical preference for terms of reward. \downarrow each action a .
 Larger the preference more often the action is taken

It is relative and is determined using soft-max distribution.

$$P_{t+1}\{A_{t+1} = a\} = \frac{e^{H_t(a)}}{\sum_{b=1}^n e^{H_t(b)}} = \pi_t(a)$$

Probability of taking action a at time t .

Initially all preferences are same : For eg: $H_1(a) = 0 \forall a$

Natural learning algorithm based on stochastic gradient ascent.

$A_t \rightarrow$ action taken on time step t

$R_t \rightarrow$ reward received

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

and

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

$\bar{R}_t \in R \rightarrow$ avg. of all rewards upto and including time t .
 serves as a

baseline. If reward is higher $\rightarrow P(A_t) \uparrow$
 else $P(A_t) \downarrow$

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

$$E[R_t] = \sum_b \pi_t(b) q_t(b)$$

$$\frac{\partial E[R_t]}{\partial H_t(a)} = \frac{\partial}{\partial H_t(a)} \left[\sum_b \pi_t(b) q_t(b) \right]$$

$$= \sum_b q_t(b) \frac{\partial \pi_t(b)}{\partial H_t(a)}$$

$$= \sum_b (q_t(b) - x_t) \frac{\partial \pi_t(b)}{\partial H_t(a)}$$

any scalar
that doesn't depend on b

\sum of this must be 0 as the sum of all probabilities will always remain 0.

$$= \sum_b \pi_t(b) (q_t(b) - x_t) \frac{\partial \pi_t(b)}{\partial H_t(a)} / \pi_t(b)$$

this can be written as

$$(-A_t + \pi_t - 1)(+\bar{q}_t - \bar{x}_t + A_t - H) = -H$$

$$(-A_t + \bar{q}_t)(\pi_t - \bar{x}_t) = E \left[(q_t(A_t) - x_t) \frac{\partial \pi_t(A_t)}{\partial H_t(a)} \mid \pi_t(A_t) \right]$$

$$= E \left[(R_t - \bar{R}_t) \frac{\partial \pi_t(A_t)}{\partial H_t(a)} \mid \pi_t(A_t) \right]$$

Shortly we establish that $\frac{\partial \pi_t(b)}{\partial H_t(a)} = \pi_t(b) / (I_{a=b} - \pi_t(a))$

Assuming this to be true

$$= E[(R_t - \bar{R}_t) I_{t(A_t)} (I_{a=A_t} - \pi_t(b)) | \mathcal{F}_t(A_t)]$$

$$= E[(R_t - \bar{R}_t) (I_{a=A_t} - \pi_t(b))]$$

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

Proving $\frac{\partial \pi_t(b)}{\partial H_t(a)} = \pi_t(b) / (I_{a=b} - \pi_t(a))$

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

$$\frac{\partial \pi_t(b)}{\partial H_t(a)} = \frac{\partial \left[\frac{e^{H_t(b)}}{\sum_i^n e^{H_t(c)}} \right]}{\partial H_t(a)}$$

$$= \frac{\partial H_t(b)}{\partial H_t(a)} \left[\sum_1^n e^{H_t(c)} \right] - \frac{e^{H_t(b)} \sum_1^n e^{H_t(c)}}{\partial H_t(a)}$$

$$= \frac{I_{a=b} e^{H_t(a)} \sum_1^n e^{H_t(c)}}{\left(\sum_1^n e^{H_t(c)} \right)^2} - \frac{e^{H_t(b)} e^{H_t(a)}}{\left(\sum_1^n e^{H_t(c)} \right)^2}$$

$$= \frac{I_{a=b} e^{H_t(b)}}{\sum_1^n e^{H_t(c)}} - \frac{e^{H_t(b)} e^{H_t(a)}}{\left(\sum_1^n e^{H_t(c)} \right)^2}$$

$$= \boxed{I_{a=b} \pi_t(b) - \pi_t(b) I_{a=b}}$$

$$= \pi_t(b) [I_{a=b} - \pi_t(a)]$$

$$= \pi_t(b) [I_{a=b} - \pi_t(a)]$$

$$\frac{(a+b)(a-b)}{(b-a)} = \frac{(a+b)(a-b)}{(a-b)}$$

$$\frac{(a+b)(a-b)}{(a-b)} = \frac{(a+b)(a-b)}{(a-b)}$$

$$\frac{(a+b)(a-b)}{(a-b)} = (a+b)$$

$$\frac{(a+b)(a-b)}{(a-b)} = \frac{(a+b)(a-b)}{(a-b)}$$

$$(a+b)(a-b) = (a+b)(a-b)$$

$$(a+b)(a-b) = (a+b)(a-b)$$