

$$V_{i+1}(s) = \max_a \left(\sum_{s'} T(s, a, s') (R(s, a, s') + \gamma V_i(s')) \right) \quad (1)$$

s	(1,1)	(1,2)	(1,3)	(2,1)	(2,2)	(2,3)
V_0	0	0	-5	0	0	5
V_1	0	0	-5	0	$0.1 \times 5 = 0.5$	5
V_2	0	$0.9 \times 0.1 \times 5 + 0.1 \times -5 = 0.4$	-5	$0.1 \times 0.9 \times 5 = 0.45$	$0.1 \times 5 = 0.5$	5

s	(1,1)	(1,2)	(1,3)	(2,1)	(2,2)	(2,3)
$\pi^*(s)$	بالا	چپ	-	بالا	بالا	-

۱) $C(1,1) - C(1,2) - C(1,3)$

۱۱) $C(1,1) - C(1,2) - C(2,3)$

۱۱۱) $C(1,1) - C(2,1) - C(2,2) - C(2,3)$

برای محاسبه تخمین‌های Monte Carlo

میانگین Reward هایی دریافتی را در مسیرهای مختلف میانگین می‌گیریم.

$$V(1,1) = \frac{-5 + 5 + 5}{3} = 1.66$$

$$V(2,2) = \frac{5 + 5}{2} = 5$$

TD-learning شکل کلی: $V(s) = V(s) + \alpha (r + \gamma V(s') - V(s))$

پس از مرحله نخست

همه آپدیت فایز این مورد
صفر خواهد بود
 $V(1,2) = 0 + 0.1 (-5 + (0.9 \times 0) - 0) = -0.5$

پس از مرحله دوم به این شکل خواهد بود:

$$\begin{cases} V(1,1) = 0 + 0.1 (0 + (0.9 \times -0.5) - 0) = -0.045 \\ V(1,2) = -0.5 + 0.1 (0 + (0.9 \times 0) + 0.5) = -0.45 \\ V(2,2) = 0 + 0.1 (5 + (0.9 \times 0) - 0) = 0.5 \end{cases}$$

DQN

۱) DQN به اختصار Deep Reinforcement Learning که همانطور که از نامش پیداست ترکیب شده یادگیری عمیق و RL است. RL با آزمون و خطا تقسیم گرفته و بخش Deep آن یادگیری عمیق در حال و عمق از داده های ورودی بدون ساختار و مدلهای فیکس فضاها را تقسیم گیری کند.

ورودی های بسیار بزرگ می گیرند و تصمیمات بهینه سازی آنها انجام می دهد در فضاها NLP ریاضیات بازی های ویدیویی و بینایی ماشین استفاده می شود.

در فضای از مسئله های تصمیم گیری، حالات و از MDP دارای ابعاد بالا هستند مثل تسلیف دور بین، که RL های سنتی در حال آنها مشکل دارند و DQN ها این چنین MDP های را حل می کنند (Sutton).

الگوریتم های مختلف در مطالعاتی در نمونه های DQN هکلو، یادگیری مشروط به هدف و یادگیری چند عاملی از جمله این پژوهشی ها هستند.