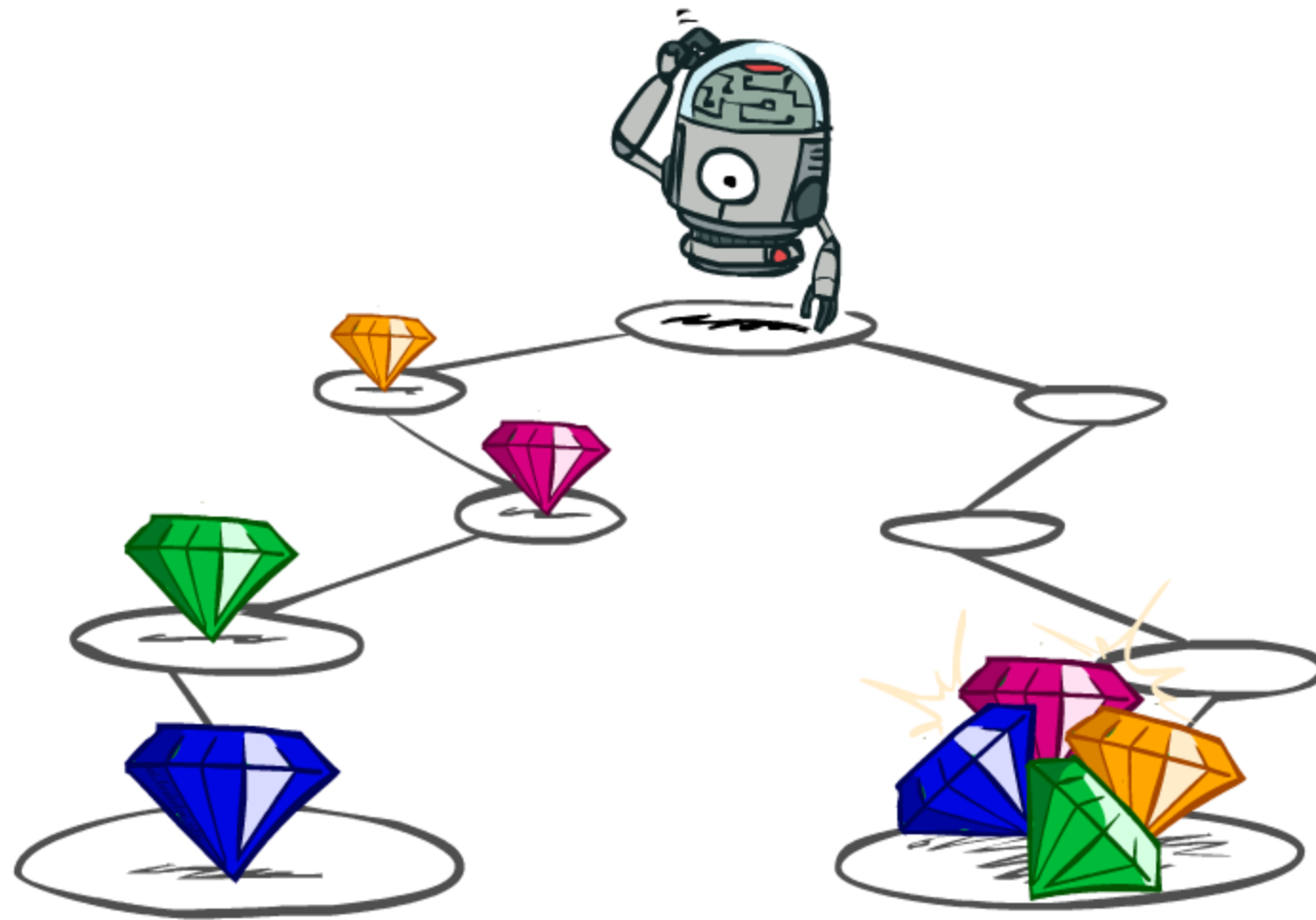


# یادگیری تقویتی (Reinforcement Learning)

1. فرایند تصمیم مارکوف (Markov Decision Processes)
2. الگوریتم تکرار مقدار (Value Iteration)
3. الگوریتم تکرار سیاست (Policy Iteration)
4. یادگیری تقویتی (Reinforcement Learning)

