# Reinforcement Learning

## Chapter 3: Model-free RL

### Ali Bereyhi

ali.bereyhi@utoronto.ca

Department of Electrical and Computer Engineering
University of Toronto

Fall 2025

# Control Loop via Temporal Difference

+ *But still we are not fully online! We need to wait till end of each episode!*
– Well! That's right! But, we could use TD!

*Using TD in the control loop will make our algorithm fully online*

- *We update values after each state-action pair*
- *We then improve the policy*

> *We should yet use $\epsilon$-greedy improvement to keep exploration*

# SARSA: *State-Action-Reward State-Action*

## SARSA $\equiv$ *State-Action Reward State-Action*

*SARSA algorithms use TD along with $\epsilon$-greedy update for the control loop*

In general, we can develop various forms of SARSA

- *We may use TD-0 for updating action-values*
    - ↳ *This is the basic SARSA*
- *We may use TD-$n$ for updating action-values*
    - ↳ *This is $n$-SARSA*
- *We may use TD$_\lambda$ for updating action-values*
    - ↳ *This is SARSA($\lambda$)*

# SARSA: *First Try*

*Let's try to make a simple TD-based control loop*

---

```
TD_Control():
```
*1: Initiate estimator as $\hat{q}_\pi(s, a) = 0$ for all states and actions*
*2: **for** episode $= 1 : K$ or until $\pi$ stops changing **do***
*3:     Initiate with a random state-action pair $(S_0, A_0)$*
*4:     **for** $t = 0 : T - 1$ that is either terminal or terminated **do***
*5:         Act $A_t$ and observe*

$$S_t, A_t \xrightarrow{R_{t+1}} S_{t+1}$$

*6:         Update policy to $\pi \leftarrow \epsilon\text{-Greedy}(\hat{q}_\pi)$*
*7:         Draw the new action $A_{t+1}$ from $\pi(\cdot|S_{t+1})$*
*8:         Compute $\hat{v}_\pi(S_{t+1})$ from $\hat{q}_\pi(S_{t+1}, a)$ and $\pi(\cdot|S_{t+1})$*
*9:         Set $G \leftarrow R_{t+1} + \gamma\hat{v}_\pi(S_{t+1})$*
*10:        Update $\hat{q}_\pi(S_t, A_t) \leftarrow \hat{q}_\pi(S_t, A_t) + \alpha(G - \hat{q}_\pi(S_t, A_t))$*
*11:    **end for***
*12: **end for***

---

# SARSA: *Going On-Policy*

In *line 8* of our control algorithm: *we compute $\hat{v}_\pi (S_{t+1})$ as*

$$\hat{v}_\pi (S_{t+1}) = \sum_{m=1}^{M} \pi (a^m | S_{t+1}) \, \hat{q}_\pi (S_{t+1}, a^m)$$

*But, we do know that*

1. *our estimates $\hat{q}_\pi (S_{t+1}, a^m)$ are not that good, and also*
2. *our policy has led use to next action $A_{t+1}$*

*So, we could move on our policy and write*

$$\pi (a | S_{t+1}) = \begin{cases} 1 & a = A_{t+1} \\ 0 & a \neq A_{t+1} \end{cases} \rightsquigarrow \hat{v}_\pi (S_{t+1}) = \hat{q}_\pi (S_{t+1}, A_{t+1})$$

*We call this approach on-policy, since move on our policy*

# SARSA: *Basic Algorithm*

SARSA():
 *1: Initiate estimator as $\hat{q}_\pi(s, a) = 0$ for all states and actions*
 *2: for episode $= 1 : K$ or until $\pi$ stops changing do*
 *3:     Initiate with a random state-action pair $(S_0, A_0)$*
 *4:     for $t = 0 : T - 1$ that is either terminal or terminated do*
 *5:         Act $A_t$ and observe*

