# Supervised Learning

**Data**: (x, y)
x is data, y is label

**Goal**: Learn a *function* to map x -> y

**Examples**: Classification,  regression, object detection,  semantic segmentation, image  captioning, etc.
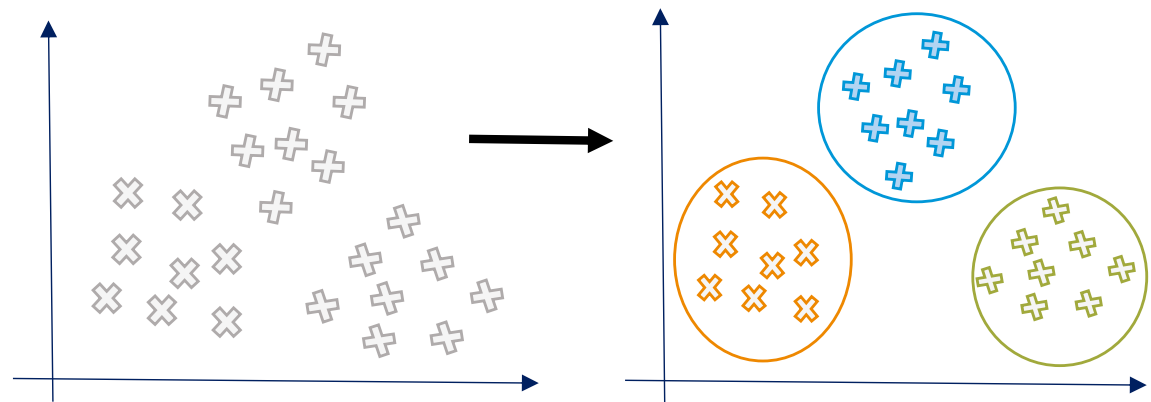


Cat

Classification

ITsMOre than a
UNIVERSITY

# Unsupervised Learning

**Data**: x
Just data, no labels!

**Goal**: Learn some underlying
hidden *structure* of the data

**Examples**: Clustering,  dimensionality
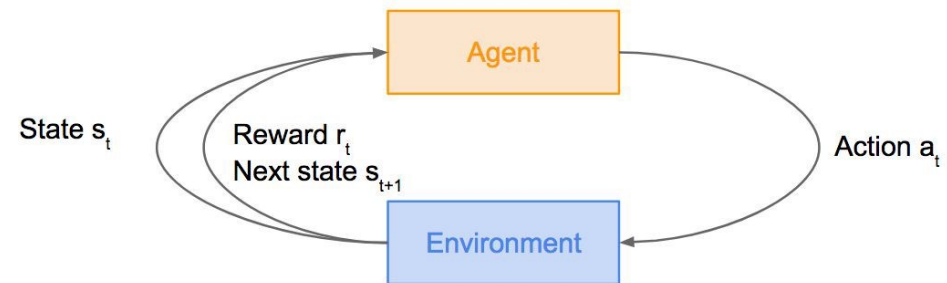reduction, feature  learning, density
estimation, etc.

density estimation

# Reinforcement Learning

Problems involving an **agent** interacting with an **environment**, which provides numeric **reward** signals
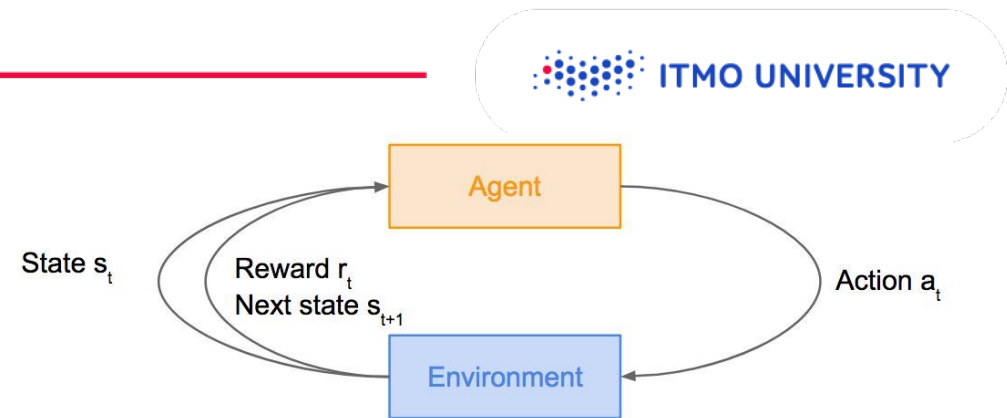
**Goal**: Learn how to take actions in order to maximize reward

# Reinforcement Learning

Problems involving an **agent** interacting with an **environment**, which provides numeric **reward** signals

**Goal**: Learn how to take actions in order to maximize reward



A person learns how to walk not by looking through a huge number of teaching examples, but by trying new things and making mistakes, so that, by getting both positive and negative experience

# Overview

- What is Reinforcement Learning?

