**Chapter 18 – Reinforcement Learning**

*This notebook contains all the sample code and solutions to the exercises in chapter 18.*

`CO` Open in Colab    `k` Open in Kaggle

## ˅ Setup

This project requires Python 3.7 or above:

```
import sys

assert sys.version_info >= (3, 7)
```

And TensorFlow ≥ 2.8:

```
from packaging import version
import tensorflow as tf

assert version.parse(tf.__version__) >= version.parse("2.8.0")
```

As we did in earlier chapters, let's define the default font sizes to make the figures prettier. We will also display some Matplotlib animations, and there are several possible options to do that: we will use the Javascript option.

```
import matplotlib.animation
import matplotlib.pyplot as plt

plt.rc('font', size=14)
plt.rc('axes', labelsize=14, titlesize=14)
plt.rc('legend', fontsize=14)
plt.rc('xtick', labelsize=10)
plt.rc('ytick', labelsize=10)
plt.rc('animation', html='jshtml')
```

And let's create the `images/rl` folder (if it doesn't already exist), and define the `save_fig()` function which is used through this notebook to save the figures in high-res for the book:

```
from pathlib import Path

IMAGES_PATH = Path() / "images" / "rl"
IMAGES_PATH.mkdir(parents=True, exist_ok=True)

def save_fig(fig_id, tight_layout=True, fig_extension="png", resolution=300):
    path = IMAGES_PATH / f"{fig_id}.{fig_extension}"
    if tight_layout:
        plt.tight_layout()
    plt.savefig(path, format=fig_extension, dpi=resolution)
```

This chapter can be very slow without a GPU, so let's make sure there's one, or else issue a warning:

```
if not tf.config.list_physical_devices('GPU'):
    print("No GPU was detected. Neural nets can be very slow without a GPU.")
    if "google.colab" in sys.modules:
        print("Go to Runtime > Change runtime and select a GPU hardware "
              "accelerator.")
    if "kaggle_secrets" in sys.modules:
        print("Go to Settings > Accelerator and select GPU.")

    No GPU was detected. Neural nets can be very slow without a GPU.
    Go to Runtime > Change runtime and select a GPU hardware accelerator.
```

Let's install the Gymnasium library, which provides many environments for Reinforcement Learning. We'll also install the extra libraries needed for classic control environments (including CartPole, which we will use shortly), as well as for Box2D and Atari environments, which are needed for the exercises.

**Important notes:**

- OpenAI has handed over the maintenance and developmnent of the Gym library to the Farama foundation (see the [announcement](#)), and the library was renamed to Gymnasium. It's a drop-in replacement for OpenAI Gym: you can just install `gymnasium` instead of `gym`, and `import gymnasium as gym`, and everything should work fine.
- by running the following cell, you accept the Atari ROM license.

```
if "google.colab" in sys.modules or "kaggle_secrets" in sys.modules:
    %pip install -q -U gymnasium swig
    %pip install -q -U gymnasium[classic_control,box2d,atari,accept-rom-license]
```

```
          ──────────────────────────────────────  953.9/953.9 kB 6.2 MB/s eta 0:00:00
          ──────────────────────────────────────  1.9/1.9 MB 33.1 MB/s eta 0:00:00
          ──────────────────────────────────────  374.4/374.4 kB 3.1 MB/s eta 0:00:00
       Preparing metadata (setup.py) ... done
          ──────────────────────────────────────  434.7/434.7 kB 20.1 MB/s eta 0:00:00
       Installing build dependencies ... done
       Getting requirements to build wheel ... done
       Preparing metadata (pyproject.toml) ... done
          ──────────────────────────────────────  1.7/1.7 MB 27.1 MB/s eta 0:00:00
       Building wheel for box2d-py (setup.py) ... done
       Building wheel for AutoROM.accept-rom-license (pyproject.toml) ... done
```

## ⌄ Introduction to ~~OpenAI gym~~ Gymnasium

In this notebook we will be using [gymnasium](#), a great toolkit for developing and comparing Reinforcement Learning algorithms. It provides many environments for your learning *agents* to interact with. Let's import Gym and make a new CartPole environment:

```
import gymnasium as gym

env = gym.make("CartPole-v1", render_mode="rgb_array")
```

The CartPole (version 1) is a very simple environment composed of a cart that can move left or right, and pole placed vertically on top of it. The agent must move the cart left or right to keep the pole upright.

**Tip**: `gym.envs.registry` is a dictionary containing all available environments:

```
# extra code – shows the first few environments
envs = gym.envs.registry
sorted(envs.keys())[:5] + ["..."]
```

```
       /usr/local/lib/python3.10/dist-packages/ipykernel/ipkernel.py:283: DeprecationWarning: `should_run_async` will not call `transform_cell`
         and should_run_async(code)
       ['ALE/Adventure-ram-v5',
        'ALE/Adventure-v5',
        'ALE/AirRaid-ram-v5',
        'ALE/AirRaid-v5',
        'ALE/Alien-ram-v5',
        '...']
```

◀ ▭▭▭▭▭▭▭▭▭▭▭▭▭                                                                     ▶

The registry values are environment specifications:

```
# extra code – shows the specification for the CartPole-v1 environment
envs["CartPole-v1"]
```

```
       EnvSpec(id='CartPole-v1', entry_point='gymnasium.envs.classic_control.cartpole:CartPoleEnv', reward_threshold=475.0,
       nondeterministic=False, max_episode_steps=500, order_enforce=True, autoreset=False, disable_env_checker=False,
       apply_api_compatibility=False, kwargs={}, namespace=None, name='CartPole', version=1, additional_wrappers=(),
       vector_entry_point='gymnasium.envs.classic_control.cartpole:CartPoleVectorEnv')
```

Let's initialize the environment by calling is `reset()` method. This returns an observation, as well as a dictionary that may contain extra information. Both are environment-specific.

```
obs, info = env.reset(seed=42)
obs
```

```
array([ 0.0273956 , -0.00611216,  0.03585979,  0.0197368 ], dtype=float32)
```

```
info
```

```
{}
```

In the case of the CartPole, each observation is a 1D NumPy array composed of 4 floats: they represent the cart's horizontal position, its velocity, the angle of the pole (0 = vertical), and the angular velocity.

An environment can be visualized by calling its `render()` method. If you set `render_mode` to `"rgb_array"` when creating the environment, then this will return a NumPy array.
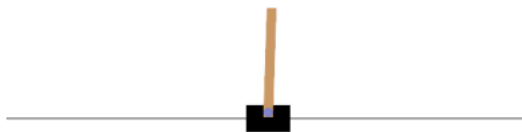
```
img = env.render()
img.shape  # height, width, channels (3 = Red, Green, Blue)
```

```
(400, 600, 3)
```

```
# extra code – creates a little function to render and plot an environment

def plot_environment(env, figsize=(5, 4)):
    plt.figure(figsize=figsize)
    img = env.render()
    plt.imshow(img)
    plt.axis("off")
    return img

plot_environment(env)
plt.show()
```



Let's see how to interact with an environment. Your agent will need to select an action from an "action space" (the set of possible actions). Let's see what this environment's action space looks like:

```
env.action_space
```

```
Discrete(2)
```

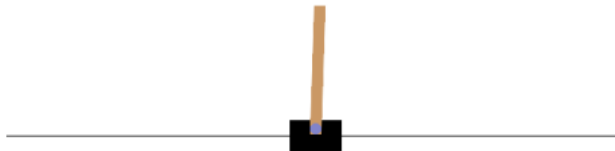Yep, just two possible actions: accelerate towards the left (0) or towards the right (1).

Since the pole is leaning toward the right ( `obs[2] > 0` ), let's accelerate the cart toward the right:

```
action = 1  # accelerate right
obs, reward, done, truncated, info = env.step(action)
obs
```

```
array([ 0.02727336,  0.18847767,  0.03625453, -0.26141977], dtype=float32)
```

Notice that the cart is now moving toward the right (`obs[1] > 0`). The pole is still tilted toward the right (`obs[2] > 0`), but its angular velocity is now negative (`obs[3] < 0`), so it will likely be tilted toward the left after the next step.

```
# extra code – displays the environment
plot_environment(env)
save_fig("cart_pole_plot")
plt.show()
```



Looks like it's doing what we're telling it to do!

The environment also tells the agent how much reward it got during the last step:

```
reward
```

```
    1.0
```

When the game is over, the environment returns `done=True`. In this case, it's not over yet:

```
done
```

```
    False
```

Some environment wrappers may want to interrupt the environment early. For example, when a time limit is reached or when an object goes out of bounds. In this case, `truncated` will be set to `True`. In this case, it's not truncated yet:

```
truncated
```

```
    False
```

Finally, `info` is an environment-specific dictionary that can provide some extra information that you may find useful for debugging or for training. For example, in some games it may indicate how many lives the agent has.

```
info
```

```
    {}
```

The sequence of steps between the moment the environment is reset until it is done or truncated is called an "episode". At the end of an episode (i.e., when `step()` returns `done=True` or `truncated=True`), you should reset the environment before you continue to use it.

```
if done or truncated:
    obs, info = env.reset()
```

Now how can we make the poll remain upright? We will need to define a *policy* for that. This is the strategy that the agent will use to select an action at each step. It can use all the past actions and observations to decide what to do.

## ⌄ A simple hard-coded policy

Let's hard code a simple strategy: if the pole is tilting to the left, then push the cart to the left, and *vice versa*. Let's see if that works:

```python
def basic_policy(obs):
    angle = obs[2]
    return 0 if angle < 0 else 1

totals = []
for episode in range(500):
    episode_rewards = 0
    obs, info = env.reset(seed=episode)
    for step in range(200):
        action = basic_policy(obs)
        obs, reward, done, truncated, info = env.step(action)
        episode_rewards += reward
        if done or truncated:
            break

    totals.append(episode_rewards)
```

```python
import numpy as np

np.mean(totals), np.std(totals), min(totals), max(totals)
```

```
    (41.698, 8.389445512070509, 24.0, 63.0)
```

Well, as expected, this strategy is a bit too basic: the best it did was to keep the poll up for only 63 steps. This environment is considered solved when the agent keeps the poll up for 200 steps.

Let's visualize one episode. You can learn more about Matplotlib animations in the [Matplotlib tutorial notebook](#).
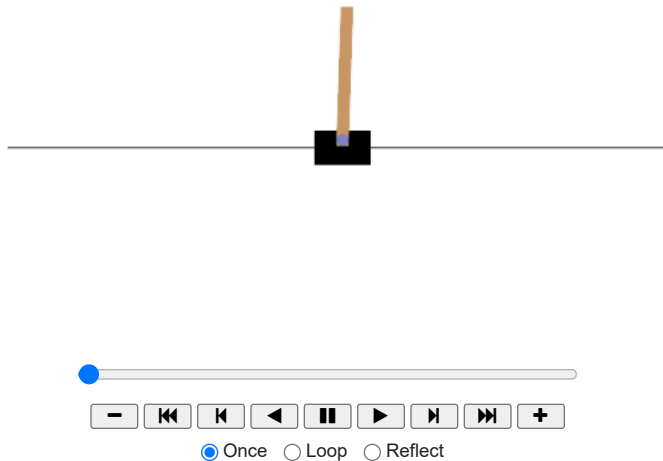
```python
# extra code – this cell displays an animation of one episode

def update_scene(num, frames, patch):
    patch.set_data(frames[num])
    return patch,

def plot_animation(frames, repeat=False, interval=40):
    fig = plt.figure()
    patch = plt.imshow(frames[0])
    plt.axis('off')
    anim = matplotlib.animation.FuncAnimation(
        fig, update_scene, fargs=(frames, patch),
        frames=len(frames), repeat=repeat, interval=interval)
    plt.close()
    return anim

def show_one_episode(policy, n_max_steps=200, seed=42):
    frames = []
    env = gym.make("CartPole-v1", render_mode="rgb_array")
    np.random.seed(seed)
    obs, info = env.reset(seed=seed)
    for step in range(n_max_steps):
        frames.append(env.render())
        action = policy(obs)
        obs, reward, done, truncated, info = env.step(action)
        if done or truncated:
            break
    env.close()
    return plot_animation(frames)

show_one_episode(basic_policy)
```

Clearly the system is unstable and after just a few wobbles, the pole ends up too tilted: game over. We will need to be smarter than that!

## ⌄ Neural Network Policies

Let's create a neural network that will take observations as inputs, and output the probabilities of actions to take for each observation. To choose an action, the network will estimate a probability for each action, then we will select an action randomly according to the estimated probabilities. In the case of the CartPole environment, there are just two possible actions (left or right), so we only need one output neuron: it will output the probability `p` of the action 0 (left), and of course the probability of action 1 (right) will be `1 - p`.

```
import tensorflow as tf

tf.random.set_seed(42)  # extra code – ensures reproducibility on the CPU

model = tf.keras.Sequential([
    tf.keras.layers.Dense(5, activation="relu"),
    tf.keras.layers.Dense(1, activation="sigmoid"),
])
```

In this particular environment, the past actions and observations can safely be ignored, since each observation contains the environment's full state. If there were some hidden state then you may need to consider past actions and observations in order to try to infer the hidden state of the environment. For example, if the environment only revealed the position of the cart but not its velocity, you would have to consider not only the current observation but also the previous observation in order to estimate the current velocity. Another example is if the observations are noisy: you may want to use the past few observations to estimate the most likely current state. Our problem is thus as simple as can be: the current observation is noise-free and contains the environment's full state.
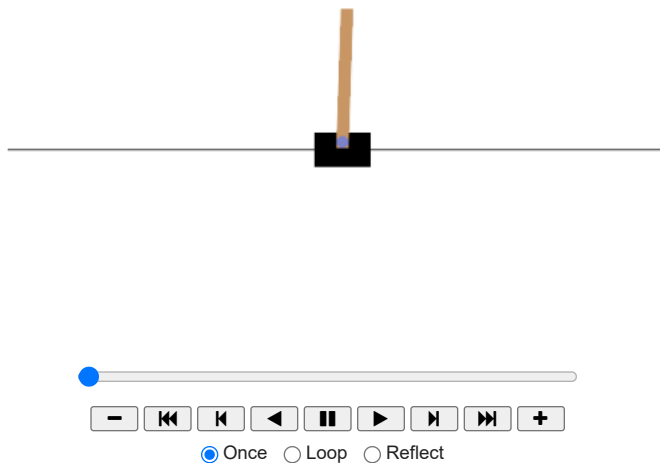
You may wonder why we plan to pick a random action based on the probability given by the policy network, rather than just picking the action with the highest probability. This approach lets the agent find the right balance between *exploring* new actions and *exploiting* the actions that are known to work well. Here's an analogy: suppose you go to a restaurant for the first time, and all the dishes look equally appealing so you randomly pick one. If it turns out to be good, you can increase the probability to order it next time, but you shouldn't increase that probability to 100%, or else you will never try out the other dishes, some of which may be even better than the one you tried.

Let's write a small policy function that will use the neural net to get the probability of moving left, then let's use it to run one episode:

```
# extra code – a function that creates an animation for a given policy model

def pg_policy(obs):
    left_proba = model.predict(obs[np.newaxis], verbose=0)[0][0]
    return int(np.random.rand() > left_proba)

np.random.seed(42)
show_one_episode(pg_policy)
```



Yeah... pretty bad. The neural network will have to learn to do better. First let's see if it is capable of learning the basic policy we used earlier: go left if the pole is tilting left, and go right if it is tilting right.

Let's see if it can learn a better policy on its own. One that does not wobble as much.

## ˅ Policy Gradients

To train this neural network we will need to define the target probabilities **y**. If an action is good we should increase its probability, and conversely if it is bad we should reduce it. But how do we know whether an action is good or bad? The problem is that most actions have delayed effects, so when you win or lose points in an episode, it is not clear which actions contributed to this result: was it just the last action? Or the last 10? Or just one action 50 steps earlier? This is called the *credit assignment problem*.

The *Policy Gradients* algorithm tackles this problem by first playing multiple episodes, then making the actions near positive rewards slightly more likely, while actions near negative rewards are made slightly less likely. First we play, then we go back and think about what we did.

Let's start by creating a function to play a single step using the model. We will also pretend for now that whatever action it takes is the right one, so we can compute the loss and its gradients. We will just save these gradients for now, and modify them later depending on how good or bad the action turned out to be.

```
def play_one_step(env, obs, model, loss_fn):
    with tf.GradientTape() as tape:
        left_proba = model(obs[np.newaxis])
        action = (tf.random.uniform([1, 1]) > left_proba)
        y_target = tf.constant([[1.]]) - tf.cast(action, tf.float32)
        loss = tf.reduce_mean(loss_fn(y_target, left_proba))

    grads = tape.gradient(loss, model.trainable_variables)
    obs, reward, done, truncated, info = env.step(int(action))
    return obs, reward, done, truncated, grads
```

If `left_proba` is high, then `action` will most likely be `False` (since a random number uniformly sampled between 0 and 1 will probably not be greater than `left_proba`). And `False` means 0 when you cast it to a number, so `y_target` would be equal to 1 - 0 = 1. In other words, we set the target to 1, meaning we pretend that the probability of going left should have been 100% (so we took the right action).

Now let's create another function that will rely on the `play_one_step()` function to play multiple episodes, returning all the rewards and gradients, for each episode and each step:

```
def play_multiple_episodes(env, n_episodes, n_max_steps, model, loss_fn):
    all_rewards = []
    all_grads = []
    for episode in range(n_episodes):
        current_rewards = []
        current_grads = []
        obs, info = env.reset()
        for step in range(n_max_steps):
            obs, reward, done, truncated, grads = play_one_step(
                env, obs, model, loss_fn)
            current_rewards.append(reward)
            current_grads.append(grads)
            if done or truncated:
                break

        all_rewards.append(current_rewards)
        all_grads.append(current_grads)

    return all_rewards, all_grads
```

The Policy Gradients algorithm uses the model to play the episode several times (e.g., 10 times), then it goes back and looks at all the rewards, discounts them and normalizes them. So let's create couple functions for that: the first will compute discounted rewards; the second will normalize the discounted rewards across many episodes.

```
def discount_rewards(rewards, discount_factor):
    discounted = np.array(rewards)
    for step in range(len(rewards) - 2, -1, -1):
        discounted[step] += discounted[step + 1] * discount_factor
    return discounted

def discount_and_normalize_rewards(all_rewards, discount_factor):
    all_discounted_rewards = [discount_rewards(rewards, discount_factor)
                              for rewards in all_rewards]
    flat_rewards = np.concatenate(all_discounted_rewards)
    reward_mean = flat_rewards.mean()
    reward_std = flat_rewards.std()
    return [(discounted_rewards - reward_mean) / reward_std
            for discounted_rewards in all_discounted_rewards]
```

Say there were 3 actions, and after each action there was a reward: first 10, then 0, then -50. If we use a discount factor of 80%, then the 3rd action will get -50 (full credit for the last reward), but the 2nd action will only get -40 (80% credit for the last reward), and the 1st action will get 80% of -40 (-32) plus full credit for the first reward (+10), which leads to a discounted reward of -22:

```
discount_rewards([10, 0, -50], discount_factor=0.8)
```

```
    array([-22, -40, -50])
```

To normalize all discounted rewards across all episodes, we compute the mean and standard deviation of all the discounted rewards, and we subtract the mean from each discounted reward, and divide by the standard deviation:

```
      discount_and_normalize_rewards([[10, 0, -50], [10, 20]],
                                     discount_factor=0.8)

      [array([-0.28435071, -0.86597718, -1.18910299]),
       array([1.26665318, 1.0727777 ])]
```

```python
n_iterations = 150
n_episodes_per_update = 10
n_max_steps = 200
discount_factor = 0.95
```

```python
# extra code – let's create the neural net and reset the environment, for
#              reproducibility

tf.random.set_seed(42)

model = tf.keras.Sequential([
    tf.keras.layers.Dense(5, activation="relu"),
    tf.keras.layers.Dense(1, activation="sigmoid"),
])

obs, info = env.reset(seed=42)
```

```python
optimizer = tf.keras.optimizers.Nadam(learning_rate=0.01)
loss_fn = tf.keras.losses.binary_crossentropy
```

```python
for iteration in range(n_iterations):
    all_rewards, all_grads = play_multiple_episodes(
        env, n_episodes_per_update, n_max_steps, model, loss_fn)

    # extra code – displays some debug info during training
    total_rewards = sum(map(sum, all_rewards))
    print(f"\rIteration: {iteration + 1}/{n_iterations},"
          f" mean rewards: {total_rewards / n_episodes_per_update:.1f}", end="")

    all_final_rewards = discount_and_normalize_rewards(all_rewards,
                                                       discount_factor)
    all_mean_grads = []
    for var_index in range(len(model.trainable_variables)):
        mean_grads = tf.reduce_mean(
            [final_reward * all_grads[episode_index][step][var_index]
             for episode_index, final_rewards in enumerate(all_final_rewards)
                 for step, final_reward in enumerate(final_rewards)], axis=0)
        all_mean_grads.append(mean_grads)

    optimizer.apply_gradients(zip(all_mean_grads, model.trainable_variables))
```

```
    /usr/local/lib/python3.10/dist-packages/tensorflow/python/framework/ops.py:291: DeprecationWarning: Conversion of an array with ndim > 0
      return int(self._numpy())
    Iteration: 150/150, mean rewards: 177.8
```
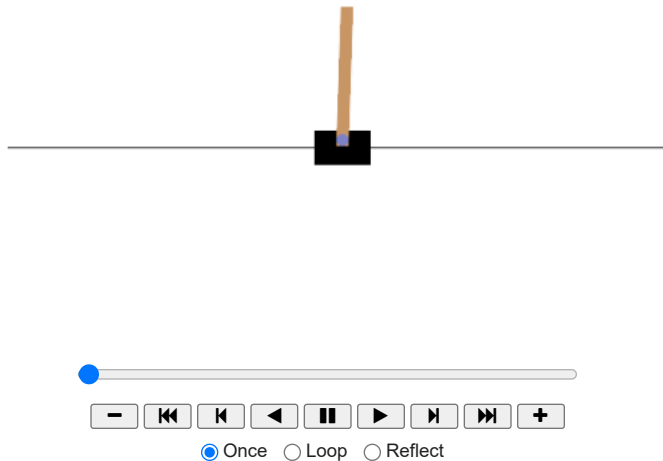
```python
np.mean(total_rewards), np.std(total_rewards), min(total_rewards), max(total_rewards)
```

```
    (135.0, 42.26345939461179, 73.0, 200.0)
```

```python
# extra code – displays the animation
np.random.seed(42)
show_one_episode(pg_policy)
```

## ⌄ Extra Material – Markov Chains

The following transition probabilities correspond to the Markov Chain represented in Figure 18–7. Let's run this stochastic process a few times to see what it looks like:

```python
np.random.seed(42)

transition_probabilities = [ # shape=[s, s']
        [0.7, 0.2, 0.0, 0.1],  # from s0 to s0, s1, s2, s3
        [0.0, 0.0, 0.9, 0.1],  # from s1 to s0, s1, s2, s3
        [0.0, 1.0, 0.0, 0.0],  # from s2 to s0, s1, s2, s3
        [0.0, 0.0, 0.0, 1.0]]  # from s3 to s0, s1, s2, s3

n_max_steps = 1000  # to avoid blocking in case of an infinite loop
terminal_states = [3]

def run_chain(start_state):
    current_state = start_state
    for step in range(n_max_steps):
        print(current_state, end=" ")
        if current_state in terminal_states:
            break
        current_state = np.random.choice(
            range(len(transition_probabilities)),
            p=transition_probabilities[current_state]
        )
    else:
        print("...", end="")

    print()

for idx in range(10):
    print(f"Run #{idx + 1}: ", end="")
    run_chain(start_state=0)
```

```
 Run #1: 0 0 3
 Run #2: 0 1 2 1 2 1 2 1 2 1 3
 Run #3: 0 1 2 1 2 1 2 1 2 1 2 1 2 1 2 1 2 1 2 1 2 1 3
 Run #4: 0 3
 Run #5: 0 1 2 1 2 1 2 1 2 1 2 1 2 1 3
 Run #6: 0 1 3
 Run #7: 0 1 2 1 2 1 2 1 2 1 2 1 2 1 2 1 2 1 2 1 2 1 2 1 2 1 2 1 2 1 2 1 2 1 2 1 2 1 2 1 2 1 2 1 2 1 2 1 2 1 2 1 3
```

```
Run #8: 0 0 0 1 2 1 2 1 3
Run #9: 0 1 2 1 2 1 2 1 2 1 2 1 2 1 2 1 2 1 2 1 3
Run #10: 0 0 0 1 2 1 3
```

## Markov Decision Process

Let's define some transition probabilities, rewards and possible actions. For example, in state s0, if action a0 is chosen then with proba 0.7 we will go to state s0 with reward +10, with probability 0.3 we will go to state s1 with no reward, and with never go to state s2 (so the transition probabilities are `[0.7, 0.3, 0.0]`, and the rewards are `[+10, 0, 0]`):

```python
transition_probabilities = [  # shape=[s, a, s']
    [[0.7, 0.3, 0.0], [1.0, 0.0, 0.0], [0.8, 0.2, 0.0]],
    [[0.0, 1.0, 0.0], None, [0.0, 0.0, 1.0]],
    [None, [0.8, 0.1, 0.1], None]
]
rewards = [  # shape=[s, a, s']
    [[+10, 0, 0], [0, 0, 0], [0, 0, 0]],
    [[0, 0, 0], [0, 0, 0], [0, 0, -50]],
    [[0, 0, 0], [+40, 0, 0], [0, 0, 0]]
]
possible_actions = [[0, 1, 2], [0, 2], [1]]
```

## Q-Value Iteration

```python
possible_actions = env.action_space
```

```python
from sklearn.preprocessing import KBinsDiscretizer
import time, math, random
from typing import Tuple
```

```python
n_bins = ( 6 , 12 )
lower_bounds = [ env.observation_space.low[2], -math.radians(50) ]
upper_bounds = [ env.observation_space.high[2], math.radians(50) ]

def discretizer( _ , __ , angle, pole_velocity ) -> Tuple[int,...]:
    """Convert continues state intro a discrete state"""
    est = KBinsDiscretizer(n_bins=n_bins, encode='ordinal', strategy='uniform')
    est.fit([lower_bounds, upper_bounds ])
    return tuple(map(int,est.transform([[angle, pole_velocity]])[0]))
```

```python
Q_values = np.full((n_bins[0], n_bins[1], possible_actions.n ), 0.0)

Q_values.shape
```

```
    (6, 12, 2)
```

```python
def nstep(action):
    obs, reward, done, truncated, info = env.step(action)
    # Do discretization
    next_state = discretizer(_, _, obs[2], obs[3])
    return next_state, reward, done, truncated
```

```python
def exploration_policy():
    return np.random.choice([0,1])
```

Haz doble clic (o pulsa Intro) para editar

```
gamma = 0.90  # the discount factor

history1 = []  # extra code – needed for the figure below
for iteration in range(50):
    Q_prev = Q_values.copy()
    history1.append(Q_prev) # extra code
    for s in range(3):
        for a in possible_actions[s]:
            Q_values[s, a] = np.sum([
                    transition_probabilities[s][a][sp]
                    * (rewards[s][a][sp] + gamma * Q_prev[sp].max())
                for sp in range(3)])

history1 = np.array(history1)  # extra code
```

```
Q_values = np.full((3, 3), -np.inf)  # -np.inf for impossible actions
for state, actions in enumerate(possible_actions):
    Q_values[state, actions] = 0.0  # for all possible actions

Q_values
```

```
    array([[18.91891892, 17.02702702, 13.62162162],
           [ 0.         ,        -inf, -4.87971488],
           [       -inf, 50.13365013,        -inf]])
```

```
Q_values.argmax(axis=1)  # optimal action for each state
```

```
    array([0, 0, 1])
```

The optimal policy for this MDP, when using a discount factor of 0.90, is to choose action a0 when in state s0, and choose action a0 when in state s1, and finally choose action a1 (the only possible action) when in state s2. If you try again with a discount factor of 0.95 instead of 0.90, you will find that the optimal action for state s1 becomes a2. This is because the discount factor is larger so the agent values the future more, and it is therefore ready to pay an immediate penalty in order to get more future rewards.

## ⌄ Q-Learning

Q-Learning works by watching an agent play (e.g., randomly) and gradually improving its estimates of the Q-Values. Once it has accurate Q-Value estimates (or close enough), then the optimal policy consists in choosing the action that has the highest Q-Value (i.e., the greedy policy).

We will need to simulate an agent moving around in the environment, so let's define a function to perform some action and get the new state and a reward:

```
def step(state, action):
    probas = transition_probabilities[state][action]
    next_state = np.random.choice([0, 1, 2], p=probas)
    reward = rewards[state][action][next_state]
    return next_state, reward
```

```
    /usr/local/lib/python3.10/dist-packages/ipykernel/ipkernel.py:283: DeprecationWarning: `should_run_async` will not call `transform_cell`
      and should_run_async(code)
```
◀                                                                                                                                          ▶

We also need an exploration policy, which can be any policy, as long as it visits every possible state many times. We will just use a random policy, since the state space is very small:

```
def exploration_policy(state):
    return np.random.choice(possible_actions[state])
```

Now let's initialize the Q-Values like earlier, and run the Q-Learning algorithm:

```
# extra code – initializes the Q-Values, just like earlier
np.random.seed(42)
Q_values = np.full((3, 3), -np.inf)
for state, actions in enumerate(possible_actions):
    Q_values[state][actions] = 0


alpha0 = 0.05  # initial learning rate
decay = 0.005  # learning rate decay
gamma = 0.90  # discount factor
state = 0  # initial state
history2 = []  # extra code – needed for the figure below

for iteration in range(10_000):
    history2.append(Q_values.copy())  # extra code
    action = exploration_policy(state)
    next_state, reward = step(state, action)
    next_value = Q_values[next_state].max()  # greedy policy at the next step
    alpha = alpha0 / (1 + iteration * decay)
    Q_values[state, action] *= 1 - alpha
    Q_values[state, action] += alpha * (reward + gamma * next_value)
    state = next_state

history2 = np.array(history2)  # extra code
```
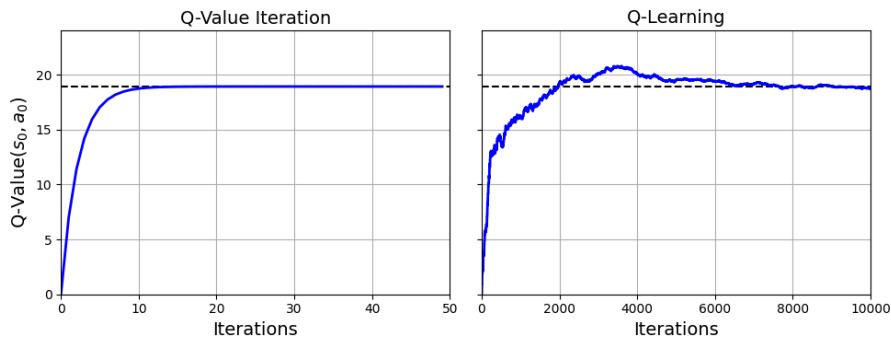
```
# extra code – this cell generates and saves Figure 18–9

true_Q_value = history1[-1, 0, 0]

fig, axes = plt.subplots(1, 2, figsize=(10, 4), sharey=True)
axes[0].set_ylabel("Q-Value$(s_0, a_0)$", fontsize=14)
axes[0].set_title("Q-Value Iteration", fontsize=14)
axes[1].set_title("Q-Learning", fontsize=14)
for ax, width, history in zip(axes, (50, 10000), (history1, history2)):
    ax.plot([0, width], [true_Q_value, true_Q_value], "k--")
    ax.plot(np.arange(width), history[:, 0, 0], "b-", linewidth=2)
    ax.set_xlabel("Iterations", fontsize=14)
    ax.axis([0, width, 0, 24])
    ax.grid(True)

save_fig("q_value_plot")
plt.show()
```

```
totals3 = []
for episode in range(episodes):
    # Print the current episode number
    if episode % 100 == 0:
        print(f"Episode {episode}/{episodes}")
    obs, info = env.reset(seed=episode)
    _, _, pole_angle, pole_velocity = obs
    state, done, truncated = discretizer(_,_,pole_angle, pole_velocity), False, False
    total_reward = 0
    done = False

    while not done:
      action = exploration_policy() # exploration policy
      next_state, reward, done, truncated = nstep(action)
      total_reward += reward
      next_value = Q_values[next_state].max()

      alpha = alpha0 / (1 + episode * decay)
      Q_values[state][action] *= 1 - alpha
      Q_values[state][action] += alpha * (reward + gamma * next_value)
      state = next_state

    totals3.append(total_reward)
    history2.append(Q_values.copy())
```

```
    Episode 0/10000
    Episode 100/10000
    Episode 200/10000
    Episode 300/10000
    Episode 400/10000
    Episode 500/10000
    Episode 600/10000
    Episode 700/10000
    Episode 800/10000
    Episode 900/10000
    Episode 1000/10000
    Episode 1100/10000
    Episode 1200/10000
    Episode 1300/10000
    Episode 1400/10000
    Episode 1500/10000
    Episode 1600/10000
    Episode 1700/10000
    Episode 1800/10000
    Episode 1900/10000
    Episode 2000/10000
    Episode 2100/10000
    Episode 2200/10000
    Episode 2300/10000
    Episode 2400/10000
    Episode 2500/10000
    Episode 2600/10000
    Episode 2700/10000
    Episode 2800/10000
    Episode 2900/10000
    Episode 3000/10000
    Episode 3100/10000
    Episode 3200/10000
    Episode 3300/10000
    Episode 3400/10000
    Episode 3500/10000
    Episode 3600/10000
    Episode 3700/10000
    Episode 3800/10000
    Episode 3900/10000
    Episode 4000/10000
    Episode 4100/10000
    Episode 4200/10000
    Episode 4300/10000
    Episode 4400/10000
    Episode 4500/10000
    Episode 4600/10000
    Episode 4700/10000
    Episode 4800/10000
    Episode 4900/10000
    Episode 5000/10000
    Episode 5100/10000
    Episode 5200/10000
    Episode 5300/10000
    Episode 5400/10000
```

```
        Episode 5500/10000
        Episode 5600/10000
        Episode 5700/10000
```

```python
np.mean(totals3), np.std(totals3), min(totals3), max(totals3)
```

```
    (22.291, 11.874818693352752, 8.0, 206.0)
```

```python
def exploration_rate(n : int, min_rate= 0.1 ) -> float :
    """Decaying exploration rate"""
    return max(min_rate, min(1, 1.0 - math.log10((n  + 1) / 25)))
```

```python
Q_values = np.full((n_bins[0], n_bins[1], possible_actions.n ), 0.0)
```

```python
totals2 = []
for episode in range(episodes):
    # Print the current episode number
    if episode % 100 == 0:
        print(f"Episode {episode}/{episodes}")
    # Reset
    obs, info = env.reset(seed=episode)
    _, _, pole_angle, pole_velocity = obs
    state, done, truncated = discretizer(_,_,pole_angle, pole_velocity), False, False
    total_reward = 0
    done = False

    while not done:
      # Exploitation policy
      action = np.argmax(Q_values[state])  # exploitation policy  greedly selects

      # Exploration policy
      if (np.random.random() < exploration_rate(episode)):
        action = exploration_policy() # exploration policy

       # Set next state and value
      next_state, reward, done, truncated = nstep(action)
      total_reward += reward
      next_value = Q_values[next_state].max()

      # Fill Q table with TD
      alpha = alpha0 / (1 + episode * decay)
      Q_values[state][action] *= 1 - alpha
      Q_values[state][action] += alpha * (reward + gamma * next_value)

      state = next_state

    totals2.append(total_reward)
    history2.append(Q_values.copy())
```

```
    Episode 0/10000
    Episode 100/10000
    Episode 200/10000
    Episode 300/10000
    Episode 400/10000
    Episode 500/10000
    Episode 600/10000
    Episode 700/10000
    Episode 800/10000
    Episode 900/10000
    Episode 1000/10000
    Episode 1100/10000
    Episode 1200/10000
    Episode 1300/10000
    Episode 1400/10000
    Episode 1500/10000
    Episode 1600/10000
    Episode 1700/10000
    Episode 1800/10000
    Episode 1900/10000
    Episode 2000/10000
    Episode 2100/10000
    Episode 2200/10000
    Episode 2300/10000
    Episode 2400/10000
    Episode 2500/10000
    Episode 2600/10000
    Episode 2700/10000
```

```
Episode 2800/10000
Episode 2900/10000
Episode 3000/10000
Episode 3100/10000
Episode 3200/10000
Episode 3300/10000
Episode 3400/10000
Episode 3500/10000
Episode 3600/10000
Episode 3700/10000
Episode 3800/10000
Episode 3900/10000
Episode 4000/10000
Episode 4100/10000
Episode 4200/10000
Episode 4300/10000
Episode 4400/10000
Episode 4500/10000
Episode 4600/10000
Episode 4700/10000
Episode 4800/10000
Episode 4900/10000
Episode 5000/10000
Episode 5100/10000
Episode 5200/10000
Episode 5300/10000
Episode 5400/10000
Episode 5500/10000
Episode 5600/10000
Episode 5700/10000
```

```
np.mean(totals2), np.std(totals2), min(totals2), max(totals2)
```

```
(83.2468, 75.87192425238733, 8.0, 747.0)
```

```
totals = []
for episode in range(episodes):
    # Print the current episode number
    if episode % 100 == 0:
        print(f"Episode {episode}/{episodes}")
    # Reset
    obs, info = env.reset(seed=episode)
    _, _, pole_angle, pole_velocity = obs
    state, done, truncated = discretizer(_,_,pole_angle, pole_velocity), False, False
    total_reward(0)

    while not done:
      # Exploitation policy
      action = np.argmax(Q_values[state])  # exploitation policy  greedly selects

      # Exploration policy
      if (np.random.random() < exploration_rate(episode)):
        action = exploration_policy() # exploration policy

      # Set next state and value
      next_state, reward, done, truncated = nstep(action)
      total_reward += reward

      # Fill Q table
      alpha = max(alpha0, min(1.0, 1.0 - math.log10((episode + 1) / 25)))
      Q_values[state][action] *= 1 - alpha
      Q_values[state][action] += alpha*(reward + gamma * np.max(Q_values[next_state]))
      state = next_state
    totals.append(total_reward)
    history2.append(Q_values.copy())
```

```
Episode 0/10000
Episode 100/10000
Episode 200/10000
Episode 300/10000
Episode 400/10000
Episode 500/10000
Episode 600/10000
Episode 700/10000
Episode 800/10000
Episode 900/10000
Episode 1000/10000
Episode 1100/10000
```

```
Episode 1200/10000
Episode 1300/10000
Episode 1400/10000
Episode 1500/10000
Episode 1600/10000
Episode 1700/10000
Episode 1800/10000
Episode 1900/10000
Episode 2000/10000
Episode 2100/10000
Episode 2200/10000
Episode 2300/10000
Episode 2400/10000
Episode 2500/10000
Episode 2600/10000
Episode 2700/10000
Episode 2800/10000
Episode 2900/10000
Episode 3000/10000
Episode 3100/10000
Episode 3200/10000
Episode 3300/10000
Episode 3400/10000
Episode 3500/10000
Episode 3600/10000
Episode 3700/10000
Episode 3800/10000
Episode 3900/10000
Episode 4000/10000
Episode 4100/10000
Episode 4200/10000
Episode 4300/10000
Episode 4400/10000
Episode 4500/10000
Episode 4600/10000
Episode 4700/10000
Episode 4800/10000
Episode 4900/10000
Episode 5000/10000
Episode 5100/10000
Episode 5200/10000
Episode 5300/10000
Episode 5400/10000
Episode 5500/10000
Episode 5600/10000
Episode 5700/10000
```

```python
np.mean(totals), np.std(totals), min(totals), max(totals)
```

```
(137017.28518629508, 73124.97786538936, 12.0, 262570.0)
```

```python
Q_values = np.full((n_bins[0], n_bins[1], possible_actions.n ), 0.0)
```

## Deep Q-Network

Let's build the DQN. Given a state, it will estimate, for each possible action, the sum of discounted future rewards it can expect after it plays that action (but before it sees its outcome):

```python
tf.random.set_seed(42)  # extra code – ensures reproducibility on the CPU

input_shape = [4]  # == env.observation_space.shape
n_outputs = 2  # == env.action_space.n

model = tf.keras.Sequential([
    tf.keras.layers.Dense(32, activation="elu", input_shape=input_shape),
    tf.keras.layers.Dense(32, activation="elu"),
    tf.keras.layers.Dense(n_outputs)
])
```

To select an action using this DQN, we just pick the action with the largest predicted Q-value. However, to ensure that the agent explores the environment, we choose a random action with probability `epsilon`.

```
def epsilon_greedy_policy(state, epsilon=0):
    if np.random.rand() < epsilon:
        return np.random.randint(n_outputs)  # random action
    else:
        Q_values = model.predict(state[np.newaxis], verbose=0)[0]
        return Q_values.argmax()  # optimal action according to the DQN
```

We will also need a replay buffer. It will contain the agent's experiences, in the form of tuples: `(obs, action, reward, next_obs, done)`. We can use the `deque` class for that:

```
from collections import deque

replay_buffer = deque(maxlen=2000)
```

**Note**: for very large replay buffers, you may want to use a circular buffer instead, as random access time will be O(1) instead of O(N). Or you can check out DeepMind's [Reverb library](#).

```
# extra code – A basic circular buffer implementation

class ReplayBuffer:
    def __init__(self, max_size):
        self.buffer = np.empty(max_size, dtype=object)
        self.max_size = max_size
        self.index = 0
        self.size = 0

    def append(self, obj):
        self.buffer[self.index] = obj
        self.size = min(self.size + 1, self.max_size)
        self.index = (self.index + 1) % self.max_size

    def sample(self, batch_size):
        indices = np.random.randint(self.size, size=batch_size)
        return self.buffer[indices]
```

And let's create a function to sample experiences from the replay buffer. It will return 6 NumPy arrays: `[obs, actions, rewards, next_obs, dones, truncateds]`.

```
def sample_experiences(batch_size):
    indices = np.random.randint(len(replay_buffer), size=batch_size)
    batch = [replay_buffer[index] for index in indices]
    return [
        np.array([experience[field_index] for experience in batch])
        for field_index in range(6)
    ]  # [states, actions, rewards, next_states, dones, truncateds]
```

Now we can create a function that will use the DQN to play one step, and record its experience in the replay buffer:

```
def play_one_step(env, state, epsilon):
    action = epsilon_greedy_policy(state, epsilon)
    next_state, reward, done, truncated, info = env.step(action)
    replay_buffer.append((state, action, reward, next_state, done, truncated))
    return next_state, reward, done, truncated, info
```

Lastly, let's create a function that will sample some experiences from the replay buffer and perform a training step:

```
# extra code – for reproducibility, and to generate the next figure
env.reset(seed=42)
np.random.seed(42)
tf.random.set_seed(42)
rewards = []
best_score = 0
```

```
batch_size = 32
discount_factor = 0.95
optimizer = tf.keras.optimizers.Nadam(learning_rate=1e-2)
loss_fn = tf.keras.losses.mean_squared_error

def training_step(batch_size):
    experiences = sample_experiences(batch_size)
    states, actions, rewards, next_states, dones, truncateds = experiences
    next_Q_values = model.predict(next_states, verbose=0)
    max_next_Q_values = next_Q_values.max(axis=1)
    runs = 1.0 - (dones | truncateds)  # episode is not done or truncated
    target_Q_values = rewards + runs * discount_factor * max_next_Q_values
    target_Q_values = target_Q_values.reshape(-1, 1)
    mask = tf.one_hot(actions, n_outputs)
    with tf.GradientTape() as tape:
        all_Q_values = model(states)
        Q_values = tf.reduce_sum(all_Q_values * mask, axis=1, keepdims=True)
        loss = tf.reduce_mean(loss_fn(target_Q_values, Q_values))

    grads = tape.gradient(loss, model.trainable_variables)
    optimizer.apply_gradients(zip(grads, model.trainable_variables))
```

And now, let's train the model!

```
for episode in range(600):
    obs, info = env.reset()
    for step in range(200):
        epsilon = max(1 - episode / 500, 0.01)
        obs, reward, done, truncated, info = play_one_step(env, obs, epsilon)
        if done or truncated:
            break

    # extra code – displays debug info, stores data for the next figure, and
    #               keeps track of the best model weights so far
    print(f"\rEpisode: {episode + 1}, Steps: {step + 1}, eps: {epsilon:.3f}",
          end="")
    rewards.append(step)
    if step >= best_score:
        best_weights = model.get_weights()
        best_score = step

    if episode > 50:
        training_step(batch_size)

model.set_weights(best_weights)  # extra code – restores the best model weights
```
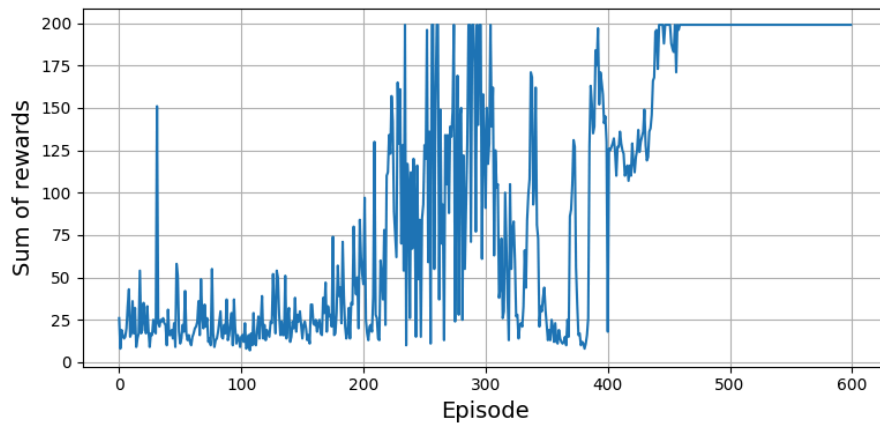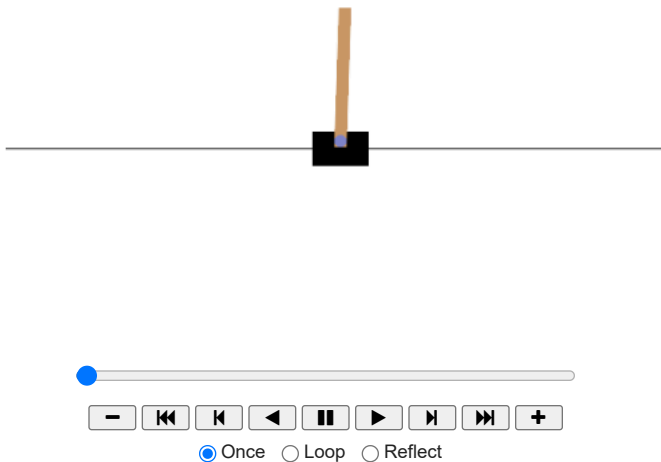
```
    Episode: 600, Steps: 146, eps: 0.010
```

```
# extra code – this cell generates and saves Figure 18–10
plt.figure(figsize=(8, 4))
plt.plot(rewards)
plt.xlabel("Episode", fontsize=14)
plt.ylabel("Sum of rewards", fontsize=14)
plt.grid(True)
save_fig("dqn_rewards_plot")
plt.show()
```

```
# extra code – shows an animation of the trained DQN playing one episode
show_one_episode(epsilon_greedy_policy)
```



Not bad at all! 😜

## ⌄  Fixed Q-Value Targets

Let's create the online DQN:

```
# extra code – creates the same DQN model as earlier

tf.random.set_seed(42)

model = tf.keras.Sequential([
    tf.keras.layers.Dense(32, activation="elu", input_shape=input_shape),
    tf.keras.layers.Dense(32, activation="elu"),
    tf.keras.layers.Dense(n_outputs)
])
```

```
/usr/local/lib/python3.10/dist-packages/ipykernel/ipkernel.py:283: DeprecationWarning: `should_run_async` will not call `transform_cell`
  and should_run_async(code)
```

Now create the target DQN: it's just a clone of the online DQN:

```
target = tf.keras.models.clone_model(model)  # clone the model's architecture
target.set_weights(model.get_weights())  # copy the weights
```

Next, we use the same code as above except for the line marked with # <= CHANGED:

```
env.reset(seed=42)
np.random.seed(42)
tf.random.set_seed(42)
rewards = []
best_score = 0

batch_size = 32
discount_factor = 0.95
optimizer = tf.keras.optimizers.Nadam(learning_rate=1e-2)
loss_fn = tf.keras.losses.mean_squared_error

replay_buffer = deque(maxlen=2000)  # resets the replay buffer

def training_step(batch_size):
    experiences = sample_experiences(batch_size)
    states, actions, rewards, next_states, dones, truncateds = experiences
    next_Q_values = target.predict(next_states, verbose=0)  # <= CHANGED
    max_next_Q_values = next_Q_values.max(axis=1)
    runs = 1.0 - (dones | truncateds)  # episode is not done or truncated
    target_Q_values = rewards + runs * discount_factor * max_next_Q_values
    target_Q_values = target_Q_values.reshape(-1, 1)
    mask = tf.one_hot(actions, n_outputs)
    with tf.GradientTape() as tape:
        all_Q_values = model(states)
        Q_values = tf.reduce_sum(all_Q_values * mask, axis=1, keepdims=True)
        loss = tf.reduce_mean(loss_fn(target_Q_values, Q_values))

    grads = tape.gradient(loss, model.trainable_variables)
    optimizer.apply_gradients(zip(grads, model.trainable_variables))
```

Again, this is the same code as earlier, except for the lines marked with # <= CHANGED:

```
for episode in range(600):
    obs, info = env.reset()
    for step in range(200):
        epsilon = max(1 - episode / 500, 0.01)
        obs, reward, done, info, truncated = play_one_step(env, obs, epsilon)
        if done or truncated:
            break

    # extra code – displays debug info, stores data for the next figure, and
    #               keeps track of the best model weights so far
    print(f"\rEpisode: {episode + 1}, Steps: {step + 1}, eps: {epsilon:.3f}",
          end="")
    rewards.append(step)
    if step >= best_score:
        best_weights = model.get_weights()
        best_score = step

    if episode > 50:
        training_step(batch_size)
        if episode % 50 == 0:                       # <= CHANGED
            target.set_weights(model.get_weights())  # <= CHANGED

    # Alternatively, you can do soft updates at each step:
    #if episode > 50:
        #training_step(batch_size)
        #target_weights = target.get_weights()
        #online_weights = model.get_weights()
        #for index, online_weight in enumerate(online_weights):
        #    target_weights[index] = (0.99 * target_weights[index]
        #                             + 0.01 * online_weight)
        #target.set_weights(target_weights)

model.set_weights(best_weights)  # extra code – restores the best model weights
```
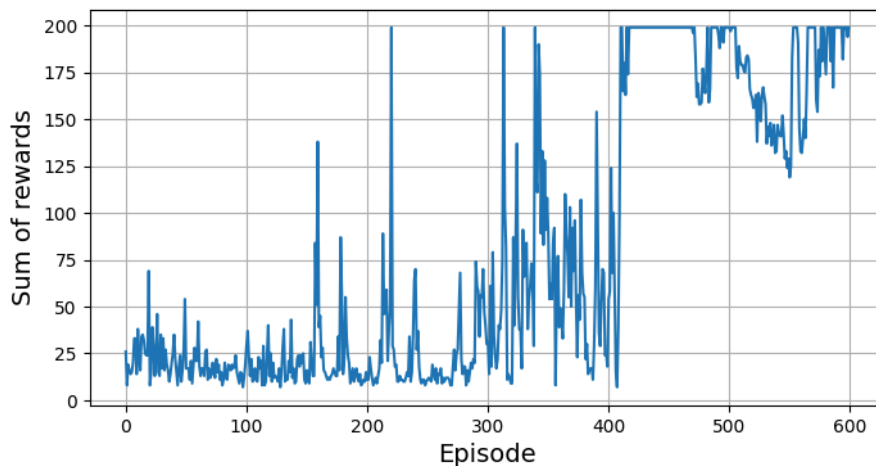
```
Episode: 600, Steps: 200, eps: 0.010
```

```
# extra code – this cell plots the learning curve
plt.figure(figsize=(8, 4))
plt.plot(rewards)
plt.xlabel("Episode", fontsize=14)
plt.ylabel("Sum of rewards", fontsize=14)
plt.grid(True)
plt.show()
```
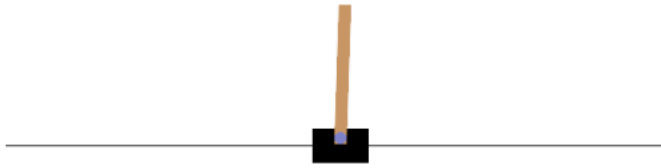


```
# extra code – shows an animation of the trained DQN playing one episode
show_one_episode(epsilon_greedy_policy)
```

## Double DQN

The code is exactly the same as for fixed Q-Value targets, except for the section marked as changed in the `training_step()` function:

```python
tf.random.set_seed(42)

model = tf.keras.Sequential([
    tf.keras.layers.Dense(32, activation="elu", input_shape=input_shape),
    tf.keras.layers.Dense(32, activation="elu"),
    tf.keras.layers.Dense(n_outputs)
])

target = tf.keras.models.clone_model(model)  # clone the model's architecture
target.set_weights(model.get_weights())  # copy the weights

env.reset(seed=42)
np.random.seed(42)
tf.random.set_seed(42)
rewards = []
best_score = 0

batch_size = 32
discount_factor = 0.95
optimizer = tf.keras.optimizers.Nadam(learning_rate=1e-2)
loss_fn = tf.keras.losses.mean_squared_error

def training_step(batch_size):
    experiences = sample_experiences(batch_size)
    states, actions, rewards, next_states, dones, truncateds = experiences

    #################### CHANGED SECTION ####################
    next_Q_values = model.predict(next_states, verbose=0)  # ≠ target.predict()
    best_next_actions = next_Q_values.argmax(axis=1)
    next_mask = tf.one_hot(best_next_actions, n_outputs).numpy()
    max_next_Q_values = (target.predict(next_states, verbose=0) * next_mask
                         ).sum(axis=1)
    ########################################################

    runs = 1.0 - (dones | truncateds)  # episode is not done or truncated
    target_Q_values = rewards + runs * discount_factor * max_next_Q_values
    target_Q_values = target_Q_values.reshape(-1, 1)
    mask = tf.one_hot(actions, n_outputs)
    with tf.GradientTape() as tape:
        all_Q_values = model(states)
        Q_values = tf.reduce_sum(all_Q_values * mask, axis=1, keepdims=True)
        loss = tf.reduce_mean(loss_fn(target_Q_values, Q_values))

    grads = tape.gradient(loss, model.trainable_variables)
    optimizer.apply_gradients(zip(grads, model.trainable_variables))

replay_buffer = deque(maxlen=2000)

for episode in range(600):
    obs, info = env.reset()
    for step in range(200):
        epsilon = max(1 - episode / 500, 0.01)
        obs, reward, done, info, truncated = play_one_step(env, obs, epsilon)
        if done or truncated:
            break

    print(f"\rEpisode: {episode + 1}, Steps: {step + 1}, eps: {epsilon:.3f}",
          end="")
    rewards.append(step)
    if step >= best_score:
        best_weights = model.get_weights()
        best_score = step

    if episode > 50:
        training_step(batch_size)
        if episode % 50 == 0:
            target.set_weights(model.get_weights())

model.set_weights(best_weights)

max_steps_per_episode = 200  # Máxima cantidad de pasos permitidos en un episodio
learning_threshold = max_steps_per_episode - 5  # Umbral para considerar que el modelo ha aprendido

if best_score >= learning_threshold:
    print("\n¡El modelo ha aprendido a mantener el poste en posición vertical!")
else:
```
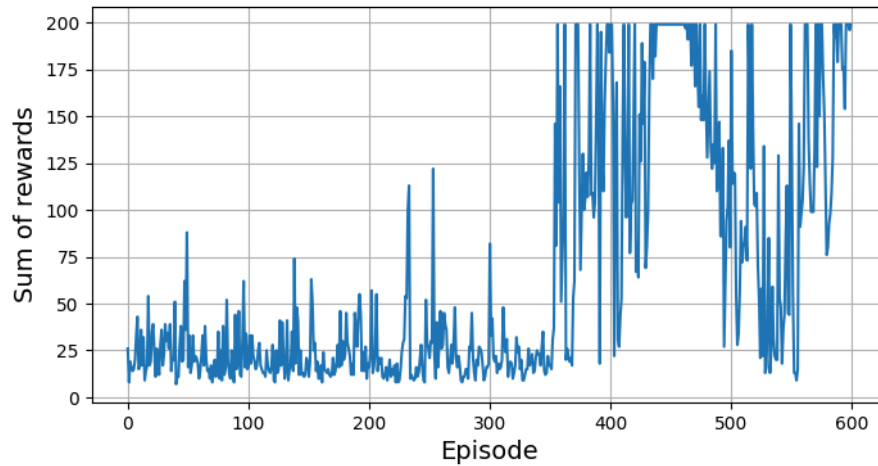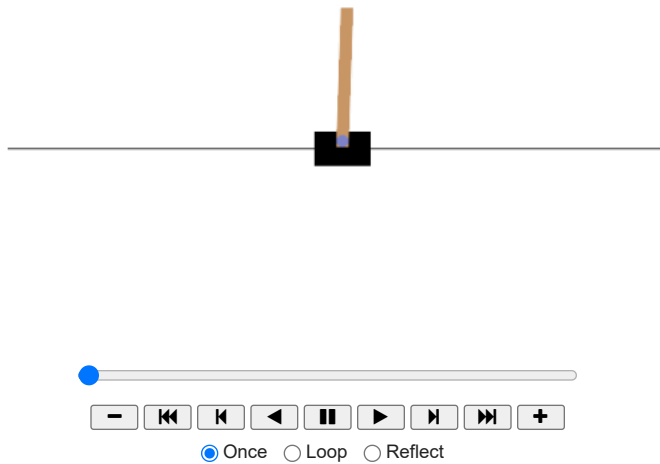
```
    print("\nEl modelo aún no ha aprendido a mantener el poste en posición vertical.")
```

```
    Episode: 481, Steps: 200, eps: 0.040
```

```python
# extra code – this cell plots the learning curve
plt.figure(figsize=(8, 4))
plt.plot(rewards)
plt.xlabel("Episode", fontsize=14)
plt.ylabel("Sum of rewards", fontsize=14)
plt.grid(True)
plt.show()
```



```python
# extra code – shows an animation of the trained DQN playing one episode
show_one_episode(epsilon_greedy_policy)
```



## Dueling Double DQN

```
tf.random.set_seed(42)  # extra code – ensures reproducibility on the CPU

input_states = tf.keras.layers.Input(shape=[4])
hidden1 = tf.keras.layers.Dense(32, activation="elu")(input_states)
hidden2 = tf.keras.layers.Dense(32, activation="elu")(hidden1)
state_values = tf.keras.layers.Dense(1)(hidden2)
raw_advantages = tf.keras.layers.Dense(n_outputs)(hidden2)
advantages = raw_advantages - tf.reduce_max(raw_advantages, axis=1,
                                            keepdims=True)
Q_values = state_values + advantages
model = tf.keras.Model(inputs=[input_states], outputs=[Q_values])
```

```
/usr/local/lib/python3.10/dist-packages/ipykernel/ipkernel.py:283: DeprecationWarning: `should_run_async` will not call `transform_cell`
  and should_run_async(code)
```

◀                                                                                                          ▶

The rest is the same code as earlier:

```
# extra code – trains the model

batch_size = 32
discount_factor = 0.95
optimizer = tf.keras.optimizers.Nadam(learning_rate=5e-3)
loss_fn = tf.keras.losses.mean_squared_error

target = tf.keras.models.clone_model(model)  # clone the model's architecture
target.set_weights(model.get_weights())  # copy the weights

env.reset(seed=42)
replay_buffer = deque(maxlen=2000)
rewards = []
best_score = 0

for episode in range(600):
    obs, info = env.reset()
    for step in range(200):
        epsilon = max(1 - episode / 500, 0.01)
        obs, reward, done, info, truncated = play_one_step(env, obs, epsilon)
        if done or truncated:
            break

    print(f"\rEpisode: {episode + 1}, Steps: {step + 1}, eps: {epsilon:.3f}",
          end="")
    rewards.append(step)
    if step >= best_score:
        best_weights = model.get_weights()
        best_score = step

    if episode > 50:
        training_step(batch_size)
        if episode % 50 == 0:
            target.set_weights(model.get_weights())

model.set_weights(best_weights)
```
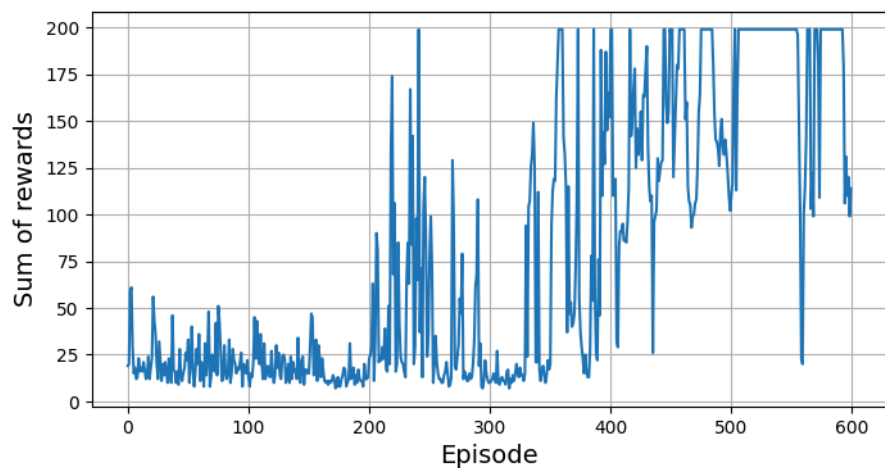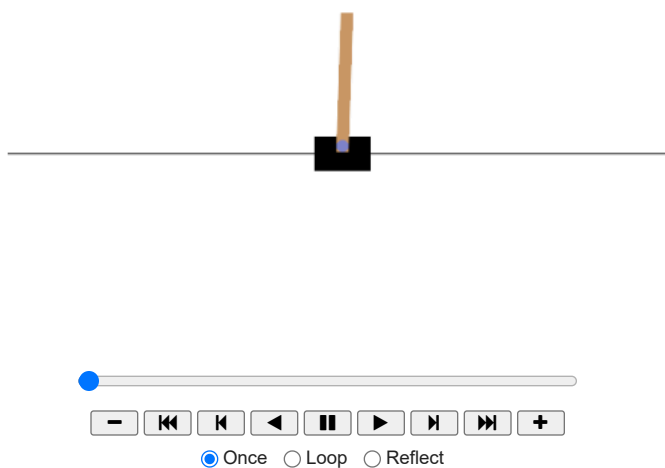
```
    Episode: 600, Steps: 115, eps: 0.010
```

```
# extra code – this cell plots the learning curve
plt.figure(figsize=(8, 4))
plt.plot(rewards)
plt.xlabel("Episode", fontsize=14)
plt.ylabel("Sum of rewards", fontsize=14)
plt.grid(True)
plt.show()
```
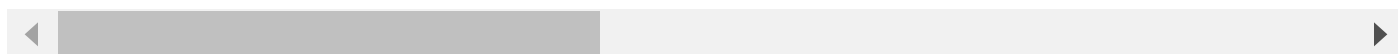
```
# extra code – shows an animation of the trained DQN playing one episode
show_one_episode(epsilon_greedy_policy)
```



