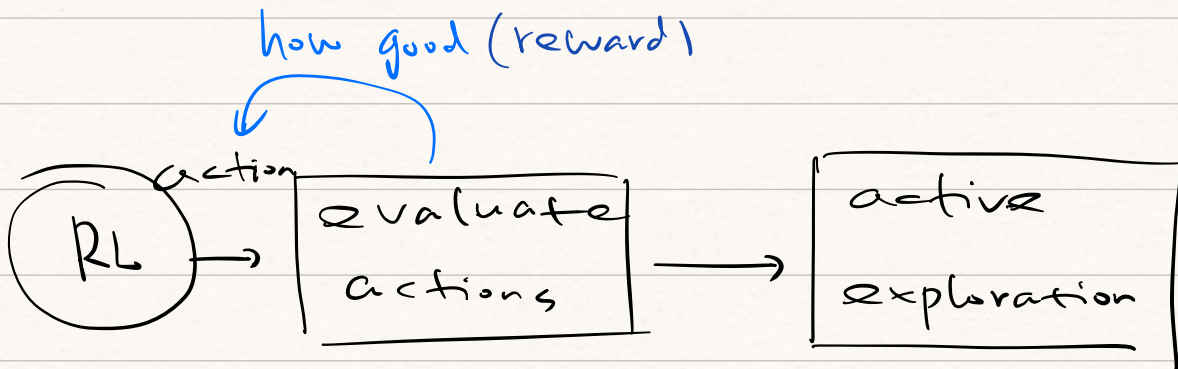


RL:

evaluative feedback: depends on action taken
(how good action)

Others:

instructive feedback: independent of the action taken
(correct action)



k-armed bandit problem

k different options (actions)

action selections (time steps)

Value (of the action): expected/mean reward
given that action

$$Q_*(a) = \mathbb{E}[R_t | A_t = a]$$

$$Q_t(a) \xrightarrow[\text{greedy actions}]{\frac{1}{2} \cdot \frac{1}{2}} Q_*(a)$$

current knowledge

Select one action: exploit current knowledge
of the values of the actions
(on one step)

Select non-greedy action: exploring

improve your estimates of
the nongreedy action's
value.

(in the long run)

Action-Value Method

Value:

$$Q_t(a) \doteq \frac{\text{sum of rewards when } a \text{ taken prior to } t}{\text{number of times } a \text{ taken prior to } t} = \frac{\sum_{i=1}^{t-1} R_i \cdot \mathbb{1}_{A_i=a}}{\sum_{i=1}^{t-1} \mathbb{1}_{A_i=a}}$$

$$t \rightarrow \infty, \quad \sum_{i=1}^{t-1} \mathbb{1}_{A_i=a} \rightarrow \infty, \quad Q_t(a) \rightarrow Q_*(a)$$

(Sample-average method)

Greedy action selection method:

$$A_t \doteq \operatorname{argmax}_a Q_t(a)$$

exploit current knowledge to maximize immediate

reward

small probability ϵ ; (ϵ -greedy)

Select randomly from among all the actions
with equal probability, independently of
the action-value estimates.

10-armed Testbed: $\epsilon = 0.01$ slower but better

"Get stuck performing suboptimal actions"

Incremental Implementation

Action value:

$$Q_n = \frac{R_1 + R_2 + \dots + R_n}{n} = \frac{1}{n} \sum_{i=1}^n R_i$$

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

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

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

$$= \frac{1}{n} (R_n + nQ_n - Q_n)$$

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

$$\text{New Estimate} \leftarrow \text{old Estimate} + \text{Step Size} \times [\text{Target} - \text{old Estimate}]$$

Error

Non-stationary: (Exponential Recency-weighted Average)

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

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

Conditions:

$$\left\{ \begin{array}{l} \sum_{n=1}^{\infty} \alpha_n(a) = \infty \\ \sum_{n=1}^{\infty} \alpha_n^2(a) < \infty \end{array} \right.$$

