

A deep dive into reinforcement learning

prof. dr. Pieter Libin
Artificial intelligence lab
Department of computer science
Vrije Universiteit Brussel



06/02/2023

▶ Introduction

✉ Pieter.Libin@vub.be
🐦 @PieterLibin



- PhD in Computer science, AI lab, VUB
- Postdoc at Data science institute, UHasselt
- Assistant Professor, AI lab, VUB
- Research: decision making using detailed models and RL
- Interests: computational biology, big data, engineering

■ Program

- 10h00-12h00: Lecture on Reinforcement Learning
- 13h00-15h00: RL Programming exercise in Python
- **Interactive session:** ask questions + participate!

| 5

■ Attendance

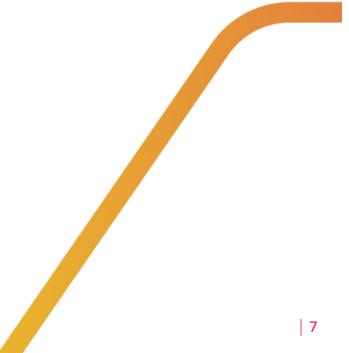


| 6

■ Course objectives

- **Intuition** about reinforcement learning
- How can you use reinforcement learning
- **Understanding** important reinforcement learning **algorithms**

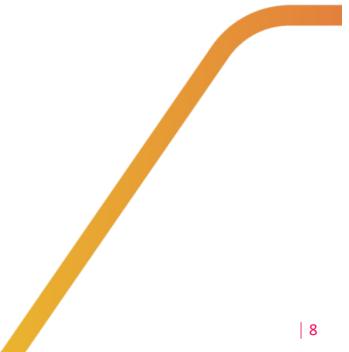
| 7

A thick yellow S-shaped curve, resembling a sigmoid function, starts at the bottom left and ends at the top right, indicating a transition or process.

■ Reinforcement learning, a bit different

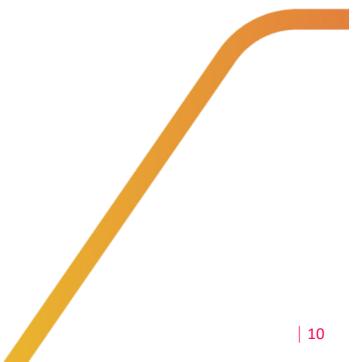
- Supervised vs unsupervised machine learning
 - Learn from labelled vs unlabelled data

| 8

A thick yellow S-shaped curve, resembling a sigmoid function, starts at the bottom left and ends at the top right, indicating a transition or process.

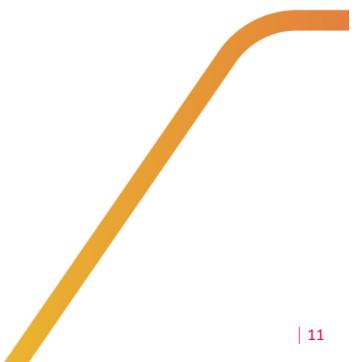
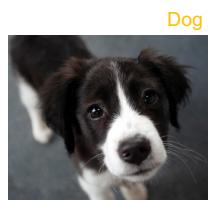
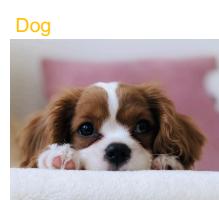
► Reinforcement learning, a bit different

- **Supervised** vs unsupervised machine learning
 - Learn from labelled vs unlabelled data