- Markov Decision Processes

- Q-Learning

- Policy Gradients

*partly based on CS231n Stanford RL lecture*

# Reinforcement Learning

Agent

Environment

State $s_t$

Reward $r_t$
Next state $s_{t+1}$

Action $a_t$

# Cart-Pole Problem

**Objective**: Balance a pole on top of a movable cart

**State:** angle, angular speed, position, horizontal velocity
**Action:** horizontal force applied on the cart
**Reward:** 1 at each time step if the pole is upright

ITsMOre than a
UNIVERSITY

# Robot motions

**Objective**: Make the robot move forward

**State:** Angle and position of the joints
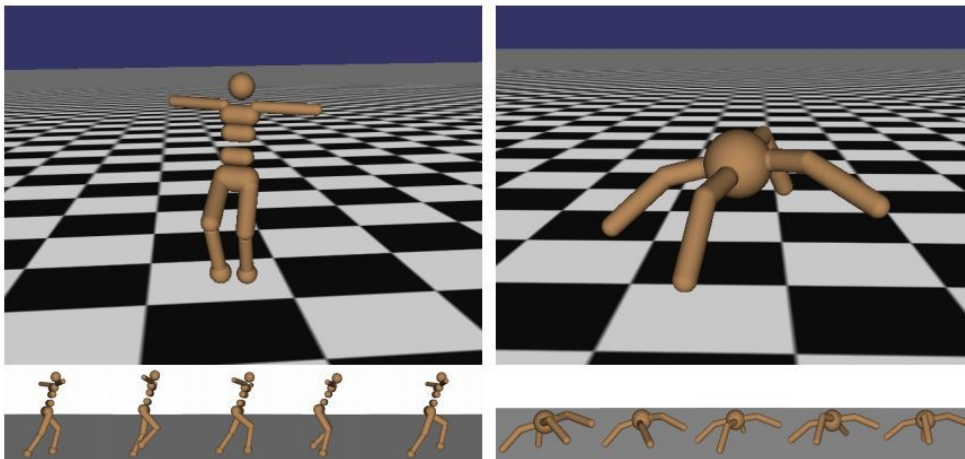**Action:** Torques applied on joints
**Reward:** 1 at each time step upright + forward movement

# Atari Games

**Objective**: Complete the game with the highest score

**State:** Raw pixel inputs of the game state
**Action:** Game controls e.g. Left, Right, Up, Down
**Reward:** Score increase/decrease at each time step

# Go

**Objective**: Win the game!

**State:** Position of all pieces
**Action:** Where to put the next piece down
**Reward:** 1 if win at the end of the game, 0 otherwise

# How can we mathematically formalize the RL problem?

State $s_t$

Reward $r_t$
Next state $s_{t+1}$

Agent

Environment

Action $a_t$

# Markov Decision Process

- Mathematical formulation of the RL problem
- **Markov property**: Current state completely characterises the state of the world

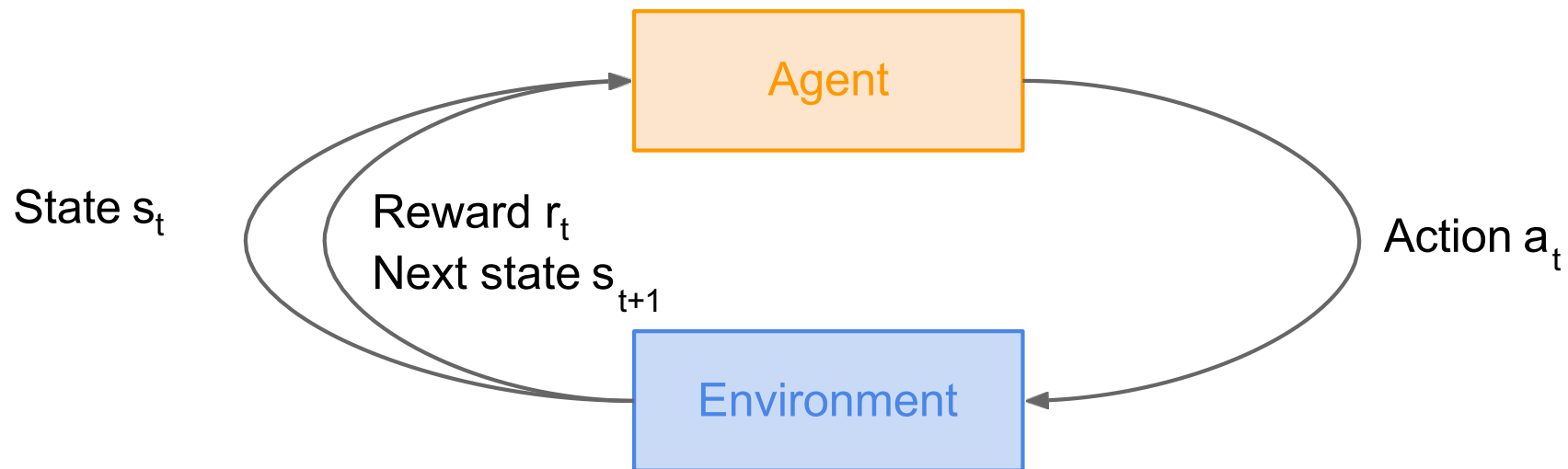Process is defined by: $(\mathcal{S}, \mathcal{A}, \mathcal{R}, \mathbb{P}, \gamma)$

$\mathcal{S}$ : set of possible states
$\mathcal{A}$ : set of possible actions
$\mathcal{R}$ : distribution of reward given (state, action) pair
$\mathbb{P}$ : transition probability i.e. distribution over next state given (state, action) pair
$\gamma$ : discount factor – the less it is the less objective depends on further rewards

ITsMOre than a
UNIVERSITY

# Markov Decision Process

- At time step t=0, environment samples initial state $s_0 \sim p(s_0)$
- Then, for t=0 until done:
    - Agent selects action $a_t$
    - Environment samples reward $r_t \sim R( \, . \mid s_t, a_t)$
    - Environment samples next state $s_{t+1} \sim P( \, . \mid s_t, a_t)$
    - Agent receives reward $r_t$ and next state $s_{t+1}$

- A policy $\pi$ is a function from S to A that specifies what action to take in  each state
- **Objective**: find policy $\pi^*$ that maximizes cumulative discounted reward:  $\sum_{t \geq 0} \gamma^t r_t$

# A simple MDP: Grid World

actions = {

1. right

2. left

3. up

4. down

}
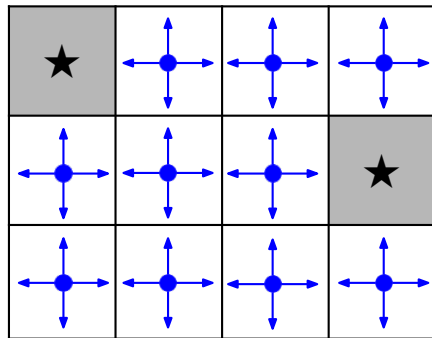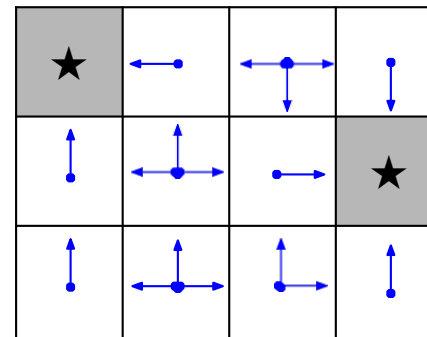
states

Set a negative "reward"
for each transition
(e.g. $r = -1$)

**Objective:** reach one of terminal states (greyed out) in
least number of actions

ITMOre than a
UNIVERSITY

# A simple MDP: Grid World



Random Policy



Optimal Policy

# "Multi-armed bandits"

Assume that we have the following setup:

- Several automates with different reward probability (give money)

- Push the arms to earn more money (Our goal is to maximize **reward** with t$\rightarrow \infty$)

- No different states for agent – only actions with reward;

- How to choose the best one? (say, «the best policy")



Pure exploration    Pure exploitation

**Greedy algorithm**: on every step we choose the arm with the highest reward.

**Exploration algorithm**: on every step choose another arm in a hope that it will give even more reward.

**ε-greedy algorithm:** with probability **(1-ε)** we perform in greedy way and with **ε we explore the system.**



Accuray of epsilon-Greedy alg.     Avg. Rewards of epsilon-Greedy alg.     Cumulative Rewards of epsilon-Greedy alg.

# The optimal policy п*

We want to find optimal policy п* that maximizes the sum of rewards.

How do we handle the randomness (initial state, transition probability…)?
Maximize the **expected sum of rewards!**

Formally: $\pi^* = \arg\max_{\pi} \mathbb{E}\left[\sum_{t \geq 0} \gamma^t r_t \mid \pi\right]$ with $s_0 \sim p(s_0), a_t \sim \pi(\cdot|s_t), s_{t+1} \sim p(\cdot|s_t, a_t)$

# Definitions: Value function and Q-value function

Following a policy produces sample trajectories (or paths) $s_0$, $a_0$, $r_0$, $s_1$, $a_1$, $r_1$, …

How to estimate how good the current state is?

# Chess

Reward is known only in the end of a game….

How good the current state is?
How to choose the next step?



ITsMOre than a
UNIVERSITY

# Definitions: Value function and Q-value function

Following a policy produces sample trajectories (or paths) $s_0$, $a_0$, $r_0$, $s_1$, $a_1$, $r_1$, …

How good is a state?
The **value function** at state s, is the expected cumulative reward from following the policy from state s:

$$V^\pi(s) = \mathbb{E}\left[\sum_{t\geq 0} \gamma^t r_t | s_0 = s, \pi\right]$$

How good is a state-action pair?
The **Q-value function** at state s and action a, is the expected cumulative reward from taking action a in state s and then following the policy:

$$Q^\pi(s,a) = \mathbb{E}\left[\sum_{t\geq 0} \gamma^t r_t | s_0 = s, a_0 = a, \pi\right]$$

# Bellman equation

The optimal Q-value function Q* is the maximum expected cumulative reward achievable from a given (state, action) pair:

$$Q^*(s, a) = \max_{\pi} \mathbb{E}\left[\sum_{t \geq 0} \gamma^t r_t | s_0 = s, a_0 = a, \pi\right]$$

Q* satisfies the following **Bellman equation**:

$$Q^*(s, a) = \mathbb{E}_{s' \sim \mathcal{E}}\left[r + \gamma \max_{a'} Q^*(s', a') | s, a\right]$$

**Intuition:** if the optimal state-action values for the next time-step Q*(s',a') are known, then the optimal strategy is to take the action that maximizes the expected value of $r + \gamma Q^*(s', a')$

The optimal policy π* corresponds to taking the best action in any state as specified by Q*

ITMOre than a
UNIVERSITY

# Solving for the optimal policy

**Value iteration** algorithm: Use Bellman equation as an iterative update

$$Q_{i+1}(s, a) = \mathbb{E}\left[r + \gamma \max_{a'} Q_i(s', a')|s, a\right]$$

$Q_i$ will converge to $Q^*$ as i -> infinity

What's the problem with this?
Not scalable. Must compute Q(s,a) for every state-action pair. If state is e.g. current game state pixels, computationally infeasible to compute for entire state space!

Solution: use a function approximator to estimate Q(s,a). E.g. a neural network!

IT'sMO$re$ $than$ $a$
UNIVERSITY

# Solving for the optimal policy: Q-learning

Q-learning: Use a function approximator to estimate the action-value function

$$Q(s, a; \theta) \approx Q^*(s, a)$$

function parameters (weights)

If the function approximator is a deep neural network => **deep q-learning**!