Upper - confidence - bound Action Selection

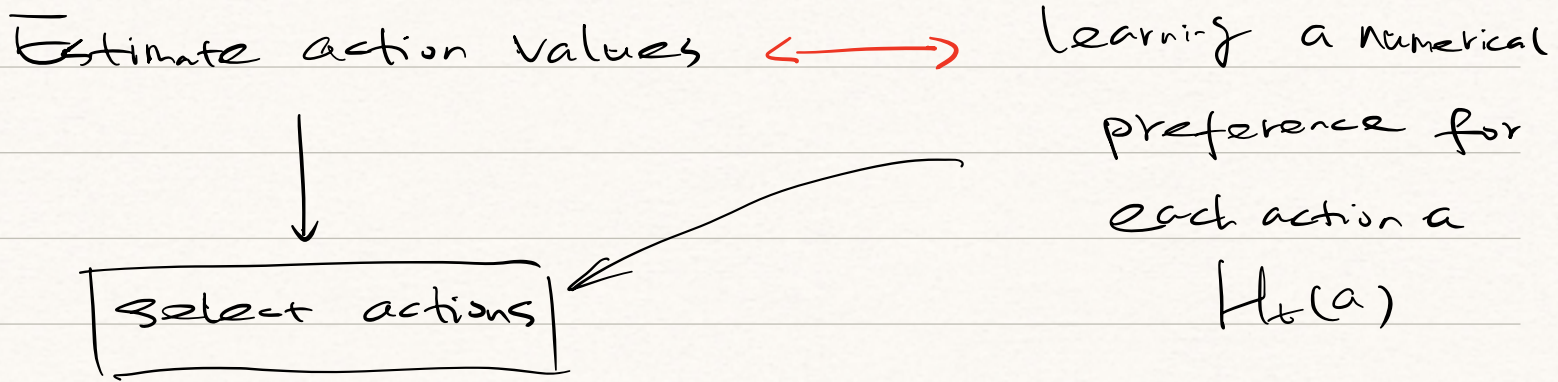
(UCB)

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

$N_t(a)$: number of times that action a has been selected prior to time t

$c > 0$: control the degree of exploration

Gradient Bandit Algorithm



Soft-max distribution (Gibbs / Boltzmann distribution)

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

$$\begin{cases} H_{t+1}(a) \doteq H_t(a) + \alpha \frac{\partial \mathbb{E}[R_t]}{\partial H_t(a)} \\ \mathbb{E}[R_t] = \sum_x \pi_t(x) q_*(x) \end{cases}$$

Associative Search

(Contextual Bandits)

Learn a policy : a mapping from situations to actions
that are best in those situations

trial-and-error learning to search for the best actions

association of these actions with the situations
in which they are best

parameter study :

perform best

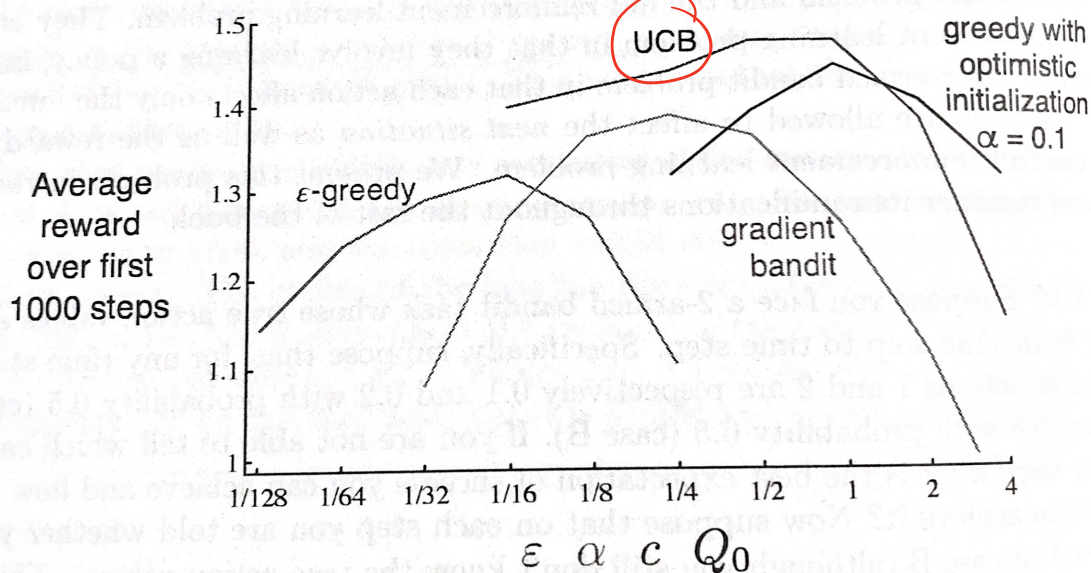


Figure 2.6: A parameter study of the various bandit algorithms presented in this chapter. Each point is the average reward obtained over 1000 steps with a particular algorithm at a particular setting of its parameter.