

# **Week 8**

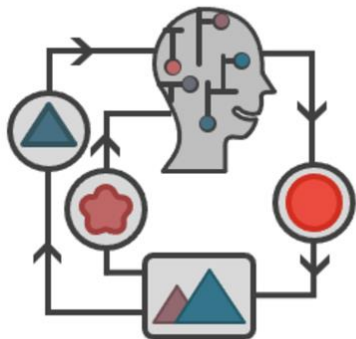
## **Reinforcement Learning**

## Content

- Introduction to Reinforcement Learning
- Fundamentals of Reinforcement Learning
- Markov Decision Processes (MDPs)
- Q-Learning and Deep Q-Networks (DQNs)
- Policy Gradient Methods
- Python: Building a Simple Reinforcement Learning Agent

# Introduction to Reinforcement Learning

- Learning through experience/data to make good decisions under uncertainty
- Essential part of intelligence
- Builds strongly from theory and ideas starting in the 1950s with Richard Bellman
- A number of impressive successes in the last decade



“  
Reinforcement Learning (RL) is an application of Machine Learning where suitable actions are taken to maximize rewards based on rigorous testing, aiming to find the best behavior or path in a specific condition.  
”

# Introduction to Reinforcement Learning

- Beyond Human Performance on the Board Game Go

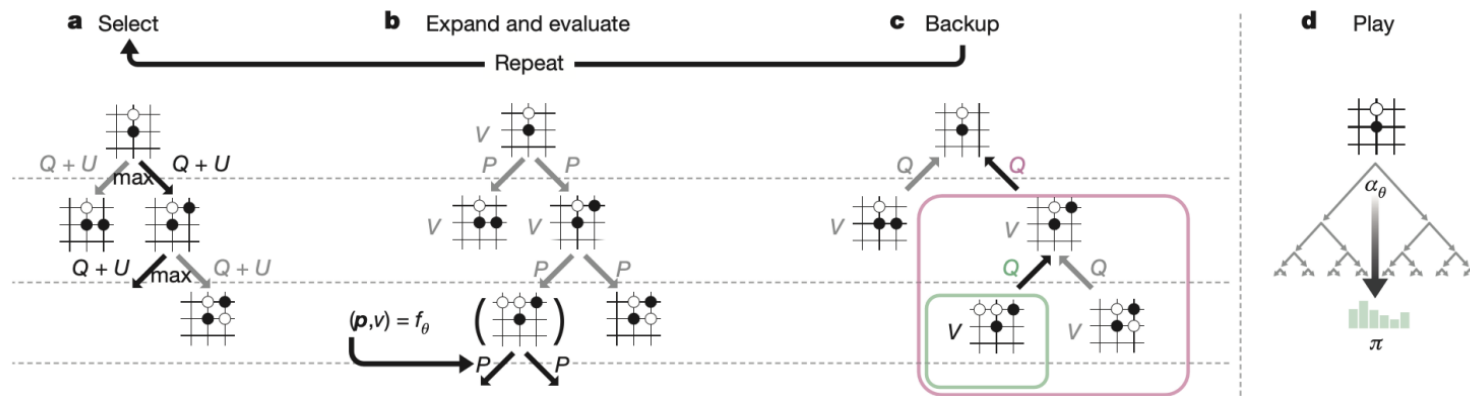
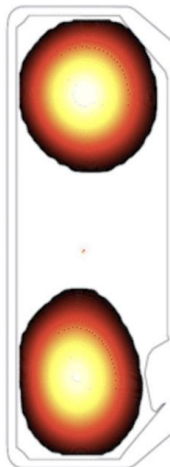


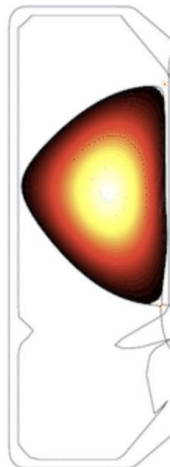
Image credits: Silver et al. Nature 2017  
<https://www.nature.com/articles/nature24270>

# Introduction to Reinforcement Learning

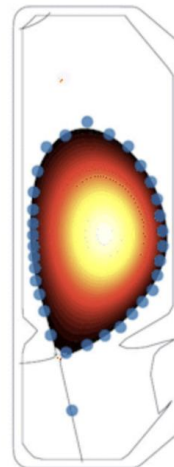
- **Learning Plasma Control for Fusion Science**



