# Function Approximation

- Function approximation in reinforcement learning
- Estimating the state-value function from on-policy data
- Approximating $v_\pi$ from experience generated using a known policy

# Approximation example

Gridworld:

| | | | | | 10 |
|---|---|---|---|---|---|
| 1 | 2 | 3 | 4 | 5 | 6 |

# Approximation example

Gridworld:

| | | | | | 10 |
|---|---|---|---|---|---|
| 1 | 2 | 3 | 4 | 5 | 6 |

State-Values:

| 5 | 6 | 7 | 8 | 9 | 10 |
|---|---|---|---|---|---|
| 1 | 2 | 3 | 4 | 5 | 6 |

# Approximation example

**Gridworld:**

| 1 | 2 | 3 | 4 | 5 | 10 6 |
|---|---|---|---|---|---|

**State-Values:**

| 5 1 | 6 2 | 7 3 | 8 4 | 9 5 | 10 6 |
|---|---|---|---|---|---|

**State-Value Approximate Function:**

$$\mathbf{w} \in \mathbb{R}^d$$

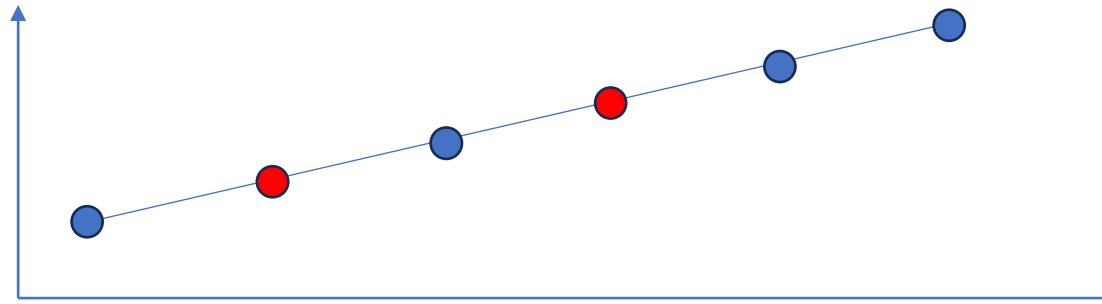$$\hat{v}(s, \mathbf{w}) \approx v_\pi(s)$$

# Approximation example

**Gridworld:**

| 1 | 2 | 3 | 4 | 5 | 10 6 |
|---|---|---|---|---|---|

**State-Values:**

| 5 1 | ? 2 | 7 3 | ? 4 | 9 5 | 10 6 |
|---|---|---|---|---|---|

**State-Value Approximate Function:**



the approximate value function is represented not as a table but as a parameterized functional form with weight vector

$$\mathbf{w} \in \mathbb{R}^d$$
$$\hat{v}(s, \mathbf{w}) \approx v_\pi(s)$$

# Why use Function Approximation?

- The dimensionality of $w$ is much less than the number of states ($d<<S$)
- When a single state is updated, the change generalizes to affect the values of many other states
- Generalization makes the learning potentially more powerful
- Can be applicable to partially observable problems
- Potentially more difficult to manage and understand
- What function approximation can't do, however, is augment the state representation with memories of past observations

ELTE | FACULTY OF INFORMATICS

# Trade-off

**Tabular based approach:**

- A continuous measure of prediction quality was not necessary because the learned value function could come to equal the true value function exactly

- The learned values at each state were decoupled (an update at one state affected no other)

**Approximation based approach:**

- By assumption there are far more states than weights, so making one state's estimate more accurate invariably means making others' less accurate

# Objective function

- Mean Squared Value Error

$$\overline{\text{VE}}(\mathbf{w}) \doteq \sum_{s \in \mathcal{S}} \mu(s) \Big[ v_\pi(s) - \hat{v}(s, \mathbf{w}) \Big]^2$$

- On-policy distribution

$$\mu(s) \geq 0, \ \sum_s \mu(s) = 1$$

# Stochastic-gradient and Semi-gradient Methods

- Weight update

$$\mathbf{w}_{t+1} \doteq \mathbf{w}_t - \frac{1}{2}\alpha\nabla\Big[v_\pi(S_t) - \hat{v}(S_t, \mathbf{w}_t)\Big]^2$$

$$= \mathbf{w}_t + \alpha\Big[v_\pi(S_t) - \hat{v}(S_t, \mathbf{w}_t)\Big]\nabla\hat{v}(S_t, \mathbf{w}_t)$$

- Partial derivate

$$\nabla f(\mathbf{w}) \doteq \left(\frac{\partial f(\mathbf{w})}{\partial w_1}, \frac{\partial f(\mathbf{w})}{\partial w_2}, \ldots, \frac{\partial f(\mathbf{w})}{\partial w_d}\right)^\top$$

- General update rule

$$\mathbf{w}_{t+1} \doteq \mathbf{w}_t + \alpha\Big[U_t - \hat{v}(S_t, \mathbf{w}_t)\Big]\nabla\hat{v}(S_t, \mathbf{w}_t)$$

# Stochastic-gradient and Semi-gradient Methods

- Weight update

$$\mathbf{w}_{t+1} \doteq \mathbf{w}_t - \frac{1}{2}\alpha\nabla\left[v_\pi(S_t) - \hat{v}(S_t,\mathbf{w}_t)\right]^2$$

$$= \mathbf{w}_t + \alpha\left[v_\pi(S_t) - \hat{v}(S_t,\mathbf{w}_t)\right]\nabla\hat{v}(S_t,\mathbf{w}_t)$$

- Partial derivate

$$\nabla f(\mathbf{w}) \doteq \left(\frac{\partial f(\mathbf{w})}{\partial w_1}, \frac{\partial f(\mathbf{w})}{\partial w_2}, \ldots, \frac{\partial f(\mathbf{w})}{\partial w_d}\right)^\top$$

- General update rule

$$\mathbf{w}_{t+1} \doteq \mathbf{w}_t + \alpha\left[U_t - \hat{v}(S_t,\mathbf{w}_t)\right]\nabla\hat{v}(S_t,\mathbf{w}_t)$$

- General update

$$s \mapsto u$$

- MC

$$S_t \mapsto G_t$$

- TD(0)

$$S_t \mapsto R_{t+1} + \gamma\hat{v}(S_{t+1},\mathbf{w}_t)$$

- N-step TD

$$S_t \mapsto G_{t:t+n}$$

ELTE | FACULTY OF INFORMATICS

# Pseudocode

**Gradient Monte Carlo Algorithm for Estimating $\hat{v} \approx v_\pi$**

Input: the policy $\pi$ to be evaluated
Input: a differentiable function $\hat{v} : \mathcal{S} \times \mathbb{R}^d \to \mathbb{R}$
Algorithm parameter: step size $\alpha > 0$
Initialize value-function weights $\mathbf{w} \in \mathbb{R}^d$ arbitrarily (e.g., $\mathbf{w} = \mathbf{0}$)

Loop forever (for each episode):
    Generate an episode $S_0, A_0, R_1, S_1, A_1, \ldots, R_T, S_T$ using $\pi$
    Loop for each step of episode, $t = 0, 1, \ldots, T-1$:
    $\mathbf{w} \leftarrow \mathbf{w} + \alpha \big[ G_t - \hat{v}(S_t, \mathbf{w}) \big] \nabla \hat{v}(S_t, \mathbf{w})$

Not biased

ELTE | FACULTY OF INFORMATICS

# Semi-gradient (bootstrapping) methods

- Do not converge as robustly as gradient methods

- Converge reliably in important cases such as the linear case

- Typically enable significantly faster learning

- Enable learning to be continual and online
(without waiting for the end of an episode)

- This enables them to be used on continuing problems and
provides computational advantages

# Pseudocode

**Semi-gradient TD(0) for estimating $\hat{v} \approx v_\pi$**

Input: the policy $\pi$ to be evaluated
Input: a differentiable function $\hat{v} : \mathcal{S}^+ \times \mathbb{R}^d \to \mathbb{R}$ such that $\hat{v}(\text{terminal}, \cdot) = 0$
Algorithm parameter: step size $\alpha > 0$
Initialize value-function weights $\mathbf{w} \in \mathbb{R}^d$ arbitrarily (e.g., $\mathbf{w} = \mathbf{0}$)

Loop for each episode:
    Initialize $S$
    Loop for each step of episode:
        Choose $A \sim \pi(\cdot|S)$
        Take action $A$, observe $R, S'$
        $\mathbf{w} \leftarrow \mathbf{w} + \alpha \big[R + \gamma\hat{v}(S',\mathbf{w}) - \hat{v}(S,\mathbf{w})\big]\nabla\hat{v}(S,\mathbf{w})$
        $S \leftarrow S'$
    until $S'$ is terminal

