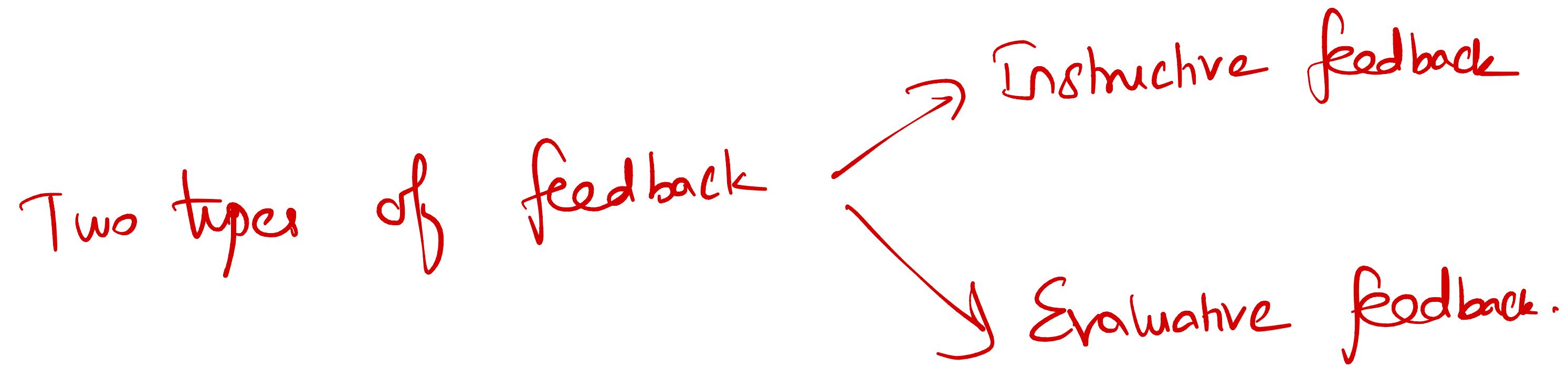


## Lecture -02

### Immediate Reinforcement Learning:-

#### Key Characteristics of an RL Problem:-

- ① Learn to act in many situations. X
- ② Delayed rewards and credit assignment. X
- ③ Exploration / Exploitation dilemma.



## Instructive feedback

- instruct the right action  $a^*$ .
- ignores the action taken.
- SL

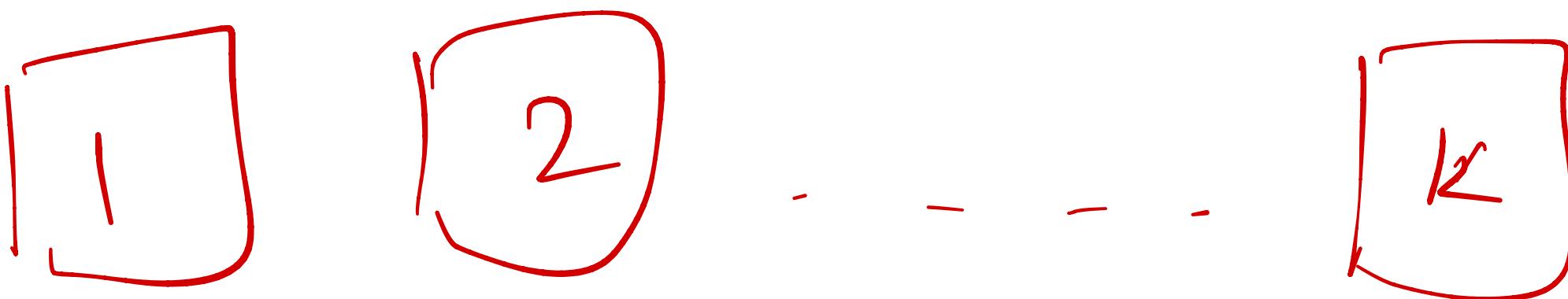
## Evaluative feedback

- evaluates the action taken At by giving some reward.
- Completely depends on the At.
- RL  
“Exploration Vs. Exploitation”

Multi-armed bandits:-

{ - Same state for all 't'.

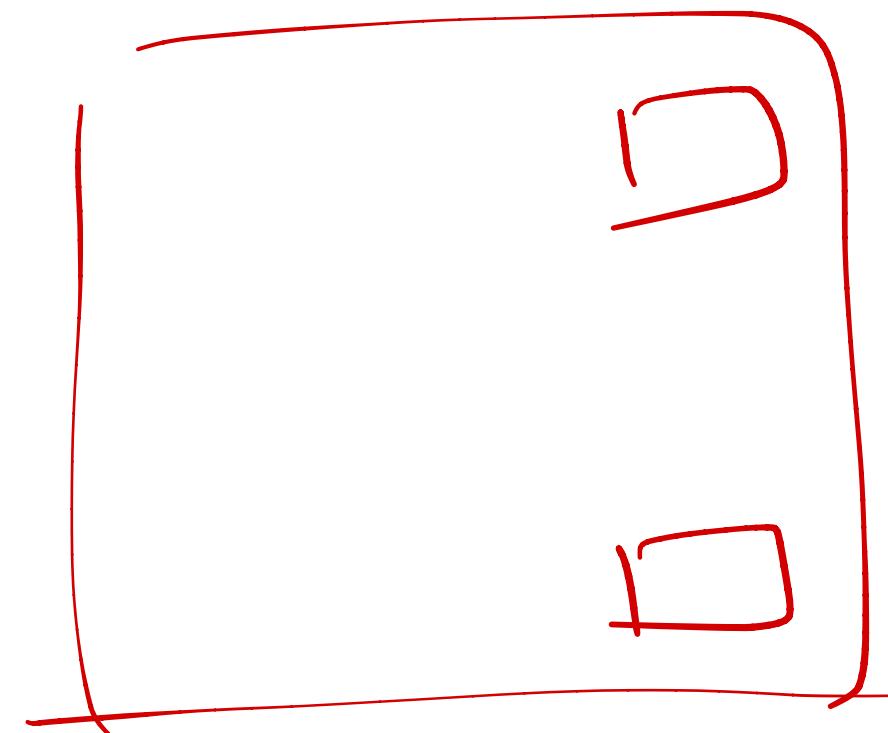
- rewards are  
immediate



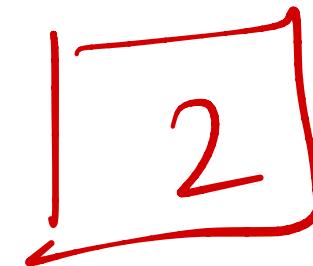
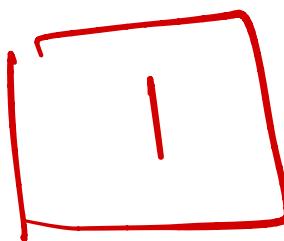
"T" "Obj!" to maximize - the expected  
total reward over some  
time period T.

