

بِسْمِ اللَّهِ الرَّحْمَنِ الرَّحِيمِ



هوش مصنوعی

فصل ۱۷

# اتخاذ تصمیم‌های پیچیده

**Making Complex Decisions**

کاظم فولادی قلعه  
دانشکده مهندسی، دانشکدگان فارابی  
دانشگاه تهران

<http://courses.fouladi.ir/ai>

## انواع تصمیم

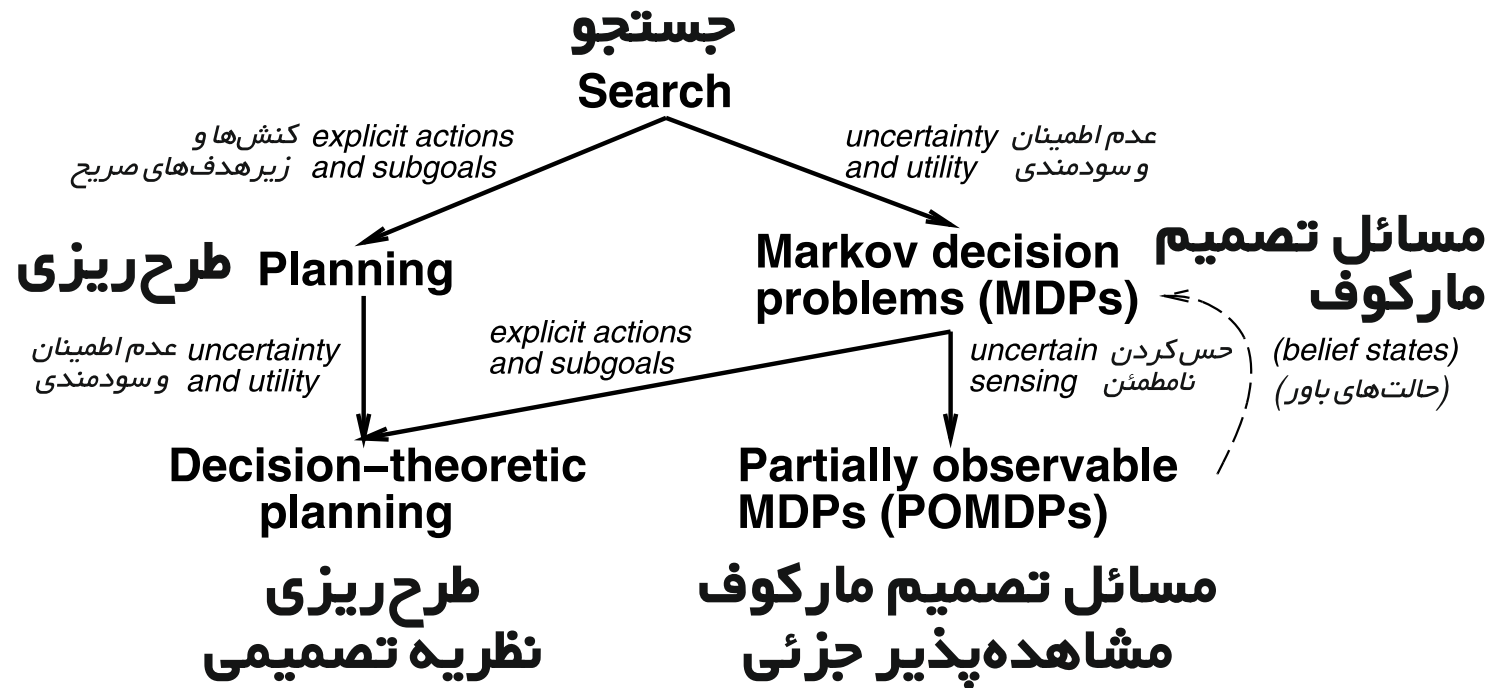
انواع تصمیم	
تصمیم پیچیده <i>Complex Decision</i>	تصمیم ساده <i>Simple Decision</i>
تصمیم‌های چندمرحله‌ای <i>Multi-Stage</i>	تصمیم‌های تک‌ضرب <i>One-Shot</i>
تصمیم بر روی یک دنباله از کنش‌ها	تصمیم بر روی یک کنش