Bootstrapped

# State Aggregation

- A simple form of generalizing function approximation in which states are grouped together, with one estimated value
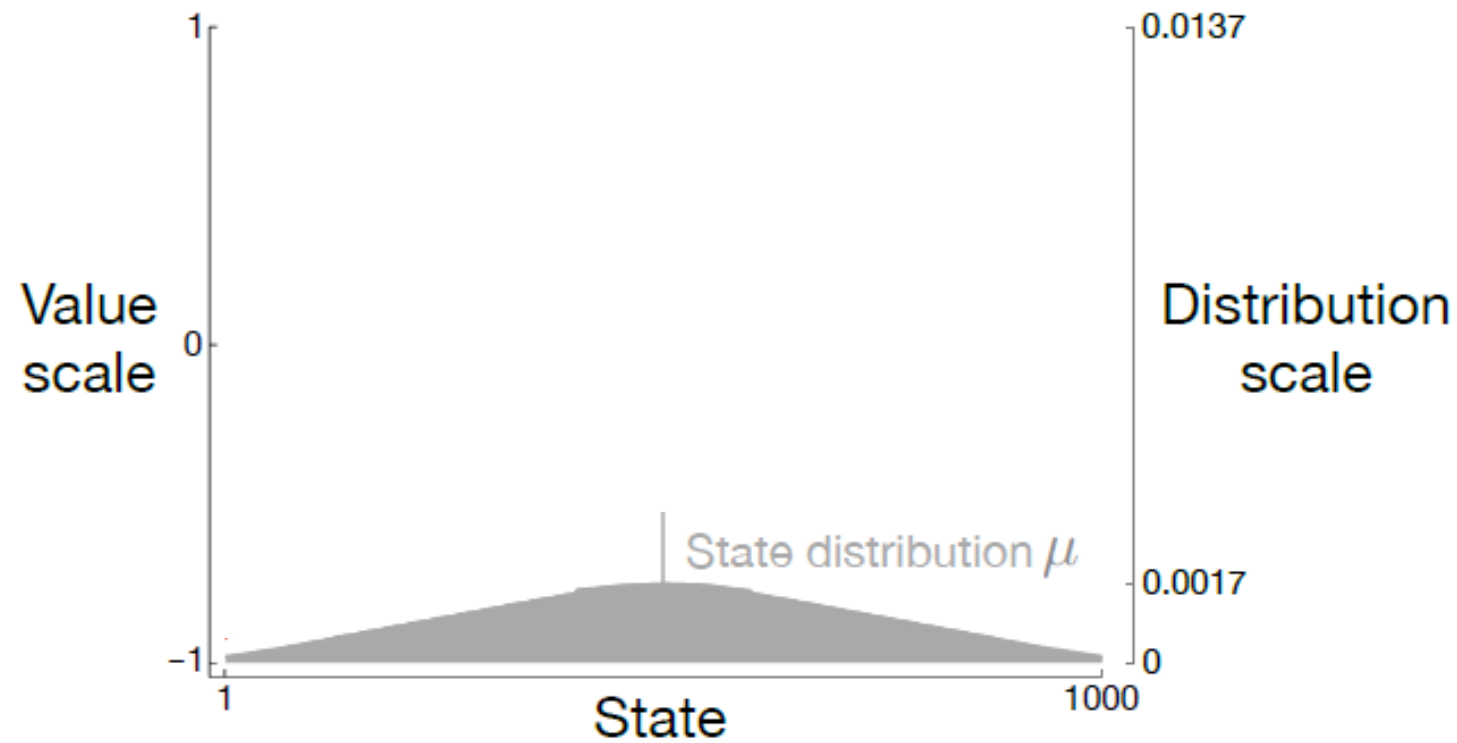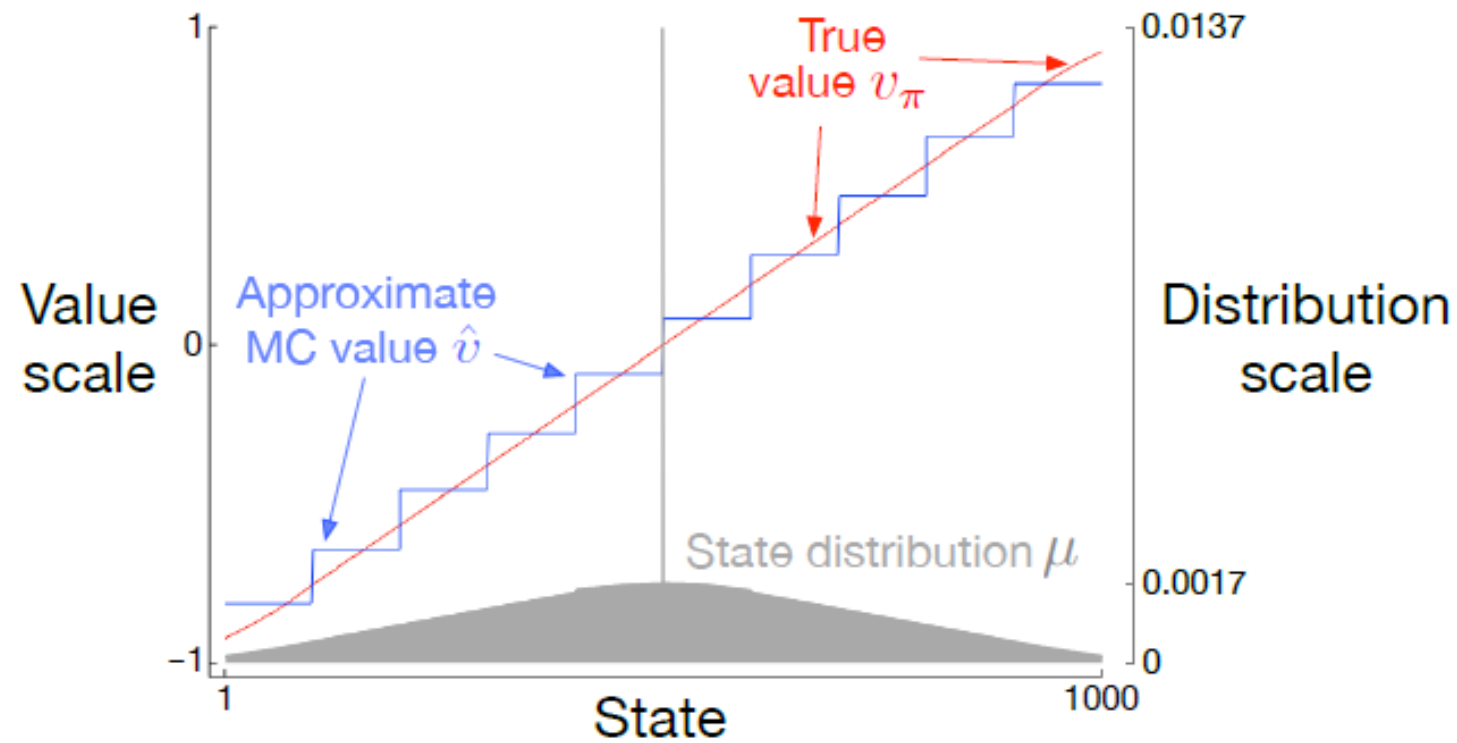
# 1000-state Random Walk



State aggregation:
10 groups
of 100 states

# 1000-state Random Walk



State aggregation:
10 groups
of 100 states

2024. 07. 23.

# Linear function approximation

- special cases in which the approximate function is a linear function of the weight vector

$$\hat{v}(s, \mathbf{w}) \doteq \mathbf{w}^\top \underline{\mathbf{x}(s)} \doteq \sum_{i=1}^{d} w_i x_i(s)$$

feature vector

# Linear function approximation

- Use SGD updates with linear function approximation
- The gradient of the approximate value function

$$\nabla \hat{v}(s, \mathbf{w}) = \mathbf{x}(s)$$

- SGD update in Linear form:

$$\mathbf{w}_{t+1} \doteq \mathbf{w}_t + \alpha \left[ U_t - \hat{v}(S_t, \mathbf{w}_t) \right] \mathbf{x}(S_t)$$

# Feature Construction for Linear Methods

- Depends critically on how the states are represented in terms of features

- Features appropriate to the task is an important way of adding prior domain knowledge

- Limitation of the linear form:
cannot take into account any interactions between features

# Coarse coding

- Features corresponding to circles in state space

- Binary features
  - state is inside a circle: the corresponding feature value = 1 (otherwise 0)

