

ECE276B: Planning & Learning in Robotics

Lecture 16: Model-free Control

Lecturer:

Nikolay Atanasov: natanasov@ucsd.edu

Teaching Assistants:

Tianyu Wang: tiw161@eng.ucsd.edu

Yongxi Lu: yol070@eng.ucsd.edu

UC San Diego

JACOBS SCHOOL OF ENGINEERING

Electrical and Computer Engineering

Generalized Policy Iteration

- **Policy Evaluation:** Given π , compute V^π :

$$V^\pi(x) = g(x, \pi(x)) + \gamma \mathbb{E}_{x' \sim p_f(\cdot|x, \pi(x))} [V^\pi(x')], \quad \forall x \in \mathcal{X}$$

- **Policy Improvement:** Given V^π obtain a new policy π' :

$$\pi'(x) = \arg \min_{u \in \mathcal{U}(x)} \{g(x, u) + \gamma \mathbb{E}_{x' \sim p_f(\cdot|x, u)} [V^\pi(x')]\}, \quad \forall x \in \mathcal{X}$$

Policy Improvement Theorem

Let π and π' be deterministic policies such that $V^\pi(x) \geq Q^\pi(x, \pi'(x))$ for all $x \in \mathcal{X}$. Then, π' is at least as good as π , i.e., $V^\pi(x) \geq V^{\pi'}(x)$ for all $x \in \mathcal{X}$

- **Proof:**

$$\begin{aligned} V^\pi(x) &\geq Q^\pi(x, \pi'(x)) = g(x, \pi'(x)) + \gamma \mathbb{E}_{x' \sim p_f(\cdot|x, \pi'(x))} V^\pi(x') \\ &\geq g(x, \pi'(x)) + \gamma \mathbb{E}_{x' \sim p_f(\cdot|x, \pi'(x))} Q^\pi(x', \pi'(x')) \\ &= g(x, \pi'(x)) + \gamma \mathbb{E}_{x' \sim p_f(\cdot|x, \pi'(x))} \{g(x', \pi'(x')) + \gamma \mathbb{E}_{x'' \sim p_f(\cdot|x', \pi'(x'))} V^\pi(x'')\} \\ &\geq \dots \geq \mathbb{E} \left[\sum_{t=0}^{\infty} \gamma^t g(x_t, \pi'(x_t)) \middle| x_0 = x \right] = V^{\pi'}(x) \end{aligned}$$

Model-free Generalized Policy Iteration

- **Policy Evaluation:** given π , GPI iterates \mathcal{T}_π to compute V^π :

$$DP : \mathcal{T}_\pi[V](x) = g(x, \pi(x)) + \gamma \mathbb{E}_{x' \sim p_f(\cdot|x, \pi(x))} [V(x')]$$

$$TD : \mathcal{T}_\pi[V](x_t) \approx V(x_t) + \alpha [g(x_t, u_t) + \gamma V(x_{t+1}) - V(x_t)]$$

$$MC : \mathcal{T}_\pi[V](x_t) \approx V(x_t) + \alpha \left[\sum_{k=0}^{T-t-1} \gamma^k g(x_{t+k}, u_{t+k}) + \gamma^{T-t} g_T(x_T) - V(x_t) \right]$$

or alternatively Q^π :

$$DP : \mathcal{T}_\pi[Q](x, u) = g(x, u) + \gamma \mathbb{E}_{x' \sim p_f(\cdot|x, u)} [Q(x', \pi(x'))]$$

$$TD : \mathcal{T}_\pi[Q](x_t, u_t) \approx Q(x_t, u_t) + \alpha [g(x_t, u_t) + \gamma Q(x_{t+1}, u_{t+1}) - Q(x_t, u_t)]$$

$$MC : \mathcal{T}_\pi[Q](x_t, u_t) \approx Q(x_t, u_t) + \alpha \left[\sum_{k=0}^{T-t-1} \gamma^k g(x_{t+k}, u_{t+k}) + \gamma^{T-t} g_T(x_T) - Q(x_t, u_t) \right]$$

- **Policy Improvement:** given V^π or Q^π compute improved π' :

$$Q^\pi(x, u) = g(x, u) + \gamma \mathbb{E}_{x' \sim p_f(\cdot|x, u)} [V^\pi(x')]$$

$$\pi'(x) = \arg \min_{u \in \mathcal{U}(x)} Q^\pi(x, u)$$

Model-free Generalized Policy Iteration

- ▶ **Policy Evaluation:** use MC or TD to estimate Q^π instead of V^π so that the policy improvement step is model-free, i.e., can compute $\min_u Q^\pi(x, u)$ without knowing p_f
- ▶ **Policy Improvement:** the fact that Q^π is approximation still causes some problems:
 - ▶ Picking the “best” control according to the current estimate Q^π might not be the actual best control
 - ▶ If a deterministic policy is used for Evaluation/Improvement, one will observe returns for only one of the possible controls at each state and also might not visit many states. Hence, estimating Q^π will not be possible at those never-visited states and controls

Example: Greedy Control Selection (David Silver)

- ▶ There are two doors in front of you
- ▶ You open the left door and get reward 0
 $V(\textit{left}) = +0$
- ▶ You open the right door and get reward +1
 $V(\textit{right}) = +1$
- ▶ You open the right door and get reward +3
 $V(\textit{right}) = +3$
- ▶ You open the right door and get reward +2
 $V(\textit{right}) = +2$
- ▶ Are you sure you have chosen the best door?



"Behind one door is tenure - behind the other is flipping burgers at McDonald's."

Model-free Control

- ▶ Two ideas to ensure that you do not commit to the (wrong) controls too early and continue exploring the state and control spaces:
 1. **Exploring Starts**: in each episode $\rho^{(k)} \sim \pi$, choose initial state-control pairs with non-zero probability among all possible pairs $\mathcal{X} \times \mathcal{U}$
 2. ϵ -**Soft Policy**: use a **stochastic policy** under which every control has a non-zero probability of being chosen and hence every reachable state will have non-zero probability of being encountered

First-visit MC Policy Iteration with Exploring Starts

Algorithm 1 MC Policy Iteration with Exploring Starts

- 1: **Init:** $Q(x, u), \pi(x)$ for all $x \in \mathcal{X}$ and $u \in \mathcal{U}$
 - 2: **loop**
 - 3: Choose $(x_0, u_0) \in \mathcal{X} \times \mathcal{U}$ randomly ▷ exploring starts!
 - 4: Generate an episode $\rho = x_0, u_0, x_1, u_1, \dots, x_{T-1}, u_{t-1}, x_T$ from π
 - 5: **for** each x, u in ρ **do**
 - 6: $G \leftarrow$ return following the first occurrence of x, u
 - 7: $Q(x, u) \leftarrow Q(x, u) + \alpha(G - Q(x, u))$
 - 8: **for** each x in ρ **do**
 - 9: $\pi(x) \leftarrow \arg \min_u Q(x, u)$
-

ϵ -Greedy Exploration

- ▶ As an alternative to exploring starts, to ensure continual exploration it must be possible to encounter all $|\mathcal{U}(x)|$ controls at state x with non-zero probability
- ▶ **ϵ -Greedy Policy** a stochastic policy that picks the best control according to $Q(x, u)$ in the policy improvement step but ensures that all other controls are selected with a small (non-zero) probability:

$$\pi(u | x) = \begin{cases} 1 - \epsilon + \frac{\epsilon}{|\mathcal{U}(x)|} & \text{if } u = \arg \min_{u' \in \mathcal{U}(x)} Q(x, u') \\ \frac{\epsilon}{|\mathcal{U}(x)|} & \text{otherwise} \end{cases}$$

ϵ -Greedy Policy Improvement

Theorem: ϵ -Greedy Policy Improvement

For any ϵ -soft policy π , the ϵ -greedy policy π' with respect to Q^π is an improvement, i.e., $V^{\pi'}(x) \leq V^\pi(x)$ for all $x \in \mathcal{X}$

