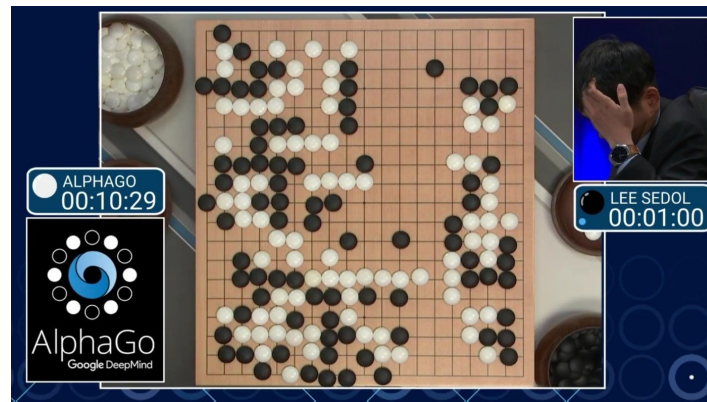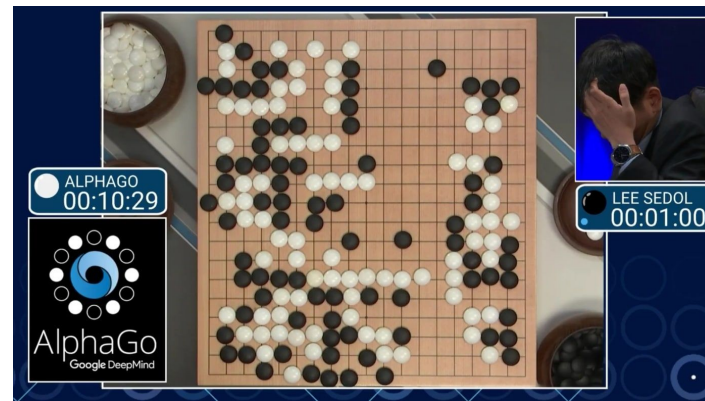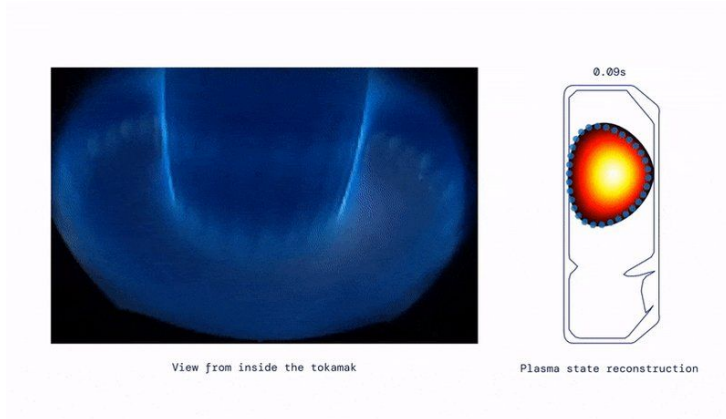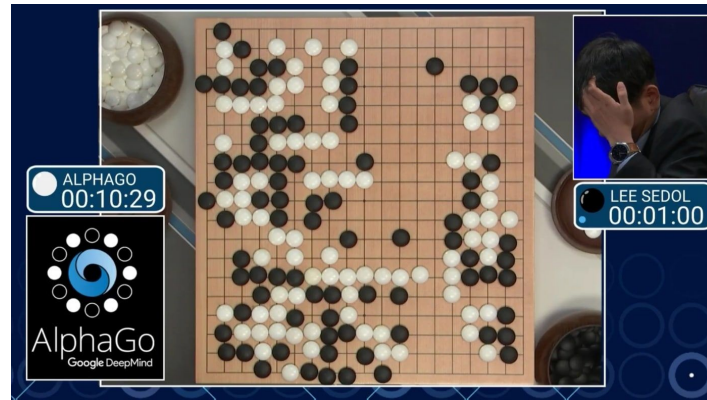# A Brief Introduction To

# Reinforcement Learning

Thomas Moerland

View from inside the tokamak

Plasma state reconstruction

0.09s

The international journal of science / 6 October 2022

**nature**

**MATRIX GAMES**
Deep reinforcement learning opens route to faster algorithms for matrix multiplication

**Protecting Peru**
Can technology help Indigenous groups preserve the Amazon?

**Invisible touch**
How marine clouds are affected by aerosols emitted from shipping

**Preferential practice**
US universities favour prestige in faculty hiring and retention



ALPHAGO
00:10:29

AlphaGo
Google DeepMind

LEE SEDOL
00:01:00



0.09s

View from inside the tokamak

Plasma state reconstruction



**a** Chemical representation of the synthesis plan
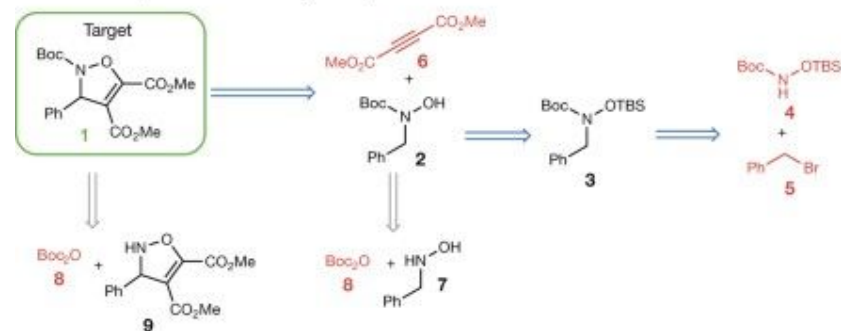
View from inside the tokamak

0.09s
Plasma state reconstruction

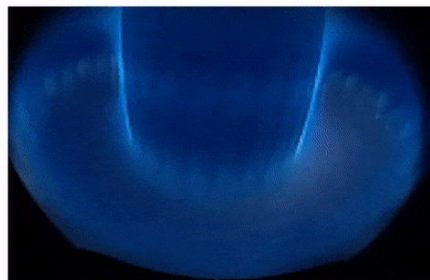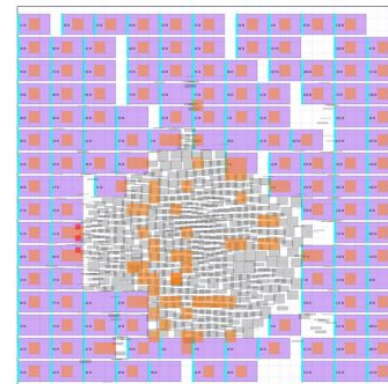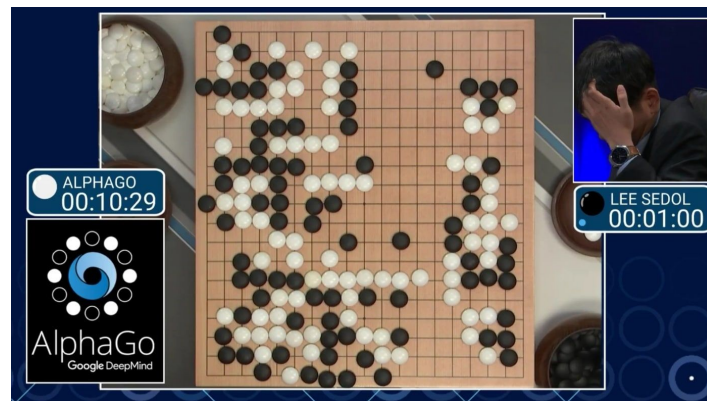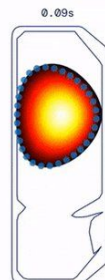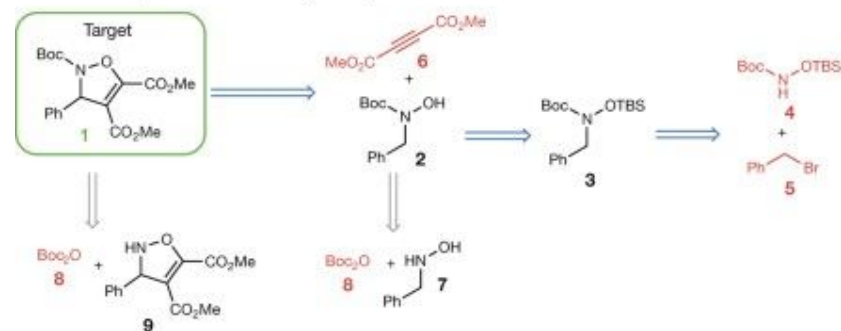a Chemical representation of the synthesis plan

View from inside the tokamak

Plasma state reconstruction

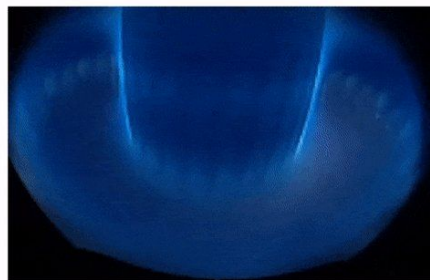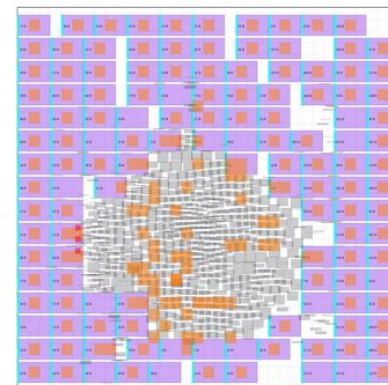**a** Chemical representation of the synthesis plan

Silver, David, et al. "Mastering the game of Go with deep neural networks and tree search." *Nature* 529.7587 (2016): 484-489.

Fawzi, Alhussein, et al. "Discovering faster matrix multiplication algorithms with reinforcement learning." *Nature* 610.7930 (2022): 47-53.

Degrave, Jonas, et al. "Magnetic control of tokamak plasmas through deep reinforcement learning." *Nature* 602.7897 (2022): 414-419.

Segler, Marwin HS et al. "Planning chemical syntheses with deep neural networks and symbolic AI." *Nature* 555.7698 (2018): 604-610.

Mirhoseini, Azalia, et al. "A graph placement methodology for fast chip design." *Nature* 594.7862 (2021): 207-212.

Many (real-world) problems can be formulated as a

**sequential decision-making problem**

Many (real-world) problems can be formulated as a

**sequential decision-making problem**

which may be solved through <u>reinforcement learning</u>.

# Content

# Part I

# Introduction

# Biology

# Biology

# Biology



*Skinner box*



*B.F. Skinner (1904 – 1990)*

# Biology



*Skinner box*



*B.F. Skinner (1904 – 1990)*

*Instrumental conditioning*:
Learning behaviour based on reward and punishment (trial and error)

# Biology



*Skinner box*



*B.F. Skinner (1904 – 1990)*

*Instrumental conditioning*:
Learning behaviour based on reward and punishment (trial and error)

RL is the computational specification of this idea

# Supervised versus reinforcement learning

# Supervised versus reinforcement learning

|  | **Supervised learning** | **Reinforcement learning** |
|---|---|---|
| Dataset |  |  |
| Feedback |  |  |

# Supervised versus reinforcement learning

| | Supervised learning | Reinforcement learning |
| --- | --- | --- |
| <u>Dataset</u> | Given | |
| <u>Feedback</u> | | |

# Supervised versus reinforcement learning

|  | Supervised learning | Reinforcement learning |
|---|---|---|
| <u>Dataset</u> | Given | Active collection |
| <u>Feedback</u> |  |  |

# Supervised versus reinforcement learning

| | Supervised learning | Reinforcement learning |
|---|---|---|
| Dataset | Given | Active collection |
| Feedback | Full<br><br>(x with correct y) | |

# Supervised versus reinforcement learning

|  | **Supervised learning** | **Reinforcement learning** |
|---|---|---|
| <u>Dataset</u> | Given | Active collection |
| <u>Feedback</u> | Full | Partial |
|  | (x with correct y) | (~~state with correct action~~) (feedback on some outcomes) |

# Benefits of Reinforcement Learning

# Benefits of Reinforcement Learning



**Autonomous behaviour/learning**
(only specify goals)

# Benefits of Reinforcement Learning



**Autonomous behaviour/learning**
(only specify goals)



**Solve tasks that you can't label**
(only need to label the outcome)

# Benefits of Reinforcement Learning



**Autonomous behaviour/learning**
(only specify goals)



**Solve tasks that you can't label**
(only need to label the outcome)



**Outperform human solution**
(only need to label the outcome)

# Part II

# Problem Formulation

# Agent-Environment loop

# Agent-Environment loop

# Agent-Environment loop



action  $a$

# Agent-Environment loop



action   $a$

# Agent-Environment loop

state
$s$

Agent

action $a$

Environment

# Agent-Environment loop



Agent

Environment

state
*s*

reward
*r*

action   *a*

# Agent-Environment loop



state
$s$

action $a$

reward
$r$

can be negative

Agent

Environment

# Agent-Environment loop
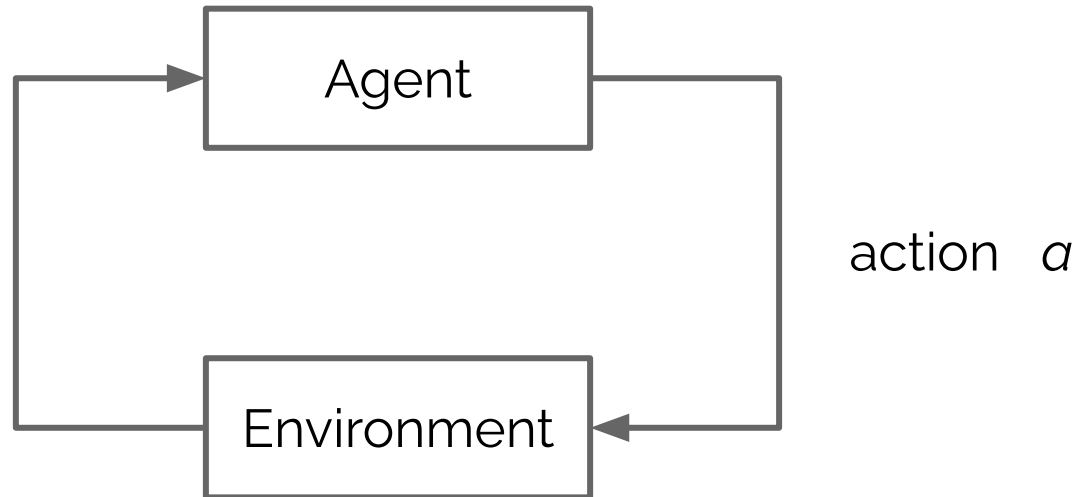
# Agent-Environment loop



Agent

Environment

state $s$

reward $r$

action $a$

# Agent-Environment loop

# Agent-Environment loop

state
*s*

reward
*r*

Find an action selection strategy...
( *policy* π(a|s) )

Agent

action *a*

Environment

... that gets as much reward as possible!

