k                    Learning

# Outline

**1** Reinforcement Learning

**2** Markov Decision Process

**3** Learning Algorithms

**4** Q-Learning in R

**5** Wrap-Up

# Outline

1 **Reinforcement Learning**

2 Markov Decision Process

3 Learning Algorithms

4 Q-Learning in R

5 Wrap-Up
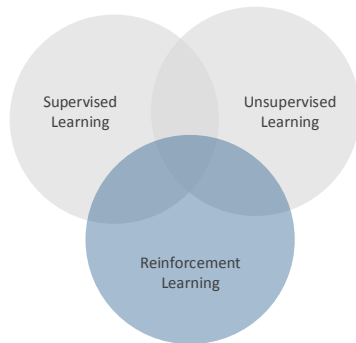
# Branches of Machine Learning

**Supervised Learning**

► Learns from pairs of input and desired outcome (i. e. labels)

**Unsupervised Learning**

► Tries to find hidden structure in unlabeled data



Supervised Learning

Unsupervised Learning

Reinforcement Learning

**Reinforcement Learning**

► Learning from interacting with the environment
► No need for pairs of input and correct outcome
► Feedback restricted to a reward signal
► Mimics human-like learning in actual environments

# Example: Backgammon

Reinforcement learning can reach a level similar to the top three human players in backgammon

## Learning task

- ► Select best move at arbitrary board states
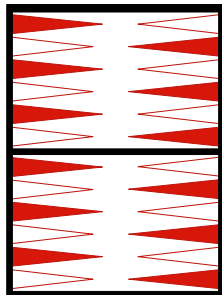  - → i. e. with highest probability to win

## Training signal

- ► Win or loss of overall game

## Training

- ► 300,000 games played against the system itself
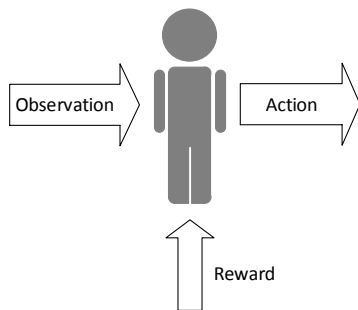
## Algorithm

- ► Reinforcement learning (plus neural network)



→ Tesauro (1995): Temporal Difference Learning and TD-Gammon. In: Comm. of the ACM, 38:3, pp. 58–68
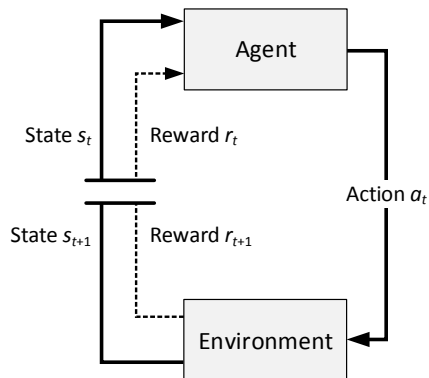
# Reinforcement Learning

- An agent interacts with its environment
- Agent takes actions that affect the state of the environment
- Feedback is limited to a reward signal that indicates how well the agent is performing
- Goal: improve the behavior given only this limited feedback



**Examples**

- Defeat the world champions at backgammon or Go
- Manage an investment portfolio
- Make a humanoid robot walk

# Agent and Environment



At each step $t$, the agent:

- Executes action $a_t$
- Receives observation $s_t$
- Receives scalar reward $r_t$

The environment:

- Changes upon action $a_t$
- Emits observation $s_{t+1}$
- Emits scalar reward $r_{t+1}$

- Time step $t$ is incremented after each iteration

# Agent and Environment

**Example**

$(1)$ ENVIRONMENT    ► You are in state 3 with 4 possible actions

$(2)$ AGENT             ► I'll take action 2

$(3)$ ENVIRONMENT    ► You received a reward of 5 units

                             ► You are in state 1 with 2 possible actions

⋮   ⋮                ⋮

**Formalization**

# Reinforcement Learning Problem

**Finding an optimal behavior**

► Learn optimal behavior $\pi$ based on past actions

► Maximize the expected cumulative reward over time

**Challenges**

► Feedback is delayed, not instantaneous

► Agent must reason about the long-term consequences of its actions

**Illustration**

► In order to maximize one's future income, one has to study now

► However, the immediate monetary reward from this might be negative

$\Rightarrow$ How do we learn optimal behavior?

