

# Reinforcement learning

# Reinforcement learning

- **Model Based RL:** Have the transition probabilities  $P(s', s, a)$  and solve the Markov Decision Process (MDP).
- **Model-free RL:** Don't have transition probability & directly learn what action to do when.
  - *Gradient free:*
  - *Gradient based:*

# Reinforcement learning

- **Model Based RL:** Have the transition probabilities  $P(s', s, a)$  and solve the Markov Decision Process (MDP).
  - Policy Iteration and Value Iteration algorithms
- **Model-free RL:** Don't have transition probability & directly learn what action to do when.
  - *Gradient based:*

# Reinforcement learning

- **Model Based RL:** Have the transition probabilities  $P(s', s, a)$  and solve the Markov Decision Process (MDP).
  - Policy Iteration and Value Iteration algorithms
- **Model-free RL:** Don't have transition probability & directly learn what action to do when.
  - *Gradient free:*
    - Off Policy: different policy to select the next action and value evaluation, e.g., Q-Learning
    - On Policy: same policy to select the next action and value evaluation
      - SARSA (State–action–reward–state–action)
      - Monte Carlo (Model-free): play many games to estimate policy
  - *Gradient based:*

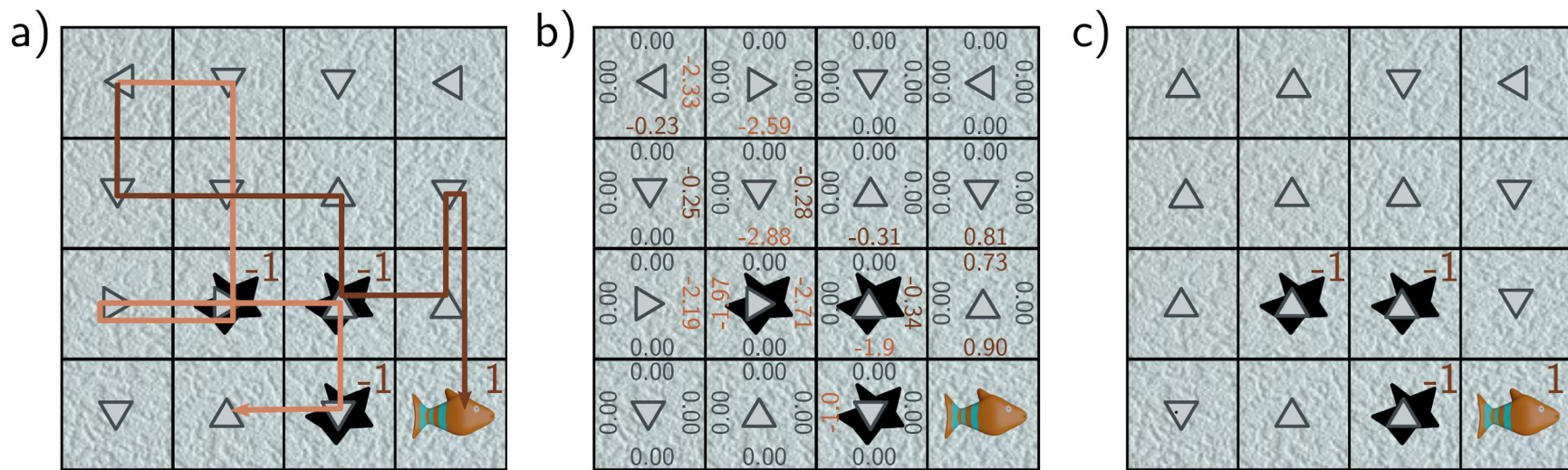
# Reinforcement learning

- **Model Based RL:** Have the transition probabilities  $P(s', s, a)$  and solve the Markov Decision Process (MDP).
  - Policy Iteration and Value Iteration algorithms
- **Model-free RL:** Don't have transition probability & directly learn what action to do when.
  - *Gradient free:*
    - Off Policy: different policy to select the next action and value evaluation, e.g., Q-Learning
    - On Policy: same policy to select the next action and value evaluation
      - SARSA (State–action–reward–state–action)
      - Monte Carlo (Model-free): play many games to estimate policy
  - *Gradient based:*
    - Policy Gradient Optimization, e.g., Deep reinforcement learning for policy network

