



Reinforcement Learning

CSE 415: Introduction to Artificial Intelligence
University of Washington
Winter 2019

Presented by S. Tanimoto, University of Washington, based on material by Dan Klein and Pieter Abbeel - University of California.



Outline

- Planning vs Learning
- Sample-Based Policy Evaluation
- Temporal Difference Learning
- Active Reinforcement Learning
- Q-Learning
- Exploration vs Exploitation
- Regret

2



Planning vs. Learning

- **Planning:**
 - The **agent knows the whole MDP (including S, T and R)**.
 - If there is no noise, compute a plan (sequence of actions) from start to goal that maximizes total reward. (**Unif. Cost Search** can often do it).
 - If there is noise, compute an optimal policy, combining **Value Iteration** with one step of argmaxing from Q-values.
- **Learning:**
 - The **agent does NOT know T and R at all**. The agent doesn't even know S, but can recognize states it has visited before.
 - Model-based learning tries to build representations of S, T, and R, but can be very inefficient.
 - Sample-based methods (**Temporal-Difference Learning**, and **Q-Learning**) focus on learning good approximations to V and Q values directly.

Agent does not even know state space! Slowly gathering information about the system.

Approach 1: try to approximate T and R by making tons of transitions. Takes a lot of time, and some irrelevant work; some areas of state space we rarely ever visit! Bad!

we don't know T, so instead think of the 'RHS as the EXPECTATION: average of a bunch of samples



Sample-Based Policy Evaluation?

- We want to improve our estimate of V by computing these averages:

$$V_{k+1}^\pi(s) \leftarrow \sum_{s'} T(s, \pi(s), s') [R(s, \pi(s), s') + \gamma V_k^\pi(s')]$$

FIXED policy that we have given agent

- Idea: Take samples of outcomes s' (by doing the action!) and average

$$\text{sample}_1 = R(s, \pi(s), s'_1) + \gamma V_k^\pi(s'_1)$$

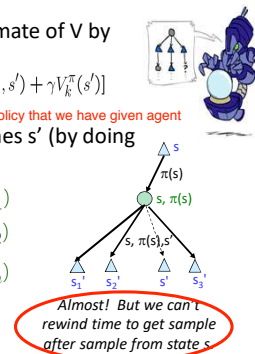
$$\text{sample}_2 = R(s, \pi(s), s'_2) + \gamma V_k^\pi(s'_2)$$

$$\dots$$

$$\text{sample}_n = R(s, \pi(s), s'_n) + \gamma V_k^\pi(s'_n)$$

$$V_{k+1}^\pi(s) \leftarrow \frac{1}{n} \sum_i \text{sample}_i$$

NO! We will average a different way

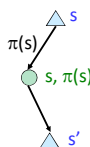


passive learning: agent does not choose actions; we dictate a set of actions and it learns about state space from it
active learning: agent chooses actions



Temporal Difference Learning

- **Big idea: learn from every experience!**
 - Update V(s) each time we experience a transition (s, a, s', r)
 - Likely outcomes s' will contribute updates more often **rather than take mean of samples after many trials, continuously update state value and learn!**
- Temporal difference learning of values
 - Policy still fixed, still doing evaluation!
 - Move values toward value of whatever successor occurs: running average



$$\text{Sample of } V(s): \text{ sample} = R(s, \pi(s), s') + \gamma V^\pi(s')$$

$$\text{Update to } V(s): V^\pi(s) \leftarrow (1 - \alpha)V^\pi(s) + (\alpha)\text{sample}$$

$$\text{Same update: } V^\pi(s) \leftarrow V^\pi(s) + \alpha(\text{sample} - V^\pi(s))$$

alpha is our learning rate!

5



Exponential Moving Average

- Exponential moving average
 - The running interpolation update: $\bar{x}_n = (1 - \alpha) \cdot \bar{x}_{n-1} + \alpha \cdot x_n$
 - **Makes recent samples more important:**

$$\bar{x}_n = \frac{x_n + (1 - \alpha) \cdot x_{n-1} + (1 - \alpha)^2 \cdot x_{n-2} + \dots}{1 + (1 - \alpha) + (1 - \alpha)^2 + \dots}$$

- Forgets about the past (distant past values were wrong anyway)

- **Decreasing learning rate (alpha) can give converging averages**

6

Example: Temporal Difference Learning

States

Observed Transitions

B, east, C, -2 C, east, D, -2

Assume: $\gamma = 1$, $\alpha = 1/2$

$$V(\text{Bnew}) = (1 - 1/2)V(\text{B_old}) + 1.2(-2 + 1 \cdot V(\text{C})) = -1$$

$$V^\pi(s) \leftarrow (1 - \alpha)V^\pi(s) + \alpha [R(s, \pi(s), s') + \gamma V^\pi(s')]$$

Problems with TD Value Learning

- TD value learning is a model-free way to do policy evaluation, mimicking Bellman updates with running sample averages
- However, if we want to turn values into a (new) policy, we don't have the needed Q values:

$$\pi(s) = \arg \max_a Q(s, a)$$

$$Q(s, a) = \sum_{s'} T(s, a, s') [R(s, a, s') + \gamma V(s')]$$

Idea: learn Q-values, not values

Makes action selection model-free too!