Droplets



Negative  
Triangularity

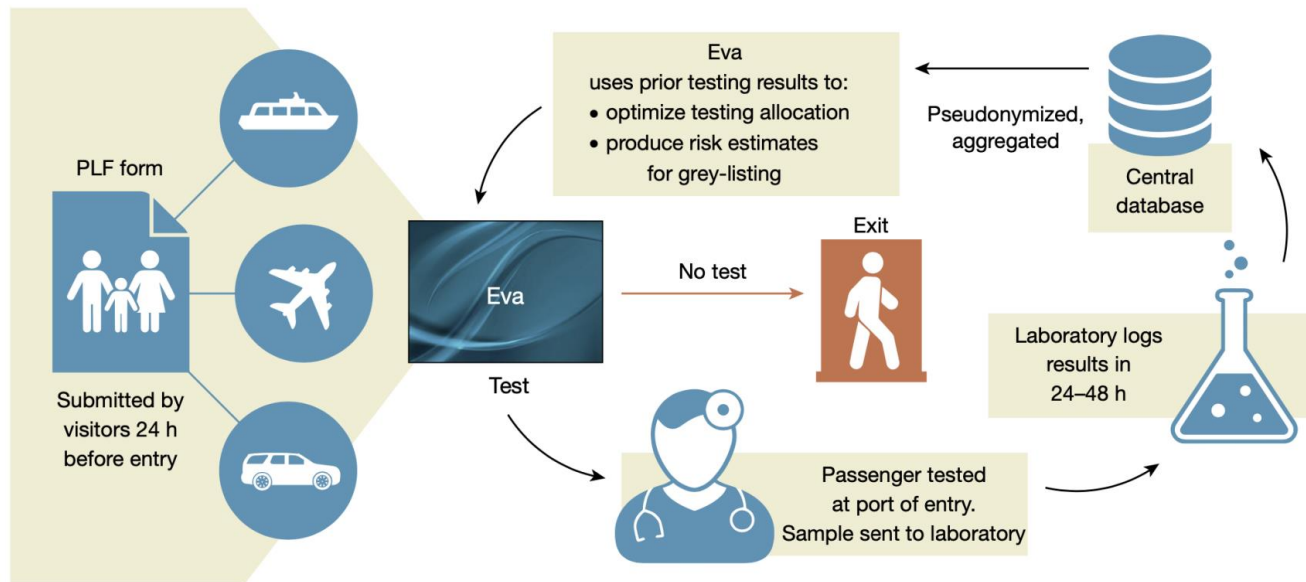


ITER-like  
shape

*Image credits: left Alain Herzog / EPFL, right DeepMind & SPC/EPFL. Degraeve et al.  
Nature 2022 <https://www.nature.com/articles/s41586-021-04301-9>*

# Introduction to Reinforcement Learning

- Efficient and targeted COVID-19 border testing via RL



# Introduction to Reinforcement Learning

- ChatGPT (<https://openai.com/blog/chatgpt/>)

Step 1

Collect demonstration data and train a supervised policy.

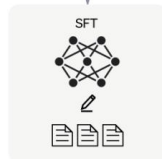
A prompt is sampled from our prompt dataset.



A labeler demonstrates the desired output behavior.



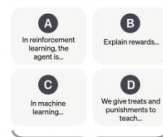
This data is used to fine-tune GPT-3.5 with supervised learning.



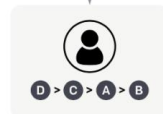
Step 2

Collect comparison data and train a reward model.

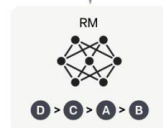
A prompt and several model outputs are sampled.



A labeler ranks the outputs from best to worst.



This data is used to train our reward model.



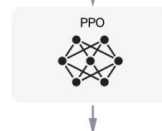
Step 3

Optimize a policy against the reward model using the PPO reinforcement learning algorithm.

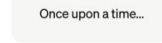
A new prompt is sampled from the dataset.



The PPO model is initialized from the supervised policy.



The policy generates an output.



The reward model calculates a reward for the output.



The reward is used to update the policy using PPO.



# Fundamentals of Reinforcement Learning

- **Basic Concepts**

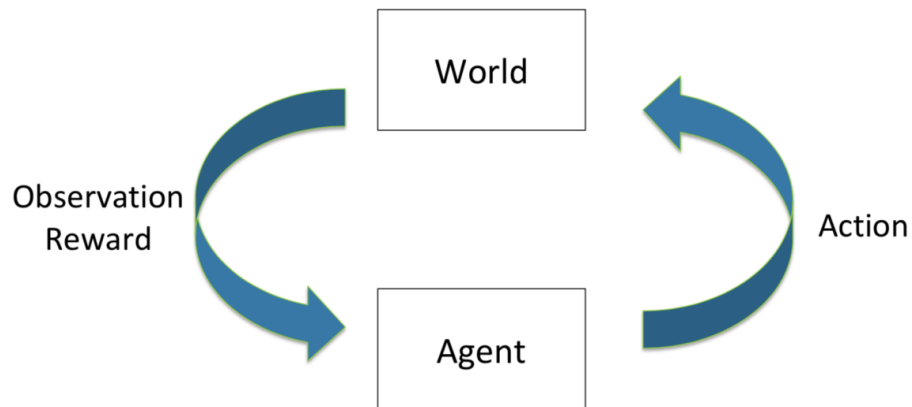
- **Agent:** The entity that learns and makes decisions.
- **Environment:** The system within which the agent operates.
- **State (S):** A representation of the environment at a specific point in time.
- **Action (A):** Choices available to the agent to interact with the environment.
- **Reward (R):** A scalar feedback signal received from the environment after taking an action.
- **Policy ( $\pi$ ):** The strategy the agent uses to decide which action to take in a given state.
- $\pi(S, A)$