# Reward

# Reward

Immediate reward

$$r_t$$

# Reward

~~Immediate reward~~

Cumulative reward

$$r_t + r_{t+1} + r_{t+2} + \ldots$$

# Reward

~~Immediate reward~~

~~Cumulative reward~~

Expected cumulative reward

$$\mathbb{E}\big[r_t + r_{t+1} + r_{t+2} + \ldots \quad \big]$$

# Reward

Immediate reward

Cumulative reward

Expected cumulative reward

$$\mathbb{E}\big[r_t + r_{t+1} + r_{t+2} + \ldots \qquad \big]$$

Average over stochasticity in 1) environment and 2) own policy.

# Reward

Immediate reward

Cumulative reward

Expected cumulative reward

= Value

$$Q^{\pi}(s, a) = \mathbb{E}\left[r_t + r_{t+1} + r_{t+2} + ... \middle| s_t = s, a_t = a\right]$$

# Reward

~~Immediate reward~~

~~Cumulative reward~~

Expected cumulative reward

= Value

$$Q^{\pi}(s, a) = \mathbb{E}\big[r_t + r_{t+1} + r_{t+2} + ... \big| s_t = s, a_t = a\big]$$

**Q-value**: total reward we get on average after taking action a in state s.

# Reward

~~Immediate reward~~

~~Cumulative reward~~

Expected cumulative reward

= Value

$$Q^{\pi}(s, a) = \mathbb{E}[r_t + r_{t+1} + r_{t+2} + ... | s_t = s, a_t = a]$$

**Q-value**: total reward we get on average after taking action a in state s.
- Depends on our own future behaviour π (if we act stupid, reward will be low)

# Reward

~~Immediate reward~~

~~Cumulative reward~~

Expected cumulative reward

= Value

$$Q^{\pi}(s, a) = \mathbb{E}[r_t + r_{t+1} + r_{t+2} + \ldots | s_t = s, a_t = a]$$

**Q-value**: total reward we get on average after taking action a in state s.
- Depends on our own future behaviour π (if we act stupid, reward will be low)

Can show each state-action has one optimal value, denoted by Q*(s,a).
- These are the quantities we want to know!

# Illustration: Optimal Value

# Illustration: Optimal Value

# Illustration: Optimal Value



**Q\* = ?**

Home
Go Out
To University
r = +2.0  Bar
Study  r = -1.0
Go Out
Take Exam
r = +2.0  Bar
Pass  r = +10.0

**Question**:     What is Q*(Home, Go Out)?

# Illustration: Optimal Value



**Question**: What is Q*(Home, Go Out)?

**Answer**: 2.0

# Illustration: Optimal Value



**Question**:    What is Q*(Home, To University)?

# Illustration: Optimal Value



**Q\* = 2.0**     Go Out          To University    **Q\* = 9.0**

Home

Go Out

To University

r = +2.0   Bar

Study   r = -1.0

Go Out

Take Exam

r = +2.0   Bar

Pass   r = +10.0

**Question**:      What is Q\*(Home, To University)?

**Answer**:        -1.0 + 10.0 = 9.0

# Illustration: Optimal Value



**Q* = 2.0**     Go Out     To University     **Q* = 9.0**

Home

r = +2.0  Bar

Study  r = -1.0

Go Out     Take Exam

r = +2.0  Bar

Pass  r = +10.0

**Question**:     What should you do at Home?

# Illustration: Optimal Value



**Question**: What should you do at Home?

**Answer**: Come to University!

# Illustration: Optimal Value



**Q\* = 2.0**  Go Out

**Q\* = 9.0**  To University

Once we know the optimal values we also know how to act optimally

Home

r = +2.0   Bar

Study   r = -1.0

Go Out

Take Exam

r = +2.0   Bar

Pass   r = +10.0

**Question**:      What should you do at Home?

**Answer**:       Come to University!

# Part III


# The Reinforcement Learning Cycle

# Challenge

# Challenge

**Problem:** In practice we don't know the problem structure and optimal Q-values.

# Challenge

**Problem:** In practice we don't know the problem structure and optimal Q-values.
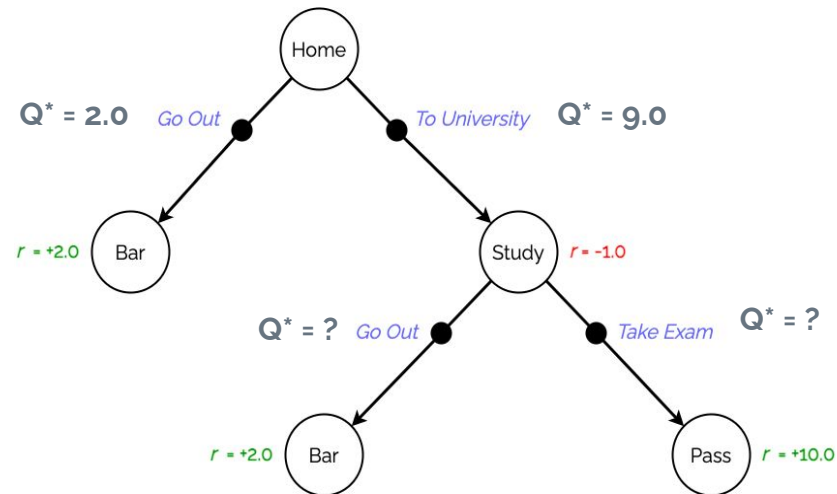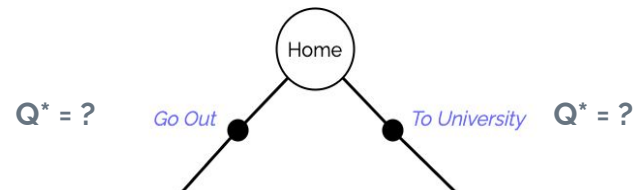
# Challenge

**Problem:** In practice we don't know the problem structure and optimal Q-values.

# Challenge

**Problem:** In practice we don't know the problem structure and optimal Q-values.



**Solution**: Learn through trial and error.

# The Reinforcement Learning Cycle

# The Reinforcement Learning Cycle

*Pseudocode*

# The Reinforcement Learning Cycle

*Pseudocode*

Initialize **Q(s,a)** solution estimates for all states and actions (e.g. to 0)

# The Reinforcement Learning Cycle

*Pseudocode*

Initialize **Q(s,a)** solution estimates for all states and actions (e.g. to 0)

Repeat:

# The Reinforcement Learning Cycle

*Pseudocode*

Initialize **Q(s,a)** solution estimates for all states and actions (e.g. to 0)

Repeat:

1) Exploration: Sample a sequence of actions.

# The Reinforcement Learning Cycle

*Pseudocode*

Initialize **Q(s,a)** solution estimates for all states and actions (e.g. to 0)

Repeat:

   1) Exploration: Sample a sequence of actions.

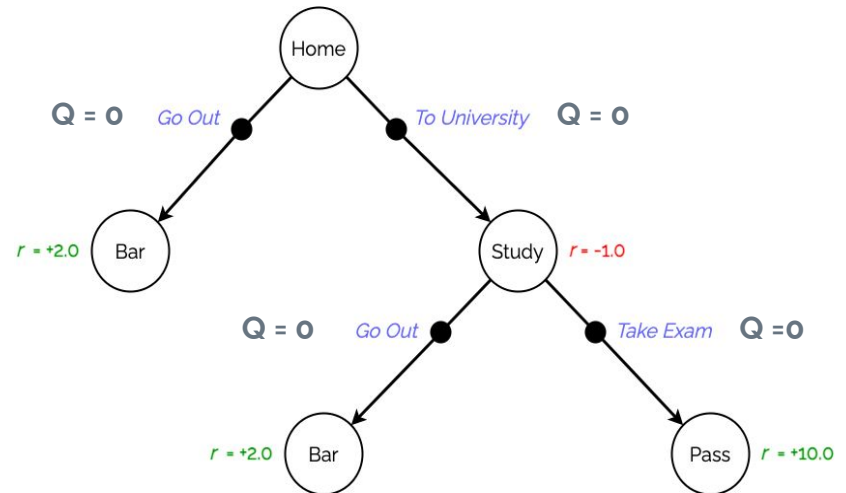   2) Credit assignment: Compute new value estimates **Q^{back-up}(s,a)** for all actions along the path.

# The Reinforcement Learning Cycle

*Pseudocode*

Initialize **Q(s,a)** solution estimates for all states and actions (e.g. to 0)

Repeat:

1) Exploration: Sample a sequence of actions.

2) Credit assignment: Compute new value estimates $Q^{back-up}(s,a)$ for all actions along the path.

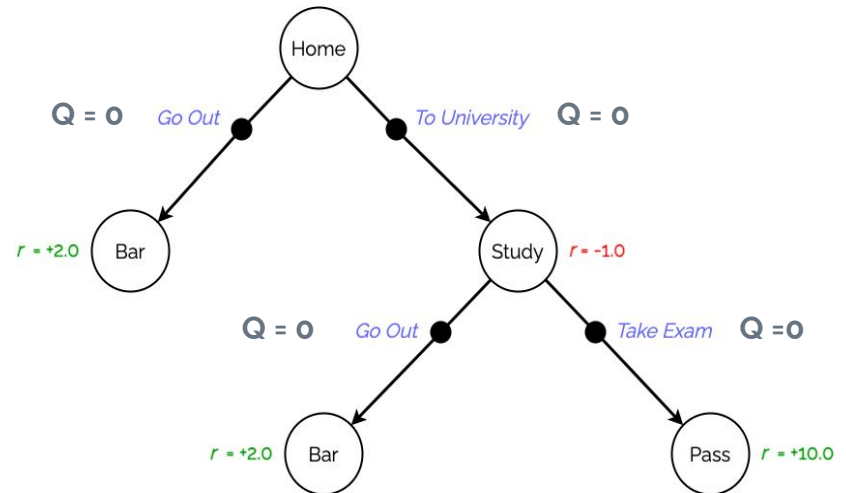3) Learning update: Adjust our **Q(s,a)** solution based on the back-up estimates $Q^{back-up}(s,a)$.
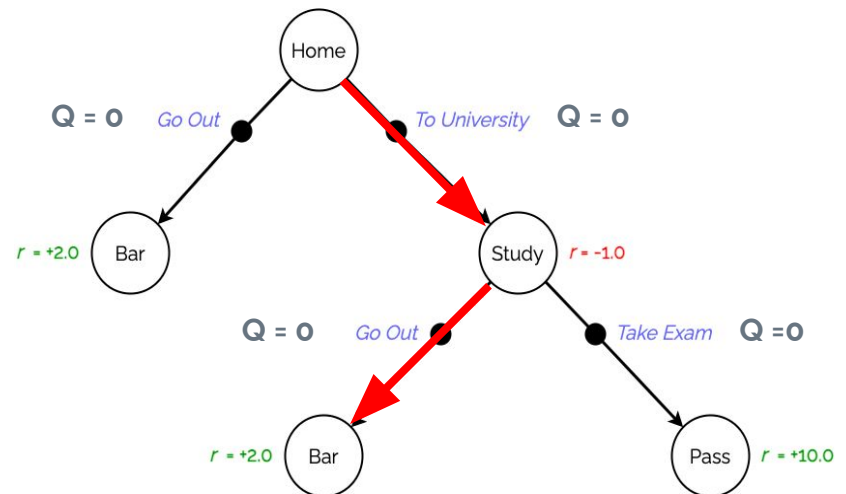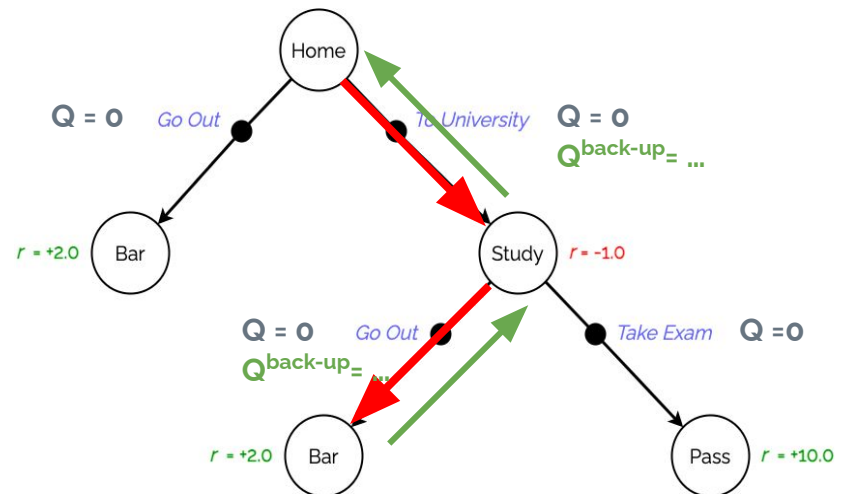
# The Reinforcement Learning Cycle

*Pseudocode*

Initialize **Q(s,a)** solution estimates for all states and actions (e.g. to 0)

Repeat:

   1) **<u>Exploration</u>**: Sample a sequence of actions.

   2) <u>Credit assignment</u>: Compute new value estimates **Q<sup>back-up</sup>(s,a)** for all actions along the path.

   3) <u>Learning update</u>: Adjust our **Q(s,a)** solution based on the back-up estimates **Q<sup>back-up</sup>(s,a)**.
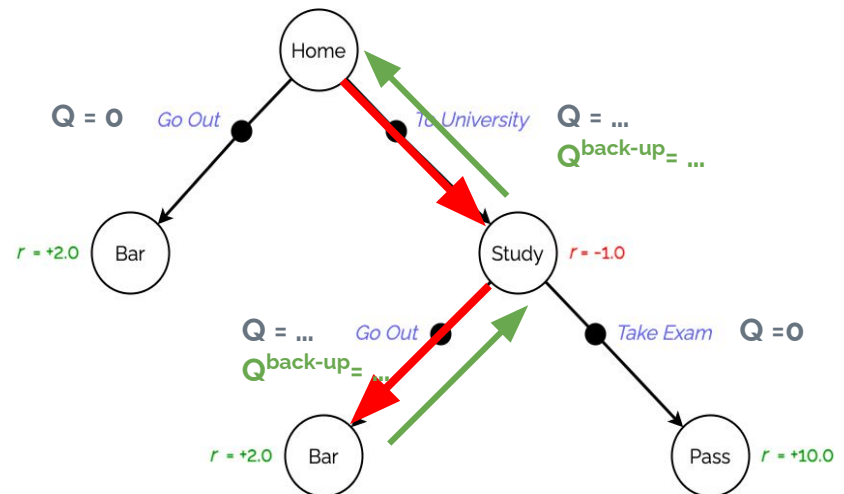
# The Reinforcement Learning Cycle

*Pseudocode*

Initialize **Q(s,a)** solution estimates for all states and actions (e.g. to 0)

Repeat:

1) <u>Exploration</u>: Sample a sequence of actions.

2) **Credit assignment**: Compute new value estimates **Q$^{back-up}$(s,a)** for all actions along the path.

3) <u>Learning update</u>: Adjust our **Q(s,a)** solution based on the back-up estimates **Q$^{back-up}$(s,a)**.



**Q = 0**   *Go Out*      *To University*   **Q = ...**
**Q$^{back-up}$= ...**

*r = +2.0*   Bar                    Study   *r = -1.0*

Home

**Q = ...**   *Go Out*      *Take Exam*   **Q =0**

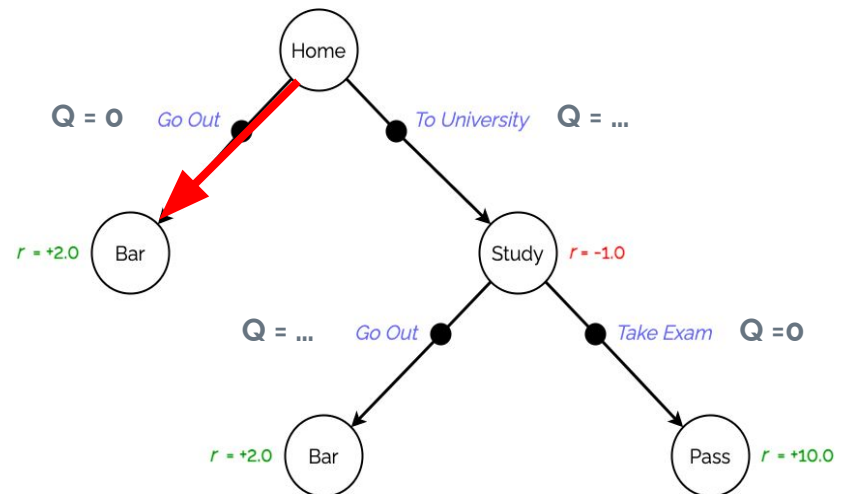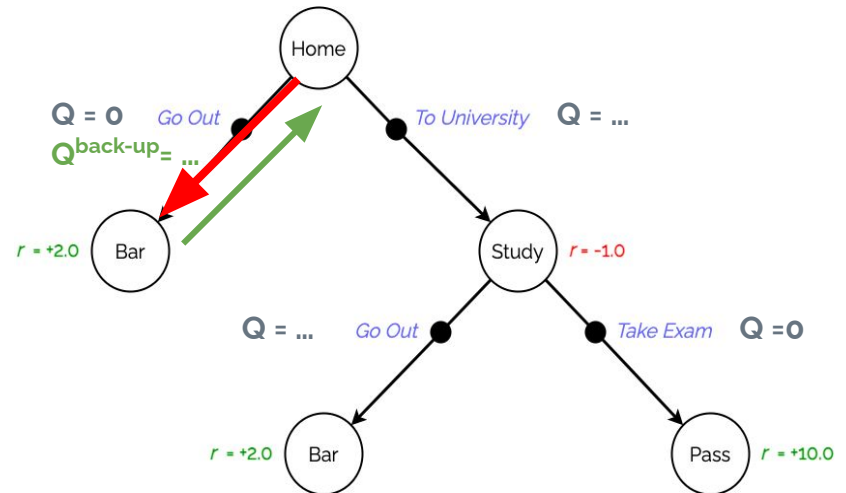*r = +2.0*   Bar                    Pass   *r = +10.0*

# The Reinforcement Learning Cycle

*Pseudocode*

Initialize **Q(s,a)** solution estimates for all states and actions (e.g. to 0)

Repeat:

  1) Exploration: Sample a sequence of actions.

  2) Credit assignment: Compute new value estimates $Q^{back-up}(s,a)$ for all actions along the path.

  3) **Learning update**: Adjust our **Q(s,a)** solution based on the back-up estimates $Q^{back-up}(s,a)$.



Home

**Q = ...**  *Go Out*      *To University*  **Q = ...**

$Q^{back-up}=$ ...

$r$ = +2.0   Bar      Study   $r$ = -1.0

**Q = ...**  *Go Out*     *Take Exam*  **Q =0**

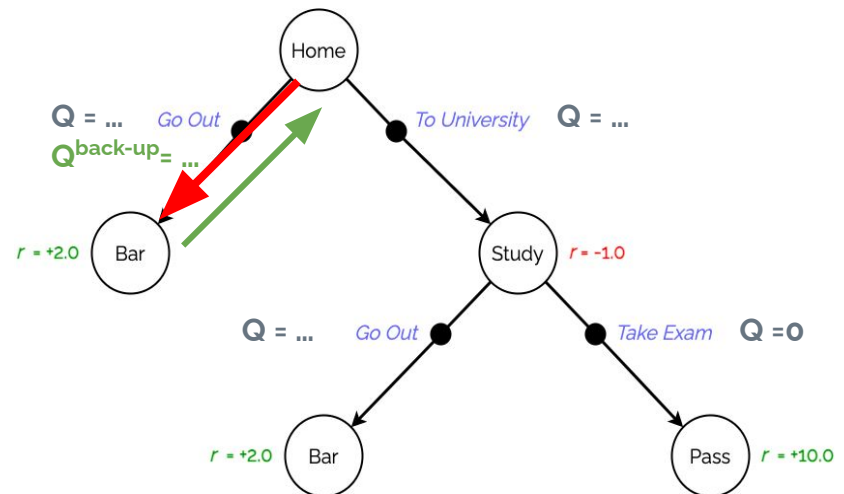$r$ = +2.0   Bar      Pass   $r$ = +10.0

# The Reinforcement Learning Cycle

*Pseudocode*

Initialize **Q(s,a)** solution estimates for all states and actions (e.g. to 0)

Repeat:

1) <u>Exploration</u>: Sample a sequence of actions.

2) <u>Credit assignment</u>: Compute new value estimates **Q$^{\text{back-up}}$(s,a)** for all actions along the path.

3) **Learning update**: Adjust our **Q(s,a)** solution based on the back-up estimates **Q$^{\text{back-up}}$(s,a)**.

etc.

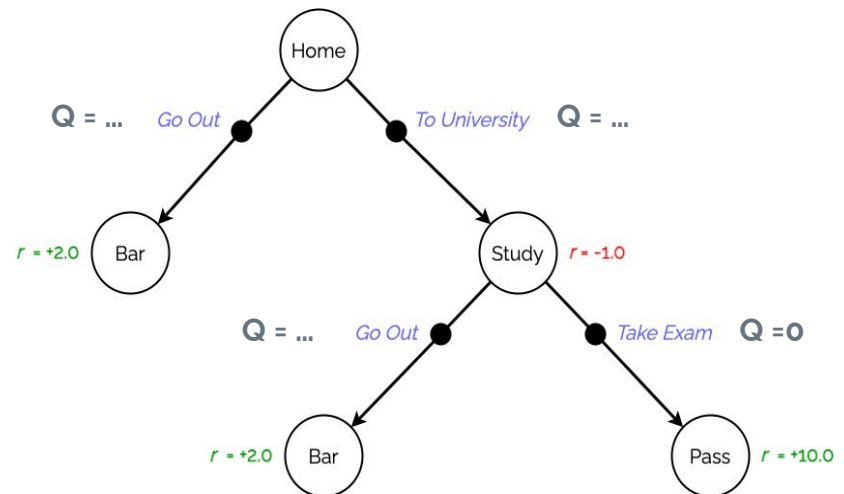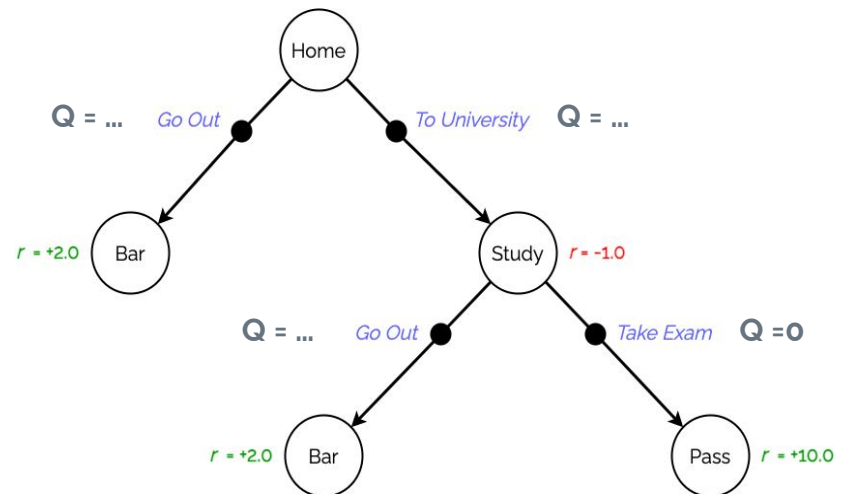# The Reinforcement Learning Cycle

*Pseudocode*

Initialize **Q(s,a)** solution estimates for all states and actions (e.g. to 0)

Repeat:

1) <u>Exploration</u>: Sample a sequence of actions.

2) <u>Credit assignment</u>: Compute new value estimates $Q^{back-up}(s,a)$ for all actions along the path.

3) **Learning update**: Adjust our **Q(s,a)** solution based on the back-up estimates $Q^{back-up}(s,a)$.

Home

**Q = ...** *Go Out*    *To University* **Q = ...**

*r* = +2.0  Bar    Study  *r* = -1.0

**Q = ...** *Go Out*    *Take Exam* **Q =0**

*r* = +2.0  Bar    Pass  *r* = +10.0

<u>We will discuss all three steps, but in reverse order</u>

# Part III A


# Learning Update

# Learning update (tabular)

# Learning update (tabular)

(Really a supervised learning topic, we will briefly discuss this in one slide.)

# Learning update (tabular)

$$Q(s, a) \leftarrow Q(s, a) + \eta \cdot \left( Q^{\text{back-up}}(s, a) - Q(s, a) \right)$$

# Learning update (tabular)

$$Q(s,a) \leftarrow Q(s,a) + \boxed{\eta} \cdot \left( Q^{\text{back-up}}(s,a) - Q(s,a) \right)$$

learning rate: $\quad \eta \in [0, 1]$

# Learning update (tabular)

$$Q(s, a) \leftarrow Q(s, a) + \eta \cdot \left( Q^{\text{back-up}}(s, a) - Q(s, a) \right)$$

To update our solution…

# Learning update (tabular)

$$Q(s, a) \leftarrow \boxed{Q(s, a)} + \eta \cdot \left( Q^{\text{back-up}}(s, a) - Q(s, a) \right)$$

To update our solution we take the current solution…

# Learning update (tabular)

$$Q(s, a) \leftarrow Q(s, a) \boxed{+ \eta \cdot} \left( Q^{\text{back-up}}(s, a) - Q(s, a) \right)$$

To update our solution we take the current solution and move it a (small) step...

# Learning update (tabular)

$$Q(s, a) \leftarrow Q(s, a) + \eta \cdot \boxed{\left( Q^{\text{back-up}}(s, a) - Q(s, a) \right)}$$

To update our solution we take the current solution and move it a (small) step…
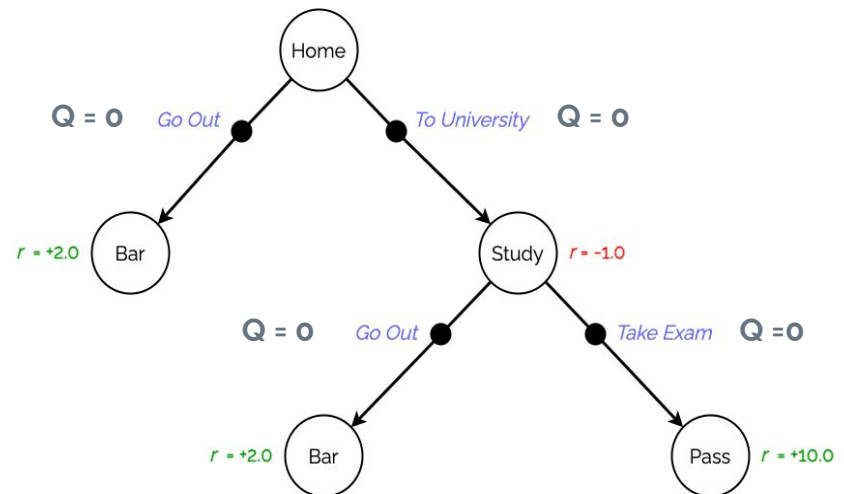
…in the direction of the back-up estimate.

# Learning update (tabular)

'training target'

$$Q(s, a) \leftarrow Q(s, a) + \eta \cdot \left( Q^{\text{back-up}}(s, a) - Q(s, a) \right)$$

To update our solution we take the current solution and move it a (small) step…

…in the direction of the back-up estimate.

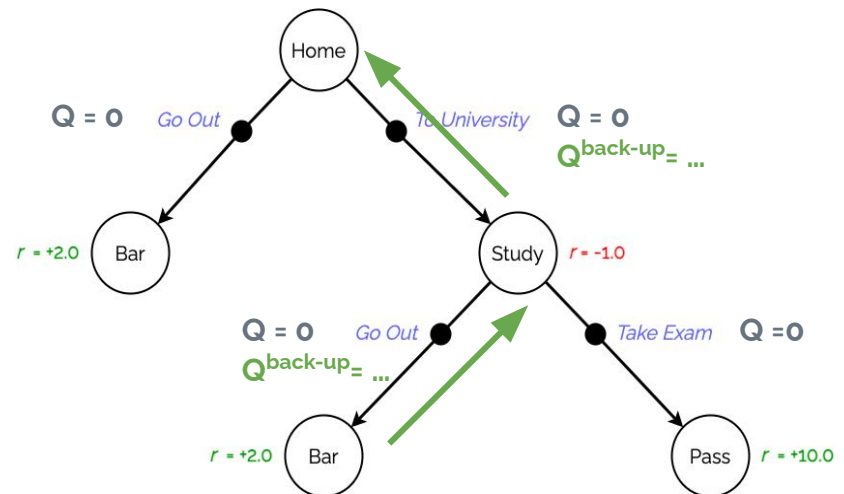# The Reinforcement Learning Cycle

*Pseudocode*

Initialize **Q(s,a)** solution estimates for all states and actions (e.g. to 0)

Repeat:

1) <span style="color:red">Exploration</span>: Sample a sequence of actions.

2) <span style="color:green">Credit assignment</span>: Compute new value estimates **Q^back-up(s,a)** for all actions along the path.

3) Learning update: Adjust our **Q(s,a)** solution based on the back-up estimates **Q^back-up(s,a)**.

**Q = 0**  *Go Out*     *To University*  **Q = 0**

Home

*r = +2.0*  Bar     Study  *r = -1.0*

**Q = 0**  *Go Out*     *Take Exam*  **Q = 0**

*r = +2.0*  Bar     Pass  *r = +10.0*

# The Reinforcement Learning Cycle

*Pseudocode*

Initialize **Q(s,a)** solution estimates for all states and actions (e.g. to 0)

Repeat:

1) <u>Exploration</u>: Sample a sequence of actions.

2) **Credit assignment**: Compute new value estimates **Q$^{back-up}$(s,a)** for all actions along the path.

3) <u>Learning update</u>: Adjust our **Q(s,a)** solution based on the back-up estimates **Q$^{back-up}$(s,a)**.
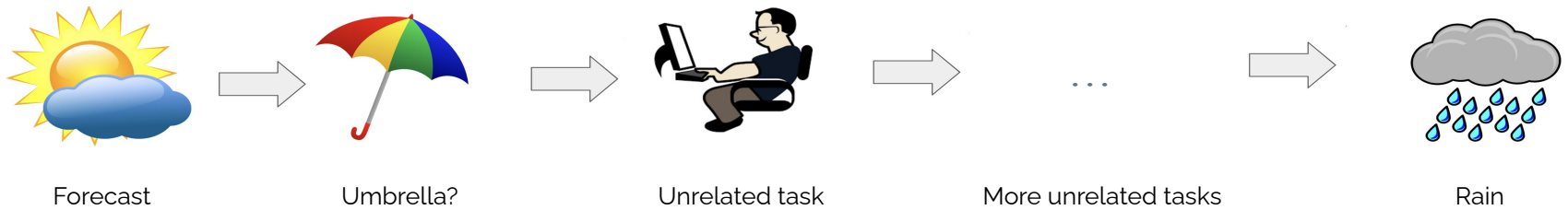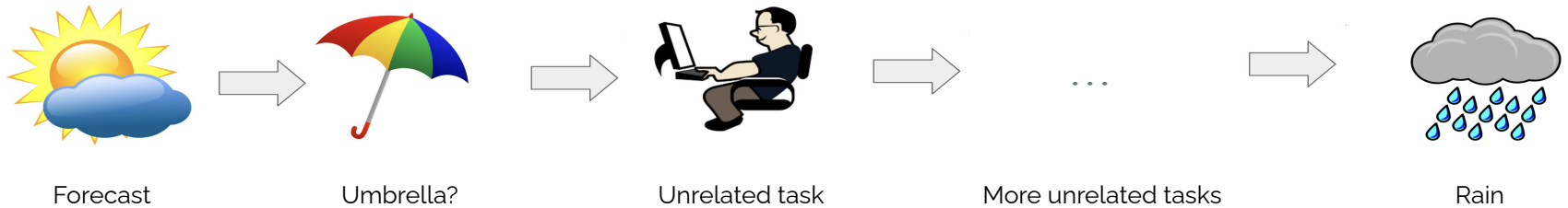
# Part III B

# Credit Assignment

# Credit assignment

# Credit assignment

(Also a topic in neural network training, but then the objective is differentiable)

# Credit assignment



Forecast       Umbrella?       Unrelated task       More unrelated tasks       Rain

# Credit assignment



Forecast      Umbrella?      Unrelated task      More unrelated tasks      Rain

**Question**: You get the reward (not soaked), but which of your previous actions deserve credit?

# Credit assignment

You think this is easy, but humans actually also struggle:

# Credit assignment

You think this is easy, but humans actually also struggle:

# Credit assignment

You think this is easy, but humans actually also struggle:

# Credit assignment

You think this is easy, but humans actually also struggle:

# Credit assignment

You think this is easy, but humans actually also struggle:

# Credit assignment

You think this is easy, but humans actually also struggle:

# Credit assignment

You think this is easy, but humans actually also struggle:

# Credit assignment

You think this is easy, but humans actually also struggle:



Superstition
=
failed credit assignment

# Credit assignment

# Credit assignment

# Credit assignment



align bottles
(r= –1)

$S_{t-3}$

$S_{t-2}$

$S_{t-1}$

$S_t$

$S_{t+1}$

Time

Nadal wins the tennis game
(r = +100)

# Credit assignment



He needs to determine how much credit each previous chosen action gets

align bottles
(r= –1)

$S_{t-3}$  $S_{t-2}$  $S_{t-1}$  $S_t$  $S_{t+1}$

Time

Nadal wins the tennis game
(r = +100)

# Credit assignment

One extreme: each action along the way gets full credit

align bottles
(r= –1)

$S_{t-3}$

$S_{t-2}$

$S_{t-1}$

$S_t$

$S_{t+1}$

Time

Nadal wins the tennis game
(r = +100)

# Credit assignment

One extreme: each action along the way gets full credit



align bottles
(r= –1)

Time

$S_{t-3}$ $S_{t-2}$ $S_{t-1}$ $S_t$ $S_{t+1}$

Nadal wins the tennis game
(r = +100)

$$Q^{back-up}(s_t, a_t) = \sum_{i=0}^{\infty} r_{t+i}$$

# Credit assignment

One extreme: each action along the way gets full credit

align bottles
(r= –1)

$Q^{back-up} = 99$

$S_{t-3}$

$S_{t-2}$

$S_{t-1}$

$S_t$

$S_{t+1}$

Time

$Q^{back-up} = 100$

Nadal wins the tennis game
(r = +100)

# Credit assignment

One extreme: each action along the way gets full credit



$S_{t-3}$

align bottles
(r= –1)

$Q^{back-up}$ = 99

$S_{t-2}$

$S_{t-1}$

$S_t$

Time

$Q^{back-up}$ = 100

$S_{t+1}$

Nadal wins the tennis game
(r = +100)

Monte Carlo back-up

+   Fast propagation.
-   High variance (action may seem better or worse than it really is)

# Credit assignment

Other extreme: only last action gets credit (for now).



align bottles
(r= –1)

Time

Nadal wins the tennis game
(r = +100)

# Credit assignment

Other extreme: only last action gets credit (for now).



align bottles
(r= –1)

Time

Nadal wins the tennis game
(r = +100)

$$Q^{back-up}(s_t, a_t) = r_t + Q(s_{t+1}, a_{t+1})$$

# Credit assignment

Other extreme: only last action gets credit (for now).



$S_{t-3}$

align bottles
(r= –1)

$S_{t-2}$

$Q^{back-up}$ = -1

$S_{t-1}$

$S_t$

Time

$Q^{back-up}$ = 100

$S_{t+1}$

Nadal wins the tennis game
(r = +100)

# Credit assignment

Other extreme: only last action gets credit (for now).



align bottles
(r= –1)

$Q^{back-up}$ = –1

$S_{t-3}$  $S_{t-2}$  $S_{t-1}$  $S_t$  $S_{t+1}$

Time

$Q^{back-up}$ = 100

Nadal wins the tennis game
(r = +100)

# Credit assignment

Other extreme: only last action gets credit (for now).



$Q^{back\text{-}up}$ = 100

align bottles
(r= -1)

$Q^{back\text{-}up}$ = -1

$S_{t-3}$

$S_{t-2}$

$S_{t-1}$

$S_t$

Time

$Q^{back\text{-}up}$ = 100

$S_{t+1}$

Nadal wins the tennis game
(r = +100)

# Credit assignment

Other extreme: only last action gets credit (for now).



$Q^{back-up}$ = 100

align bottles
(r= –1)

$Q^{back-up}$ = –1

$S_{t-3}$

$S_{t-2}$

$S_{t-1}$

$S_t$

Time

$Q^{back-up}$ = 100

$S_{t+1}$

Nadal wins the tennis game
(r = +100)

One-step (temporal difference) back-up

+ Low variance.
- Slow propagation.

# Credit assignment

# Credit assignment

Spectrum of back-up estimators

# Credit assignment

Spectrum of back-up estimators

# The Reinforcement Learning Cycle

*Pseudocode*

Initialize **Q(s,a)** estimates for all states,actions (e.g. to 0)

Repeat:

   1) <u>Exploration</u>: Sample a sequence of actions.

   2) <u>Credit assignment</u>: Compute new value estimates $Q^{back-up}(s,a)$ for all actions along the path.

   3) <u>Update</u>: Adjust our **Q(s,a)** solution based on the new back-up estimates $Q^{back-up}(s,a)$.

# The Reinforcement Learning Cycle

*Pseudocode*

Initialize **Q(s,a)** estimates for all states,actions (e.g. to 0)

Repeat:

1) **<u>Exploration</u>**: Sample a sequence of actions.

2) <u>Credit assignment</u>: Compute new value estimates **Q<sup>back-up</sup>(s,a)** for all actions along the path.

3) <u>Update</u>: Adjust our **Q(s,a)** solution based on the new back-up estimates **Q<sup>back-up</sup>(s,a)**.

# Part III C


# Exploration

# Exploration versus exploitation

# Exploration versus exploitation

# Exploration versus exploitation



**Question**: The usual place scores a 7/10 on average. Which place would you choose?

# Exploration versus exploitation



$Q = 7$    $Q = ?$

$a_1$    $a_2$

**Question**: The usual place scores a 7/10 on average. Which place would you choose?

# Exploration versus exploitation

Exploitation

(commit to the current best option)

Exploration

(try something which is new or – currently – seems suboptimal)

# Exploration versus exploitation



Exploitation

(commit to the current best option)

Exploration

(try something which is new or – currently – seems suboptimal)

We actually need to balance both

# Exploration/Exploitation trade-off

We need **exploration** because actions may look worse than they are.

# Exploration/Exploitation trade-off

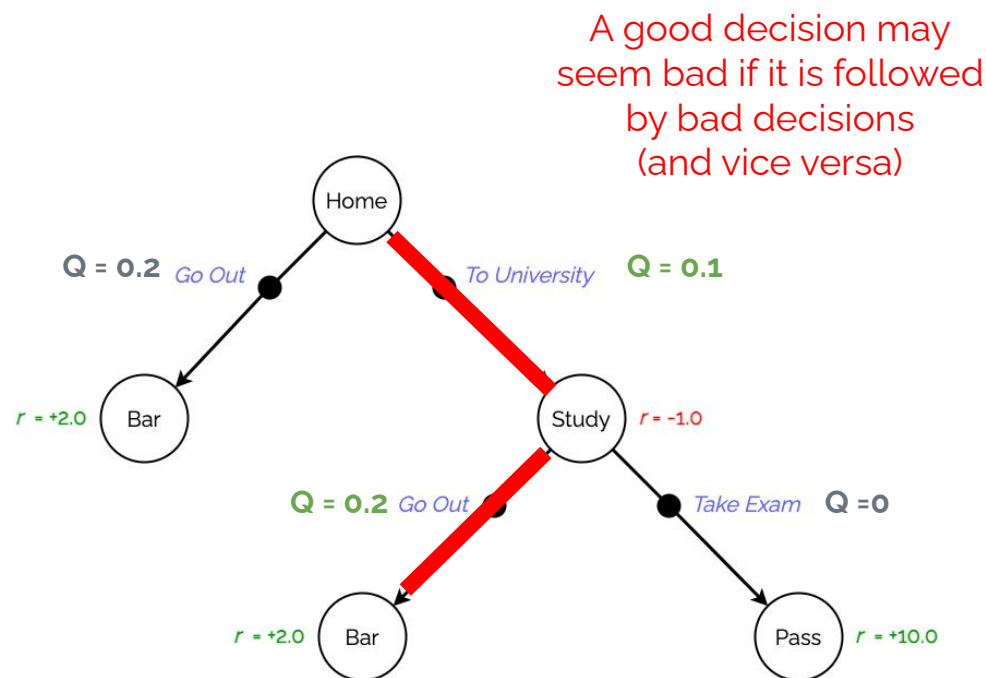We need **exploration** because actions may look worse than they are.

Reasons:

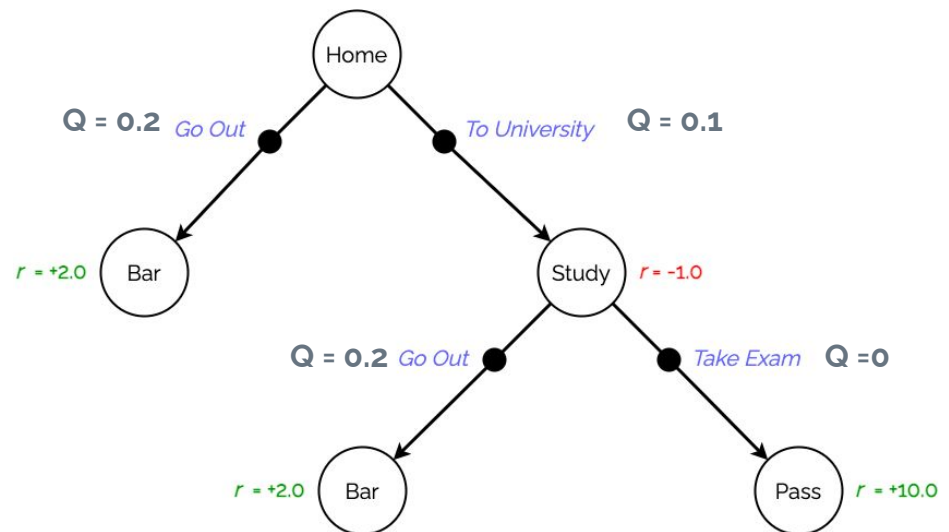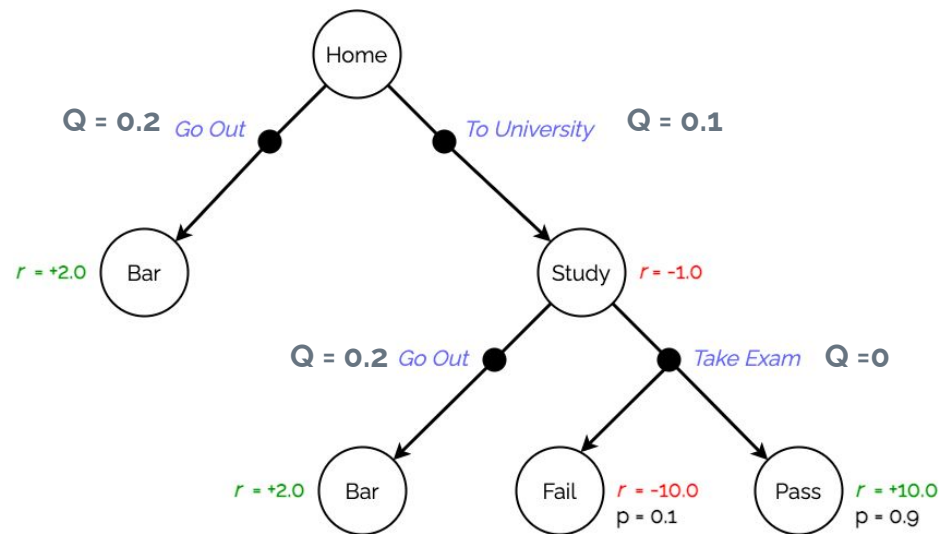# Exploration/Exploitation trade-off

We need **exploration** because actions may look worse than they are.

<u>Reasons</u>:

We will use Monte Carlo back-ups and a learning rate of 0.1

# Exploration/Exploitation trade-off

We need **exploration** because actions may look worse than they are.

Reasons:

# Exploration/Exploitation trade-off

We need **exploration** because actions may look worse than they are.

Reasons:

# Exploration/Exploitation trade-off

We need **exploration** because actions may look worse than they are.

