

# Chapter 9

## Q-Learning

### Reading

1. Sutton & Barto, Chapter 6, 11
2. Human-level control through deep reinforcement learning  
<https://www.nature.com/articles/nature14236>
3. Deterministic Policy Gradient Algorithms,  
<http://proceedings.mlr.press/v32/silver14.html>
4. Addressing Function Approximation Error in Actor-Critic Methods  
<https://arxiv.org/abs/1802.09477>
5. An Application of Reinforcement Learning to Aerobatic Helicopter Flight, <https://papers.nips.cc/paper/3151-an-application-of-reinforcement-learning-to-aerobatic-helicopter-flight>

In the previous chapter, we looked at what are called “on-policy” methods, these are methods where the current controller  $u_{\theta^k}$  is used to draw fresh data from the dynamical system and used to update to parameters  $\theta^k$ . *The key inefficiency in on-policy methods is that this data is thrown away in the next iteration.* We need to draw a fresh set of trajectories from the system for  $u_{\theta^{k+1}}$ . This lecture will discuss off-policy methods which are a way to reuse past data. These methods require much fewer data than on-policy methods (in practice, about 10–100× less).

### 9.1 Tabular Q-Learning

Recall the value iteration algorithm for discrete (and finite) state and control spaces; this is also called “tabular” Q-Learning in the RL literature because we can store the Q-function  $q(x, u)$  as a large table with number of rows being the number of states and number of columns being the number of controls, with

16 each entry in this table being the value  $q(x, u)$ . Value iteration when written  
 17 using the Q-function at the  $k^{\text{th}}$  iteration for the tabular setting looks like

$$\begin{aligned} q^{(k+1)}(x, u) &= \sum_{x' \in X} \mathbf{P}(x' | x, u) \left( r(x, u) + \gamma \max_{u'} q^{(k)}(x', u') \right) \\ &= \mathbf{E}_{x' \sim \mathbf{P}(\cdot | x, u)} \left[ r(x, u) + \gamma \max_{u'} q^{(k)}(x', u') \right]. \end{aligned}$$

In the simplest possible instantiation of Q-learning, the expectation in the value iteration above (which we can only compute if we know a model of the dynamics) is replaced by samples drawn from the environment.

18 We will imagine the robot as using an *arbitrary* controller

$$u_e(\cdot | x)$$

19 that has a fairly large degree of randomness in how it picks actions. We call  
 20 such a controller an “exploratory controller”. Conceptually, its goal is to lead  
 21 the robot to diverse states in the state-space so that we get a faithful estimate  
 22 of the expectation in value iteration. We maintain the value  $q^{(k)}(x, u)$  for all  
 23 states  $x \in X$  and controls  $u \in U$  and update these values to get  $q^{(k+1)}$  after  
 24 *each step* of the robot.

25 From the results on Bellman iteration, we know that any Q-function that  
 26 satisfies the above equation is the optimal Q-function; we would therefore like  
 27 our Q-function to satisfy

$$q^*(x_k, u_k) \approx r(x_k, u_k) + \gamma \max_{u'} q^*(x_{k+1}, u').$$

28 over samples  $(x_k, u_k, x_{k+1})$  collected as the robot explores the environment.

**Tabular Q-Learning** Let us imagine the robot travels for  $n$  trajectories each of  $T$  time-steps each. We can now solve for  $q^*$  by minimizing the objective

$$\min_q \frac{1}{n(T+1)} \sum_{i=1}^n \sum_{k=0}^T \|q(x_k^i, u_k^i) - r(x_k^i, u_k^i) - \gamma \max_{u' \in U} q(x_{k+1}^i, u')\|_2^2. \quad (9.1)$$

on the data collected by the robot. The variable of optimization here are all values  $q^*(x, u)$  for  $x \in X$  and  $u \in U$ .

29 Notice a few important things about the above optimization problem. First,  
 30 the last term is a maximization over  $u' \in U$ , it is  $\max_{u' \in U} q(x_{k+1}^i, u')$  and  
 31 not  $q(x_{k+1}^i, u_{k+1}^i)$ . In practice, you should imagine a robot performing Q-  
 32 Learning in a grid-world setting where it seeks to find the optimal trajectory  
 33 to go from a source location to a target location. If at each step, the robot  
 34 has 4 controls to choose from, computing this last term involves taking the  
 35 maximum of 4 different values (4 columns in the tabular Q-function).

36 Notice that for finite-horizon dynamic programming we initialized the  
 37 Q-function at the terminal time to a known value (the terminal cost). Similarly,  
 38 for infinite-horizon value iteration, we discussed how we can converge to the  
 39 optimal Q-function with any initialization. In the above case, we do not impose  
 40 any such constraint upon the Q-function, but there is an implicit constraint. All  
 41 values  $q(x, u)$  have to be consistent with each other and ideally, the residual

$$\|q(x_k^i, u_k^i) - r(x_k^i, u_k^i) - \gamma \max_{u' \in U} q(x_{k+1}^i, u')\|_2^2 = 0$$

42 for all trajectories  $i$  and all timesteps  $T$ .