# سودمندی (Utility) یک دنباله از عملیات

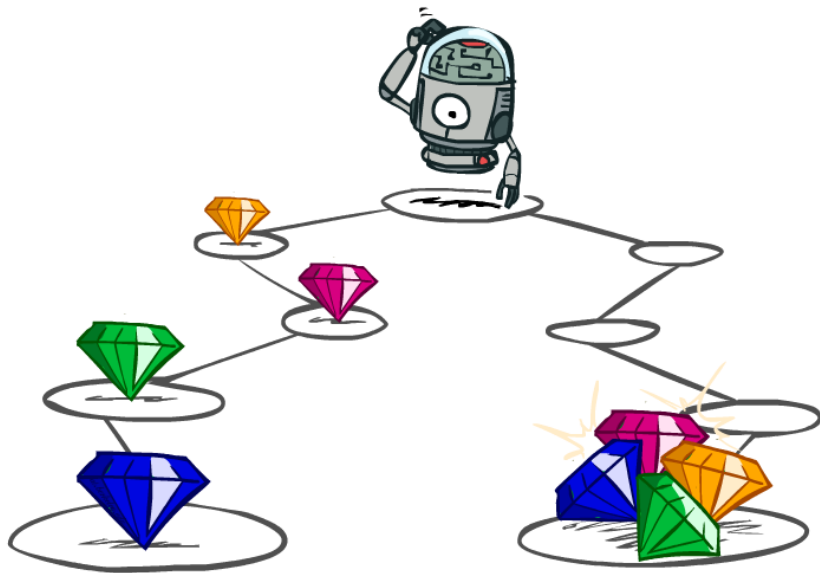


# سودمندی یک دنباله از عملیات

□ عامل باید کدام دنباله از پاداش‌ها را ترجیح دهد؟

□ کمتر یا بیشتر؟  
[1, 2, 2]      or      [2, 3, 4]

□ حالا یا بعدا؟  
[0, 0, 1]      or      [1, 0, 0]



# کاهش پاداش‌ها (discounting)

- منطقی است که عامل بخواهد مجموع پاداش‌ها رو بیشینه کند.
- همین‌طور منطقی است که پاداش‌های الان رو به پاداش‌های آینده ترجیح دهد.
- یکی از راه‌حل‌ها این است که ارزش پاداش‌ها با گذشت زمان به‌صورت نمایی کاهش پیدا کند.



1

ارزش کنونی



$\gamma$

ارزش پس از یک  
مرحله



$\gamma^2$

ارزش پس از دو مرحله

✓ ضریب کاهش =  $\gamma$

# کاهش پاداش‌ها

## چگونگی کاهش پاداش‌ها

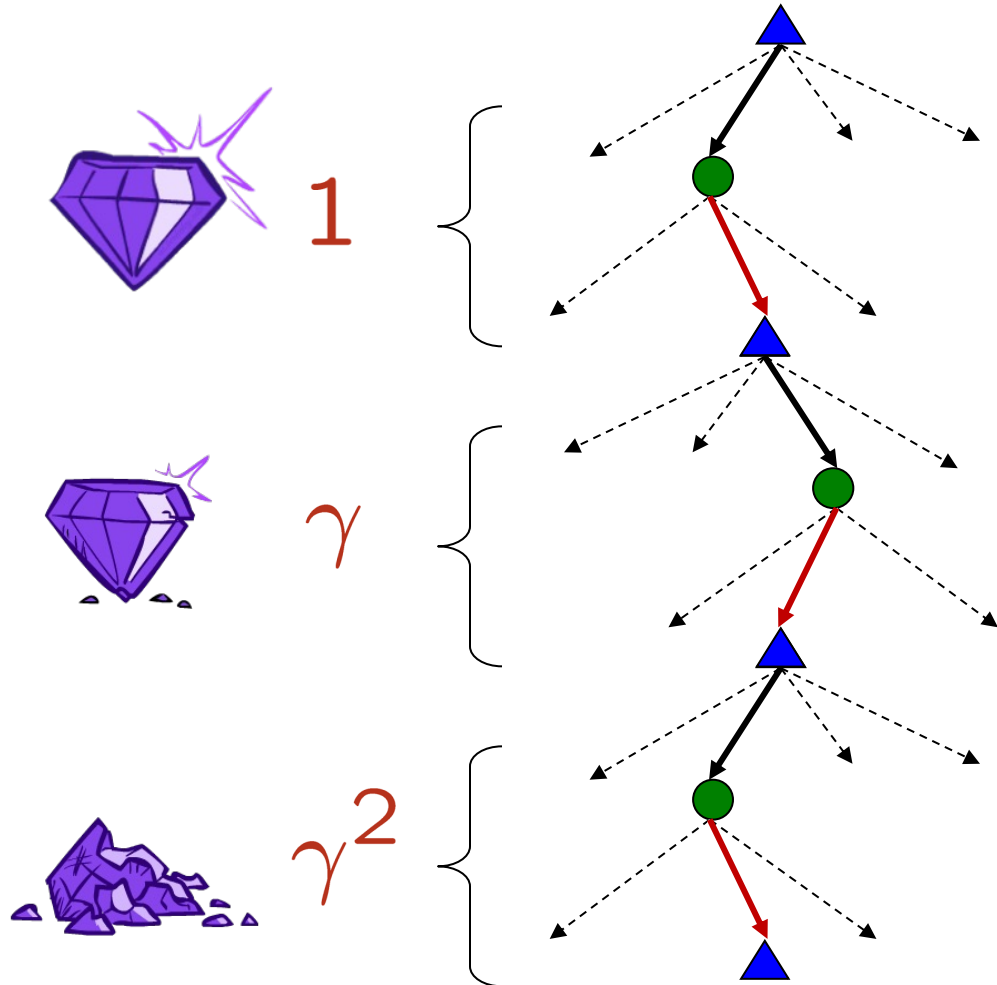
- هر بار که یک سطح پایین‌تر می‌رویم (یعنی یک گام جلوتر در زمان)، یک بار ضریب کاهش را در پاداش ضرب می‌کنیم.

