



Iran University
of Science and
Technology

به نام خدا

یادگیری تقویتی در کنترل

دکتر سعید شمقدری

دانشکده مهندسی برق
گروه کنترل

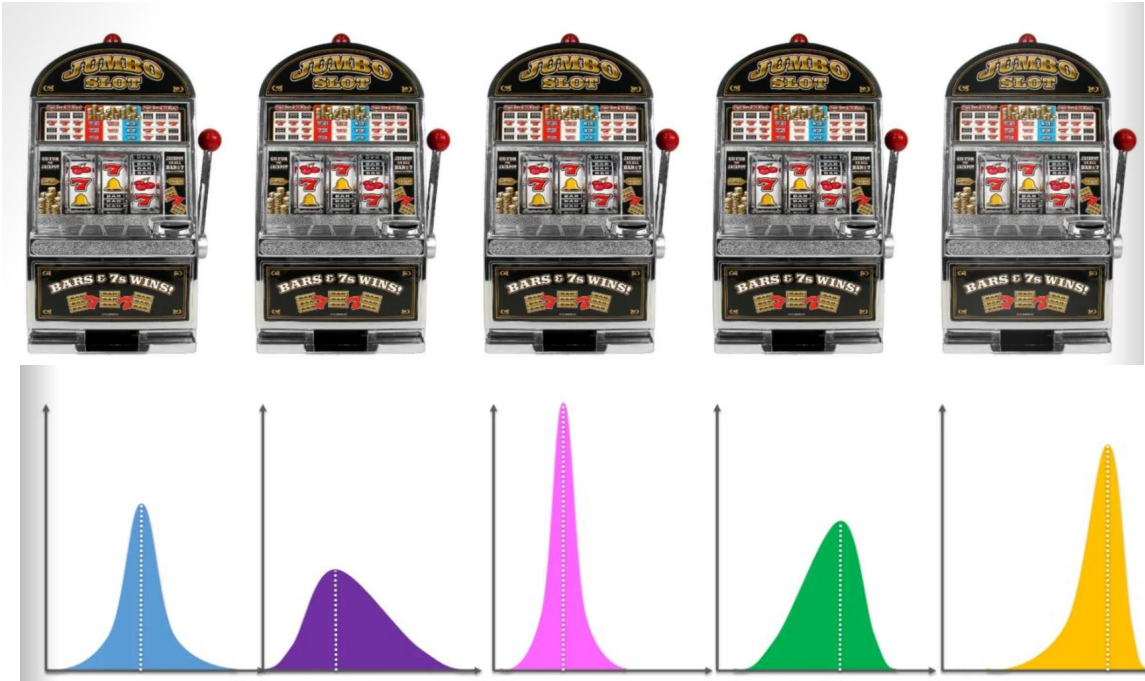
نیمسال اول ۱۴۰۵-۱۴۰۴

Multi-Armed Bandit

مسئله Multi-Armed Bandits



I Multi-Armed Bandit



مسئله Multi-Armed Bandits

Action: $a_t \in A$

Reward: $r_t \sim R$ $R_a(r) = P(r|a)$

Goal: maximize $\sum r_t$

Value: $Q(a) \approx E[r|a]$

Exact Value: $q_*(a) = E[r|a]$

Optimal Value: $\max_{a \in A} Q(a)$

Greedy and ϵ -Greedy Actions

Greedy Action: Exploitation only

Exploration: Selecting a Non-Greedy Action

Exploitation or Exploration??

Greedy Action:
 ϵ -Greedy Action

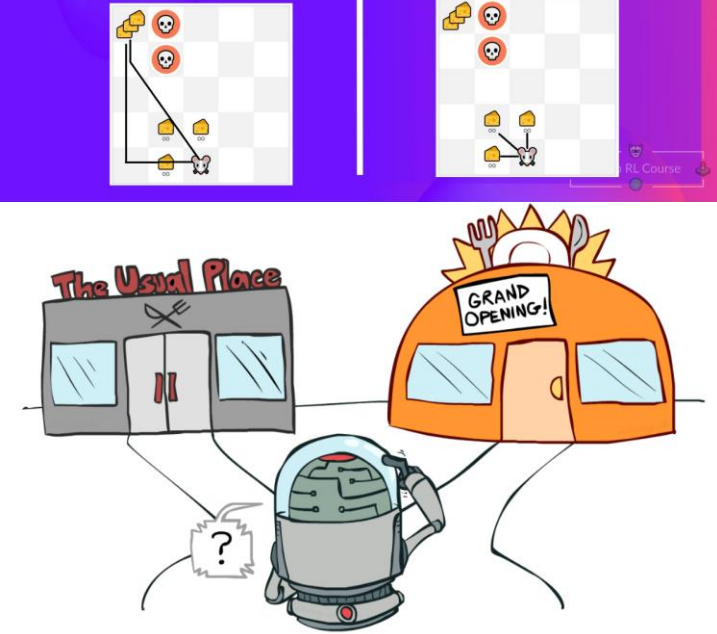
$$A_t = \arg \max Q_t(a)$$

انتخاب Action:

Exploration/ Exploitation tradeoff

Exploration: trying random actions in order to find more information about the environment.

Exploitation: using known information to maximize the reward.