## انواع تصمیم

## تصمیم‌های پیچیده



## مسائل تصمیم‌گیری

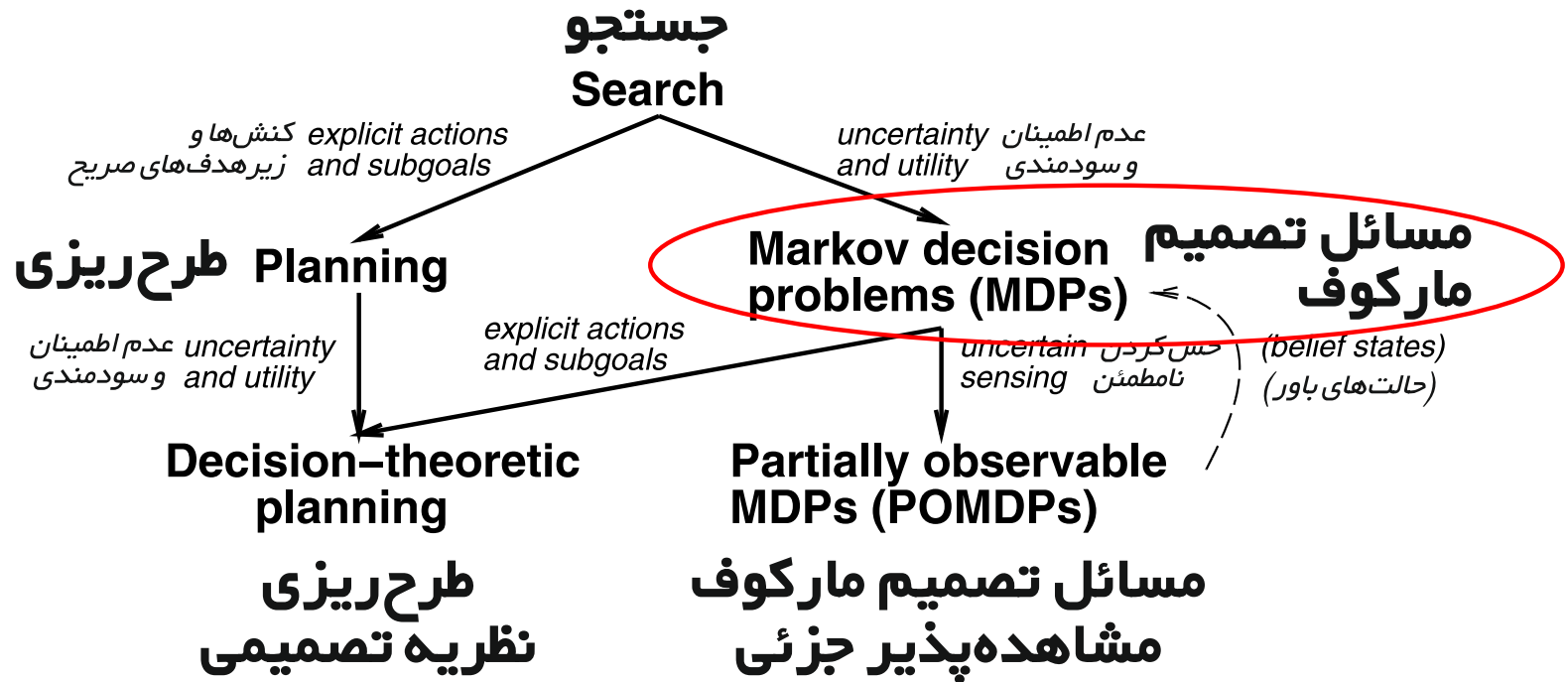


۱

# مسائل تصمیم‌گیری ترتیبی

## مسائل تصمیم‌گیری

مسائل تصمیم‌گیری مارکوف

MARKOV DECISION PROBLEMS (MDPs)

## مسائل / فرآیند تصمیم‌گیری مارکوف

مثال

MARKOV DECISION PROBLEMS / PROCESS (MDPs)

مؤلفه‌های تعریف یک MDP		
تابع پاداش <i>Reward Function</i>	مدل گذر <i>Transition Model</i>	حالت اولیه <i>Initial State</i>
$R(s)$ $R(s, a)$ $R(s, a, s')$	$T(s, a, s')$	$s_0$

States  $s \in S$ , actions  $a \in A$ Model  $T(s, a, s') \equiv P(s'|s, a)$  = probability that  $a$  in  $s$  leads to  $s'$ 

فرض می‌شود که گذرها مارکوف باشند، یعنی:

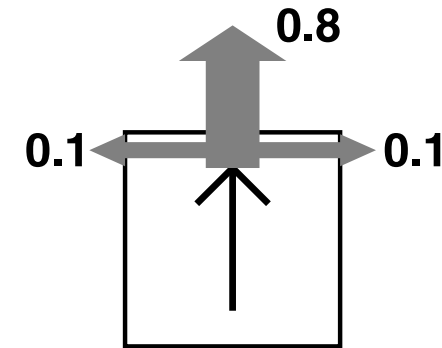
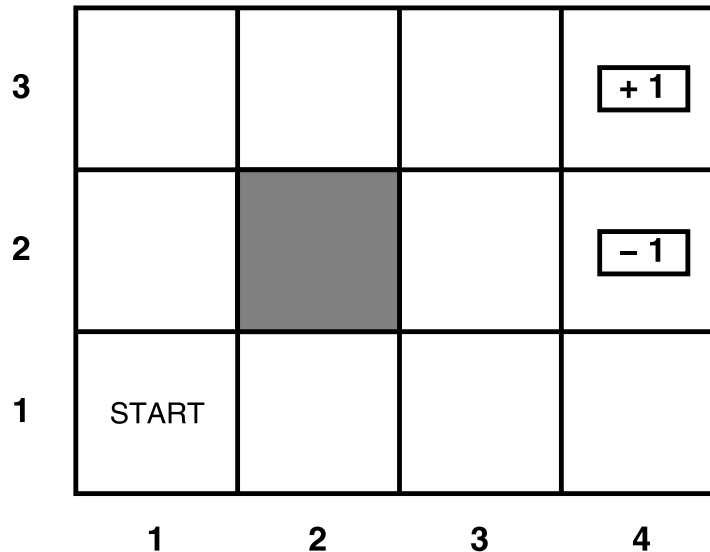
احتمال رسیدن از حالت فعلی به حالت دیگر، فقط به حالت فعلی وابسته است و به تاریخچه‌ی حالت‌های اخیر وابسته نیست.