# Fundamentals of Reinforcement Learning

- **Sequential Decision Making**

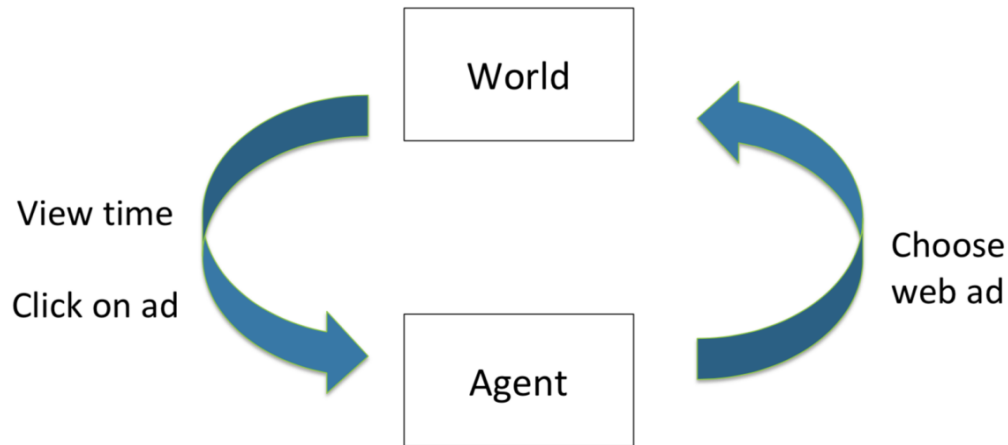
- **Goal:** Select actions to maximize total expected future reward.
- May require balancing immediate & long term rewards



# Fundamentals of Reinforcement Learning

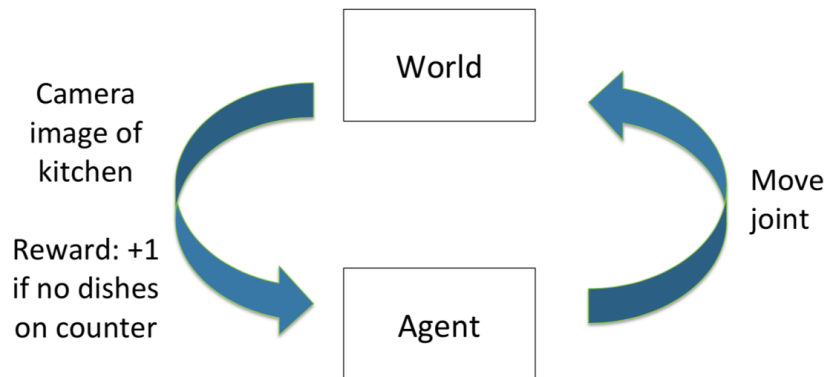
- **Example: Web Advertising**

- **Goal:** Select actions to maximize total expected future reward.
- May require balancing immediate & long term rewards



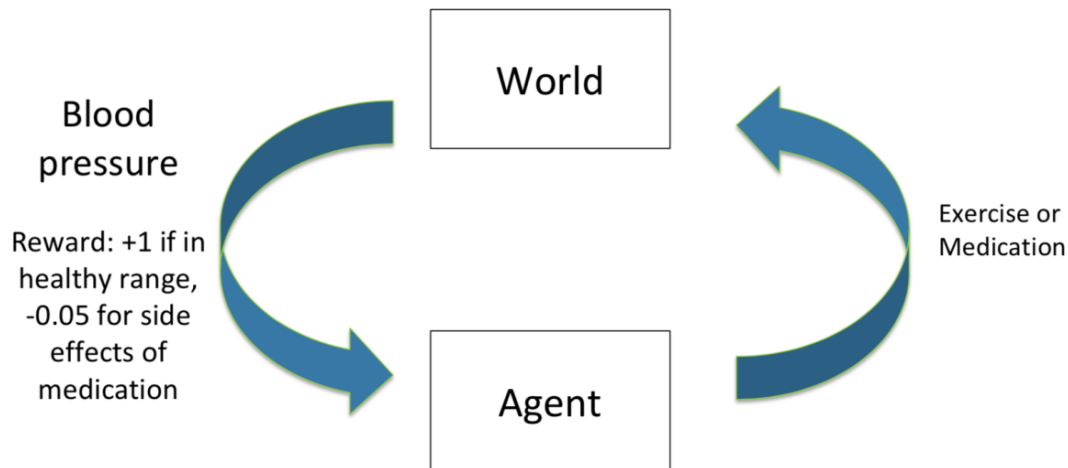
# Fundamentals of Reinforcement Learning

