

Chapter 1

n -Armed bandits

1.1 Notes

1.1.1 n -Armed Bandit Problem

We have n different options (actions) representing n different slot machines. Each action has a given reward, sampled from a stationary probability $q(a)$ only dependent on the chosen action a . We want to maximize the (expected) total reward over a given (large) time T : $\sum_{t=1}^T R_t$. To do that, we estimate the value $Q_t(a)$ of each action given what we have seen so far. Let R_t the reward at time t and $N_t(a)$ the number of times the action a has been chosen so far.

1.1.2 Estimating value

We estimate the value with:

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

with $Q_t(a) = Q_1(a)$ a default value. With $N_t(a) \rightarrow \infty$ we have $Q_t(a) \rightarrow q(a)$.

Step-by-step, this can be calculated using incremental implementation to save computation time:

$$\begin{aligned} Q_{k+1} &= \frac{1}{k} \sum_{i=1}^k R_i \\ &\dots \\ Q_{k+1} &= Q_k + \frac{1}{k} (R_k - Q_k) \end{aligned}$$

This looks like $NewEstimate \leftarrow OldEstimate + StepSize (Target - OldEstimate)$, with $StepSize = \frac{1}{k}$ here.

For tasks that never stop this estimation diverges, plus we may be interested in tracking a nonstationary problem. To achieve this, we can introduce a constant step size, that effectively weights recent rewards more heavily:

$$\begin{aligned} Q_{k+1} &= Q_k + \alpha (R_k - Q_k) \\ &\dots \\ Q_{k+1} &= (1 - \alpha)^k Q_1 + \alpha \sum_{i=1}^k (1 - \alpha)^{k-i} R_i \end{aligned}$$

As it turns out, this defines a weighted average with weights $(1 - \alpha)^k, \alpha(1 - \alpha)^{k-1}, \dots, \alpha(1 - \alpha)^0$ (they sum to 1).

By denoting $\alpha_k(a)$ the weight (step-size) used for the k -th selection of action a , we need to have two conditions:

1. $\sum_{k=1}^{\infty} \alpha_k(a) = \infty$, to guarantee that we overcome initial estimate, and
2. $\sum_{k=1}^{\infty} \alpha_k^2(a) < \infty$, to guarantee convergence.

Choosing actions

To choose the action, the *greedy* way is to select the one with the highest value: $A_t = \operatorname{argmax}_a Q_t(a)$. Problem: this does not spend any time to sample other actions to refine the estimates $Q_t(a)$.

First solution: ϵ -greedy algorithms, where $A_t = \operatorname{argmax}_a Q_t(a)$ $1 - \epsilon$ of the times and $A_t = \operatorname{uniform}(a)$ the other ϵ of the times.

Second solution: optimistic initial values, to preferentially select unsampled actions.

Third solution: *Upper Confidence Bound (UCB)* action selection, with

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

Rationale: the square-root term is a measure of the uncertainty (or variance) in the estimate of the value of a . The higher this term (and c), the higher the chance the action will be taken instead of the optimal action. UCB is often hard to transpose outside of the n -armed bandit problem.

Gradient bandits

Instead of estimating the value of each action, we can estimate the relative preference $H_t(a)$ of one action over others. We then compute the probability of taking action a using the softmax distribution (softmax "normalizes" its inputs so that any constant added to all preferences has no effect on the probability):

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

Initially, $H_1(a) = 0$, so every action has the same probability to be chosen. We then use a variation on stochastic gradient ascent to update the preference after each step (A_t is the action taken at step t and R_t the corresponding reward):

$$\begin{aligned} H_{t+1}(A_t) &= H_t(A_t) + \alpha(R_t - \bar{R}_t)(1 - \pi_t(A_t)) && \text{and} \\ H_{t+1}(a) &= H_t(a) + \alpha(R_t - \bar{R}_t)\pi_t(A_t) && \forall a \neq A_t \end{aligned}$$

Explanation: if the reward (R_t) is better than the average reward observed thus far (\bar{R}_t), we increase the probability to select A_t and decrease the probabilities of all other actions in proportion to that difference ($R_t - \bar{R}_t$). On the contrary, if the reward is lower than the current average reward, we decrease $\pi_t(A_t)$ and increase all others. We increase/decrease more if the selected action had lower probability (thus the $1 - \pi_t(A_t)$ term).