43 **Solving tabular Q-Learning** How should we solve the optimization problem  
 44 in (9.1)? This is easy, every entry  $q(x, u)$  for  $x \in U$  and  $u \in U$  is a  
 45 variable of this objective and each  $\|\cdot\|_2^2$  term in the objective simply represents  
 46 a constraint that ties these different values of the Q-function together. We can  
 47 solve for all  $q(x, u)$  iteratively as

$$\begin{aligned} q(x, u) &\leftarrow q(x, u) - \eta \nabla_{q(x, u)} \ell(q) \\ &= (1 - \eta) q(x, u) - \eta \left( r(x, u) + \gamma \max_{u'} q(x', u') \right) \end{aligned} \quad (9.2)$$

48 where  $\ell(q)$  is the entire objective  $\frac{1}{n(T+1)} \sum_i \sum_k \dots$  above and  $(x, u, x') \equiv$   
 49  $(x_k^i, u_k^i, x_{k+1}^i)$  in the second equation. An important point to note here is that  
 50 although the robot collects a finite number of data

$$D = \{(x_k^i, u_k^i)_{k=0,1,\dots,T}\}_{i=1}^n$$

51 we have an estimate for the value  $q(x, u)$  at all states  $x \in X$ . Intuitively,  
 52 tabular Q-learning looks at the returns obtained by the robot after starting from  
 53 a state  $x$  (the reward-to-come  $J(x)$ ) and patches the returns from nearby states  
 54  $x, x'$  using the constraints in the objective (9.1).

55 **Terminal state** One must be very careful about the terminal state in such  
 56 implementations of Q-learning. Typically, most research papers imagine that  
 57 they are solving an infinite horizon problem but use simulators that have an  
 58 explicit terminal state, i.e., the simulator does not proceed to the next timestep  
 59 after the robot reaches the goal. A workaround for using such simulators (this  
 60 applies for essentially all simulators) is to modify (9.2) as

$$q(x, u) = (1 - \eta) q(x, u) - \eta \left( r(x, u) + \gamma (1 - \mathbf{1}_{\{x' \text{ is terminal}\}}) \max_{u'} q(x', u') \right).$$

61 Effectively, we are setting  $q(x', u) = 0$  for all  $u \in U$  if  $x'$  is a terminal state  
 62 of problem. This is a very important point to remember and Q-Learning will  
 63 never work if you forget to include the term  $\mathbf{1}_{\{x' \text{ is terminal}\}}$  in your expression.

64 **What is the controller in tabular Q-Learning?** The controller in tabular Q-  
 65 Learning is easy to get after we solve (9.1). At test time, we use a deterministic

66 controller given by

$$u^*(x) = \underset{u'}{\operatorname{argmax}} q^*(x, u').$$

### 67 9.1.1 How to perform exploration in Q-Learning

68 The exploratory controller used by the robot  $u_e(\cdot | x)$  is critical to perform  
 69 Q-Learning well. If the exploratory controller does not explore much, we do  
 70 not get states from all parts of the state-space. This is quite bad, because in this  
 71 case the estimates of Q-function at *all states* will be bad, not just at the states  
 72 that the robot did not visit. To make this intuitive, imagine if we cordoned off  
 73 some nodes in the graph for the backward version of Dijkstra’s algorithm and  
 74 never used them to update the dist variable. We would never get to the optimal  
 75 cost-to-go for *all* states in this case because there could be trajectories that  
 76 go through these cordoned off states that lead to a smaller cost-to-go. So it is  
 77 quite important to pick the right exploratory controller.

78 It turns out that a random exploratory controller, e.g., a controller  $u_e(\cdot | x)$   
 79 that picks controls uniformly randomly is pretty good. We can show that our  
 80 tabular Q-Learning will converge to the optimal Q-function  $q^*(x, u)$  as the  
 81 amount of data drawn from the random controller goes to infinity, even if we  
 82 initialize the table to arbitrary values. In other words, if we are guaranteed that  
 83 the robot visits each state in the finite MDP infinitely often, it is a classical  
 84 result that updates of the form (9.2) for minimizing the objective in (9.1)  
 85 converge to the optimal Q-function.