$$S_t, A_t \xrightarrow{R_{t+1}} S_{t+1}$$

 *6:         Update policy to $\pi \leftarrow \epsilon\text{-Greedy}(\hat{q}_\pi)$*
 *7:         Draw the new action $A_{t+1}$ from $\pi(\cdot|S_{t+1})$ and move on policy*

$$S_t, A_t \xrightarrow{R_{t+1}} S_{t+1}, A_{t+1}$$

 *8:         Set $G \leftarrow R_{t+1} + \gamma\hat{q}_\pi(S_{t+1}, A_{t+1})$*
 *9:         Update $\hat{q}_\pi(S_t, A_t) \leftarrow \hat{q}_\pi(S_t, A_t) + \alpha(G - \hat{q}_\pi(S_t, A_t))$*
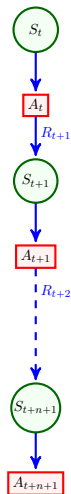 *10:     end for*
 *11: end for*

# SARSA: *Deeper Return Samples*



*We can use a longer trajectory while we learn on-policy, i.e.,*

$$G^n = \sum_{i=0}^{n} R_{t+i+1} + \gamma \hat{q}_\pi \left( S_{t+n+1}, A_{t+n+1} \right)$$

*This will however add extra delay!*

*As a practice, you could*

*re-write the basic SARSA with $n$-return* ☺

# SARSA($\lambda$):*Tracing Eligibility of State-Action Pairs*

We can extend SARSA to the case with $\lambda$-return: *we have two options*

- *the case with forward-view*
  - ↳ *We know this is not practical! So, let's skip the details*
- *the case with backward-view and eligibility tracing*
  - ↳ *Let's look into this one*

---

*We first extend eligibility tracing to the case with state-action pairs*

```
ElgTrace(S_t, A_t, E(·) |λ):
  1: Eligibility tracing function has NM components, i.e., E(s,a) for all state-action pairs
  2: for all state-action pairs (s,a) do
  3:     Update E(s,a) ← γλE(s,a)
  4: end for
  5: Update E(S_t, A_t) ← E(S_t, A_t) + 1
```

# SARSA: *Alternative via TD-$\lambda$*

---

SARSA($\lambda$):

  *1:* Initiate $\hat{q}_\pi(s, a) = 0$ and $E(s, a) = 0$ for *all states* and *actions*

  *2:* **for** episode $= 1 : K$ or until $\pi$ stops changing **do**

  *3:*     Initiate with a *random state-action pair* $(S_0, A_0)$

  *4:*     **for** $t = 0 : T - 1$ that is either terminal or terminated **do**

  *5:*         $E(\cdot) \leftarrow \texttt{ElgTrace}(S_t, A_t, E(\cdot) \mid \lambda)$

  *6:*         Act $A_t$ and observe $R_{t+1}$ and $S_{t+1}$

  *7:*         Update policy to $\pi \leftarrow \epsilon\texttt{-Greedy}(\hat{q}_\pi)$

  *8:*         Draw the new action $A_{t+1}$ from $\pi(\cdot | S_{t+1})$ and move on policy

$$S_t, A_t \xrightarrow{R_{t+1}} S_{t+1}, A_{t+1}$$

  *9:*         Set $\Delta \leftarrow R_{t+1} + \gamma\hat{q}_\pi(S_{t+1}, A_{t+1}) - \hat{q}_\pi(S_t, A_t)$

*10:*         **for** all state-*action* pairs $(s, a)$ **do**

*11:*            Update $\hat{q}_\pi(s, a) \leftarrow \hat{q}_\pi(s, a) + \alpha\Delta E(s, a)$

*12:*         **end for**

*13:*     **end for**

*14:* **end for**

---

# Going Off-Policy

Let's think about a fundamental question: *while sampling the environment with a specific policy $\pi$, can we estimate the values of another policy $\bar{\pi}$?*