► **Proof:**

$$\begin{aligned}\mathbb{E}_{u' \sim \pi'(\cdot|x)} [Q^\pi(x, u')] &= \sum_{u' \in \mathcal{U}(x)} \pi'(u' | x) Q^\pi(x, u') \\ &= \frac{\epsilon}{|\mathcal{U}(x)|} \sum_{u' \in \mathcal{U}(x)} Q^\pi(x, u') + (1 - \epsilon) \min_{u \in \mathcal{U}(x)} Q^\pi(x, u) \\ &\leq \frac{\epsilon}{|\mathcal{U}(x)|} \sum_{u' \in \mathcal{U}(x)} Q^\pi(x, u') + (1 - \epsilon) \sum_{u \in \mathcal{U}(x)} \frac{\pi(u | x) - \frac{\epsilon}{|\mathcal{U}(x)|}}{1 - \epsilon} Q^\pi(x, u) \\ &= \sum_{u \in \mathcal{U}(x)} \pi(u | x) Q^\pi(x, u) = V^\pi(x)\end{aligned}$$

► Then, from the policy improvement theorem, $V^{\pi'}(x) \leq V^\pi(x)$, $\forall x \in \mathcal{X}$

First-visit MC Policy Iteration with ϵ -Greedy Improvement

Algorithm 2 First-visit MC Policy Iteration with ϵ -Greedy Improvement

- 1: **Init:** $Q(x, u)$, $\pi(u|x)$ (ϵ -soft policy) for all $x \in \mathcal{X}$ and $u \in \mathcal{U}$
 - 2: **loop**
 - 3: Generate an episode $\rho := x_0, u_0, x_1, u_1, \dots, x_{T-1}, u_{t-1}, x_T$ from π
 - 4: **for** each x, u in ρ **do**
 - 5: $G \leftarrow$ return following the first occurrence of x, u
 - 6: $Q(x, u) \leftarrow Q(x, u) + \alpha (G - Q(x, u))$
 - 7: **for** each x in ρ **do**
 - 8: $u^* \leftarrow \arg \min_u Q(x, u)$
 - 9:
$$\pi(u|x) \leftarrow \begin{cases} 1 - \epsilon + \frac{\epsilon}{|\mathcal{U}(x)|} & \text{if } u = u^* \\ \frac{\epsilon}{|\mathcal{U}(x)|} & \text{if } u \neq u^* \end{cases}$$
-

Temporal-Difference Control

- ▶ TD prediction has several advantages over MC prediction:
 - ▶ Works with incomplete episodes
 - ▶ Lower variance
 - ▶ Can perform online updates after every transition
- ▶ To use TD instead of MC in a complete policy iteration algorithm, we still need to trade-off exploration and exploitation:
 - ▶ Apply TD to $Q(x, u)$ for policy evaluation
 - ▶ Can update $Q(x, u)$ after every transition within an episode
 - ▶ Use an ϵ -greedy policy for policy improvement

TD Policy Iteration with ϵ -Greedy Improvement (SARSA)

- ▶ **SARSA**: estimates the action-value function Q^π using TD updates after every $S_t, A_t, R_{t+1}, S_{t+1}, A_{t+1}$ transition:

$$Q(x_t, u_t) \leftarrow Q(x_t, u_t) + \alpha [g(x_t, u_t) + \gamma Q(x_{t+1}, u_{t+1}) - Q(x_t, u_t)]$$

- ▶ Ensures exploration via an ϵ -greedy policy in the policy improvement step

Algorithm 3 SARSA

- 1: **Init**: $Q(x, u)$ for all $x \in \mathcal{X}$ and all $u \in \mathcal{U}$
 - 2: **loop**
 - 3: $\pi \leftarrow \epsilon$ -greedy policy derived from Q
 - 4: Generate episode $\rho := (x_{0:T}, u_{0:T-1})$ from π
 - 5: **for** $(x, u, x', u') \in \rho$ **do**
 - 6: $Q(x, u) \leftarrow Q(x, u) + \alpha [g(x, u) + \gamma Q(x', u') - Q(x, u)]$
-

Convergence of Model-free Policy Iteration

► Greedy in the Limit with Infinite Exploitation (GLIE):

- All state-control pairs are explored infinitely many times: $\lim_{k \rightarrow \infty} N(x, u) = \infty$
- The ϵ -greedy policy converges to a greedy policy

$$\lim_{k \rightarrow \infty} \pi_k(u | x) = \mathbb{1}\{u = \arg \min_{u' \in \mathcal{U}(x)} Q(x, a')\}$$