- **Example: Robot Unloading Dishwasher**
  - **Goal:** Select actions to maximize total expected future reward.
  - May require balancing immediate & long term rewards



# Fundamentals of Reinforcement Learning

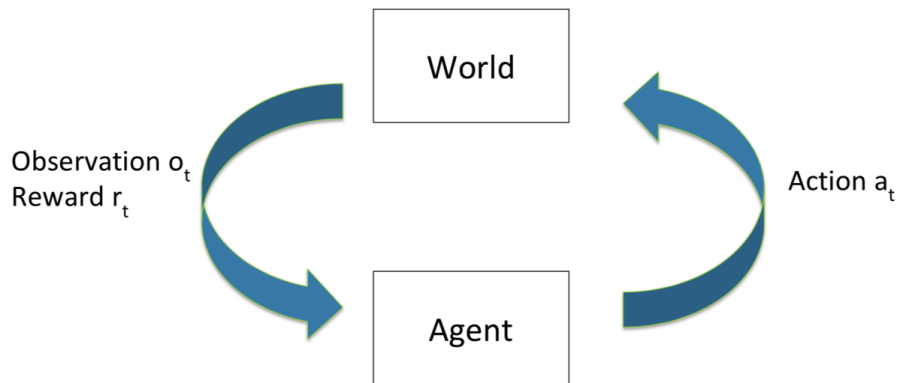
- **Example: Blood Pressure Control**
  - **Goal:** Select actions to maximize total expected future reward.
  - May require balancing immediate & long term rewards



# Fundamentals of Reinforcement Learning

- **Sequential Decision Process: Agent & the World (Discrete Time)**

- Each time step  $t$ :
  - Agent takes an action  $a_t$
  - World updates given action  $a_t$ , emits observation  $o_t$  and reward  $r_t$
  - Agent receives observation  $o_t$  and reward  $r_t$



# Markov Decision Processes (MDPs)

- **Markov Assumption**

- Information state: sufficient statistic of history
- State  $s_t$  is Markov if and only if:

$$p(s_{t+1}|s_t, a_t) = p(s_{t+1}|h_t, a_t)$$

- Future is independent of past given present

# Markov Decision Processes (MDPs)

- **Why is Markov Assumption Popular?**

- Simple and often can be satisfied if include some history as part of the state
- In practice often assume most recent observation is sufficient statistic of history

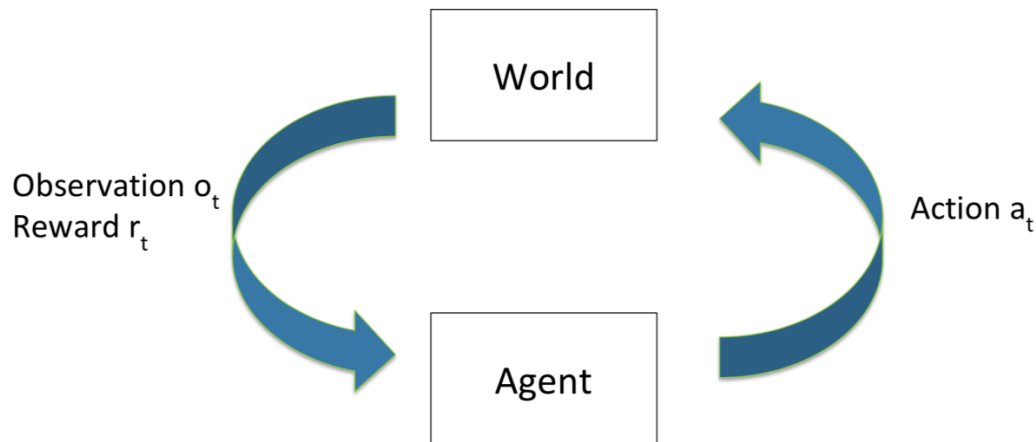
$$s_t = o_t$$

- State representation has big implications for:
  - Computational complexity
  - Data required
  - Resulting performance

# Markov Decision Processes (MDPs)

- **Why is Markov Assumption Popular?**

- Is state Markov? Is world partially observable? (POMDP)
- Are dynamics deterministic or stochastic?
- Do actions influence only immediate reward (bandits) or reward and next state ?