یادآوری از Lecture 1:
Exploitation در مقابل Exploration

یک روش تخمین تابع Value

Sample Average:

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

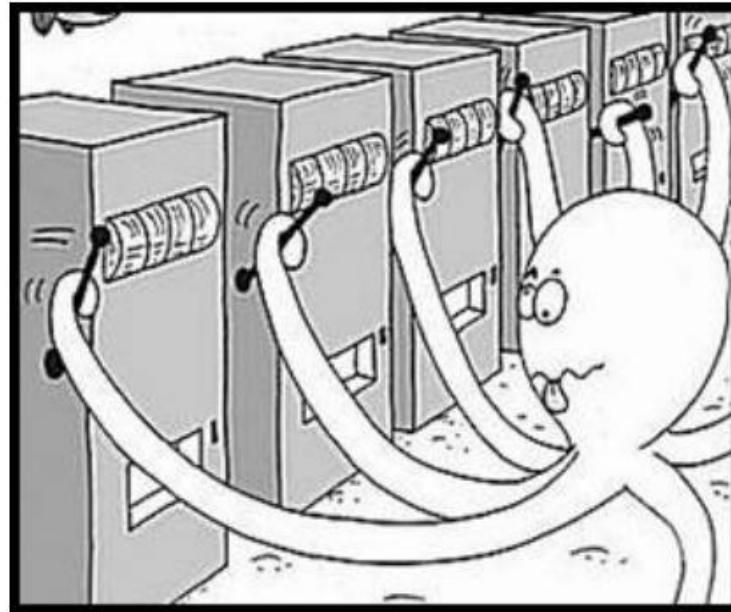
Convergence: (by the law of large numbers...)