# Trial-and-Error Learning

The agent should discover optimal behavior via trial-and-error learning

**1 Exploration**

- ▶ Try new or non-optimal actions to learn their reward
- ▶ Gain a better understanding of the environment

**2 Exploitation**

- ▶ Use current knowledge
- ▶ This might not be optimal yet, but should deviate only slightly

**Examples**

**1** Restaurant selection

- ▶ **Exploitation:** go to your favorite restaurant
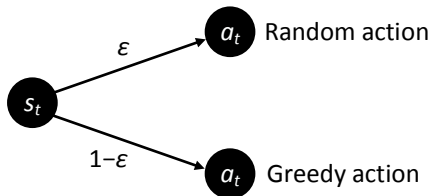- ▶ **Exploration:** try a new restaurant

**2** Game playing

- ▶ **Exploitation:** play the move you believe is best
- ▶ **Exploration:** play an experimental move

# $\varepsilon$-Greedy Action Selection

**Idea**

- ► Provide a simple heuristic to choose between exploitation and exploration
- ► Implemented via a random number $0 \leq \varepsilon \leq 1$
    - ► With probability $\varepsilon$, try a random action
    - ► With probability $1 - \varepsilon$, choose the current best



- ► Typical choice is e. g. $\varepsilon = 0.1$
- ► Other variants decrease this value over time
    - $\rightarrow$ i. e. agent gains confidence and thus needs less exploration

# Outline

# Markov Decision Process

- A Markov decision process (MDP) specifies a setup for reinforcement learning
- MDPs allow to model decision making in situations where outcomes are partly random and partly under the control of a decision maker

**Definition**

1. A Markov Decision Process is a 4-tuple $(S, A, R, T)$ with
   - A set of possible world states $S$
   - A set of possible actions $A$
   - A real-valued reward function $R$
   - Transition probabilities $T$

2. A MDP must fulfill the so-called Markov property
   - The effects of an action taken in a state depend only on that state and not on the prior history

# Markov Decision Process

**State**

- A state $s_t$ is a representation of the environment at time step $t$
- Can be directly observable to the agent or hidden

**Actions**

- At each state, the agent is able to perform an action $a_t$ that affects the subsequent state of the environment $s_{t+1}$
- Actions can be any decisions which one wants to learn

**Transition probabilities**

- Given a current state $s$, a possible subsequent state $s'$ and an action $a$
- The transition probability $T_{ss'}^a$ from $s$ to $s'$ is defined by
$$T_{ss'}^a = P\left[s_{t+1} = s' \mid s_t = s, a_t = a\right]$$

# Rewards

- A reward $r_{t+1}$ is a scalar feedback signal emitted by the environment
- Indicates how well agent is performing when reaching step $t+1$
- The expected reward $R_{ss'}^a$ when moving from state $s$ to $s'$ via action $a$ is given by

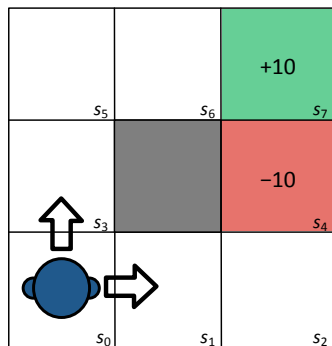$$R_{ss'}^a = E\left[r_{t+1} \mid s_t = s, a_t = a, s_{t+1} = s'\right]$$

**Examples**

1. Playing backgammon or Go
   - Zero reward after each move
   - A positive/negative reward for winning/losing a game
2. Managing an investment portfolio
   - A positive reward for each dollar left in the bank

**Goal:** maximize the expected cumulative reward over time

# Markov Decision Process

**Example**: Moving a pawn to a destination on a grid



- States $S = \{s_0, s_1, \ldots, s_7\}$
- Actions $A = \{\text{up}, \text{down}, \text{left}, \text{right}\}$
- Transition probabilities
  - $T_{s_0, s_3}^{\text{up}} = 0.9$
  - $T_{s_0, s_1}^{\text{right}} = 0.1$
  - $\ldots$
- Rewards
  - $R_{s_6, s_7}^{\text{right}} = +10$
  - $R_{s_2, s_4}^{\text{up}} = -10$
  - Otherwise $R = 0$
- Start in $s_0$
- Game over when reaching $s_7$

$\rightarrow$ available actions $A(s)$
depend on current state $s$

# Policy

**Learning task of an agent**

- Execute actions in the environment and observe results, i. e. rewards
- Learn a policy $\pi : S \to A$ that works as a selection function of choosing an action given a state
- A policy fully defines the behavior of an agent, i. e. its actions
- MDP policies depend only on the current state and not its history
- Policies are stationary (i. e. time-independent)