- Coarsely code for the location of the state
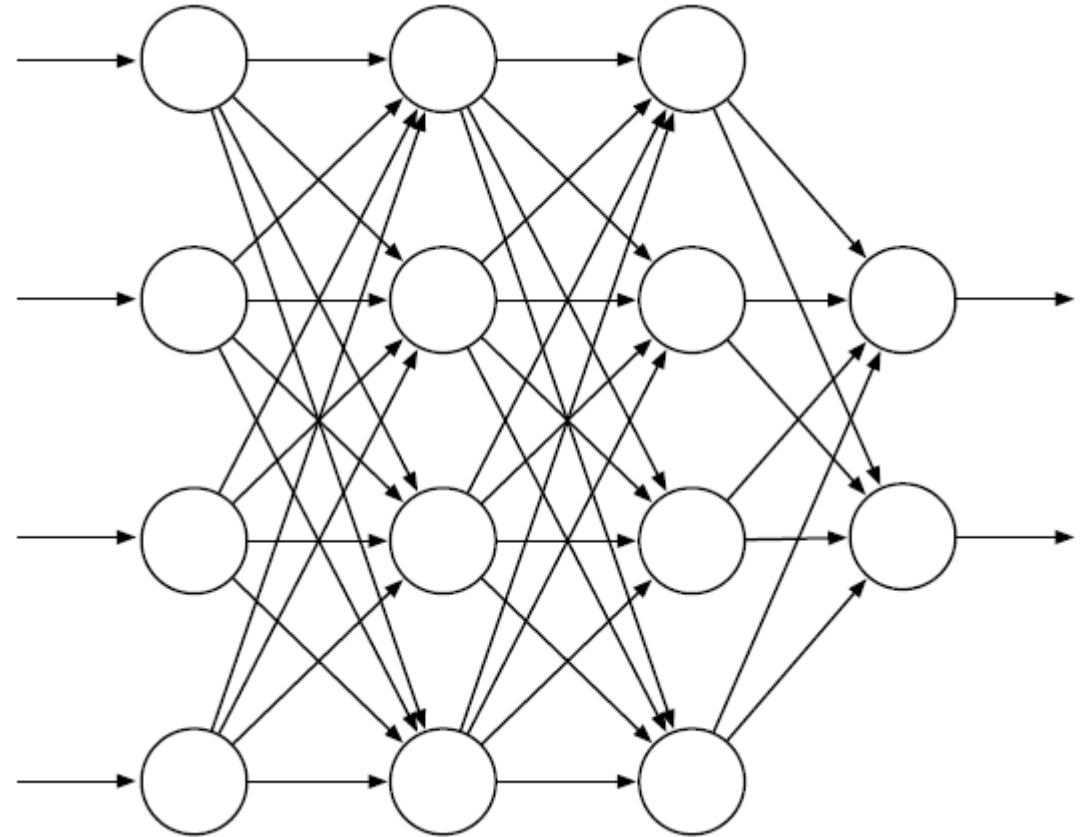
# Effect of Coarse Coding

# Tile coding

- A form of coarse coding for multi-dimensional continuous spaces
- Flexible and computationally efficient
- Practical feature representation for sequential data

# Nonlinear function approximation

**Universal approximation theorem**: any continuous function can be approximated arbitrarily well by a **neural network** with at least 1 hidden layer with a finite number of weights

## Value Learning

Find $Q(s,a)$
$a = \underset{a}{\mathrm{argmax}}\ Q(s,a)$

## Policy Learning

Find $\pi(s)$
Sample $a \sim \pi(s)$

# Deep Reinforcement Learning Algorithms

**Value Learning**

Find $Q(s,a)$

$a = \underset{a}{\mathrm{argmax}}\, Q(s,a)$

**Policy Learning**
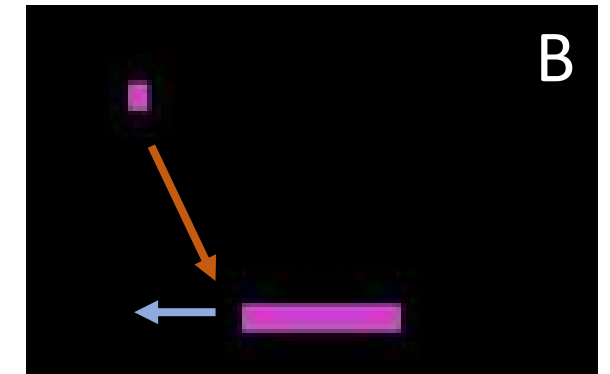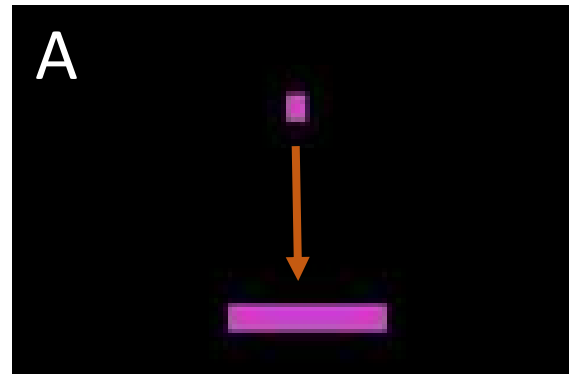
Find $\pi(s)$
Sample $a \sim \pi(s)$