86 **Epsilon-greedy exploration** Instead of the robot using a arbitrary controller  
 87  $u_e(\cdot | x)$  to gather data, we can use the current estimate of the Q-function  
 88 with some added randomness to ensure that the robot visits all states in the  
 89 state-space. This is a key idea in Q-Learning and is known as “epsilon-greedy”  
 90 exploration. We set

$$u_e(u | x) = \begin{cases} \operatorname{argmax}_u q(x, u) & \text{with probability } 1 - \epsilon \\ \operatorname{uniform}(U) & \text{with probability } \epsilon. \end{cases} \quad (9.3)$$

91 for some user-chosen value of  $\epsilon$ . Effectively, the robot repeats the controls it  
 92 took in the past with probability  $1 - \epsilon$  and uniformly samples from the entire  
 93 control space with probability  $\epsilon$ . The former ensures that the robot moves  
 94 towards the parts of the state-space where states have a high return-to-come  
 95 (after all, that is the what the Q-function  $q(x, u)$  indicates). The latter ensures  
 96 that even if the robot’s estimate of the Q-function is bad, it is still visiting  
 97 every state in the state-space infinitely often.

98 **A different perspective on Q-Learning** Conceptually, we can think of  
 99 tabular Q-learning as happening in two stages. In the first stage, the robot  
 100 gathers a large amount of data

$$D = \{(x_k^i, u_k^i)_{k=0,1,\dots,T}\}_{i=1}^n$$

101 using the exploratory controller  $u_e(\cdot | x)$ ; let us consider the case when we are  
 102 using an arbitrary exploratory controller, not epsilon-greedy exploration. Using

❶ This is again the power of dynamic programming at work. The Bellman equation guarantees the convergence of value iteration provided we compute the expectation exactly. But if the robot does give us lots of data from the environment, then Q-Learning also inherits this property of convergence to the optimal Q-function from any initialization.

103 this data, the robot fits a model for the system, i.e., it learns the underlying  
104 MDP

$$P(x' | x, u);$$

105 this is very similar to the step in the Baum-Welch algorithm that we saw for  
106 learning the Markov state transition matrix of the HMM in Chapter 2. We  
107 simply take frequency counts to estimate this probability

$$P(x' | x, u) \approx \frac{1}{N} \sum_i \mathbf{1}_{\{x' \text{ was reached from } x \text{ using control } u\}}$$

108 where  $N$  is the number of the times the robot took control  $u$  at state  $x$ . Given  
109 this transition matrix, we can now perform value iteration on the MDP to learn  
110 the Q-function

$$q^{(k+1)}(x, u) = \mathbb{E}_{x' \sim P(\cdot | x, u)} \left[ r(x, u) + \gamma \max_{u'} q^{(k)}(x', u) \right].$$

111 The success of this two-stage approach depends upon how accurate our esti-  
112 mate of  $P(x' | x, u)$  is. This in turn depends on how much the robot explored  
113 the domain and the size of the dataset it collected, both of these need to be  
114 large. We can therefore think of Q-learning as interleaving these two stages  
115 in a single algorithm, it learns the dynamics of the system and the Q-function  
116 for that dynamics simultaneously. But the Q-Learning algorithm does not  
117 really maintain a representation of the dynamics, i.e., at the end of running  
118 Q-Learning, we do not know what  $P(x' | x, u)$  is.

## 119 9.2 Function approximation (Deep Q Networks)

120 Tabular methods are really nice but they do not scale to large problems. The  
121 grid-world in the homework problem on policy iteration had 100 states, a  
122 typical game of Tetris has about  $10^{60}$  states. For comparison, the number of  
123 atoms in the known universe is about  $10^{80}$ . The number of different states  
124 in a typical Atari game is more than  $10^{300}$ . These are all problems with a  
125 discrete number of states and controls, for continuous state/control-space, the  
126 number of distinct states/controls is infinite. So it is essentially impossible to  
127 run the tabular Q-Learning method from the previous section for most real-  
128 world problems. In this section, we will look at a powerful set of algorithms  
129 that parameterize the Q-function using a neural network to work around this  
130 problem.

131 We use the same idea from the previous chapter, that of parameterizing the  
132 Q-function using a deep network. We will denote

$$q_\varphi(x, u) : X \times U \mapsto \mathbb{R}$$

133 as the Q-function and our goal is to fit the deep network to obtain the weights  $\hat{\varphi}$ ,  
134 instead of maintaining a very large table of size  $|X| \times |U|$  for the Q-function.  
135 Fitting the Q-function is quite similar to the tabular case: given a dataset

136  $D = \{(x_t^i, u_t^i)_{t=0,1,\dots,T}\}_{i=1}^n$  from the system, we want to enforce

$$q_\varphi(x_t^i, u_t^i) = r(x_t^i, u_t^i) + \gamma \max_{u'} q_\varphi(x_{t+1}^i, u')$$

137 for all tuples  $(x_t^i, u_t^i, x_{t+1}^i)$  in the dataset. Just like the previous section, we  
138 will solve