## Applications:-

- Medical trials.
- A/B testing
- Ads. [contented badils]  
↑



## Multi-armed bandits :-



A<sub>1</sub> : 1

A<sub>2</sub> : 5

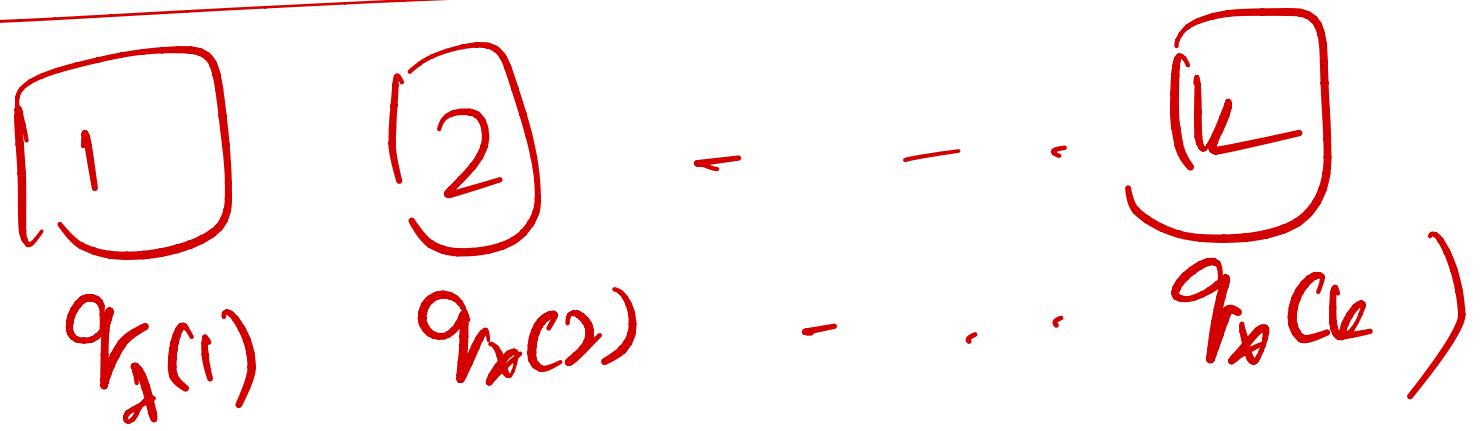
A<sub>3</sub> : 0

A<sub>4</sub> : 0

A<sub>5</sub> : 10

A<sub>6</sub> : 1

- Reward distributions are stationary.



Value of Action  $a = q_{r_k}(a) = \mathbb{E}[R_t | A_t = a]$

Optimal action  $a^* = \operatorname{argmax}_{a \in \{a_1, a_2, \dots, a_k\}} q_{r_k}(a)$

$Q_t(a)$  = estimated value of action ' $a'$   
at time ' $t$ '.

$q_{\pi}(a)$  = true expected reward.

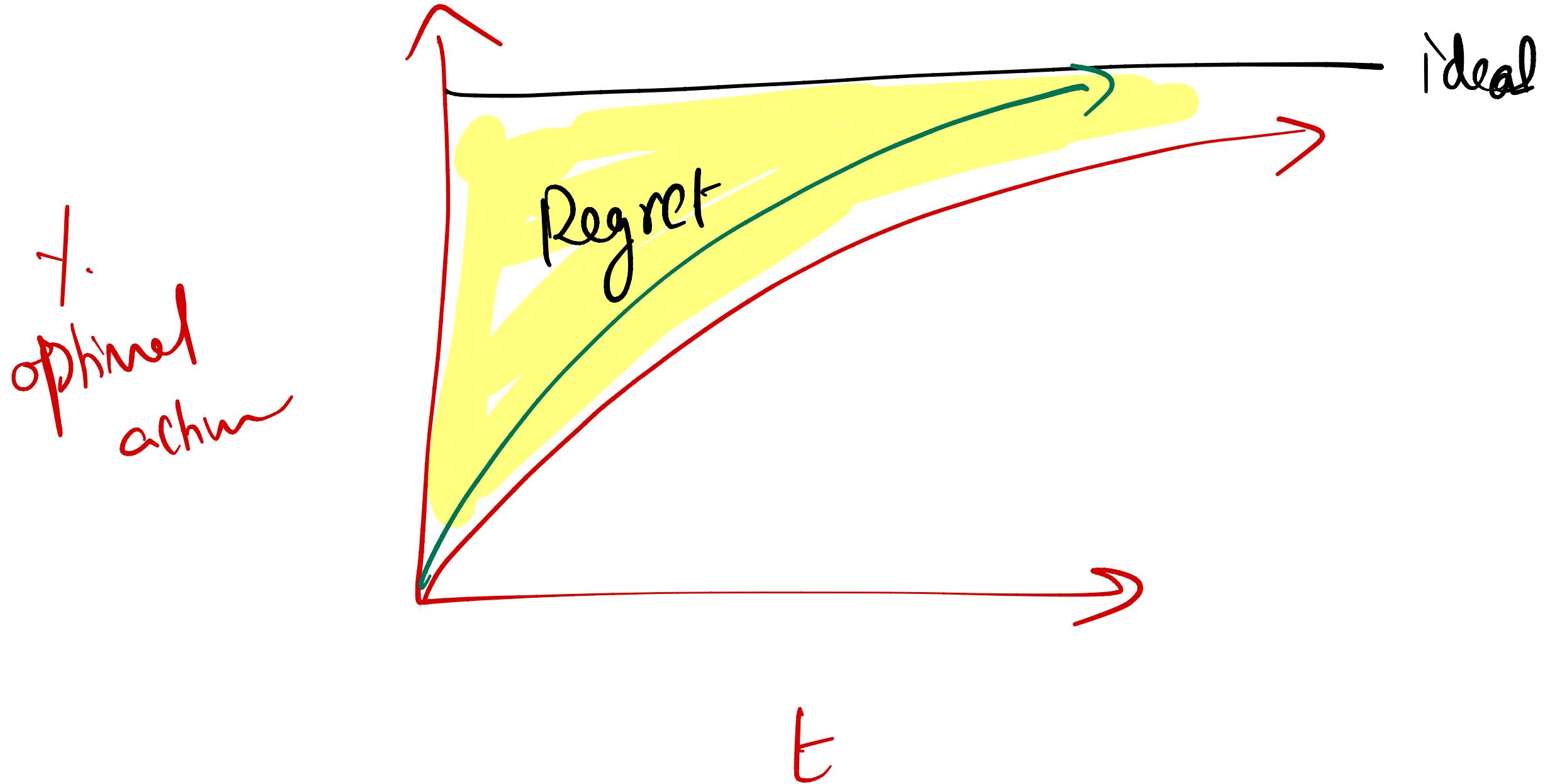
---

Regret after ' $K$ ' steps!