# Q function intuition
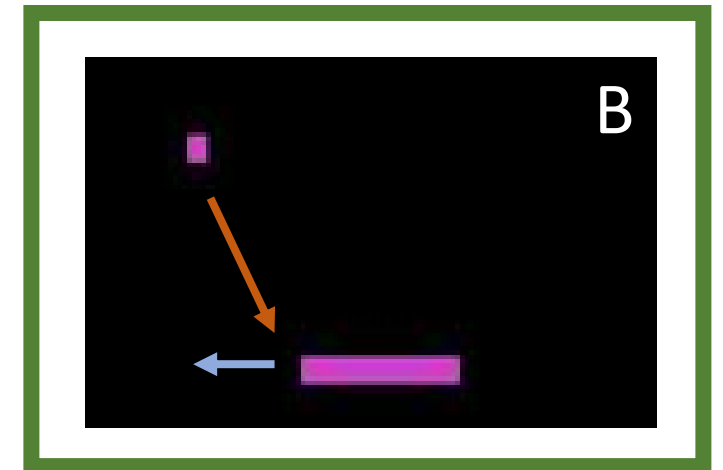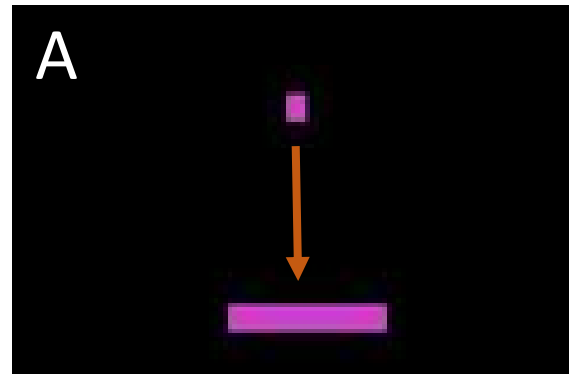


Atari game - Breakout

# Q function intuition



Atari game - Breakout

It can be very difficult for humans to accurately estimate Q-values



A



B

Which (s, a) pair has a higher Q-value? 🤔

# Q function intuition



Atari game - Breakout

It can be very difficult for humans to accurately estimate Q-values



A



B

Which (s, a) pair has a higher Q-value? 🤔

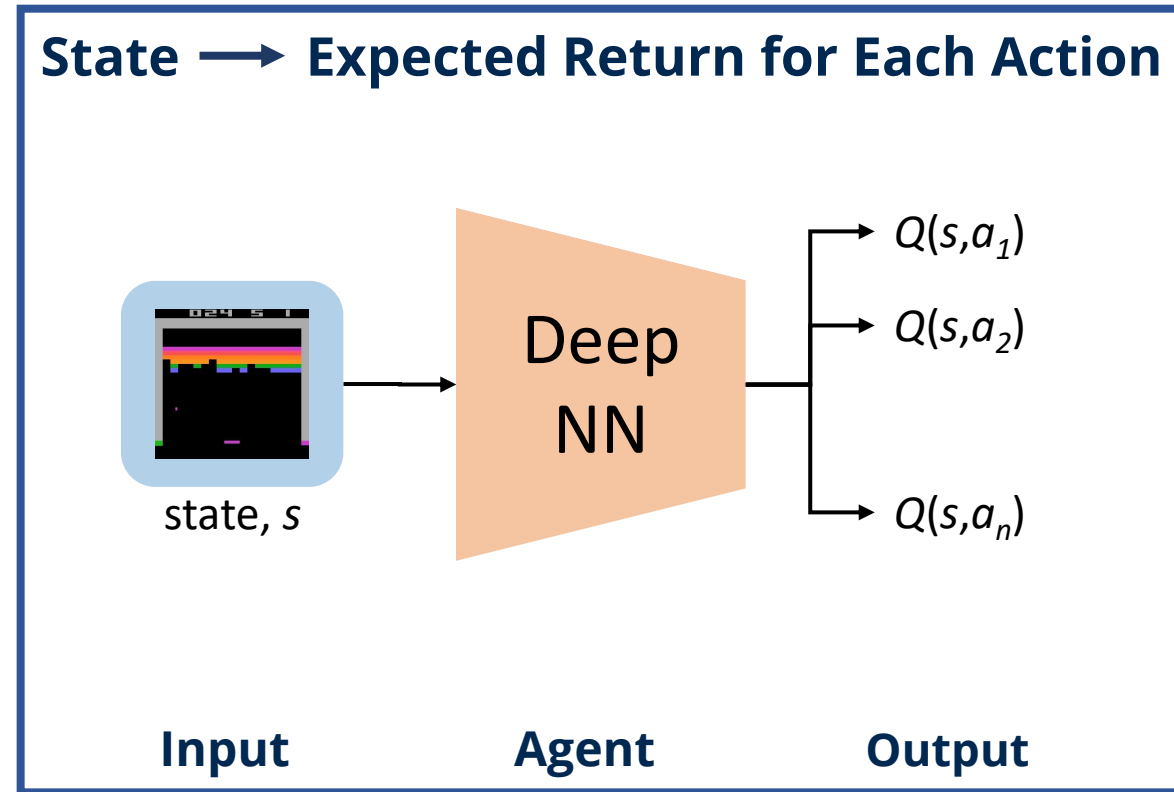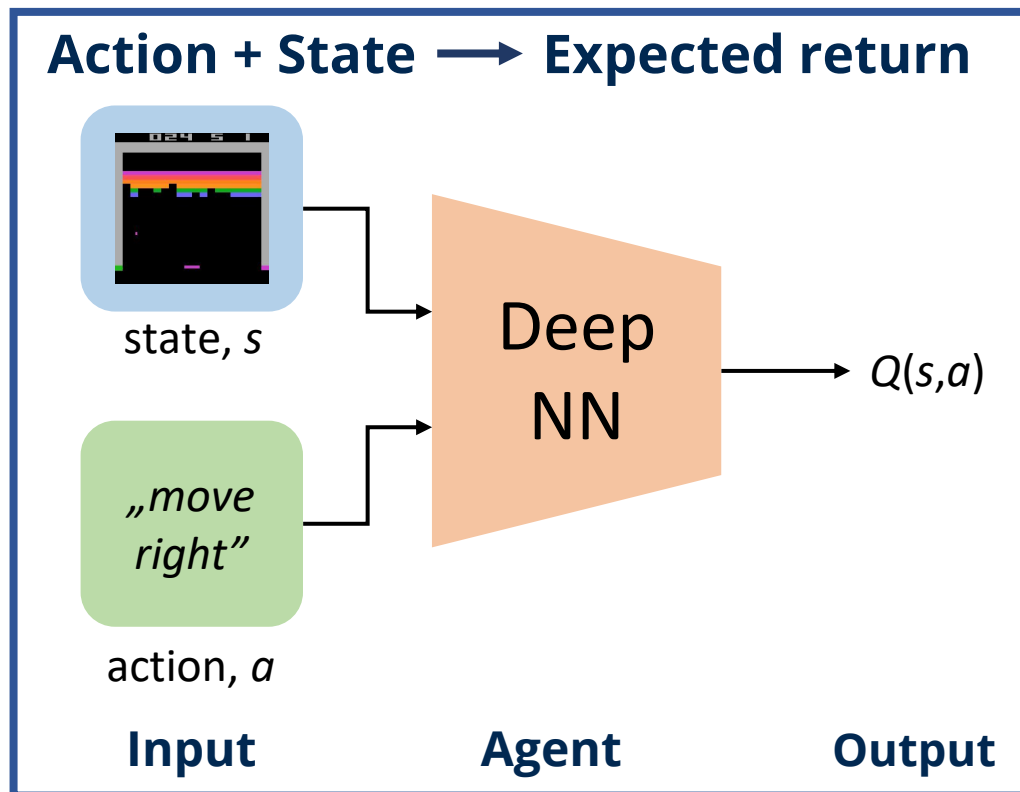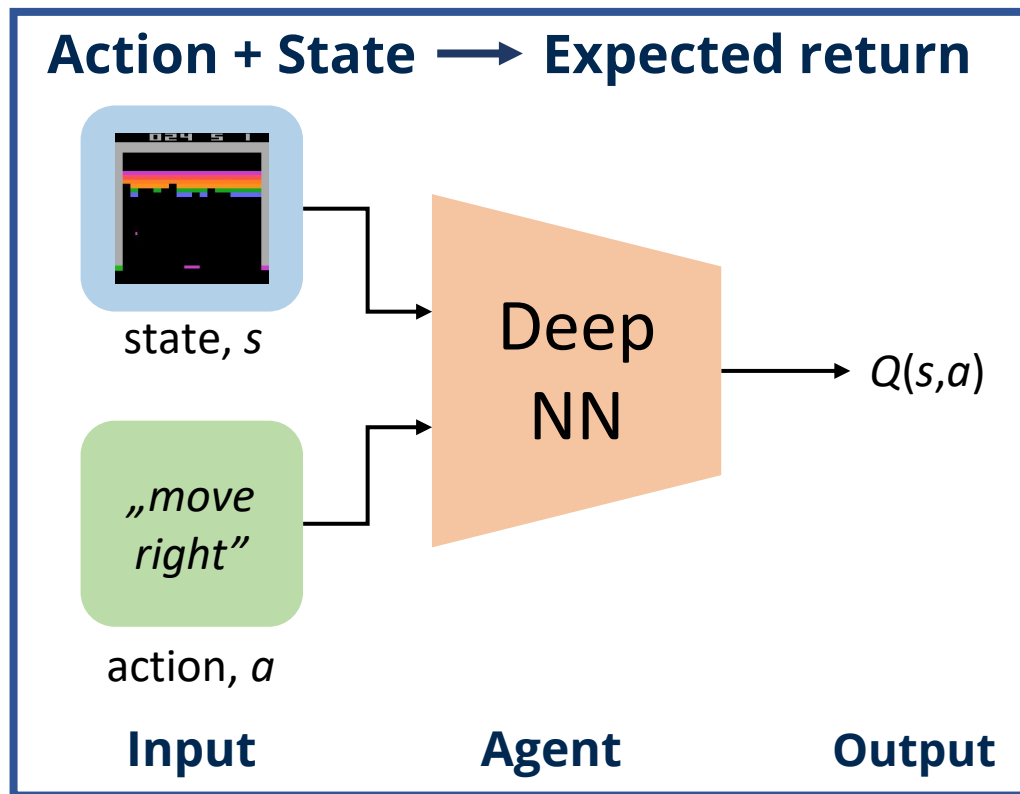ELTE | FACULTY OF INFORMATICS

# Deep Q Networks (DQN): Training

How can we use deep neural networks to model Q-functions?

# Deep Q Networks (DQN): Training

How can we use deep neural networks to model Q-functions?



