Date created: 2025-07-16-Wed

# 2.4 Incremental Implementation

We estimate action values using sample averages, but as the number of samples grow, storing these sample values and computing these averages can become expensive. How can we make this computationally efficient? We are looking for constant memory and constant per-time-step computation.

## Notation

For simplicity we consider the case with only one action.

- $R_i$: reward received after the $i$th iteration of *this* action
- $Q_n$: estimate of its action value after it has been selected $n - 1$ times.

It is easy to see that

$$Q_n := \frac{R_1 + R_2 + ... + R_{n-1}}{n - 1}$$

## Incremental formula

Given $Q_n$ and the $n$th reward, $R_n$, the new average of all $n$ rewards can be computed by

$$
\begin{aligned}
Q_{n+1} &= \frac{1}{n} \sum_{i=1}^{n} R_i \\
&= \frac{1}{n} \left( R_n + \sum_{i=1}^{n-1} R_i \right) \\
&= \frac{1}{n} \left( R_n + \frac{n-1}{n-1} \sum_{i=1}^{n-1} R_i \right) \\
&= \frac{1}{n} \left( R_n + (n-1)Q_n \right) \\
&= \frac{1}{n} \left( R_n + nQ_n - Q_n \right) \\
&= Q_n + \frac{1}{n} \left[ R_n - Q_n \right]
\end{aligned}
$$

## General form

This update rule is of a form that occurs frequently in RL:

$$\text{NewEstimate} \leftarrow \text{OldEstimate} + \text{StepSize} \left[ \text{Taget} - \text{OldEstimate} \right]$$

Where $[Target - OldEstimate]$ is an *error* in the estimate, which we reduce by taking a step towards the Target. The Target is a desirable direction in which to move, which in our example above is the $n$th reward. Obviously, the Target may be noisy.

Note that StepSize can be variable. In this book we denote this by $\alpha_t(a)$.

# A simple bandit algorithm

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

$\qquad Q(a) \leftarrow 0$

$\qquad N(a) \leftarrow 0$

*Loop* :

$$A \leftarrow \begin{cases} \text{argmax}_a Q(a) & \text{with probability } 1 - \epsilon, \text{(ties broken randomly)} \\ \text{a random action} & \text{with probability } \epsilon \end{cases}$$

$\qquad R \leftarrow bandit(A)$

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

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