



دانشگاه صنعتی شریف

دانشکده مهندسی کامپیوتر

## هوش مصنوعی

بهار ۱۴۰۰

استاد: محمدحسین رهبان

گردآورندگان: نیما فتاحی - سروش وفایی تبار

### Markov Decision Process

مهلت ارسال: ۱۱ خرداد

تمرین هفتم بخش اول

- مهلت ارسال پاسخ تا ساعت ۲۳:۵۹ روز مشخص شده است.
- همکاری و همفکری شما در انجام تمرین مانعی ندارد اما پاسخ ارسالی هر کس حتما باید توسط خود او نوشته شده باشد.
- در صورت همفکری و یا استفاده از هر منابع خارج درسی، نام همفکران و آدرس منابع مورد استفاده برای حل سوال مورد نظر را ذکر کنید.
- لطفا تصویری واضح از پاسخ سوالات نظری بارگذاری کنید. در غیر این صورت پاسخ شما تصحیح نخواهد شد.

### سوالات نظری (۱۰۰ نمره)

۱. (۵۰ نمره) همانطور که می دانید در الگوریتم value iteration پس از هر تکرار شما می توانید بهترین سیاست (policy) در آن لحظه را پیدا کنید.  
حال اگر در value iteration بهترین سیاست در مرحله  $i$ م برابر با بهترین سیاست در مرحله  $i+1$ م باشد آیا ممکن است که در ادامه سیاست ما تغییر کند؟  
در صورت درستی عبارت بالا اثبات کنید در غیر این صورت مثال نقضی بیاورید.
۲. (۵۰ نمره) یک MDP با سه حالت و دو کنش ( $L, R$ ) در نظر بگیرید که احتمال انتقالات به صورت زیر است:

Action L:

	State 1	State 2	State 3
In state 1:	0	1/4	3/4
In state 2:	3/4	0	1/4
In state 3:	1/4	3/4	0

Action R:

	State 1	State 2	State 3
In state 1:	0	3/4	1/4
In state 2:	1/4	0	3/4
In state 3:	3/4	1/4	0

امتیاز حالت ۲ برابر ۱ و بقیه ی حالت ها ۰ است. ( $\gamma = 0.5$ )

- (a)  $V(s)$  تخمینی را به کمک روش Value Iteration برای دو مرحله محاسبه نمایید. (حالات اولیه را صفر در نظر بگیرید)
- (b) با در نظر گرفتن  $V_1 = V_3 = 0.5$  و  $V_2 = 0.25$  بهترین کنش برای حالت اول را محاسبه کنید و توضیح دهید.