## مسائل تصمیم‌گیری مارکوف

مثال

MARKOV DECISION PROBLEMS (MDPS)

حرکت از START و رسیدن به هدف +1



توزیع احتمال نتیجه‌ی کنش حرکت مستقیم

States  $s \in S$ , actions  $a \in A = \{\text{Right, Left, Down, Up}\}$  حالت‌ها و کنش‌هاModel  $T(s, a, s') \equiv P(s'|s, a)$  = probability that  $a$  in  $s$  leads to  $s'$  مدل

Reward function  $R(s)$  (or  $R(s, a)$ ,  $R(s, a, s')$ ) تابع پاداش

$$= \begin{cases} -0.04 & \text{جریمه / پناستی کوچک برای حالت‌های ناپایانی} \\ \pm 1 & \text{برای حالت‌های پایانی} \end{cases}$$

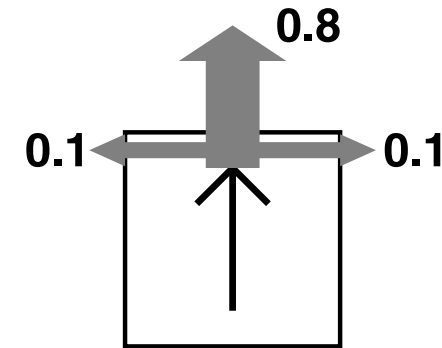
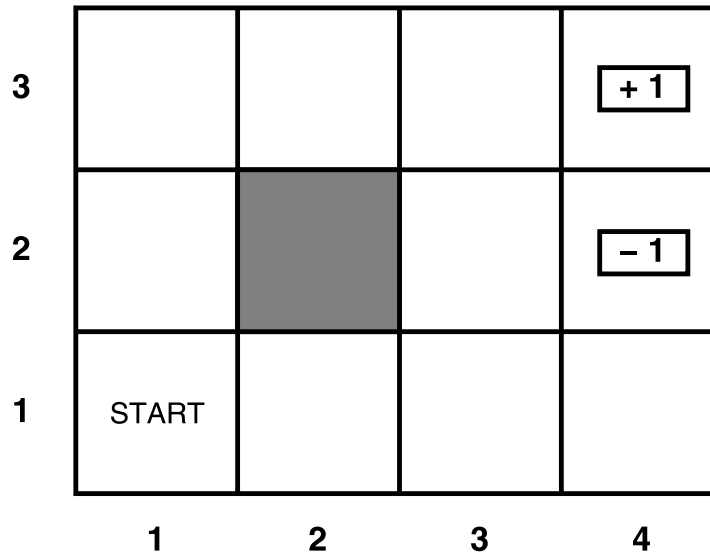


## مسائل تصمیم‌گیری مارکوف

مثال

## MARKOV DECISION PROBLEMS (MDPS)

حرکت از START و رسیدن به هدف +1



توزیع احتمال نتیجه‌ی کنش حرکت مستقیم

[Up, Up, Right, Right, Right]	راه‌حل برای محیط قطعی
?	راه‌حل برای محیط اتفاقی

وقتی محیط اتفاقی است، اثر کنش‌ها قطعی نخواهد بود.



یک ترتیب ثابت از کنش‌ها مسئله را حل نمی‌کند؛

(زیرا ممکن است عامل وارد حالتی شود که نتیجه‌ی کنش آن نبوده است)



در راه‌حل باید: کنش عامل به‌ازای همه‌ی حالت‌هایی که ممکن است به آنها برسد، تعیین شده باشد (سیاست).

## مسائل تصمیم‌گیری مارکوف

سیاست

POLICY

$$\pi(s)$$

کنش پیشنهادی برای هر حالت ممکن  $s$ سیاست  
*Policy*

وقتی عامل دارای یک سیاست کامل باشد،  
بدون توجه به نتیجه‌ی هر کنش، همیشه می‌داند برای مرحله‌ی بعدی باید چه کاری انجام دهد.

$$\pi^*(s)$$

بهترین کنش برای هر حالت ممکن  $s$ سیاست بهینه  
*Optimal Policy*

سیاست بهینه سیاستی است که به بالاترین مقدار امید سودمندی برسد.

## مسائل تصمیم‌گیری مارکوف

مسائل جستجو در مقابل مسائل تصمیم‌گیری مارکوف

### SEARCH PROBLEMS VS. MDPs

#### مسئله‌ی تصمیم‌گیری مارکوف

*Markov Decision Problem (MDP)*

هدف: یافتن یک **سیاست** بهینه  $\pi(s)$

یعنی: بهترین کنش برای هر حالت ممکن  $s$

(زیرا نمی‌توانیم پیش‌بینی کنیم که  
هر کنش قطعاً به کدام حالت منجر می‌شود)