# Monte Carlo methods

- To estimate the action values  $q[s, a]$ , a series of episodes are run.
- Each starts with a given state and action and thereafter follows the current policy, producing a series of actions, states, and returns (figure 19.11a).
- The action value for a given state-action pair under the current policy is estimated as the average of the empirical returns that follow after each time this pair is observed (figure 19.11b).
- Then the policy is updated by choosing the action with the maximum value at every state

$$\pi[a \mid s] \leftarrow \operatorname{argmax}_a [q[s, a]].$$



**Figure 19.11** Monte Carlo methods. a) The policy (arrows) is initialized randomly. The MDP is repeatedly simulated, and the trajectories of these episodes are stored (orange and brown paths represent two trajectories). b) The action values are empirically estimated based on the observed returns averaged over these trajectories. In this case, the action values were all initially zero and have been updated where an action was observed. c) The policy can then be updated according to the action which received the best (or least bad) reward.

# Monte Carlo methods: on-policy

- This is an on-policy method; the current best policy is used to guide the agent through the environment.
- This policy is based on the observed action values in every state, but it's not possible to estimate the value of actions that haven't been used.
- A different approach is to use an **epsilon greedy policy**, in which a random action is taken with probability  $\epsilon$ , and the optimal action is allotted the remaining probability.
- The choice of  $\epsilon$  trades off exploitation and exploration.



# Monte Carlo methods: off-policy

- In off-policy methods, the optimal policy  $\pi$  (the target policy) is learned based on episodes generated by a different behavior policy  $\pi'$ .
- Some off-policy methods explicitly use importance sampling to estimate the action value under policy  $\pi$  using samples from  $\pi'$ .
- Others, such as Q-learning (described in the next section), estimate the values based on the greedy action

# Temporal-difference (TD) learning

Temporal-difference methods are similar to

**Dynammic programming** update estimates based in part on other learned estimates, without waiting for the final outcome (they bootstrap)

**Monte Carlo methods** learn directly from raw experience without a model of the environment's dynamics

# Temporal difference methods

- ▶ TD methods only wait until the next time step to update the value estimates.
- ▶ At time  $t + 1$  they immediately form a target and make an update using the observed reward  $r_{t+1}$  and the current estimate  $V(s_{t+1})$ .

$$V(s_t) \leftarrow V(s_t) + \alpha (r_{t+1} + \gamma V(s_{t+1}) - V(s_t)),$$

where  $\alpha > 0$  is a step-size parameter.

# Q- Learning - Off Policy TD Learning

By contrast, Q-Learning is an off-policy algorithm with update

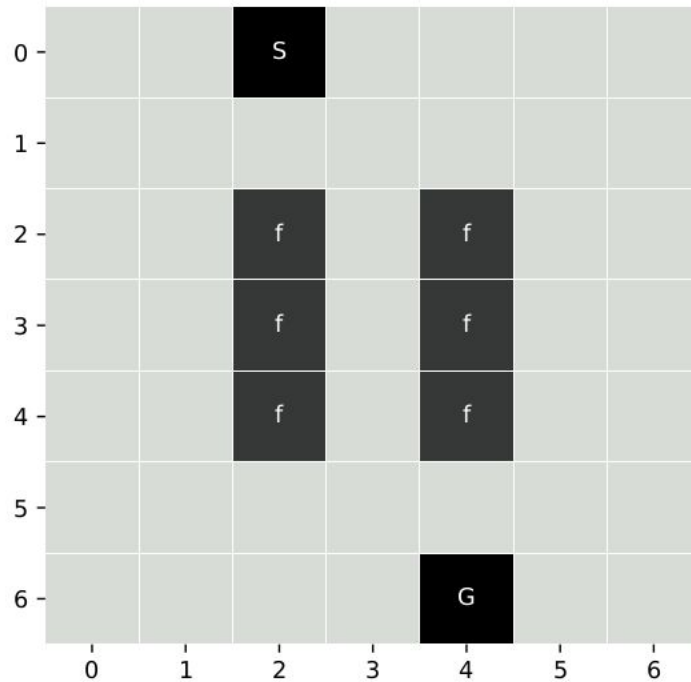
$$Q[s_t, a_t] \leftarrow Q[s_t, a_t] + \alpha \left( \overbrace{r[s_t, a_t] + \gamma \cdot \max_a [Q[s_{t+1}, a]] - Q[s_t, a_t]}^{\text{TD Error}} \right)$$

where now the choice of action at each step is derived from a different behavior policy  $\pi'$ .

- In both cases, the policy is updated by taking the maximum of the action values at each state
- It can be shown the action values will eventually converge, assuming that every state-action pair is visited an infinite number of times.