ITsMOre than a
UNIVERSITY

# Solving for the optimal policy: Q-learning

Remember: want to find a Q-function that satisfies the Bellman Equation:

$$Q^*(s, a) = \mathbb{E}_{s' \sim \mathcal{E}} \left[ r + \gamma \max_{a'} Q^*(s', a') | s, a \right]$$

Forward Pass

Loss function: $L_i(\theta_i) = \mathbb{E}_{s, a \sim \rho(\cdot)} \left[ (y_i - Q(s, a; \theta_i))^2 \right]$

where $y_i = \mathbb{E}_{s' \sim \mathcal{E}} \left[ r + \gamma \max_{a'} Q(s', a'; \theta_{i-1}) | s, a \right]$

Backward Pass

Gradient update (with respect to Q-function parameters θ):

$$\nabla_{\theta_i} L_i(\theta_i) = \mathbb{E}_{s, a \sim \rho(\cdot); s' \sim \mathcal{E}} \left[ r + \gamma \max_{a'} Q(s', a'; \theta_{i-1}) - Q(s, a; \theta_i)) \nabla_{\theta_i} Q(s, a; \theta_i) \right]$$

ITₛMOre than a
UNIVERSITY

# Case Study: Playing Atari Games



**Objective**: Complete the game with the highest score

*[Mnih et al. NIPS Workshop 2013; Nature 2015]*

**State:** Raw pixel inputs of the game state
**Action:** Game controls e.g. Left, Right, Up, Down
**Reward:** Score increase/decrease at each time step

# Q-network Architecture

ITMO UNIVERSITY

$Q(s, a; \theta):$
neural network
with weights $\theta$



FC-4 (Q-values)

FC-256

32 4x4 conv, stride 2

16 8x8 conv, stride 4

**Current state $s_t$: 84x84x4 stack of last 4 frames**
(after RGB->grayscale conversion, downsampling, and cropping)

ITsMO*re than a*
UNIVERSITY

# Q-network Architecture

ITMO UNIVERSITY

$Q(s, a; \theta)$ :
neural network
with weights $\theta$

| FC-4 (Q-values) |
| --- |

| FC-256 |
| --- |

| 32 4x4 conv, stride 2 |
| --- |

| 16 8x8 conv, stride 4 |
| --- |



← Input: state $s_t$

**Current state $s_t$: 84x84x4 stack of last 4 frames**
(after RGB->grayscale conversion, downsampling, and cropping)

IT₅MOre than a
UNIVERSITY

# Q-network Architecture

ITMO UNIVERSITY

$Q(s, a; \theta)$ :
neural network
with weights $\theta$

FC-4 (Q-values)

FC-256

32 4x4 conv, stride 2

16 8x8 conv, stride 4

← Familiar conv layers,
FC layer
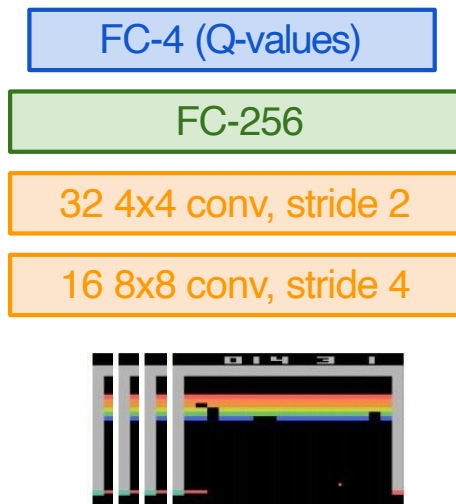
**Current state $s_t$: 84x84x4 stack of last 4 frames**
(after RGB->grayscale conversion, downsampling, and cropping)

ITsMOre than a
UNIVERSITY

# Q-network Architecture

ITMO UNIVERSITY

$Q(s, a; \theta)$:
neural network
with weights $\theta$

FC-4 (Q-values)

FC-256

32 4x4 conv, stride 2

16 8x8 conv, stride 4



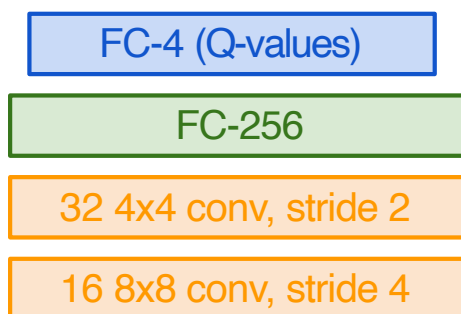**Last FC layer has 4-d output (if 4 actions), corresponding to $Q(s_t, a_1)$, $Q(s_t, a_2)$, $Q(s_t, a_3)$, $Q(s_t, a_4)$**

**Current state $s_t$: 84x84x4 stack of last 4 frames**
(after RGB->grayscale conversion, downsampling, and cropping)

IT₅MOre than a
UNIVERSITY

# Q-network Architecture

ITMO UNIVERSITY

$Q(s, a; \theta)$ :
neural network
with weights $\theta$

FC-4 (Q-values)

FC-256

32 4x4 conv, stride 2

16 8x8 conv, stride 4

Last FC layer has 4-d output (if 4 actions), corresponding to $Q(s_t, a_1)$, $Q(s_t, a_2)$, $Q(s_t, a_3)$, $Q(s_t, a_4)$

Number of actions between 4-20 depending on Atari game

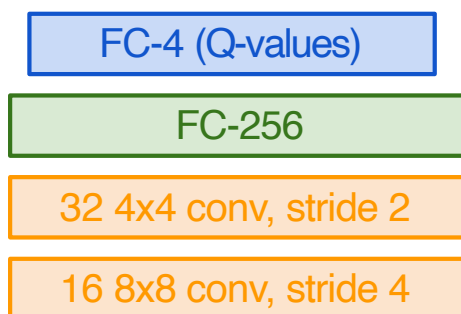**Current state $s_t$: 84x84x4 stack of last 4 frames**
(after RGB->grayscale conversion, downsampling, and cropping)

IT$_s$MO$_{re\ than\ a}$ UNIVERSITY

# Q-network Architecture

ITMO UNIVERSITY

$Q(s, a; \theta)$ :
neural network
with weights $\theta$

A single feedforward pass
to compute Q-values for all
actions from the current
state => efficient!