#### مسئله‌ی جستجو

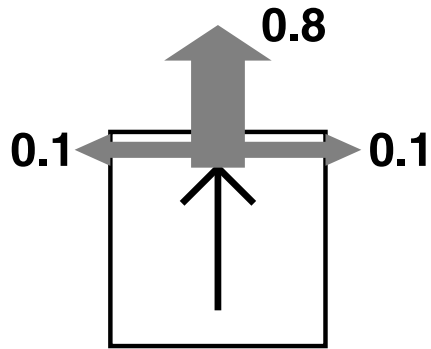
*Search Problem*

هدف: یافتن یک **دنباله‌ی** بهینه

## مسائل تصمیم‌گیری مارکوف

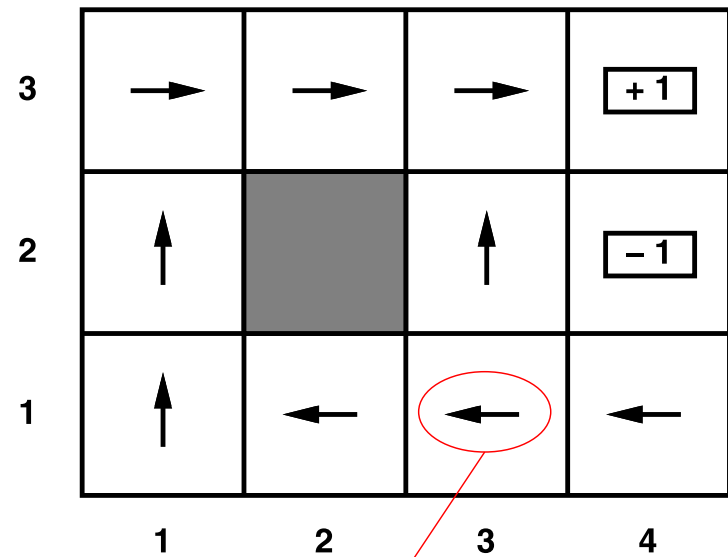
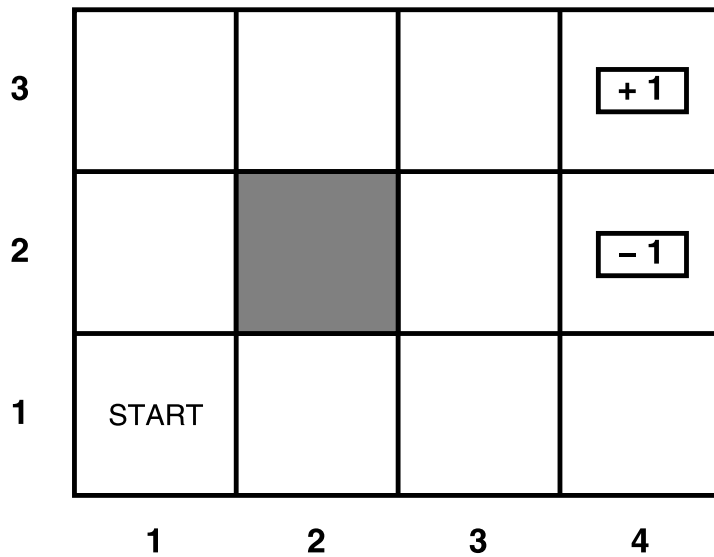
حل کردن مسائل تصمیم‌گیری مارکوف: مثال

### SOLVING MDPs



توزیع احتمال نتیجه‌ی کنش حرکت مستقیم

Optimal policy when state penalty  $R(s)$  is  $-0.04$ :



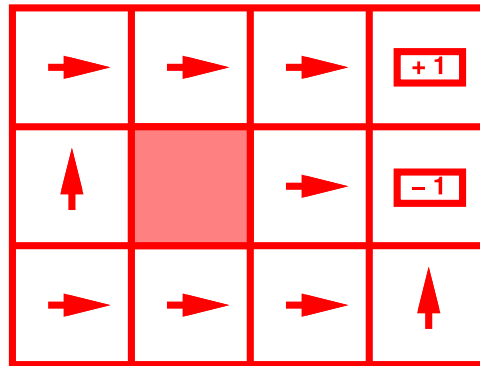
تصمیم محافظه‌کارانه: در مقابل راه میان‌بر و ریسکی ورود به 1- راه طولانی‌تر را انتخاب می‌کند.

## مسائل تصمیم‌گیری مارکوف

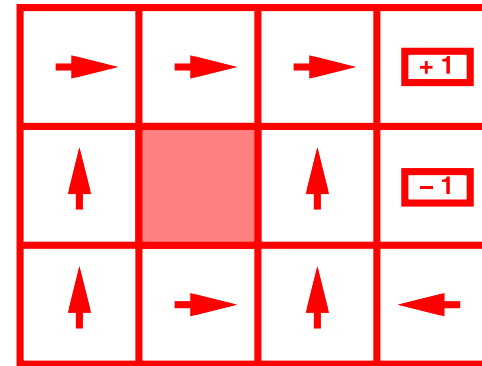
### ریسک و پاداش

#### RISK AND REWARD

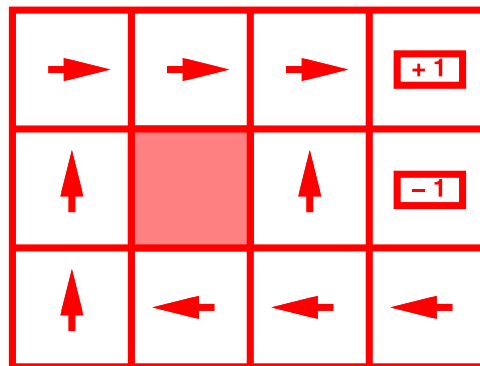
توازن دقیق ریسک و پاداش، از خصوصیات MDP است.



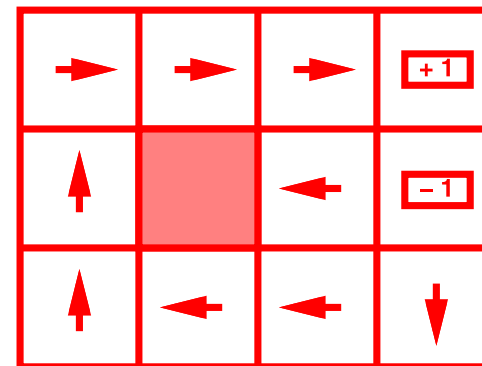
$$r = [-\infty : -1.6284]$$



$$r = [-0.4278 : -0.0850]$$



$$r = [-0.0480 : -0.0274]$$



$$r = [-0.0218 : 0.0000]$$

سیاست‌های بهینه با بازه‌های مختلف برای پاداش تغییر می‌کند.