# Q Learning - Off Policy TD Learning

- 1: Initialize  $Q(s, a), \forall s \in S, a \in A$   $t = 0$ , initial state  $s_t = s_0$
- 2: Set  $\pi_b$  to be  $\epsilon$ -greedy w.r.t.  $Q$
- 3: loop
- 4:     Take  $a_t \sim \pi_b(s_t)$  // Sample action from policy
- 5:     Observe  $(r_t, s_{t+1})$
- 6:     Update  $Q[s_t, a_t] \leftarrow Q[s_t, a_t] + \alpha (r[s_t, a_t] + \gamma \cdot \max_a [Q[s_{t+1}, a]] - Q[s_t, a_t])$
- 7:     Policy improvement  $\pi(s_t) = \arg \max_a Q(s_t, a)$  w.prob  $1 - \epsilon$ , else random
- 8:  $t = t + 1$  (until terminal state not reached)
- 9: end loop



**Figure 4:** A simple grid world. The start position is annotated with an *S* and the goal with a *G*. The squares marked *f* denote fire, which is very damaging to an agent.

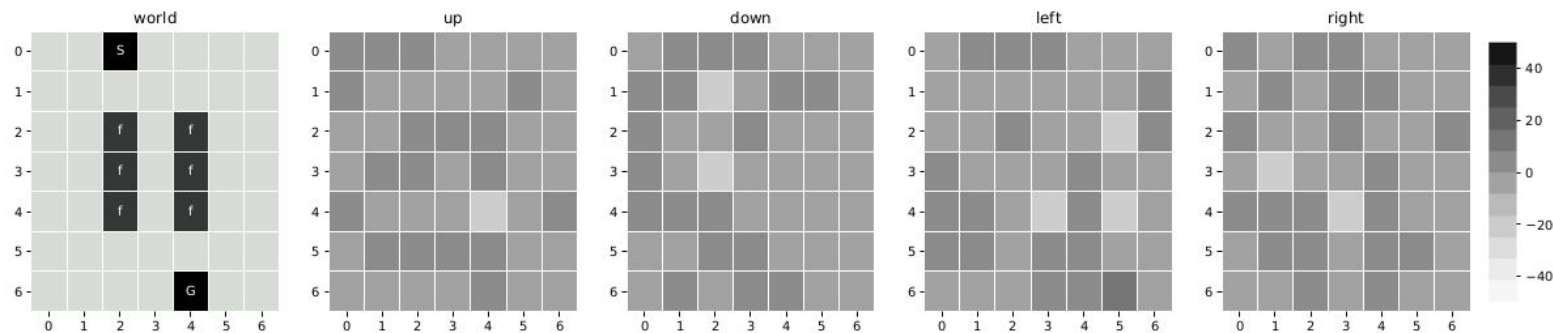
**Table 3:** A portion of the action-value table for the grid world example at its first initialization.

State	Action	Value	State	Action	Value	State	Action	Value
0-0	<i>up</i>	0.933		...			...	
0-0	<i>down</i>	-0.119	2-0	<i>left</i>	-0.691	6-2	<i>right</i>	0.201
0-0	<i>left</i>	-0.985	2-0	<i>right</i>	0.668	6-3	<i>up</i>	-0.588
0-0	<i>right</i>	0.822	2-1	<i>up</i>	-0.918	6-3	<i>down</i>	0.038
0-1	<i>up</i>	0.879	2-1	<i>down</i>	-0.228	6-3	<i>left</i>	0.859
0-1	<i>down</i>	0.164	2-1	<i>left</i>	-0.301	6-3	<i>right</i>	-0.085
0-1	<i>left</i>	0.343	2-1	<i>right</i>	-0.317	6-4	<i>up</i>	0.000
0-1	<i>right</i>	-0.832	2-2	<i>up</i>	0.633	6-4	<i>down</i>	0.000
0-2	<i>up</i>	0.223	2-2	<i>down</i>	-0.048	6-4	<i>left</i>	0.000
0-2	<i>down</i>	0.582	2-2	<i>left</i>	0.566	6-4	<i>right</i>	0.000
0-2	<i>left</i>	0.672	2-2	<i>right</i>	-0.058	6-5	<i>up</i>	0.321
0-2	<i>right</i>	0.084	2-3	<i>up</i>	0.635	6-5	<i>down</i>	-0.793
0-3	<i>up</i>	-0.308	2-3	<i>down</i>	0.763	6-5	<i>left</i>	-0.267
0-3	<i>down</i>	0.247	2-3	<i>left</i>	-0.121	6-5	<i>right</i>	0.588
0-3	<i>left</i>	0.963	2-3	<i>right</i>	0.562	6-6	<i>up</i>	-0.870
0-3	<i>right</i>	0.455	2-4	<i>up</i>	0.629	6-6	<i>down</i>	-0.720
0-4	<i>up</i>	-0.634	2-4	<i>down</i>	-0.409	6-6	<i>left</i>	0.811
	...			...		6-6	<i>right</i>	0.176