FC-4 (Q-values)

FC-256

32 4x4 conv, stride 2

16 8x8 conv, stride 4

Last FC layer has 4-d
output (if 4 actions),
corresponding to $Q(s_t, a_1)$, $Q(s_t, a_2)$, $Q(s_t, a_3)$,
$Q(s_t, a_4)$

Number of actions between 4-20
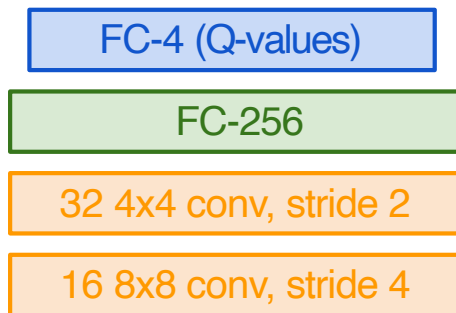depending on Atari game

**Current state $s_t$: 84x84x4 stack of last 4 frames**
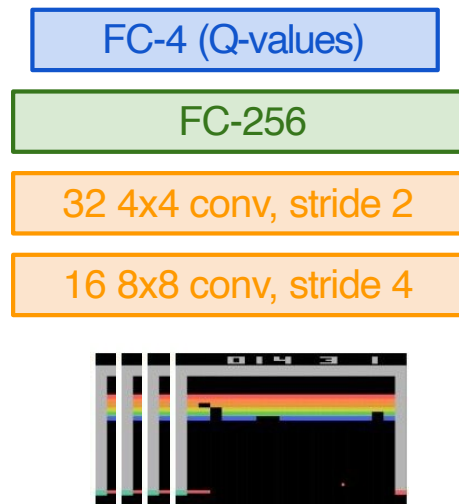(after RGB->grayscale conversion, downsampling, and cropping)

ITsMOre than a
UNIVERSITY

# Solving for the optimal policy: Q-learning

Remember: want to find a Q-function that satisfies the Bellman Equation:

$$Q^*(s,a) = \mathbb{E}_{s' \sim \mathcal{E}} \left[ r + \gamma \max_{a'} Q^*(s',a')|s,a \right]$$

**Forward Pass**

Loss function: $\quad L_i(\theta_i) = \mathbb{E}_{s,a \sim \rho(\cdot)} \left[ (y_i - Q(s,a;\theta_i))^2 \right]$

where $\quad y_i = \mathbb{E}_{s' \sim \mathcal{E}} \left[ r + \gamma \max_{a'} Q(s',a';\theta_{i-1})|s,a \right]$

**Backward Pass**

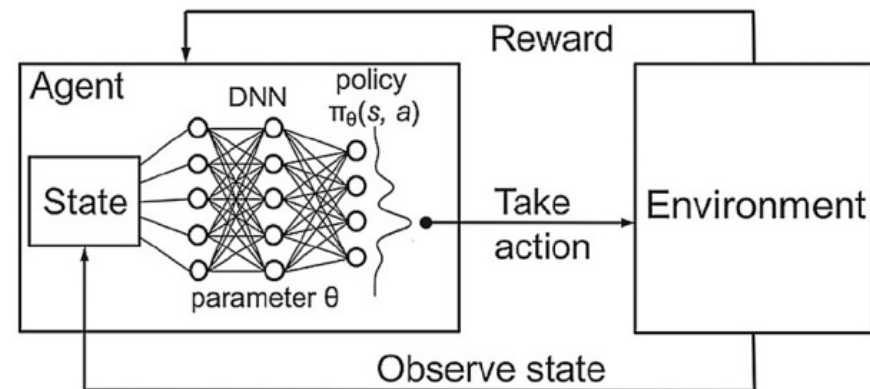Gradient update (with respect to Q-function parameters θ):

$$\nabla_{\theta_i} L_i(\theta_i) = \mathbb{E}_{s,a \sim \rho(\cdot); s' \sim \mathcal{E}} \left[ r + \gamma \max_{a'} Q(s',a';\theta_{i-1}) - Q(s,a;\theta_i)) \nabla_{\theta_i} Q(s,a;\theta_i) \right]$$

# Q-learning and Deep Q-learning



**a** Q-learning

**b** Deep Q-learning

# Training the Q-network: Experience Replay

Learning from batches of consecutive samples is problematic:
- Samples are correlated => inefficient learning
- Current Q-network parameters determines next training samples (e.g. if maximizing action is to move left, training samples will be dominated by samples from left-hand size) => can lead to bad feedback loops

Address these problems using **experience replay**
- Continually update a **replay memory** table of transitions ($s_t$, $a_t$, $r_t$, $s_{t+1}$) as game (experience) episodes are played
- Train Q-network on random minibatches of transitions from the replay memory, instead of consecutive samples

Each transition can also contribute to multiple weight updates => greater data efficiency

# Putting it together: Deep Q-Learning with Experience Replay

**Algorithm 1** Deep Q-learning with Experience Replay

Initialize replay memory $\mathcal{D}$ to capacity $N$
Initialize action-value function $Q$ with random weights
**for** episode $= 1, M$ **do**
    Initialise state $s_t$
    **for** $t = 1, T$ **do**
        With probability $\epsilon$ select a random action $a_t$
        otherwise select $a_t = \max_a Q^*(s_t, a; \theta)$
        Execute action $a_t$ and observe reward $r_t$ and state $s_{t+1}$
        Store transition $(s_t, a_t, r_t, s_{t+1})$ in $\mathcal{D}$
        Set $s_{t+1} = s_t$
        Sample random minibatch of transitions $(s_t, a_t, r_t, s_{t+1})$ from $\mathcal{D}$
        Set $y_j = \begin{cases} r_j & \text{for terminal } s_{t+1} \\ r_j + \gamma \max_{a'} Q(s_{t+1}, a'; \theta) & \text{for non-terminal } s_{t+1} \end{cases}$
        Perform a gradient descent step on $(y_j - Q(s_t, a_j; \theta))^2$
    **end for**
**end for**