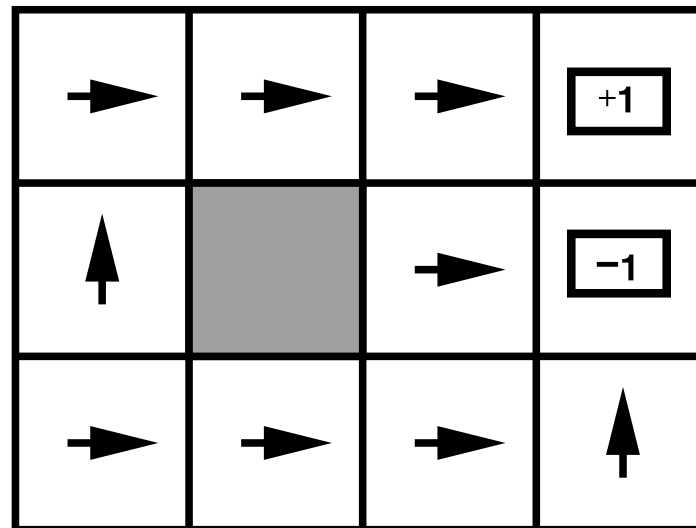
توازن بین تغییرات پاداش و ریسک، به مقدار پاداش حالت‌های غیرپایانی بستگی دارد.

## مسائل تصمیم‌گیری مارکوف

توازن ریسک و پاداش: مثال ۱

### RISK AND REWARD

سیاست بهینه برای پاداش‌های منفی (جریمه) بزرگ



$$R(s) < -1.6284$$

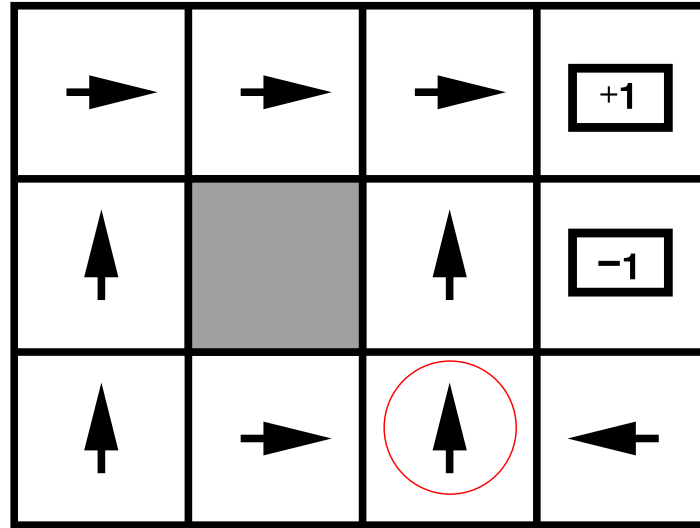
شرایط آن قدر سخت می‌شود که عامل برای یافتن نزدیک‌ترین خروجی تلاش می‌کند.  
حتی اگر خروجی دارای مقدار -1 باشد.

## مسائل تصمیم‌گیری مارکوف

توازن ریسک و پاداش: مثال ۲

RISK AND REWARD

سیاست بهینه برای پاداش‌های منفی (جریمه) متوسط



$$-0.4278 < R(s) < -0.0850$$

زندگی برای عامل ناخوشایند است و عامل کوتاه‌ترین مسیر برای رسیدن به +1 را انتخاب می‌کند؛  
و ریسک (خطر) افتادن ناگهانی در -1 را نیز می‌پذیرد.

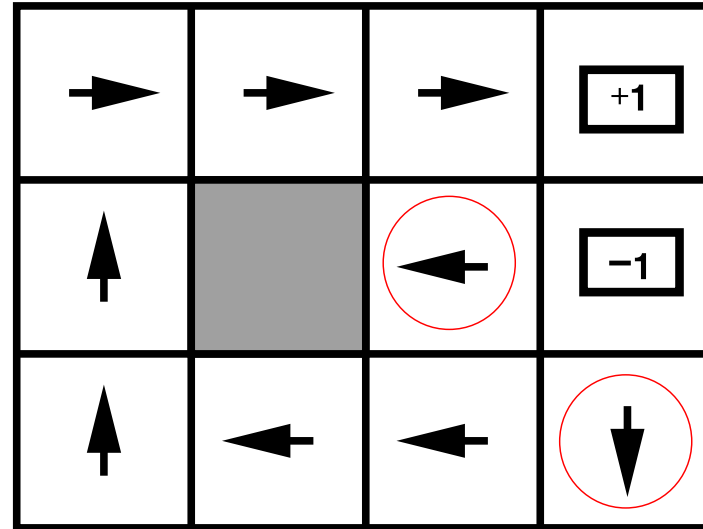
(به‌طور خاص برای زمانی که عامل راه میان‌بری را از (3,1) انتخاب کند.)

## مسائل تصمیم‌گیری مارکوف

توازن ریسک و پاداش: مثال ۳

### RISK AND REWARD

سیاست بهینه برای پاداش‌های منفی (جریمه) کوچک



$$-0.0221 < R(s) < 0$$

زندگی عامل با کمی افسردگی همراه است:  
در این حالت سیاست بهینه هیچ ریسکی را نمی‌پذیرد.

(در (4,1) و (3,2) عامل سعی می‌کند با دور شدن از حالت 1- به‌طور ناگهانی در حالت 1- گرفتار نشود؛ حتی اگر به دیوار بخورد.)