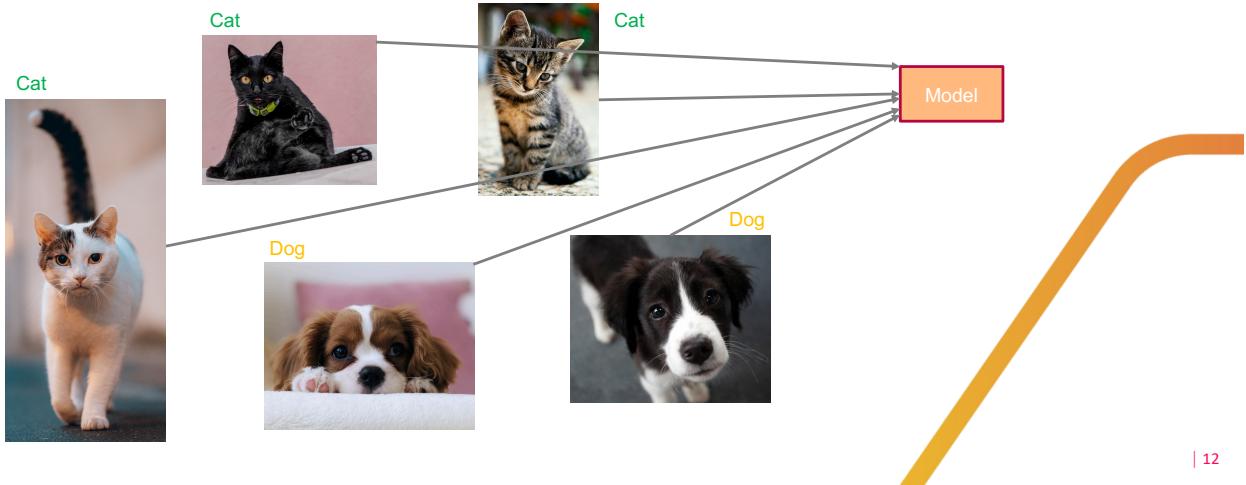
► Reinforcement learning, a bit different

- **Supervised** vs unsupervised machine learning
 - Learn from labelled vs unlabelled data



Reinforcement learning, a bit different

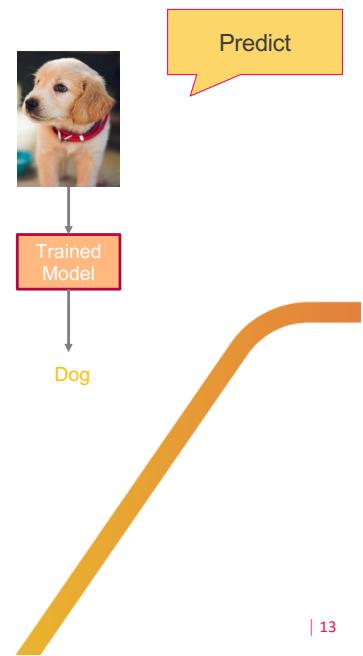
- **Supervised** vs unsupervised machine learning
 - Learn from labelled vs unlabelled data



| 12

Reinforcement learning, a bit different

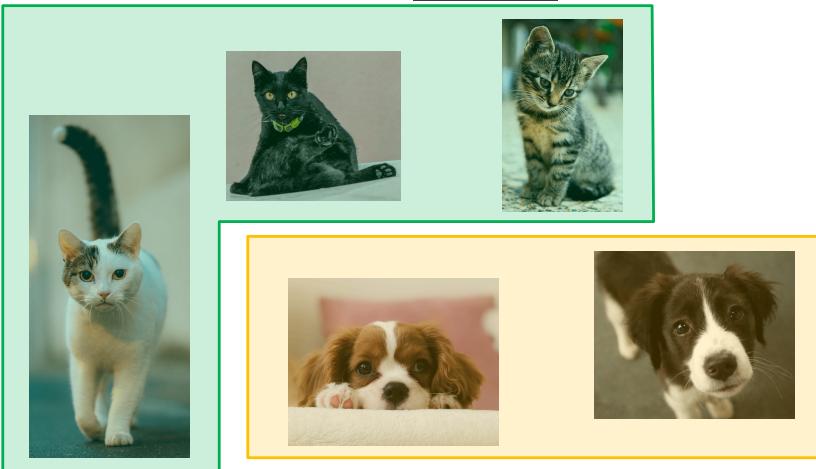
- **Supervised** vs unsupervised machine learning
 - Learn from labelled vs unlabelled data