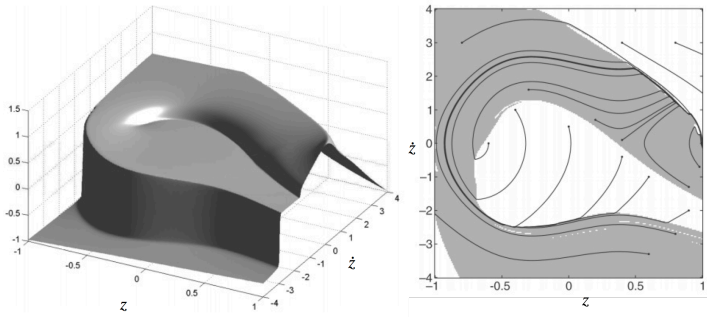
$$\begin{aligned} \hat{\varphi} &= \operatorname{argmin}_{\varphi} \frac{1}{n(T+1)} \sum_{i=1}^n \sum_{t=1}^T \left( \underbrace{q_\varphi(x_t^i, u_t^i) - r(x_t^i, u_t^i) - \gamma \left(1 - \mathbf{1}_{\{x_{t+1}^i \text{ is terminal}\}}\right) \max_{u'} q_\varphi(x_{t+1}^i, u')}_{\operatorname{target}(x'; \varphi)} \right)^2 \\ &\equiv \operatorname{argmin}_{\varphi} \frac{1}{n(T+1)} \sum_{i=1}^n \sum_{t=1}^T (q_\varphi(x_t^i, u_t^i) - \operatorname{target}(x_{t+1}^i; \varphi))^2 \end{aligned} \quad (9.4)$$

139 The last two terms in this expression above are together called the “target”  
140 because the problem is very similar to least squares regression, except that the  
141 targets also depend on the weights  $\varphi$ . This is what makes it challenging to  
142 solve.

143 As discussed above, Q-Learning with function approximation is known  
144 as “Fitted Q Iteration”. Remember that very important point that the robot  
145 collects data using the exploratory controller  $u_e(\cdot | x)$  but the Q-function that  
146 we fit is the *optimal* Q-function.

147 **Fitted Q-Iteration with function approximation may not converge to**  
148 **the optimal Q-function** It turns out that (9.4) has certain mathematical  
149 intricacies that prevent it from converging to the optimal Q-function. We  
150 will first give the intuitive reason. In the tabular Q-Learning setting, if we  
151 modify some entry  $q(x, u)$  for an  $x \in X$  and  $u \in U$ , the other entries (which  
152 are tied together using the Bellman equation) are all modified. This is akin  
153 to you changing the dist value of one node in Dijkstra’s algorithm; the dist  
154 values of *all* other nodes will have to change to satisfy the Bellman equation.  
155 This is what (9.2) achieves if implemented with a decaying step-size  $\eta$ ; see  
156 <http://users.isr.ist.utl.pt/~mtjspaam/readingGroup/ProofQlearning.pdf> for the  
157 proof. This does not hold for (9.4). Even if the objective in (9.4) is zero  
158 on our collected dataset, i.e., the Q-function fits data collected by the robot  
159 perfectly, the Q-function may not be the *optimal* Q-function. An intuitive  
160 way of understanding this problem is that even if the Bellman error is zero on  
161 samples in the dataset, the optimization objective says nothing about states  
162 that are not present in the dataset; the Bellman error on them is completely  
163 dependent upon the smoothness properties of the function expressed by the  
164 neural architecture. Contrast this comment with the solution of the HJB  
165 equation in Chapter 6 where the value function was quite non-smooth at some  
166 places. If our sampled dataset does not contain those places, there is no way  
167 the neural network can know the optimal form of the value function.

❗ The mathematical reason behind this is that the Bellman operator, i.e., the update to the Q/value-function is a contraction for the tabular setting, this is not the case for Fitted Q-Iteration unless the function approximation has some technical conditions imposed upon it.



168

## 169 9.2.1 Embellishments to Q-Learning

170 We next discuss a few practical aspects of implementing Q-Learning. Each of  
 171 the following points is extremely important to understand how to get RL to  
 172 work on real-world problems, so you should internalize these.

173 **Pick mini-batches from different trajectories in SGD** . In practice, we fit  
 174 the Q-function using stochastic gradient descent. At each iteration we sample  
 175 a mini-batch of inputs  $(x_t^i, u_t^i, x_{t+1}^i)$  from different trajectories  $i \in \{1, \dots, n\}$   
 176 and update the weights  $\varphi$  in the direction of the negative gradient.

$$\varphi^{k+1} = \varphi^k - \eta \nabla_{\varphi} (q_{p^k}(x, u) - \text{target}(x'; \varphi^k))^2.$$

177 The mini-batch is picked to have samples from different trajectories because  
 178 samples from the same trajectory are correlated to each other (after all, the  
 179 robot obtains the next tuple  $(x', u', x'')$  from the previous tuple  $(x, u, x')$ ).

180 **Replay buffer** The dataset  $D$  is known as the replay buffer.