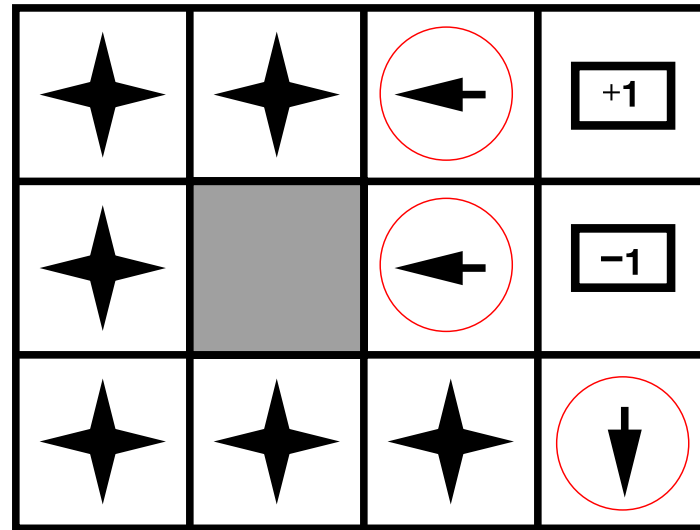


## مسائل تصمیم‌گیری مارکوف

توازن ریسک و پاداش: مثال ۴

### RISK AND REWARD

سیاست بهینه برای پاداش‌های مثبت



$$R(s) > 0$$

زندگی عامل لذت‌بخش است:  
عامل از هر دو حالت خروج امتناع می‌کند.

(تازمانی که کنش‌ها در (4,1) و (3,2) و (3,3) مطابق شکل فوق باشد، هر سیاستی بهینه است و عامل به خاطر عدم ورود به حالت‌های پایانی، در کل پاداشی نامتناهی به دست می‌آورد.)

## سودمندی دنباله‌های حالت

### UTILITY OF STATE SEQUENCES

برای تعیین ترجیح‌ها بین دنباله‌های حالت‌ها

ترجیح‌های ایستان را بر روی دنباله‌های پاداش در نظر می‌گیریم:  
**stationary preferences**

$$[r, r_0, r_1, r_2, \dots] \succ [r, r'_0, r'_1, r'_2, \dots] \Leftrightarrow [r_0, r_1, r_2, \dots] \succ [r'_0, r'_1, r'_2, \dots]$$

**قضیه:**

با فرض ایستان بودن  
 برای انتساب سودمندی به دنباله‌های حالت  
 (ترکیب پاداش‌ها در طول زمان) فقط دو راه وجود دارد:

- |                          |  |
|--------------------------|--|
| تابع سودمندی جمعی        | 1) <i>Additive</i> utility function:<br>$U([s_0, s_1, s_2, \dots]) = R(s_0) + R(s_1) + R(s_2) + \dots$   |
| تابع سودمندی تخفیف‌یافته | 2) <i>Discounted</i> utility function:<br>$U([s_0, s_1, s_2, \dots]) = R(s_0) + \gamma R(s_1) + \gamma^2 R(s_2) + \dots$ <p>where <math>\gamma</math> is the <u>discount factor</u> فاکتور تخفیف</p> |

## سودمندی حالت‌ها

UTILITY OF STATES

امید مجموع (تخفیف‌یافته) پاداش‌ها (تا رسیدن به پایان) با فرض کنش‌های بهینه  $U(s)$

ارزش یک حالت  
*Value of a state*

سودمندی یک حالت  
*Utility of a state*

با داشتن سودمندی حالت‌ها، انتخاب بهترین کنش صرفاً یک MEU است:  
امید سودمندی مابعد‌های بی‌واسطه را ماکزیمم کنید.

$$\pi^*(s) = \operatorname{argmax}_{a \in A(s)} \sum_{s'} P(s' | s, a) U(s')$$

## سودمندی حالتها

مثال

UTILITY OF STATES

Utility of States:

3	0.812	0.868	0.912	$+1$
2	0.762		0.660	$-1$
1	0.705	0.655	0.611	0.388
	1	2	3	4

Optimal policy when state penalty  $R(s)$  is  $-0.04$ :

3	→	→	→	$+1$
2	↑		↑	$-1$
1	↑	←	←	←
	1	2	3	4

## سودمندی حالتها

برخورد با مشکل طول عمر نامتناهی

### UTILITY OF STATES

Problem: infinite lifetimes  $\Rightarrow$  additive utilities are infinite

- 1) **Finite horizon**: termination at a *fixed time*  $T$   
 $\Rightarrow$  **nonstationary** policy:  $\pi(s)$  depends on time left
- 2) **Absorbing state(s)**: w/ prob. 1, agent eventually “dies” for any  $\pi$   
 $\Rightarrow$  expected utility of every state is finite
- 3) **Discounting**: assuming  $\gamma < 1$ ,  $R(s) \leq R_{\max}$ ,

$$U([s_0, \dots s_\infty]) = \sum_{t=0}^{\infty} \gamma^t R(s_t) \leq R_{\max} / (1 - \gamma)$$

Smaller  $\gamma \Rightarrow$  shorter horizon

- 4) Maximize **system gain** = average reward per time step

Theorem: optimal policy has constant gain after initial transient

E.g., taxi driver's daily scheme cruising for passengers

۲

تکرار  
ارزش

## برنامه ریزی پویا

معادله‌ی بلمن

### DYNAMIC PROGRAMMING: THE BELLMAN EQUATION

Definition of utility of states leads to a simple relationship among utilities of neighboring states:

expected sum of rewards

= current reward

+  $\gamma \times$  expected sum of rewards after taking best action

Bellman equation (1957):

$$U(s) = R(s) + \gamma \max_a \sum_{s'} U(s') T(s, a, s')$$

$$U(1, 1) = -0.04$$

$$+ \gamma \max \{ 0.8U(1, 2) + 0.1U(2, 1) + 0.1U(1, 1), \\ 0.9U(1, 1) + 0.1U(1, 2) \\ 0.9U(1, 1) + 0.1U(2, 1) \\ 0.8U(2, 1) + 0.1U(1, 2) + 0.1U(1, 1) \}$$