**Action + State → Expected return**
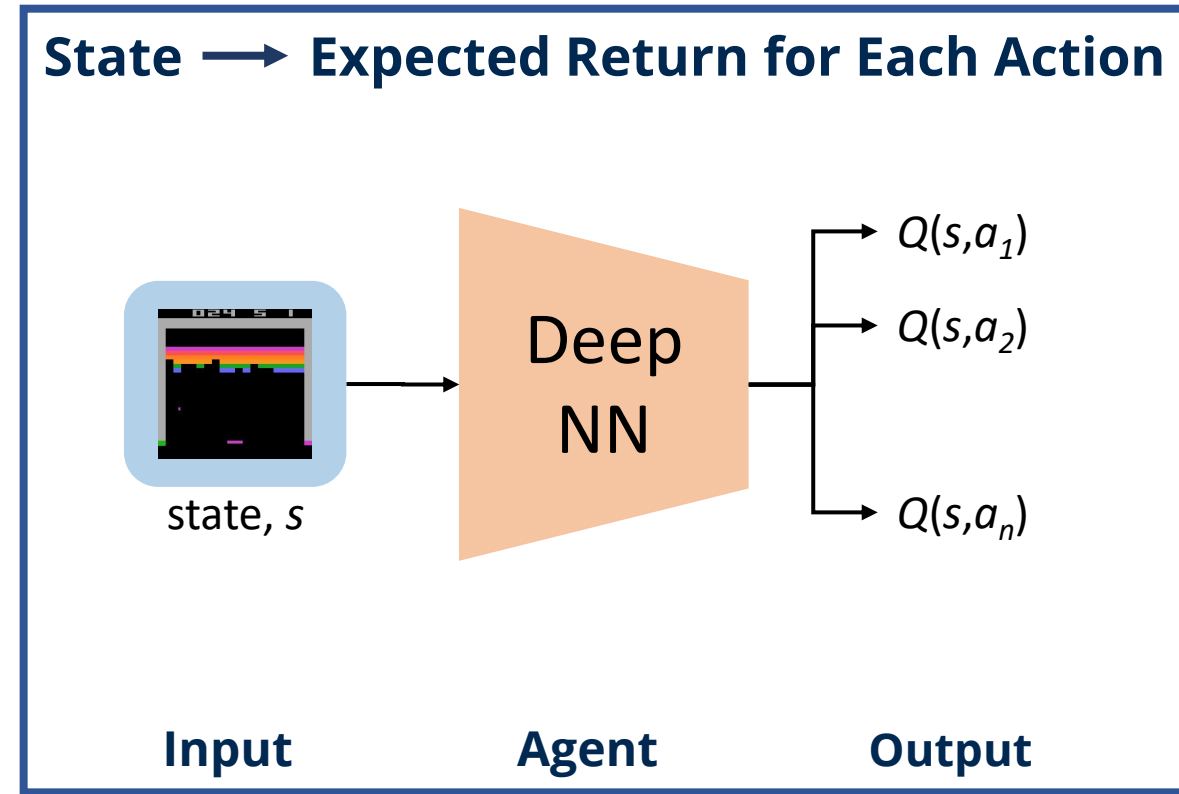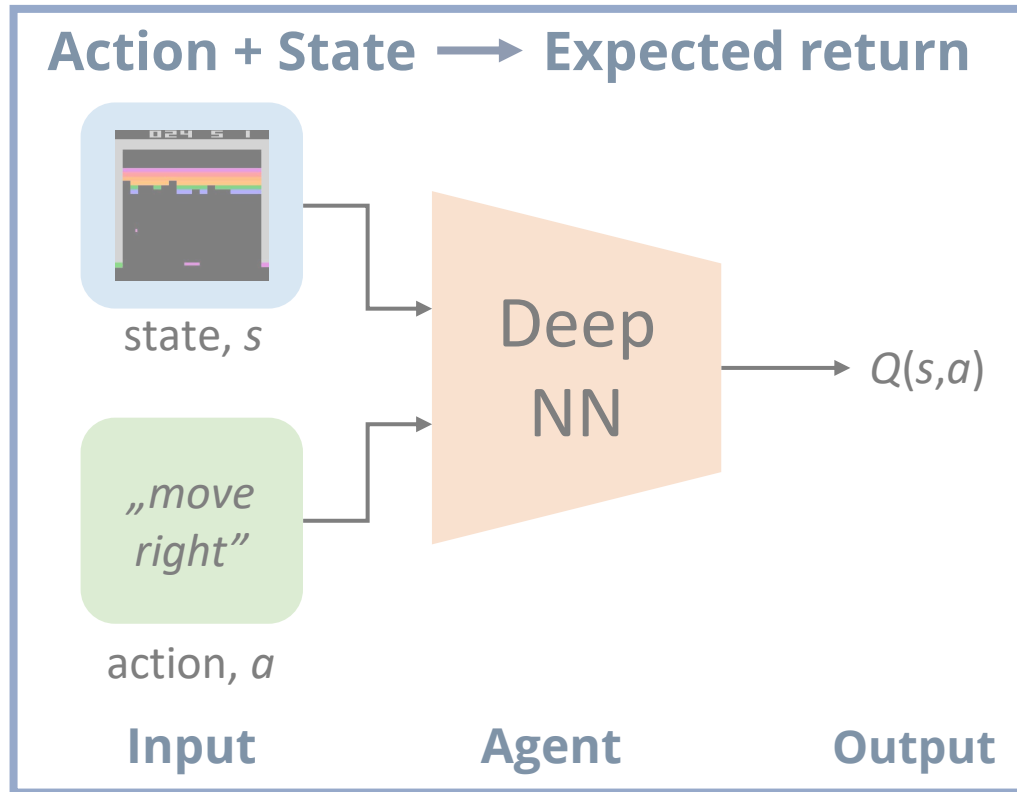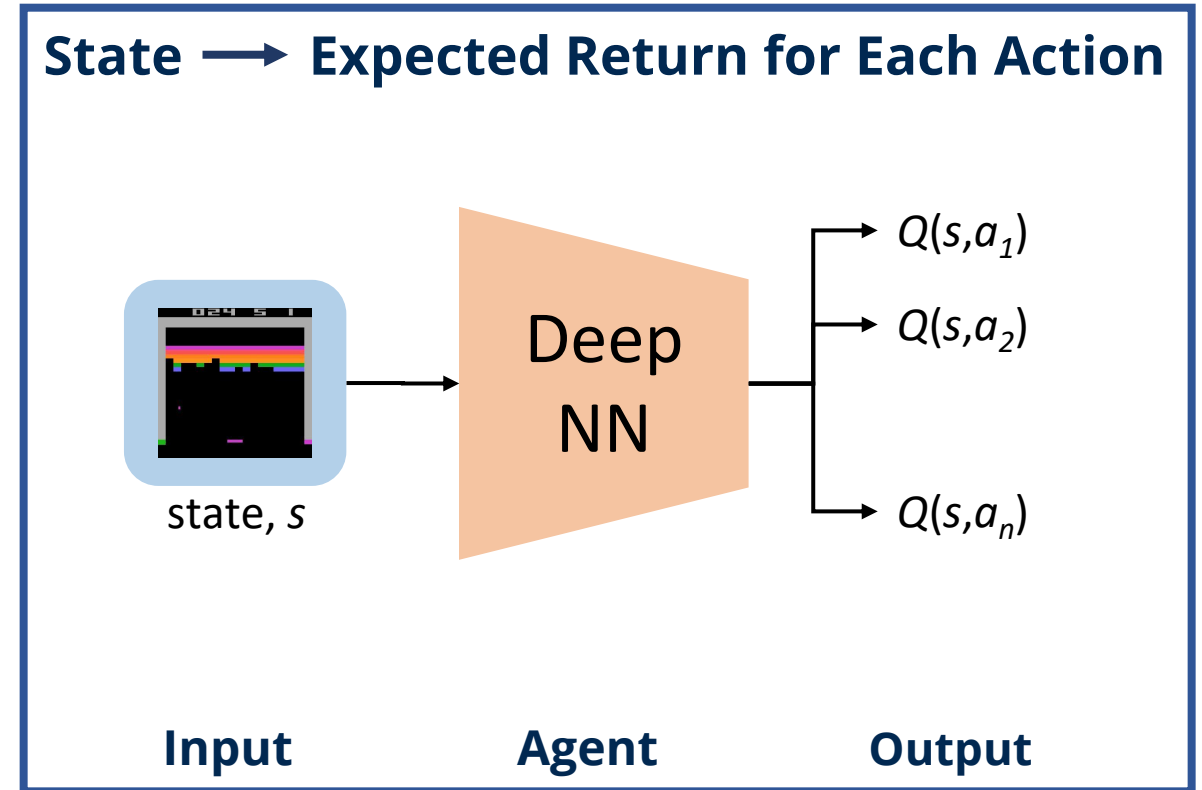
state, $s$

„move right"

action, $a$

Deep NN

$Q(s,a)$

**Input**     **Agent**     **Output**

**State → Expected Return for Each Action**

state, $s$

Deep NN

$Q(s,a_1)$

$Q(s,a_2)$

$Q(s,a_n)$

**Input**     **Agent**     **Output**

# Deep Q Networks (DQN): Training

## How can we use deep neural networks to model Q-functions?



**Action + State ⟶ Expected return**

state, $s$

„move right"

action, $a$

Deep NN

$Q(s,a)$

**Input**       **Agent**       **Output**

**Multiple evaluation needed**

**State ⟶ Expected Return for Each Action**

state, $s$

Deep NN

$Q(s,a_1)$

$Q(s,a_2)$

$Q(s,a_n)$

**Input**       **Agent**       **Output**

**More efficient**

# Deep Q Networks (DQN): Training

How can we use deep neural networks to model Q-functions?