- Example: If $\epsilon_k = \frac{1}{k}$, then ϵ -greedy is GLIE

Theorem: Convergence of MC Policy Iteration

GLIE MC Policy Iteration converges to the optimal action-value function, $Q(x, u) \rightarrow Q^*(x, u)$ as the number of episodes $k \rightarrow \infty$

Theorem: Convergence of TD Policy Iteration

SARSA converges to the optimal action-value function, $Q(x, u) \rightarrow Q^*(x, u)$ as $k \rightarrow \infty$ as long as:

- The sequence of ϵ -greedy policies $\pi_k(u | x)$ is GLIE
- The sequence of step sizes α_k is Robbins-Monro

On-Policy vs Off-Policy Learning

- ▶ **On-policy Prediction:** estimate V^π or Q^π using experience from π
- ▶ **Off-policy Prediction:** estimate V^π or Q^π using experience from μ
- ▶ On-policy methods:
 - ▶ evaluate or improve the policy that is used to make decisions
 - ▶ require well-designed exploration functions
 - ▶ empirically successful with function approximation
- ▶ Off-policy methods:
 - ▶ evaluate or improve a different policy from the (behavior) policy used to generate data
 - ▶ can use an effective exploratory policy to generate data while learning about an optimal policy
 - ▶ can learn from observing humans or other agents
 - ▶ can re-use experience from old policies $\pi_1, \pi_2, \dots, \pi_{n-1}$
 - ▶ can learn about multiple policies while following one policy
 - ▶ have problems with function approximation and eligibility traces

Importance Sampling for Off-policy Learning

- ▶ To use returns generated from μ to evaluate π , we need to re-weight the stage-costs according to the similarity (i.e., likelihood) of states encountered by the two different policies
- ▶ **Importance Sampling**: estimates the expectation of a function with respect to a different distribution:

$$\begin{aligned}\mathbb{E}_{x \sim p(\cdot)}[f(x)] &= \int p(x)f(x)dx \\ &= \int q(x)\frac{p(x)}{q(x)}f(x)dx = \mathbb{E}_{x \sim q(\cdot)}\left[\frac{p(x)}{q(x)}f(x)\right]\end{aligned}$$

Importance Sampling for Off-policy MC Learning

- ▶ To use returns generated from μ to evaluate π via MC, weight the long-term cost G_t via importance-sampling corrections along the whole episode:

$$G_t^{\pi/\mu} = \frac{\pi(u_t|x_t)}{\mu(u_t|x_t)} \frac{\pi(u_{t+1}|x_{t+1})}{\mu(u_{t+1}|x_{t+1})} \dots \frac{\pi(u_{T-1}|x_{T-1})}{\mu(u_{T-1}|x_{T-1})} G_t$$

- ▶ Update the value estimate towards the *corrected return*:

$$V(x_t) \leftarrow V(x_t) + \alpha \left(G_t^{\pi/\mu} - V(x_t) \right)$$

- ▶ Importance sampling in MC can dramatically increase the variance and cannot be used if μ is zero when π is non-zero

Importance Sampling for Off-policy TD Learning

- ▶ To use returns generated from μ to evaluate π via TD, weight the TD target $g(x, u) + \gamma V(x')$ by importance sampling:

$$V(x_t) \leftarrow V(x_t) + \alpha \left(\frac{\pi(u_t | x_t)}{\mu(u_t | x_t)} (g(x_t, u_t) + \gamma V(x_{t+1})) - V(x_t) \right)$$

- ▶ Importance sampling in TD is much lower variance than in MC and the policies need to be similar (i.e., μ should not be zero when π is non-zero) over a single step only

Off-policy TD Control without Importance Sampling

- ▶ **Q-Learning** (Watkins, 1989): one of the early breakthroughs in reinforcement learning was the development of an off-policy TD algorithm that does not use importance sampling
- ▶ Q-Learning approximates $\mathcal{T}_*[Q](x, u)$ directly using samples:

$$Q(x_t, u_t) \leftarrow Q(x_t, u_t) + \alpha \left[g(x_t, u_t) + \gamma \min_{u \in \mathcal{U}(x_{t+1})} Q(x_{t+1}, u) - Q(x_t, u_t) \right]$$

- ▶ The learned Q function eventually approximates Q^* **regardless of the policy being followed!**

Theorem

Q-Learning converges almost surely to Q^* assuming that all state-control pairs continue to be updated and the learning rate α is chosen via the usual Robbins-Monro stochastic approximation condition

- ▶ C. J. Watkins and P. Dayan. "Q-learning," Machine learning, 1992.