## چرا باید مقدار پاداش را کاهش دهیم؟

- چون گرفتن پاداش در حال حاضر بهتر از گرفتن همان پاداش در آینده است.
- همچنین می‌توان این‌طور نگاه کرد که در هر گام، با احتمال  $1-\gamma$  ممکن است فرآیند تمام شود.
- کاهش مقدار پاداش به همگرایی الگوریتم‌ها کمک می‌کند

## مثال: ضریب کاهش پاداش = 0.5

- $U([1,2,3]) = 1*1 + 0.5*2 + 0.25*3$
- $U([1,2,3]) < U([3,2,1])$



$$U([r_0, r_1, r_2, \dots]) = r_0 + \gamma r_1 + \gamma^2 r_2 \dots$$

# کوئیز: کاهش پاداش‌ها

10				1
a	b	c	d	e

❑ داده‌های مسئله:

- اعمال: شرق، غرب و خروج (فقط در حالت‌های a و e قابل انجام است)
- با انجام یک عمل قطعا به حالت مشخصی می‌رویم (Transitions deterministic)

10	<-	<-	<-	1
----	----	----	----	---

➤ کوئیز ۱: برای  $\gamma = 1$  سیاست بهینه چیست؟

10	<-	<-	->	1
----	----	----	----	---

➤ کوئیز ۲: برای  $\gamma = 0.1$  سیاست بهینه چیست؟

➤ کوئیز ۳: در حالت d،  $\gamma$  چه مقداری باشد تا به سمت شرق یا غرب حرکت شود سودمندی یکسان باشد؟

$$1 = 10 \gamma^2 \longrightarrow \gamma = \sqrt{\frac{1}{10}} \approx 0.316$$

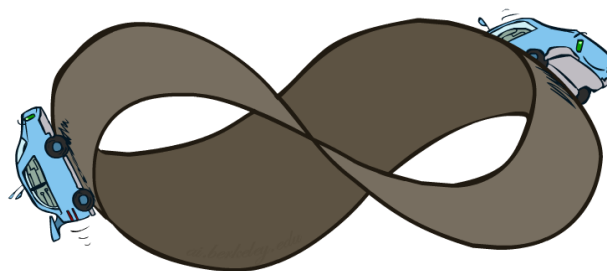
# سودمندی نامتناهی؟!!

□ مسئله: اگر بازی تا ابد ادامه پیدا کند، آیا پاداش بی نهایت دریافت خواهیم کرد؟

■ راه حل ها:

1. افق متناهی

- مشابه جستجوی عمقی محدود
- عامل فقط برای تعداد مشخصی از مراحل (مثلاً 100 مرحله) تصمیم گیری می کند.



2. استفاده از ضریب کاهش پاداش:  $0 < \gamma < 1$

• ضریب  $\gamma$  کوچکتر = افق محدودتر = تمرکز بر پاداش های نزدیکتر

$$U([r_0, \dots, r_\infty]) = \sum_{t=0}^{\infty} \gamma^t r_t \leq R_{\max} / (1 - \gamma)$$

سری هندسی با قدر نسبت  $\gamma$ :  $1 + \gamma + \gamma^2 + \dots = \frac{1}{1 - \gamma}$

3. اطمینان حاصل می کنیم که تحت هر سیاستی، در نهایت به یک حالت پایانی خواهیم رسید.

- مانند حالت «جوش» در مسابقه ای اتومبیل رانی

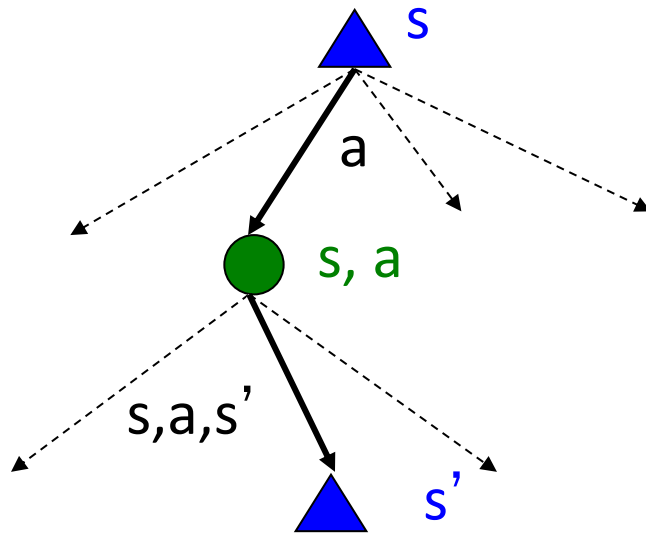
# مرور: فرایند تصمیم مارکوف

□ یک MDP به صورت زیر تعریف می شود:

- یک مجموعه از حالت ها  $S$
- یک مجموعه از اعمال  $A$
- یک تابع تغییر حالت  $T(s, a, s')$  (یا  $P(s'|s, a)$ )
- یک تابع پاداش  $R(s, a, s')$  (و ضریب کاهش پاداش  $\gamma$ )
- یک حالت شروع  $s_0$