$$\text{regret} = k q_{\pi}(a^*) - \underbrace{\sum_{t=1}^K R_t}_{\text{To}}$$

100

To



## Action-value methods:-

$Q_T(a)$  = Sum of rewards when 'a' was taken prior to  $t'$ .

$$1_{A_i=a} = \begin{cases} 1 & \text{if 'a' was taken at } i^{\text{th}} \text{ step} \\ 0 & \text{o.w.} \end{cases}$$

$\# \text{ of times 'a' was taken prior to } t'$ .

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

$Q_T(a) \rightarrow q_{T,x}(a)$   
as  $t \rightarrow \infty$ .

## Alg 1: Explore then Commit (ETC) :

"m"

for each arm

Pull that arm "m" times

$Q(\cdot)$  = Average of "m" rewards.

best action =  $\arg\max_a Q(a)$

pull the best arm for the rest of the trials -

Alg 2:

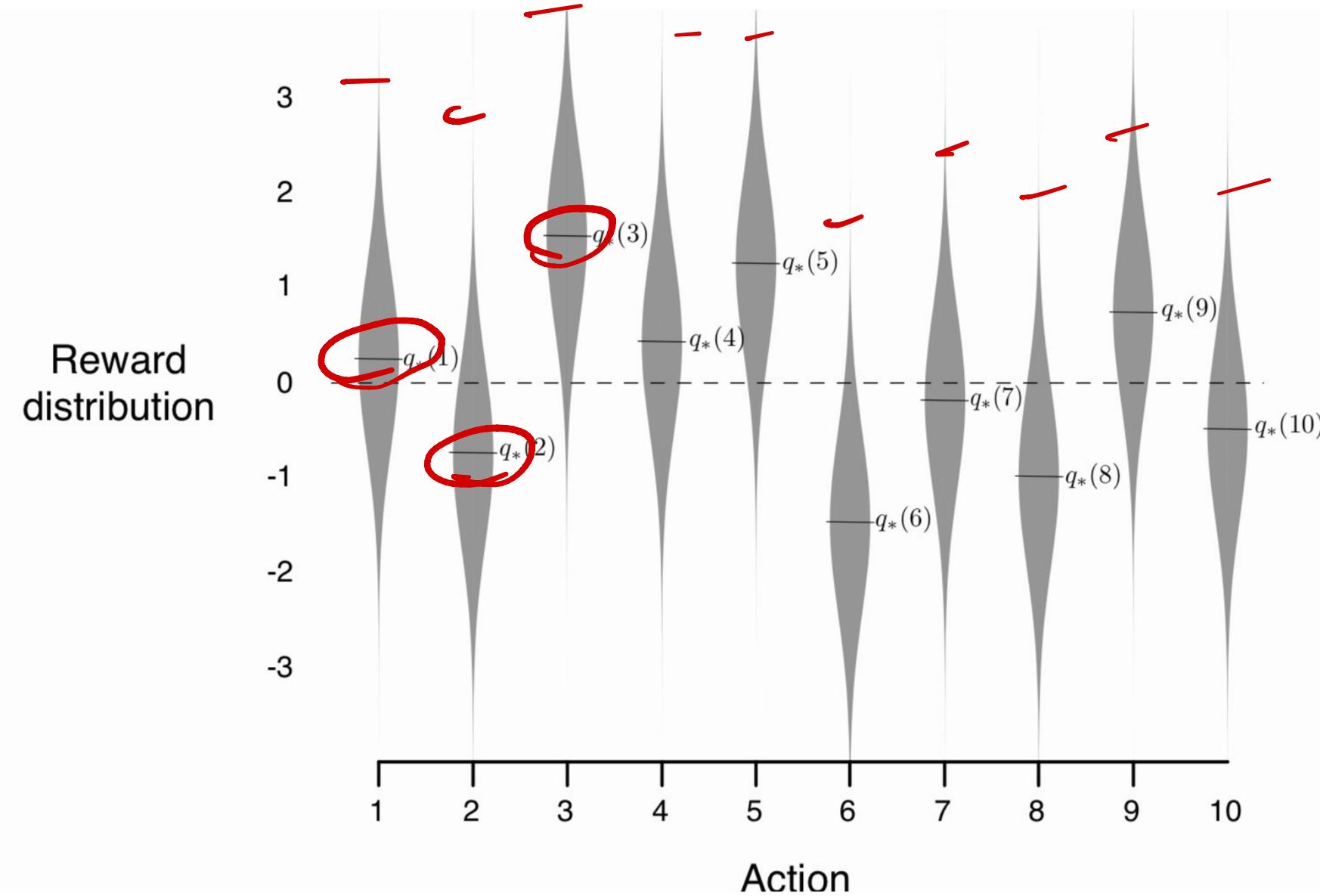
$\epsilon$ -greedy

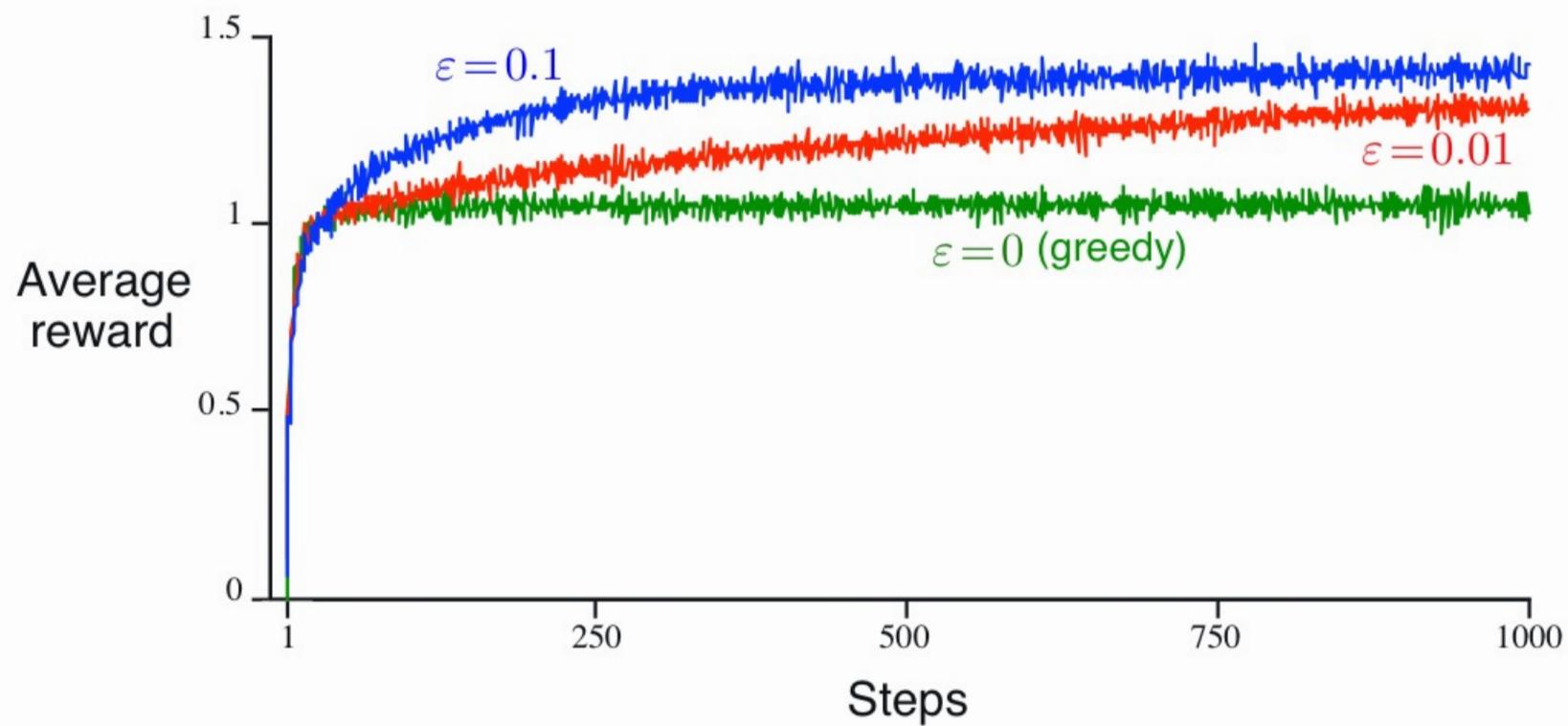
$$A_t = \begin{cases} \underset{a}{\operatorname{arg\max}} Q_t(a) & \text{w.p. } 1-\epsilon \\ \text{random other} & \text{w.p. } \epsilon \end{cases}$$

$$\underline{\epsilon = 0.01}$$

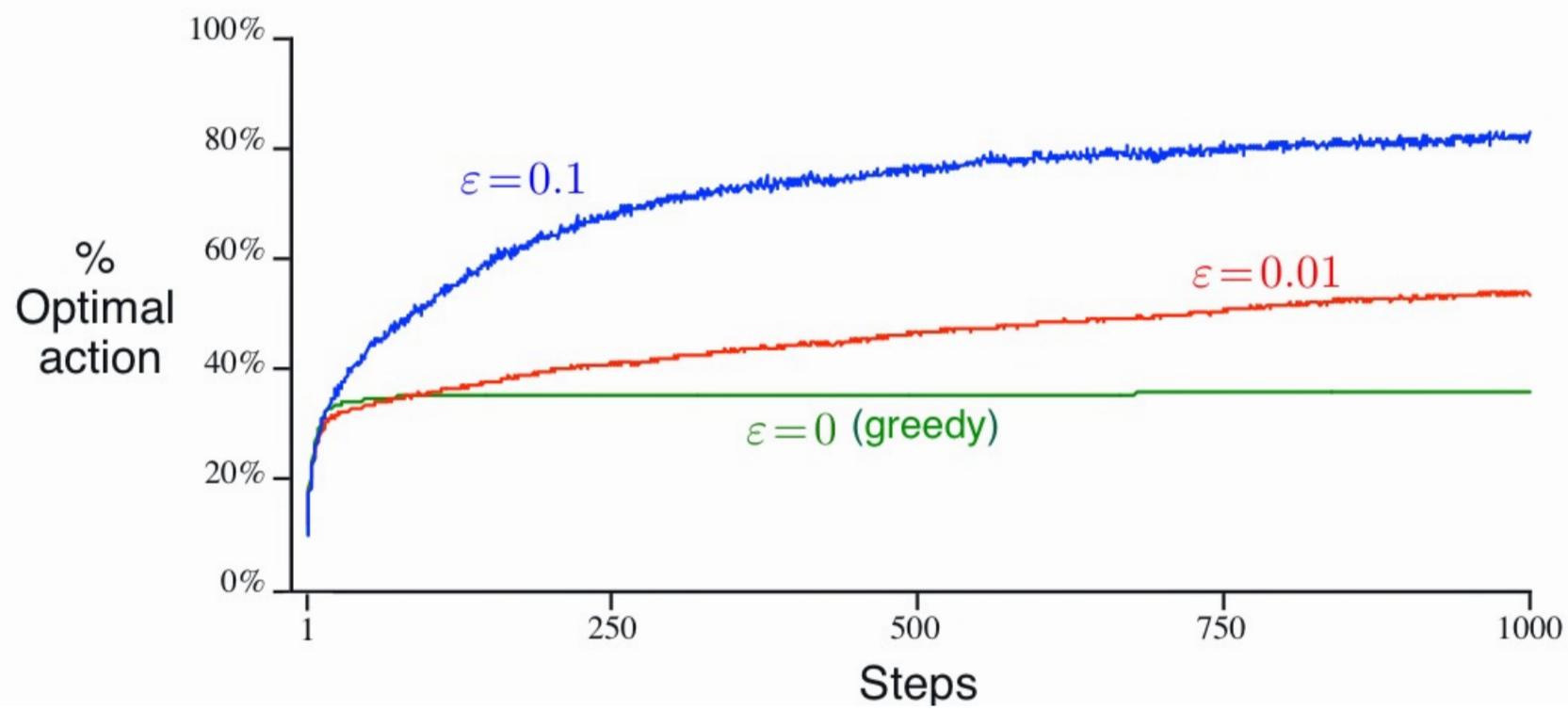
Note : If denom. is 0,  $Q_k(c) = \infty$