up  
left  
down  
right

One equation per state =  $n$  **nonlinear** equations in  $n$  unknowns

## تکرار ارزش

## الگوریتم

VALUE ITERATION ALGORITHM

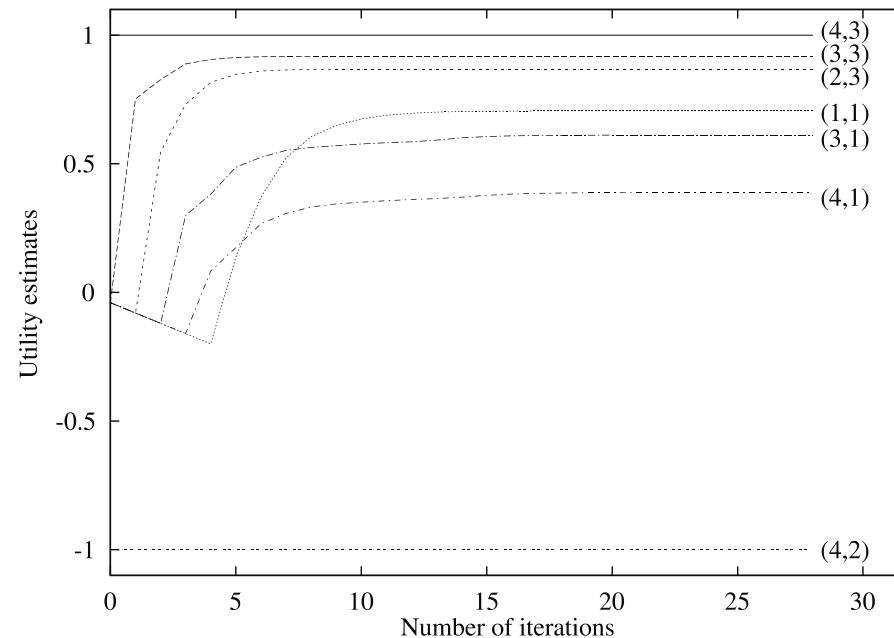
Idea: Start with arbitrary utility values

Update to make them **locally consistent** with Bellman eqn.

Everywhere locally consistent  $\Rightarrow$  global optimality

Repeat for every  $s$  simultaneously until “no change”

$$U(s) \leftarrow R(s) + \gamma \max_a \sum_{s'} U(s') T(s, a, s') \quad \text{for all } s$$





## تکرار ارزش

قضیه‌ی همگرایی

### VALUE ITERATION ALGORITHM: CONVERGENCE

Define the **max-norm**  $\|U\| = \max_s |U(s)|$ ,  
so  $\|U - V\|$  = maximum difference between  $U$  and  $V$

Let  $U^t$  and  $U^{t+1}$  be successive approximations to the true utility  $U$

**Theorem:** For any two approximations  $U^t$  and  $V^t$

$$\|U^{t+1} - V^{t+1}\| \leq \gamma \|U^t - V^t\|$$

I.e., any distinct approximations must get closer to each other  
so, in particular, any approximation must get closer to the true  $U$   
and value iteration converges to a unique, stable, optimal solution

**Theorem:** if  $\|U^{t+1} - U^t\| < \epsilon$ , then  $\|U^{t+1} - U\| < 2\epsilon\gamma/(1 - \gamma)$

I.e., once the change in  $U^t$  becomes small, we are almost done.

MEU policy using  $U^t$  may be optimal long before convergence of values

۳

تکرار  
سیاست

## تکرار سیاست

## الگوریتم

POLICY ITERATION ALGORITHM

Howard, 1960: search for optimal policy and utility values simultaneously

Algorithm:

$\pi \leftarrow$  an arbitrary initial policy

repeat until no change in  $\pi$

    compute utilities given  $\pi$

    update  $\pi$  as if utilities were correct (i.e., local MEU)

To compute utilities given a fixed  $\pi$  (**value determination**):

$$U(s) = R(s) + \gamma \sum_{s'} U(s') T(s, \pi(s), s') \quad \text{for all } s$$

i.e.,  $n$  simultaneous **linear** equations in  $n$  unknowns, solve in  $O(n^3)$

## تکرار سیاست

شکل اصلاح شده

### MODIFIED POLICY ITERATION ALGORITHM

Policy iteration often converges in few iterations, but each is expensive

Idea: use a few steps of value iteration (but with  $\pi$  fixed) starting from the value function produced the last time to produce an approximate value determination step.

Often converges much faster than pure VI or PI

Leads to much more general algorithms where Bellman value updates and Howard policy updates can be performed locally in any order

**Reinforcement learning** algorithms operate by performing such updates based on the observed transitions made in an initially unknown environment

# ۴

## MDPهای مشاهده‌پذیر جزئی (POMDPs)

## MDPهای مشاهده‌پذیر جزئی (POMDPs)

POMDP has an **observation model**  $O(s, e)$  defining the probability that the agent obtains evidence  $e$  when in state  $s$

Agent does not know which state it is in

$\Rightarrow$  makes no sense to talk about policy  $\pi(s)!!$

**Theorem** (Astrom, 1965): the optimal policy in a POMDP is a function  $\pi(b)$  where  $b$  is the **belief state** (probability distribution over states)

Can convert a POMDP into an MDP in belief-state space, where

$T(b, a, b')$  is the probability that the new belief state is  $b'$  given that the current belief state is  $b$  and the agent does  $a$ .  
I.e., essentially a filtering update step

## MDPهای مشاهده‌پذیر جزئی (POMDPs)

ویژگی‌ها

Solutions automatically include information-gathering behavior

If there are  $n$  states,  $b$  is an  $n$ -dimensional real-valued vector  
 $\Rightarrow$  solving POMDPs is very (actually, PSPACE-) hard!