$$Q(0-3, left) \leftarrow Q(0-3, left) + \alpha \times (R(0-3, left) + \gamma \times Q(\quad) - Q(0-3, left))$$

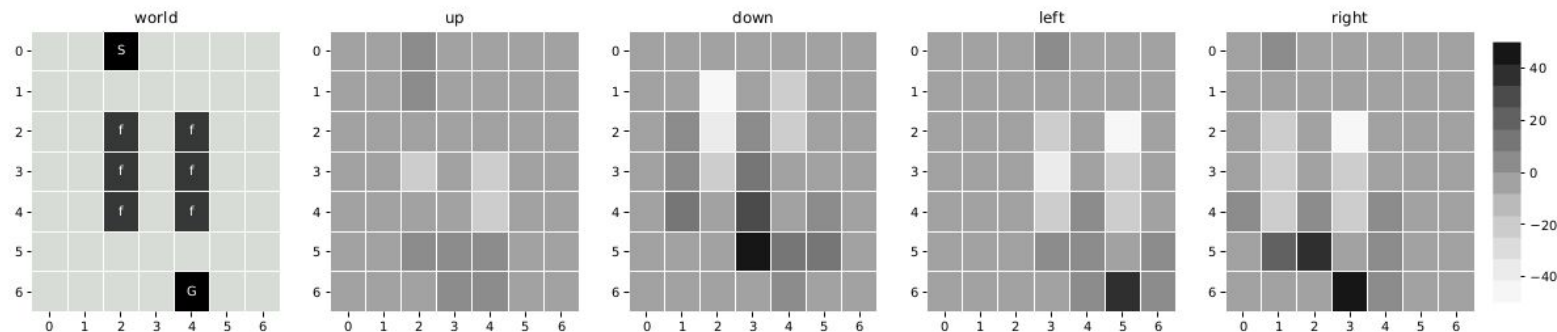


$$\begin{aligned}
 Q(0-3, left) &\leftarrow Q(0-3, left) + \alpha \times (R(0-3, left) + \gamma \times Q(0-2, left) - Q(0-3, left)) \\
 &0.963 + 0.2 \times (-1 + 0.9 \times 0.672 - 0.963) \\
 &0.691
 \end{aligned}$$



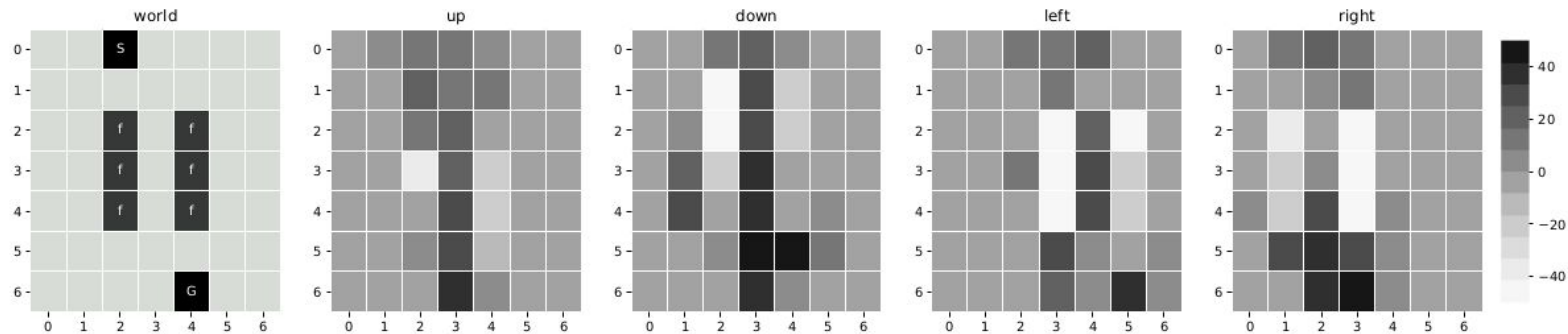
(a) Action-value table after 1 episode

**Figure 5:** (a)–(c) The evolution of the entries in the action-value table over episodes of Q-learning off-policy temporal-difference learning across the grid world.