# Markov Decision Processes (MDPs)

- **Example: Mars Rover as a Markov Decision Process**

- States: Location of rover ( $s_1, \dots, s_7$ )
- Actions: TryLeft or TryRight
- Rewards:
  - +1 in state  $s_1$
  - +10 in state  $s_7$
  - 0 in all other states


$s_1$	$s_2$	$s_3$	$s_4$	$s_5$	$s_6$	$s_7$
						

Figure: Mars rover image: NASA/JPL-Caltech

# Markov Decision Processes (MDPs)

- **MDP Model**

- Agent's representation of how world changes given agent's action
- Transition / dynamics model predicts next agent state

$$p(s_{t+1} = s' | s_t = s, a_t = a)$$

- Reward model predicts immediate reward

$$r(s_t = s, a_t = a) = \mathbb{E}[r_t | s_t = s, a_t = a]$$

# Markov Decision Processes (MDPs)

- **Example: Mars Rover Stochastic Markov Model**

- Numbers above show RL agent's reward model
- Part of agent's transition model:

$$\begin{aligned} 0.5 &= P(s_1 | s_1, \text{TryRight}) = P(s_2 | s_1, \text{TryRight}) \\ 0.5 &= P(s_2 | s_2, \text{TryRight}) = P(s_3 | s_2, \text{TryRight}) \dots \end{aligned}$$

- Model may be wrong

[illegible]

# Markov Decision Processes (MDPs)

- **Policy**

- Policy  $\pi$  determines how the agent chooses actions
- $\pi : S \rightarrow A$ , mapping from states to actions
- Deterministic policy:

$$\pi(s) = a$$

- Stochastic policy:

$$\pi(a|s) = Pr(a_t = a | s_t = s)$$

# Markov Decision Processes (MDPs)

- **Example: Mars Rover Policy**

- $\pi(s_1) = \pi(s_2) = \dots = \pi(s_7) = \text{TryRight}$
- Q: is this a deterministic policy or a stochastic policy?

# Markov Decision Processes (MDPs)

- **Evaluation and Control**
  - Evaluation
    - Estimate/predict the expected rewards from following a given policy
  - Control
    - Optimization: find the best policy

# Markov Decision Processes (MDPs)

- **Making Sequences of Good Decisions Given a Model of the World**
  - Assume finite set of states and actions
  - Given models of the world (dynamics and reward)
  - Evaluate the performance of a particular decision policy
  - Compute the best policy
  - This can be viewed as an AI planning problem

# Markov Decision Processes (MDPs)

- **Making Sequences of Good Decisions Given a Model of the World**
  - Markov Processes Markov
  - Reward Processes (MRPs)
  - Markov Decision Processes (MDPs)
  - Evaluation and Control in MDPs



# Q-Learning and Deep Q-Networks (DQNs)

- Reinforcement Learning (Generally)  
Involves
- Optimization
- Delayed consequences
- Exploration
- Generalization

# Q-Learning and Deep Q-Networks (DQNs)

- **Reinforcement Learning (Generally) Involves**
  - Optimization
    - Goal is to find an optimal way to make decisions
    - Explicit notion of decision utility
  - Delayed consequences
    - Decisions now can impact things much later...
    - Saving for retirement Finding a key in video game Montezuma's revenge
  - Exploration
    - Learning about the world by making decisions
    - Decisions impact what we learn about
  - Generalization
    - Policy is mapping from past experience to action

# Q-Learning and Deep Q-Networks (DQNs)

- **Generalized Policy Improvement**

- Model-free Policy Iteration

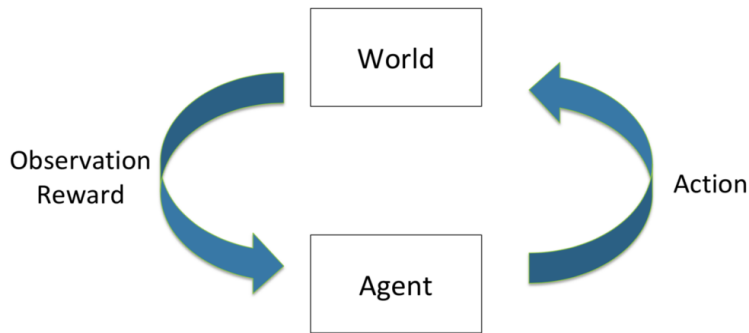
- Initialize policy  $\pi$
    - Repeat:
      - Policy evaluation: compute  $Q\pi$
      - Policy improvement: update  $\pi$  given  $Q\pi$
    - May need to modify policy evaluation:
      - If  $\pi$  is deterministic, can't compute  $Q(s, a)$  for any  $a \neq \pi(s)$
    - How to interleave policy evaluation and improvement?
      - Policy improvement is now using an estimated  $Q$

# Q-Learning and Deep Q-Networks (DQNs)

- **Generalized Policy Improvement**

- The Problem of Exploration

- Goal: Learn to select actions to maximize total expected future reward
    - Problem: Can't learn about actions without trying them (need to explore)
    - Problem: But if we try new actions, spending less time taking actions that our past experience suggests will yield high reward (need to exploit knowledge of domain to achieve high rewards)