**Action + State → Expected return**

state, $s$

„move right"

action, $a$

Deep NN

$Q(s,a)$

**Input**          **Agent**          **Output**

**Multiple evaluation needed**

**State → Expected Return for Each Action**

state, $s$

Deep NN

$Q(s,a_1)$

$Q(s,a_2)$

$Q(s,a_n)$

**Input**          **Agent**          **Output**
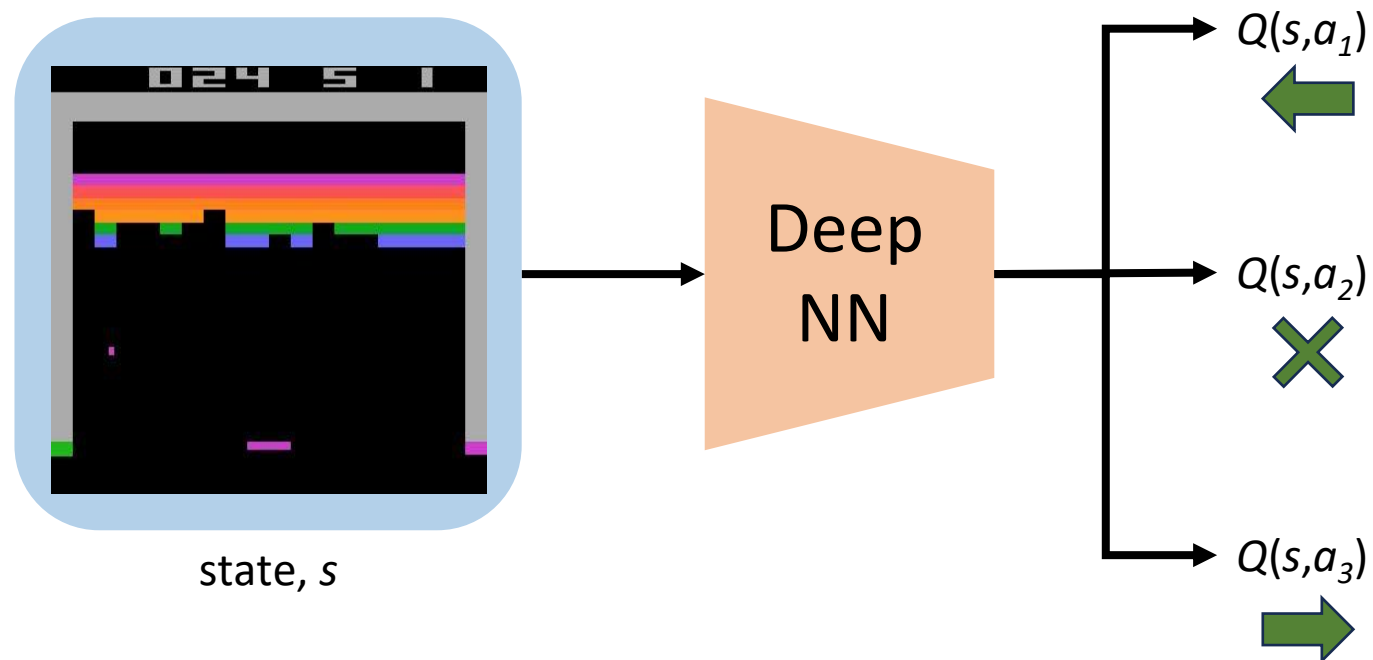
**More efficient**

# Deep Q Networks Summary

Use NN to learn Q-function and then use to infer the optimal policy, $\pi(s)$



state, $s$

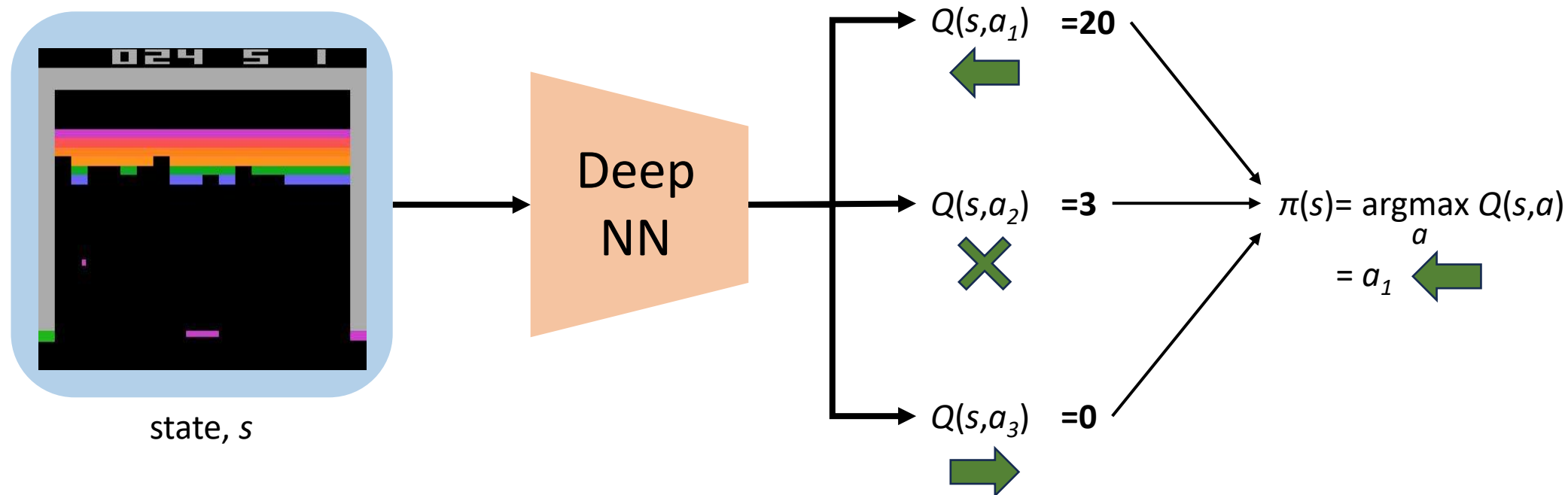$Q(s,a_1)$

$Q(s,a_2)$

$Q(s,a_3)$

# Deep Q Networks Summary

Use NN to learn Q-function and then use to infer the optimal policy, $\pi(s)$



state, $s$

Deep NN

$Q(s,a_1)$ **=20**

$Q(s,a_2)$ **=3**

$Q(s,a_3)$ **=0**

# Deep Q Networks Summary

Use NN to learn Q-function and then use to infer the optimal policy, $\pi(s)$



state, $s$

Deep NN

$Q(s,a_1)$ =**20**

$Q(s,a_2)$ =**3**

$Q(s,a_3)$ =**0**

$\pi(s)= \underset{a}{\text{argmax }} Q(s,a)$

$= a_1$

# Deep Q Networks Summary

Use NN to learn Q-function and then use to infer the optimal policy, $\pi(s)$



state, $s$

$Q(s,a_1)$ =20

$Q(s,a_2)$ =3

$Q(s,a_3)$ =0

$\pi(s)= \underset{a}{\mathrm{argmax}}\, Q(s,a)$

$= a_1$

Send action back to the environment and receive next state

