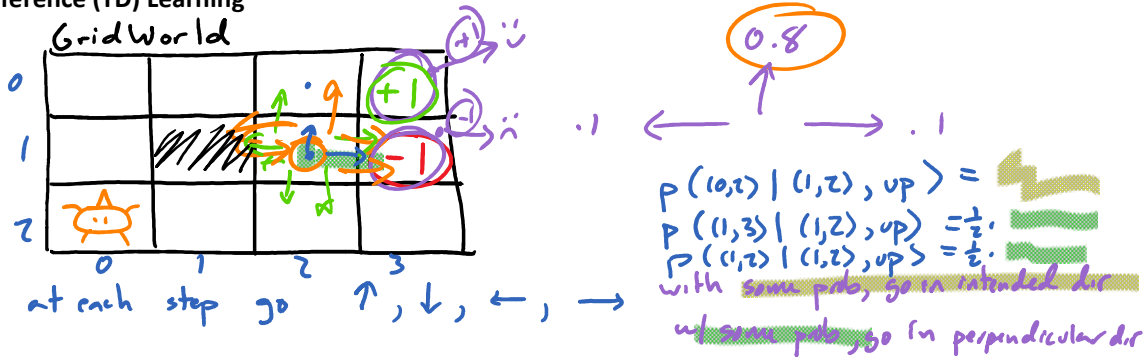


Temporal Difference (TD) Learning



model

$p(s' | s, a)$ = prob that from state s , choosing action a puts you in s'
 $r(s, a, s')$ = reward for going $s \xrightarrow{a} s'$
 $r(s, a, (0,3)) = 1$
 $r(s, a, (1,3)) = -1$
 $r(s, a, s') = 0$ for other s'

policy: function $\pi: \text{states} \rightarrow \text{action}$

$$v_{\pi}(s) = \begin{cases} 0 & \text{for terminal states} \\ \sum_{s'} p(s' | s, \pi(s)) \cdot (r(s, a, s') + \gamma \cdot v_{\pi}(s')) & \end{cases}$$

π^* : policy that maximizes $v_{\pi}(s)$ for each s

$$v^*(s) = v_{\pi^*}(s)$$

Value iteration

initialize $v(s)$ to 0 for each s
 until v has converged

for each state s

$$v(s) = \sum_a P(s' | s, a) \cdot (r(s, a, s') + \gamma \cdot v_{\pi}(s'))$$

$$v^*(s) = \max_a g^*(s, a)$$

$g(s, a)$ the value of taking action a from s

$$\pi(s) = \operatorname{argmax}_a \sum_{s'} P(s' | s, a) \cdot (r(s, a, s') + \gamma \cdot v_{\pi}(s'))$$

TD Value Learning (TD(0))

temporal difference

episode

learn $v_{\pi}(s)$ w/o knowing model (model-free)
 until done

$s \leftarrow s_0$ initial state

while s not terminal

$a \leftarrow \pi(s)$ δ -learning: choose an action (E-greedy)

converges if for all states s ,
 $\sum_{t=1}^{\infty} \alpha_{s,t} = \infty$ $\sum_{t=1}^{\infty} \alpha_{s,t}^2 < \infty$
 α applied to t^{th} update to state s

$$\alpha_{s,t} = \frac{1}{t}$$

episode

$a \leftarrow \pi(s)$ ϵ -learning: choose an action (ϵ -greedy)

$$\alpha_{s,t} = \frac{1}{t}$$

observe reward r , new state s'

$$v(s) = (1-\alpha) \cdot v(s) + \alpha \cdot (r + \gamma \cdot v(s'))$$

$$v(s) + \alpha \cdot (r + \gamma \cdot v(s') - v(s))$$

learning rate surprise

$$v(s') = \max_{a'} g(s', a')$$

observe $s \rightarrow s'$

$$g(s, a) + \alpha \cdot (r + \gamma \cdot \max_{a'} g(s', a) - g(s, a))$$

sample

Q-learning

sampled value of $v(s) =$

update estimate $v(s) \leftarrow v(s)$

$$= (1-\alpha) v(s) + \alpha R(s, s', a) + \gamma V(s')$$