# Q-Learning and Deep Q-Networks (DQNs)

- **Generalized Policy Improvement**

- $\epsilon$ -greedy Policies

- Simple idea to balance exploration and achieving rewards
    - Let  $|A|$  be the number of actions
    - Then an  $\epsilon$ -greedy policy w.r.t. a state-action value  $Q(s, a)$  is  $\pi(a|s) =$

$$\begin{aligned} & \arg \max_a Q(s, a), \text{ w. prob } 1 - \epsilon + \frac{\epsilon}{|A|} \\ & a' \neq \arg \max Q(s, a) \text{ w. prob } \frac{\epsilon}{|A|} \end{aligned}$$

- In words: select argmax action with probability  $1 - \epsilon$ , else select action uniformly at random

# Q-Learning and Deep Q-Networks (DQNs)

- **Generalized Policy Improvement**

- Policy Improvement with  $\epsilon$ -greedy policies

- Recall we proved that policy iteration using given dynamics and reward models, was guaranteed to monotonically improve
    - That proof assumed policy improvement output a deterministic policy
    - Same property holds for  $\epsilon$ -greedy policies
    - Monotonic  $\epsilon$ -greedy Policy Improvement
      - Theorem

For any  $\epsilon$ -greedy policy  $\pi_i$ , the  $\epsilon$ -greedy policy w.r.t.  $Q^{\pi_i}$ ,  $\pi_{i+1}$  is a monotonic improvement  $V^{\pi_{i+1}} \geq V^{\pi_i}$

$$\begin{aligned} Q^{\pi_i}(s, \pi_{i+1}(s)) &= \sum_{a \in A} \pi_{i+1}(a|s) Q^{\pi_i}(s, a) \\ &= (\epsilon/|A|) \left[ \sum_{a \in A} Q^{\pi_i}(s, a) \right] + (1 - \epsilon) \max_a Q^{\pi_i}(s, a) \end{aligned}$$

# Q-Learning and Deep Q-Networks (DQNs)

- Recall Monte Carlo Policy Evaluation, Now for Q

---

```
1: Initialize  $Q(s, a) = 0, N(s, a) = 0 \forall (s, a), k = 1$ , Input  $\epsilon = 1, \pi$ 
2: loop
3:   Sample  $k$ -th episode  $(s_{k,1}, a_{k,1}, r_{k,1}, s_{k,2}, \dots, s_{k,T})$  given  $\pi$ 
3:   Compute  $G_{k,t} = r_{k,t} + \gamma r_{k,t+1} + \gamma^2 r_{k,t+2} + \dots \gamma^{T_i-1} r_{k,T_i} \forall t$ 
4:   for  $t = 1, \dots, T$  do
5:     if First visit to  $(s,a)$  in episode  $k$  then
6:        $N(s, a) = N(s, a) + 1$ 
7:        $Q(s_t, a_t) = Q(s_t, a_t) + \frac{1}{N(s,a)} (G_{k,t} - Q(s_t, a_t))$ 
8:     end if
9:   end for
10:   $k = k + 1$ 
11: end loop
```

---

# Q-Learning and Deep Q-Networks (DQNs)

- Monte Carlo Online Control / On Policy Improvement

---

```
1: Initialize  $Q(s, a) = 0, N(s, a) = 0 \forall (s, a)$ , Set  $\epsilon = 1, k = 1$ 
2:  $\pi_k = \epsilon$ -greedy( $Q$ ) // Create initial  $\epsilon$ -greedy policy
3: loop
4:   Sample  $k$ -th episode  $(s_{k,1}, a_{k,1}, r_{k,1}, s_{k,2}, \dots, s_{k,T})$  given  $\pi_k$ 
4:    $G_{k,t} = r_{k,t} + \gamma r_{k,t+1} + \gamma^2 r_{k,t+2} + \dots + \gamma^{T-t} r_{k,T}$ 
5:   for  $t = 1, \dots, T$  do
6:     if First visit to  $(s, a)$  in episode  $k$  then
7:        $N(s, a) = N(s, a) + 1$ 
8:        $Q(s_t, a_t) = Q(s_t, a_t) + \frac{1}{N(s,a)}(G_{k,t} - Q(s_t, a_t))$ 
9:     end if
10:  end for
11:   $k = k + 1, \epsilon = 1/k$ 
12:   $\pi_k = \epsilon$ -greedy( $Q$ ) // Policy improvement
13: end loop
```

