

ACTIVE REINFORCEMENT LEARNING

- ① Generally, we don't just want to evaluate the expected utility of a particular policy, but rather we want to find the policy with the best expected utility. One step in this direction is to rewrite TDLEARNING in terms of the expected utility $U(q, \sigma)$ of taking action σ in state q , rather than the expected utility of policy π in state q :

TDLEARNING (Q, Σ, R, π):

set $U(q, \sigma) = 0$, $n_q = 1$ for all $(q, \sigma) \in Q \times \Sigma$

repeat:

observe transition $q \xrightarrow{\pi(q)} q'$

update $U(q, \pi(q)) = U(q, \pi(q)) + \alpha(n_q) \left(\begin{array}{l} R_o(q) \\ + \gamma \cdot U(q', \pi(q')) \\ - U^\pi(q) \end{array} \right)$

$n_q += 1$

ACTIVE REINFORCEMENT LEARNING

② Rather than averaging

$$R_o(q) + \gamma \cdot U(q', \pi(q'))$$

for each transition $q \xrightarrow{\pi(q)} q'$ we observe, we now want to average

$$R_o(q) + \gamma \max_{\sigma' \in \Sigma} U(q', \sigma')$$

reward for
state q

utility of the observed next
state, given we perform the
best action in that state

for each transition $q \xrightarrow{\sigma} q'$ we observe.

③ In TDLEARNING, we collect observations by rigidly following a fixed policy π . Now we have the flexibility to choose different actions ("experimentation") to observe their outcomes. This gives us the following "active" algorithm:

set $U(q, \sigma) = 0$, $n_{q, \sigma} = 1$ for all $(q, \sigma) \in Q \times \Sigma$
 $q = q_0$
repeat:
 choose action σ
 observe transition $q \xrightarrow{\sigma} q'$
 $U(q, \sigma) += \alpha(n_{q, \sigma}) \left[\left(R_o(q) + \gamma \max_{\sigma'} U(q', \sigma') \right) - U(q, \sigma) \right]$
 $n_{q, \sigma} += 1$

ACTIVE REINFORCEMENT LEARNING

④ This algorithm is called Q-learning. The main piece we left unspecified is how to "choose action σ ".

A simple approach is:

with prob p :

choose random action σ

("exploration")

otherwise:

choose $\underset{\sigma}{\operatorname{argmax}} U(q, \sigma)$

("exploitation")

ACTIVE REINFORCEMENT LEARNING

⑤ Generally in blackjack, the player gets to see one of the dealer's two cards (called the "up card"). This gives a player additional information on which to base a decision. For instance, if the player has a total of 14, but the dealer's up card is a 5, then it might be good to STAND, knowing that the dealer is likely to draw a 10 and then be forced to HIT on 15 (which is likely to result in a bust). We could make a state machine for this more typical version of blackjack:

$$Q = \{ \text{START, WIN, LOSE, DRAW} \} \cup \{ A_{k,u} \mid k \in [2, 21], u \in [1, 10] \}$$

$$\Sigma = \{ \text{DEAL, HIT, STAND} \}$$

$$\Delta = \{ (\text{START, DEAL, } A_{k,u} \mid k \in [2, 20], u \in [1, 10]) \} \\ \cup \{ (A_{k,u}, \text{STAND, } q) \mid k \in [2, 20], u \in [1, 10], q \in \{ \text{WIN, LOSE, DRAW} \} \} \\ \cup \{ (A_{21,u}, \text{STAND, } q) \mid u \in [1, 10], q \in \{ \text{WIN, DRAW} \} \} \\ \cup \{ (A_{k,u}, \text{HIT, } A_{k+m,u}) \mid k \in [2, 20], u, m \in [1, 10], k+m \leq 21 \} \\ \cup \{ (A_{k,u}, \text{HIT, LOSE}) \mid k \in [12, 21], u \in [1, 10] \}$$

we assume aces
always count as 1



$$q_0 = \text{START}$$

$$F = \{ \text{WIN, LOSE, DRAW} \}$$

here, $[i, j]$ denotes the
integers from i to j (inclusive)



ACTIVE REINFORCEMENT LEARNING

- ⑥ The trouble with this new state machine is that it has many more states than the "up card ignorant" formulation:

"upcard ignorant"

$$4 + 20$$

$$= 24$$

"upcard informed"

$$4 + 20 \cdot 10$$

$$= 204$$

-
- ⑦ This means that QLEARNING is now trying to estimate $U(q, \sigma)$ for approximately 400 state-action pairs (~ 200 states, ~ 2 actions/state) rather than approximately 40 state-action pairs, which intuitively will take much more experimentation and exploitation.

For instance, to get a good estimate of $U(A_{6,7}, \text{Hit})$, we have to reach a card total of 6 with a dealer upcard of 7 enough times to get a reliable sense of its utility.