(b) Action-value table after 35 episodes

**Figure 6:** (a)–(c) The evolution of the entries in the action-value table over episodes of Q-learning off-policy temporal-difference learning across the grid world.



(c) Action-value table after 350 episodes

**Figure 7:** (a)–(c) The evolution of the entries in the action-value table over episodes of Q-learning off-policy temporal-difference learning across the grid world.



# SARSA - On Policy TD Learning

*SARSA* (State-Action-Reward-State-Action) is an on-policy algorithm with update:

$$Q[s_t, a_t] \leftarrow Q[s_t, a_t] + \alpha \left( \overbrace{r[s_t, a_t] + \gamma \cdot Q[s_{t+1}, a_{t+1}] - Q[s_t, a_t]}^{\text{TD Error}} \right)$$

where  $\alpha \in \mathbb{R}^+$  is the learning rate.

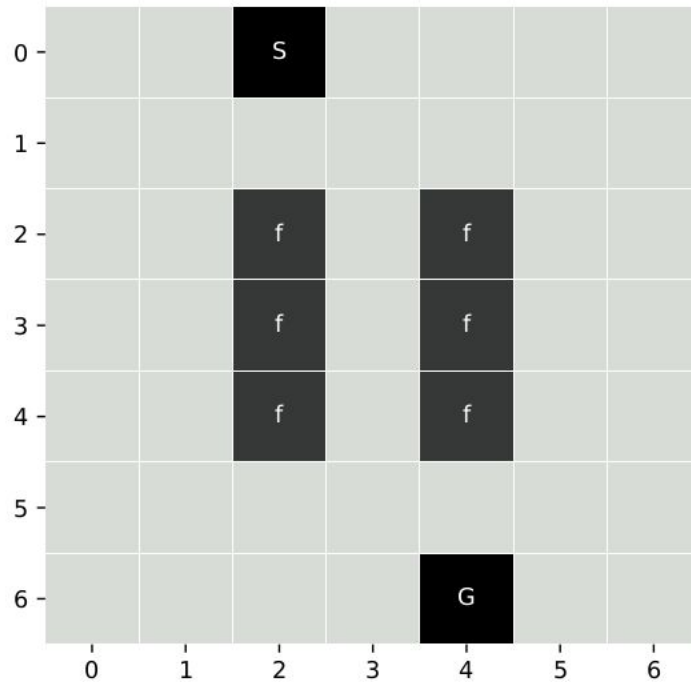
- The bracketed term is called the TD error and measures the consistency between the estimated action value  $Q[s_t, a_t]$  and the estimate  $r[s_t, a_t] + \gamma \cdot Q[s_{t+1}, a_{t+1}]$  after taking a single step.

# SARSA - On Policy TD Learning

Initialize all  $Q(,)$  arbitrarily

For all episodes

- Initialize  $s_t$
- Choose  $a_t$  using policy derived from  $Q$ , e.g.,  $\epsilon$ -greedy
- Repeat
  - Take action  $a_t$ , observe  $r_t$  and  $s_{t+1}$
  - Choose  $a_{t+1}$  using policy derived from  $Q$ , e.g.,  $\epsilon$ -greedy
  - Update  $Q[s_t, a_t] \leftarrow Q[s_t, a_t] + \alpha (r[s_t, a_t] + \gamma \cdot Q[s_{t+1}, a_{t+1}] - Q[s_t, a_t])$
  - $s_t \leftarrow s_{t+1}, a_t \leftarrow a_{t+1}$
- Until  $s_t$  is terminal state