---



# Q-Learning and Deep Q-Networks (DQNs)

- **Greedy in the Limit of Infinite Exploration (GLIE)**

- Definition of GLIE

- All state-action pairs are visited an infinite number of times

$$\lim_{i \rightarrow \infty} N_i(s, a) \rightarrow \infty$$

- Behavior policy (policy used to act in the world) converges to greedy policy

- GLIE Monte-Carlo Control using Tabular Representations

- GLIE Monte-Carlo control converges to the optimal state-action value function

$$Q(s, a) \rightarrow Q^*(s, a)$$

# Q-Learning and Deep Q-Networks (DQNs)

- **Model-free Policy Iteration with TD Methods**

- Initialize policy  $\pi$
- Repeat:
  - Policy evaluation: compute  $Q\pi$  using temporal difference updating with  $\epsilon$ -greedy policy
  - Policy improvement: Same as Monte carlo policy improvement, set  $\pi$  to  $\epsilon$ -greedy ( $Q\pi$ )
- On policy: SARSA computes an estimate  $Q$  of policy used to act

# Q-Learning and Deep Q-Networks (DQNs)

- General Form of SARSA(**What is SARSA**) Algorithm

---

- 1: Set initial  $\epsilon$ -greedy policy  $\pi$  randomly,  $t = 0$ , initial state  $s_t = s_0$
- 2: Take  $a_t \sim \pi(s_t)$
- 3: Observe  $(r_t, s_{t+1})$
- 4: **loop**
- 5:   Take action  $a_{t+1} \sim \pi(s_{t+1})$  // Sample action from policy
- 6:   Observe  $(r_{t+1}, s_{t+2})$
- 7:   Update Q given  $(s_t, a_t, r_t, s_{t+1}, a_{t+1})$ :
- 8:   Perform policy improvement:
- 9:    $t = t + 1$
- 10: **end loop**

---

# Q-Learning and Deep Q-Networks (DQNs)

- **Q-Learning: Learning the Optimal State-Action Value**
  - SARSA is an on-policy learning algorithm
  - SARSA estimates the value of the current behavior policy (policy using to take actions in the world)
  - And then updates that (behavior) policy
  - Alternatively, can we directly estimate the value of  $\pi^*$  while acting with another behavior policy  $\pi_b$ ?
  - Yes! Q-learning, an off-policy RL algorithm

# Q-Learning and Deep Q-Networks (DQNs)

- **Q-Learning: Learning the Optimal State-Action Value**

- SARSA is an on-policy learning algorithm
  - Estimates the value of behavior policy (policy using to take actions in the world)
  - And then updates the behavior policy Q-learning estimate the Q value of  $\pi$  \* while acting with another behavior policy  $\pi_b$
- Key idea: Maintain Q estimates and bootstrap for best future value
- Recall SARSA

$$Q(s_t, a_t) \leftarrow Q(s_t, a_t) + \alpha((r_t + \gamma Q(s_{t+1}, a_{t+1})) - Q(s_t, a_t))$$

- Q-learning:

$$Q(s_t, a_t) \leftarrow Q(s_t, a_t) + \alpha((r_t + \gamma \max_{a'} Q(s_{t+1}, a')) - Q(s_t, a_t))$$

# Q-Learning and Deep Q-Networks (DQNs)

- Q-Learning with  $\epsilon$ -greedy Exploration

---

```
1: Initialize  $Q(s, a), \forall s \in S, a \in A$   $t = 0$ , initial state  $s_t = s_0$ 
2: Set  $\pi_b$  to be  $\epsilon$ -greedy w.r.t.  $Q$ 
3: loop
4:   Take  $a_t \sim \pi_b(s_t)$  // Sample action from policy
5:   Observe  $(r_t, s_{t+1})$ 
6:    $Q(s_t, a_t) \leftarrow Q(s_t, a_t) + \alpha(r_t + \gamma \max_a Q(s_{t+1}, a) - Q(s_t, a_t))$ 
7:    $\pi(s_t) = \arg \max_a Q(s_t, a)$  w.prob  $1 - \epsilon$ , else random
8:    $t = t + 1$ 
9: end loop
```

---

# Q-Learning and Deep Q-Networks (DQNs)

- **Q-Learning with Neural Networks**
  - Q-learning converges to optimal  $Q^*(s, a)$  using tabular representation
  - In value function approximation Q-learning minimizes MSE loss by stochastic gradient descent using a target Q estimate instead of true Q
  - But Q-learning with VFA can diverge
  - Two of the issues causing problems:
    - Correlations between samples
    - Non-stationary targets
  - Deep Q-learning (DQN) addresses these challenges by using
    - Experience replay
    - Fixed Q-targets

# Q-Learning and Deep Q-Networks (DQNs)

- **DQNs: Experience Replay**

- To help remove correlations, store dataset (called a replay buffer)  $\mathcal{D}$  from prior experience
- To perform experience replay, repeat the following:

$(s, a, r, s') \sim \mathcal{D}$ : sample an experience tuple from the dataset

Compute the target value for the sampled  $s$ :  $r + \gamma \max_{a'} \hat{Q}(s', a'; \mathbf{w})$

