

A simple bandit algorithm

Initialize, for $a = 1$ to k :

$$Q(a) \leftarrow 0$$
$$N(a) \leftarrow 0$$

Loop forever:

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

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

$$\underline{N(A) \leftarrow N(A) + 1}$$

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

Lecture -02

Immediate RL / Bandits

$$q_x(1) \quad q_x(2) \quad \dots \quad q_x(k)$$

$\boxed{1} \quad \boxed{2} \quad \dots \quad \boxed{k}$

$$Q_t(1) \quad Q_t(2) \quad \dots \quad Q_t(k)$$

$Q_t(a) = \frac{\text{Sum of rewards when } a \text{ taken proto}}{\text{# of times 'a' taken proto}}$

$$= \frac{\sum_{i=1}^{t-1} R_i \frac{1}{A_i = a}}{\sum_{i=1}^t \frac{1}{A_i = a}}$$

Greedy $A_t = \underset{a}{\operatorname{argmax}} Q_t(a)$
 exploratory action for 'E' times
 You pick a random action.

Any action a

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

Naive implementation: Store all rewards

& Compute average every time.

Incremental implementation: -

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

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

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

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

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

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

$$Q_{n+1} = \underline{Q_n} + \frac{1}{n} \boxed{R_n - \underline{Q_n}}$$

↑
new estimate
↑
old estimate
↑
Current reward

$$\text{for } n=1, \quad Q_2 = R_1$$

$$\text{New estimate} = \text{old estimate} + \text{stepsize}$$

$$[\text{target} - \text{old estmat}]$$

Non-stationary problems:

Stationary bandits: reward prob. do not change over time.

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



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

$$\alpha \in [0, 1]$$

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

$$= \alpha R_n + (1-\alpha) Q_n$$

$$= \alpha R_n + (1-\alpha) [\alpha R_{n-1} + (1-\alpha) Q_{n-1}]$$

$$= \alpha R_n + (1-\alpha) \alpha R_{n-1} + (1-\alpha)^2 Q_{n-1}$$

$$\begin{aligned}
 &= \alpha R_n + (1-\alpha) \alpha R_{n-1} + (1-\alpha)^2 \alpha R_{n-2} + \\
 &\quad \cdots + (1-\alpha)^{n-1} \alpha R_1 + \\
 &\quad (1-\alpha)^n Q_1
 \end{aligned}$$

$$(1-\alpha)^n + \sum_{j=1}^n \alpha (1-\alpha)^{n-j-1} = 1$$

$\alpha_n(a)$ - Step size param. used
to process the reward
received after the n^{th}
selection of a .

$$\alpha_n(a) = \frac{1}{n}$$

✓ $\sum_{n=1}^{\infty} \alpha_n(a) = \infty$

✓ $\sum_{n=1}^{\infty} \alpha_n^2(a) < \infty$.

α

Initial values :-

bias

$$Q_1(1) = 0$$

$$Q_1(2) = 0 \dots Q_1(k) = 0$$

$$Q_2 = R_1$$

:

:

:

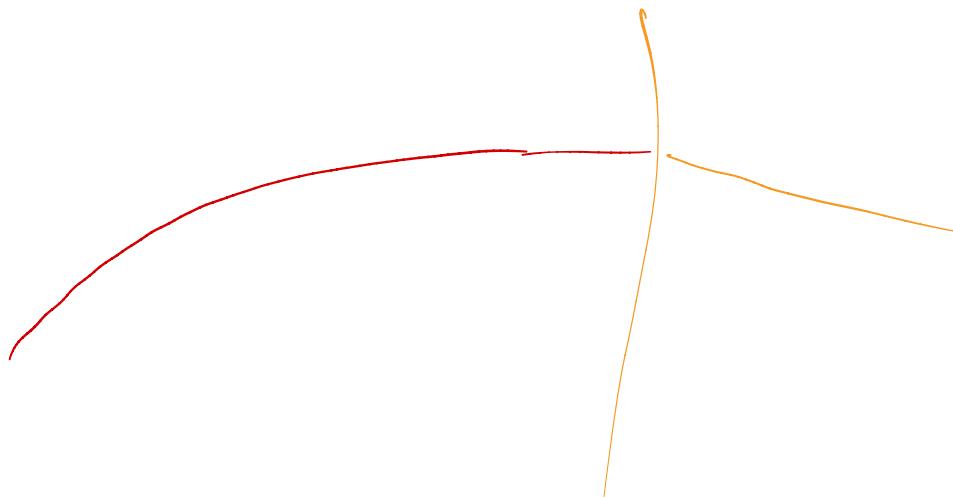
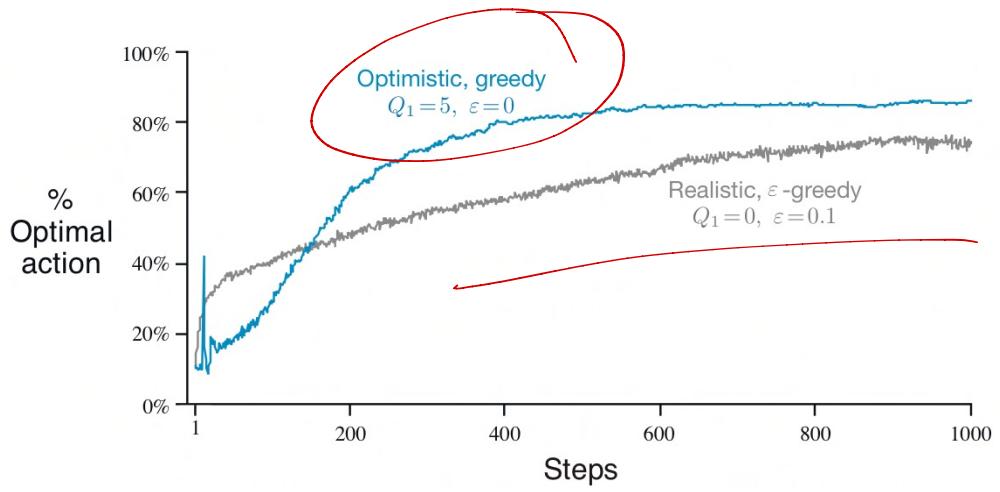
$$Q_1(1) = 5$$

$$Q_1(2) = 5 \dots Q_1(k) = 5$$

1

2

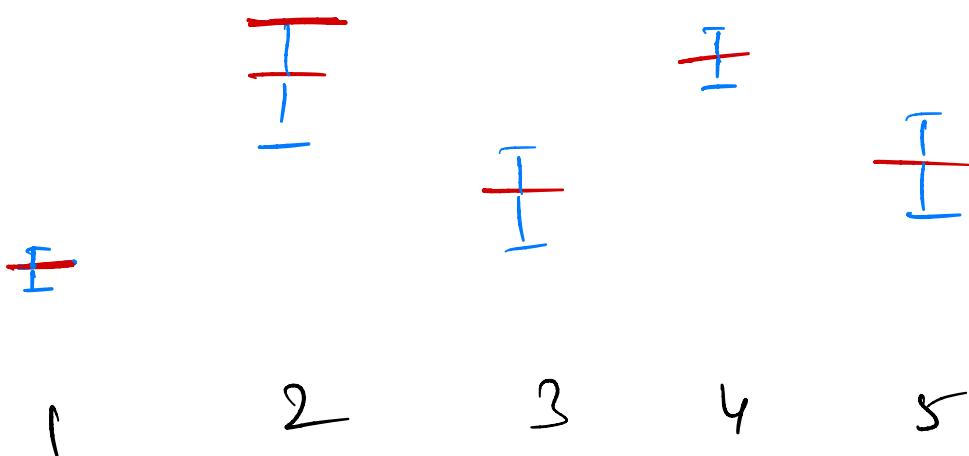
1,5



Upper-Confidence bound (UCB) action

Selection

- There are other promising actions.
- there are other actions with high uncertainty.



$$A_L = \underset{a}{\operatorname{argmax}} \left[Q_L(a) + C \sqrt{\frac{\ln t}{N_L(a)}} \right]$$

$C > 0$ Controls the degree of exploration

- ✓ E-Greedy
 - ✓ UCB
- }
- ~~With~~ action-value estimation.
-

Gradient bandit algorithm ;—

$H_t(a)$ — numerical preference for choosing action a

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

Gibbs or Boltzmann distribution.