| 13

► Reinforcement learning, a bit different

- Supervised vs unsupervised machine learning
 - Learn from labelled vs unlabelled data



| 15

► Reinforcement learning, a bit different

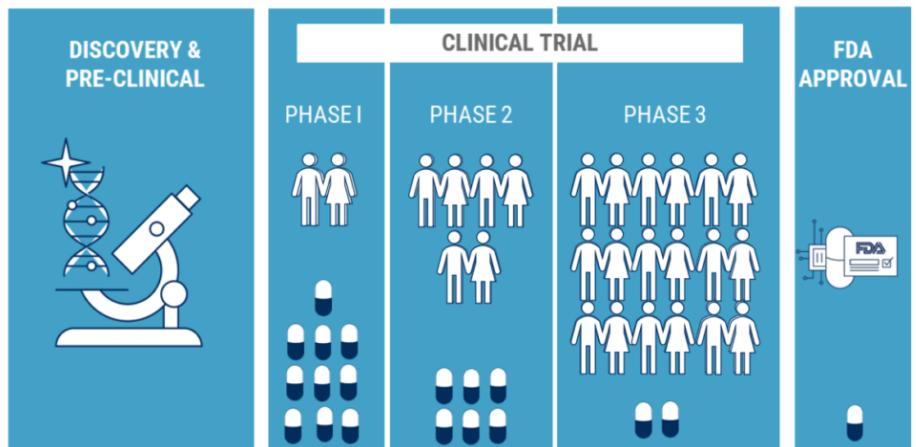
- Reinforcement learning:
 - sequential decision making problems
 - optimize behaviour by interacting with an environment
 - get to know an environment by trial-and-error
 - maximise long-term reward
- We have an **agent** that wants to learn an optimal **policy**, by **interacting** with an **environment**, by maximising a **reward** signal



| 17

■ Reinforcement learning, some applications

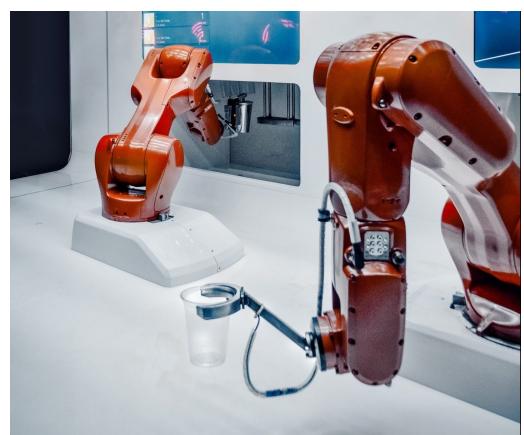
- Design of clinical trials, Williamson (2019)



| 19

■ Reinforcement learning, some applications

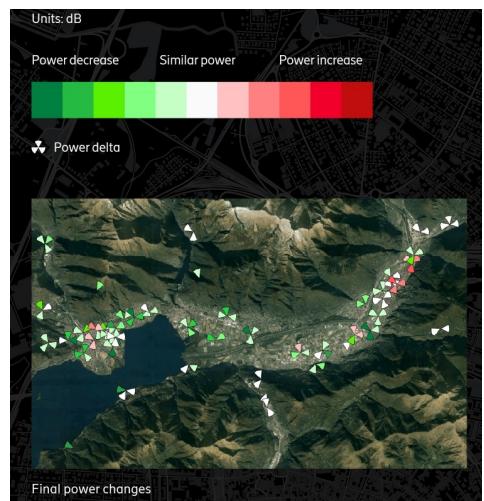
- Design of clinical trials, Williamson (2019)
- Robot control, Kober (2013)



| 20

► Reinforcement learning, some applications

- Design of clinical trials, Williamson (2019)
- Robot control, Kober (2013)
- Telecom, Ericsson (2020)