# 10-armed testbed :-





$\epsilon$   
= hyperparameter  
for the algo



Incremental estimation of  $Q_T(a)$ ! -

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

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

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

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

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

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

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

$$\text{New estimate} = \text{old estimate} + \\ \text{stepsize} [\text{target} - \text{old est.}]$$

# E-greedy:

## A simple bandit algorithm

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

$$Q(a) \leftarrow 0$$

$$N(a) \leftarrow 0$$

Loop forever:

$$\left\{ \begin{array}{ll} A \leftarrow \begin{cases} \arg \max_a Q(a) & \text{with probability } 1 - \varepsilon \\ \text{a random action} & \text{with probability } \varepsilon \end{cases} & \text{(breaking ties randomly)} \\ R \leftarrow \text{bandit}(A) & \\ N(A) \leftarrow N(A) + 1 & \\ Q(A) \leftarrow Q(A) + \frac{1}{N(A)} [R - Q(A)] & \end{array} \right.$$

Non-stationary problems! -

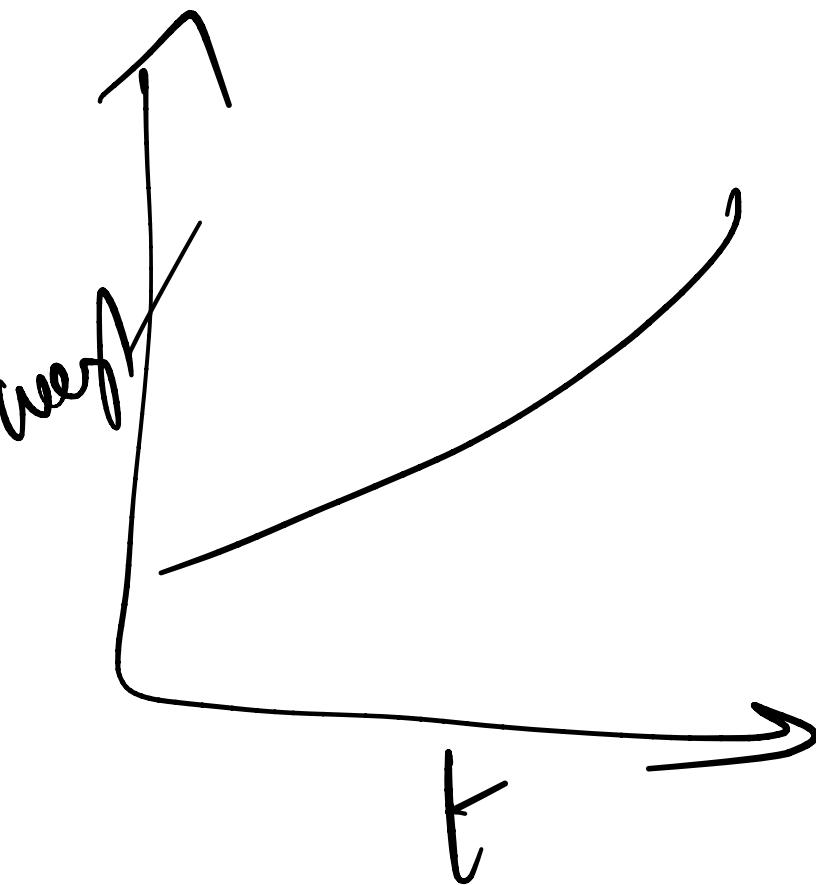
Stationary problems! -  $Q_{n+1} = Q_n + \frac{1}{n} [R_n - Q_n]$

Not ideal in non-stat. settings.

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

$\alpha \in (0, 1]$  is constant.

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



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

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

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

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

$$Q_{n+1} = (-\alpha)^n Q_1 + \sum_{i=1}^n \alpha (-\alpha)^{n-i} R_i + (-\alpha)^{n-2} Q_{n-2}$$

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

will  
not  
converge

$$d_n(a) = a$$

Stationary  
problems.

} Sample  
average

Non-stati.  
prob

} recent  
weighted  
average

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



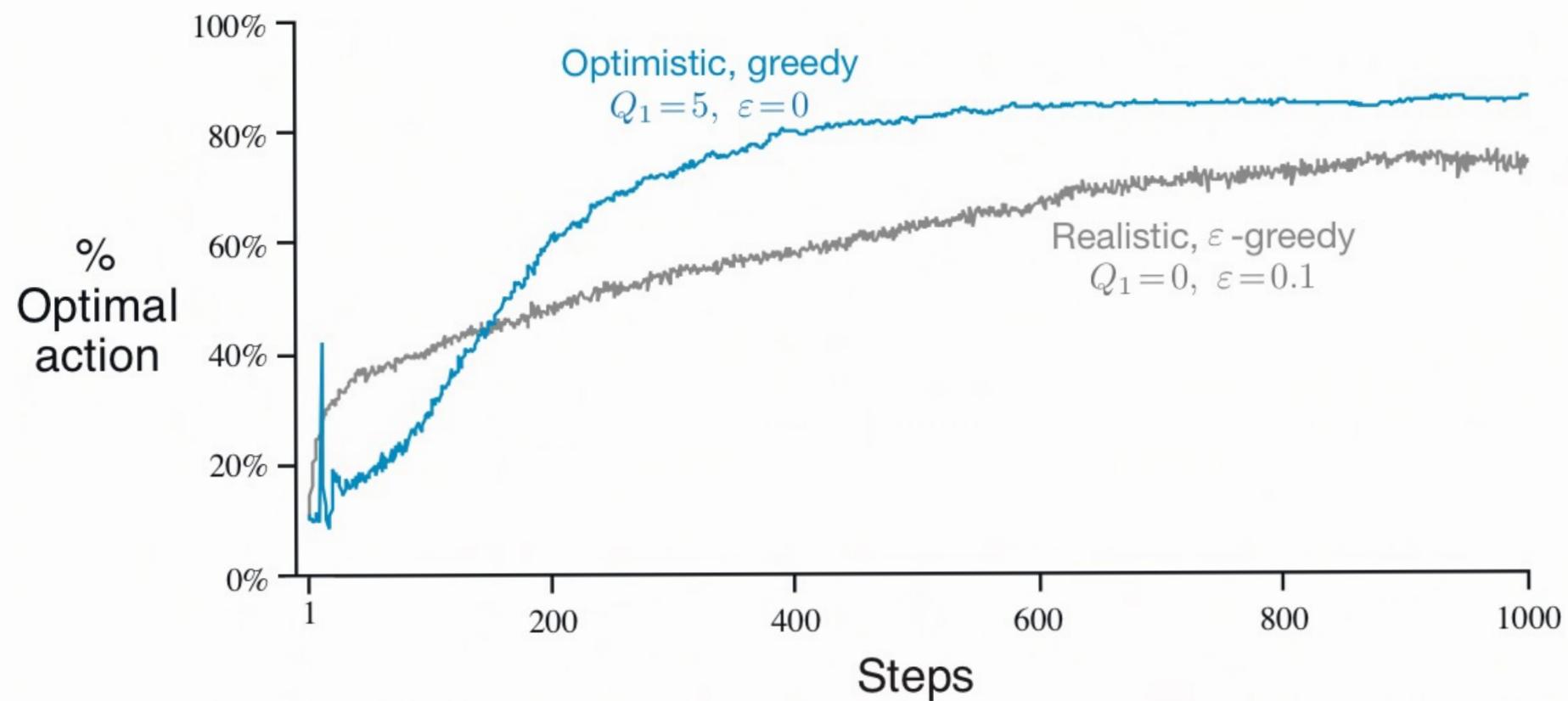
$$\sum_{n=1}^{\infty} d_n^2(a) < \infty$$