$$Q_t(a) \rightarrow q_*(a) \quad \text{for all } a \in A$$

Q: در شرایط همگرایی برای روش ϵ -Greedy Action، احتمال انتخاب Action بهینه؟

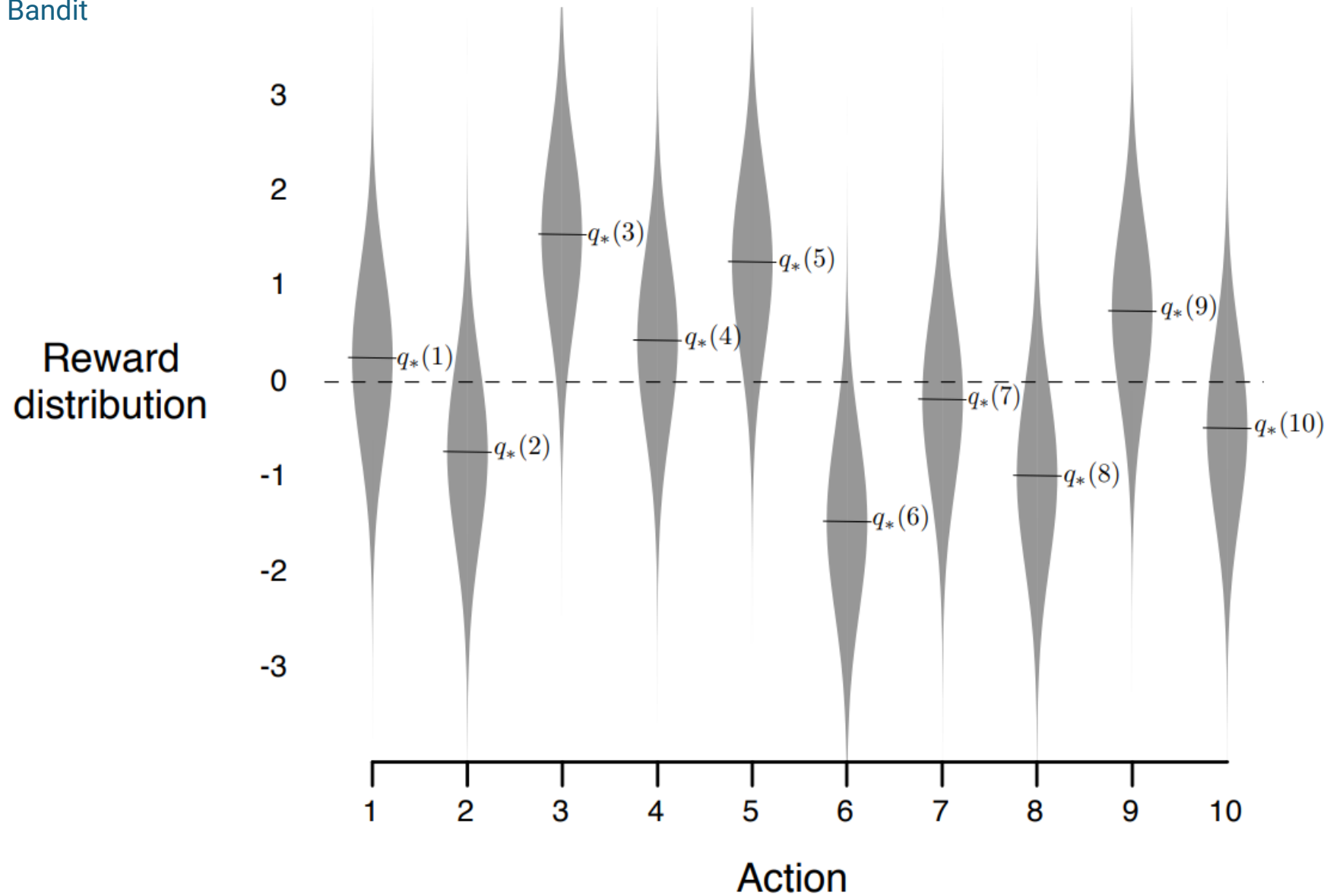
The 10-Armed Testbed

One Run? 1000 sample
Average Behavior?
Average in 2000 problem
solution

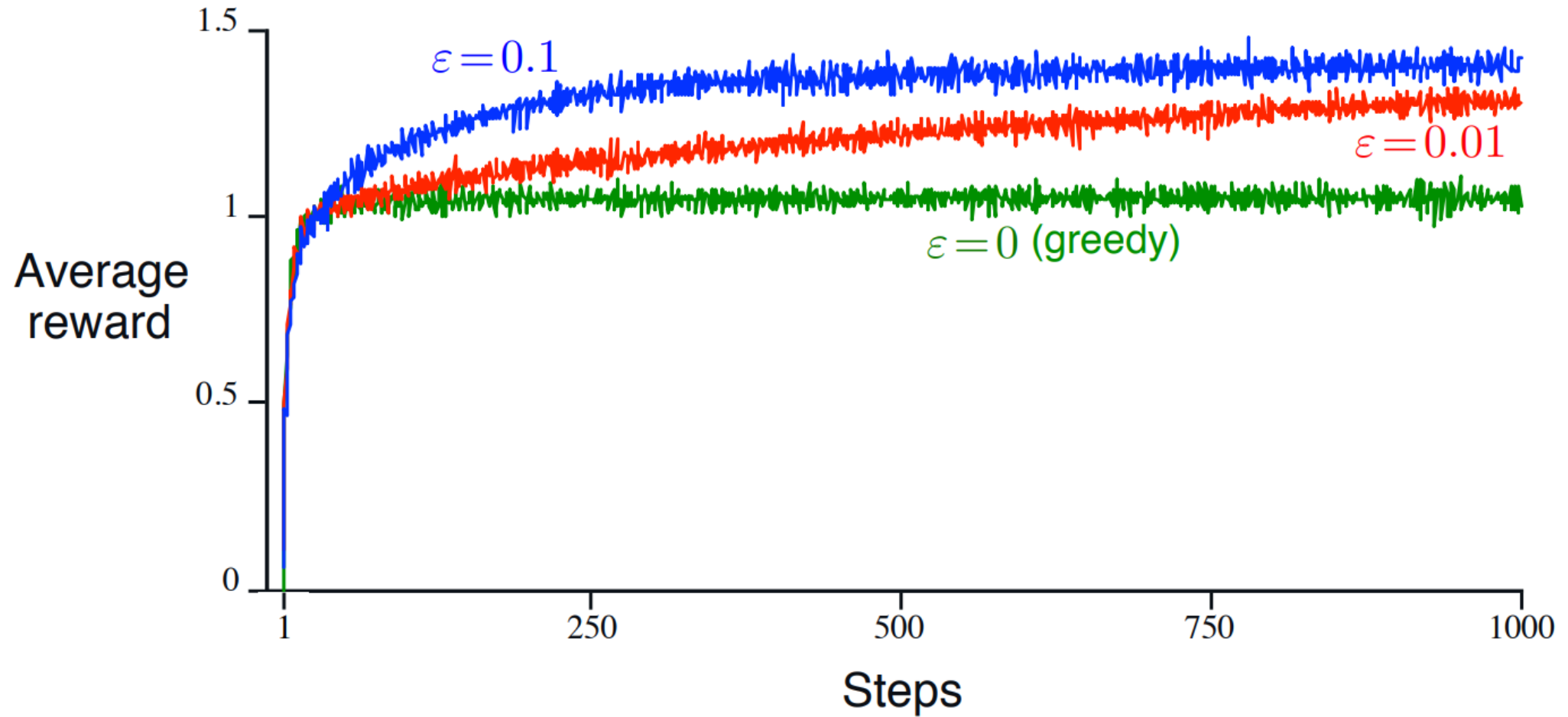


تولید مقادیر $q_*(a)$ برای هر a : توزیع نرمال با میانگین صفر و واریانس ۱
مقادیر پاداش برای هر a : توزیع نرمال با میانگین $q_*(a)$ و واریانس ۱

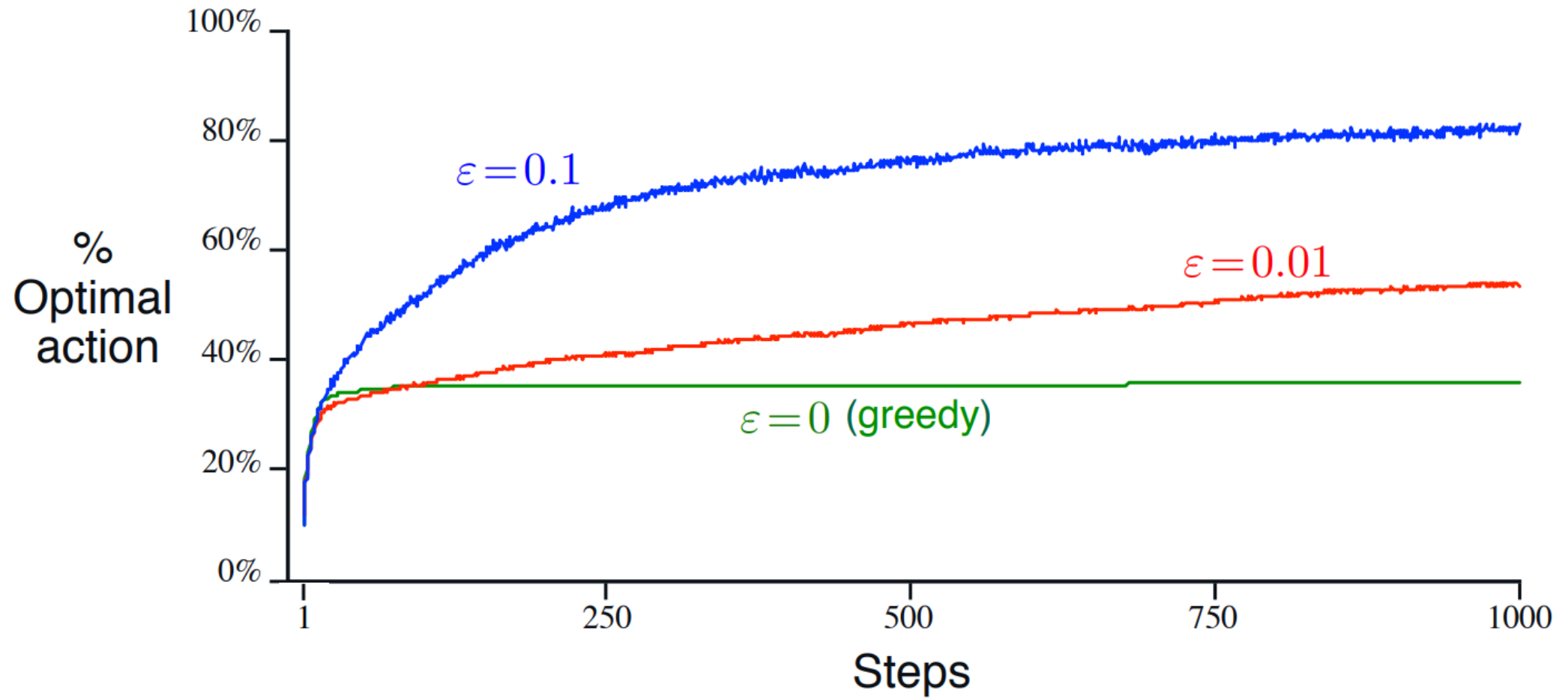
I Multi-Armed Bandit



I Greedy and ϵ -Greedy Actions



I Greedy and ϵ -Greedy Actions




Incremental Implementation

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

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

Incremental Implementation

$$NewEstimate \leftarrow OldEstimate + StepSize \underbrace{[Target - OldEstimate]}_{Error}$$