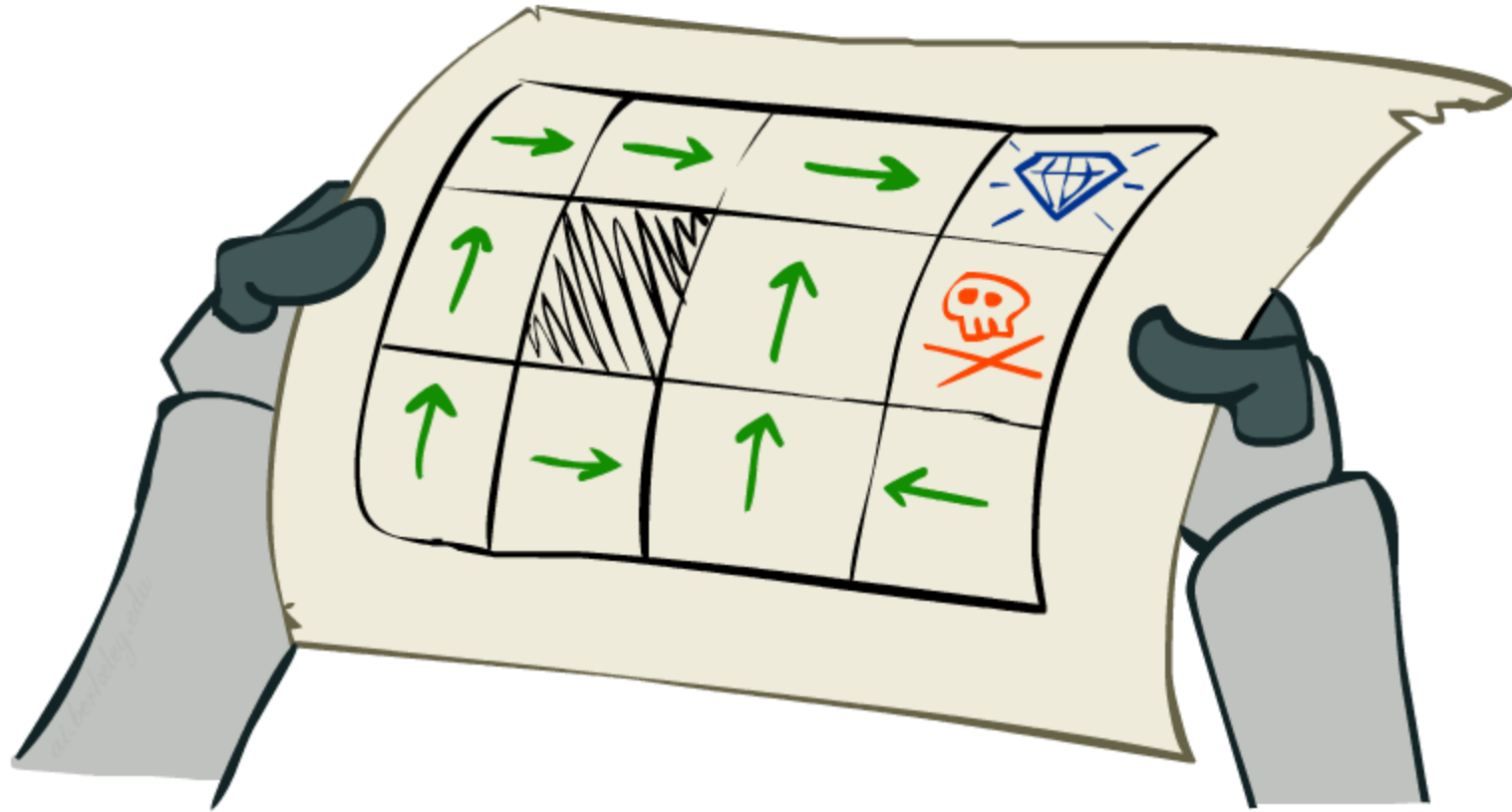
□ دو کمیت مهم MDP تا اینجا:

- سیاست (Policy): انتخاب یک عمل برای هر حالت
- سودمندی (Utility): مجموع (کاهش یافته) پاداش ها





# حل مسائل MDP



# کمیت‌های بهینه

■ ارزش (یا سودمندی) یک حالت  $s$

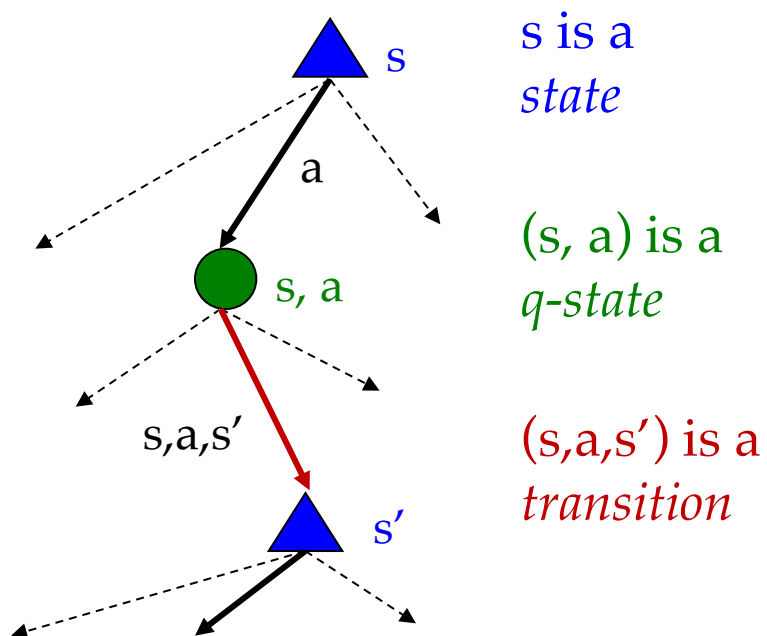
$V^*(s)$  = سودمندی مورد انتظار با شروع از  $s$  و بهینه عمل کردن.

