



## سوال ۱: (۲۰ نمره)

مشخص کنید کدام یک از گزاره های زیر برقرار است:

- (آ) اضافه کردن مقدار ثابتی به تابع پاداش، policy بهینه را تغییر نمی دهد.
- (ب) در هر MDP متناهی، حداقل یک policy بهینه deterministic وجود دارد.
- (ج) زمانی که Value Iteration همگرا می شود، همواره به policy بهینه میل می کند.

## سوال ۲: (۴۰ نمره)

در یک MDP به ازای یک policy مشخص مانند  $\pi$ ، گذار زیر بین حالت ها و پاداش آن ها نمایش داده شده است. حالت  $s_3$  حالت پایانی است، به این معنا که با وارد شدن به آن، اپیزود خاتمه میابد. فرض کنید  $\gamma = 1$  است. معادلات bellman را برای value حالت ها نوشته، با حل آن ها  $v_\pi(s_i)$  را به ازای  $i = 0, 1, 2$  بدست آوردید. (توجه:  $v_\pi(s_3) = 0$  فرض شود).