So it will take longer to train. This problem only worsens as we add more information about the current state (e.g. card-counting statistics, "hard" or "soft" 16, etc.)

ACTIVE REINFORCEMENT LEARNING

- ⑤ One way to get reinforcement learning to scale to larger state spaces is to assume that each utility is a sum of weights:

$$\begin{aligned} U(A_{14,5}, \text{HIT}) &= \theta_{14, \text{HIT}}^{\text{total}} + \theta_{5, \text{HIT}}^{\text{up}} \\ U(A_{16,5}, \text{HIT}) &= \theta_{16, \text{HIT}}^{\text{total}} + \theta_{5, \text{HIT}}^{\text{up}} \\ U(A_{16,10}, \text{HIT}) &= \theta_{16, \text{HIT}}^{\text{total}} + \theta_{10, \text{HIT}}^{\text{up}} \end{aligned}$$

Annotations: "this is shared!" with arrows pointing to $\theta_{5, \text{HIT}}^{\text{up}}$ in the first two equations, and "this is shared!" with an arrow pointing to $\theta_{16, \text{HIT}}^{\text{total}}$ in the third equation.

For example, we are assuming that the expected reward of HITTING given a card total of 16 and an upcard of 5 can be expressed as the sum of

- the "value" of HITTING given a card total of 16 ($\theta_{16, \text{HIT}}^{\text{total}}$)
- the "value" of HITTING given an upcard of 5. ($\theta_{5, \text{HIT}}^{\text{up}}$)

Generally:

$$U(A_{k,u}, \sigma) = \sum_{k', u'} \theta_{k'}^{\text{total}} \delta_{k'}^{\text{total}}(A_{k,u}) + \theta_{u'}^{\text{up}} \delta_{u'}^{\text{up}}(A_{k,u})$$

where:

$$\delta_{k'}^{\text{total}}(A_{k,u}) = \begin{cases} 1 & \text{if } k=k' \\ 0 & \text{otherwise} \end{cases}$$

$$\delta_{u'}^{\text{up}}(A_{k,u}) = \begin{cases} 1 & \text{if } u=u' \\ 0 & \text{otherwise} \end{cases}$$

ACTIVE REINFORCEMENT LEARNING

- ⑨ The advantage of this assumption is that we can generalize learning across states. For instance, it is generally bad to HIT on 20 (regardless of the upcard). Even if we don't have many examples of state $A_{20,8}$ in our data so far, if we've seen several states of the form $A_{20,u}$, then we can still know not to HIT on $A_{20,8}$ because the weight $\theta_{20, \text{HIT}}^{\text{total}}$ will be low.

Rather than learning ~ 400 independent utilities $U(q, \sigma)$, we need to estimate 40 $\theta_{k, \sigma}^{\text{total}}$ weights and 20 $\theta_{u, \sigma}^{\text{up}}$ weights. So convergence (if it happens) should be faster.

ACTIVE REINFORCEMENT LEARNING

⑩ The modification to QLEARNING is minimal:

QLEARNING

set $\theta_i = 0$ for all $i \in [1, I]$

set $n_{q,\sigma} = 1$ for all $(q, \sigma) \in Q \times \Sigma$

$q = q_0$

repeat:

choose action σ

observe transition $q \xrightarrow{\sigma} q'$

$$\boxed{\theta_i} += \alpha(n_{q,\sigma}) \left[(R_0(q) + \gamma \max_{\sigma'} U(q', \sigma')) - U(q, \sigma) \right]$$

$$\boxed{\text{if } \sigma_i(q) = 1}$$

$$n_{q,\sigma} += 1$$

⑪ Why would such a modification work? Well, it doesn't always. But the rationale is this. If the weights θ_i have been set perfectly, then:

$$U(q, \sigma) = R_0(q) + \gamma \max_{\sigma'} U(q', \sigma')$$

$$\Rightarrow (R_0(q) + \gamma \max_{\sigma'} U(q', \sigma')) - U(q, \sigma) = 0$$

Thus, we won't update the weights if they're correct.

ACTIVE REINFORCEMENT LEARNING

12) If the weights that compose $U(q, \sigma)$ are wrong, then either $U(q, \sigma)$ is underestimated:

$$\sum_{i: \phi_i(q)=1} \theta_i < R_0(q) + \gamma \max_{\sigma'} U(q', \sigma')$$

$$\Rightarrow (R_0(q) + \gamma \max_{\sigma'} U(q', \sigma')) - \sum_{i: \phi_i(q)=1} \theta_i \text{ is positive.}$$

So we increase the weights θ_i , or $U(q, \sigma)$ is overestimated, so we (similarly) decrease the weights θ_i .

However, there are no convergence guarantees. Sometimes however, it works spectacularly.