# n-step temporal difference learning

**Discounted future rewards**
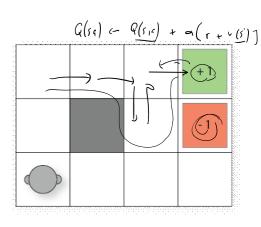
$$G_t = r_t + V(s')$$
$$= r_t + \gamma r_{t+1} + \gamma^2 r_{t+1} \cdots + U(s_{t+n})$$

**Truncated future rewards**



TD (1-step)   2-step   3-step   $n$-step   Monte Carlo

*SARSA: $Q(s,a) := Q(s,a) + \alpha [r + \gamma Q(s',a') - Q(s,a)]$*

Change the update:

$$G \leftarrow \sum_{i=t+1}^{t+n} r_i \qquad n=1$$

$$Q(s,a) \leftarrow Q(s,a) + \alpha \left( G + \gamma^n Q(s_{t+1}, a_{t+1}) - Q(s,a) \right)$$



| With 1-step learning | | | | |
|---|---|---|---|---|
| **State** | **Action** | | | |
| | North | South | East | West |
| (0,0) | 0 | 0 | 0 | 0 |
| (0,1) | 0 | 0 | 0 | 0 |
| (0,2) | 0 | 0 | 0 | 0 |
| . . . | | | | |
| (1,2) | 0 | 0 | 0 | 0 |
| (2,1) | 0 | 0 | 0 | 0 |
| (2,2) | 0 | 0 | 0.45 | 0 |
| (2,3) | 0 | 0 | 0 | 0 |
| . . . | | | | |

$$G(s_q) \leftarrow Q(s_{1s}) + a\left( r + v(s') \right)$$





n-step Sarsa for estimating $Q \approx q_*$, or $Q \approx q_\pi$ for a given $\pi$

Initialize $Q(s,a)$ arbitrarily, for all $s \in S$, $a \in A$
Initialize $\pi$ to be $\varepsilon$-greedy with respect to $Q$, or to a fixed given policy
Parameters: step size $\alpha \in (0,1]$, small $\varepsilon > 0$, a positive integer $n$
All store and access operations (for $S_t$, $A_t$, and $R_t$) can take their index mod $n$

Repeat (for each episode):
  Initialize and store $S_0 \neq$ terminal
  Select and store an action $A_0 \sim \pi(\cdot|S_0)$
  $T \leftarrow \infty$
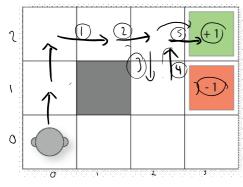  For $t = 0, 1, 2, \ldots$ :
    If $t < T$, then:
      Take action $A_t$
      Observe and store the next reward as $R_{t+1}$ and the next state as $S_{t+1}$
      If $S_{t+1}$ is terminal, then:
        $T \leftarrow t+1$
      else:
        Select and store an action $A_{t+1} \sim \pi(\cdot|S_{t+1})$
    $\tau \leftarrow t - n + 1$   ($\tau$ is the time whose estimate is being updated)
    If $\tau \geq 0$:
      $G \leftarrow \sum_{i=\tau+1}^{\min(\tau+n,T)} \gamma^{i-\tau-1} R_i$
      If $\tau + n < T$, then $G \leftarrow G + \gamma^n Q(S_{\tau+n}, A_{\tau+n})$   $(G_{\tau:\tau+n})$
      $Q(S_\tau, A_\tau) \leftarrow Q(S_\tau, A_\tau) + \alpha [G - Q(S_\tau, A_\tau)]$
      If $\pi$ is being learned, then ensure that $\pi(\cdot|S_\tau)$ is $\varepsilon$-greedy wrt $Q$
  Until $\tau = T-1$

## Exercise: Grid World



Compute 5-step SARSA update
$\alpha = 0.5$
$\gamma = 0.9$

G_5 = $\gamma \cdot 1$
G_4 = $c + \gamma^2 \cdot 1$
G_3 = $c' + 0 + \gamma^3 \cdot 1$
G_2 = $\vdots$
G_1 = $\gamma^{5} \cdot 1$

$Q(s_{(2,2)},E) = 0 + 0.5 \left[ 1\frac{0.9}{5} + G - G \right] = 0.45$

$Q(s_{(1,2)},N) = 0 + 0.5 \left[ 0.81 + 0 - 0 \right] = 0.405$

...

$Q(s_{(2,2)},S) = $

$Q(s_{(2,2)},E) = 0 + 0.5 [0.45 + 0 - 0] = 0.45$

$Q(s_{(1,2)},N) = 0 + 0.5 [0.81 + 0 - 0] = 0.405$

...

$Q(s_{(2,2)},S) = $  "

$Q(s_{(2,1)},E) = $  "

$$Q(s,s) = Q(r,a) + \alpha [\xi_t + \gamma^n Q(s',a') - Q(s,a)]$$

| State | With 1-step learning |  |  |  |
|-------|-------|-------|------|------|
|  | North | South | East | West |
| (0,0) | 0 | 0 | 0 | 0 |
| (0,1) | 0 | 0 | 0 | 0 |
| (0,2) | 0 | 0 | 0 | 0 |
| ... |  |  |  |  |
| (1,2) | 0 | 0 | 0 | 0 |
| (2,1) | 0 | 0 | 0 | 0 |
| (2,2) | 0 | 0 | 0.45 | 0 |
| (2,3) | 0 | 0 | 0 | 0 |
| ... |  |  |  |  |

| State | With 5-step learning |  |  |  |
|-------|-------|-------|------|------|
|  | North | South | East | West |
| (0,0) | 0 | 0 | 0 | 0 |
| (0,1) | 0 | 0 | 0 | 0 |
| (0,2) | 0 | 0 | 0.2953 | 0 |
| ... |  |  |  |  |
| (1,2) | 0 | 0 | 0.3281 | 0 |
| (2,1) | 0.405 | 0 | 0 | 0 |
| (2,2) | 0 | 0.3645 | 0.45 | 0 |
| (2,3) | 0 | 0 | 0 | 0 |
| ... |  |  |  |  |

Example: Random walk





MCTS + Reinforcement learning

AlphaGoZero: