

50.039 Theory and Practice of Deep Learning

W12S2 Introduction to Reinforcement Learning

Matthieu De Mari, Berrak Sisman



About this week (Week 12)

1. What is **Reinforcement Learning**?
2. What are the key ideas behind **reinforcement learning** and its **framework**?
3. What is the **exploration vs. exploitation tradeoff**?
4. How do we **train** an **RL agent** by exploring, then progressively exploiting?
5. What are some **advanced strategies** in **multi-arm bandit problems**?
6. What are the **Q** and **V functions** for a RL problem?
7. What is **Q-learning** and how can it be implemented in RL problems?

About this week (Week 12)

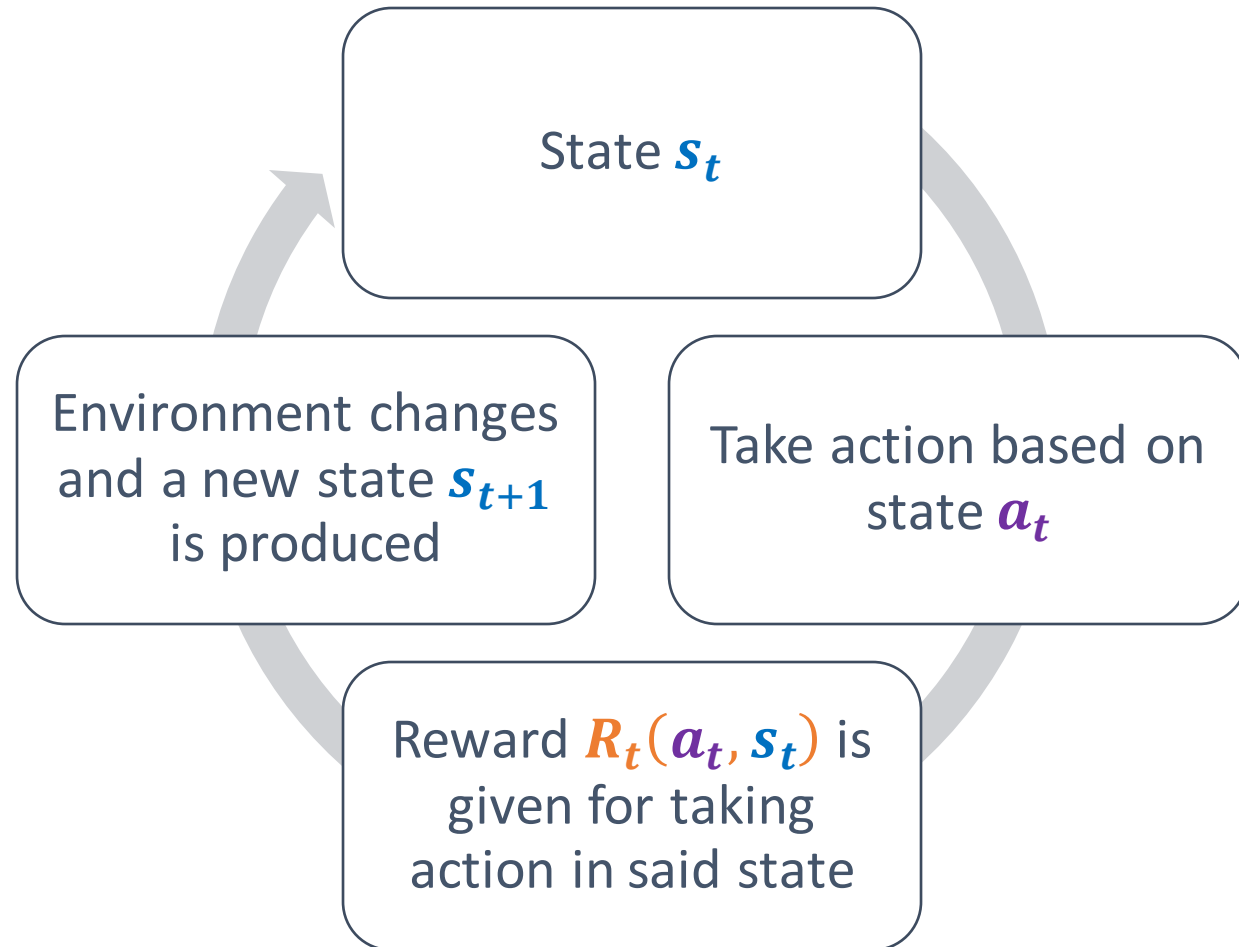
8. What is **Deep Q Learning**? And which problem does it address?
9. What is a **trainer** and a **main Q-network**? What is the interleaved training for Deep Q learning?
10. What are **actor-critic** learning methods? And which problems do these approaches address?
11. What are more advanced problems in RL?
 - Markov states
 - Partially observable environment
 - SARSA
 - Non-stationary problems

Reusing the RL formalism

Definition (**agent**):

In RL, we refer to the AI, as an **agent**. At each step, the **agent**:

- Looks at the current **state** s_t ,
- Then takes an **action**, in this given state, a_t .
- The **action** has an effect, which is eventually measured in terms of **reward**, R_t .
- And a **new state** s_{t+1} is produced.



Exploration and exploitation tradeoff

Definition (**exploration** and **exploitation**):

The phase during which, you try out things to acquire knowledge about the problem is called **exploration**.

The second phase, where you rely on your acquired knowledge, and play what you feel is the best move, is called **exploitation**.

Definition (**exploration vs. exploitation tradeoff**):

A good RL-based AI, needs to smartly combine **exploration** and **exploitation** phases.

- **Too much exploration?** You have wasted coins trying out bad machines.
- **No enough exploration?** You might end up choosing the wrong machine as the “best” one.

A second toy problem

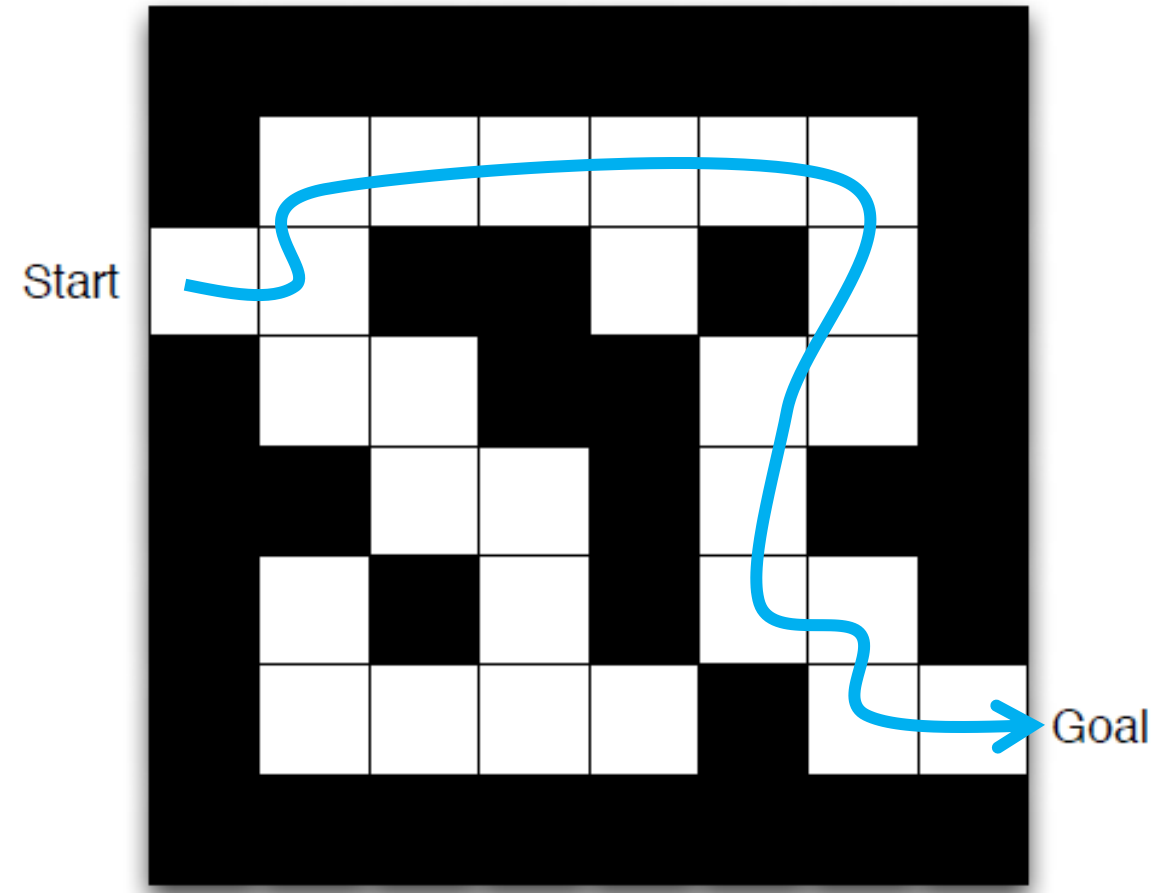
The Maze problem

The maze problem

Let us consider the maze problem,
described on the right.

The objective is

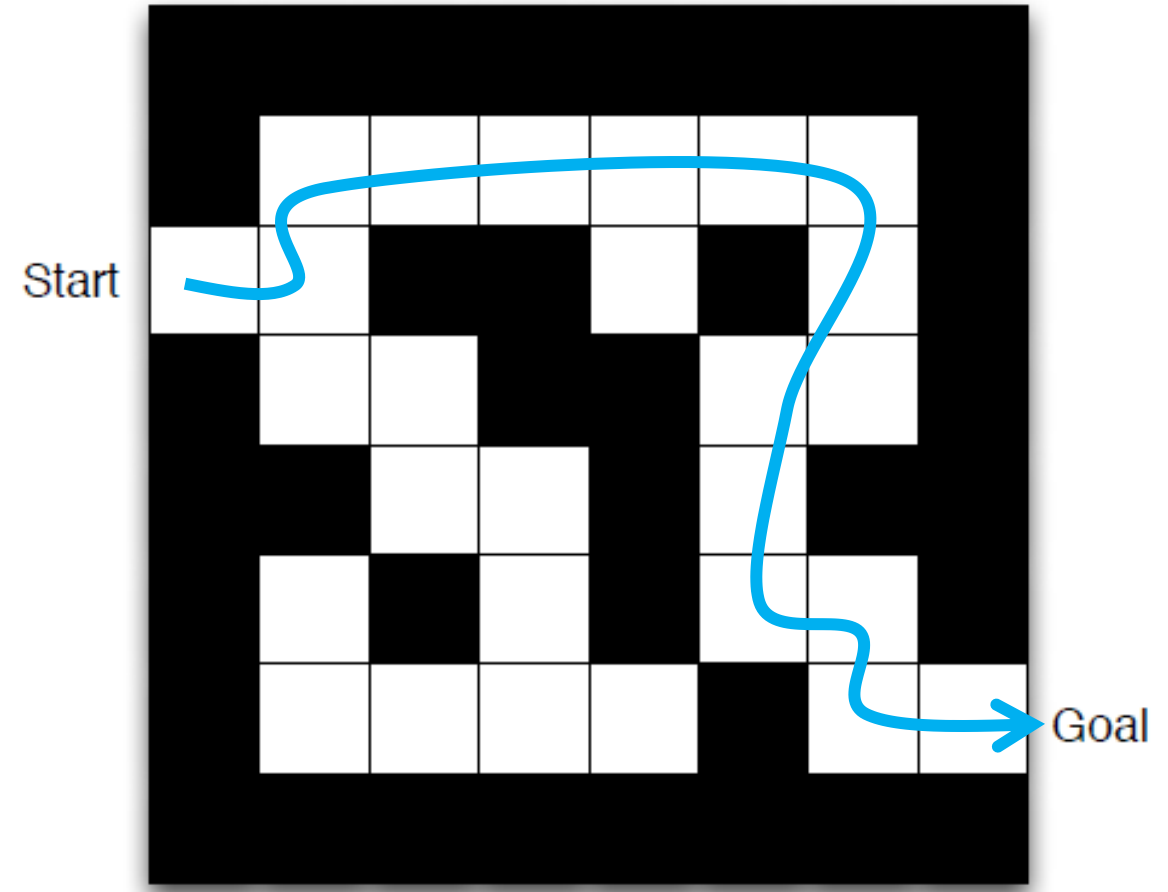
- to find the shortest path between the start and goal squares,
- without any knowledge of what the maze is beforehand.



The maze problem

Agent state:

- Our agent will start at the starting position, with coordinates **(3,1)**. This is the initial **state s_1** .



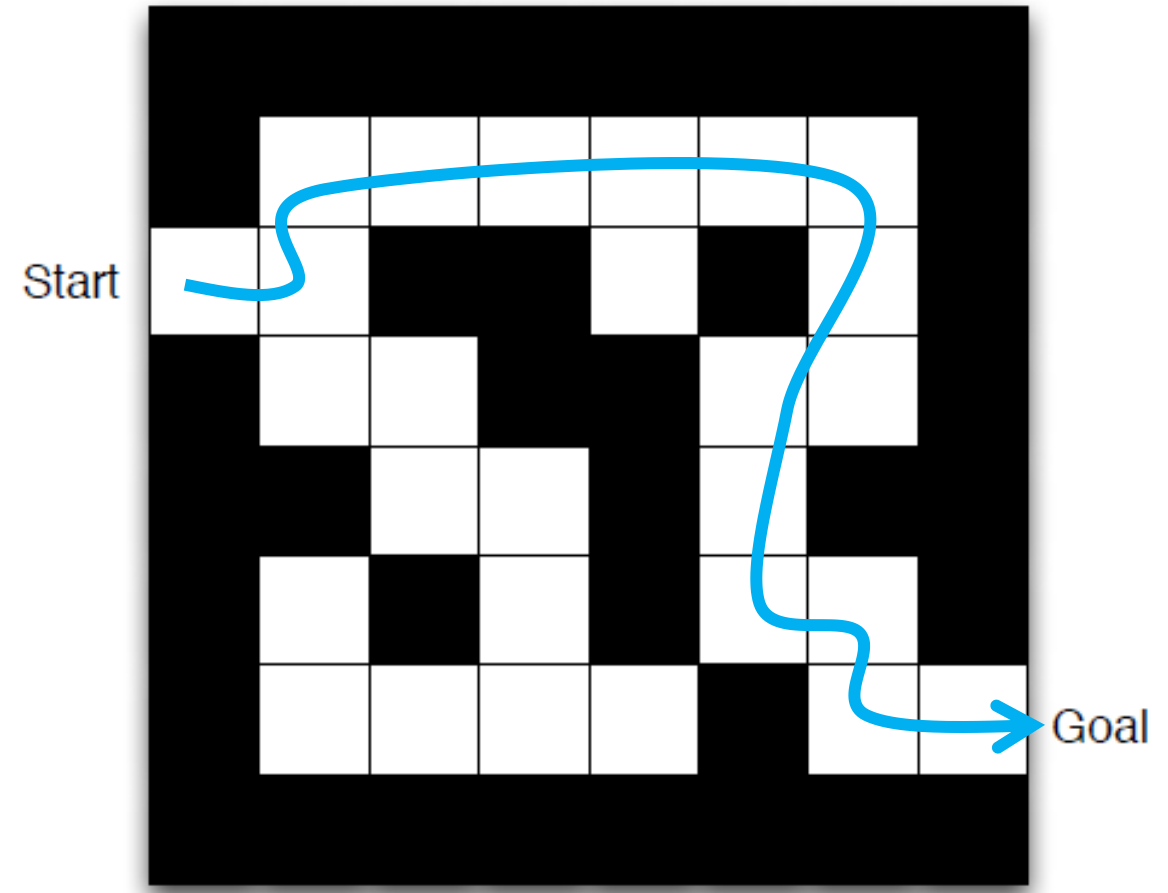
The maze problem

Agent state:

- Our agent will start at the starting position, with coordinates **(3,1)**. This is the initial **state s_1** .

Action space:

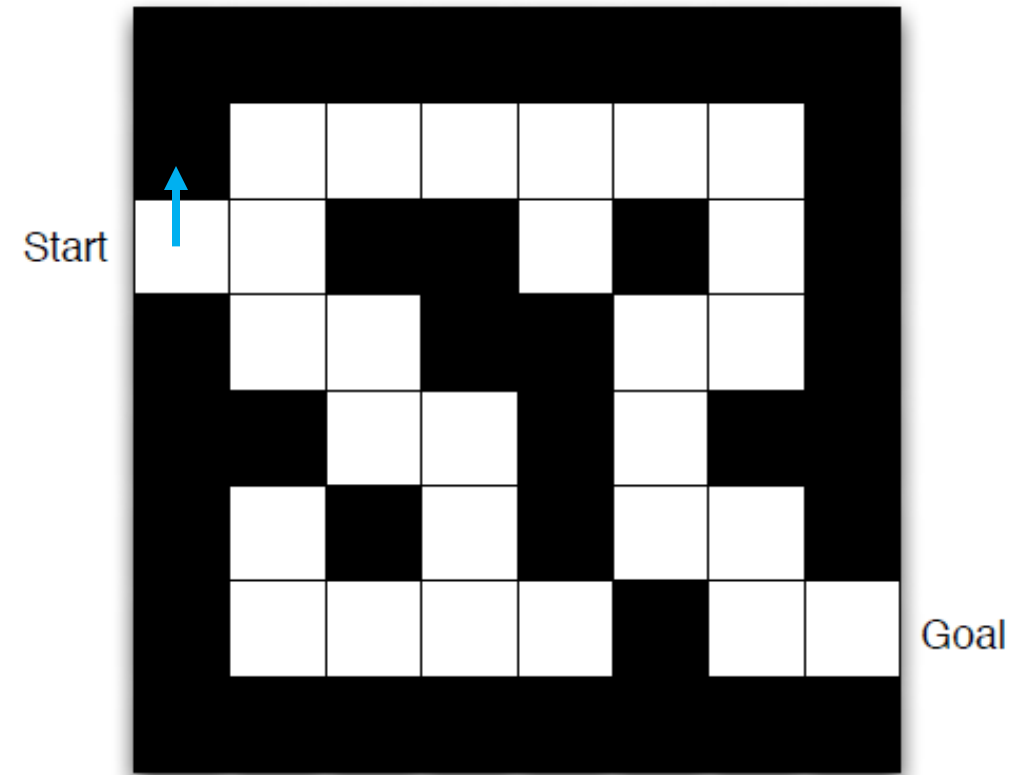
- In every position of the maze, four **actions** are possible
 $A = \{Go\ North, Go\ South, Go\ East, Go\ West\}$



The maze problem

State update:

- For a given **state** s_t , and an **action** a_t , the **state** will be updated following these rules.
- If the **action** a_t moves the user in a wall or left of the start point, the **new state** s_{t+1} is the same as the **previous state** s_t .

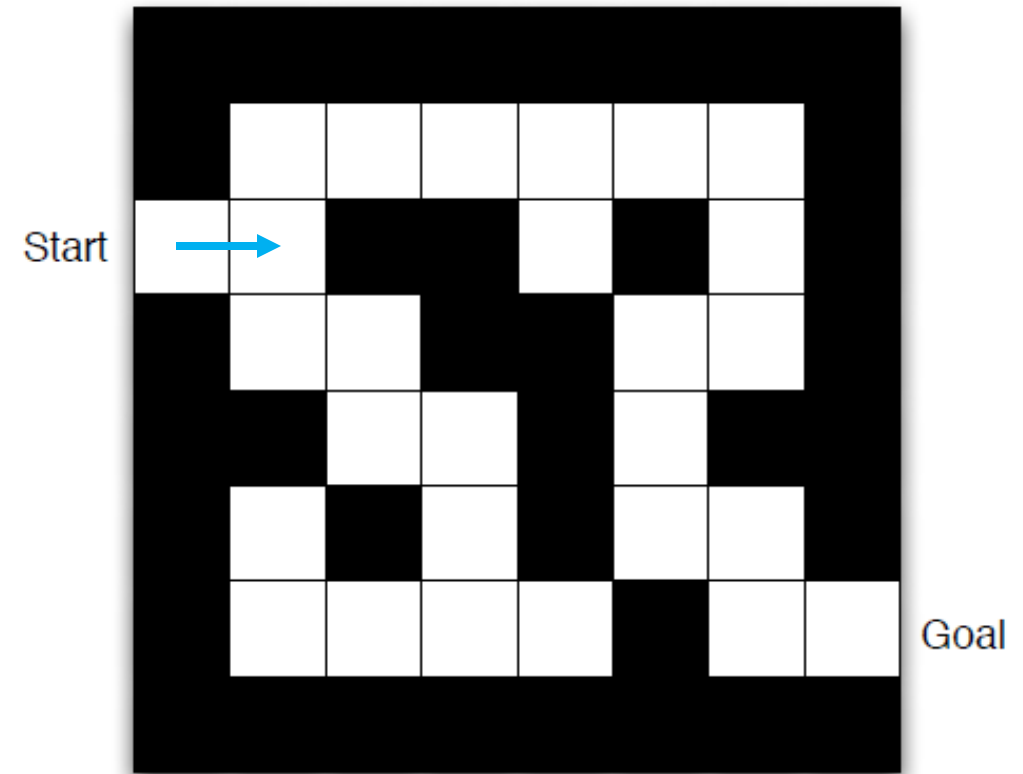


The maze problem

State update:

- For a given **state** s_t , and an **action** a_t , the **state** will be updated following these rules.
- If the **action** a_t moves the user in a wall or left of the start point, the **new state** s_{t+1} is the same as the **previous state** s_t .
- Otherwise, move the use to next square and this becomes s_{t+1} .

$$s_t = (3, 1) + a_t = \text{East} \rightarrow s_{t+1} = (3, 2)$$



The maze problem

Reward value:

- For every **state** and every **action** of the agent, the **reward** R_t is set to -1.
- The game stops when the agent reaches the Goal square, i.e. **state becoming (8, 7)**.
- The **cumulated reward/gain** $G_t = \sum_t R_t$ will be equal to minus 1 times the number of steps taken to get from the start square to the goal square.
- Maximizing $G_t = \sum_t R_t$ is then strictly equivalent to **finding the shortest path out of the maze**.
- This will be our **RL framework** for this task.

Value function V

Definition (**value function V**):

The **value function V** is a prediction/estimation of the future cumulated reward/gain, if in state s_t at time t , for the policy π .

$$\forall t, \forall s_t, \quad V_t^\pi(s_t) = E_\pi[R_t + R_{t+1} + R_{t+2} + \dots | s_t]$$

In general, we add a parameter $\gamma \in [0,1]$, which gives more or less importance to the future or present rewards.

$$\forall t, \forall s_t, \quad V_t^\pi(s_t) = E_\pi[R_t + \gamma R_{t+1} + \gamma^2 R_{t+2} \dots | s_t]$$

Value function V

Definition (**value function V**):

The **value function V** is a prediction/estimation of the future cumulated reward/gain, if in state s_t at time t , for the policy π .

In general, we add a parameter $\gamma \in [0,1]$, which gives more or less importance to the future or present rewards, in the same way as regularization.

$$\forall t, \forall s_t, \quad V_t^\pi(s_t) = E_\pi[R_t + \gamma R_{t+1} + \gamma^2 R_{t+2} \dots | s_t]$$

This function is used to **evaluate the goodness/badness of a state s_t** , and **can help to identify the best action to use in a given state**.

Optimization with value function V

The **value function V** can be used to rewrite our optimization problem, one step at a time.

Let us assume we are current at time t , in **state s_t** .

The best **action a_t** to use in this current state will simply maximize

- The **immediate reward** we will get after playing this action, that is R_t
- Plus the **expected reward** we will get in the future, if we end up in a new **state s_{t+1}** after playing **action a_t** .

$$a_t = \arg \max_{a \in A} [R_t(a, s_t) + V_{t+1}^\pi(s_{t+1})]$$

The state-action function Q

Definition (state-action function Q):

The **state-action function Q** is a function, which is used to quantify the goodness/badness of taking **action a_t** in **state s_t** at time t , according to our policy π .

It is closely related to the value function V , as it consists of our previous optimization function term, which was combining

- The **immediate reward** we will get after playing this action, that is R_t
- Plus the **expected reward** we will get in the future, if we end up in a **new state s_{t+1}** after playing **action a_t** .

$$Q_t^\pi(a_t, s_t) = R_t(a_t, s_t) + V_{t+1}^\pi(s_{t+1})$$

On the policy, V and Q functions relationship

Theorem (Bellman principle of optimality):

The policy π , V and Q functions are closely related. For instance, we have a clear relationship between Q and V . The value of Q can be used to compute the value of V immediately, and vice-versa.

$$Q_t^\pi(a_t, s_t) = R_t(a, s_t) + V_{t+1}^\pi(s_{t+1})$$

The policy π , can then be derived from Q , as:

$$a_t = \pi_t(s_t) = \arg \max_{a \in A} Q_t^\pi(a, s_t)$$

Defining a Q-table

- In our problem, we have a finite number of **actions** and **states**.
- It is then simpler to write the Q function as a table, which is then referred to as the **Q -table**.
- Later on, we will let our agent explore the maze and update the values of the Q -table.
- Finding the converged Q -table, later gives us the best policy π to use in any state.

Initialized

Q-Table		Actions			
		South (0)	North (1)	East (2)	West (3)
States	0	0	0	0	0

	327	0	0	0	0

	499	0	0	0	0

Training

Q-Table		Actions			
		South (0)	North (1)	East (2)	West (3)
States	0	0	0	0	0

	328	-2.30108105	-1.97092096	-2.30357004	-2.20591839

	499	9.96984239	4.02706992	12.96022777	29

Updating the Q-table via exploration

- To update the Q -table, we will let our agent play a few rounds of the game and “explore” the maze.
- At the end of each round, we obtain a finite **history of actions, states and rewards**, i.e. a sequence of values for run k

$$h_k = \{s_1, a_1, R_1, s_2, a_2, R_2, \dots\}$$

Updating the Q-table via exploration

- This sequence can then be used, to update the Q table, according to

$$Q^{new}(s_t, a_t) \leftarrow \underbrace{Q(s_t, a_t)}_{\text{old value}} + \underbrace{\alpha}_{\text{learning rate}} \cdot \underbrace{\left(\underbrace{r_t}_{\text{reward}} + \underbrace{\gamma}_{\text{discount factor}} \cdot \underbrace{\max_a Q(s_{t+1}, a)}_{\text{estimate of optimal future value}} - \underbrace{Q(s_t, a_t)}_{\text{old value}} \right)}_{\text{new value (temporal difference target)}}$$

temporal difference

- In a sense, this formula is the “**equivalent**” of our **gradient descent update**, but in the case of reinforcement learning, and is commonly referred to as **Q-learning**.

Updating the Q table via exploration

- A quick note on its parameters...

$$Q^{new}(s_t, a_t) \leftarrow \underbrace{Q(s_t, a_t)}_{\text{old value}} + \underbrace{\alpha}_{\text{learning rate}} \cdot \underbrace{\left(\underbrace{r_t}_{\text{reward}} + \underbrace{\gamma}_{\text{discount factor}} \cdot \underbrace{\max_a Q(s_{t+1}, a)}_{\text{estimate of optimal future value}} - \underbrace{Q(s_t, a_t)}_{\text{old value}} \right)}_{\text{new value (temporal difference target)}}$$

temporal difference

- The learning rate α determines to what extent newly acquired information overrides old information in the Q-table.
- A value 0 makes the agent learn nothing (exclusively exploiting prior knowledge), while a factor of 1 makes the agent consider only the most recent information (ignoring prior knowledge to explore possibilities).

Updating the Q table via exploration

- A quick note on its parameters...

$$Q^{new}(s_t, a_t) \leftarrow \underbrace{Q(s_t, a_t)}_{\text{old value}} + \underbrace{\alpha}_{\text{learning rate}} \cdot \underbrace{\left(\underbrace{r_t}_{\text{reward}} + \underbrace{\gamma}_{\text{discount factor}} \cdot \underbrace{\max_a Q(s_{t+1}, a)}_{\text{estimate of optimal future value}} - \underbrace{Q(s_t, a_t)}_{\text{old value}} \right)}_{\text{new value (temporal difference target)}}$$

temporal difference

- The discount factor γ determines the importance of future rewards.
- A value 0 will make the agent "myopic" (or short-sighted) by only considering current rewards, i.e. R_t (in the update rule above).
- A value approaching 1 will make it strive for a long-term high reward.

Updating the Q table via exploration

- A quick note on its parameters...

$$Q^{new}(s_t, a_t) \leftarrow \underbrace{Q(s_t, a_t)}_{\text{old value}} + \underbrace{\alpha}_{\text{learning rate}} \cdot \underbrace{\left(\underbrace{r_t}_{\text{reward}} + \underbrace{\gamma}_{\text{discount factor}} \cdot \underbrace{\max_a Q(s_{t+1}, a)}_{\text{estimate of optimal future value}} - \underbrace{Q(s_t, a_t)}_{\text{old value}} \right)}_{\text{new value (temporal difference target)}}$$

temporal difference

- Q-learning implicitly assumes some initial values in the Q -table before the first update occurs.
- High initial values, also known as "optimistic initial conditions" tend to encourage exploration.

Training an agent

- Training an agent then requires to define a policy π , which will smoothly transition from exploration into exploitation, as in the random candy machines before.
- Upon seeing “convergence” on the values of the Q -table, we can then claim that the agent has been “properly” trained.
- From there, exploiting should then give the best strategy.

Initialized

Q-Table		Actions			
		South (0)	North (1)	East (2)	West (3)
States	0	0	0	0	0

	327	0	0	0	0

	499	0	0	0	0

Training

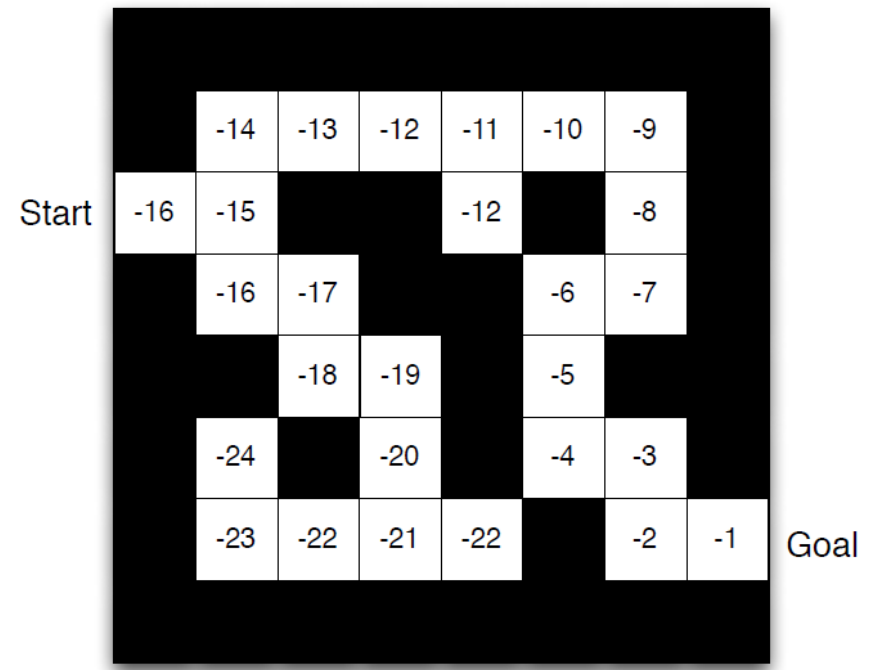
Q-Table		Actions			
		South (0)	North (1)	East (2)	West (3)
States	0	0	0	0	0

	328	-2.30108105	-1.97092096	-2.30357004	-2.20591839

	499	9.96984239	4.02706992	12.96022777	29

Training an agent

- By successfully updating V and Q at the end of each round, we will eventually obtain a good estimate of the “true” V (or alternatively Q) values for our problem.
- That matches the minimal number of steps which should be taken from any given square of the maze to reach the end.



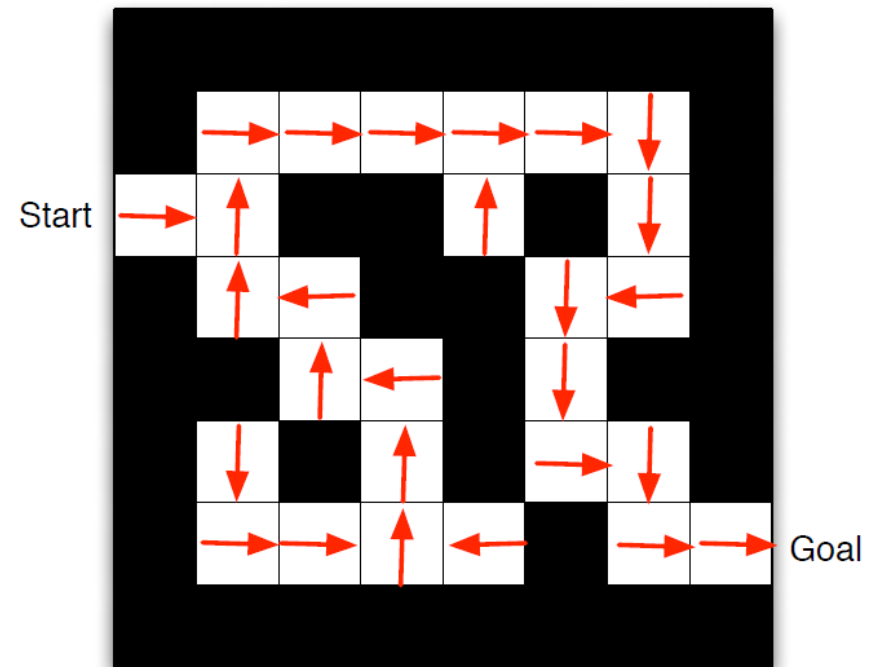
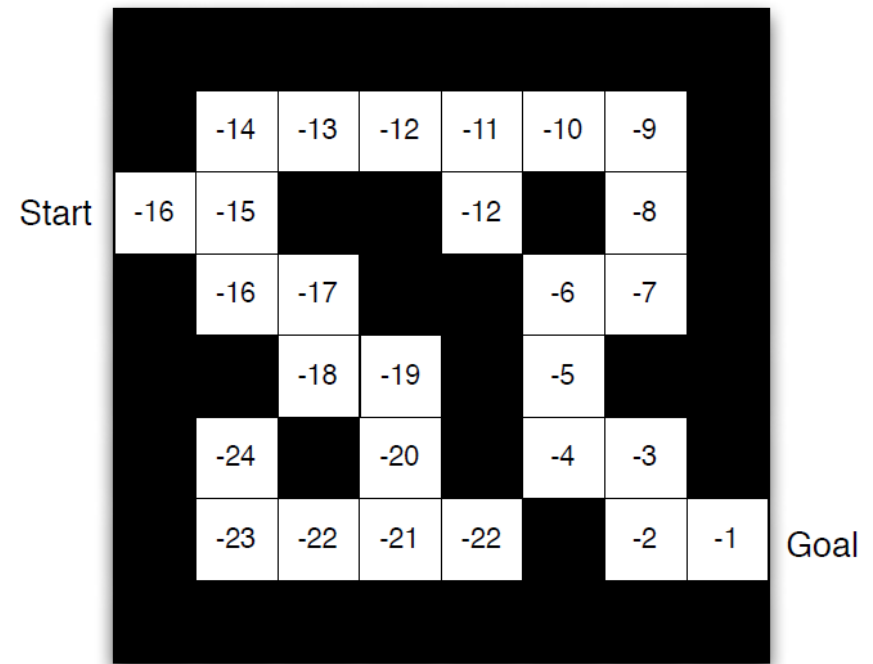
Note: the values above are the V ones, literally telling us the number of steps needed to reach goal from any position.

Training an agent

- Reusing the V (or Q) value allows to define the optimal policy π to use in any given state.
- That is, the optimal direction in which we should go for any given square.

$$a_t = \pi_t(s_t) = \arg \max_{a \in A} Q_t^\pi(a, s_t)$$

$$Q_t^\pi(a_t, s_t) = R_t(a, s_t) + V_{t+1}^\pi(s_{t+1})$$



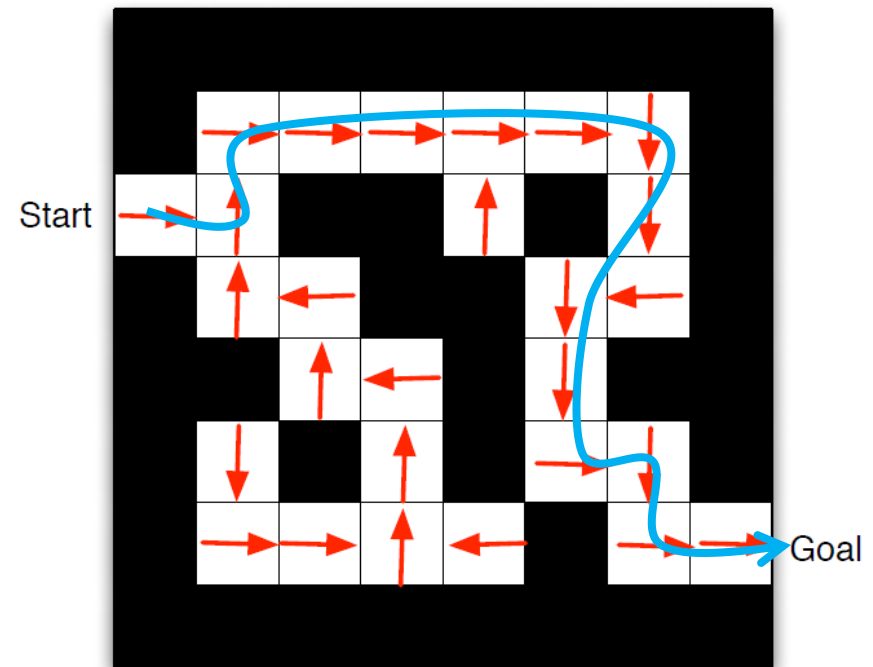
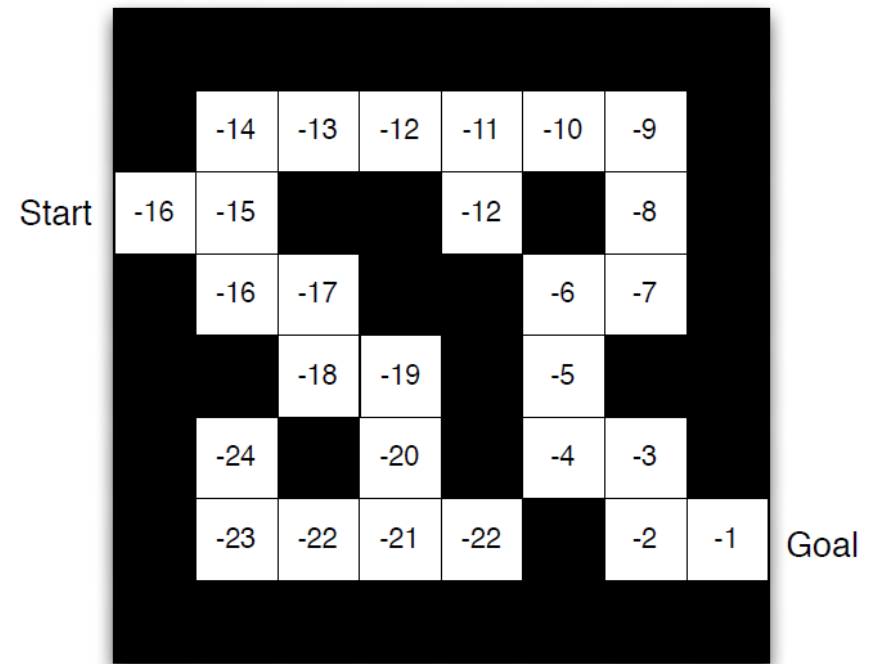
Training an agent

Later on, we can then replay the game,

- in full exploitation mode, i.e. by playing the best action according to our policy π every time,
- without any exploration moves.

This gives us the shortest path.

Our agent has then learnt to recognize the maze and figured the shortest path through trial and error!



A representation problem

- **Problem:** in many RL problems, the states and/or actions sets are not necessarily finite.
- In that case, it is impossible to represent the Q and V functions as tables.
- And, even worse, for these problems, coming up with a closed form expression of the V and Q functions might prove challenging.

Initialized

Q-Table		Actions			
		South (0)	North (1)	East (2)	West (3)
States	0	0	0	0	0

	327	0	0	0	0

	499	0	0	0	0

Training

Q-Table		Actions			
		South (0)	North (1)	East (2)	West (3)
States	0	0	0	0	0

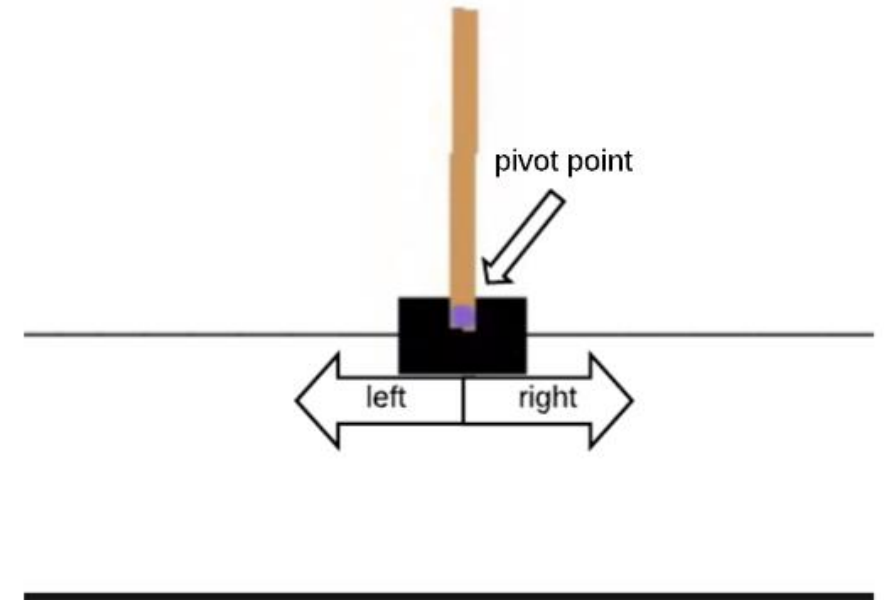
	328	-2.30108105	-1.97092096	-2.30357004	-2.20591839

	499	9.96984239	4.02706992	12.96022777	29

Toy example #3 (shown, but implementation out of the scope of this class)

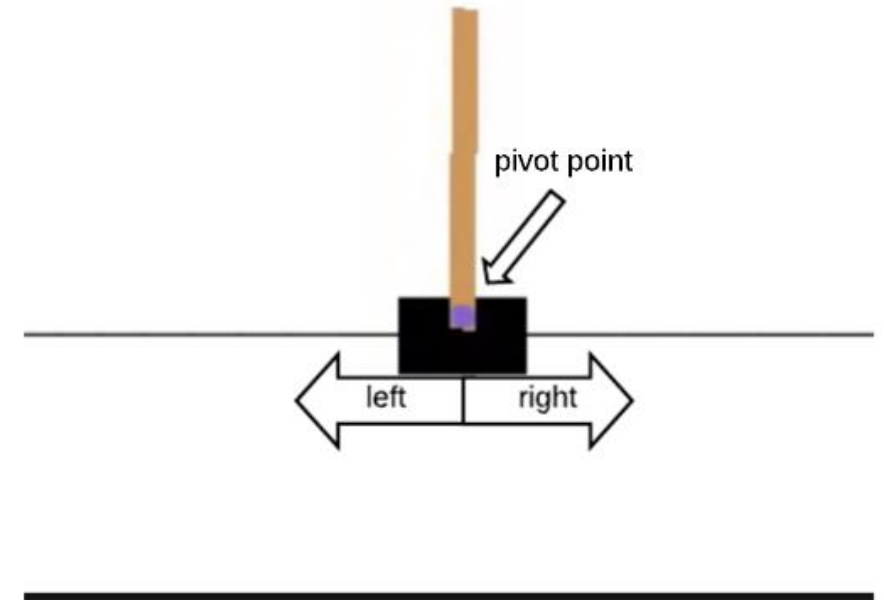
Example: the cart-pole problem.

- **State:** our current visualization of the cart (i.e. an image).
- **Actions:** 2 of them, go left or go right at a fixed speed.
- **Reward:** +1 for each unit of time where the cart does not leave the screen and the pole does not fall below a certain angle.
- **Next state generation:** cart and pole both follow simple programmed rules of physics.



Toy example #3 (shown, but implementation out of the scope of this class)

- **Problem:** in many RL problems, the states and/or actions sets are not necessarily finite.
- In that case, it is impossible to represent the Q and V functions as tables.
- How do we address this issue?
- Give it to an AI! (as usual)



Toy example #3 (shown, but implementation out of the scope of this class)

- **Solution:** replace the Q -table with a Deep Neural Network, whose job is to estimate the value of each action (left/right) in the current state.
- The objective is then to train, just like before with our Q -table.
- **However, we are no longer changing the table values but the Neural Net parameters!**

```
1 class DQN(nn.Module):
2
3     def __init__(self, h, w, outputs):
4         super(DQN, self).__init__()
5         self.conv1 = nn.Conv2d(3, 16, kernel_size=5, stride=2)
6         self.bn1 = nn.BatchNorm2d(16)
7         self.conv2 = nn.Conv2d(16, 32, kernel_size=5, stride=2)
8         self.bn2 = nn.BatchNorm2d(32)
9         self.conv3 = nn.Conv2d(32, 32, kernel_size=5, stride=2)
10        self.bn3 = nn.BatchNorm2d(32)
11
12        # Number of Linear input connections depends on output of conv2d layers
13        # and therefore the input image size, so compute it.
14        def conv2d_size_out(size, kernel_size = 5, stride = 2):
15            return (size - (kernel_size - 1) - 1) // stride + 1
16        convw = conv2d_size_out(conv2d_size_out(conv2d_size_out(w)))
17        convh = conv2d_size_out(conv2d_size_out(conv2d_size_out(h)))
18        linear_input_size = convw * convh * 32
19        self.head = nn.Linear(linear_input_size, outputs)
20
21        # Called with either one element to determine next action, or a batch
22        # during optimization. Returns tensor([[left0exp,right0exp]...]).
23        def forward(self, x):
24            x = F.relu(self.bn1(self.conv1(x)))
25            x = F.relu(self.bn2(self.conv2(x)))
26            x = F.relu(self.bn3(self.conv3(x)))
27            return self.head(x.view(x.size(0), -1))
```


Toy example #3 (shown, but implementation out of the scope of this class)

- On each round of the game, use the Q network to compute the Q -value of both actions (left/right) in the current state.
- Use the one with the maximal value (**exploitation**) or a randomly chosen action (**exploration**).
- Use **ϵ -greedy policy** to decide how to explore/exploit.

```
1 def select_action(state):
2     global steps_done
3     sample = random.random()
4     eps_threshold = EPS_END + (EPS_START - EPS_END) * \
5         math.exp(-1. * steps_done / EPS_DECAY)
6     steps_done += 1
7     if sample > eps_threshold:
8         with torch.no_grad():
9             # Here, t.max(1) will return largest column value of each row.
10            # Second column on max result is index of where max element was
11            # found, so we pick action with the larger expected reward.
12            return policy_net(state).max(1)[1].view(1, 1)
13     else:
14         return torch.tensor([random.randrange(n_actions)], device=device, dtype=torch.long)
```

Toy example #3

To train this DNN, we need a dataset of some sort.

- Do so by playing the game multiple times and keeping a history of the (state, action, rewards, next_state, done) tuples.
- Here, done indicates that the game has ended (out of screen or low angle on pole).
- Structure is roughly similar to our dataloaders?

```
1 # Define namedtuples for transitions and history
2 Transition = namedtuple('Transition', ('state', 'action', 'next_state', 'reward'))
```

```
1 class ReplayMemory(object):
2
3     def __init__(self, capacity):
4         self.capacity = capacity
5         self.memory = []
6         self.position = 0
7
8     def push(self, *args):
9         """
10         Saves a transition to memory.
11         """
12         if len(self.memory) < self.capacity:
13             self.memory.append(None)
14         self.memory[self.position] = Transition(*args)
15         self.position = (self.position + 1) % self.capacity
16
17     def sample(self, batch_size):
18         """
19         Get sample from history.
20         """
21         return random.sample(self.memory, batch_size)
22
23     def __len__(self):
24         """
25         Get length of history (number of samples).
26         """
27         return len(self.memory)
```

Toy example #3 (shown, but implementation out of the scope of this class)

- **Core idea for memory replay:** we are trying to approximate a complex, nonlinear function Q , with a Neural Network.
- To do this, we must calculate targets using the **Bellman equation** and then consider that we have a **supervised learning** problem at hand.
- **Important:** However, one of the fundamental requirements for SGD optimization is that the training data is independent and identically distributed and when the Agent interacts with the game, the sequence of experience tuples can be highly correlated.
- The naive Q -learning algorithm that learns from each of these experiences tuples in sequential order runs the risk of getting swayed by the effects of this correlation.

Toy example #3 (shown, but implementation out of the scope of this class)

Definition (**experience buffer in RL**):

- We can prevent action values from oscillating or diverging catastrophically using a large buffer of our past experience and sample training data from it, instead of using our latest experience.
- This is called an **experience buffer**.
- The experience buffer contains a collection of experience tuples (state, action, rewards, next_state).
- The tuples are gradually added to the buffer as the agents keep on interacting with the game.

Toy example #3 (shown, but implementation out of the scope of this class)

Definition (**experience replay**):

- The simplest implementation is a buffer of fixed size, with new data added to the end of the experience buffer, so that it pushes the oldest experience out of it.
- The act of sampling a small batch of tuples from the experience buffer in order to learn is known as **experience replay**.
- In addition to breaking harmful correlations, experience replay allows us to learn more from individual tuples multiple times, recall rare occurrences, and in general make better use of our experience.

Toy example #3 (shown, but implementation out of the scope of this class)

To train this DNN, we need a dataset of some sort.

- Do so by playing the game multiple times and keeping a history of the (state, action, rewards, next_state, done) tuples.
- Here, done indicates that the game has ended (out of screen or low angle on pole).
- Structure is roughly similar to our dataloaders?

```
1 # Define namedtuples for transitions and history
2 Transition = namedtuple('Transition', ('state', 'action', 'next_state', 'reward'))
```

```
1 class ReplayMemory(object):
2
3     def __init__(self, capacity):
4         self.capacity = capacity
5         self.memory = []
6         self.position = 0
7
8     def push(self, *args):
9         """
10         Saves a transition to memory.
11         """
12         if len(self.memory) < self.capacity:
13             self.memory.append(None)
14         self.memory[self.position] = Transition(*args)
15         self.position = (self.position + 1) % self.capacity
16
17     def sample(self, batch_size):
18         """
19         Get sample from history.
20         """
21         return random.sample(self.memory, batch_size)
22
23     def __len__(self):
24         """
25         Get length of history (number of samples).
26         """
27         return len(self.memory)
```

Toy example #3 (shown, but implementation out of the scope of this class)

To train this DNN, we need a loss function and weight update procedure of some sort, as well.

Our previous Q -learning was using this iterative update formula.

$$Q^{new}(s_t, a_t) \leftarrow \underbrace{Q(s_t, a_t)}_{\text{old value}} + \underbrace{\alpha}_{\text{learning rate}} \cdot \underbrace{\left(\underbrace{r_t}_{\text{reward}} + \underbrace{\gamma}_{\text{discount factor}} \cdot \underbrace{\max_a Q(s_{t+1}, a)}_{\text{estimate of optimal future value}} - \underbrace{Q(s_t, a_t)}_{\text{old value}} \right)}_{\text{temporal difference}}$$

new value (temporal difference target)

Toy example #3 (shown, but implementation out of the scope of this class)

To train this DNN, we need a loss function and weight update procedure of some sort, as well.

Our previous Q -learning was using this iterative update formula.

$$Q^{new}(s_t, a_t) \leftarrow \underbrace{Q(s_t, a_t)}_{\text{old value}} + \underbrace{\alpha}_{\text{learning rate}} \cdot \underbrace{\left(\underbrace{r_t}_{\text{reward}} + \underbrace{\gamma}_{\text{discount factor}} \cdot \underbrace{\max_a Q(s_{t+1}, a)}_{\text{estimate of optimal future value}} - \underbrace{Q(s_t, a_t)}_{\text{old value}} \right)}_{\text{temporal difference}}$$

new value (temporal difference target)

- **Problem:** update $Q(s_t, a)$ via $Q(s_{t+1}, a)$. However, both states have only one step between them. This makes them very similar, and it is very hard for a Neural Network to distinguish between them.

Toy example #3 (shown, but implementation out of the scope of this class)

Solution: use **two Neural Networks**, one for **training** $Q(s_t, a)$ and one for producing **targets** $Q(s_{t+1}, a)$, or **evaluating** the other.

- That is, the predicted Q values of this second Q -network called the target network, are used to backpropagate through and train the main Q -network.
- **Note:** the target network's parameters are not trained, but they are periodically synchronized with the parameters of the main Q -network.
- The idea is that using the target network's Q values to train the main Q -network will improve the stability of the training.

```

1 def optimize_model():
2     if len(memory) < BATCH_SIZE:
3         return
4     transitions = memory.sample(BATCH_SIZE)
5     # Transpose the batch (see https://stackoverflow.com/a/19343/3343043 for
6     # detailed explanation). This converts batch-array of Transitions
7     # to Transition of batch-arrays.
8     batch = Transition(*zip(*transitions))
9
10    # Compute a mask of non-final states and concatenate the batch elements
11    # (a final state would've been the one after which simulation ended)
12    non_final_mask = torch.tensor(tuple(map(lambda s: s is not None,
13                                           batch.next_state)), device=device, dtype=torch.bool)
14    non_final_next_states = torch.cat([s for s in batch.next_state
15                                      if s is not None])
16    state_batch = torch.cat(batch.state)
17    action_batch = torch.cat(batch.action)
18    reward_batch = torch.cat(batch.reward)
19
20    # Compute Q(s_t, a) - the model computes Q(s_t), then we select the
21    # columns of actions taken. These are the actions which would've been taken
22    # for each batch state according to policy_net
23    state_action_values = policy_net(state_batch).gather(1, action_batch)
24
25    # Compute V(s_{t+1}) for all next states.
26    # Expected values of actions for non_final_next_states are computed based
27    # on the "older" target_net; selecting their best reward with max(1)[0].
28    # This is merged based on the mask, such that we'll have either the expected
29    # state value or 0 in case the state was final.
30    next_state_values = torch.zeros(BATCH_SIZE, device=device)
31    next_state_values[non_final_mask] = target_net(non_final_next_states).max(1)[0].detach()
32    # Compute the expected Q values
33    expected_state_action_values = (next_state_values * GAMMA) + reward_batch
34
35    # Compute Huber loss
36    loss = F.smooth_l1_loss(state_action_values, expected_state_action_values.unsqueeze(1))
37
38    # Optimize the model
39    optimizer.zero_grad()
40    loss.backward()
41    for param in policy_net.parameters():
42        param.grad.data.clamp_(-1, 1)
43    optimizer.step()

```

Toy example #3 (shown, but implementation out of the scope of this class)

To train this DNN, we need a loss function and weight update procedure of some sort.

To create a loss function, let us first recall that

$$Q_t^\pi(s_t, a_t) = R_t(s_t, a_t) + \gamma Q_{t+1}^\pi(s_{t+1}, \pi(s_{t+1}))$$

Let us denote the error δ as

$$\delta = Q_t^\pi(s_t, a_t) - \left(R_t(s_t, a_t) + \gamma \max_a Q_{t+1}^\pi(s_{t+1}, a) \right)$$

Toy example #3 (shown, but implementation out of the scope of this class)

To train this DNN, we need a loss function and weight update procedure of some sort.

To train our DNN, we want to minimize this error δ .

We will use the L1 norm on delta to do so.

$$L(\delta) = |\delta|$$

Toy example #3 (shown, but implementation out of the scope of this class)

To train this DNN, we need a loss function and weight update procedure of some sort.

To train our DNN, we want to minimize this error δ .

We will use the L1 norm on delta to do so.

Note: we can also use a slightly different loss function known as the **Huber loss**, which is slightly more robust to outliers.

$$L_d(\delta) = \begin{cases} \frac{1}{2} \delta^2 & \text{if } |\delta| \leq d \\ d \left(|\delta| - \frac{1}{2} d \right) & \text{else} \end{cases}$$

```

1 def optimize_model():
2     if len(memory) < BATCH_SIZE:
3         return
4     transitions = memory.sample(BATCH_SIZE)
5     # Transpose the batch (see https://stackoverflow.com/a/19343/3343043 for
6     # detailed explanation). This converts batch-array of Transitions
7     # to Transition of batch-arrays.
8     batch = Transition(*zip(*transitions))
9
10    # Compute a mask of non-final states and concatenate the batch elements
11    # (a final state would've been the one after which simulation ended)
12    non_final_mask = torch.tensor(tuple(map(lambda s: s is not None,
13                                           batch.next_state)), device=device, dtype=torch.bool)
14    non_final_next_states = torch.cat([s for s in batch.next_state
15                                      if s is not None])
16    state_batch = torch.cat(batch.state)
17    action_batch = torch.cat(batch.action)
18    reward_batch = torch.cat(batch.reward)
19
20    # Compute Q(s_t, a) - the model computes Q(s_t), then we select the
21    # columns of actions taken. These are the actions which would've been taken
22    # for each batch state according to policy_net
23    state_action_values = policy_net(state_batch).gather(1, action_batch)
24
25    # Compute V(s_{t+1}) for all next states.
26    # Expected values of actions for non_final_next_states are computed based
27    # on the "older" target_net; selecting their best reward with max(1)[0].
28    # This is merged based on the mask, such that we'll have either the expected
29    # state value or 0 in case the state was final.
30    next_state_values = torch.zeros(BATCH_SIZE, device=device)
31    next_state_values[non_final_mask] = target_net(non_final_next_states).max(1)[0].detach()
32    # Compute the expected Q values
33    expected_state_action_values = (next_state_values * GAMMA) + reward_batch
34
35    # Compute Huber loss
36    loss = F.smooth_l1_loss(state_action_values, expected_state_action_values.unsqueeze(1))
37
38    # Optimize the model
39    optimizer.zero_grad()
40    loss.backward()
41    for param in policy_net.parameters():
42        param.grad.data.clamp_(-1, 1)
43    optimizer.step()

```

Trainer function

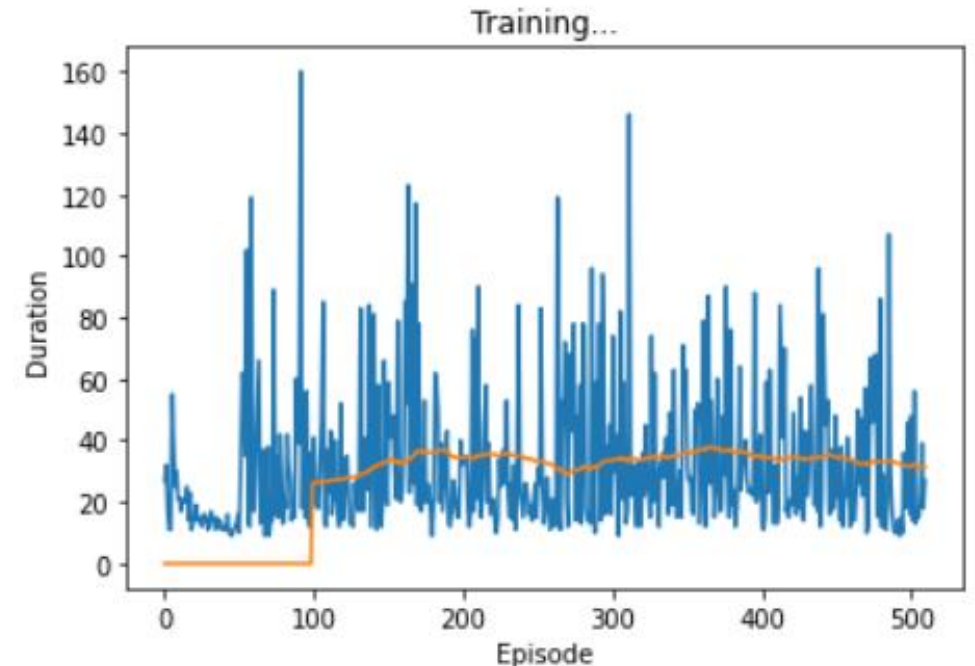
- Our trainer function will play the game 500 times.
- Keep track of different histories over the 500 games.
- Sample from history to train our main Q -Network.
- Backpropagate with mixed main and target Q -networks values.
- Occasionally update the target network.

```
1 """
2 Full trainer on 500 iteration (for meaningful improvements)
3 """
4 num_episodes = 500
5 for i_episode in range(num_episodes):
6     print("Episode:", i_episode)
7     # Initialize the environment and state
8     env.reset()
9     last_screen = get_screen()
10    current_screen = get_screen()
11    state = current_screen - last_screen
12    for t in count():
13        # Select and perform an action
14        action = select_action(state)
15        _, reward, done, _ = env.step(action.item())
16        reward = torch.tensor([reward], device=device)
17
18        # Observe new state
19        last_screen = current_screen
20        current_screen = get_screen()
21        if not done:
22            next_state = current_screen - last_screen
23        else:
24            next_state = None
25
26        # Store the transition in memory
27        memory.push(state, action, next_state, reward)
28
29        # Move to the next state
30        state = next_state
31
32        # Perform one step of the optimization (on the policy network)
33        optimize_model()
34        if done:
35            episode_durations.append(t + 1)
36            plot_durations()
37            break
38
39    # Update the target network, copying all weights and biases in DQN
40    if i_episode % TARGET_UPDATE == 0:
41        target_net.load_state_dict(policy_net.state_dict())
```

Training results

- Our RL agent will learn to balance the pole on the cart, by playing the game.
- Can display the length of each game/episode to see the progression!
- This RL approach of training some DNNs to replace the Q functions is commonly referred to as **Deep Q -learning**.

```
1 def plot_durations():
2     """
3     Show episode durations for each episode.
4     """
5     plt.figure(2)
6     plt.clf()
7     durations_t = torch.tensor(episode_durations, dtype=torch.float)
8     plt.title('Training...')
9     plt.xlabel('Episode')
10    plt.ylabel('Duration')
11    plt.plot(durations_t.numpy())
12    # Take 100 episode averages and plot them too
13    if len(durations_t) >= 100:
14        means = durations_t.unfold(0, 100, 1).mean(1).view(-1)
15        means = torch.cat((torch.zeros(99), means))
16        plt.plot(means.numpy())
17
```



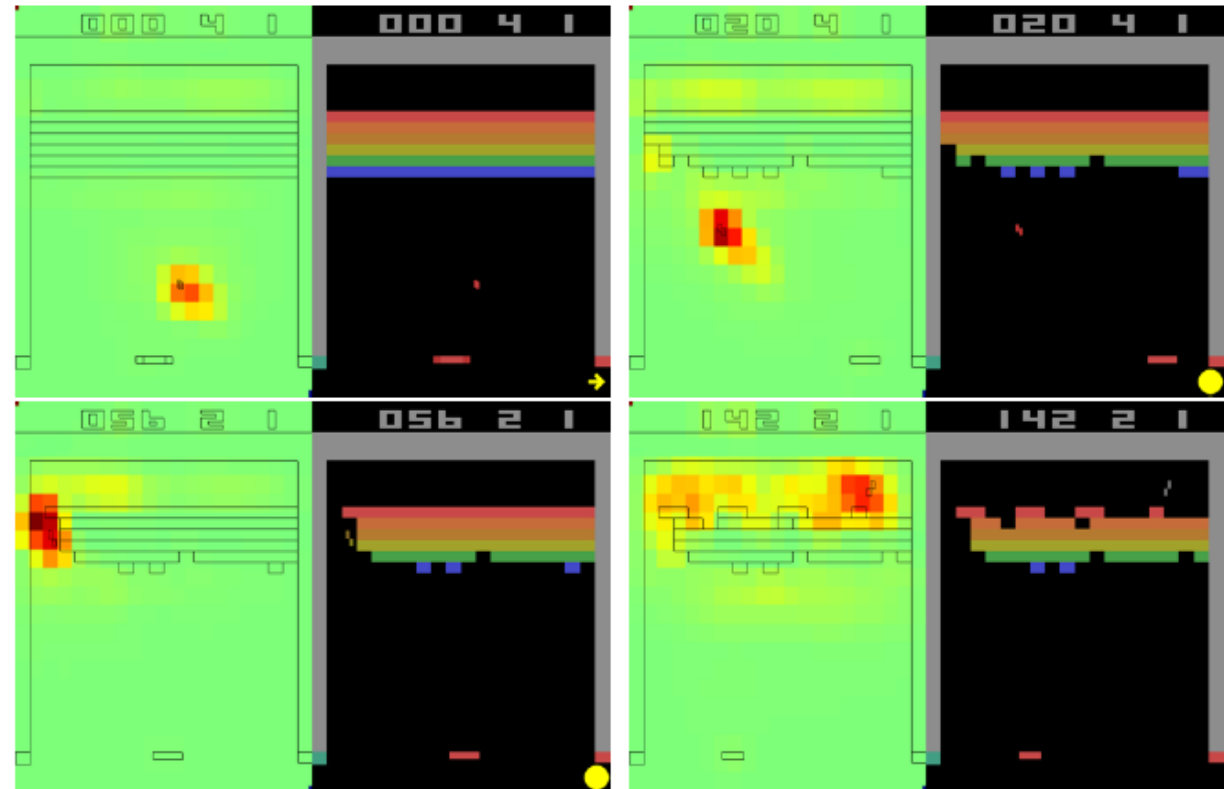
Following this cart-pole balance idea...

- Train an AI to keep a robot on its feet, despite some “*minor environment perturbations*” (a polite way of saying you kick the hell out of the robot for fun).
- Video:
<https://www.youtube.com/watch?v=NR32ULxbjYc>
- **BostonDynamics** blog:
<https://blog.bostondynamics.com/>



Following this idea of using computer vision to identify state and act...

- Train an AI to play video games with Deep Reinforcement Learning (Mnih, 2013)!
- Paper:
<https://www.cs.toronto.edu/~vmnih/docs/dqn.pdf>
- Video:
<https://www.youtube.com/watch?v=TmPfTpjtdgg>



Conclusion

1. What is **Deep Q Learning**? And which problem does it address?
2. What is a **trainer** and a **main Q-network**? What is the interleaved training for Deep Q learning?
3. What are **actor-critic** learning methods? And which problems do these approaches address?
4. What are more advanced problems in RL?
 - Markov states
 - Partially observable environment
 - SARSA
 - Non-stationary problems

Learn more about these topics

Out of class, for those of you who are curious

- [TheBibleOfRL] R. **Sutton** et al., “Reinforcement learning: An Introduction, 2nd edition”, 2018.
<http://www.incompleteideas.net/book/RLbook2020.pdf>
- [Mnih2018] **Mnih** et al., “Playing Atari with Deep Reinforcement Learning”, 2018.
<https://arxiv.org/abs/1312.5602>
- The **BostonDynamics** blog (some more robots updates!)
<https://blog.bostondynamics.com/>

Learn more about these topics

Tracking important names (Track their works and follow them on Scholar, Twitter, or whatever works for you!)

- **Richard Sutton: Professor at University of Alberta, also DeepMind.** Co-author of the Bible of RL (possibly most influential professor in the field of RL).

<https://scholar.google.ca/citations?user=6m4wv6gAAAAJ&hl=en>

<https://www.ualberta.ca/admissions-programs/online-courses/reinforcement-learning/index.html>

<http://www.incompleteideas.net/book/RLbook2020.pdf>

- **Andrew Barto: Professor at University of Massachusetts.** Co-author of the Bible of RL.

<https://people.cs.umass.edu/~barto/>

<https://scholar.google.com/citations?user=CMlgrCgAAAAJ&hl=en>

Learn more about these topics

Tracking important names (Track their works and follow them on Scholar, Twitter, or whatever works for you!)

- **David Silver**: Researcher at DeepMind, Adjunct Professor at University College London (?). Inventor of Alpha GO AIs, has a fantastic course on RL as well.
<https://www.davidsilver.uk/>
<https://scholar.google.com/citations?user=-8DNE4UAAAAJ&hl=en>
<https://www.davidsilver.uk/teaching/>
- **Volodymyr Mnih**: Researcher at DeepMind.
<http://www.cs.toronto.edu/~vmnih>
<https://scholar.google.com/citations?user=rLdfJ1gAAAAJ&hl=en>

Learn more about these topics

Out of class, for those of you who are curious

- [AlphaGo] How AlphaGo was created using DL and RL.
<https://jonathan-hui.medium.com/alphago-zero-a-game-changer-14ef6e45eba5>
- Full Alpha Go movie on DeepMind's Youtube channel. Award-winning movie, pretty cool.
<https://www.youtube.com/watch?v=WXuK6gekU1Y>

Out-of-scope stuff: used to be the third lecture of RL, but no longer is...

In case you are curious about more advanced concepts in RL!

A quick word on actor-critic methods (more advanced stuff, out of scope)

- In Deep Q-learning, we realized that most RL problems cannot have their Q functions computed easily.
- We then replaced the Q function with a Deep Neural Network to approximate this function.

A quick word on actor-critic methods (more advanced stuff, out of scope)

- In Deep Q-learning, we realized that most RL problems cannot have their Q functions computed easily.
- We then replaced the Q function with a Deep Neural Network to approximate this function.
- **Additional suggestion:** Can the reward function always be computed?



A quick word on actor-critic methods (more advanced stuff, out of scope)

- In Deep Q-learning, we realized that most RL problems cannot have their Q functions computed easily.
- We then replaced the Q function with a Deep Neural Network to approximate this function.
- **Additional suggestion:** Can the reward function always be computed?



→ **Replace more elements of the RL system with Deep Neural Networks.**

A quick word on actor-critic methods (more advanced stuff, out of scope)

Definition (actor-critic):

Actor-critic algorithms consist of two components.

- **Actor:** a DNN, whose purpose is to produce actions in response to given states, i.e. a policy. Can be trained as in Deep Q-learning.

A quick word on actor-critic methods (more advanced stuff, out of scope)

Definition (actor-critic):

Actor-critic algorithms consist of two components.

- **Actor:** a DNN, whose purpose is to produce actions in response to given states, i.e. a policy. Can be trained as in Deep Q-learning.
- **Critic:** a DNN, whose purpose is to evaluate the quality of the selected actions and suggesting directions for improvement, by defining a reward function or a Q function.

A quick word on actor-critic methods (more advanced stuff, out of scope)

Definition (actor-critic):

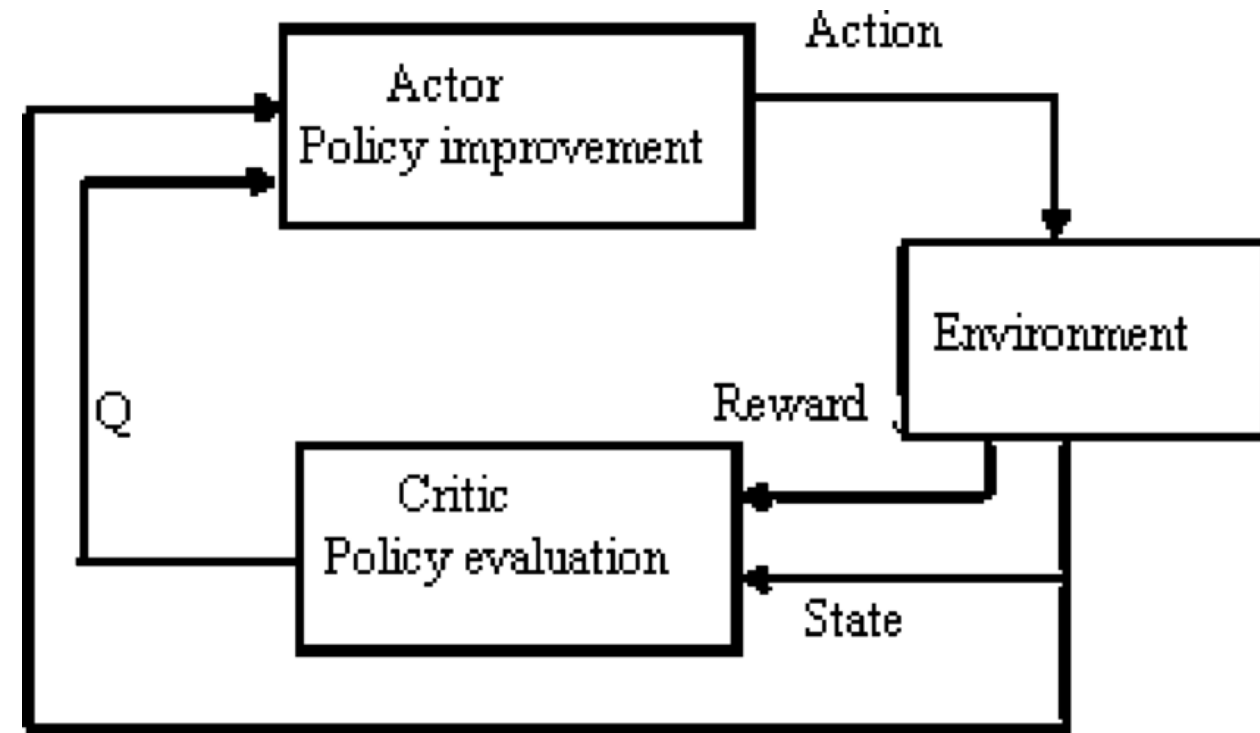
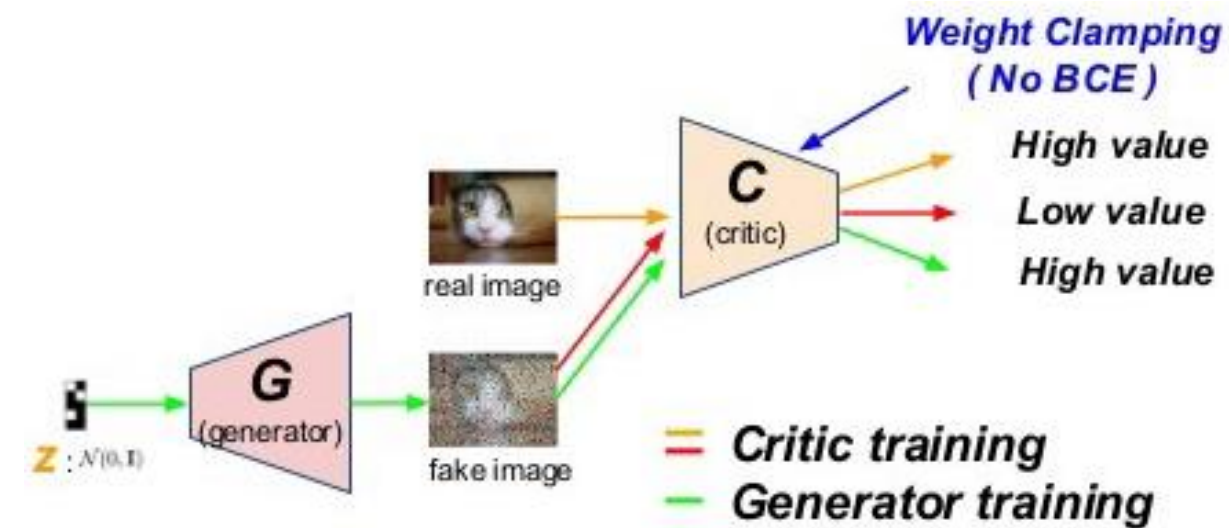
Actor-critic algorithms consist of two components.

- **Actor:** a DNN, whose purpose is to produce actions in response to given states, i.e. a policy. Can be trained as in Deep Q-learning.
- **Critic:** a DNN, whose purpose is to evaluate the quality of the selected actions and suggesting directions for improvement, by defining a reward function or a Q function.

In a sense, **similar to the Generator-Critic pair** of the Wasserstein GANs!

- **Generator:** produce fake images
- **Critic:** evaluate said images

From Deep Q learning to actor-critic



Markov Decision Processes

(more advanced stuff, out of scope)

- Sometimes, the transition from state s_t to next state s_{t+1} , following from action a_t , will not always be **deterministic**.
- In that case, we have to define some system dynamics, as below.

$$p(s', r | s, a) = P(s_{t+1} = s', r_t = r \mid s_t = s, a_t = a)$$

- Similar to our previous problem, with a stochastic twist.
- It is called a Markov Decision Process (MDP) problem.

Markov Decision Processes

(more advanced stuff, out of scope)

- In MDPs, all the previous formulas have to be reworked to account for the stochastic aspect of the problem.

$$V_t^\pi(s) = E[G_t \mid s_t = s]$$

$$Q_t^\pi(s, a) = E[G_t \mid s_t = s, a_t = a]$$

- In MDPs, the Q and V functions can still be learned from experience, but their Bellman equations change slightly, to account for the stochasticity.

Markov Decision Processes

(more advanced stuff, out of scope)

- In MDPs, the Q and V functions can still be learned from experience, but their Bellman equations change slightly, to account for the stochasticity.

$$\begin{aligned} V_t^\pi(s) &= E[G_t \mid s_t = s] \\ V_t^\pi(s) &= E[R_t + \gamma G_{t+1} \mid s_t = s] \\ V_t^\pi(s) &= \sum_a \pi(a \mid s) \sum_{s'} \sum_r p(s', r \mid s, a) [r + \gamma E[G_{t+1} \mid s_{t+1}=s']] \end{aligned}$$

$$V_t^\pi(s) = \sum_a \pi(a \mid s) \sum_{s'} \sum_r p(s', r \mid s, a) [r + \gamma V_{t+1}^\pi(s')]$$

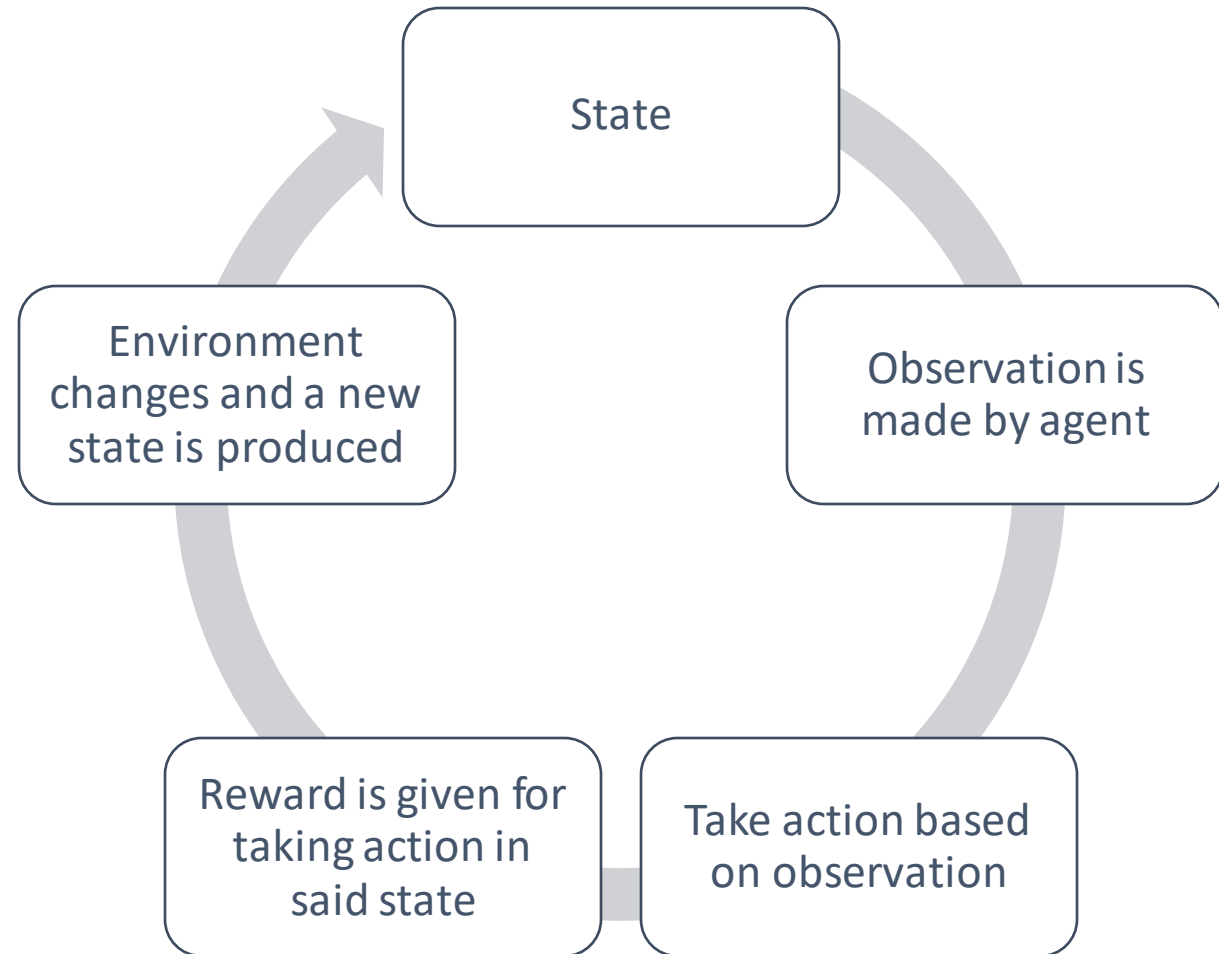
Partially observable MDP (more advanced stuff, out of scope)

Definition (partially observable Markov Decision Process):

In our original RL framework, we assumed that the agent was seeing the exact state of the game at each time t .

This is also an assumption, which can be challenged.

State s_t is ground truth, agent received observation o_t and uses to decide on action.



Partially observable MDP

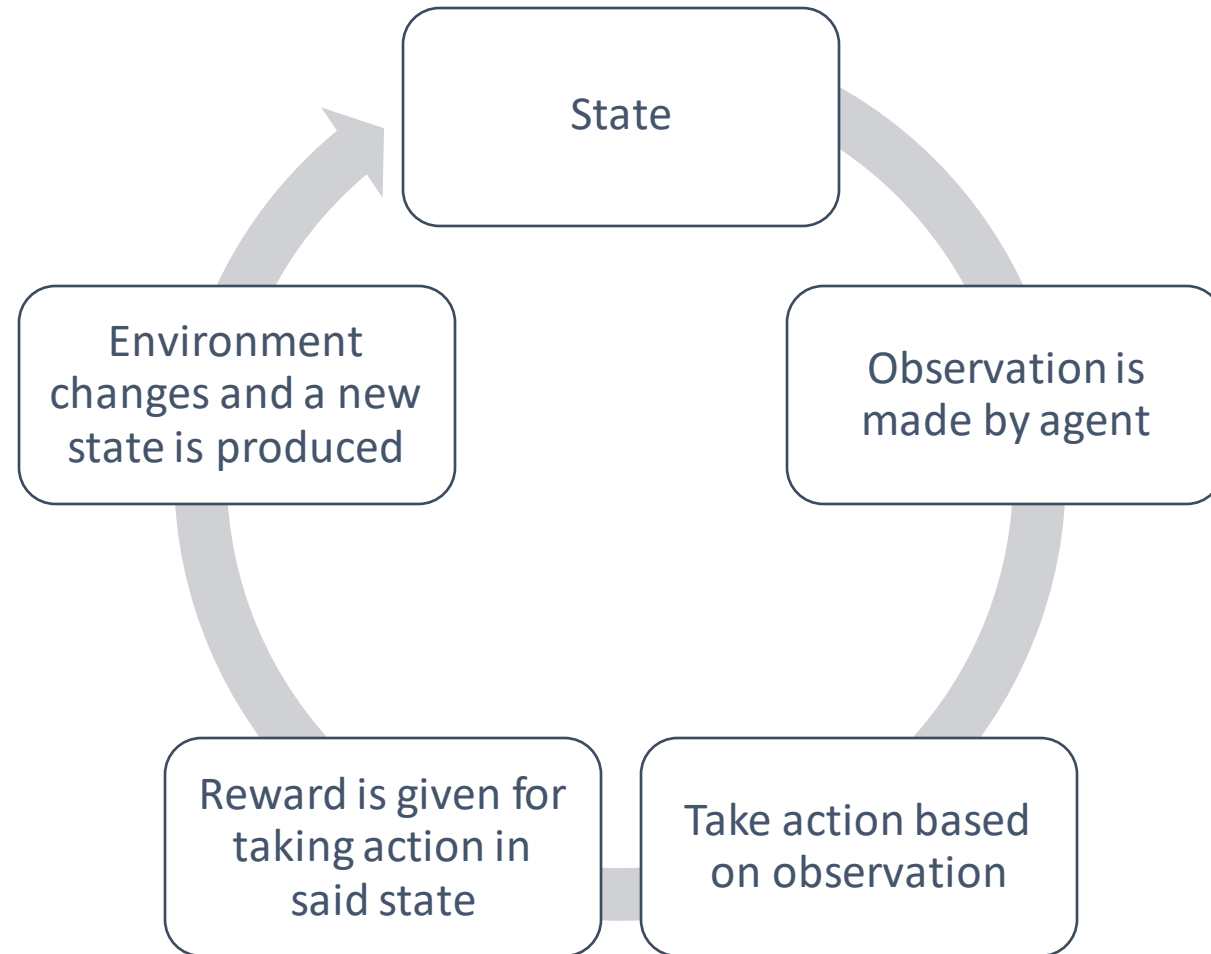
(more advanced stuff, out of scope)

Definition (partially observable Markov Decision Process):

The agent then decides on an observation made of the actual (hidden) state.

This is called a **partially observable Markov Decision Process**.

Agent will have to learn to observe on top of acting properly (e.g. card game, with opponent hiding his/her hand).



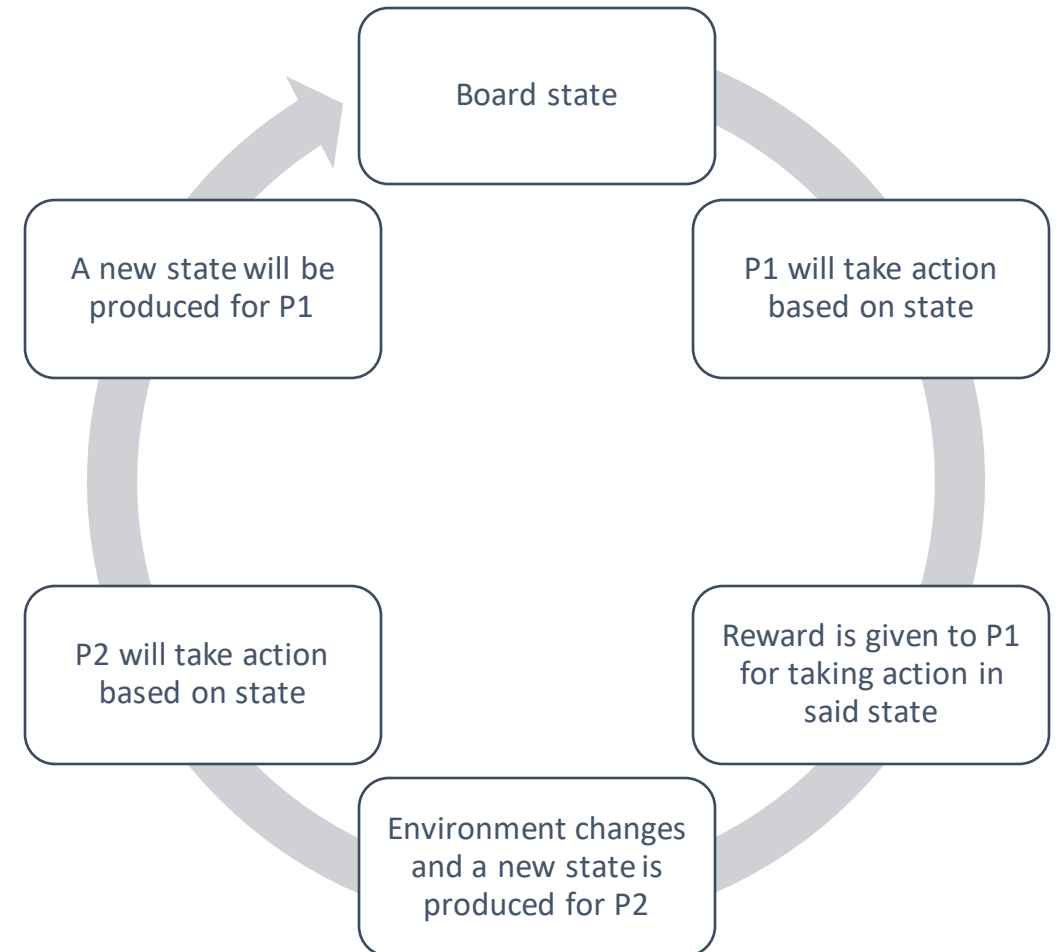
State-Action-Reward-State-Action or SARSA (more advanced stuff, out of scope)

- In many problems, e.g. Go, the new state seen by a given player is not the result of the action of the said player.
- Instead, another player has to act first, before a new state is produced.



State-Action-Reward-State-Action or SARSA (more advanced stuff, out of scope)

- In many problems, e.g. Go, the new state seen by a given player is not the result of the action of the said player.
- Instead, another player has to act first, before a new state is produced.
- This adds steps to the cycle, which becomes (state, action, reward, state_P2, action_P2).



State-Action-Reward-State-Action or SARSA (more advanced stuff, out of scope)

- This is called a **State-Action-Reward-State-Action (SARSA)** type of problem.
- In that case the agent has to learn how to play, but also has to learn how another player might respond to its actions.
- **RL meets game theory!**