Reasons:



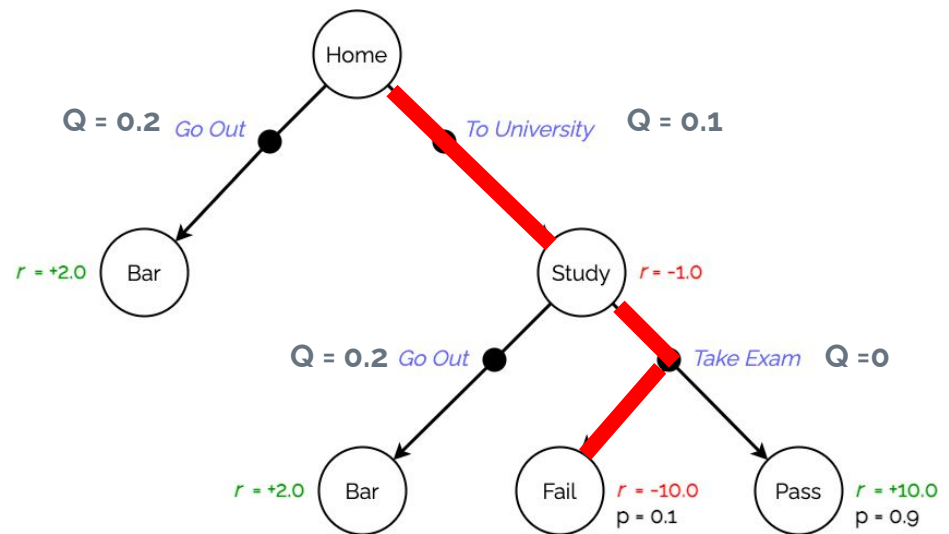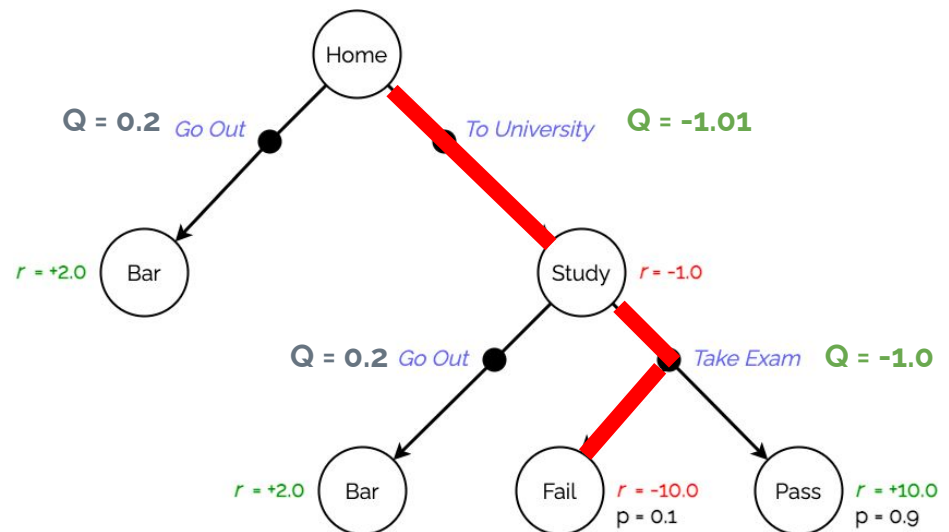If you never go to university, you will never find out the pay-off

# Exploration/Exploitation trade-off

We need **exploration** because actions may look worse than they are.

Reasons:

1. We need to collect our own data

# Exploration/Exploitation trade-off

We need **exploration** because actions may look worse than they are.

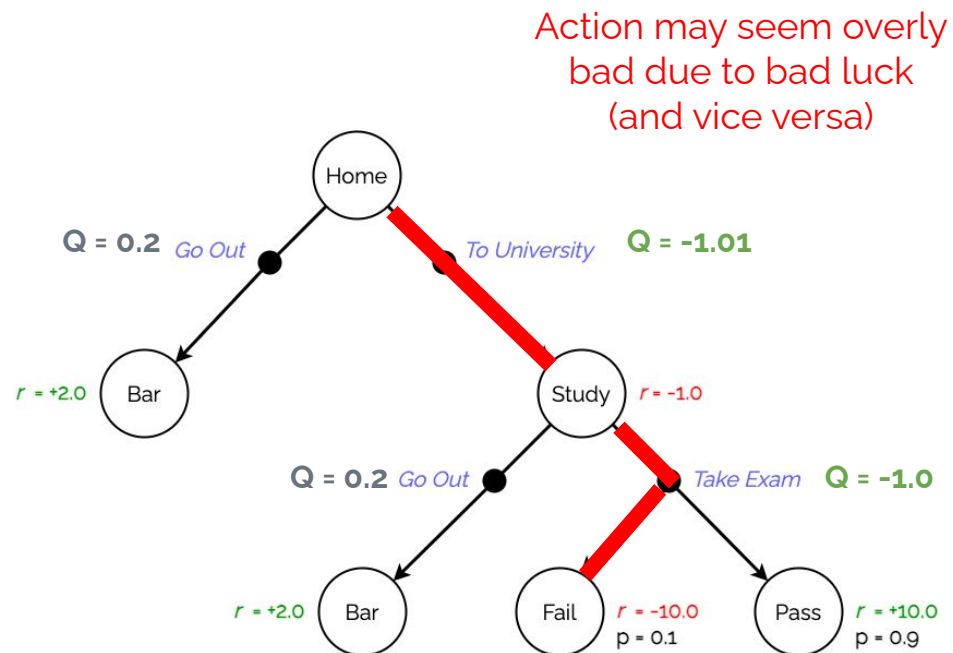Reasons:
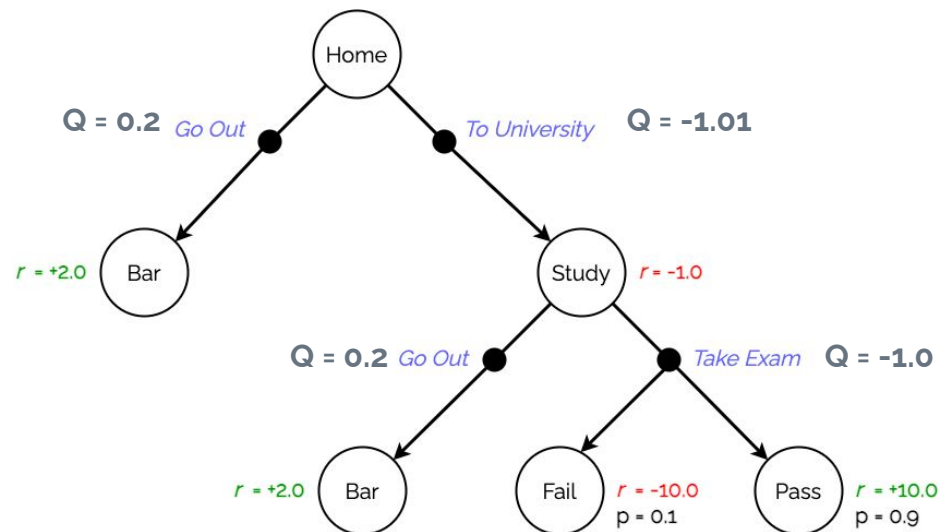
1. We need to collect our own data

# Exploration/Exploitation trade-off

We need **exploration** because actions may look worse than they are.

Reasons:

1. We need to collect our own data

# Exploration/Exploitation trade-off

We need **exploration** because actions may look worse than they are.

Reasons:

1. We need to collect our own data

A good decision may seem bad if it is followed by bad decisions (and vice versa)
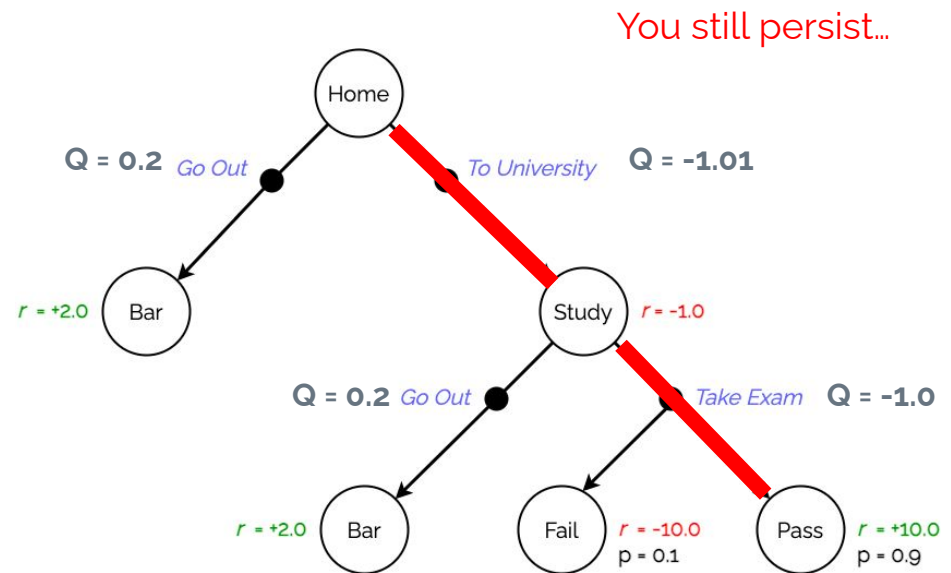
# Exploration/Exploitation trade-off

We need **exploration** because actions may look worse than they are.

Reasons:

1. We need to collect our own data
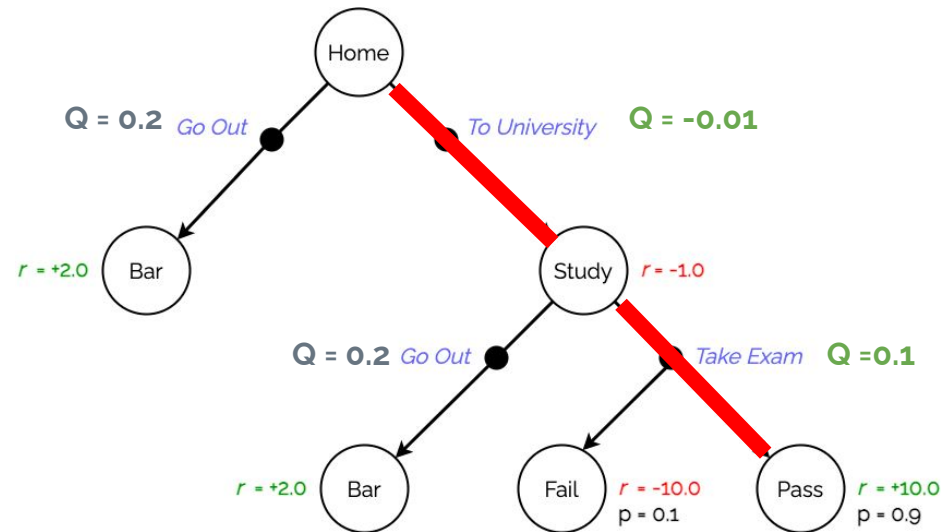2. Good action may seem bad if followed by bad actions

# Exploration/Exploitation trade-off

We need **exploration** because actions may look worse than they are.

Reasons:

1. We need to collect our own data
2. Good action may seem bad if followed by bad actions

# Exploration/Exploitation trade-off

We need **exploration** because actions may look worse than they are.

Reasons:

1. We need to collect our own data
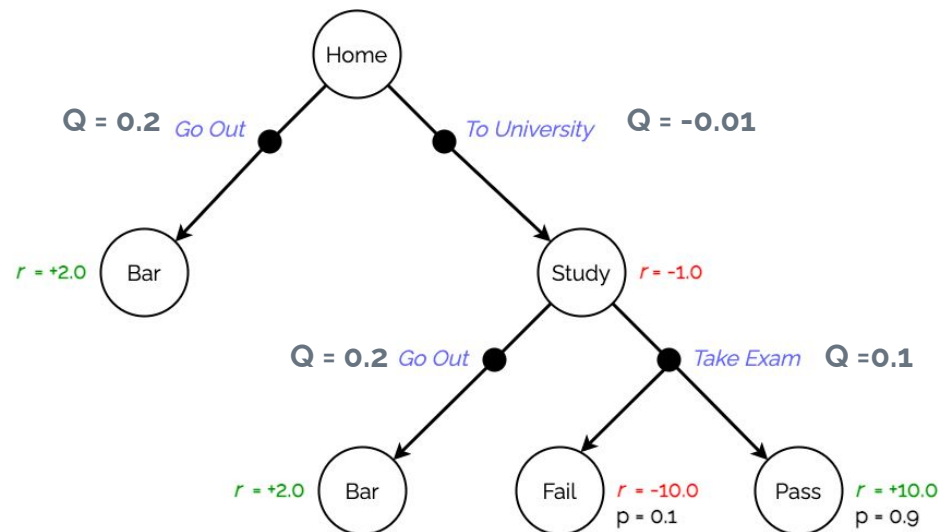2. Good action may seem bad if followed by bad actions

# Exploration/Exploitation trade-off

We need **exploration** because actions may look worse than they are.

Reasons:

1. We need to collect our own data
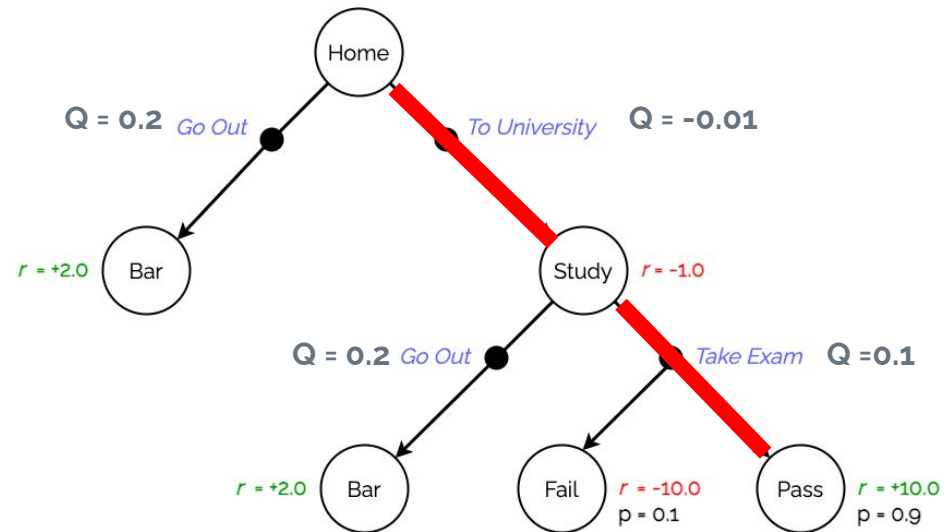2. Good action may seem bad if followed by bad actions

# Exploration/Exploitation trade-off

We need **exploration** because actions may look worse than they are.

