# Reinforcement Learning

Mathematical Foundation

Dr. Alireza Aghamohammadi

## Transition Function

❖ We need a way to model the environment and its dynamics.

❖ The agent interacts with the environment through actions, which may cause the environment to change states.

❖ The function that describes this state change is called the **transition function** $T(s, a, s')$.

❖ The transition function takes a current state $s$, an action $a$, and a next state $s'$, returning the probability of moving from $s$ to $s'$ after taking action $a$:

$$T : S \times A \times S \to [0, 1]$$

❖ Formally, the transition function is defined as the probability of reaching state $s'$ at time $t$, given the state $s$ and action $a$ at time $t - 1$:

$$P(S_t = s' \mid S_{t-1} = s, A_{t-1} = a)$$

❖ In many Reinforcement Learning (RL) algorithms, we assume this transition probability is **stationary**, meaning it does not change during training or evaluation.

# Reward Function

- ❖ The environment may provide a reward signal as feedback to the agent's actions.
- ❖ The **reward function** $R(s, a)$ assigns a numerical value indicating the desirability of a state-action pair:

$$R : S \times A \to \mathbb{R}$$

- ❖ The reward function can be defined as the expected reward at time step $t$, given the state $s$ and action $a$ at time $t - 1$:
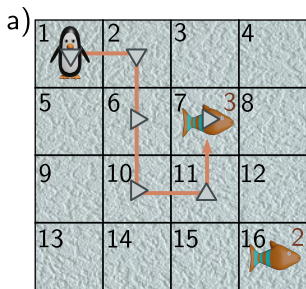
$$r(s, a) = \mathbb{E}\left[R_t \mid S_{t-1} = s, A_{t-1} = a\right]$$

- ❖ Together, the transition function and the reward function form the **model of the environment**.
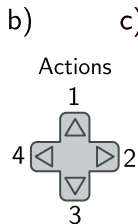
# Markov Decision Process (MDP)

❖ In RL, we often model the environment as a **Markov Decision Process (MDP)**.

❖ An MDP assumes that the probability of transitioning to the next state depends only on the current state and action, not on the full history of past states and actions. This is known as the **Markov Property**:

$$P(S_{t+1} \mid S_t, A_t) = P(S_{t+1} \mid S_t, A_t, S_{t-1}, A_{t-1}, \dots)$$



a)

b) Actions

c)

$$\tau = [1, 3, 0, 2, 3, 0, 6, 2, 0, 10, 2, 0, 11, 1, 0, 7, 2, 3]$$
$s_1\, a_1\, r_2\, s_2\, a_2\, r_3\, s_3\, a_3\, r_4\, s_4\, a_4\, r_5\, s_5\, a_5\, r_6\, s_6\, a_6\, r_7$

$$\tau = [1, 0, 3, 2, 0, 3, 6, 0, 2, 10, 0, 2, 11, 0, 1, 7, 1, 2]$$
$s_1\, r_1\, a_1\, s_2\, r_2\, a_2\, s_3\, r_3\, a_3\, s_4\, r_4\, a_4\, s_5\, r_5\, a_5\, s_6\, r_6\, a_6$

$Pr(s_{t+1}|s_t=6, a_t=1)$
$Pr(s_{t+1}|s_t=6, a_t=2)$
$Pr(s_{t+1}|s_t=6, a_t=3)$
$Pr(s_{t+1}|s_t=6, a_t=4)$

## Discount Factor

- ❖ The **discount factor** $\gamma \in [0, 1)$ determines the importance of future rewards.
- ❖ We can use the discount factor to compute the **discounted return**, which gives a measure of the total reward accumulated over time:
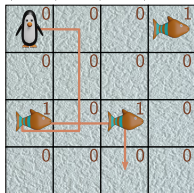
$$G_t = R_{t+1} + \gamma R_{t+2} + \gamma^2 R_{t+3} + \ldots + \gamma^{T-t-1} R_T$$

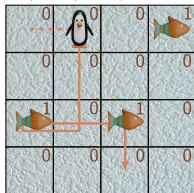- ❖ This can be simplified and generalized:

$$G_t = \sum_{k=0}^{\infty} \gamma^k R_{t+k+1}$$
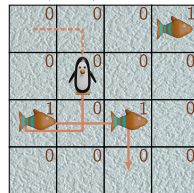
$$G_t = R_{t+1} + \gamma G_{t+1}$$

a) $G_1 = 0 + \gamma \cdot 0 + \gamma^2 \cdot 0 + \gamma^3 \cdot 0$
$+ \gamma^4 \cdot 1 + \gamma^5 \cdot 0 + \gamma^6 \cdot 1 + \gamma^7 \cdot 0 = 1.19$

b) $G_2 = 0 + \gamma \cdot 0 + \gamma^2 \cdot 0 + \gamma^3 \cdot 1$
$+ \gamma^4 \cdot 0 + \gamma^5 \cdot 1 + \gamma^6 \cdot 0 = 1.31$

c) $G_3 = 0 + \gamma \cdot 0 + \gamma^2 \cdot 1 + \gamma^3 \cdot 0$
$+ \gamma^4 \cdot 1 + \gamma^5 \cdot 0 = 1.47$



$s_1 \; r_2 \; s_2 \; r_3 \; s_3 \; r_4 \; s_4 \; r_5 \; s_5 \; r_6 \; s_6 \; r_7 \; s_7 \; r_8 \; s_8 \; r_9$
$\tau = [1, 0, 2, 0, 6, 0, 10, 0, 9, 1, 10, 0, 11, 1, 15, 0]$
$s_1, r_1, s_2, r_2, s_3, r_3, s_4, r_4, s_5, r_5, s_6, r_6, s_7, r_7, s_8, r_8$