(exploration! —————) α 

حل مسئله RL برای فرایند **غیرایستا**

وزن بیشتر به پاداش های اخیر

استفاده از Weighted Average بجای Sample Average:

فرض مقدار ثابت برای $\alpha \in (0,1]$:

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

حل مسئله RL برای فرایند **غیر ایستا**

$$\begin{aligned}
 Q_{n+1} &= Q_n + \alpha [R_n - Q_n] \\
 &= \alpha R_n + (1 - \alpha) Q_n \\
 &= \alpha R_n + (1 - \alpha) [\alpha R_{n-1} + (1 - \alpha) Q_{n-1}] \\
 &= \alpha R_n + (1 - \alpha) \alpha R_{n-1} + (1 - \alpha)^2 Q_{n-1} \\
 &= \alpha R_n + (1 - \alpha) \alpha R_{n-1} + (1 - \alpha)^2 \alpha R_{n-2} + \\
 &\quad \dots + (1 - \alpha)^{n-1} \alpha R_1 + (1 - \alpha)^n Q_1 \\
 &= (1 - \alpha)^n Q_1 + \sum_{i=1}^n \alpha (1 - \alpha)^{n-i} R_i.
 \end{aligned}$$

حل مسئله RL برای فرایند غیرایستا

$$\begin{aligned}
 Q_{n+1} &= Q_n + \alpha [R_n - Q_n] \\
 &= \alpha R_n + (1 - \alpha) Q_n \\
 &= \alpha R_n + (1 - \alpha) [\alpha R_{n-1} + (1 - \alpha) Q_{n-1}] \\
 &= \alpha R_n + (1 - \alpha) \alpha R_{n-1} + (1 - \alpha)^2 Q_{n-1} \\
 &= \alpha R_n + (1 - \alpha) \alpha R_{n-1} + (1 - \alpha)^2 \alpha R_{n-2} + \\
 &\quad \dots + (1 - \alpha)^{n-1} \alpha R_1 + (1 - \alpha)^n Q_1 \\
 &= \underbrace{(1 - \alpha)^n Q_1}_{\text{Weighted Average}} + \sum_{i=1}^n \alpha (1 - \alpha)^{n-i} R_i.
 \end{aligned}$$

Q: تاثیر افزایش α در واریانس تخمین؟؟

Weighted Average

Unbiased Estimation: if $n \rightarrow \infty$ then $Q = E[R(a)] = q_*(a)$

حل مسئله RL برای فرایند **غیرایستا**

Weighted Average با α متغیر با زمان : $\alpha_n(a)$

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

and

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

عدم کوچک شدن سری α_n ها
کاهش اثر شرط اولیه و نوسانات
Exploration کافی

شرط همگرایی

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

در حالت Sample Average:

توجه: برای هر ضریب ثابت برقرار نیست!

حل مسئله RL برای فرایند **غیرایستا**

در حالت Sample Average:

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

توجه: برای هر ضریب ثابت برقرار نیست!

$$\begin{aligned} 1 + \frac{1}{2} + \frac{1}{3} + \frac{1}{4} + \dots &\geq 1 + \frac{1}{2} + \frac{1}{4} + \frac{1}{4} + \left(\frac{1}{8} + \frac{1}{8} + \frac{1}{8} + \frac{1}{8} \right) + \dots + \left(\frac{1}{16} + \dots \right) \\ &\geq 1 + 1 + 1 + \dots = \infty \end{aligned}$$

حل مسئله RL برای فرایند **غیرایستا**


تخمینگر Unbiased:

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

$$\begin{aligned} Q_n &= Q_{n-1} + \alpha_{n-1} [R_{n-1} - Q_{n-1}] \\ &= (1 - \alpha_{n-1}) Q_{n-1} + \alpha_{n-1} R_{n-1} \end{aligned}$$

$$\begin{aligned} Q_{n+1} &= \alpha_n R_n + (1 - \alpha_n) \alpha_{n-1} R_{n-1} + (1 - \alpha_n) (1 - \alpha_{n-1}) Q_{n-1} \\ &= \alpha_n R_n + (1 - \alpha_n) \alpha_{n-1} R_{n-1} + (1 - \alpha_n) (1 - \alpha_{n-1}) \cdot \alpha_{n-2} R_{n-2} + \dots \end{aligned}$$

$$\begin{aligned} q_* &= E[R] (\alpha_n + (1 - \alpha_n) \alpha_{n-1} + (1 - \alpha_n) (1 - \alpha_{n-1}) \alpha_{n-2} \\ &\quad + \dots + (1 - \alpha_n) (1 - \alpha_{n-1}) \dots (1 - \alpha_1)) \end{aligned}$$

$$q_* = E[R] (\alpha_n + (1 - \alpha_n) [\alpha_{n-1} + (1 - \alpha_{n-1}) \alpha_{n-2} + \dots])$$


حل مسئله RL برای فرایند **غیرایستا**

تخمینگر Unbiased:

$$q_* = E[R](\alpha_n + (1 - \alpha_n)\alpha_{n-1} + (1 - \alpha_n)(1 - \alpha_{n-1})\alpha_{n-2} \\ + \dots + (1 - \alpha_n)(1 - \alpha_{n-1}) \dots (1 - \alpha_1))$$