■ ارزش (یا سودمندی) یک حالت  $q$

$Q^*(s,a)$  = سودمندی مورد انتظار با شروع از  $s$  و انتخاب عمل  $a$  و بهینه عمل کردن

■ سیاست بهینه

$\pi^*(s)$  = عمل بهینه در حالت  $s$

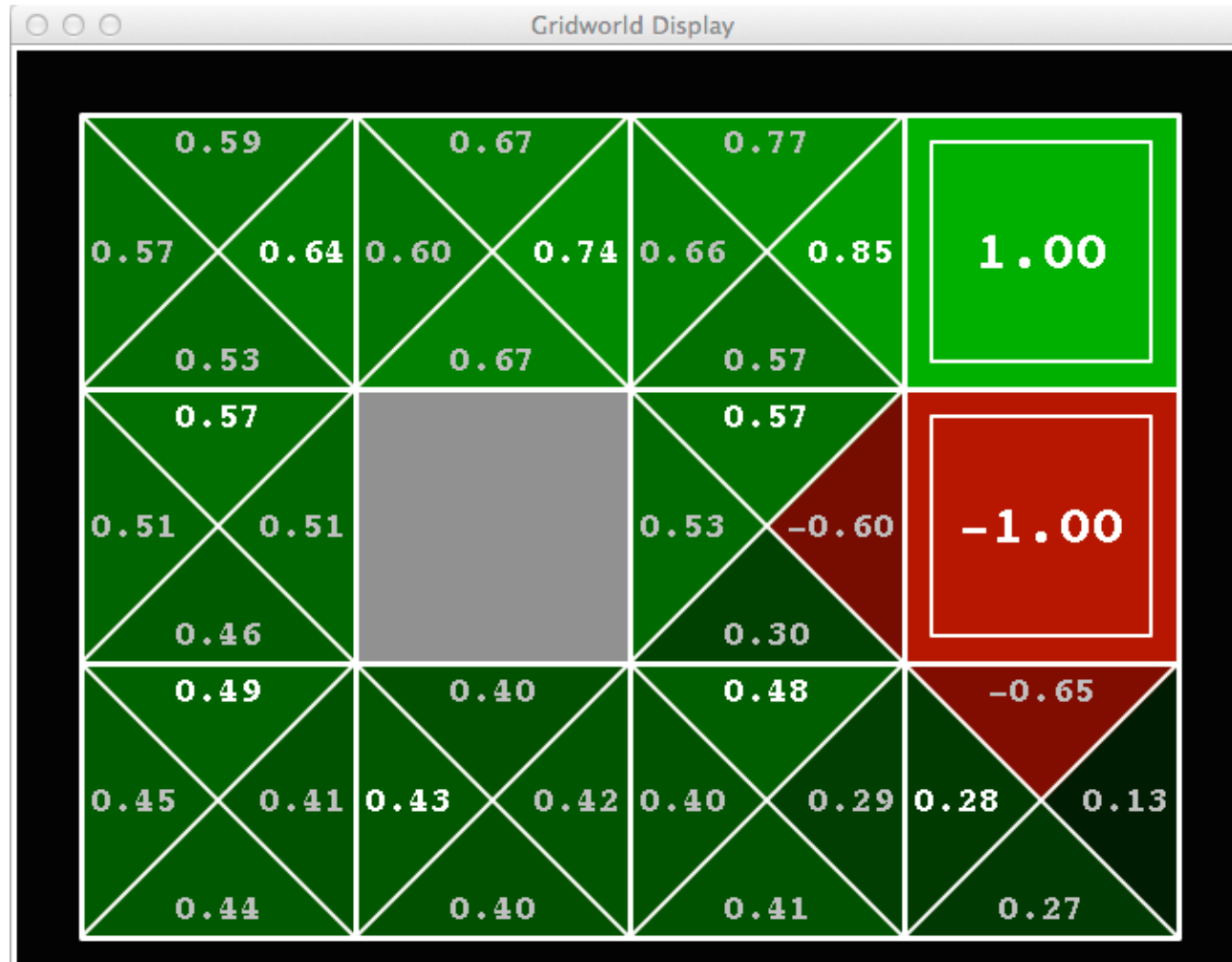


# V\* Values



- /۲ = نويز
- /۹ = ضريب کاهش
- = پاداش مرحله‌ای

# Q\* Values



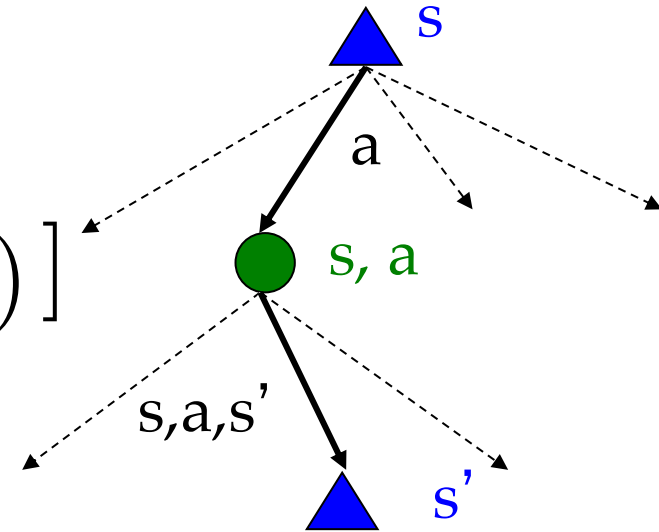
•/۲ = نويز  