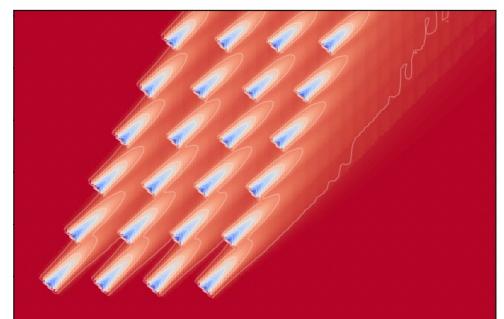


| 21

► Reinforcement learning, some applications

- Design of clinical trials, Williamson (2019)
- Robot control, Kober (2013)
- Telecom, Ericsson (2020)
- Wind farm control, Verstraeten (2020)

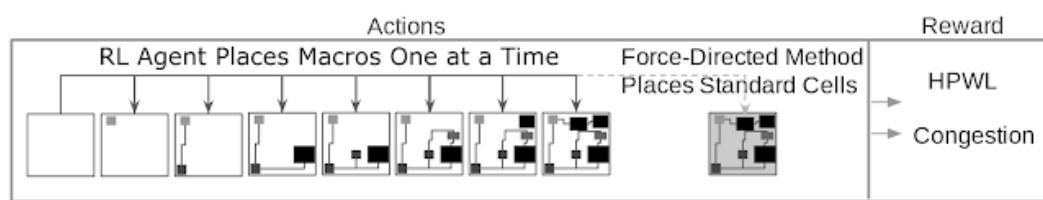
Simulator



| 22

■ Reinforcement learning, some applications

- Design of clinical trials, Williamson (2019)
- Robot control, Kober (2013)
- Telecom, Ericsson (2020)
- Wind farm control, Verstraeten (2020)
- Chip design, Mirhoseini (2020)



| 23

■ Reinforcement learning, some applications

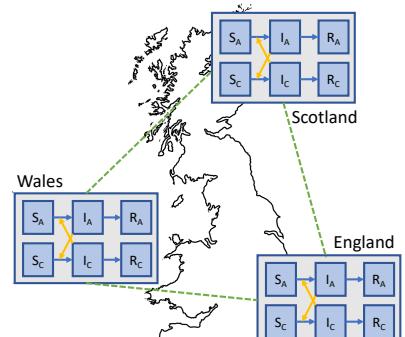
- Design of clinical trials, Williamson (2019)
- Robot control, Kober (2013)
- Telecom, Ericsson (2020)
- Wind farm control, Verstraeten (2020)
- Chip design, Mirhoseini (2020)
- Electric vehicle charging, Energis/VUB (2020)



| 24

Reinforcement learning, some applications

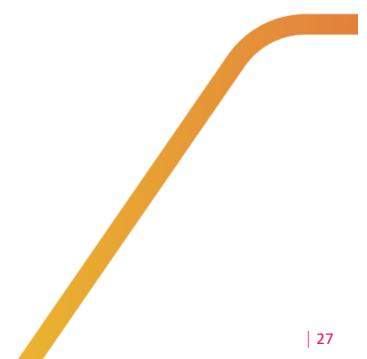
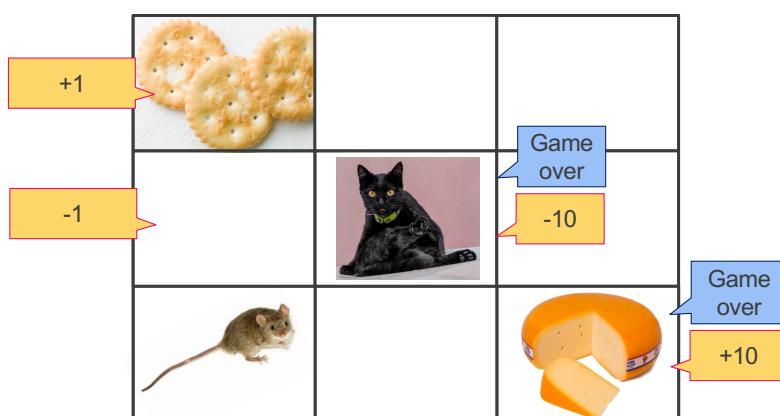
- Design of clinical trials, Williamson (2019)
- Robot control, Kober (2013)
- Telecom, Ericsson (2020)
- Wind farm control, Verstraeten (2020)
- Chip design, Mirhoseini (2020)
- Electric vehicle charging, Energis/VUB (2020)
- Epidemic control, Libin (2020)



