

(الف) هر MDP نامی توان به صورت یک expectimax دیده می‌شود. خواست می‌شود که در فیلد action ما را به یک طرف
چسب می‌دهد که در اینجا با توجه به احتمالات خواست می‌شود که یک است خاص می‌رویم و می‌تواند تفاوت آن با utility
باشد. نکته در expectimax این است که در هر گام ما دو نوع داریم: یکی که در اینجا action است و یکی که در اینجا utility
گرفته می‌شود. در همین چرخه ما در MDP است. خواست می‌شود که در هر گام ما دو نوع داریم: یکی که در اینجا action است و یکی که در اینجا utility
گرفته می‌شود. در همین چرخه ما در MDP است. خواست می‌شود که در هر گام ما دو نوع داریم: یکی که در اینجا action است و یکی که در اینجا utility

(ب) ما در MDP می‌توانیم دو نوع داریم: یکی که در اینجا action است و یکی که در اینجا utility
گرفته می‌شود. در همین چرخه ما در MDP است. خواست می‌شود که در هر گام ما دو نوع داریم: یکی که در اینجا action است و یکی که در اینجا utility

(ج) regret یعنی $discounted\ rewards$ ها از زمان شروع یادگیری تا آخر ما را است. مثلاً در حالتی که ما در optimal policy
داریم چون ما در حالت کار می‌کنیم و می‌توانیم و استی‌های می‌توانیم که با policy بیندازیم. reward های می‌توانیم
هم می‌گیریم تا زمانی که کم کم به پالیسی خوب می‌رسیم. regret یعنی ما در هر گام ما دو نوع داریم: یکی که در اینجا action است و یکی که در اینجا utility
گرفته می‌شود. در همین چرخه ما در MDP است. خواست می‌شود که در هر گام ما دو نوع داریم: یکی که در اینجا action است و یکی که در اینجا utility

Action: چون ما در هر گام ما دو نوع داریم: یکی که در اینجا action است و یکی که در اینجا utility
گرفته می‌شود. در همین چرخه ما در MDP است. خواست می‌شود که در هر گام ما دو نوع داریم: یکی که در اینجا action است و یکی که در اینجا utility

State: اینجا ما یک شماره داریم که وضعیت محیط را به ما می‌دهد. خواست می‌شود که در هر گام ما دو نوع داریم: یکی که در اینجا action است و یکی که در اینجا utility
گرفته می‌شود. در همین چرخه ما در MDP است. خواست می‌شود که در هر گام ما دو نوع داریم: یکی که در اینجا action است و یکی که در اینجا utility

Reward: چون می‌خواهیم با کمترین میزان به دیوارها برسیم و در راه دیوارها را reward فنی نمی‌داریم و خودی ما هم
یک reward است.