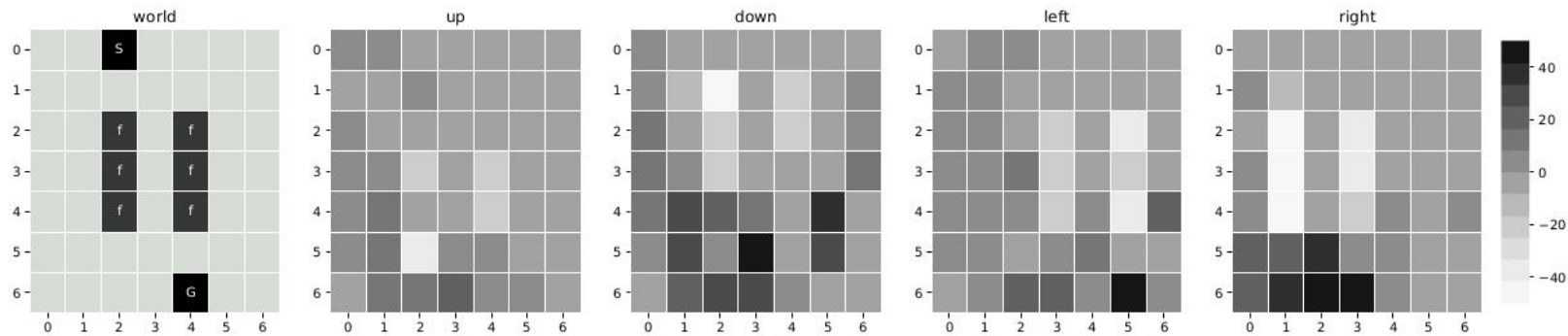
**Figure 4:** A simple grid world. The start position is annotated with an *S* and the goal with a *G*. The squares marked *f* denote fire, which is very damaging to an agent.



**Table 3:** A portion of the action-value table for the grid world example at its first initialization.

State	Action	Value	State	Action	Value	State	Action	Value
0-0	<i>up</i>	0.933		...			...	
0-0	<i>down</i>	-0.119	2-0	<i>left</i>	-0.691	6-2	<i>right</i>	0.201
0-0	<i>left</i>	-0.985	2-0	<i>right</i>	0.668	6-3	<i>up</i>	-0.588
0-0	<i>right</i>	0.822	2-1	<i>up</i>	-0.918	6-3	<i>down</i>	0.038
0-1	<i>up</i>	0.879	2-1	<i>down</i>	-0.228	6-3	<i>left</i>	0.859
0-1	<i>down</i>	0.164	2-1	<i>left</i>	-0.301	6-3	<i>right</i>	-0.085
0-1	<i>left</i>	0.343	2-1	<i>right</i>	-0.317	6-4	<i>up</i>	0.000
0-1	<i>right</i>	-0.832	2-2	<i>up</i>	0.633	6-4	<i>down</i>	0.000
0-2	<i>up</i>	0.223	2-2	<i>down</i>	-0.048	6-4	<i>left</i>	0.000
0-2	<i>down</i>	0.582	2-2	<i>left</i>	0.566	6-4	<i>right</i>	0.000
0-2	<i>left</i>	0.672	2-2	<i>right</i>	-0.058	6-5	<i>up</i>	0.321
0-2	<i>right</i>	0.084	2-3	<i>up</i>	0.635	6-5	<i>down</i>	-0.793
0-3	<i>up</i>	-0.308	2-3	<i>down</i>	0.763	6-5	<i>left</i>	-0.267
0-3	<i>down</i>	0.247	2-3	<i>left</i>	-0.121	6-5	<i>right</i>	0.588
0-3	<i>left</i>	0.963	2-3	<i>right</i>	0.562	6-6	<i>up</i>	-0.870
0-3	<i>right</i>	0.455	2-4	<i>up</i>	0.629	6-6	<i>down</i>	-0.720
0-4	<i>up</i>	-0.634	2-4	<i>down</i>	-0.409	6-6	<i>left</i>	0.811
	...			...		6-6	<i>right</i>	0.176

$$\begin{aligned}
 Q(0-3, left) &\leftarrow Q(0-3, left) + \alpha \times (R(0-3, left) + \gamma \times Q(0-2, down) - Q(0-3, left)) \\
 &0.963 + 0.2 \times (-1 + 0.9 \times 0.582 - 0.963) \\
 &0.675
 \end{aligned}$$



(a) Action-value table after 350 episodes

**Figure 9:** (a) A visualization of the final action-value table for an agent trained using SARSA on-policy temporal-difference learning across the grid world after 350 episodes.

# Summary

Algorithm	Description	Policy	Action Space	State Space	Operator
Monte Carlo (Model-free)	Estimate $Q_{\pi}^{\text{avg}}(s, a)$ average of all the trajectories	Either	Discrete	Discrete	Sample-means
Q-learning	Aim to learn $Q_{\pi}(s, a)$	Off-policy	Discrete	Discrete	Q-value
SARSA	Aim to learn $Q_{\pi}(s, a)$	On-policy	Discrete	Discrete	Q-value