Q-Learning: Off-policy TD Learning

Algorithm 4 Q-Learning

- 1: **Init:** $Q(x, u)$ for all $x \in \mathcal{X}$ and all $u \in \mathcal{U}$
 - 2: **loop**
 - 3: $\pi \leftarrow \epsilon$ -greedy policy derived from Q
 - 4: Generate episode $\rho := (x_{0:T}, u_{0:T-1})$ from π
 - 5: **for** $(x, u, x') \in \rho$ **do**
 - 6: $Q(x, u) \leftarrow Q(x, u) + \alpha [g(x, u) + \gamma \min_{u'} Q(x', u') - Q(x, u)]$
-

Relationship Between DP and TD

| Full Backups (DP) | Sample Backups (TD) |
|---|--|
| Policy Evaluation $V(x) \leftarrow \mathcal{T}_\pi[V](x) = g(x, \pi(x)) + \gamma \mathbb{E}_{x'} [V(x')]$ | TD Prediction $V(x) \leftarrow V(x) + \alpha(g(x, u) + \gamma V(x') - V(x))$ |
| Policy Q-Evaluation $Q(x, u) \leftarrow \mathcal{T}_\pi[Q](x, u) = g(x, u) + \gamma \mathbb{E}_{x'} [Q(x', \pi(x'))]$ | SARSA $Q(x, u) \leftarrow Q(x, u) + \alpha(g(x, u) + \gamma Q(x', u') - Q(x, u))$ |
| Value Iteration $V(x) \leftarrow \mathcal{T}_*[V](x) = \min_u \{g(x, u) + \gamma \mathbb{E}_{x'} [V(x')]\}$ | N/A |
| Q-Value Iteration $Q(x, u) \leftarrow \mathcal{T}_*[Q](x, u) = g(x, u) + \gamma \mathbb{E}_{x'} \left[\min_{u'} Q(x', u') \right]$ | Q-Learning $Q(x, u) \leftarrow Q(x, u) + \alpha \left(g(x, u) + \gamma \min_{u'} Q(x', u') - Q(x, u) \right)$ |

Batch Sampling-based Q-Value Iteration

Algorithm 5 Batch Sampling-based Q-Value Iteration

1: **Init:** $Q^{(0)}(x, u)$ for all $x \in \mathcal{X}$ and all $u \in \mathcal{U}$

2: **loop**

3: $\pi \leftarrow \epsilon$ -greedy policy derived from Q

4: Generate episodes $\{\rho^{(k)}\}_{k=1}^K$ from π

5: **for** $(x, u) \in \mathcal{X} \times \mathcal{U}$ **do**

6:
$$Q^{(i+1)} = \frac{1}{K} \sum_{k=1}^K \frac{\sum_{t=0}^T \mathcal{T}_*[Q^{(i)}](x_t^{(k)}, u_t^{(k)}, x_{t+1}^{(k)}) \mathbb{1}\{(x_t^{(k)}, u_t^{(k)}) = (x, u)\}}{\sum_{t=0}^T \mathbb{1}\{(x_t^{(k)}, u_t^{(k)}) = (x, u)\}}$$

- ▶ Batch Sampling-based Q-Value Iteration behaves like $Q^{(i+1)} = \mathcal{T}_*[Q^{(i)}] + \text{noise}$. Does it actually converge?