$\pi_t(a)$ = prob. of taking action a .

$H_t(a) = 0$ for all a .

How to learn the preferences?

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

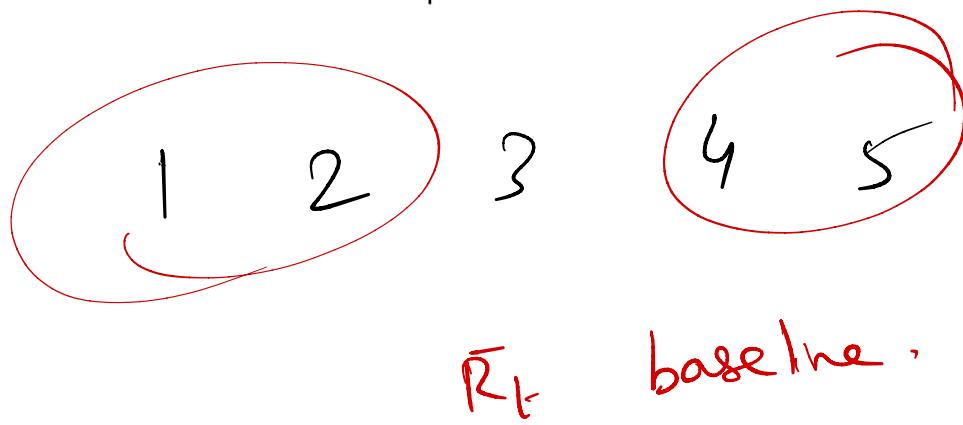
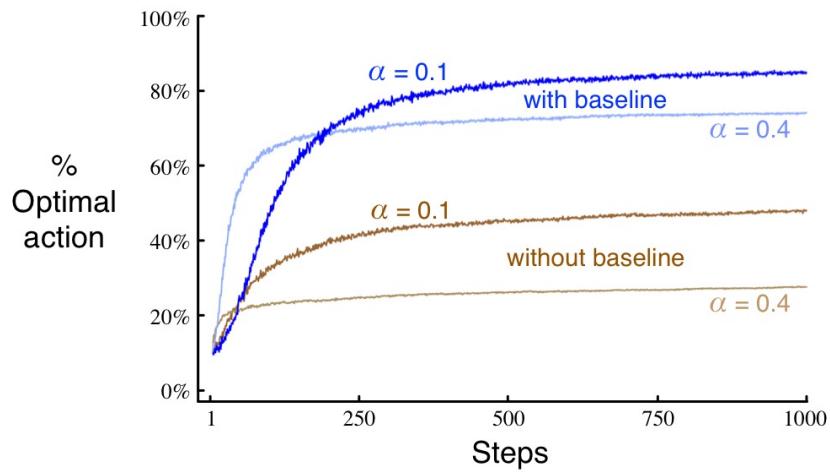
$\overbrace{\qquad\qquad\qquad}^{\text{step size } \alpha > 0}$

$\bar{R}_t \in \mathbb{R}$ is the average of all rewards up to t .

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

$\forall a \notin A_t$

mean = + 4



The bandit gradient algorithm as

Stochastic gradient ascent:

Exact gradient ascent:

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

$$\text{Performance } \Phi = E[R_t]$$

$$= \sum_n \pi_t(n) q_{t+}(n)$$

$$\frac{\partial E[R_t]}{\partial H_t(a)} = \frac{\partial}{\partial H_t(a)} \left[\sum_n \pi_t(n) q_{t+}(n) \right]$$

$$= \sum_n q_{t+}(n) \frac{\partial \pi_t(n)}{\partial H_t(a)}$$

$$\sum_n \frac{\partial \pi_t(n)}{\partial H_t(a)} = 0$$

$$= \sum_n (q_{t+}(n) - \beta_t) \frac{\partial \pi_t(n)}{\partial H_t(a)}$$

$$\frac{\partial E[R_t]}{\partial H_t(a)} = \sum_n (q_x(n) - \beta_t) \frac{\partial \pi_t(n)}{\partial H_t(a)}$$

$$= \sum_n \pi_t(n) (q_x(n) - \beta_t) \frac{\partial \pi_t(n)}{\partial H_t(a) / \pi_t(a)}$$

$$= E \left[(q_x(n) - \beta_t) \frac{\partial \pi_t(A_t)}{\partial H_t(a) / \pi_t(A_t)} \right]$$