| 25

The reinforcement learning problem

- Reinforcement learning:
 - We have an **agent** that wants to learn an optimal **policy**, by interacting with an **environment**, by maximising a **reward** signal



| 27

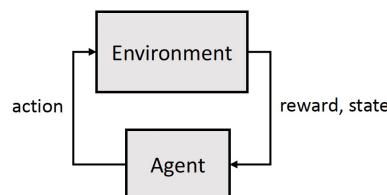
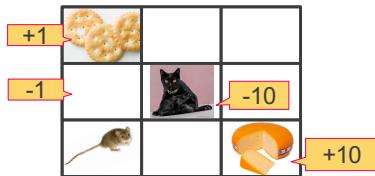
Markov decision process

Markovian



A Markov Decision Process corresponds to a tuple $\langle \mathcal{S}, \mathcal{A}, T, \gamma, R \rangle$, where

- \mathcal{S} is the set of possible states the environment can take upon
- \mathcal{A} is the set of possible actions an agent can take
- $T(s' | s, a)$ signifies the transition probability to go from state s to state s' by taking an action a
- γ is the discount factor that modulates the importance of future rewards
- $R(s, a, s')$ is the reward function that specifies which reward the agents receives upon choosing an action a in state s



| 28

Policy



A Markov Decision Process corresponds to a tuple $\langle \mathcal{S}, \mathcal{A}, T, \gamma, R \rangle$, where

- \mathcal{S} is the set of possible states the environment can take upon
- \mathcal{A} is the set of possible actions an agent can take
- $T(s' | s, a)$ signifies the transition probability to go from state s to state s' by taking an action a
- γ is the discount factor that modulates the importance of future rewards
- $R(s, a, s')$ is the reward function that specifies which reward the agents receives upon choosing an action a in state s



Given a Markov Decision Process $\langle \mathcal{S}, \mathcal{A}, T, \gamma, R \rangle$, an agent follows a policy p , that expresses the probability to take action a when in state s :

$$p : \mathcal{S} \times \mathcal{A} \rightarrow [0, 1]$$

| 30

Policy - Cumulative reward

- In the MDP, we aim to maximize **cumulative reward**

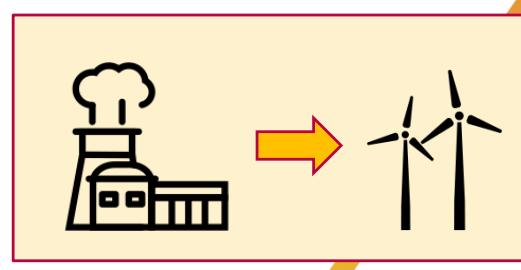
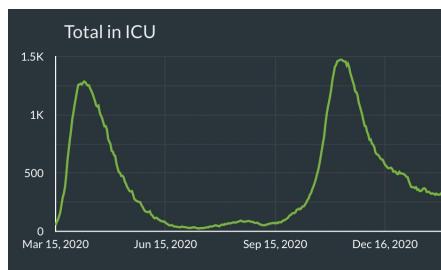


The return, or the discounted sum of rewards, starting from time step t , is defined as:

$$\sum_{i=0}^T \gamma^i R(s_{t+i}, a_{t+i}, s_{t+i+1}),$$

One trajectory through the MDP
T is horizon

where γ is the discount factor.