+ *Why should this be a fundamental question?*
- Well! There are several reasons
    ↳ *Maybe we sampled environment with our bad policy: can't we use our sample again?*
    ↳ *Maybe we are looking at other players: can't we learn something about the environment from their samples?*
        ↳ *Maybe they are good players: can't we use this fact to improve our policy?*
        ↳ *Maybe they are bad players: can't we use this fact to avoid doing mistakes?*

---

*This is the idea of off-policy control*

> *Let's start with some baiscs*

## Importance Sampling

Consider following problem: *we have random variable $X$ drawn as $X \sim p(x)$ whose mean is*

$$\mu_p = \mathbb{E}_p \{X\} = \sum_x p(x) x$$

*We want to know how would be the expectation if we had $X \sim q(x)$: we write*

$$\mu_q = \mathbb{E}_q \{X\} = \sum_x q(x) x$$
$$= \sum_x p(x) \frac{q(x)}{p(x)} x = \mathbb{E}_p \left\{ \frac{q(X)}{p(X)} X \right\}$$

*This gives us possibility to*

*estimate $\mathbb{E}_q \{X\}$ using samples drawn from $p(x)$*

# Importance Sampling

*Say we have drawn $K$ samples from $p(x)$, i.e., we have*

$$X_1, X_2, \ldots, X_K$$

*We can use Monte-Carlo to estimate $\mu_p$ as*

$$\hat{\mu}_p = \frac{1}{K} \sum_{k=1}^{K} X_k$$

*We can also use Monte-Carlo to estimate $\mu_q$ as*

$$\hat{\mu}_q = \frac{1}{K} \sum_{k=1}^{K} \frac{q(X_k)}{p(X_k)} X_k$$

*We call this method importance sampling*

# Off-Policy Control via Importance Sampling

Now, let's get back to our problem: *assume we have played with policy $\pi$ and collected $K$ sample trajectories of length $T$ all started at state $S_0 = s$, i.e.,*

$$s = S_0\left[k\right], A_0\left[k\right] \xrightarrow{R_1[k]} S_1\left[k\right], A_1\left[k\right] \xrightarrow{R_2[k]} \cdots \xrightarrow{R_T[k]} S_T\left[k\right]$$

*for $k = 1 : K$; then, we could write*

$$\hat{v}_\pi\left(s\right) = \frac{1}{K} \sum_{k=1}^{K} G\left[k\right]$$

*This is the basic Monte-Carlo*

# Off-Policy Control via Importance Sampling

*But now, we want to use samples to evaluate another policy $\bar{\pi}$*

$$s = S_0\left[k\right], A_0\left[k\right] \xrightarrow{R_1[k]} S_1\left[k\right], A_1\left[k\right] \xrightarrow{R_2[k]} \cdots \xrightarrow{R_T[k]} S_T\left[k\right]$$

*We could also use importance sampling to write*

$$\hat{v}_{\bar{\pi}}\left(s\right) = \frac{1}{K}\sum_{k=1}^{K}\frac{\mathrm{Pr}\left\{\text{same action sequence with } \bar{\pi}\right\}}{\mathrm{Pr}\left\{\text{same action sequence with } \pi\right\}}G\left[k\right]$$

$$= \frac{1}{K}\sum_{k=1}^{K}\frac{\bar{\pi}\left(A_0\left[k\right]|S_0\left[k\right]\right)\cdots\bar{\pi}\left(A_{T-1}\left[k\right]|S_{T-1}\left[k\right]\right)}{\pi\left(A_0\left[k\right]|S_0\left[k\right]\right)\cdots\pi\left(A_{T-1}\left[k\right]|S_{T-1}\left[k\right]\right)}G\left[k\right]$$

$$= \frac{1}{K}\sum_{k=1}^{K}\prod_{\ell=0}^{T-1}\frac{\bar{\pi}\left(A_\ell\left[k\right]|S_\ell\left[k\right]\right)}{\pi\left(A_\ell\left[k\right]|S_\ell\left[k\right]\right)}G\left[k\right]$$