181 **Off-policy learning** The replay buffer is typically not fixed during training.  
 182 Instead of drawing data from the exploratory controller  $u_e$ , we can think of the  
 183 following algorithm. Initialize the Q-function weights to  $\varphi^0$  and the dataset to  
 184  $D = \emptyset$ . At the  $k^{\text{th}}$  iteration,

- 185 • Draw a dataset  $D^k$  of  $n$  trajectories from the  $\epsilon$ -greedy policy

$$u_{\epsilon}(u | x) = \begin{cases} \operatorname{argmax}_u q^k(x, u) & \text{with probability } 1 - \epsilon \\ \operatorname{uniform}(U) & \text{with probability } \epsilon. \end{cases}$$

- 186 • Add new trajectories to the dataset

$$D \leftarrow D \cup D^k.$$

- 187 • Update weights to  $q^{k+1}$  using all past data  $D$  using (9.4).

188 Compare this algorithm to policy-gradient-based methods which throw away  
 189 the data from the previous iteration. Indeed, when we want to compute  
 190 the gradient  $\nabla_{\theta} \mathbb{E}_{\tau \sim p_{\theta^k}} [R(\tau)]$ , we should sample trajectories from current  
 191 weights  $\theta^k$ , we cannot use trajectories from some old weights. In contrast,

192 in Q-Learning, we maintain a cumulative dataset  $D$  that contains trajectories  
 193 from all the past  $\epsilon$ -greedy controllers and use it to find new weights of the Q-  
 194 function. We can do so because of the powerful Bellman equation, Q-Iteration  
 195 is learning the *optimal* value function and no matter what dataset (9.4) is  
 196 evaluated upon, if the error is zero, we are guaranteed that Q-function learned  
 197 is the optimal one. Policy gradients do not use the Bellman equation and that  
 198 is why they are so inefficient. This is also the reason Q-Learning with a replay  
 199 buffer is called “off-policy” learning because it learns the optimal controller  
 200 even if the data that it uses comes from some other non-optimal controller (the  
 201 exploratory controller or the  $\epsilon$ -greedy controller).

202 Using off-policy learning is an old idea, the DQN paper which demon-  
 203 strated very impressive results on Atari games using RL brought it back into  
 204 prominence.

205 **Setting a good value of  $\epsilon$  for exploration is critical** Towards the beginning  
 206 of training, we want a large value for  $\epsilon$  to gather diverse data from the envi-  
 207 ronment. As training progresses, we want to reduce  $\epsilon$  because presumably we  
 208 have a few good control trajectories that result in good returns and can focus  
 209 on searching the neighborhood of these trajectories.

210 **Prioritized experience replay** is an idea where instead of sampling from  
 211 the replay buffer  $D$  uniformly randomly when we fit the Q-function in (9.4),  
 212 we only sample data points  $(x_t^i, u_t^i)$  which have a high Bellman error

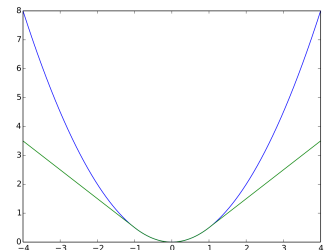
$$|q_\varphi(x_t^i, u_t^i) - r(x_t^i, u_t^i) - \gamma \left(1 - \mathbf{1}_{\{x_{t+1}^i \text{ is terminal}\}}\right) \max_{u'} q_\varphi(x_{t+1}^i, u')|$$

213 This is a reasonable idea but is not very useful in practice for two reasons.  
 214 First, if we use deep networks for parameterizing the Q-function, the network  
 215 *can* fit even very complex datasets so there is no reason to not use the data  
 216 points with low Bellman error in (9.4); the gradient using them will be small  
 217 anyway. Second, there are a lot of hyper-parameters that determine prioritized  
 218 sampling, e.g., the threshold beyond which we consider the Bellman error to be  
 219 high. These hyper-parameters are quite difficult to use in practice and therefore  
 220 it is a good idea to not use prioritized experience replay at the beginning of  
 221 development of your method on a new problem.

222 **Using robust regression to fit the Q-function** There may be states in the  
 223 replay buffer with very high Bellman error, e.g., the kinks in the value function  
 224 for the mountain car obtained from HJB above, if we happen to sample  
 225 those. For instance, these are states where the controller “switches” and is  
 226 discontinuous function of state  $x$ . In these cases, instead of these few states  
 227 dominating the gradient for the entire dataset, we can use ideas in robust  
 228 regression to reduce their effect on the gradient. A popular way to do so is to  
 229 use a Huber-loss in place of the quadratic loss in (9.4)

$$\text{huber}_\delta(a) = \begin{cases} \frac{a^2}{2} & \text{for } |a| \leq \delta \\ \delta \left(|a| - \frac{\delta}{2}\right) & \text{otherwise.} \end{cases} \quad (9.5)$$

❶ Huber loss for  $\delta = 1$  (green) compared to the squared error loss (blue).