 •/۹ = ضريب کاهش  
 • = پاداش مرحله‌ای

## محاسبه $V^*$ و $Q^*$

□ تعریف بازگشتی  $V^*$ :

$$V^*(s) = \max_a Q^*(s, a)$$

$$Q^*(s, a) = \sum_{s'} T(s, a, s') [R(s, a, s') + \gamma V^*(s')]$$

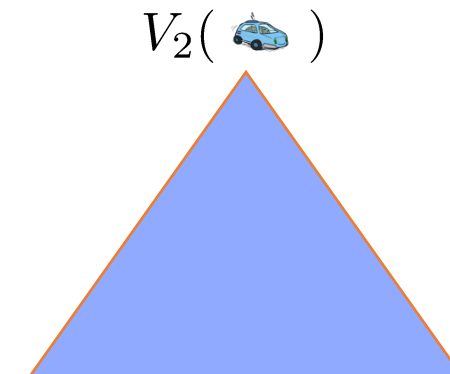
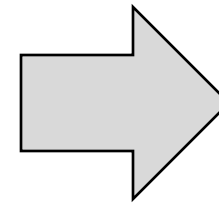
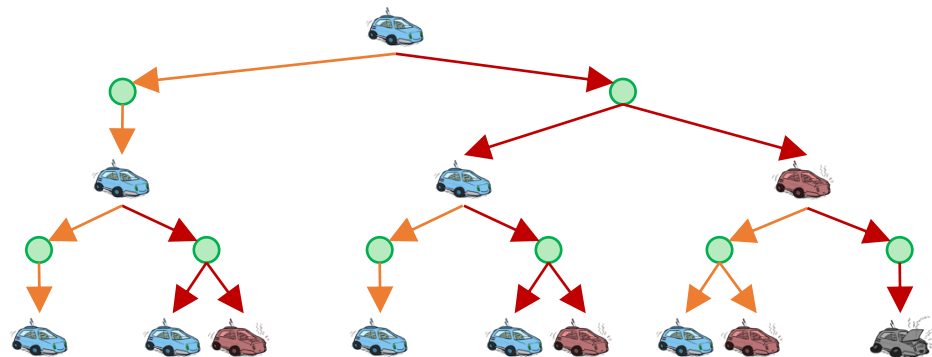
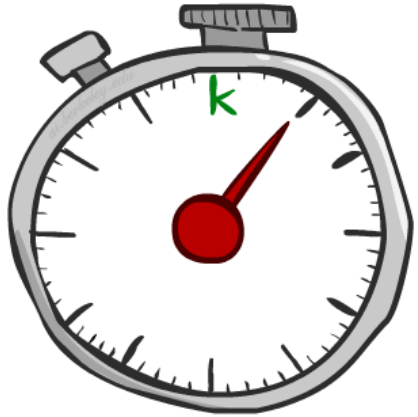