# Off-Policy Control via Importance Sampling

*We can further update the estimate in an online fashion from*

$$S_t, A_t \xrightarrow{R_{t+1}} S_{t+1}, A_{t+1} \xrightarrow{R_{t+2}} \cdots \xrightarrow{R_T} S_T$$

*by online averaging as*

$$\hat{v}_{\bar{\pi}}\left(S_t\right) \leftarrow \hat{v}_{\bar{\pi}}\left(S_t\right) + \alpha \left( \prod_{\ell=t}^{T-1} \frac{\bar{\pi}\left(A_\ell | S_\ell\right)}{\pi\left(A_\ell | S_\ell\right)} G_t - \hat{v}_{\bar{\pi}}\left(S_t\right) \right)$$

*So, we are evaluating $\bar{\pi}$ via Monte-Carlo*

*off our policy $\pi$*

*This is off-policy control*

# Off-Policy Control via Importance Sampling

*We can further apply off-policy control via TD*

$$\hat{v}_{\bar{\pi}}\left(S_t\right) \leftarrow \hat{v}_{\bar{\pi}}\left(S_t\right) + \alpha \left(\frac{\bar{\pi}\left(A_t|S_t\right)}{\pi\left(A_t|S_t\right)}\left(R_{t+1} + \gamma\hat{v}_{\bar{\pi}}\left(S_{t+1}\right)\right) - \hat{v}_{\bar{\pi}}\left(S_t\right)\right)$$

*Note that for action-values estimate*

$R_{t+1}$ *does not depend any more on policy as we know action* $A_t$

*Therefore, we have for action-value update*

$$\hat{q}_{\bar{\pi}}\left(S_t, A_t\right) \leftarrow \hat{q}_{\bar{\pi}}\left(S_t, A_t\right) + \alpha \left(R_{t+1} + \gamma\frac{\bar{\pi}\left(A_t|S_t\right)}{\pi\left(A_t|S_t\right)}\hat{v}_{\bar{\pi}}\left(S_{t+1}\right) - \hat{q}_{\bar{\pi}}\left(S_t, A_t\right)\right)$$

# Q-Learning

## Q-Learning

*Q-learning is an off-policy TD control algorithm, where we sample with $\epsilon$-greedy policy but update the action-values to evaluate greedy policy*

*This means in Q-learning $\pi$ is $\epsilon$-greedy policy and $\bar{\pi}$ is greedy. Let's consider basic TD evaluation: so, we can write*

$$\hat{q}_{\bar{\pi}}\left(S_t, A_t\right) \leftarrow \hat{q}_{\bar{\pi}}\left(S_t, A_t\right) + \alpha(G - \hat{q}_{\bar{\pi}}\left(S_t, A_t\right))$$

*where $G$ should be*

$$G = R_{t+1} + \gamma\hat{v}_{\bar{\pi}}\left(\bar{S}_{t+1}\right)$$

*Since we sample by $\epsilon$-greedy policy $\pi$, we use importance sampling and write*

$$G = R_{t+1} + \gamma\frac{\bar{\pi}\left(A_t|S_t\right)}{\pi\left(A_t|S_t\right)}\hat{v}_{\bar{\pi}}\left(S_{t+1}\right)$$

# Q-Learning

*But, we really don't need importance sampling: we can simply observe that*

$$\hat{v}_{\bar{\pi}}\left(S_{t+1}\right) = \sum_{m=1}^{M} \hat{q}_{\bar{\pi}}\left(S_{t+1}, a^m\right) \bar{\pi}\left(a^m | S_{t+1}\right) = \max_{m} \hat{q}_{\bar{\pi}}\left(S_{t+1}, a^m\right)$$

*and we do know that*

$$\frac{\bar{\pi}\left(A_t | S_t\right)}{\pi\left(A_t | S_t\right)} = \mathbf{1}\left\{A_t = \operatorname*{argmax}_{a} \hat{q}_{\bar{\pi}}\left(S_t, a\right)\right\}$$