This looks like a pretty robust agent!

```
env.close()
```

```
/usr/local/lib/python3.10/dist-packages/ipykernel/ipkernel.py:283: DeprecationWarning: `should_run_async` will not call `transform_cell`
  and should_run_async(code)
```

## ✓ Exercise Solutions

### ✓ 1. to 7.

1. Reinforcement Learning is an area of Machine Learning aimed at creating agents capable of taking actions in an environment in a way
   that maximizes rewards over time. There are many differences between RL and regular supervised and unsupervised learning. Here are a

few:

- In supervised and unsupervised learning, the goal is generally to find patterns in the data and use them to make predictions. In Reinforcement Learning, the goal is to find a good policy.
- Unlike in supervised learning, the agent is not explicitly given the "right" answer. It must learn by trial and error.
- Unlike in unsupervised learning, there is a form of supervision, through rewards. We do not tell the agent how to perform the task, but we do tell it when it is making progress or when it is failing.
- A Reinforcement Learning agent needs to find the right balance between exploring the environment, looking for new ways of getting rewards, and exploiting sources of rewards that it already knows. In contrast, supervised and unsupervised learning systems generally don't need to worry about exploration; they just feed on the training data they are given.
- In supervised and unsupervised learning, training instances are typically independent (in fact, they are generally shuffled). In Reinforcement Learning, consecutive observations are generally *not* independent. An agent may remain in the same region of the environment for a while before it moves on, so consecutive observations will be very correlated. In some cases a replay buffer (memory) is used to ensure that the training algorithm gets fairly independent observations.

2. Here are a few possible applications of Reinforcement Learning, other than those mentioned in Chapter 18:

- Music personalization: The environment is a user's personalized web radio. The agent is the software deciding what song to play next for that user. Its possible actions are to play any song in the catalog (it must try to choose a song the user will enjoy) or to play an advertisement (it must try to choose an ad that the user will be interested in). It gets a small reward every time the user listens to a song, a larger reward every time the user listens to an ad, a negative reward when the user skips a song or an ad, and a very negative reward if the user leaves.
- Marketing: The environment is your company's marketing department. The agent is the software that defines which customers a mailing campaign should be sent to, given their profile and purchase history (for each customer it has two possible actions: send or don't send). It gets a negative reward for the cost of the mailing campaign, and a positive reward for estimated revenue generated from this campaign.
- Product delivery: Let the agent control a fleet of delivery trucks, deciding what they should pick up at the depots, where they should go, what they should drop off, and so on. It will get positive rewards for each product delivered on time, and negative rewards for late deliveries.

3. When estimating the value of an action, Reinforcement Learning algorithms typically sum all the rewards that this action led to, giving more weight to immediate rewards and less weight to later rewards (considering that an action has more influence on the near future than on the distant future). To model this, a discount factor is typically applied at each time step. For example, with a discount factor of 0.9, a reward of 100 that is received two time steps later is counted as only $0.9^2 \times 100 = 81$ when you are estimating the value of the action. You can think of the discount factor as a measure of how much the future is valued relative to the present: if it is very close to 1, then the future is valued almost as much as the present; if it is close to 0, then only immediate rewards matter. Of course, this impacts the optimal policy tremendously: if you value the future, you may be willing to put up with a lot of immediate pain for the prospect of eventual rewards, while if you don't value the future, you will just grab any immediate reward you can find, never investing in the future.

4. To measure the performance of a Reinforcement Learning agent, you can simply sum up the rewards it gets. In a simulated environment, you can run many episodes and look at the total rewards it gets on average (and possibly look at the min, max, standard deviation, and so on).

5. The credit assignment problem is the fact that when a Reinforcement Learning agent receives a reward, it has no direct way of knowing which of its previous actions contributed to this reward. It typically occurs when there is a large delay between an action and the resulting reward (e.g., during a game of Atari's *Pong*, there may be a few dozen time steps between the moment the agent hits the ball and the moment it wins the point). One way to alleviate it is to provide the agent with shorter-term rewards, when possible. This usually requires prior knowledge about the task. For example, if we want to build an agent that will learn to play chess, instead of giving it a reward only when it wins the game, we could give it a reward every time it captures one of the opponent's pieces.

6. An agent can often remain in the same region of its environment for a while, so all of its experiences will be very similar for that period of time. This can introduce some bias in the learning algorithm. It may tune its policy for this region of the environment, but it will not perform well as soon as it moves out of this region. To solve this problem, you can use a replay buffer; instead of using only the most immediate experiences for learning, the agent will learn based on a buffer of its past experiences, recent and not so recent (perhaps this is why we dream at night: to replay our experiences of the day and better learn from them?).

7. An off-policy RL algorithm learns the value of the optimal policy (i.e., the sum of discounted rewards that can be expected for each state if the agent acts optimally) while the agent follows a different policy. Q-Learning is a good example of such an algorithm. In contrast, an on-policy algorithm learns the value of the policy that the agent actually executes, including both exploration and exploitation.

⌄ 8.

*Exercise: Use policy gradients to solve* ~~*OpenAI Gym*~~ *Gymnasium's LunarLander-v2 environment.*

Let's start by creating a LunarLander-v2 environment:

```
env = gym.make("LunarLander-v2", render_mode="rgb_array")
```

The inputs are 8-dimensional:

```
env.observation_space
```

```
    Box([-1.5       -1.5       -5.       -5.       -3.1415927 -5.
     -0.        -0.       ], [1.5       1.5       5.       5.       3.1415927 5.       1.
     1.       ], (8,), float32)
```

```
obs, info = env.reset(seed=42)
obs
```

```
    array([ 0.00229702,  1.4181306 ,  0.2326471 ,  0.3204666 , -0.00265488,
           -0.05269805,  0.        ,  0.        ], dtype=float32)
```

From the source code, we can see that these each 8D observation (x, y, h, v, a, w, l, r) correspond to:

- x,y: the coordinates of the spaceship. It starts at a random location near (0, 1.4) and must land near the target at (0, 0).
- h,v: the horizontal and vertical speed of the spaceship. It starts with a small random speed.
- a,w: the spaceship's angle and angular velocity.
- l,r: whether the left or right leg touches the ground (1.0) or not (0.0).

The action space is discrete, with 4 possible actions:

```
env.action_space
```

```
    Discrete(4)
```

Looking at the LunarLander-v2's description, these actions are:

- do nothing
- fire left orientation engine
- fire main engine
- fire right orientation engine

Let's create a simple policy network with 4 output neurons (one per possible action):

```
np.random.seed(42)
tf.random.set_seed(42)

n_inputs = env.observation_space.shape[0]
n_outputs = env.action_space.n

model = tf.keras.Sequential([
    tf.keras.layers.Dense(32, activation="relu", input_shape=[n_inputs]),
    tf.keras.layers.Dense(32, activation="relu"),
    tf.keras.layers.Dense(n_outputs, activation="softmax"),
])
```

Note that we're using the softmax activation function in the output layer, instead of the sigmoid activation function like we did for the CartPole-v1 environment. This is because we only had two possible actions for the CartPole-v1 environment, so a binary classification model worked fine. However, since we now how more than two possible actions, we need a multiclass classification model.

Next, let's reuse the `play_one_step()` and `play_multiple_episodes()` functions we defined for the CartPole-v1 Policy Gradient code above, but we'll just tweak the `play_one_step()` function to account for the fact that the model is now a multiclass classification model rather than a binary classification model. We'll also tweak the `play_multiple_episodes()` function to call our tweaked `play_one_step()` function rather than the original one, and we add a big penalty if the spaceship does not land (or crash) before a maximum number of steps.

```python
def lander_play_one_step(env, obs, model, loss_fn):
    with tf.GradientTape() as tape:
        probas = model(obs[np.newaxis])
        logits = tf.math.log(probas + tf.keras.backend.epsilon())
        action = tf.random.categorical(logits, num_samples=1)
        loss = tf.reduce_mean(loss_fn(action, probas))
    grads = tape.gradient(loss, model.trainable_variables)
    obs, reward, done, info, truncated = env.step(action[0, 0].numpy())
    return obs, reward, done, truncated, grads


def lander_play_multiple_episodes(env, n_episodes, n_max_steps, model, loss_fn):
    all_rewards = []
    all_grads = []
    for episode in range(n_episodes):
        current_rewards = []
        current_grads = []
        obs, info = env.reset()
        for step in range(n_max_steps):
```