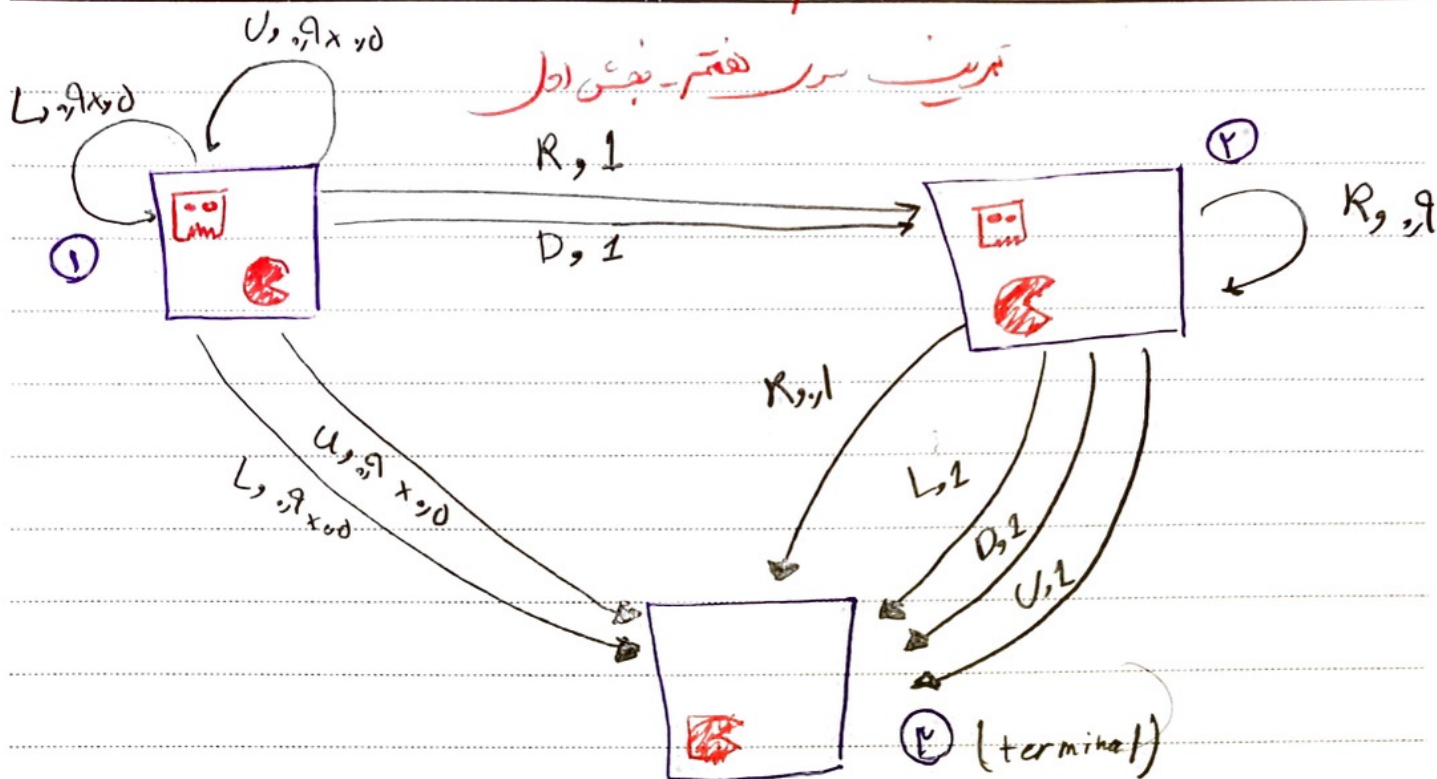


به نام خدا



* توجه ۱: در حالت ۱، روح با احتمال ۰.۵ به حالت ۲ و با احتمال ۰.۵ به پایت حرکت می‌کند. زیرا در هر دو حالت به یک دست یک خانه نزدیک تر می‌شود.

* توجه ۲: تارن حالت‌ها با حالاتی که رسم نشده‌اند، برای س فاصله روح به یک دست از پایداری است.

(I) Optimal Policy به شکل زیر تعریف می‌شود:

$$\forall s \in S : V^{\pi^*}(s) \geq V^{\pi}(s) \quad (\pi^* : \text{Optimal Policy})$$

برای باید گفتیم که معادله bellman-opt. ۱ جواب دارد
۲) یکی از این جواب‌ها از باقی مدارها بزرگتر است

۱) معادله bellman-optimization جواب دارد: صحت حرفی که در کلاس راجع به

Fixed Point زده شد، این را مل هاشد روش نوبت راضون ها به جواب بگردان شوند

۲) یک از جواب ها از بقی جواب ها بزرگتر است؛

طبق چیزی که در Policy iteration دیدیم، داریم؛

$$\pi_{i+1}(s) = \arg \max_a \sum_{s'} T(s, a, s') [R(s, a, s') + \gamma V^{\pi_i}(s')]$$

که چون در هر مرحله داریم $\arg \max$ بگیریم، قطعاً $V^{\pi_{i+1}} \geq V^{\pi_i}$ به ازای هر s ها.

← حال به دلیل اینکه $a \in A$ ، و دانه A (یا $|A|$) محدود است، این بهر

در بهترین حالت، از بدترین π به π^* می رسد و طی $|A|$ مرحله Policy iteration

حالا به π^* می رسیم که از هر سیاست دیگر بزرگتر است.

[مراجعه: [towards data science.com](https://towardsdatascience.com/)]

ب) طبق همان مدل Policy iteration که در بخش قبل معرفی کردیم، چقدر داریم برابر

$a = \arg \max$ بگیریم، پس در هر مرحله از این iteration به سیاست deterministic

داریم که $a = \pi_i(s)$ و در مرحله آخر که به π^* می رسیم $a = \pi^*(s)$

که از هر سیاست ها بهر سیاست deterministic قوی تر است.