# Putting it together: Deep Q-Learning with Experience Replay

**Algorithm 1** Deep Q-learning with Experience Replay

Initialize replay memory $\mathcal{D}$ to capacity $N$ — Initialize replay memory, Q-network
Initialize action-value function $Q$ with random weights
**for** episode $= 1, M$ **do**
$\quad$ Initialise state $s_t$
$\quad$ **for** $t = 1, T$ **do**
$\quad\quad$ With probability $\epsilon$ select a random action $a_t$
$\quad\quad$ otherwise select $a_t = \max_a Q^*(s_t, a; \theta)$
$\quad\quad$ Execute action $a_t$ and observe reward $r_t$ and state $s_{t+1}$
$\quad\quad$ Store transition $(s_t, a_t, r_t, s_{t+1})$ in $\mathcal{D}$
$\quad\quad$ Set $s_{t+1} = s_t$
$\quad\quad$ Sample random minibatch of transitions $(s_t, a_t, r_t, s_{t+1})$ from $\mathcal{D}$
$\quad\quad$ Set $y_j = \begin{cases} r_j & \text{for terminal } s_{t+1} \\ r_j + \gamma \max_{a'} Q(s_{t+1}, a'; \theta) & \text{for non-terminal } s_{t+1} \end{cases}$
$\quad\quad$ Perform a gradient descent step on $(y_j - Q(s_t, a_j; \theta))^2$
$\quad$ **end for**
**end for**

**Algorithm 1** Deep Q-learning with Experience Replay

Initialize replay memory $\mathcal{D}$ to capacity $N$
Initialize action-value function $Q$ with random weights
**for** episode $= 1, M$ **do**  ← Play M episodes (full games)
    Initialise state $s_t$
    **for** $t = 1, T$ **do**
        With probability $\epsilon$ select a random action $a_t$
        otherwise select $a_t = \max_a Q^*(s_t, a; \theta)$
        Execute action $a_t$ and observe reward $r_t$ and state $s_{t+1}$
        Store transition $(s_t, a_t, r_t, s_{t+1})$ in $\mathcal{D}$
        Set $s_{t+1} = s_t$
        Sample random minibatch of transitions $(s_t, a_t, r_t, s_{t+1})$ from $\mathcal{D}$
        Set $y_j = \begin{cases} r_j & \text{for terminal } s_{t+1} \\ r_j + \gamma \max_{a'} Q(s_{t+1}, a'; \theta) & \text{for non-terminal } s_{t+1} \end{cases}$
        Perform a gradient descent step on $(y_j - Q(s_t, a_j; \theta))^2$
    **end for**
**end for**

# Putting it together: Deep Q-Learning with Experience Replay

**Algorithm 1** Deep Q-learning with Experience Replay

Initialize replay memory $\mathcal{D}$ to capacity $N$
Initialize action-value function $Q$ with random weights
**for** episode $= 1, M$ **do**
    Initialise state $s_t$
    **for** $t = 1, T$ **do**
        With probability $\epsilon$ select a random action $a_t$
        otherwise select $a_t = \max_a Q^*(s_t, a; \theta)$
        Execute action $a_t$ and observe reward $r_t$ and state $s_{t+1}$
        Store transition $(s_t, a_t, r_t, s_{t+1})$ in $\mathcal{D}$
        Set $s_{t+1} = s_t$
        Sample random minibatch of transitions $(s_t, a_t, r_t, s_{t+1})$ from $\mathcal{D}$
        Set $y_j = \begin{cases} r_j & \text{for terminal } s_{t+1} \\ r_j + \gamma \max_{a'} Q(s_{t+1}, a'; \theta) & \text{for non-terminal } s_{t+1} \end{cases}$
        Perform a gradient descent step on $(y_j - Q(s_t, a_j; \theta))^2$
    **end for**
**end for**

Initialize state (starting game screen pixels) at the beginning of each episode

# Putting it together: Deep Q-Learning with Experience Replay

**Algorithm 1** Deep Q-learning with Experience Replay

Initialize replay memory $\mathcal{D}$ to capacity $N$
Initialize action-value function $Q$ with random weights
**for** episode $= 1, M$ **do**
    Initialise state $s_t$
    **for** $t = 1, T$ **do**
        With probability $\epsilon$ select a random action $a_t$
        otherwise select $a_t = \max_a Q^*(s_t, a; \theta)$
        Execute action $a_t$ and observe reward $r_t$ and state $s_{t+1}$
        Store transition $(s_t, a_t, r_t, s_{t+1})$ in $\mathcal{D}$
        Set $s_{t+1} = s_t$
        Sample random minibatch of transitions $(s_t, a_t, r_t, s_{t+1})$ from $\mathcal{D}$
        Set $y_j = \begin{cases} r_j & \text{for terminal } s_{t+1} \\ r_j + \gamma \max_{a'} Q(s_{t+1}, a'; \theta) & \text{for non-terminal } s_{t+1} \end{cases}$
        Perform a gradient descent step on $(y_j - Q(s_t, a_j; \theta))^2$
    **end for**
**end for**