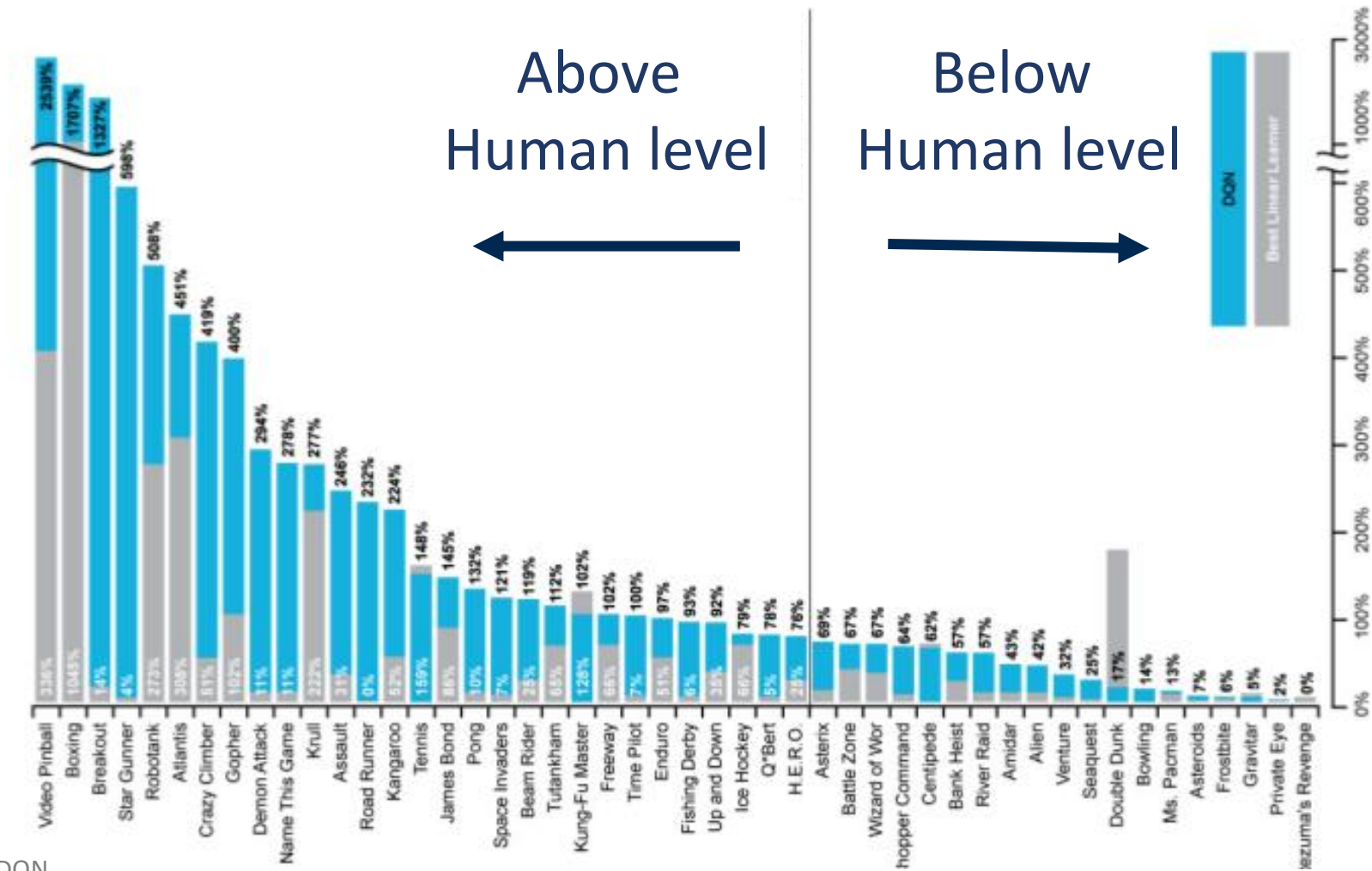ELTE | FACULTY OF INFORMATICS

# DQN Atari playing Network



Mnih, V., Kavukcuoglu, K., Silver, D. *et al.* Human-level control through deep reinforcement learning. *Nature* **518**, 529–533 (2015). https://doi.org/10.1038/nature14236

# DQN Atari Results



Above Human level ← | → Below Human level

# Downsides of Q-learning

Complexity:

- Can model scenarios where the action space is discrete and small
- Cannot handle continuous action spaces

Flexibility:

- Policy is deterministically computed from the Q function by maximizing the reward → cannot learn stochastic policies

# Downsides of Q-learning

Complexity:

• Can model scenarios where the action space is discrete and small

• Cannot handle continuous action spaces

Flexibility:

• Policy is deterministically computed from the Q function by maximizing the reward → cannot learn stochastic policies

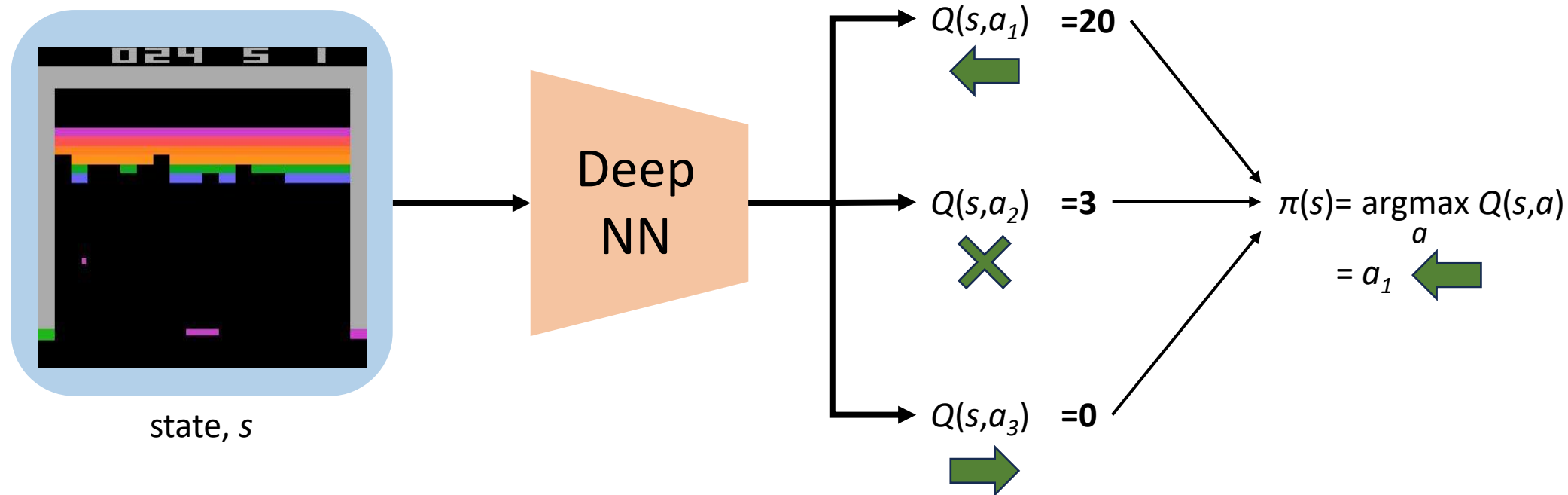**To address these, consider a new class of RL training algorithms: Policy gradient methods**

# Deep Reinforcement Learning Algorithms

**Value Learning**

Find $Q(s,a)$

$a = \underset{a}{\text{argmax}}\, Q(s,a)$

**Policy Learning**

Find $\pi(s)$
Sample $a \sim \pi(s)$
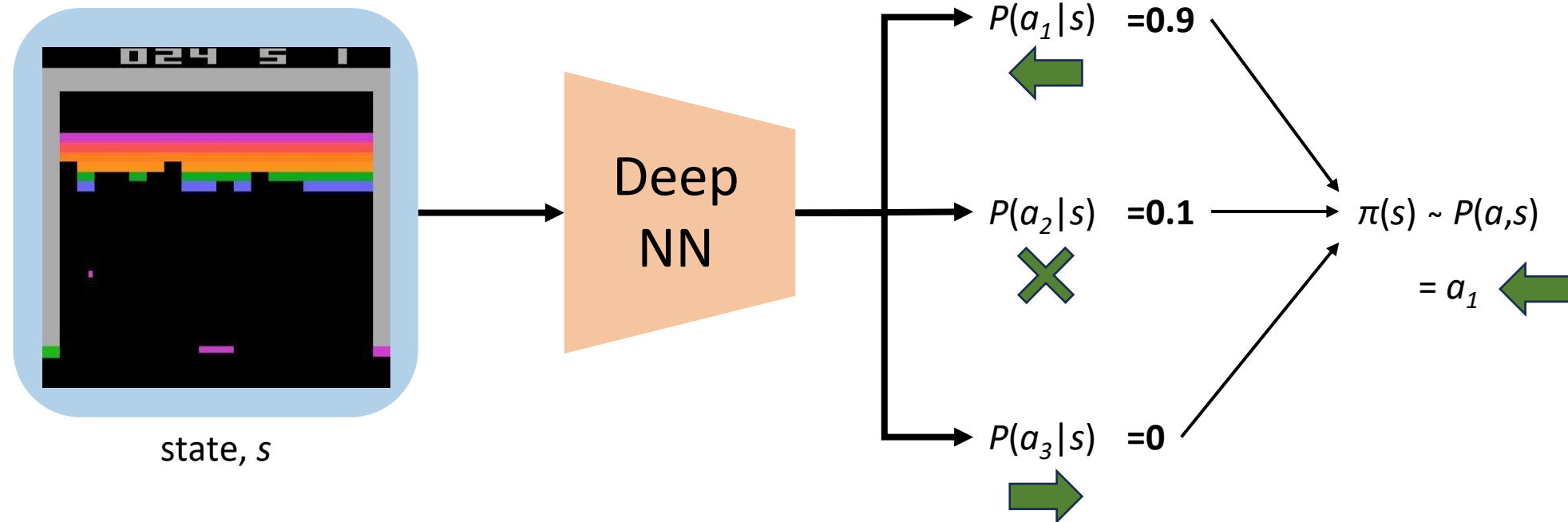
# Deep Q Networks Summary