$$V^*(s) = \max_a \sum_{s'} T(s, a, s') [R(s, a, s') + \gamma V^*(s')]$$

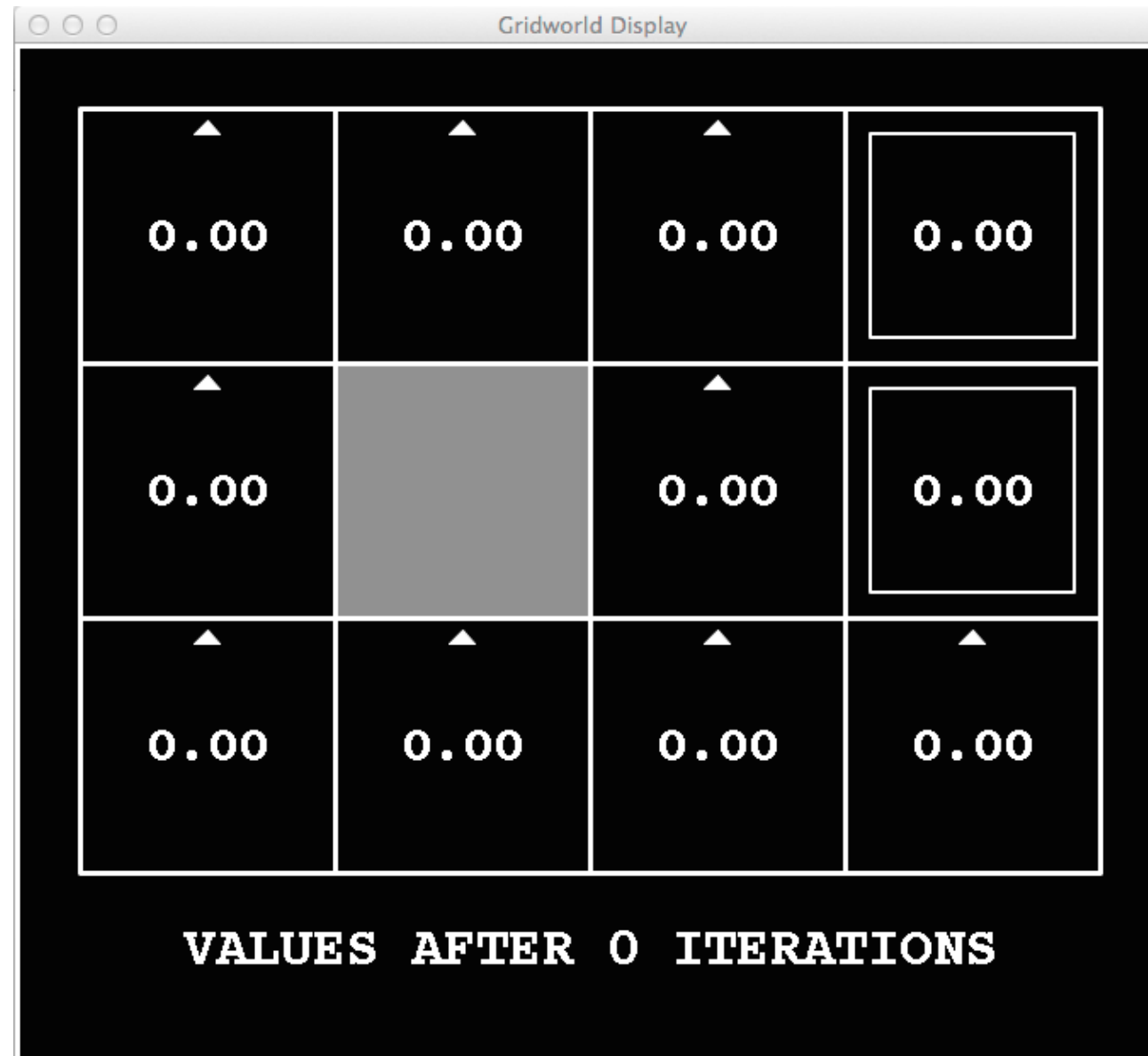
# ارزش‌های با زمان محدود

□ ایده‌ی کلیدی: ارزش‌های با زمان محدود!

□ تعریف  $V_k(s)$  به عنوان ارزش بهینه  $s$  اگر بازی پس از  $k$  مرحله‌ی دیگر خاتمه یابد.