For each timestep t
of the game

# Putting it together: Deep Q-Learning with Experience Replay

---

**Algorithm 1** Deep Q-learning with Experience Replay

Initialize replay memory $\mathcal{D}$ to capacity $N$
Initialize action-value function $Q$ with random weights
**for** episode $= 1, M$ **do**
    Initialise state $s_t$
    **for** $t = 1, T$ **do**
        With probability $\epsilon$ select a random action $a_t$
        otherwise select $a_t = \max_a Q^*(s_t, a; \theta)$
        Execute action $a_t$ and observe reward $r_t$ and state $s_{t+1}$
        Store transition $(s_t, a_t, r_t, s_{t+1})$ in $\mathcal{D}$
        Set $s_{t+1} = s_t$
        Sample random minibatch of transitions $(s_t, a_t, r_t, s_{t+1})$ from $\mathcal{D}$
        Set $y_j = \begin{cases} r_j & \text{for terminal } s_{t+1} \\ r_j + \gamma \max_{a'} Q(s_{t+1}, a'; \theta) & \text{for non-terminal } s_{t+1} \end{cases}$
        Perform a gradient descent step on $(y_j - Q(s_t, a_j; \theta))^2$
    **end for**
**end for**

With small probability, select a random action (explore), otherwise select greedy action from current policy

# Putting it together: Deep Q-Learning with Experience Replay

---

**Algorithm 1** Deep Q-learning with Experience Replay

Initialize replay memory $\mathcal{D}$ to capacity $N$
Initialize action-value function $Q$ with random weights
**for** episode $= 1, M$ **do**
    Initialise state $s_t$
    **for** $t = 1, T$ **do**
        With probability $\epsilon$ select a random action $a_t$
        otherwise select $a_t = \max_a Q^*(s_t, a; \theta)$
        Execute action $a_t$ and observe reward $r_t$ and state $s_{t+1}$
        Store transition $(s_t, a_t, r_t, s_{t+1})$ in $\mathcal{D}$
        Set $s_{t+1} = s_t$
        Sample random minibatch of transitions $(s_t, a_t, r_t, s_{t+1})$ from $\mathcal{D}$
        Set $y_j = \begin{cases} r_j & \text{for terminal } s_{t+1} \\ r_j + \gamma \max_{a'} Q(s_{t+1}, a'; \theta) & \text{for non-terminal } s_{t+1} \end{cases}$
        Perform a gradient descent step on $(y_j - Q(s_t, a_j; \theta))^2$
    **end for**
**end for**

Take the action ($a_t$), and observe the reward $r_t$ and next state $s_{t+1}$

# Putting it together: Deep Q-Learning with Experience Replay

**Algorithm 1** Deep Q-learning with Experience Replay

Initialize replay memory $\mathcal{D}$ to capacity $N$
Initialize action-value function $Q$ with random weights
**for** episode $= 1, M$ **do**
    Initialise state $s_t$
    **for** $t = 1, T$ **do**
        With probability $\epsilon$ select a random action $a_t$
        otherwise select $a_t = \max_a Q^*(s_t, a; \theta)$
        Execute action $a_t$ and observe reward $r_t$ and state $s_{t+1}$
        Store transition $(s_t, a_t, r_t, s_{t+1})$ in $\mathcal{D}$
        Set $s_{t+1} = s_t$
        Sample random minibatch of transitions $(s_t, a_t, r_t, s_{t+1})$ from $\mathcal{D}$
        Set $y_j = \begin{cases} r_j & \text{for terminal } s_{t+1} \\ r_j + \gamma \max_{a'} Q(s_{t+1}, a'; \theta) & \text{for non-terminal } s_{t+1} \end{cases}$
        Perform a gradient descent step on $(y_j - Q(s_t, a_j; \theta))^2$
    **end for**
**end for**

Store transition in replay memory

# Putting it together: Deep Q-Learning with Experience Replay

---

**Algorithm 1** Deep Q-learning with Experience Replay

Initialize replay memory $\mathcal{D}$ to capacity $N$
Initialize action-value function $Q$ with random weights
**for** episode $= 1, M$ **do**
    Initialise state $s_t$
    **for** $t = 1, T$ **do**
        With probability $\epsilon$ select a random action $a_t$
        otherwise select $a_t = \max_a Q^*(s_t, a; \theta)$
        Execute action $a_t$ and observe reward $r_t$ and state $s_{t+1}$
        Store transition $(s_t, a_t, r_t, s_{t+1})$ in $\mathcal{D}$
        Set $s_{t+1} = s_t$
        Sample random minibatch of transitions $(s_t, a_t, r_t, s_{t+1})$ from $\mathcal{D}$
        Set $y_j = \begin{cases} r_j & \text{for terminal } s_{t+1} \\ r_j + \gamma \max_{a'} Q(s_{t+1}, a'; \theta) & \text{for non-terminal } s_{t+1} \end{cases}$
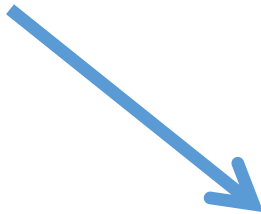        Perform a gradient descent step on $(y_j - Q(s_t, a_j; \theta))^2$
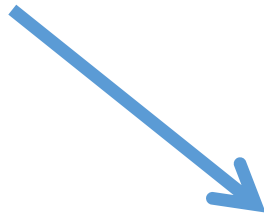    **end for**
**end for**

Experience Replay: Sample a random minibatch of transitions from replay memory and perform a gradient descent step

**Environments for experiments!**

https://www.gymlibrary.dev/environments/atari/breakout/

**Example!**

https://www.youtube.com/watch?v=V1eYniJ0Rnk

# Policy Gradients

What is a problem with Q-learning?
The Q-function can be very complicated!