**DQN**: Approximate Q function and use to infer the optimal policy, $\pi(s)$



state, $s$

Deep NN

$Q(s,a_1)$ =**20**

$Q(s,a_2)$ =**3**

$Q(s,a_3)$ =**0**

$\pi(s) = \underset{a}{\mathrm{argmax}}\ Q(s,a)$

$= a_1$

# Policy Gradient (PG): Key Idea

**DQN**: Approximate Q function and use to infer the optimal policy, $\pi(s)$

**Policy Gradient**: Directly optimize the policy $\pi(s)$

state, $s$

Deep NN

$P(a_1|s)$ **=0.9**

$P(a_2|s)$ **=0.1**

$P(a_3|s)$ **=0**

$\pi(s) \sim P(a,s)$

$= a_1$

# Policy Gradient (PG): Key Idea

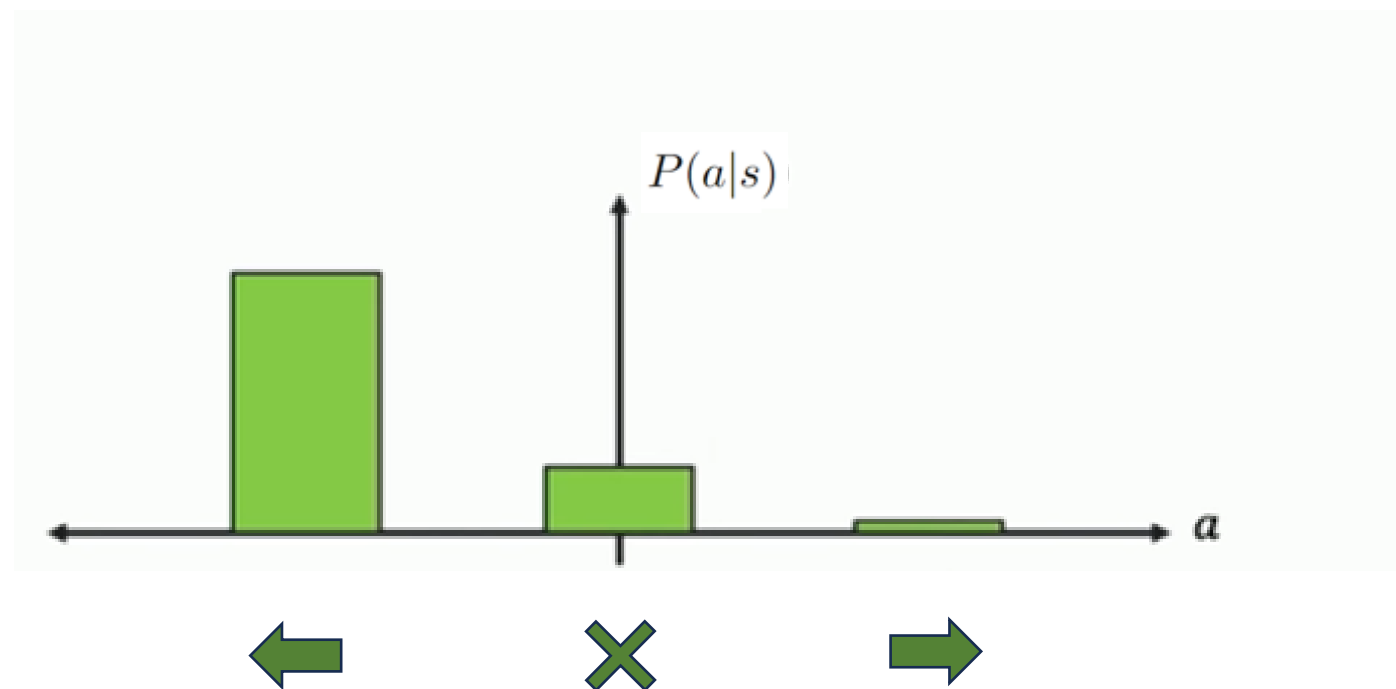**DQN**: Approximate Q function and use to infer the optimal policy, $\pi(s)$

**Policy Gradient**: Directly optimize the policy $\pi(s)$



$$\sum_{a_i \in A} P(a_i|s) = 1$$

$P(a_1|s)$ **=0.9**

$P(a_2|s)$ **=0.1**

$P(a_3|s)$ **=0**

state, $s$

Deep NN

$\pi(s) \sim P(a,s)$

$= a_1$

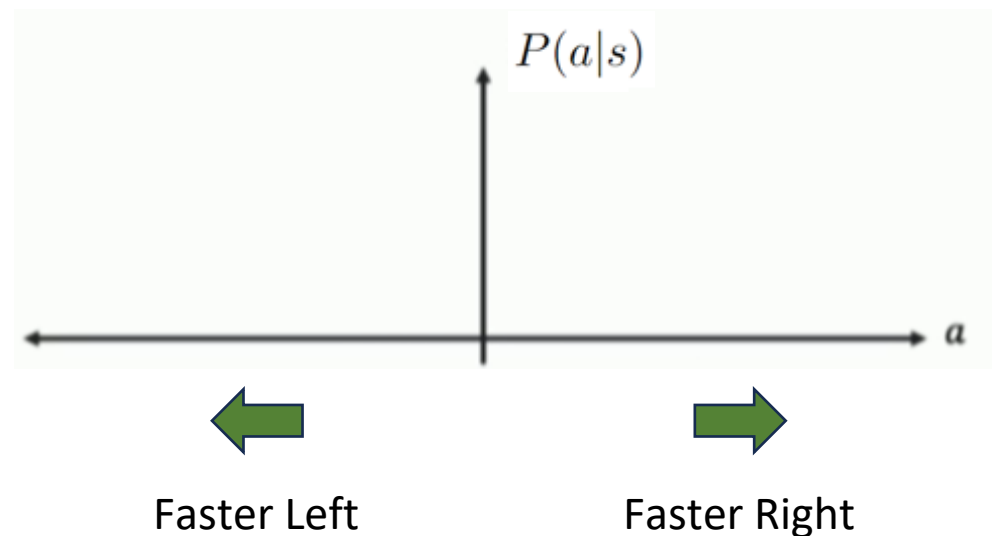**Discrete action space**: which direction should I move?  ⬅ ✖ ➡

# Discrete vs Continuous Action Spaces

**Discrete action space**: which direction should I move? ⬅ ✖ ➡

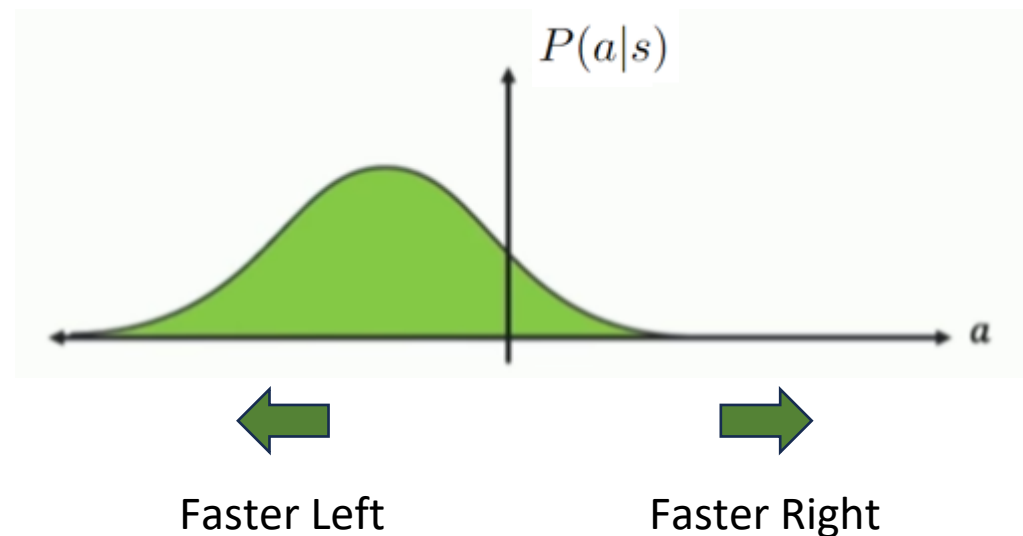**Continuous action space**: how fast should I move? ⬅ 7 m/s



$P(a|s)$

$a$

Faster Left          Faster Right

# Discrete vs Continuous Action Spaces

**Discrete action space**: which direction should I move? ⬅ ✖ ➡

**Continuous action space**: how fast should I move? ⬅ 7 m/s



$P(a|s)$

Faster Left          Faster Right

ELTE | FACULTY OF INFORMATICS

# Discrete vs Continuous Action Spaces

**Discrete action space**: which direction should I move? ⬅ ✖ ➡

**Continuous action space**: how fast should I move? ⬅ 7 m/s

Parametrization

$P(a|s)$

$a$

Faster Left          Faster Right

# Policy Gradient (PG): Key Idea

**Policy Gradient:** Enables modeling of continuous action space



state, $s$

Deep
NN

$$P(a|s) = \mathcal{N}(\mu, \sigma^2)$$

Faster Left          Faster Right

ELTE | FACULTY OF INFORMATICS

**Policy Gradient:** Enables modeling of continuous action space



state, *s*

Deep NN

Mean, **μ**　　=-1

Variance, **σ²** =-0.5

$$P(a|s) = \mathcal{N}(\mu, \sigma^2)$$

Faster Left　　　　Faster Right

# Policy Gradient (PG): Key Idea

**Policy Gradient:** Enables modeling of continuous action space



state, $s$

Deep NN

Mean, **μ** =-1

Variance, **σ²** =-0.5

$$P(a|s) = \mathcal{N}(\mu, \sigma^2)$$
$$\pi^*(s) \sim P(a|s)$$
$$= -0.8[m/s]$$

$$\int_{a=-\infty}^{\infty} P(a|s) = 1$$

$$P(a|s) = \mathcal{N}(\mu, \sigma^2)$$

Faster Left        Faster Right