<u>Reasons</u>:

1. We need to collect our own data
2. Good action may seem bad if followed by bad actions

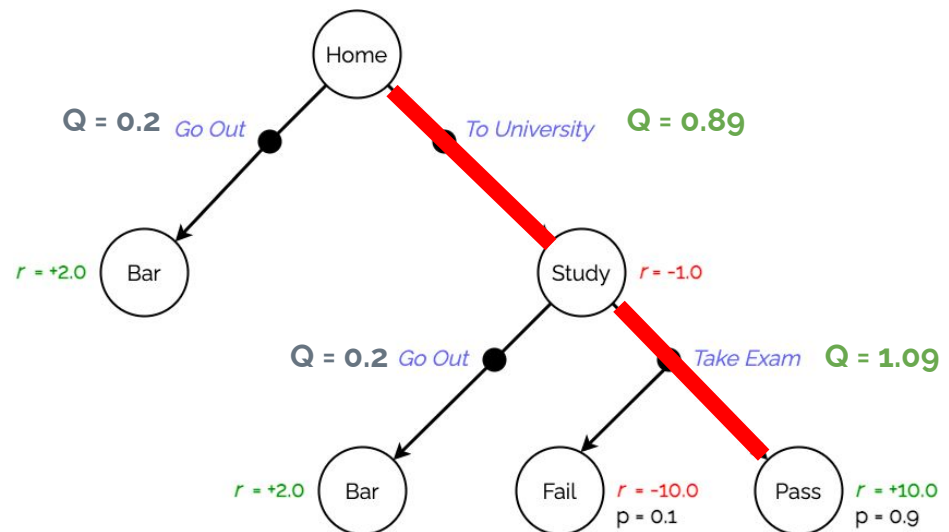Action may seem overly bad due to bad luck (and vice versa)

# Exploration/Exploitation trade-off

We need **exploration** because actions may look worse than they are.

Reasons:

1. We need to collect our own data
2. Good action may seem bad if followed by bad actions
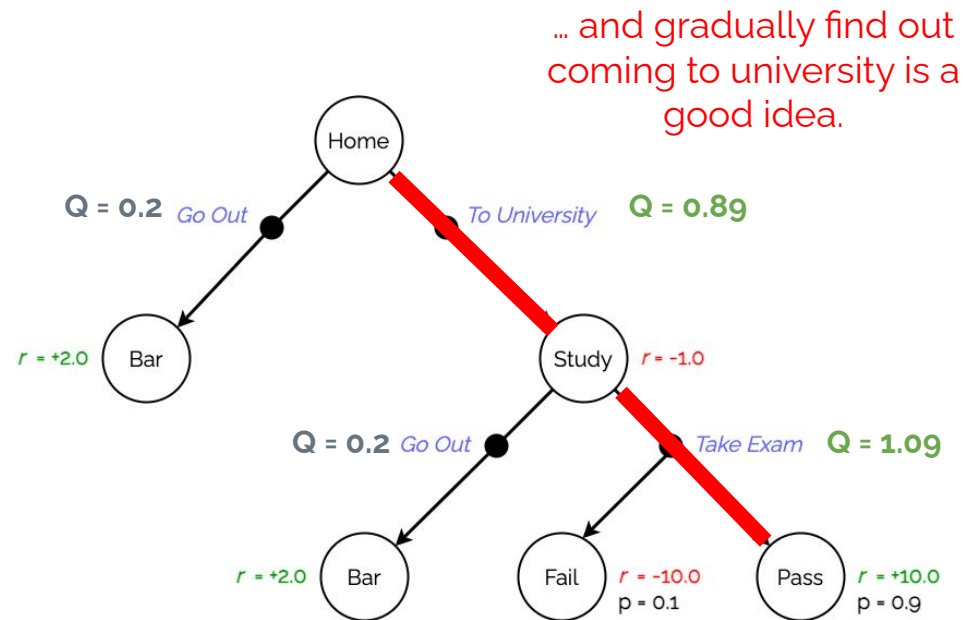3. Environment can be stochastic

# Exploration/Exploitation trade-off

We need **exploration** because actions may look worse than they are.

Reasons:

1. We need to collect our own data
2. Good action may seem bad if followed by bad actions
3. Environment can be stochastic
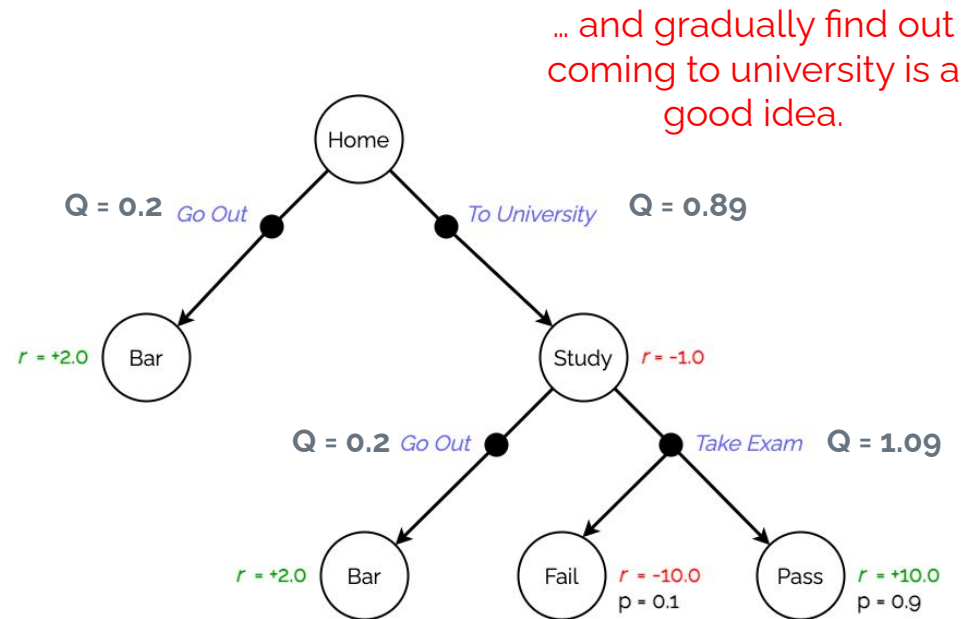
You still persist...

# Exploration/Exploitation trade-off

We need **exploration** because actions may look worse than they are.

Reasons:

1. We need to collect our own data
2. Good action may seem bad if followed by bad actions
3. Environment can be stochastic

# Exploration/Exploitation trade-off

We need **exploration** because actions may look worse than they are.
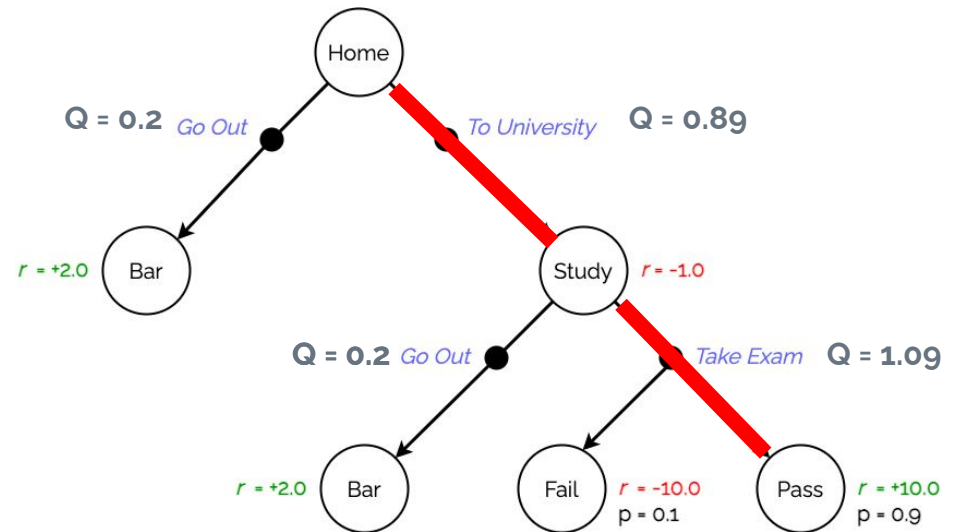
Reasons:

1. We need to collect our own data
2. Good action may seem bad if followed by bad actions
3. Environment can be stochastic

# Exploration/Exploitation trade-off

We need **exploration** because actions may look worse than they are.

Reasons:

1. We need to collect our own data
2. Good action may seem bad if followed by bad actions
3. Environment can be stochastic

# Exploration/Exploitation trade-off

We need **exploration** because actions may look worse than they are.

Reasons:

1. We need to collect our own data
2. Good action may seem bad if followed by bad actions
3. Environment can be stochastic

# Exploration/Exploitation trade-off

We need **exploration** because actions may look worse than they are.

Reasons:

1. We need to collect our own data
2. Good action may seem bad if followed by bad actions
3. Environment can be stochastic

... and gradually find out coming to university is a good idea.

# Exploration/Exploitation trade-off

We need **exploration** because actions may look worse than they are.

Reasons:

1. We need to collect our own data
2. Good action may seem bad if followed by bad actions
3. Environment can be stochastic

We also need **exploitation**.

... and gradually find out coming to university is a good idea.

# Exploration/Exploitation trade-off

We need **exploration** because actions may look worse than they are.

Reasons:

1. We need to collect our own data
2. Good action may seem bad if followed by bad actions
3. Environment can be stochastic

We also need **exploitation**.

Reasons:

1. Want to use what we learned
2. In bigger problems: move in promising directions to further explore.

# Exploration/Exploitation strategies

# Exploration/Exploitation strategies

Huge amount of strategies, we will here discuss one (simple) example:

Boltzmann (softmax) exploration

# Boltzmann (softmax) exploration

# Boltzmann (softmax) exploration



$a_1$    $a_2$    $a_3$

Q = 2     Q = -1     Q = 4

# Boltzmann (softmax) exploration



Q = 2        Q = -1        Q = 4

<u>Intuition</u>: give all actions a chance (exploration), but actions with higher Q-estimate deserve a higher probability (exploitation).

# Boltzmann (softmax) exploration

# Boltzmann (softmax) exploration

$$\pi(a_i|s) = \frac{e^{Q(s,a_i)/\tau}}{\sum_{a \in \mathcal{A}} e^{Q(s,a)/\tau}}$$

# Boltzmann (softmax) exploration

To get the probability of
selecting action $a_i$ in state s…

$$\pi(a_i | s) = \frac{e^{Q(s, a_i)/\tau}}{\sum_{a \in \mathcal{A}} e^{Q(s, a)/\tau}}$$

# Boltzmann (softmax) exploration

To get the probability of
selecting action $a_i$ in state s…

… we exponentiate its Q-value…

$$\pi(a_i|s) = \frac{e^{Q(s,a_i)/\tau}}{\sum_{a \in \mathcal{A}} e^{Q(s,a)/\tau}}$$

# Boltzmann (softmax) exploration

To get the probability of selecting action $a_i$ in state s…

… we exponentiate its Q-value…

$$\pi(a_i|s) = \frac{e^{Q(s,a_i)/\tau}}{\sum_{a\in\mathcal{A}} e^{Q(s,a)/\tau}}$$

… and normalize over the sum of exponentiated Q-values of all actions (to make it a valid probability distribution).

# Boltzmann (softmax) exploration

To get the probability of
selecting action $a_i$ in state s…

… we exponentiate its Q-value…

$$\pi(a_i|s) = \frac{e^{Q(s,a_i)\boxed{/\tau}}}{\sum_{a \in \mathcal{A}} e^{Q(s,a)\boxed{/\tau}}}$$

… and normalize over the sum of
exponentiated Q-values of all actions
(to make it a valid probability
distribution).

Temperature **τ** scales the
amount of exploration:

$\tau \to 0$  :    one-hot (exploit)

$\tau \to \infty$  :    uniform (explore)

# Boltzmann (softmax) exploration



a₁ — $a_1$

a₂ — $a_2$

a₃ — $a_3$

Q = 2    Q = -1    Q = 4

# Boltzmann (softmax) exploration



$a_1$    $a_2$    $a_3$

Q = 2      Q = -1      Q = 4

**τ=1.0**      π=0.12      π=0.01      π=0.87

# Boltzmann (softmax) exploration



| | Q = 2 | Q = -1 | Q = 4 |
|---|---|---|---|
| $\tau$=1.0 | $\pi$=0.12 | $\pi$=0.01 | $\pi$=0.87 |
| $\tau$=0.01 | $\pi$=0.0 | $\pi$=0.0 | $\pi$=1.0 |

# Boltzmann (softmax) exploration



| | Q = 2 | Q = -1 | Q = 4 | |
|---|---|---|---|---|
| **τ=1.0** | π=0.12 | π=0.01 | π=0.87 | |
| **τ=0.01** | π=0.0 | π=0.0 | π=1.0 | **full exploitation** |

# Boltzmann (softmax) exploration



|        | Q = 2 | Q = -1 | Q = 4 |  |
|--------|-------|--------|-------|--|
| $\tau$=1.0 | π=0.12 | π=0.01 | π=0.87 |  |
| $\tau$=0.01 | π=0.0 | π=0.0 | π=1.0 | **full exploitation** |
| $\tau$=100 | π=0.33 | π=0.32 | π=0.34 |  |

# Boltzmann (softmax) exploration



|  | Q = 2 | Q = -1 | Q = 4 |  |
|---|---|---|---|---|
| **τ=1.0** | π=0.12 | π=0.01 | π=0.87 |  |
| **τ=0.01** | π=0.0 | π=0.0 | π=1.0 | **full exploitation** |
| **τ=100** | π=0.33 | π=0.32 | π=0.34 | **full exploration** |

# Boltzmann (softmax) exploration



Can anneal $\tau$ during training to gradually transition from exploration to exploitation

|  | Q = 2 | Q = -1 | Q = 4 |  |
| --- | --- | --- | --- | --- |
| $\tau$=1.0 | $\pi$=0.12 | $\pi$=0.01 | $\pi$=0.87 |  |
| $\tau$=0.01 | $\pi$=0.0 | $\pi$=0.0 | $\pi$=1.0 | full exploitation |
| $\tau$=100 | $\pi$=0.33 | $\pi$=0.32 | $\pi$=0.34 | full exploration |

# Exploration

(video)

# The Reinforcement Learning Cycle

*Pseudocode*

Initialize **Q(s,a)** estimates for all states,actions (e.g. to 0)

Repeat:

  1) Exploration:

  2) Credit assignment:

  3) Update:

# The Reinforcement Learning Cycle

*Pseudocode*

Initialize **Q(s,a)** estimates for all states,actions (e.g. to 0)

Repeat:

   1) <u>Exploration</u>: Boltzmann policy with annealing temperature.

   2) <u>Credit assignment</u>:

   3) <u>Update</u>:

# The Reinforcement Learning Cycle

*Pseudocode*

Initialize **Q(s,a)** estimates for all states,actions
(e.g. to 0)


Repeat:

    1) Exploration: Boltzmann policy with
       annealing temperature.


    2) Credit assignment: Monte Carlo back-up.


    3) Update:

# The Reinforcement Learning Cycle

*Pseudocode*

Initialize **Q(s,a)** estimates for all states,actions (e.g. to 0)

Repeat:

1) <u>Exploration</u>: Boltzmann policy with annealing temperature.