*So, we could directly update as*

$$\hat{q}_{\bar{\pi}}\left(S_t, A_t\right) \leftarrow \hat{q}_{\bar{\pi}}\left(S_t, A_t\right) + \alpha\left(R_{t+1} + \gamma \max_{m} \hat{q}_{\bar{\pi}}\left(S_{t+1}, a^m\right) - \hat{q}_{\bar{\pi}}\left(S_t, A_t\right)\right)$$

*This concludes Q-learning algorithm*
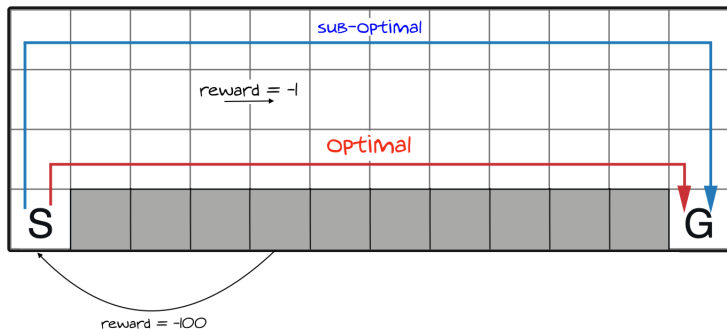
# Q-Learning: *Basic Algorithm*

Q-Learning():
 *1: Initiate estimator as $\hat{q}_\star(s, a) = 0$ for all states and actions*
 *2: **for** episode $= 1 : K$ or until $\pi$ stops changing **do***
 *3:     Initiate with a random state $S_0$*
 *4:     **for** $t = 0 : T - 1$ that is either terminal or terminated **do***
 *5:         Update policy to $\pi \leftarrow \epsilon\text{-Greedy}(\hat{q}_\star)$*
 *6:         Draw action $A_t$ from $\pi(\cdot|S_t)$ and observe*

$$S_t, A_t \xrightarrow{R_{t+1}} S_{t+1}$$
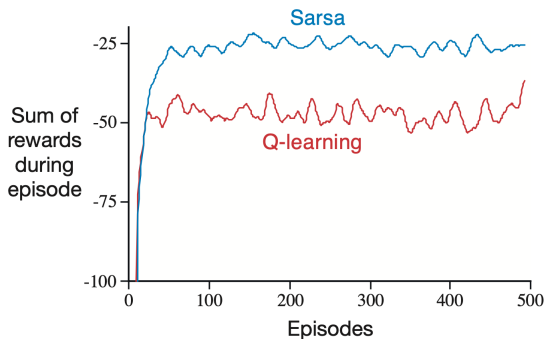
 *7:         Set $G \leftarrow R_{t+1} + \gamma \max_m \hat{q}_\star(S_{t+1}, a^m)$*
 *8:         Update $\hat{q}_\star(S_t, A_t) \leftarrow \hat{q}_\star(S_t, A_t) + \alpha(G - \hat{q}_\star(S_t, A_t))$*
 *9:     **end for***
*10: **end for***

# Example: *Cliff Walking*



*Let's compare SARSA to Q-Learning algorithm!*

# Example: *Cliff Walking*



## Don't Mistake!

*Q-learning collects less reward since it goes off-policy; however, it estimates optimal action-values: at some point it can start playing optimally*

# Convergence of SARSA

## Recall: *GLIE Algorithms*

*A GPI-type control loop is GLIE, if for any **state**-**action** pair $(s, a)$, we have the following asymptotic properties*

① *The number of visits to all **state**-**action** pair grows large*

$$\lim_{K \to \infty} \mathcal{C}_K (s, a) = \infty$$

② *The improved policy in last episode converges to greedy policy*

$$\lim_{K \to \infty} \pi_K (a^m | s) = \begin{cases} 1 & m = \underset{m}{\operatorname{argmax}} \, q_{\pi_K} (s, a^m) \\ 0 & m \neq \underset{m}{\operatorname{argmax}} \, q_{\pi_K} (s, a^m) \end{cases}$$