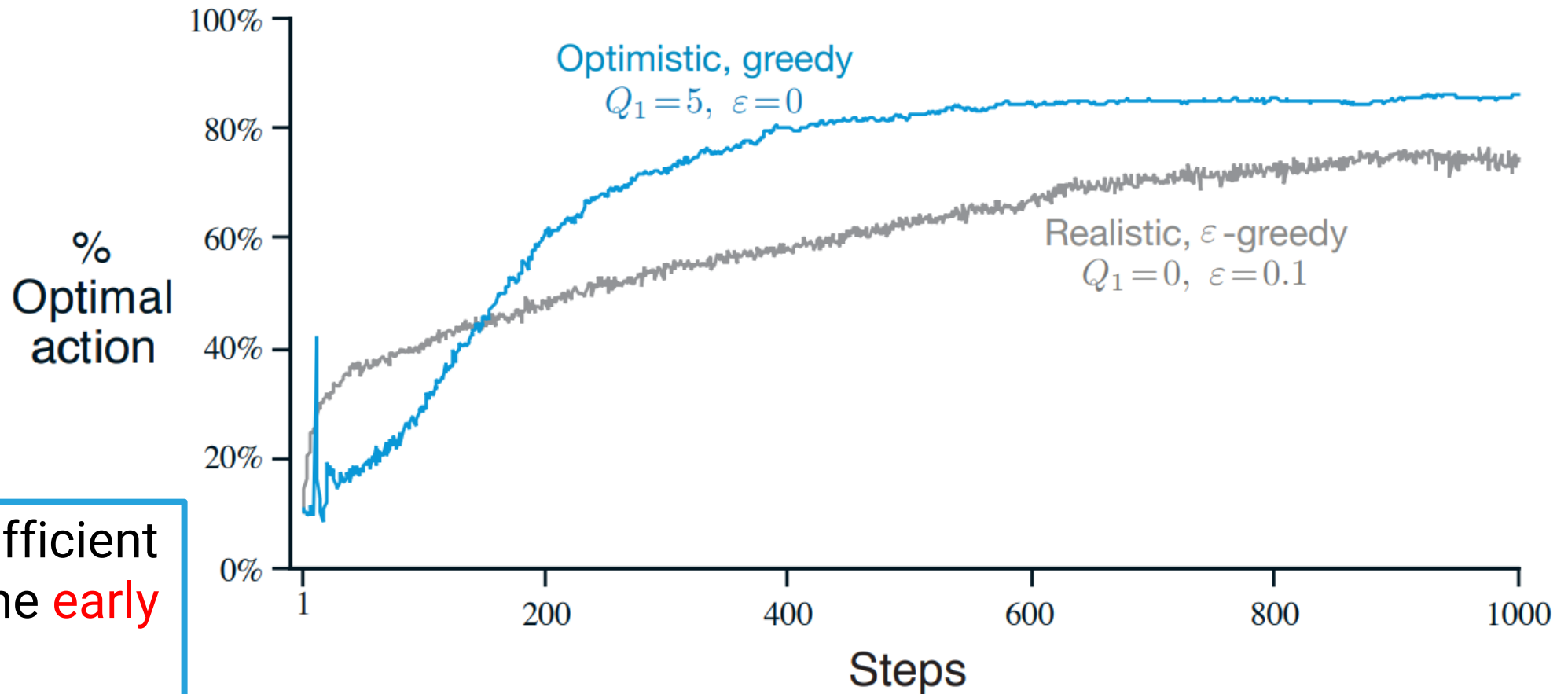
$$q_* = E[R](\alpha_n + (1 - \alpha_n)[\underbrace{\alpha_{n-1} + (1 - \alpha_{n-1})\alpha_{n-2} + \dots}_S])$$

$$S = \alpha_n + (1 - \alpha_n)S$$

$$(1 + \alpha_n)S = S + \alpha_n \rightarrow S = 1$$

انتخاب Optimistic مقادير اوليه

مثال: $Q_1 = 5$



The need for sufficient exploration in the **early** stages.

Upper Confidence Bound

عدم قطعیت در تخمین Value: نیاز به exploration

Q: ضعف روش ϵ -Greedy؟

در روش ϵ -greedy اکتشاف از بین عمل‌های غیرحریصانه به صورت تصادفی انجام می‌شود.

مناسب است که میزان احتمال بهینه بودن عمل‌های غیرحریصانه نیز در explore لحاظ شود.

Upper Confidence Bound

لحاظ کردن عدم قطعیت تخمین!
میزان عدم قطعیت

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

کنترل میزان exploration

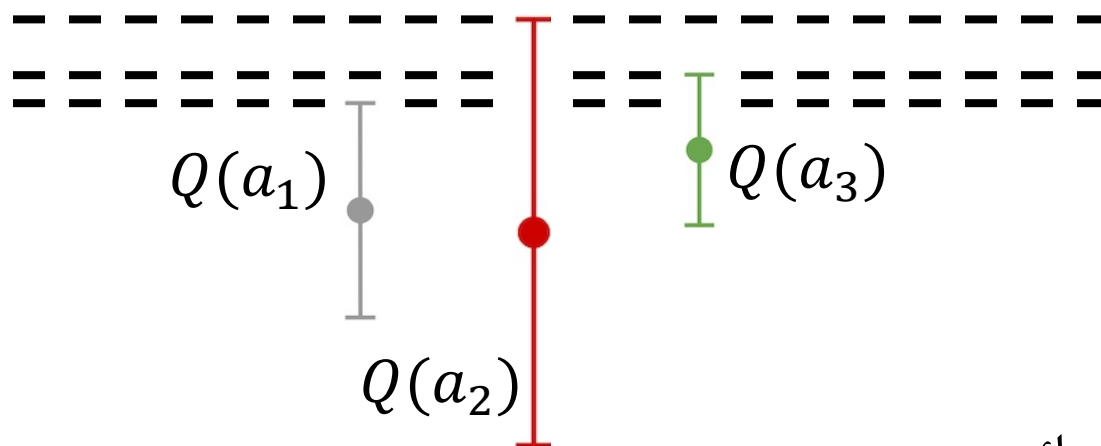
دفعاتی که قبل از t اکشن a انتخاب شده است (if zero?)