Least-squares Backup Version

- ▶ $Q^{(i+1)}(x, u) = \mathbf{mean} \left\{ \mathcal{T}_*[Q^{(i)}](x_t^{(k)}, u_t^{(k)}, x_{t+1}^{(k)}), \forall k, t \text{ such that } (x_t^{(k)}, u_t^{(k)}) = (x, u) \right\}$
- ▶ $\mathbf{mean} \left\{ x^{(k)} = \arg \min_x \sum_{k=1}^K \|x^{(k)} - x\|^2 \right\}$
- ▶ $Q^{(i+1)}(x, u) = \arg \min_q \sum_{k=1}^K \sum_{(x_t^{(k)}, u_t^{(k)})=(x, u)} \left\| \mathcal{T}_*[Q^{(i)}](x_t^{(k)}, u_t^{(k)}, x_{t+1}^{(k)}) - q \right\|^2$
- ▶ $Q^{(i+1)} = \arg \min_Q \sum_{k=1}^K \sum_{t=0}^T \left\| \mathcal{T}_*[Q^{(i)}](x_t^{(k)}, u_t^{(k)}, x_{t+1}^{(k)}) - Q(x_t^{(k)}, u_t^{(k)}) \right\|^2$

Algorithm 6 Batch Least-squares Q-Value Iteration

1: **Init:** $Q^{(0)}(x, u)$ for all $x \in \mathcal{X}$ and all $u \in \mathcal{U}$

2: **loop**

3: $\pi \leftarrow \epsilon$ -greedy policy derived from Q

4: Generate episodes $\{\rho^{(k)}\}_{k=1}^K$ from π

5: $Q^{(i+1)} = \arg \min_Q \sum_{k=1}^K \sum_{t=0}^T \left\| \mathcal{T}_*[Q^{(i)}](x_t^{(k)}, u_t^{(k)}, x_{t+1}^{(k)}) - Q(x_t^{(k)}, u_t^{(k)}) \right\|^2$

Small Steps in the Backup Direction

- ▶ Full backup: $Q^{(i+1)} \leftarrow \mathcal{T}_*[Q^{(i)}]$
- ▶ Partial backup: $Q^{(i+1)} \leftarrow \epsilon \mathcal{T}_*[Q^{(i)}] + (1 - \epsilon)Q^{(i)}$
- ▶ Equivalent to a gradient step on squared error:

$$\begin{aligned} Q^{(i+1)} \leftarrow \epsilon \mathcal{T}_*[Q^{(i)}] + (1 - \epsilon)Q^{(i)} &= Q^{(i)} - \epsilon \left(Q^{(i)} - \mathcal{T}_*[Q^{(i)}] \right) \\ &= Q^{(i)} - \epsilon \left(\frac{1}{2} \nabla_Q \|Q^{(i)} - \mathcal{T}_*[Q^{(i)}]\|^2 \Big|_{Q=Q^{(i)}} + \text{noise} \right) \end{aligned}$$

- ▶ Behaves like stochastic gradient descent on $L(Q) := \frac{1}{2} \|\mathcal{T}_*[Q] - Q\|^2$ but the objective is changing, i.e., $\mathcal{T}_*[Q^{(i)}]$ is a moving target. Does it converge?
- ▶ **Stochastic Approximation Theory:** a “partial update” to ensure contraction + appropriate step size ϵ implies convergence to the contraction fixed point: $\lim_{i \rightarrow \infty} Q^{(i)} = Q^*$
- ▶ T. Jaakkola, M. Jordan, S. Singh, “On the convergence of stochastic iterative dynamic programming algorithms,” Neural computation, 1994, 3

Least-squares Partial Backup Version

Algorithm 7 Batch Gradient Least-squares Q-Value Iteration

- 1: **Init:** $Q^{(0)}(x, u)$ for all $x \in \mathcal{X}$ and all $u \in \mathcal{U}$
 - 2: **loop**
 - 3: $\pi \leftarrow \epsilon$ -greedy policy derived from Q
 - 4: Generate episodes $\{\rho^{(k)}\}_{k=1}^K$ from π
 - 5:
$$Q^{(i+1)}(x, u) \leftarrow Q^{(i)}(x, u) + \frac{\epsilon}{2} \nabla_Q \sum_{k=1}^K \sum_{t=0}^{T-1} \|\mathcal{T}_*[Q^{(i)}](x_t^{(k)}, u_t^{(k)}, x_{t+1}^{(k)}) - Q(x_t^{(k)}, u_t^{(k)})\|^2$$
-

- ▶ Watkins Q-learning is a special case with $T = 1$