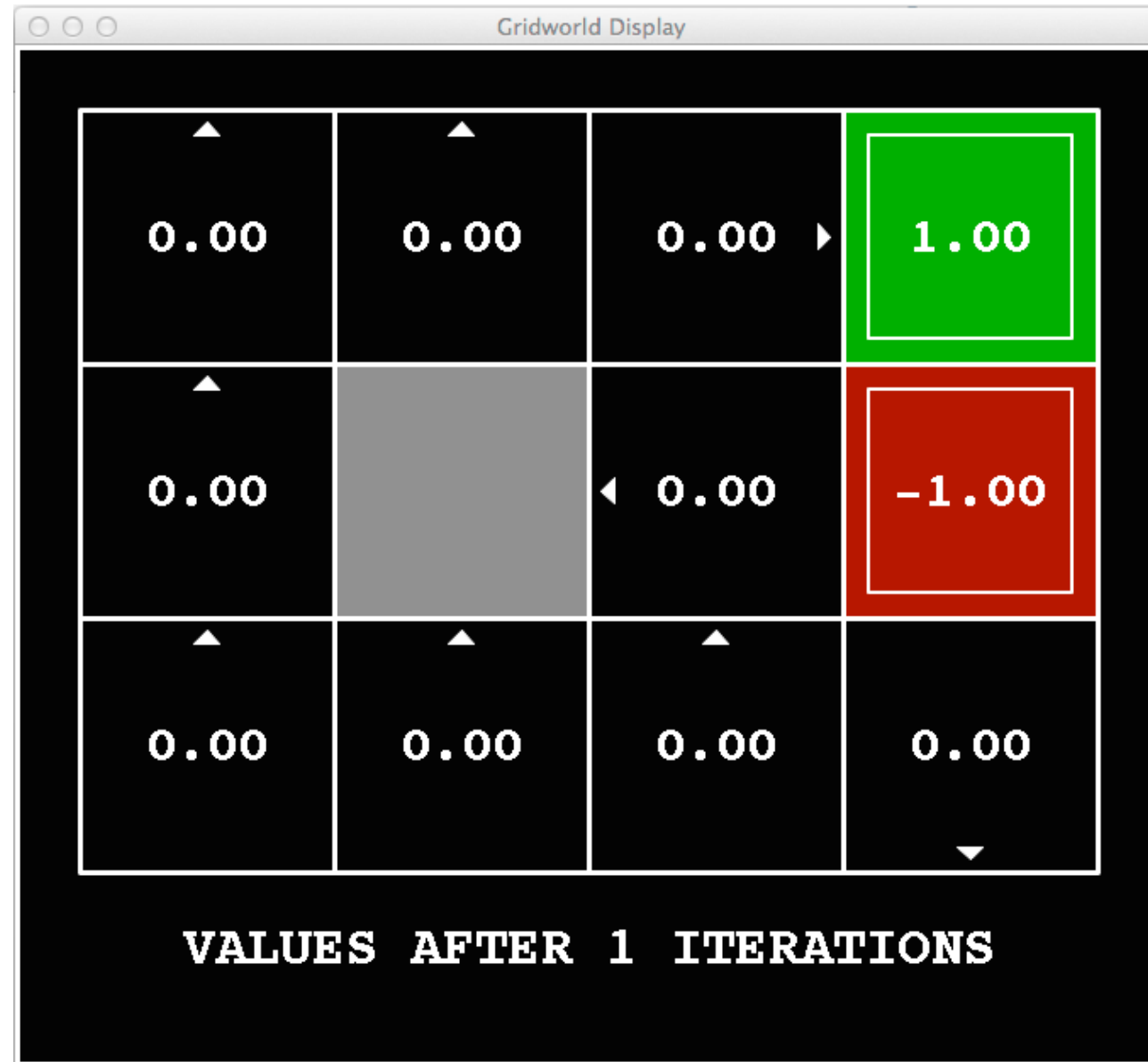


$k=0$



•/۲ = نويز  
•/۹ = ضريب کاهش  
• = پاداش مرحله‌ای

k=1



- / ۲ = نويز
- / ۹ = ضريب کاهش
- = پاداش مرحله‌ای



$$V^*(s) = \max_a \sum_{s'} T(s, a, s') [R(s, a, s') + \gamma V^*(s')]$$

$$V^{(2)}(s) = 0.8 \cdot (0 + 0.9 \cdot 1.00) + 0.1 \cdot (0 + 0.9 \cdot 0.00) + 0.1 \cdot (0 + 0.9 \cdot 0.00) = 0.72$$

k=2

For action = right



•/۲ = نويز  
•/۹ = ضريب کاهش  
• = پاداش مرحله‌ای

$$V^*(s) = \max_a \sum_{s'} T(s, a, s') [R(s, a, s') + \gamma V^*(s')]$$

$$V^{(3)}(s) = 0.8 \cdot (0 + 0.9 \cdot 1.00) + 0.1 \cdot (0 + 0.9 \cdot 0.72) + 0.1 \cdot (0 + 0.9 \cdot 0.00) = 0.7848$$

k=3

For action = right



•/۲ = نويز  
•/۹ = ضريب کاهش  
• = پاداش مرحله‌ای

k=4



•/۲ = نويز  
•/۹ = ضريب کاهش  
• = پاداش مرحله‌ای

k=5



•/۲ = نويز  
•/۹ = ضريب کاهش  
• = پاداش مرحله‌ای

k=6



•/۲ = نويز  
•/۹ = ضريب کاهش  
• = پاداش مرحله‌ای

k=7



•/۲ = نويز  
•/۹ = ضريب کاهش  
• = پاداش مرحله‌ای

k=8



•/۲ = نويز  
•/۹ = ضريب کاهش  
• = پاداش مرحله‌ای

k=9



•/۲ = نويز  
•/۹ = ضريب کاهش  
• = پاداش مرحله‌ای



k=10



•/۲ = نويز  
•/۹ = ضريب کاهش  
• = پاداش مرحله‌ای

k=11



•/۲ = نويز  
•/۹ = ضريب کاهش  
• = پاداش مرحله‌ای

k=12

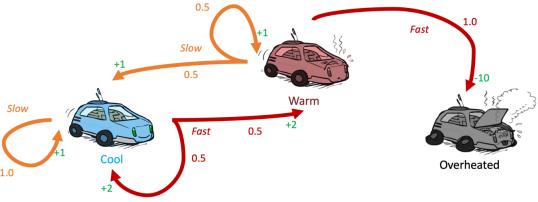


•/۲ = نويز  
•/۹ = ضريب کاهش  
• = پاداش مرحله‌ای

$k=100$

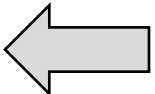


•/۲ = نويز  
•/۹ = ضريب کاهش  
• = پاداش مرحله‌ای

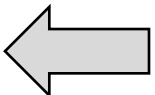


# محاسبه‌ی ارزش‌های با زمان محدود

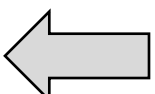
$V_4(\text{Cool}) \quad V_4(\text{Warm}) \quad V_4(\text{Overheated})$



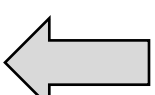
$V_3(\text{Cool}) \quad V_3(\text{Warm}) \quad V_3(\text{Overheated})$



$V_2(\text{Cool}) \quad V_2(\text{Warm}) \quad V_2(\text{Overheated})$



$V_1(\text{Cool}) \quad V_1(\text{Warm}) \quad V_1(\text{Overheated})$



$V_0(\text{Cool}) \quad V_0(\text{Warm}) \quad V_0(\text{Overheated})$