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

Note: $E[R_t | A_t] = q_x(A_t)$

$$\frac{\partial \pi_t(a)}{\partial h_t(a)} = \frac{\partial}{\partial h_t(a)} \pi_t(a)$$

$$= \frac{\partial}{\partial h_t(a)} \left[\frac{e^{H_t(n)}}{\sum_{y=1}^k e^{H_t(y)}} \right]$$

$\frac{\partial f(y)}{\sqrt{u_n - u_d}} \frac{u_d - u}{\sqrt{2}}$

$$= \frac{\sum_{y=1}^k e^{H_t(y)} \frac{\partial e^{H_t(n)}}{\partial h_t(a)} - e^{H_t(n)} \frac{\partial}{\partial h_t(a)} \sum_{y=1}^k e^{H_t(y)}}{\left(\sum_{y=1}^k e^{H_t(y)} \right)^2}$$

$$= \frac{\mathbb{1}_{a=n} e^{H_t(n)} \sum_{y=1}^k e^{H_t(y)} - e^{H_t(n)} e^{H_t(a)}}{\left(\sum_{y=1}^k e^{H_t(y)} \right)^2}$$

$$= \frac{\mathbb{1}_{a=n} e^{H_t(n)}}{\sum_{y=1}^k e^{H_t(y)}} - \frac{e^{H_t(n)} e^{H_t(a)}}{\left(\sum_{y=1}^k e^{H_t(y)} \right)^2}$$

$$= \frac{\frac{1}{a=n} e^{H_t(n)}}{\sum_{y=1}^k e^{H_t(y)}} - \frac{e^{H_t(n)} e^{H_t(a)}}{\sum_{y=1}^k e^{H_t(y)}}^2$$

$$= \frac{1}{a=n} \pi_t(n) - \pi_t(n) \cdot \pi_t(a)$$

$$\frac{\partial \pi_t(n)}{\partial H_t(a)} = \pi_t(n) \left(\frac{1}{a=n} - \pi_t(a) \right)$$

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

$$= E \left[(R_t - \bar{R}_t) \pi_t(A_t) \left(\frac{1}{a=A_t} - \pi_t(a) \right) \right]$$

$$= E \left[(R_t - \bar{R}_t) \left(\frac{1}{a=A_t} - \pi_t(a) \right) \right]$$

$$H_{tf}(a) = H_t(c) + \alpha (R_t - \bar{R}_t) \left(\frac{1}{a=a_t} - \pi_t(a) \right)$$

Contextual bandits :-

Bandits → Non - aggressive tasks.

Full RL — many situations.