Use stochastic gradient descent to update the network weights

$$\Delta \mathbf{w} = \alpha (r + \gamma \max_{a'} \hat{Q}(s', a'; \mathbf{w}) - \hat{Q}(s, a; \mathbf{w})) \nabla_{\mathbf{w}} \hat{Q}(s, a; \mathbf{w})$$

$s_1, a_1, r_2, s_2$
$s_2, a_2, r_3, s_3$
$s_3, a_3, r_4, s_4$
$\dots$
$s_t, a_t, r_{t+1}, s_{t+1}$

$\rightarrow s, a, r, s'$

- Uses target as a scalar, but function weights will get updated on the next round, changing the target value



# Q-Learning and Deep Q-Networks (DQNs)

- DQN Pseudocode

---

```
1: Input  $C, \alpha, D = \{\}$ , Initialize  $\mathbf{w}, \mathbf{w}^- = \mathbf{w}, t = 0$ 
2: Get initial state  $s_0$ 
3: loop
4:   Sample action  $a_t$  given  $\epsilon$ -greedy policy for current  $\hat{Q}(s_t, a; \mathbf{w})$ 
5:   Observe reward  $r_t$  and next state  $s_{t+1}$ 
6:   Store transition  $(s_t, a_t, r_t, s_{t+1})$  in replay buffer  $D$ 
7:   Sample random minibatch of tuples  $(s_i, a_i, r_i, s_{i+1})$  from  $D$ 
8:   for  $j$  in minibatch do
9:     if episode terminated at step  $i + 1$  then
10:       $y_i = r_i$ 
11:     else
12:       $y_i = r_i + \gamma \max_{a'} \hat{Q}(s_{i+1}, a'; \mathbf{w}^-)$ 
13:     end if
14:     Do gradient descent step on  $(y_i - \hat{Q}(s_i, a_i; \mathbf{w}))^2$  for parameters  $\mathbf{w}$ :  $\Delta \mathbf{w} = \alpha(y_i - \hat{Q}(s_i, a_i; \mathbf{w})) \nabla_{\mathbf{w}} \hat{Q}(s_i, a_i; \mathbf{w})$ 
15:   end for
16:    $t = t + 1$ 
17:   if  $\text{mod}(t, C) == 0$  then
18:      $\mathbf{w}^- \leftarrow \mathbf{w}$ 
19:   end if
20: end loop
```

---

Note there are several hyperparameters and algorithm choices. One needs to choose the neural network architecture, the learning rate, and how often to update the target network. Often a fixed size replay buffer is used for experience replay, which introduces a parameter to control the size, and the need to decide how to populate it.

# Q-Learning and Deep Q-Networks (DQNs)

- **DQNs Summary (Other RF model)**
  - DQN uses experience replay and fixed Q-targets
  - Store transition ( $s_t, a_t, r_{t+1}, s_{t+1}$ ) in replay memory  $D$
  - Sample random mini-batch of transitions ( $s, a, r, s'$ ) from  $D$
  - Compute Q-learning targets w.r.t. old, fixed parameters  $w^-$
  - Optimizes MSE between Q-network and Q-learning targets
  - Uses stochastic gradient descent

# Q-Learning and Deep Q-Networks (DQNs)

- Which Aspects of DQN were Important for Success?

Game	Linear	Deep Network	DQN w/ fixed Q	DQN w/ replay	DQN w/replay and fixed Q
Breakout	3	3	10	241	317
Enduro	62	29	141	831	1006
River Raid	2345	1453	2868	4102	7447
Seaquest	656	275	1003	823	2894
Space Invaders	301	302	373	826	1089

Replay is hugely important

Why? Beyond helping with correlation between samples, what does replaying do?

# Python: Building a Simple Reinforcement Learning Agent

```
class QLearningAgent:
    def __init__(self, env, learning_rate=0.1, discount_factor=0.99, exploration_rate=1.0, exploration_decay=0.995):
        self.env = env
        self.q_table = np.zeros((env.size, env.size, 4)) # Q-table initialized to zeros
        self.learning_rate = learning_rate
        self.discount_factor = discount_factor
        self.exploration_rate = exploration_rate
        self.exploration_decay = exploration_decay

    def choose_action(self, state):
        if random.uniform(0, 1) < self.exploration_rate:
            return random.choice([0, 1, 2, 3]) # Explore
        else:
            x, y = state
            return np.argmax(self.q_table[x, y]) # Exploit

    def learn(self, state, action, reward, next_state):
        x, y = state
        next_x, next_y = next_state
        best_next_action = np.argmax(self.q_table[next_x, next_y])
        td_target = reward + self.discount_factor * self.q_table[next_x, next_y, best_next_action]
        td_error = td_target - self.q_table[x, y, action]
        self.q_table[x, y, action] += self.learning_rate * td_error
```