The real world is a POMDP (with initially unknown  $T$  and  $O$ )

# ۵

تصمیم‌هایی  
با عامل‌های  
چندگانه:  
نظریه‌ی  
بازی



## تصمیم‌هایی با عامل‌های چندگانه

نظریه‌ی بازی

DECISIONS WITH MULTIPLE AGENTS: GAME THEORY

هوش مصنوعی

اتخاذ تصمیم‌های پیچیده

۶

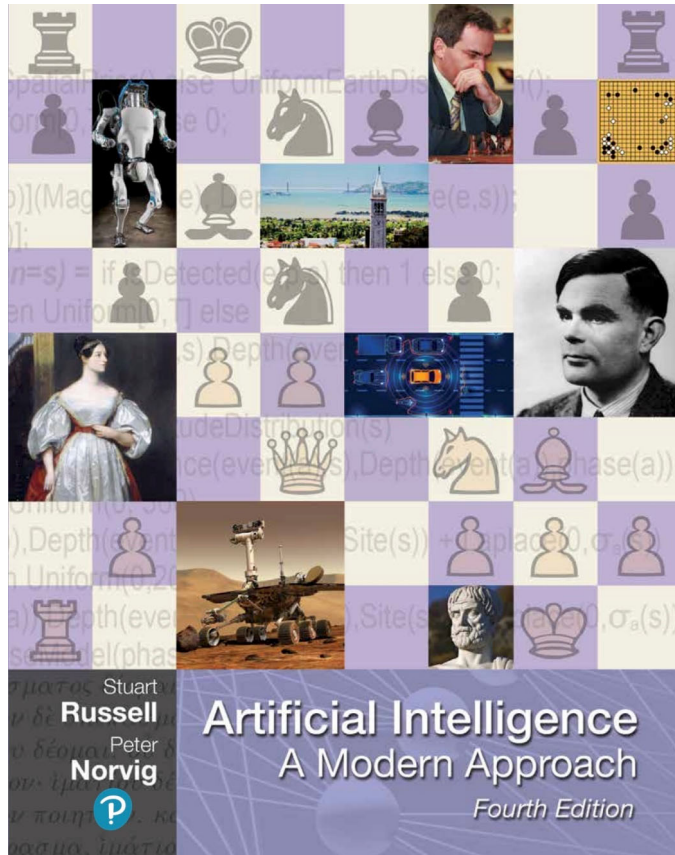
طراحی  
مکانیزم

## طراحی مکانیزم

### MECHANISM DESIGN

# ۷

منابع،  
مطالعه،  
تکلیف



Stuart Russell and Peter Norvig,  
**Artificial Intelligence: A Modern Approach**,  
 4<sup>th</sup> Edition, Prentice Hall, 2020.

## Chapter 17

## CHAPTER 17

### MAKING COMPLEX DECISIONS

*In which we examine methods for deciding what to do today, given that we may face another decision tomorrow.*

In this chapter, we address the computational issues involved in making decisions in a stochastic environment. Whereas Chapter 16 was concerned with one-shot or episodic decision problems, in which the utility of each action's outcome was well known, we are concerned here with **sequential decision problems**, in which the agent's utility depends on a sequence of decisions. Sequential decision problems incorporate utilities, uncertainty, and sensing, and include search and planning problems as special cases. Section 17.1 explains how sequential decision problems are defined, and Section 17.2 describes methods for solving them to produce behaviors that are appropriate for a stochastic environment. Section 17.3 covers **multi-armed bandit** problems, a specific and fascinating class of sequential decision problems that arise in many contexts. Section 17.4 explores decision problems in partially observable environments and Section 17.5 describes how to solve them.

#### 17.1 Sequential Decision Problems

Suppose that an agent is situated in the  $4 \times 3$  environment shown in Figure 17.1(a). Beginning in the start state, it must choose an action at each time step. The interaction with the environment terminates when the agent reaches one of the goal states, marked +1 or -1. Just as for search problems, the actions available to the agent in each state are given by  $ACTIONS(s)$ , sometimes abbreviated to  $A(s)$ ; in the  $4 \times 3$  environment, the actions in every state are *Up*, *Down*, *Left*, and *Right*. We assume for now that the environment is **fully observable**, so that the agent always knows where it is.

If the environment were deterministic, a solution would be easy: [*Up*, *Up*, *Right*, *Right*, *Right*]. Unfortunately, the environment won't always go along with this solution, because the actions are unreliable. The particular model of stochastic motion that we adopt is illustrated in Figure 17.1(b). Each action achieves the intended effect with probability 0.8, but the rest of the time, the action moves the agent at right angles to the intended direction. Furthermore, if the agent bumps into a wall, it stays in the same square. For example, from the start square (1,1), the action *Up* moves the agent to (1,2) with probability 0.8, but with probability 0.1, it moves right to (2,1), and with probability 0.1, it moves left, bumps into the wall, and stays in (1,1). In such an environment, the sequence [*Up*, *Up*, *Right*, *Right*, *Right*] goes up around the barrier and reaches the goal state at (4,3) with probability  $0.8^5 = 0.32768$ . There is also a small chance of accidentally reaching the goal by going the other way around with probability  $0.1^4 \times 0.8$ , for a grand total of 0.32776. (See also Exercise 17.MDPX.)

Sequential decision problem