230 **Delayed target** Notice that the target also depends upon the weights  $\varphi$ :

$$\text{target}(x'; \varphi) := r(x, u) + \gamma (1 - \mathbf{1}_{\{x' \text{ is terminal}\}}) \max_{u'} q_{\varphi}(x', u').$$

231 This creates a very big problem when we fit the Q-function. Effectively, both  
 232 the covariate and the target in (9.4) depend upon the weights of the Q-function.  
 233 Minimizing the objective in (9.4) is akin to performing least squares regression  
 234 where the targets keep changing every time you solve for the solution. This  
 235 is the root cause of why Q-Learning is difficult to use in practice. A popular  
 236 hack to get around this problem is to use some old weights to compute the  
 237 target, i.e., use the loss

$$\frac{1}{n(T+1)} \sum_{i,t} \left( q_{\varphi^k}(x_t^i, u_t^i) - \text{target}(x_{t+1}^i; \varphi^{k'}) \right)^2. \quad (9.6)$$

238 in place of (9.4). Here  $k'$  is an iterate much older than  $k$ , say  $k' = k - 100$ .  
 239 This trick is called “delayed target”.

240 **Exponential averaging to update the target** Notice that in order to im-  
 241 plement delayed targets as discussed above we will have to save all weights  
 242  $\varphi^k, \varphi^{k-1}, \dots, \varphi^{k-100}$ , which can be cumbersome. We can however do yet  
 243 another clever hack and initialize two copies of the weights, one for the actual  
 244 Q-function  $\varphi^k$  and another for the target, let us call it  $\varphi'^k$ . We set the target  
 245 equal to the Q-function at initialization. The target copy is updated at each  
 246 iteration to be

$$\varphi'^{k+1} = (1 - \alpha)\varphi'^k + \alpha\varphi^{k+1} \quad (9.7)$$

247 with some small value, say  $\alpha = 0.05$ . The target’s weights are therefore an  
 248 exponentially averaged version of the weights of the Q-function.

249 **Why are delayed targets essential for Q-Learning to work?** There are  
 250 many explanations given why delayed targets are essential in practice but the  
 251 correct one is not really known yet.

- 252 1. For example, one reason could be that since  $q_{\varphi^k}(x, u)$  for a given state  
 253 typically increases as we train for more iterations in Q-Learning, the old  
 254 weights inside a delayed target are an underestimate of the true target.  
 255 This might lead to some stability in situations when the Q-function’s  
 256 weights  $\varphi^k$  change too quickly when we fit (9.4) or we do not have  
 257 enough data in the replay buffer yet.
- 258 2. Another reason one could hypothesize is related to concepts like self-  
 259 distillation. For example, we may write a new objective for Q-Learning  
 260 that looks like

$$\left( q_{\varphi^k}(x_t^i, u_t^i) - \text{target}(x_{t+1}^i; \varphi^k) \right)^2 + \frac{1}{2\lambda} \|\varphi^k - \varphi^{k'}\|_2^2$$

261 where the second term is known as proximal term that prevents the  
 262 weights  $\varphi^k$  from change too much from their old values  $\varphi^{k'}$ . Proximal  
 263 objectives are more stable versions of the standard quadratic objective

264 in (9.4) and help in cases when one is solving Q-Learning using SGD  
265 updates.

266 **Double Q-Learning** Even a delayed target may not be sufficient to get  
267 Q-Learning to lead to good returns in practice. Focus on one state  $x$ . One  
268 problem arise from the max operator in (9.4). Suppose that the Q-function  
269  $q_{\varphi^k}$  corresponds to a particularly bad controller, say a controller that picks a  
270 control

$$\operatorname{argmax}_u q_{\varphi^k}(x, u)$$

271 that is very different from the optimal control

$$\operatorname{argmax}_u q^*(x, u)$$

272 then, even the delayed target  $q_{\varphi^{k'}}$  may be a similarly poor controller. The ideal  
273 target is of course the return-to-come, or the value of the optimal Q-function  
274  $\max_{u'} q^*(x', u')$ , but we do not know it while fitting the Q-function. The same  
275 problem also occurs if our Q-function (or its delayed version, the target) is too  
276 optimistic about the values of certain control inputs, it will consistently pick  
277 those controls in the max operator. One hack to get around this problem is to  
278 pick the maximizing control input using the non-delayed Q-function but use  
279 the value of the delayed target

$$\operatorname{target}_{\text{DDQN}}(x_{t+1}^i; \varphi'^k) = r(x, u) + \gamma \left(1 - \mathbf{1}_{\{x_{t+1}^i \text{ is terminal}\}}\right) q_{\varphi'^k}(x_{t+1}^i, u'). \quad (9.8)$$