$$Q_1(a) = 0 \quad \text{for all } a.$$

→ Optimise initial value.

$$Q_1(a) = 5 \quad \text{for all } a$$



## Alg 2: Boltzmann exploration

$$Q_t(a_1) \quad Q_t(a_2) \quad \dots \quad Q_t(a_5)$$

$$Q [10, 5, 20, 15, 20]$$

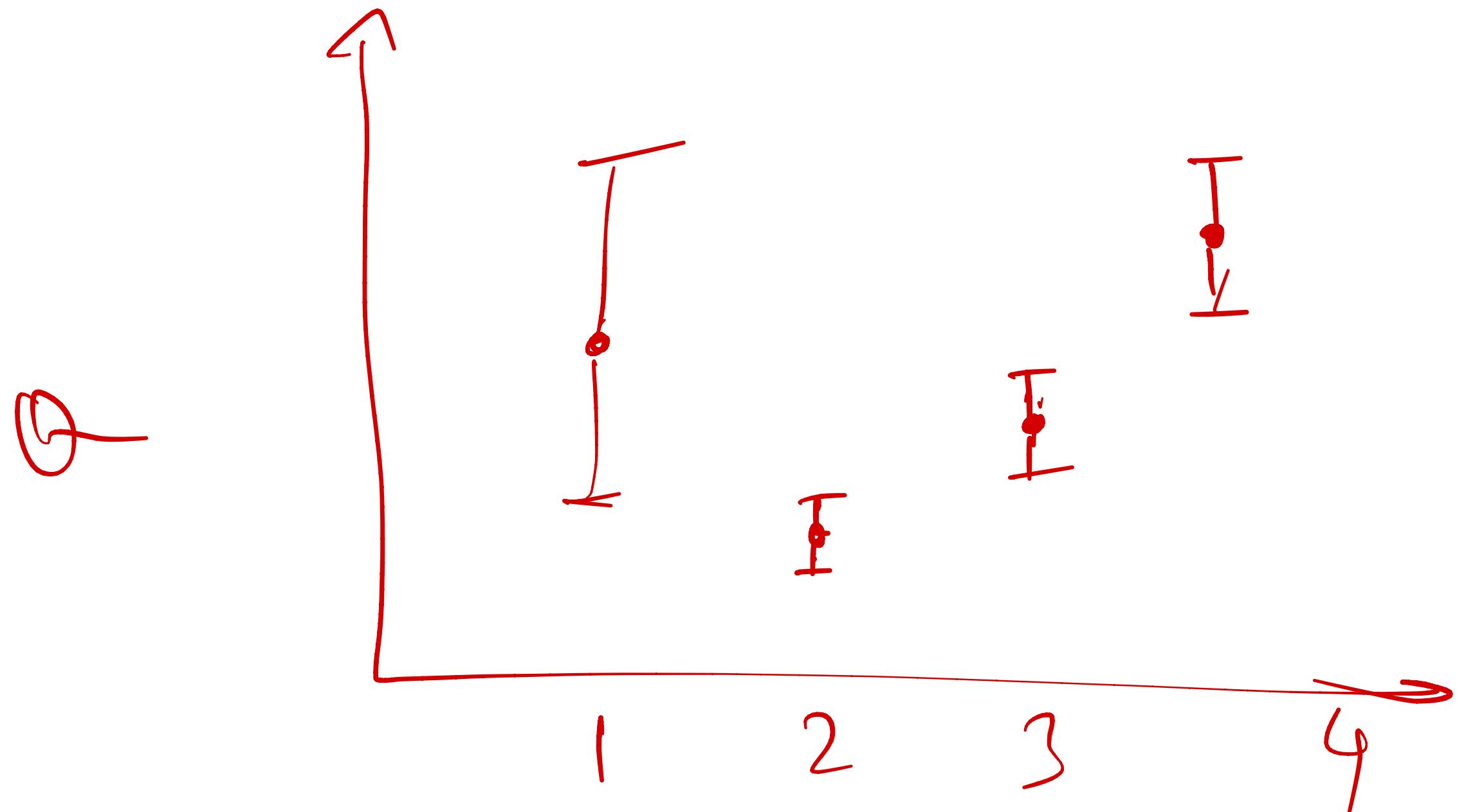
$$\text{Softmax}(Q)$$

$$P(a_i) = \frac{e^{Q_t(a_i)/\tau}}{\sum_{i=1}^k e^{Q_t(a_i)/\tau}}$$

$\tau$  temp · para.