As the Q values are what allow us to determine the optimal policy.

Active Reinforcement Learning

Let the agent decide their own moves and learn (possibly from their mistakes)

Active Reinforcement Learning

- Full reinforcement learning: optimal policies (like value iteration)
 - You don't know the transitions $T(s, a, s')$
 - You don't know the rewards $R(s, a, s')$
 - You choose the actions now
 - Goal: learn the optimal policy / values
- In this case:
 - Learner makes choices, we have little information about
 - Fundamental tradeoff: exploration vs. exploitation
 - This is NOT offline planning! You actually take actions in the world and find out what happens...

Prioritize exploration initially

Detour: Q-Value Iteration

- Value iteration: find successive (depth-limited) values
 - Start with $V_0(s) = 0$, which we know is right
 - Given V_k , calculate the depth $k+1$ values for all states:
$$V_{k+1}(s) \leftarrow \max_a \sum_{s'} T(s, a, s') [R(s, a, s') + \gamma V_k(s')]$$
- But Q-values are more useful, so compute them instead
 - Start with $Q_0(s, a) = 0$, which we know is right
 - Given Q_k , calculate the depth $k+1$ q-values for all q-states:
$$Q_{k+1}(s, a) \leftarrow \sum_{s'} T(s, a, s') [R(s, a, s') + \gamma \max_{a'} Q_k(s', a')]$$

k = "how many turns we can make" to decide value of our state. more correctly, the number of update iterations we run to approximate all our Q states

If we know T and R

Q-Learning

- Q-Learning: sample-based Q-value iteration

$$Q_{k+1}(s, a) \leftarrow \sum_{s'} T(s, a, s') [R(s, a, s') + \gamma \max_{a'} Q_k(s', a')]$$

If we know T and R
- Learn $Q(s, a)$ values as you go
 - Receive a sample (s, a, s', r)
 - Consider your old estimate: $Q(s, a)$
 - Consider your new sample estimate:

$$\text{sample} = R(s, a, s') + \gamma \max_{a'} Q(s', a')$$
 - Incorporate the new estimate into a running average:

$$Q(s, a) \leftarrow (1 - \alpha)Q(s, a) + (\alpha) \{\text{sample}\}$$

Q-VALUES AFTER 1000 EPISODES

[Demo: Q-learning - gridworld (L10D2)]
[Demo: Q-learning - crawler (L10D3)]



Video of Demo Q-Learning -- Gridworld



13



Video of Demo Q-Learning -- Crawler



14



Q-Learning Properties

- Amazing result: Q-learning converges to optimal policy -- even if you're acting suboptimally!
- This is called **off-policy learning**
continuously explore and don't starve states (visit them all)
- Caveats:
 - You have to explore enough
 - You have to eventually make the learning rate small enough
 - ... but not decrease it too quickly
 - Basically, in the limit, it doesn't matter how you select actions (!)



15



How to Explore?

- Several schemes for forcing exploration
 - Simplest: random actions (ϵ -greedy)
 - Every time step, flip a coin
 - With (small) probability ϵ , act randomly
 - With (large) probability $1-\epsilon$, act on current policy
 - Problems with random actions?
 - You do eventually explore the space, but keep thrashing around once learning is done
 - One solution: lower ϵ over time
 - Another solution: exploration functions



Video of Demo Q-learning – Manual Exploration – Bridge Grid



Video of Demo Q-learning – Epsilon-Greedy – Crawler





Exploration Functions



- When to explore?
 - Random actions: explore a fixed amount
 - Better idea: explore areas whose badness is not (yet) established, eventually stop exploring
- Exploration function
 - Takes a value estimate u and a visit count n , and returns an optimistic utility, e.g. $f(u, n) = u + k/n$

Note: this propagates the “bonus” back to states that lead to unknown states as well!

Regular Q-Update: $Q(s, a) \leftarrow \alpha R(s, a, s') + \gamma \max_{a'} Q(s', a')$

Modified Q-Update: $Q(s, a) \leftarrow \alpha R(s, a, s') + \gamma \max_{a'} f(Q(s', a'), N(s', a'))$



Video of Demo Q-learning – Exploration Function – Crawler



Regret

- Even if you learn the optimal policy you still make mistakes along the way!
- Regret is a measure of your total mistake cost: the difference between your (expected) rewards, including youthful suboptimality, and optimal (expected) rewards
- Minimizing regret goes beyond learning to be optimal – it requires optimally learning to be optimal
- Example: random exploration and exploration functions both end up optimal, but random exploration has higher regret