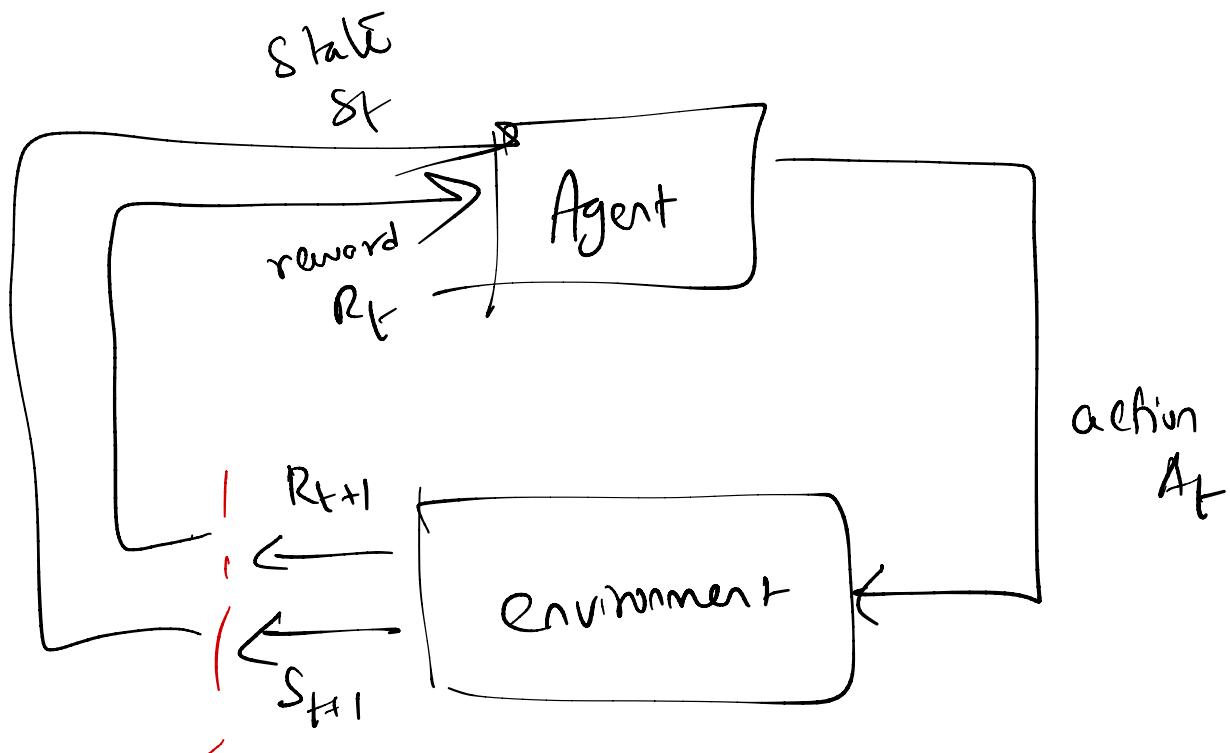
Contextual bandits

Context or id

Ad-placement:

Context: User history,

Finite Markov Decision Processes (MDP)



Sequence / trajectory: $s_0, a_0, r_1, s_1, a_1, r_2, s_2, a_2, r_3 \dots$

~~Finite MDP~~
 (S, A, R)

R_t, S_t — well defined discrete
 prob. distrib. that
 depend only on prev.
 state & action

$$P(s', r | s, a) = \Pr_{\tau} \left[\begin{array}{l} S_t = s', R_t = r \\ S_{t-1} = s, A_{t-1} = a \end{array} \right]$$

dynamics
of the MDP.

$$\forall s', s \in S, r \in R, \\ a \in A(s)$$

Note: $\sum_{s' \in S} \sum_{r \in R} P(s', r | s, a) = 1$

$$\forall s = S, a \in A(s)$$

$$p: S \times R \times S \times A \rightarrow [0, 1]$$

State-transition prob.

$$P(s' | s, a) = \sum_{r \in R} P(s', r | s, a)$$

Expected reward for s, a pair:

$$r(s, a) = \sum_{r \in R} \sum_{s' \in S} P(s', r | s, a)$$

Markov property:

$$P(S_t, R_t \mid S_{t-1}, A_{t-1})$$

$$= P(S_t, R_t \mid S_{t-1}, A_{t-1}, R_{t-1}, S_{t-2}, \\ A_{t-2}, \dots)$$

boundary between Agent & env:

The boundary represents the limit
of agent's absolute control
Not its knowledge.

MDP — ~~one step~~ the actions ✓

— the state ✓

— the reward, ✓

Goals & rewards:-

↓
What to achieve

Not how to achieve it.

Returns

Cumulative reward = return

$$G_t = R_{t+1} + R_{t+2} + \dots + R_T$$

fixed
time step

Episodes

episodic tasks

Tasks

continuing tasks

$\gamma = 0$
"myopic"

$\gamma = 1$

discount rate

$0 \leq \gamma \leq 1$

$$\begin{aligned} G_t &= R_{t+1} + \gamma R_{t+2} + \gamma^2 R_{t+3} \dots \\ &= \sum_{k=0}^{\infty} \gamma^k R_{t+k+1} \end{aligned}$$