هر بار انتخاب a ؟ کاهش عدم قطعیت
هر بار انتخاب نشدن a ؟ افزایش عدم قطعیت

Upper Confidence Bound

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

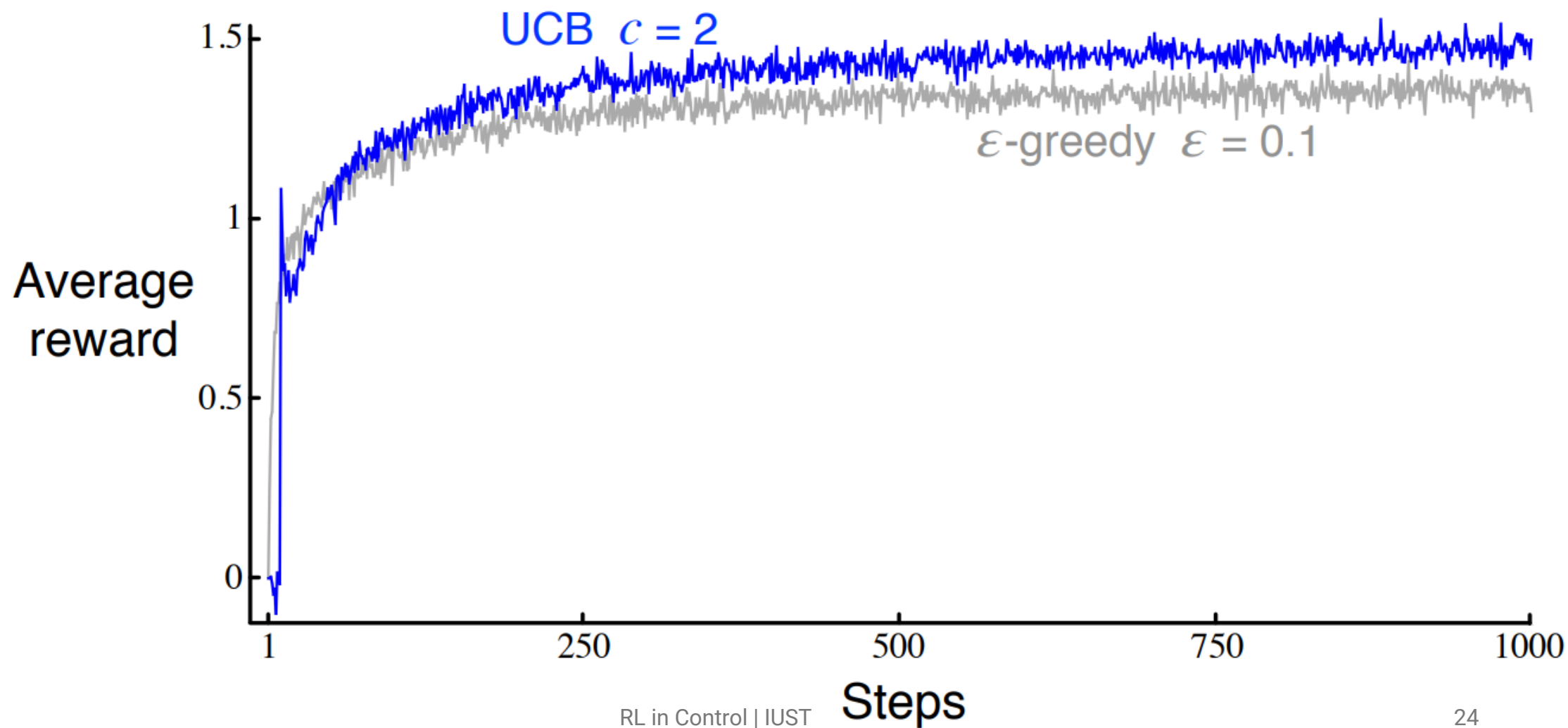
تخمین نقطه ای ← تخمین بازه ای



انتخاب تمام اکشن ها در نهایت
انتخاب با تکرار کمتر برای اکشن های با value کم
Exploration واقع بینانه تر

انجام اتوماتیک Exploration و Exploitation به صورت دائم

Upper Confidence Bound



الگوریتم Gradient Bandit

استفاده از یک معیار **preference** بجای تخمین action value $H_t(a)$

احتمال انتخاب a مطابق با توزیع soft-max:

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

الگوریتم Gradient Bandit

$$H_1(a) = 0 \text{ for all } a \in \mathcal{A} \quad :t=1$$

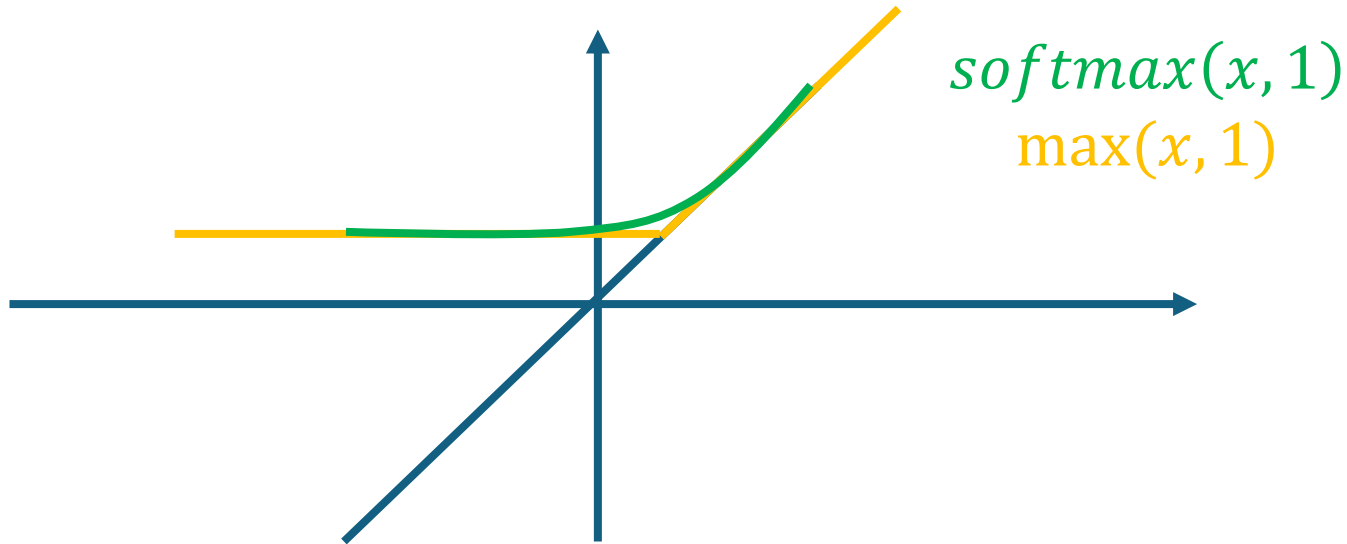
بروزرسانی در $t > 1$:

$$H_{t+1}(A_t) \doteq H_t(A_t) + \alpha(R_t - \bar{R}_t)(1 - \pi_t(A_t)), \quad \text{and}$$

$$H_{t+1}(a) \doteq H_t(a) - \alpha(R_t - \bar{R}_t)\pi_t(a), \quad \text{for all } a \neq A_t$$

متوسط پاداش ها تا زمان t (شامل t)
 بروز شدن H_t اکشنی که انتخاب نشده است

Baseline { $R_t > \bar{R}_t$ افزایش احتمال انتخاب A_t در آینده
 $R_t < \bar{R}_t$ کاهش احتمال انتخاب A_t در آینده

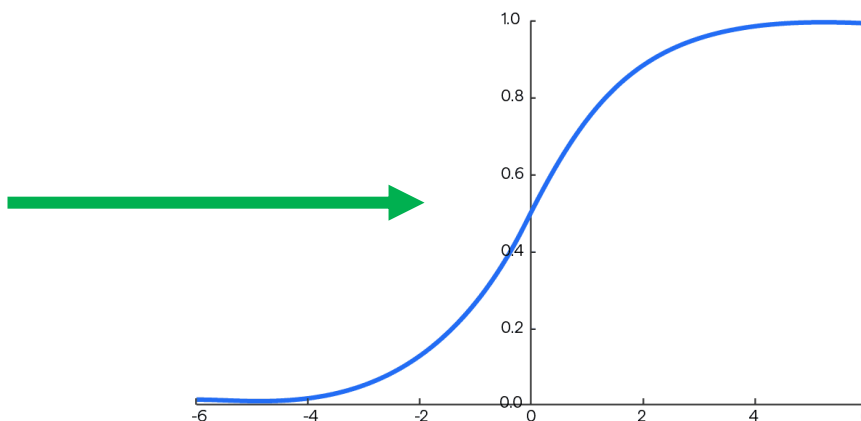


توزیع Soft-max

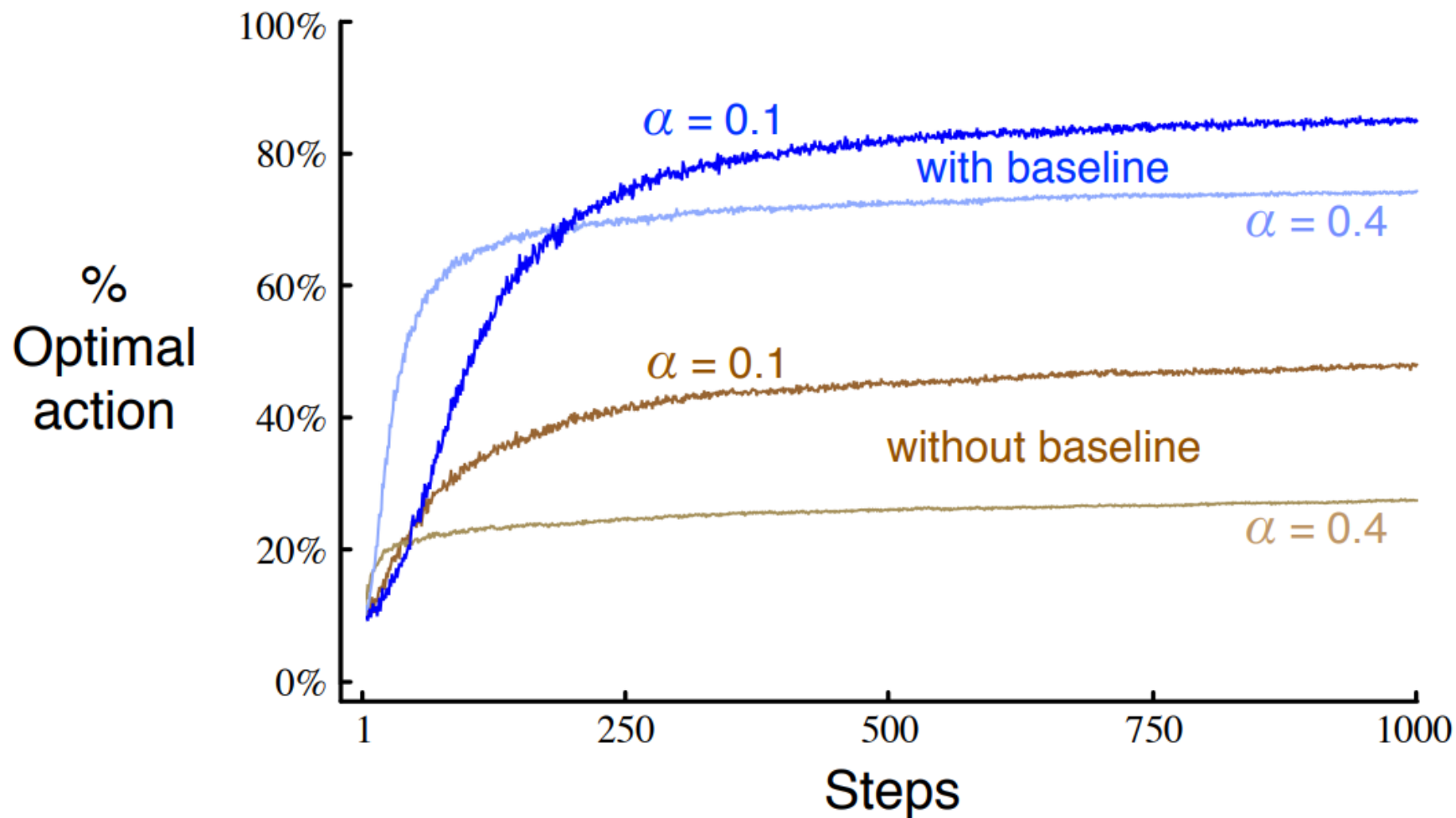
$$\text{softmax}(x_1, \dots, x_n) = \log \sum_{i=1}^n e^{x_i}$$