Example: a robot grasping an object has a very high-dimensional state => hard
to learn exact value of every (state, action) pair

# Policy Gradients

What is a problem with Q-learning?
The Q-function can be very complicated!

Example: a robot grasping an object has a very high-dimensional state => hard to learn exact value of every (state, action) pair

But the policy can be much simpler: just close your hand
Can we learn a policy directly, e.g. finding the best policy from a collection of policies?

# Policy Gradients

Formally, let's define a class of parametrized policies: $\Pi = \{\pi_\theta, \theta \in \mathbb{R}^m\}$

For each policy, define its value:

$$J(\theta) = \mathbb{E}\left[\sum_{t \geq 0} \gamma^t r_t | \pi_\theta\right]$$

# Policy Gradients

Formally, let's define a class of parametrized policies: $\Pi = \{\pi_\theta, \theta \in \mathbb{R}^m\}$

For each policy, define its value:

$$J(\theta) = \mathbb{E}\left[\sum_{t \geq 0} \gamma^t r_t | \pi_\theta\right]$$

We want to find the optimal policy $\theta^* = \arg\max_\theta J(\theta)$

How can we do this?

# Policy Gradients

Formally, let's define a class of parametrized policies: $\Pi = \{\pi_\theta, \theta \in \mathbb{R}^m\}$

For each policy, define its value:

$$J(\theta) = \mathbb{E}\left[\sum_{t \geq 0} \gamma^t r_t | \pi_\theta\right]$$

We want to find the optimal policy $\theta^* = \arg\max_\theta J(\theta)$

How can we do this?

Gradient ascent on policy parameters!

# REINFORCE algorithm

Mathematically, we can write:

$$J(\theta) = \mathbb{E}_{\tau \sim p(\tau;\theta)} \left[ r(\tau) \right]$$

$$= \int_\tau r(\tau) p(\tau; \theta) \mathrm{d}\tau$$

Where r(т) is the reward of a trajectory $\tau = (s_0, a_0, r_0, s_1, \ldots)$

ITsMOre than a
UNIVERSITY

# REINFORCE algorithm

Expected reward:

$$J(\theta) = \mathbb{E}_{\tau \sim p(\tau;\theta)}\left[r(\tau)\right]$$

$$= \int_{\tau} r(\tau)p(\tau;\theta)\mathrm{d}\tau$$

Now let's differentiate this: $\quad \nabla_{\theta}J(\theta) = \int_{\tau} r(\tau)\nabla_{\theta}p(\tau;\theta)\mathrm{d}\tau$

Intractable! Gradient of an expectation is problematic when p depends on θ

Can be estimated with Monte Carlo sampling!
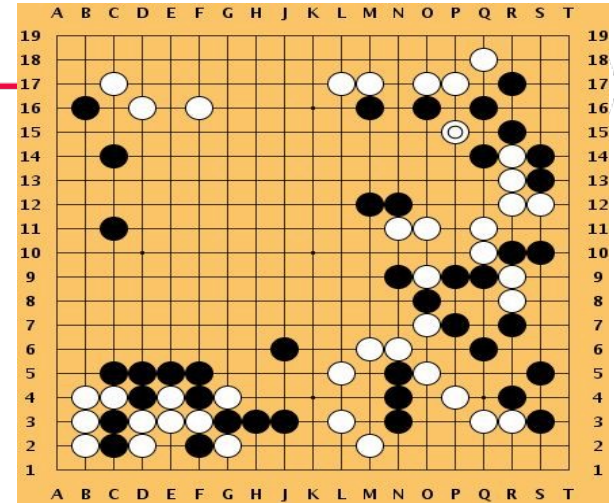
IT₃MOre than a
UNIVERSITY

# More policy gradients: AlphaGo



**AlphaGo - overview:**
- Mix of supervised learning and reinforcement learning
- Mix of old methods (Monte Carlo Tree Search) and recent ones (deep RL)
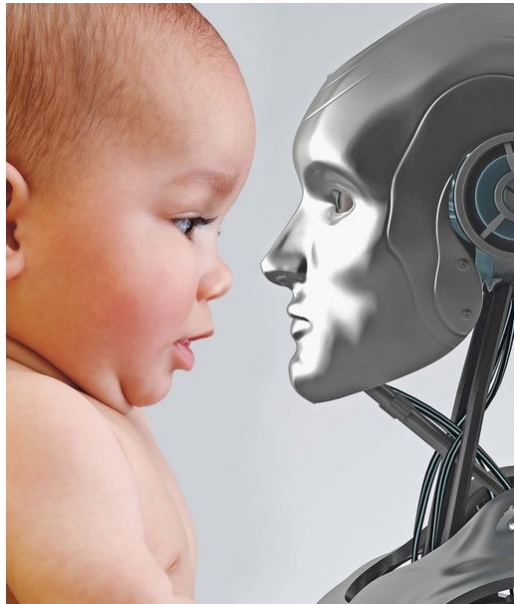
**How to beat the Go world champion:**
- Featurize the board (stone color, move legality, bias, …)
- Initialize policy network with supervised training from professional go games, then continue training using policy gradient (play against itself from random previous iterations, +1 / -1 reward for winning / losing)
- Also learn value network (critic)
- Finally, combine policy and value networks in a Monte Carlo Tree Search algorithm to select actions by lookahead search

*[Silver et al., Nature 2016]*

# Deb Roy's MIT experiment



**90 000 hours video**
**140 000 hours audio**



ITMO UNIVERSITY

ITsMOre than a UNIVERSITY

http://ai-news.ru/2018/04/kak_obuchautsya_deti_i_pochemu_iskusstvennyj_intellekt_tak_ne_mozhet.html