2) <u>Credit assignment</u>: Monte Carlo back-up.

3) <u>Update</u>: Tabular learning rule with learning rate 0.1
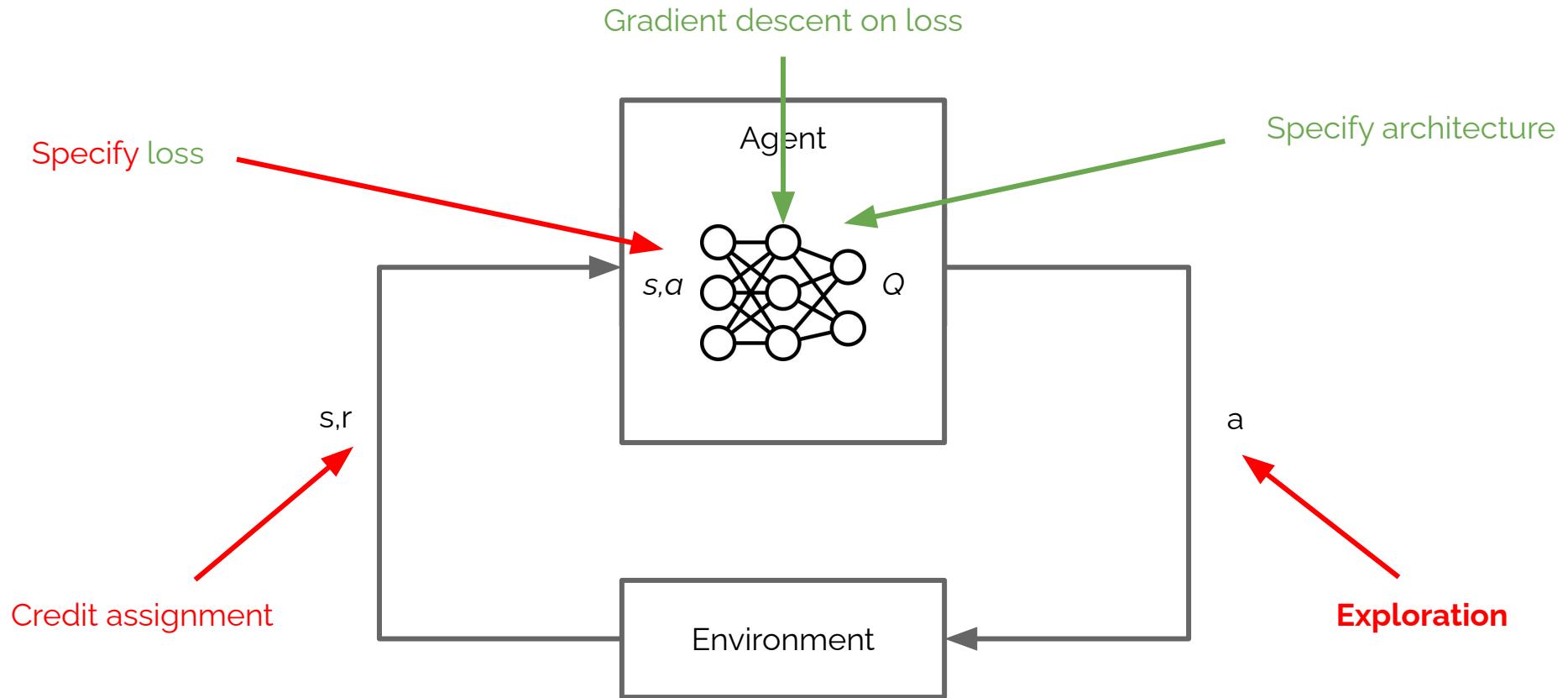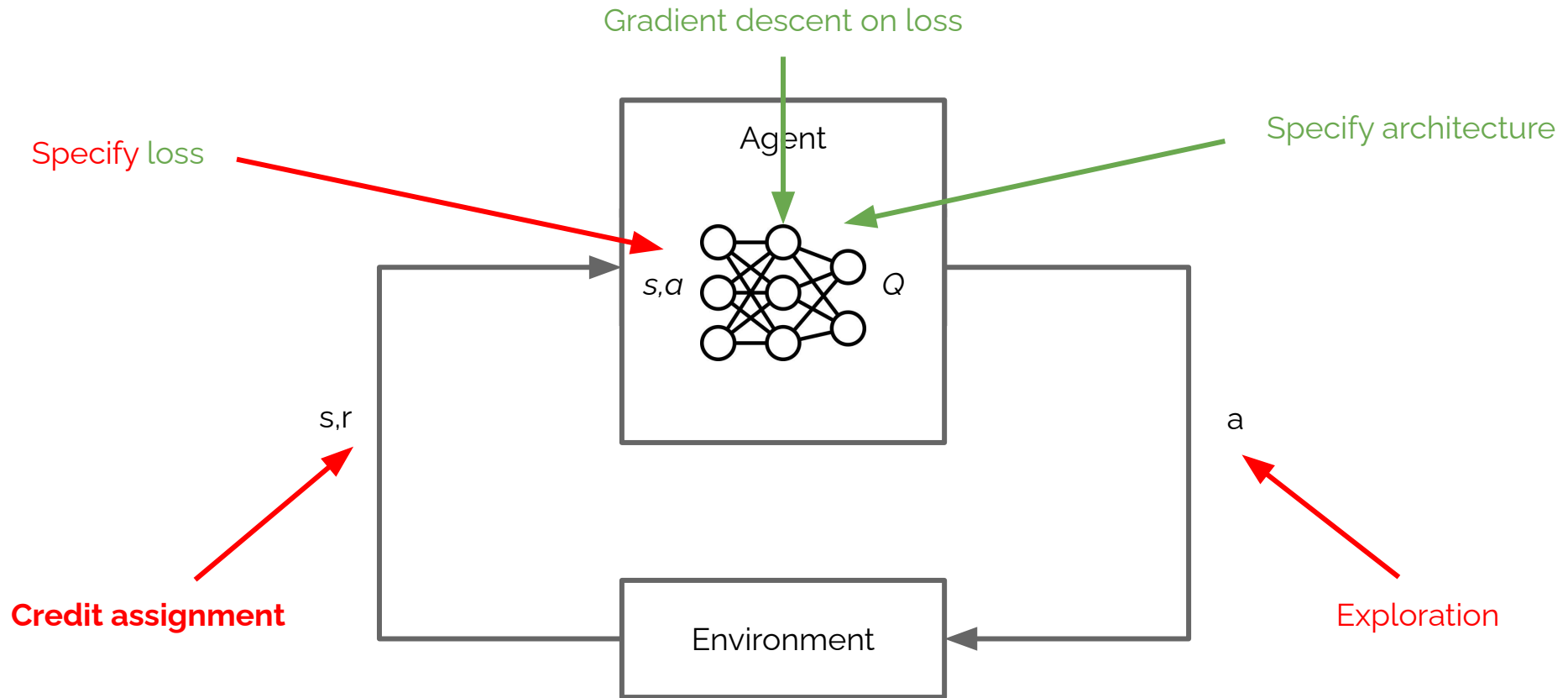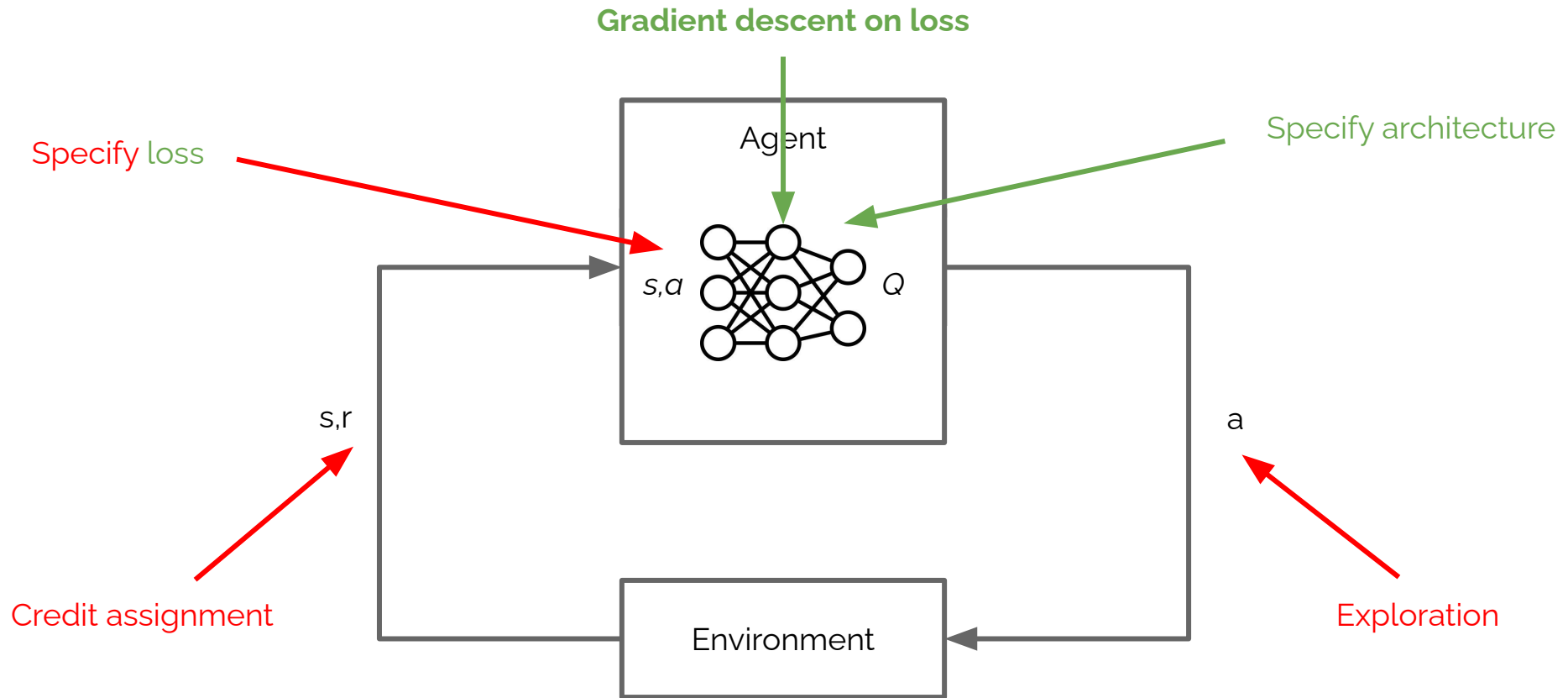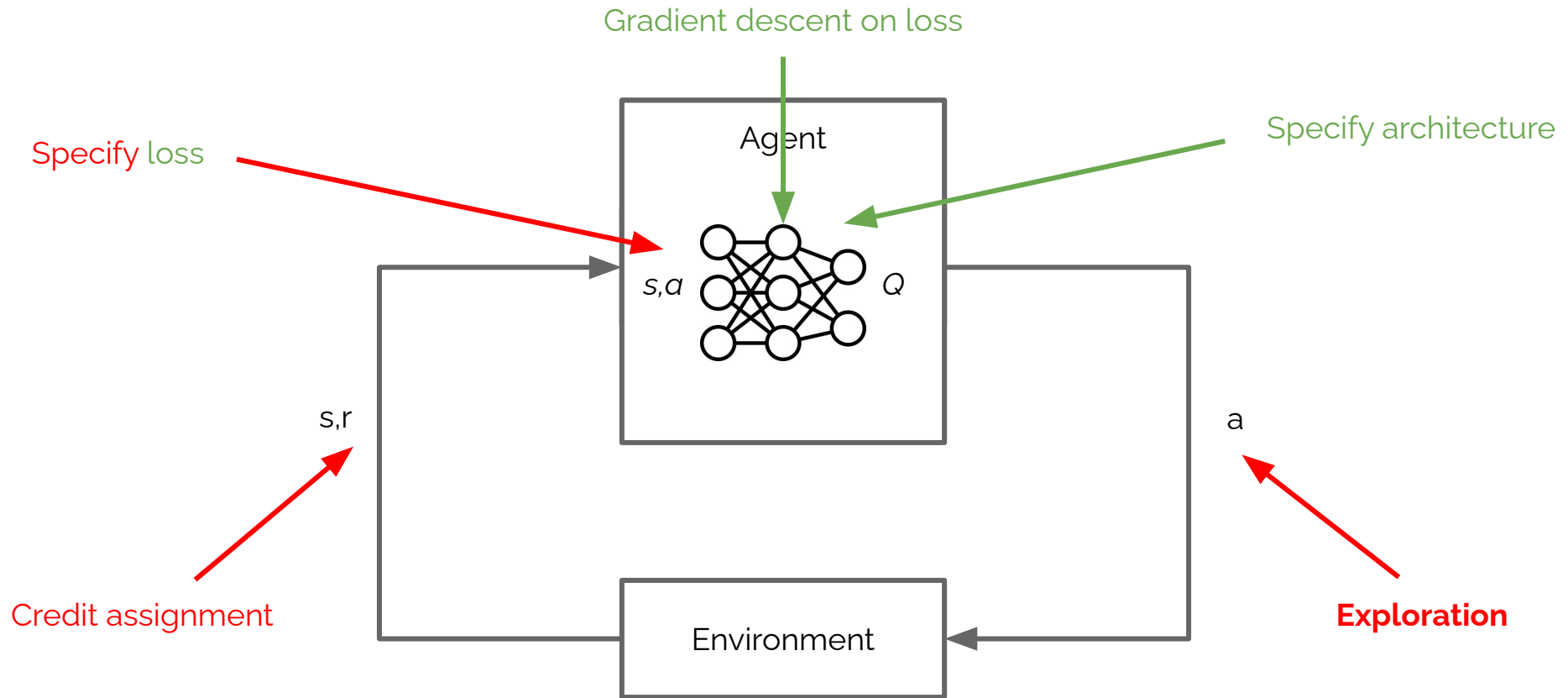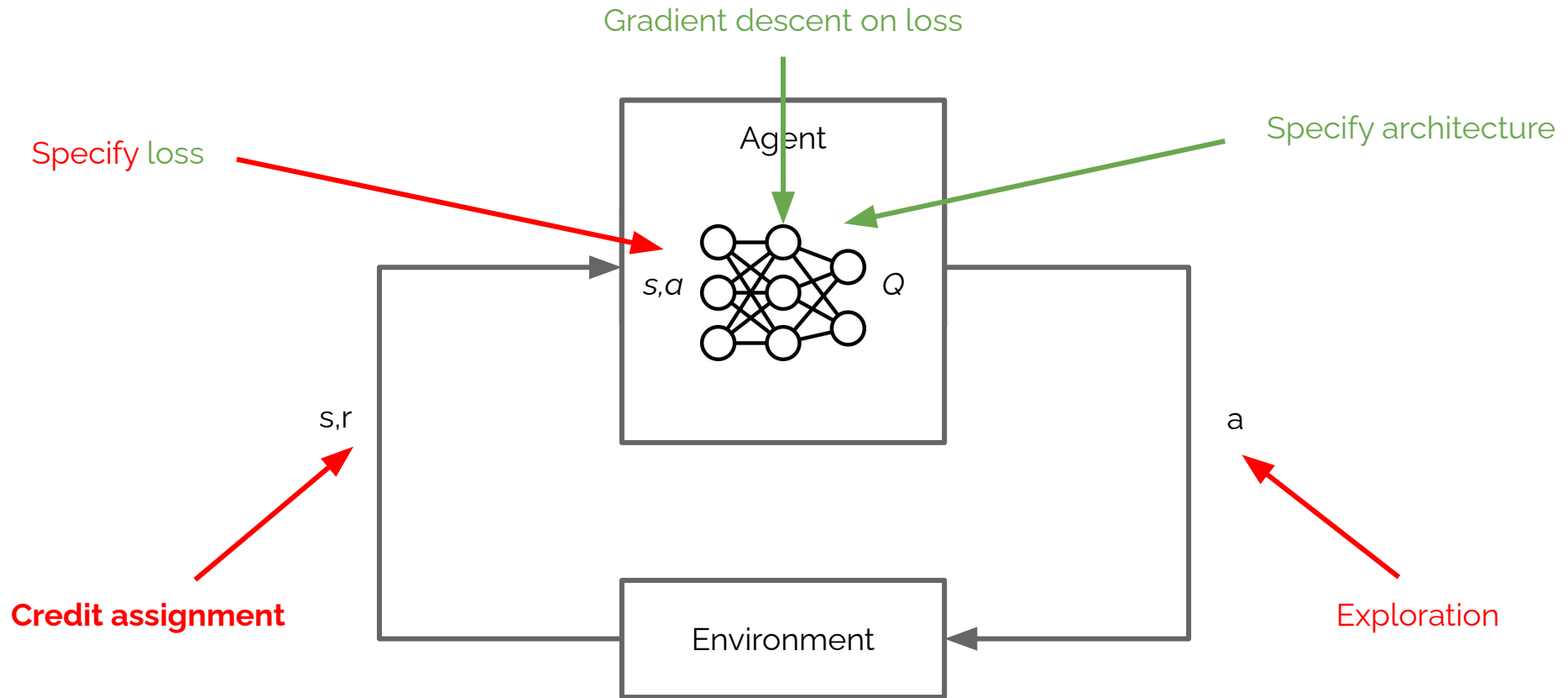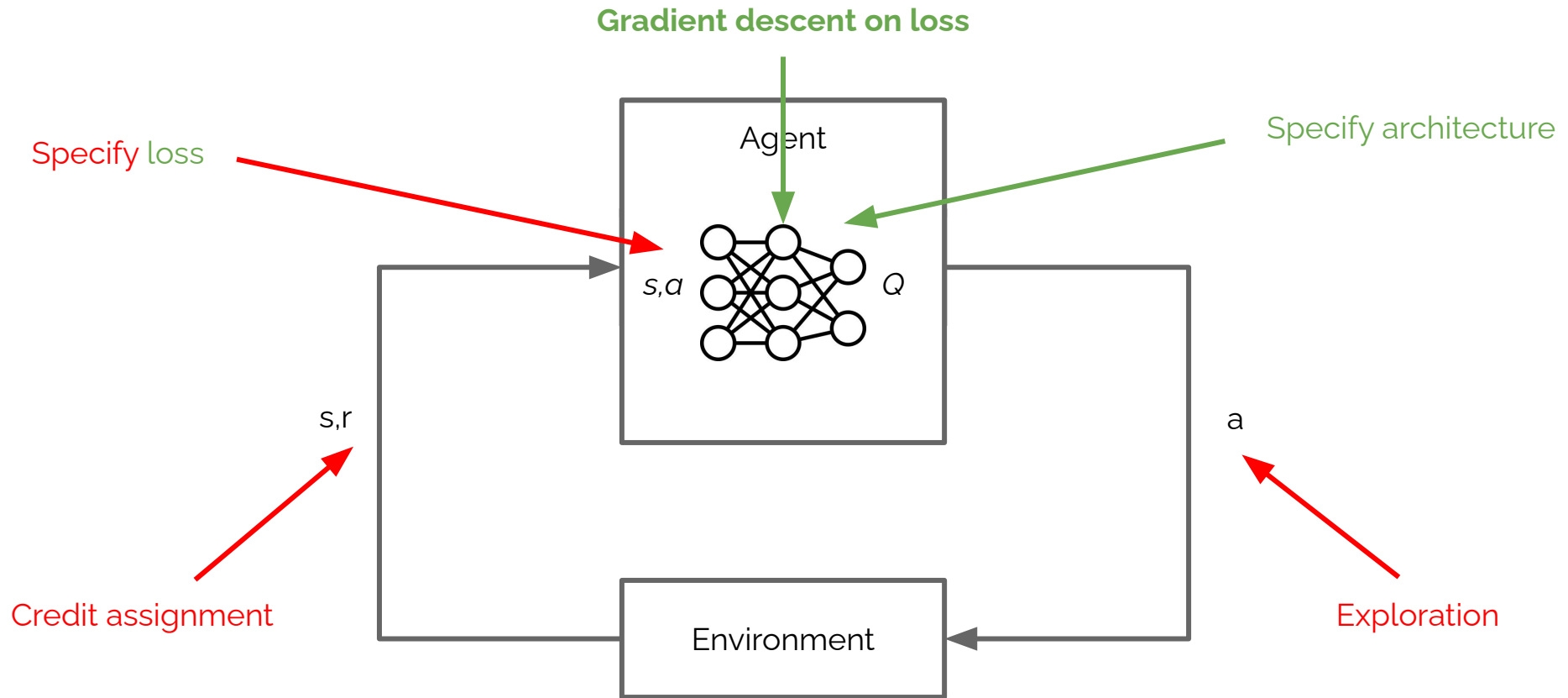
# Part IV