در یادگیری ماشین و شبکه‌های عصبی:
 برای Classification در طبقه‌ی خروجی مقادیر عددی تولید می‌شود. احتمال اینکه ورودی شبکه به کدام کلاس متعلق است از رابطه‌ی زیر بدست می‌آید.

$$\sigma(\vec{z})_i = \frac{e^{z_i}}{\sum_{j=1}^K e^{z_j}}$$

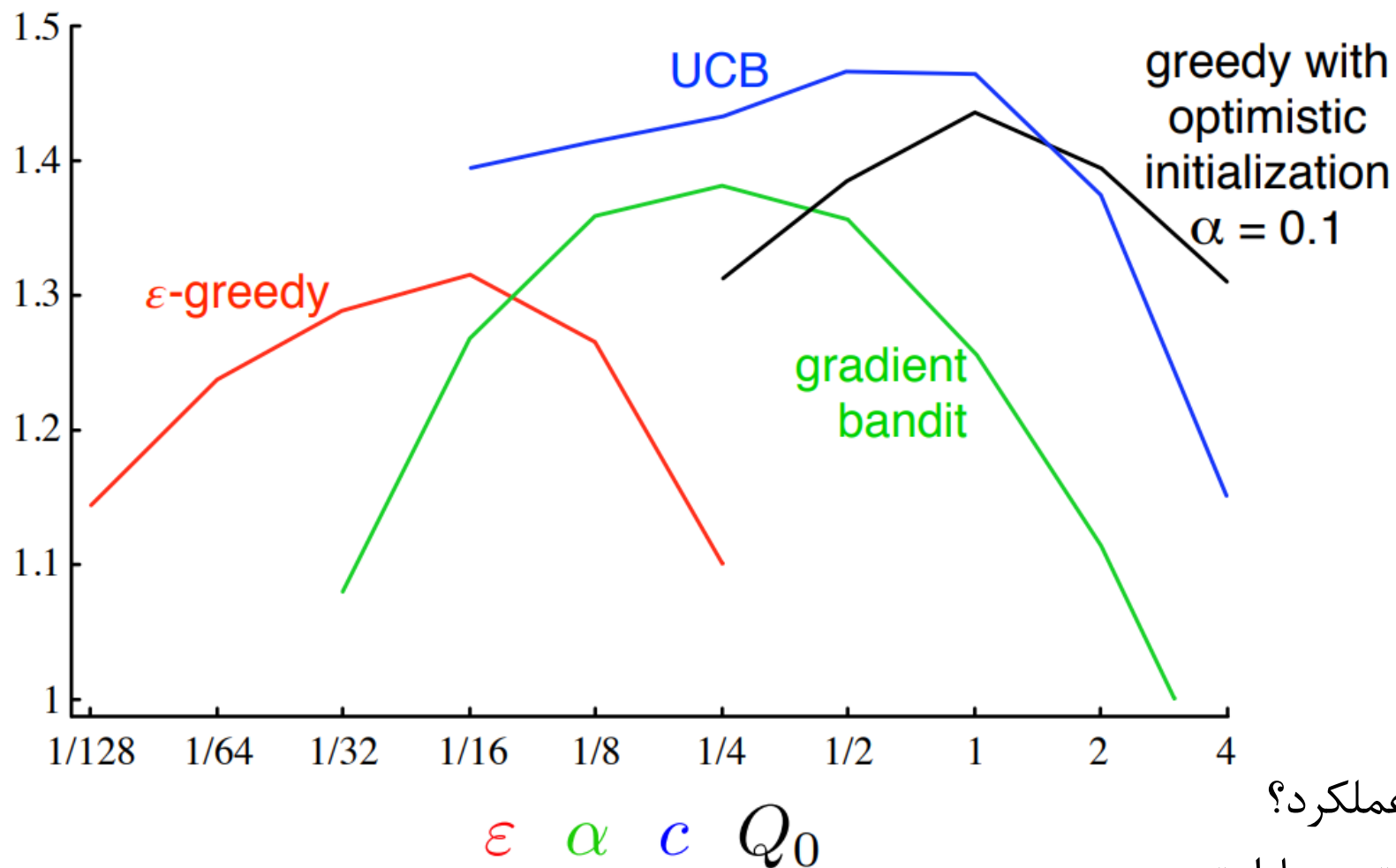


Gradient Bandit **الگوریتم**



مقایسه

Average
reward
over first
1000 steps



← Q: محدوده مناسب عملکرد؟
توجه به حساسیت تغییرات به پارامتر