280 where

$$u' = \underbrace{\operatorname{argmax}_u q_{\varphi^k}(x_{t+1}^i, u)}_{\text{control chosen by the Q-function}} .$$

281 **Training two Q-functions** We can also train two copies of the Q-function  
282 simultaneously, each with its own delayed target and mix-and-match their  
283 targets. Let  $\varphi^{(1)k}$  and  $\varphi'^{(1)k}$  be one Q-function and target pair and  $\varphi^{(2)k}$  and  
284  $\varphi'^{(2)k}$  be another pair. We update both of them using the following objective.

$$\begin{aligned} \text{For } \varphi^{(1)} : & \left( q^{(1)k}(x, u) - r(x, u) - \gamma \left(1 - \mathbf{1}_{\{x' \text{ is terminal}\}}\right) \operatorname{target}_{\text{DDQN}}(x', \varphi'^{(2)k}) \right)^2 \\ \text{For } \varphi^{(2)} : & \left( q^{(2)k}(x, u) - r(x, u) - \gamma \left(1 - \mathbf{1}_{\{x' \text{ is terminal}\}}\right) \operatorname{target}_{\text{DDQN}}(x', \varphi'^{(1)k}) \right)^2 \end{aligned} \quad (9.9)$$

285 Sometimes we also use only one target that is the minimum of the two targets  
286 (this helps because it is more pessimistic estimate of the true target)

$$\operatorname{target}(x') := \min \left\{ \operatorname{target}_{\text{DDQN}}(x', \varphi'^{(1)k}), \operatorname{target}_{\text{DDQN}}(x', \varphi'^{(2)k}) \right\} .$$

287 You will also see many papers train multiple Q-functions, many more than 2.  
288 In such cases, it is a good idea to pick the control for evaluation using all the

289 Q-functions:

$$u^*(x) := \operatorname{argmax}_u \sum_k q_{\varphi^{(k)}}(x, u).$$

290 rather than only one of them, as is often done in research papers.

291 **A remark on the various tricks used to compute the target** It may seem  
 292 that a lot of these tricks are about being pessimistic while computing the target.  
 293 This is our current understanding in RL and it is born out of the following ob-  
 294 servation: typically in practice, you will observe that the Q-function estimates  
 295 can become very large. Even if the TD error is small, the values  $q_{\varphi}(x, u)$  can  
 296 be arbitrarily large; see Figure 1 in **Continuous Doubly Constrained Batch**  
 297 **Reinforcement Learning** for an example in a slightly different setting. This  
 298 occurs because we pick the control that maximizes the Q-value of a particular  
 299 state  $x$  in (9.8). Effectively, if the Q-value  $q_{\varphi}(x', u)$  of a particular control  
 300  $u \in U$  is an over-estimate, the target will keep selecting this control as the  
 301 maximizing control, which drives up the value of the Q-function at  $q_{\varphi}(x, u)$   
 302 as well. This problem is a bit more drastic in the next section on continuous-  
 303 valued controls. It is however unclear how to best address this issue and design  
 304 mathematically sound methods that do not use arbitrary heuristics such as  
 305 “pessimism”.

### 306 9.3 Q-Learning for continuous control spaces

307 All the methods we have looked at in this chapter are for discrete control  
 308 spaces, i.e., the set of controls that the robot can take is a finite set. In this case  
 309 we can easily compute the maximizing control of the Q-function.

$$u^*(x) = \operatorname{argmax}_u q_{\varphi}(x, u).$$

310 Certainly a lot of real-world problems have continuous-valued controls and  
 311 we therefore need Q-Learning-based methods to handle this.

312 **Deterministic policy gradient** A natural way, although non-rigorous, to  
 313 think about this is to assume that we are given a Q-function  $q_{\varphi}(x, u)$  (we  
 314 will leave the controller for which this is the Q-function vague for now)  
 315 and a dataset  $D = \{(x_t^i, u_t^i)_{t=0}^T\}_{i=1}^n$ . We can find a deterministic feedback  
 316 controller that takes controls that lead to good values as

$$\theta^* = \max_{\theta} \frac{1}{n(T+1)} \sum_{i=1}^n \sum_{t=0}^T q_{\varphi}(x_t^i, u_{\theta}(x_t^i)). \quad (9.10)$$

317 Effectively we are fitting a feedback controller that takes controls  $u_{\theta^*}(x)$  that  
 318 are the maximizers of the Q-function. This is a natural analogue of the argmax  
 319 over controls for discrete/finite control spaces. Again we should think of  
 320 having a deep network that parametrizes the deterministic controller and fitting

❶ Mathematically, the fundamental problem in function-approximation-based RL is actually clear: even if the Bellman operation is a contraction for tabular RL, it need not be a contraction when we are approximating the Q-function using a neural network. Therefore minimizing TD-error which works quite well for the tabular case need not work well in the function-approximation case. There may exist other, more robust, ways of computing the Bellman fixed point  $q_{\varphi}(x, u) = r(x, u) + \max_{u'} \gamma q_{\varphi}(x', u')$  other than minimizing the the squared TD error but we do not have good candidates yet.