**Objective**

- Maximize the expected cumulative reward over time
- The expected cumulative reward from an initial state $s$ with policy $\pi$ is

$$J_\pi(s) = \sum_t R^{a_t}_{s_t, s_{t+1}} = E_\pi \left[ \sum_t r_t \mid s_0 = s \right]$$

# Value Functions

**Definition**

- The state-value function $V_\pi(s)$ of an MDP is the expected reward starting from state $s$, and then following once policy $\pi$
- $V_\pi(s) = E_\pi \left[ J_\pi(s_t) \mid s_t = s \right]$
- Quantifies how good is it to be in a particular state $s$

**Definition**

- The state-action value function $Q_\pi(s,a)$ is the expected reward starting from state $s$, taking action $a$, and then following policy $\pi$
- $Q_\pi(s,a) = E_\pi \left[ J_\pi(s_t) \mid s_t = s, a_t = a \right]$
- Quantifies how good is it to be in a particular state $s$ and apply action $a$, and afterwards follow policy $\pi$

Now, we can formalize the policy definition (with discount factor $\gamma$) via

$$\pi(s) = \arg\max_a \sum_{s'} T^a_{ss'} (R^a_{ss'} + \gamma V_\pi(s'))$$

# Optimal Value Functions

- ► While $\pi$ can be any policy, $\pi^*$ denotes the optimal one with the highest expected cumulative reward
- ► The optimal value functions specify the best possible policy
- ► A MDP is solved when the optimal value functions are known

**Definitions**

1. The optimal state-value function $V_{\pi^*}(s)$ maximizes the expected reward over all policies

$$V_{\pi^*}(s) = \max_{\pi} V_{\pi}(s)$$

2. The optimal action-value function $Q_{\pi^*}(s,a)$ maximizes the action-value function over all policies

$$Q_{\pi^*}(s,a) = \max_{\pi} Q_{\pi}(s,a)$$

# Outline

# Types of Learning Algorithms

Aim: find optimal policy and value functions

**Model-based learning**

- ▶ Aim: find optimal policy and value functions
- ▶ Model of the environment is as MDP with transition probabilities
- ▶ Approach: learn the MDP model or an approximation of it

**Model-free learning**

- ▶ Explicit model of the environment model is not available
  $\rightarrow$ i. e. transition probabilities are unknown
- ▶ Approach: derive the optimal policy without explicitly formalizing the model

# Outline

# Model-Based Learning: Policy Iteration

**Approach via policy iteration**

- Given an initial policy $\pi_0$
- Evaluate policy $\pi_i$ to find the corresponding value function $V_{\pi_i}$
- Improve policy over $V_\pi$ via greedy exploration
- Policy iteration always converges to optimal policy $\pi^*$

**Illustration**

$$\pi_0 \xrightarrow{E} V_{\pi_0} \xrightarrow{I} \pi_1 \xrightarrow{E} V_{\pi_1} \xrightarrow{I} \cdots \xrightarrow{E} V_{\pi^*} \xrightarrow{I} \pi^*$$

with

- $E$: policy evaluation
- $I$: policy improvement

# Policy Evaluation

▶ Computes the state-value function $V_\pi$ for an arbitrary policy $\pi$ via

$$V_\pi(s) = E_\pi \left[ r_{t+1} + \gamma r_{t+2} + \gamma^2 r_{t-3} + \cdots \mid s_t = s \right]$$
$$= E_\pi \left[ r_{t+1} + \gamma V_\pi(s+1) \mid s_t = s \right]$$
$$= \sum_a \pi(s,a) \sum_{s'} T_{ss'}^a \left[ R_{ss'}^a + \gamma V_\pi(s') \right]$$

▶ System of $|S|$ linear equations with $|S|$ unknowns

▶ Solvable but computational expensive if $|S|$ is large

▶ Advanced methods are available, e. g. iterative policy evaluation

**Discount factor**

▶ If $0 < \gamma < 1$, makes cumulative reward finite

▶ Necessary for setups with infinite time horizons

▶ Puts more importance on first learning steps, but less on later ones

# Iterative Policy Evaluation

- ▶ Iterative policy evaluation uses dynamic programming
- ▶ Iteratively approximate $V_\pi$
- ▶ Choose $V_0$ arbitrarily
- ▶ Then use Bellman equation as an update rule

$$V_{k+1}(s) = E_\pi \left[ r_{t+1} + \gamma V_k(s+1) \mid s_t = s \right]$$
$$= \sum_a \pi(s,a) \sum_{s'} T_{ss'}^a \left[ R_{ss'}^a + \gamma V_k(s') \right]$$