*GLIE control algorithms converge to optimal policy*

# Convergence of SARSA

+ *But do we really have large number of episodes with SARSA?*

– Not necessarily! We may have only one <span style="color:red">infinitely long</span> trajectory

+ *What should we do then?*

– We can simply treat it as a large number of episodes of length $1$

---

*In (basic) SARSA, we only need one step in the trajectory*

$$S_t, A_t \xrightarrow{R_{t+1}} S_{t+1}$$

*We could hence think of it as one episode*

1. *each time step $t$ we update the action-values*
2. *each time step we improve the policy*

# Convergence of SARSA

### Modification: *GLIE Algorithms*

*An online control loop is GLIE, if we have asymptotically in time $t$*

1. *The number of visits to all state-action pair grows large*

$$\lim_{t \to \infty} \mathcal{C}_t \left( s, a \right) = \infty$$

2. *The improved policy converges to greedy policy*

$$\lim_{t \to \infty} \pi_t \left( a^m | s \right) = \begin{cases} 1 & m = \operatorname*{argmax}_m q_{\pi_t} \left( s, a^m \right) \\ 0 & m \neq \operatorname*{argmax}_m q_{\pi_t} \left( s, a^m \right) \end{cases}$$

## Convergence of SARSA: *Make it GLIE*

+ *Can we guarantee that both conditions hold with SARSA?*

- The second one is easy: *we need to scale $\epsilon$ down with $t$, e.g., $\epsilon_t = 1/t$*

+ *What about the first condition?*

- We should scale *the step-size $\alpha$ according to Robbins-Monro*

Robbins-Monro Sequence

*Sequence $\alpha_t$ is Robbins-Monro if we have*

$$\sum_{t=0}^{\infty} \alpha_t = \infty \qquad and \qquad \sum_{t=0}^{\infty} \alpha_t^2 < \infty$$

*For instance, $\alpha_t = 1/t$ is a Robbins-Monro sequence*

# Convergence of SARSA

## Convergence of SARSA

*SARSA online control loop converges to the optimal action-values if*

1. *Step-size is scheduled by a Robbins-Monro sequence*
2. *Exploration factor $\epsilon$ decays in time*

*In practice however*

- *$\epsilon$ is a hyperparameter*
  - ↳ *We know that we should schedule it*
  - ↳ *How we should do the scheduling? This is hyperparameter tuning*
- *$\alpha$ is a hyperparameter: some people call it learning rate*
  - ↳ *Its scheduling is again hyperparameter tuning*

# Convergence of Q-Learning

## Convergence of Q-Learning

*Q-learning online control loop with exploration (non-zero $\epsilon$) converges to the optimal action-values as $t \to \infty$*

+ *That's it?*

– Yes!

> *Since we are evaluating off-policy, we don't care about behaving policy*

# Q-Learning vs SARSA

+ *So! Does it mean that Q-learning is always better?*

– *Not always!*

*In general Q-learning has several benefits*

- *Minimal convergence requirements*
- *It converges faster to the optimal policy*
  - ↳ *If we want to make SARSA that fast, we may get to a sub-optimal policy*
- *It has more flexibility and sample-efficiency*

*But, SARSA also has some benefits*

- *It is better suited for online control*
  - ↳ *Our behaving policy is the one going towards optimal one*
  - ↳ *In Q-learning, the behaving policy is not the optimal one*
- *It has lower complexity*
  - ↳ *We just deal with one policy*

# End of Story!



Model-Based RL
Bellman Equation
value iteration
policy iteration

Model-free RL
on-policy methods
temporal difference
Monte Carlo
SARSA

off-policy methods
Q-learning

*I would strongly suggest to start with programming part of Assignment 2!*

*There you solve Froozen Lake with SARSA and Q-Learning*