Alg 4:

Upper Confidence Bound (UCB)

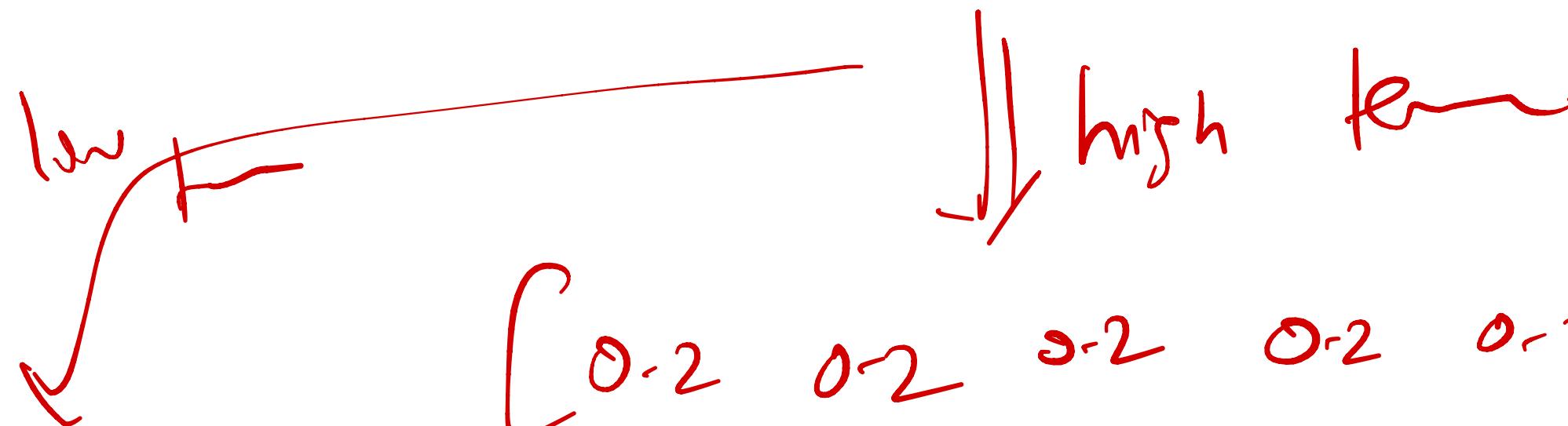


UCB1 :

$$A_T = \underset{a}{\operatorname{argmax}} \left[ Q_t(a) + C \sqrt{\frac{\ln T}{N_T(a)}} \right]$$

$[10, 5, 20, 15, 20]$

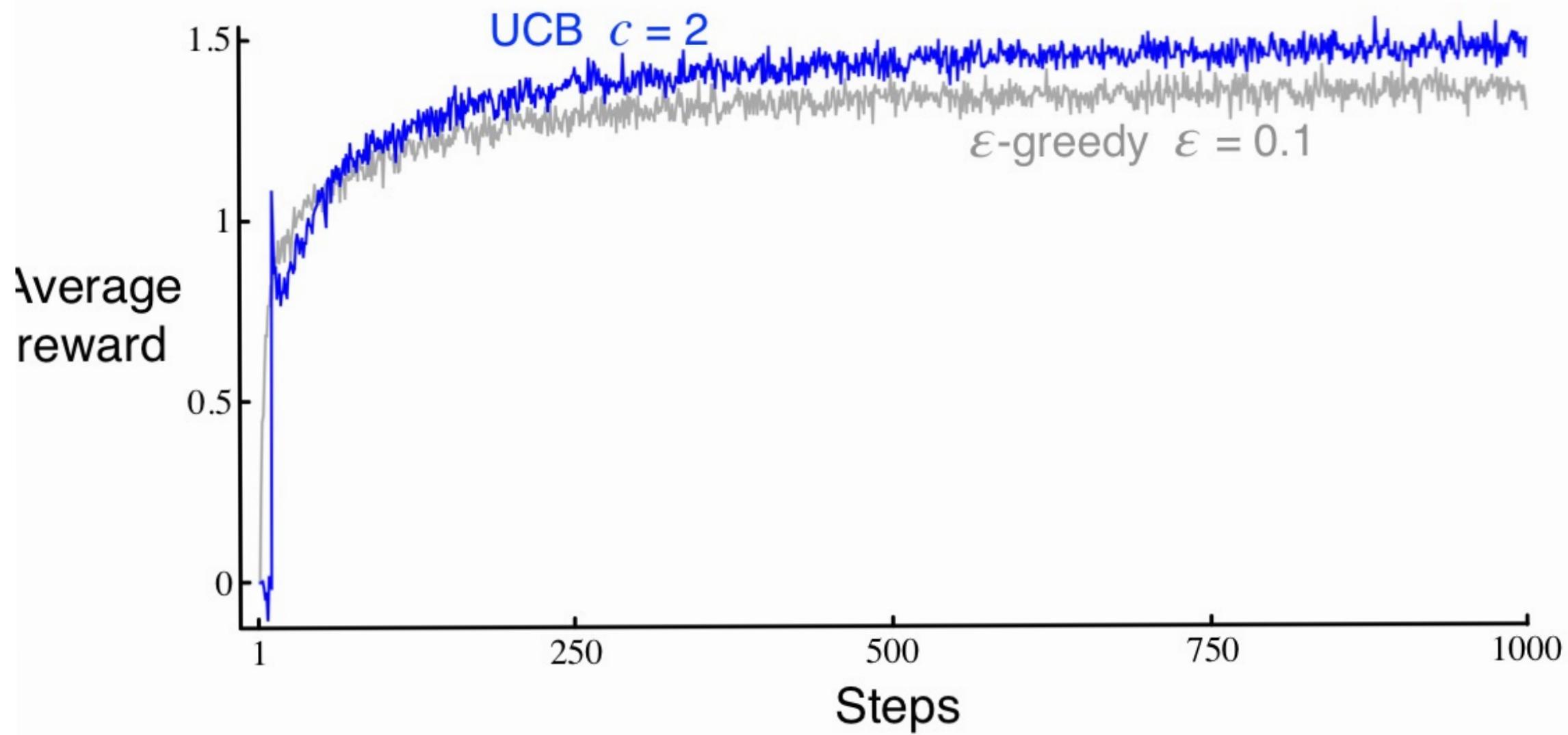
Software:  $[0.15 \ 0.05 \ 0.3 \ 0.2 \ 0.3]$



$[0.2 \ 0.2 \ 0.2 \ 0.2 \ 0.2]$

0.03 0.02 0.45 0.01 0.45

$\wedge \quad N$



Alg 5: Gradient bandit.

