

- Model-Free Control

- Last lecture, we saw Model-Free Prediction. This time we want to optimize VF
- Most interesting problems have MDPs that are either:
 - ① Unknown to us
 - ② Too big to use except with sampling
 - Model-free controls addresses those limitations

- On and Off-Policy Learning

- On-policy Learning: "Learning on the job"
 - Learning about policy π from experience sampled from π
- Off-policy Learning: "looking over someone's shoulder"
 - Learn about policy π from experience sampled from μ

- Generalized Policy Iteration (Reinforcement)

- Alternate between policy evaluation and greedy policy improvement
 - Converges on optimal VF
 - Can we swap in MC Evaluation to evaluate VF instead of using DP?
 - There are two problems:
 - Will come back to this → ② If you act greedily, there is a probability that you do not explore the entire state space and therefore miss optimal policy.
- ① We want to be model-free, but you need a model of the MDP to know state transition dynamics.
- Alternative is to use action value function

$$\pi^*(s) := \underset{a \in A}{\operatorname{argmax}} Q(s, a)$$

→ Policy Iteration with Action-Value Function

Can you discuss how this is done? → Policy Evaluation: Monte Carlo policy Evaluation $Q = q_{\pi}$

→ Take mean of all state-action pairs

Policy Improvement: Greedy policy improvement?

→ No! We may miss the optimal choice

→ Back to problem ②

→ You need to keep exploring all options to properly estimate value function at each state

→ ϵ -Greedy Exploration:

→ Take random action with probability ϵ

→ Take greedy action with probability $1-\epsilon$

→ ϵ -Greedy Guarantees a policy improvement:

$$q_{\pi}(s, \pi^*(s)) = \sum_{a \in A} \pi^*(a|s) q_{\pi}(s, a)$$

$$= \underbrace{\frac{\epsilon}{n} \sum_{a \in A} 1_{\pi}(s, a)}_{\text{Random Action}} + (1-\epsilon) \max_{a \in A} 1_{\pi}(s, a)$$

Ans is greater than any weighted sum over all actions. → $\geq \frac{\epsilon}{n} \sum_{a \in A} 1_{\pi}(s, a) + (1-\epsilon) \sum_{a \in A} \frac{\pi(a|s) + \frac{\epsilon}{n}}{1-\epsilon} 1_{\pi}(s, a)$

$$= \sum_{a \in A} \pi(a|s) q_{\pi}(s, a) = v_{\pi}(s)$$

• Monte Carlo Policy Iteration

→ Policy Evaluation: Monte Carlo policy Evaluation, $Q = q_{\pi}$

→ Mean return from every state-action pair

→ Policy Improvement: ϵ -greedy policy improvement (stochastically)

* This will converge to optimal policy, given enough time

→ Now let's try to make this more efficient.

• Monte Carlo Control: Improve policy after each episode. Always act greedily with respect to most recent updated value function.

→ Iterate between running single episodes and immediately improving policy

• How can we guarantee that we find the best possible policy?

→ We must balance exploration of new policies and exploitation of existing policies.

→ One idea: Greedy in the limit of infinite exploration (GLIE)

→ Come up with a schedule such that two conditions are met:

① All state-action pairs are explored infinitely many times

$$\lim_{k \rightarrow \infty} N_k(s, a) = \infty$$

→ E.g., ϵ -greedy policy improvement satisfies this

② The policy converges to a greedy policy:

$$\lim_{k \rightarrow \infty} \pi_k(a|s) = \mathbb{I}(a = \arg \max Q_k(s, a))$$

→ E.g. Let ϵ decay slowly such as $\epsilon_k = 1/k$

• GLIE MC Control:

• Sample k^{th} episode using π : $\{s_1, a_1, r_1, \dots, s_T\} \sim \pi$

• For each state s_t and action a_t in the episode:

$$N(s_t, a_t) \leftarrow N(s_t, a_t) + 1$$

Incremental Update to the mean

$$Q(s_t, a_t) \leftarrow Q(s_t, a_t) + \frac{1}{N(s_t, a_t)} (G_t - Q(s_t, a_t))$$

Taking returns from better and better policies into account. Past gets dominated by more informally recent q-values.

- Improve Policy based on new action-value function

$$\epsilon \leftarrow 1/k$$

$$\pi \leftarrow \epsilon\text{-Greedy}(Q)$$

- This algorithm finds optimal action-value function

↳ You can update policy after each algorithm

↳ Does it matter how you initialize Q ?

↳ For this particular algorithm, it makes no difference. In other algorithms, a good initialized Q will converge sooner, but convergence guaranteed eventually regardless of Q init.

• MC vs. TD Control:

TD has advantages over MC:

- ① Lower Variance
- ② Online
- ③ Incomplete Sequences

Natural Idea: use TD instead of MC in our policy iteration strategy

- ↳ Use TD instead of MC to estimate $Q(s, a)$
- ↳ Use ϵ -greedy improvement
- ↳ Run policy improvement after every time step

• SARSA: $s, a \xrightarrow[R]{ } s' \xrightarrow{} a'$

$$Q(s, a) \leftarrow Q(s, a) + \alpha (R + \gamma Q(s', a') - Q(s, a))$$

↳ Each step, we update q-value using SARSA

→ You've already changed your value function, so at every timestep, update policy using ϵ -greedy policy improvement

Algorithm

Initialize $Q(s,a)$

Repeat: For every episode

 Initialize s (pick random starting state)

 Choose action from s using ϵ -greedy policy

Repeat: For each step of the episode

 Take action, observe R, s'

 Choose A' from s' using ϵ -greedy policy

 Update $Q(s,a) \leftarrow Q(s,a) + \alpha [R + \gamma Q(s',A') - Q(s,a)]$

 Until s is terminal