321 its parameters  $\theta$  using stochastic gradient descent on (9.10)

$$\begin{aligned}\theta^{k+1} &= \theta^k + \eta \nabla_{\theta} q_{\varphi}(x^{\omega}, u_{\theta^k}(x^{\omega})) \\ &= \theta^k + \eta (\nabla_u q_{\varphi}(x^{\omega}, u)) (\nabla_{\theta} u_{\theta^k}(x^{\omega}))\end{aligned}\quad (9.11)$$

322 where  $\omega$  is the index of the datum in the dataset  $D$ . The equality was obtained  
323 by applying the chain rule. This result is called the “deterministic policy  
324 gradient” and we should think of it as the limit of the policy gradient for a  
325 stochastic controller as the stochasticity goes to zero. Also notice that the term

$$\nabla_u q_{\varphi}(x^{\omega}, u)$$

326 is the gradient of the output of the Q-function  $q_{\varphi} : X \times U \mapsto \mathbb{R}$  with respect  
327 to its second input  $u$ . Such gradients can also be easily computed using  
328 backpropagation in PyTorch. It is different than the gradient of the output with  
329 respect to its weights

$$\nabla_{\varphi} q_{\varphi}(x^{\omega}, u).$$

330 **On-policy deterministic actor-critic** Let us now construct an analogue of  
331 the policy gradient method for the case of a deterministic controller. The  
332 algorithm would proceed as follows. We initialize weights of a Q-function  $\varphi^0$   
333 and weights of the deterministic controller  $\theta^0$ .

- 334 1. At the  $k^{\text{th}}$  iteration, we collect a dataset from the robot using the  
335 latest controller  $u_{\theta^k}$ . Let this dataset be  $D^k$  that consists of tuples  
336  $(x, u, x', u')$ .
- 337 2. Fit a Q-function  $q^{\theta^k}$  to this dataset by minimizing the temporal differ-  
338 ence error

$$\varphi^{k+1} = \underset{\varphi}{\operatorname{argmin}} \sum_{(x, u, x', u') \in D^k} (q_{\varphi}(x, u) - r(x, u) - \gamma (1 - \mathbf{1}_{\{x' \text{ is terminal}\}}) q_{\varphi'}(x', u'))^2.$$

(9.12)

339 Notice an important difference in the expression above, instead of using  
340  $\max_u$  in the target, we are using the control that the current controller,  
341 namely  $u_{\theta^k}$  has taken. This is because we want to evaluate the controller  
342  $u_{\theta^k}$  and simply parameterize the Q-function using weights  $\varphi^{k+1}$ . More  
343 precisely, we hope that we have

$$q_{\varphi^{k+1}}(x, u_{\theta^k}(x)) \approx \max_u q^{\theta^k}(x, u).$$

- 344 3. We can now update the controller using this Q-function:

$$\theta^{k+1} = \theta^k + \eta \nabla_{\theta} q_{\varphi^{k+1}}(x^{\omega}, u_{\theta^k}(x^{\omega})) \quad (9.13)$$

345 This algorithm is called “on-policy SARSA” because at each iteration we draw  
346 fresh data  $D^k$  from the environment; this is the direct analogue of actor-critic  
347 methods that we studied in the previous chapter for deterministic controllers.

348 **Off-policy deterministic actor-critic methods** We can also run the above  
349 algorithm using data from an exploratory controller. The only difference is

**i** SARSA is an old algorithm in RL that is the tabular version of what we did here. It stands for state-action-reward-state-action ...

350 that the we now do not throw away the data  $D^k$  from older iterations

$$D = D^1 \cup \dots \cup D^k$$

351 and therefore have to change (9.12) to be

$$\varphi^{k+1} = \underset{\varphi}{\operatorname{argmin}} \sum_{(x,u,x',u') \in D} \left( q_{\varphi}(x,u) - r(x,u) - \gamma (1 - \mathbf{1}_{\{x' \text{ is terminal}\}}) q_{\varphi'}(x', \underbrace{u_{\theta^k}(x')}_{\text{notice the difference}}) \right)^2. \quad (9.14)$$

352 Effectively, we are fitting the optimal Q-function using the data  $D$  but since  
 353 we can no longer take the maximum over controls directly, we plug in the  
 354 controller in the computation of the target. This is natural; we think of the  
 355 controller as the one that maximizes the Q-function when we update (9.13).  
 356 When used with deep networks, this is called the “deep deterministic policy  
 357 gradient” algorithm, it is popular by the name DDPG.