# Deep reinforcement learning

# Deep reinforcement learning

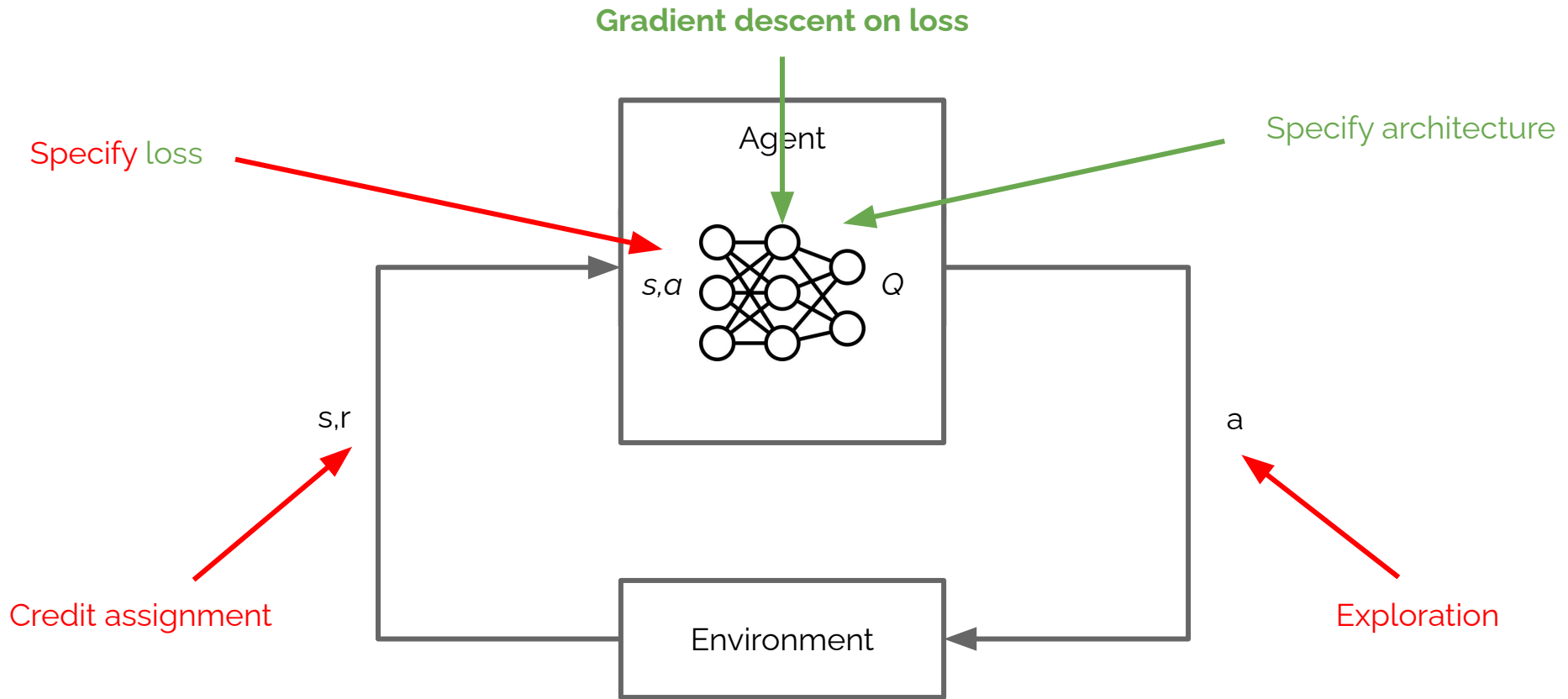# Deep reinforcement learning

Deep reinforcement learning = deep learning + reinforcement learning

# Deep reinforcement learning

Deep reinforcement learning = deep learning + reinforcement learning



Observation spaces in reinforcement are usually high-dimensional

# Deep reinforcement learning

Deep reinforcement learning = deep learning + reinforcement learning



Observation spaces in reinforcement are usually high-dimensional

We need to use *function approximation*, e.g., deep learning, to store our solution
(to fit it in memory & profit from generalization)

# Deep Reinforcement Learning

# Deep Reinforcement Learning

# Deep Reinforcement Learning

# Deep Reinforcement Learning

# Deep Reinforcement Learning

# Deep Reinforcement Learning

# Deep Reinforcement Learning

Gradient descent on loss

Specify architecture

Specify loss

Agent

s,a

Q

s,r

a

Credit assignment

Environment

Exploration

# Deep Reinforcement Learning

# Deep Reinforcement Learning

Gradient descent on loss

Specify architecture

Specify loss

Agent

s,a    Q

Credit assignment

s,r

a

Exploration

Environment

# Deep Reinforcement Learning

# Deep Reinforcement Learning

Gradient descent on loss

Agent

Specify architecture

Specify loss

*s,a*          *Q*

s,r

a

Credit assignment

Environment

**Exploration**

# Deep Reinforcement Learning



Gradient descent on loss

Specify architecture

Agent

Specify loss

s,a

Q

Credit assignment

s,r

a

Exploration

Environment

# Deep Reinforcement Learning

# Deep Reinforcement Learning

**Gradient descent on loss**

Specify architecture

Specify loss

Agent

s,a    Q

Credit assignment

s,r

a

Exploration

Environment

RL is supervised learning

# Deep Reinforcement Learning

**Gradient descent on loss**

Specify loss

Specify architecture

Agent

*s,a*     *Q*

s,r

Credit assignment

a

Exploration

Environment

RL is supervised learning on a moving target function

# Deep Reinforcement Learning

**Gradient descent on loss**

Specify loss

Agent

Specify architecture

*s,a*

*Q*

s,r

a

Credit assignment

Environment

Exploration

RL is supervised learning on a moving target function that influences which data you see.

# Conclusion

# Conclusion

# Conclusion

- Many tasks can be formulated as a sequential decision making problem, for which you can use reinforcement learning (RL).

# Conclusion

- Many tasks can be formulated as a sequential decision making problem, for which you can use reinforcement learning (RL).

- The main benefit of RL is that you can learn tasks you only label on outcomes, potentially outperforming human solutions.

# Conclusion

- Many tasks can be formulated as a sequential decision making problem, for which you can use reinforcement learning (RL).

- The main benefit of RL is that you can learn tasks you only label on outcomes, potentially outperforming human solutions.

- Key topics in RL are exploration (which action should I try next) and credit assignment (how do I process the obtained reward information).

# Conclusion

- Many tasks can be formulated as a sequential decision making problem, for which you can use reinforcement learning (RL).

- The main benefit of RL is that you can learn tasks you only label on outcomes, potentially outperforming human solutions.

- Key topics in RL are exploration (which action should I try next) and credit assignment (how do I process the obtained reward information).

- Reinforcement learning is supervised learning on a moving target that influences which data you see.

# Courses

# Courses



[irl.liacs.nl](irl.liacs.nl)

# Courses





[irl.liacs.nl](irl.liacs.nl)

[rl.liacs.nl](rl.liacs.nl)

# Courses



[irl.liacs.nl](irl.liacs.nl)



[rl.liacs.nl](rl.liacs.nl)



[arl.liacs.nl](arl.liacs.nl)

# AI & Robotics challenge



Artificial Intelligence & Robotics (AIR) Challenge Leiden

Home    Timeline    Information

The AI & Robotics Challenge is a yearly bachelor student competition that runs within the Leiden Institute of Advanced Computer Science (LIACS).

# AI & Robotics challenge



Artificial Intelligence & Robotics (AIR)

Challenge Leiden

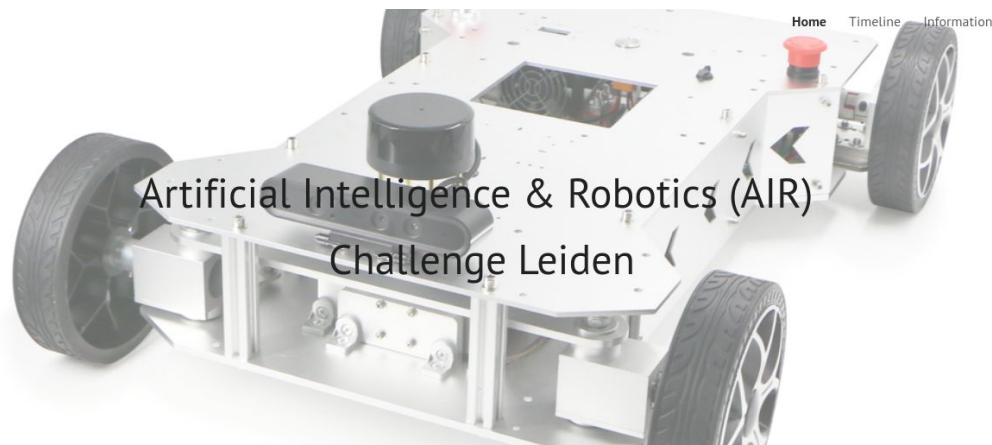Home    Timeline    Information

The AI & Robotics Challenge is a yearly bachelor student competition that runs within the Leiden Institute of Advanced Computer Science (LIACS).

Extra-curricular course (2 ECTS)

Sign-up in September 2023

# Questions?