$$H_t(a)$$

pref. for taking action  
 $a_t$

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

$$H_t(a) = 0 \text{ for all } a.$$

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

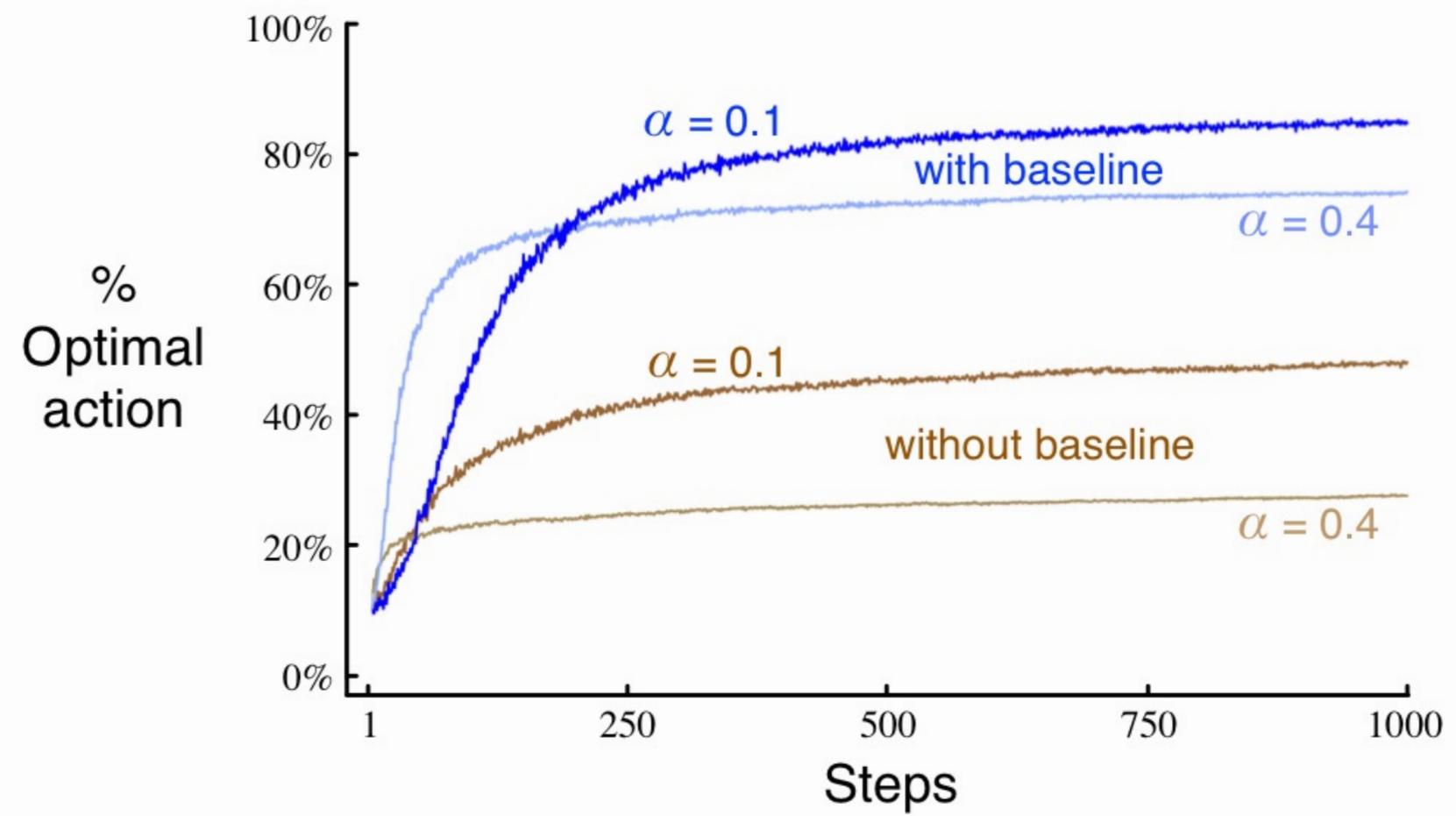
$$H_{t+1}(a) = H_t(a) + \alpha (R_t - \bar{R}_t) \pi_t(a)$$

$\forall a \neq A_t$

$\bar{R}_t$  = average reward so far.

$R_t > \bar{R}_t$  : increases the prob of  $A_t$

$R_t < \bar{R}_t$  : decreases



Gradient bandit  $\Rightarrow$  Stochastic gradient ascent.

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

$$E[R_t] = \sum_x \pi_t(x) q_x(x)$$

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

$$= \sum_n q_{rx}(x) \frac{\partial \pi_t(n)}{\partial H_t(a)}$$

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

$$\frac{\partial E[R_T]}{\partial H_T(a)} = \sum_n (q_{\pi}(n) - \beta_T) \frac{\partial \pi_T(n)}{\partial H_T(a)}$$

$$= \sum_n \pi_T(n) \cdot (q_{\pi}(n) - \beta_T) \frac{\partial \pi_T(n)}{\partial H_T(a)} / \pi_T(n)$$

$$= E \left[ (q_{\pi}(A_T) - \beta_T) \frac{\partial \pi_T(n)}{\partial H_T(a)} / \pi_T(A_T) \right]$$

$$2 E \left[ (R_T - \bar{R}_T) \frac{\partial \pi_T(n)}{\partial H_T(a)} / \pi_T(A_T) \right]$$

$$\begin{aligned}
 \frac{\partial \pi_L(x)}{\partial H_L(c)} &= \frac{\partial}{\partial H_L(c)} \pi_L(x) \\
 &= \frac{\partial}{\partial H_L(c)} \left[ \sum_{y=1}^k e^{H_L(y)} \right] \\
 &= \sum_{y=1}^k e^{H_L(y)} \frac{\partial e^{H_L(y)}}{\partial H_L(c)} - e^{H_L(n)} \frac{\partial \sum_{y=1}^k e^{H_L(y)}}{\partial H_L(c)} \\
 &\hline
 &\quad \left( \sum_{y=1}^k e^{H_L(y)} \right)^2
 \end{aligned}$$

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

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

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

$$= \mathbb{1}_{a=n} \pi_t(n) - \pi_t(n) \pi_t(c)$$

$$\frac{\partial \pi_t(n)}{\partial h_t(a)} = \pi_t(n) (1_{a=n} - \pi_t(a))$$

$$\begin{aligned} \frac{\partial E[R_t]}{\partial h_t(a)} &= E \left[ (R_t - \bar{R}_t) \frac{\partial \pi_t(A_t)}{\partial h_t(a)} \Bigg/ \pi_t(A_t) \right] \\ &= E \left[ (R_t - \bar{R}_t) \cancel{\pi_t(A_t)} (1_{a=A_t} - \pi_t(a)) \Bigg/ \cancel{\pi_t(A_t)} \right] \\ &= E \left[ (R_t - \bar{R}_t) (1_{a=A_t} - \pi_t(a)) \right] \end{aligned}$$

$$= \mathbb{E} \left[ (r_t - \bar{r}_t) \left( \mathbf{1}_{a=A_t} - \pi_t(a) \right) \right]$$

$$H_{G+T}(a) = H_T(a) + \alpha (R_t - \bar{R}_t) \left( \mathbf{1}_{a=A_t} - \pi_t(a) \right)$$

$\forall a$