- ▶ Sequence $V_k, V_{k+1}, \ldots$ converges to $V_\pi$ as $k \to \infty$

# Policy Improvement

- Policy evaluation determines the value function $V_\pi$ for a policy $\pi$
- The alternative step exploits this knowledge to select the optimal action in each state
- For that, policy improvement searches policy $\pi'$ that is as good as or better than $\pi$
- Remedy is to use state-action value function via

$$
\begin{aligned}
\pi'(s) &= \arg\max_a Q_\pi(s,a) \\
&= \arg\max_a E\left[r_{t+1} + \gamma V_k(s+1) \mid s_t = s\right] \\
&= \arg\max_a \sum_{s'} T_{ss'}^a \left[R_{ss'}^a + \gamma V_k(s')\right]
\end{aligned}
$$

- Afterwards, continue with policy evaluation and policy improvement until a desired convergence criterion is reached

# Outline

# Model-Free Learning

**Drawbacks of model-based learning**

- ► Requires MDP, i. e. explicit model of the dynamics in the environment
- ► Transition probabilities are often not available or difficult to define
- ► Model-based learning is thus often intractable even in "simple" cases

**Model-free learning**

- ► Idea: learn directly from interactions with the environment
- ► Only use experience from the sequences of states, action, and rewards

**Common approaches**

1. **Monte Carlo methods** are simple but has slow convergence
2. **Q-learning** is more efficient due to off-policy learning

# Monte Carlo Method

- ► Monte Carlo methods require no knowledge of transition as in MDPs
- ► Perform reinforcement learning from a sequence of interactions
- ► Mimic policy iteration to find optimal policy
- ► Estimate the value of each action $Q(s,a)$ instead of $V(s)$
- ► Store average rewards in state-action table

**Example**

- ► State-action table

| State | Actions | | Optimal Policy |
|:---:|:---:|:---:|:---:|
| | $a_1$ | $a_2$ | |
| $s_1$ | 2 | 1 | $a_1$ |
| $s_2$ | 1 | 3 | $a_2$ |
| $s_3$ | 2 | 4 | $a_2$ |

# Monte Carlo Method

**Algorithm**

1. Start with an arbitrary state-action table (and corresponding policies)
   $\rightarrow$ Often all rewards are initially set to zero

2. Observe first state

3. Choose an action according to $\varepsilon$-greedy action selection, i. e.
   - With probability $\varepsilon$, pick a random action
   - Otherwise, take action with highest expected reward

4. Update state-action table with new reward (averaging)

5. Observe new state

6. Go to step 3

**Disadvantage**

- High computational time and thus slow convergence
  $\rightarrow$ Method must frequently evaluate a suboptimal policy

# Q-Learning

- One of the most important breakthroughs in reinforcement learning
- Off-policy learning concept
  - Explore the environment and at the same time exploit the current knowledge
- In each step, take a look forward to the next state and observe the maximum possible reward for all available actions in that state
- Use this knowledge to update the action-value of the corresponding action in the current state
- Apply update rule with learning rate $\alpha$ ($0 < \alpha \leq 1$)

$$Q(s,a) \leftarrow \underbrace{Q(s,a)}_{\text{old value}} + \underbrace{\alpha}_{\substack{\text{learning} \\ \text{rate}}} \left[ \underbrace{r'}_{\text{reward}} + \underbrace{\gamma}_{\substack{\text{discount} \\ \text{factor}}} \underbrace{\max_{a'} Q(s',a')}_{\text{expected optimal value}} - \underbrace{Q(s,a)}_{\text{old value}} \right]$$

- Q-learning is repeated for different episodes (e. g. games, trials, etc.)

# Q-Learning

**Algorithm**

**1** Initialize the table $Q(s,a)$ to zero for all state-action pairs $(s,a)$

**2** Observe the current state $s$

**3** Repeat until convergence

- Select an action $a$ and apply it
- Receive immediate reward $r$
- Observe the new state $s'$
- Update the table entry for $Q(s,a)$ according to

$$Q(s,a) \leftarrow Q(s,a) + \alpha \left[ r + \gamma \max_{a'} Q(s',a') - Q(s,a) \right]$$

- Move to next state, i.e. $s \leftarrow s'$

# Outline

# Wrap-Up

**Summary**

- Reinforcement learning learns through trial-and-error from interactions
- The reward indicates the performance of the agent
  $\rightarrow$ But without showing how to improve its behavior
- Learning is grouped into model-based and model-free strategies
- A common and efficient model-free variant is Q-learning
- Similar to human-like learning in real-world environments
- Common for trade-offs between long-term vs. short-term benefits

**Drawbacks**

- Can be computational expensive when state-action space is large
- No R library is yet available for model-free learning