| 31

Policy

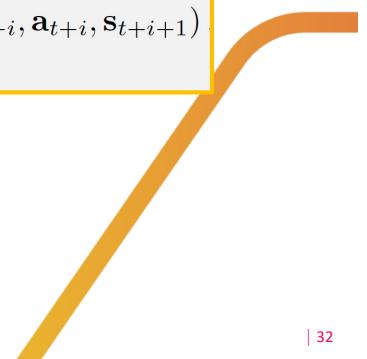


Given a Markov Decision Process $\langle S, A, T, \gamma, R \rangle$, an agent follows a policy p , that expresses the probability to take action a when in state s :

$$p : S \times A \rightarrow [0, 1]$$

+1		
-1		Game over -10
		Game over +10

$$\sum_{i=0}^T \gamma^i R(s_{t+i}, a_{t+i}, s_{t+i+1})$$



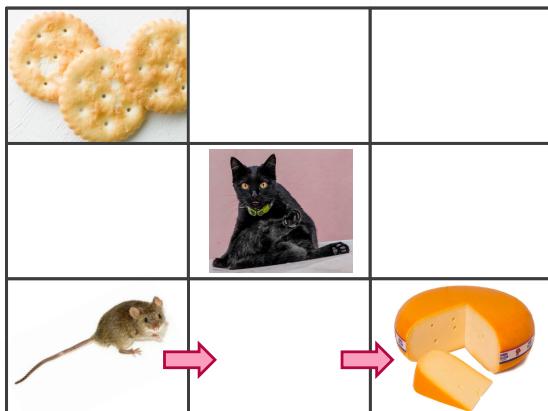
| 32

Policy

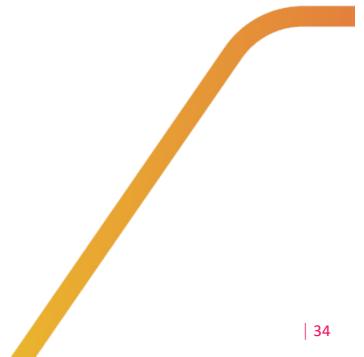


Given a Markov Decision Process $\langle \mathcal{S}, \mathcal{A}, T, \gamma, R \rangle$, an agent follows a policy p , that expresses the probability to take action a when in state s :

$$p : \mathcal{S} \times \mathcal{A} \rightarrow [0, 1]$$



$$p(., F) = 1$$



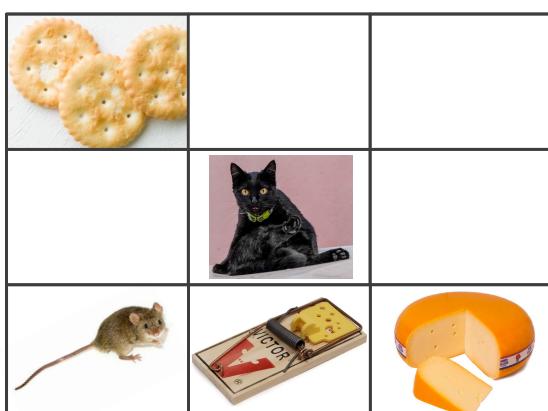
| 34

Policy

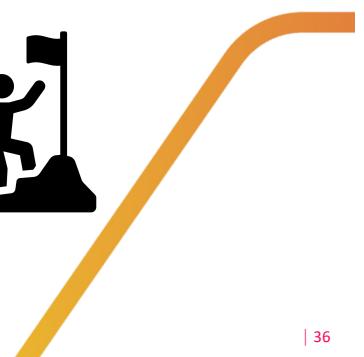


Given a Markov Decision Process $\langle \mathcal{S}, \mathcal{A}, T, \gamma, R \rangle$, an agent follows a policy p , that expresses the probability to take action a when in state s :

$$p : \mathcal{S} \times \mathcal{A} \rightarrow [0, 1]$$



$$p(., .) = \text{"Avoid cat"}$$



| 36

► Policies, some intuition

- What is a policy to you?
- How to treat a patient that is admitted to the hospital?



| 39

► Policies, some intuition

- What is a policy to you?
- How to treat a patient that is admitted to the hospital?
- How to treat an HIV patient?



| 40

► Policies, some intuition

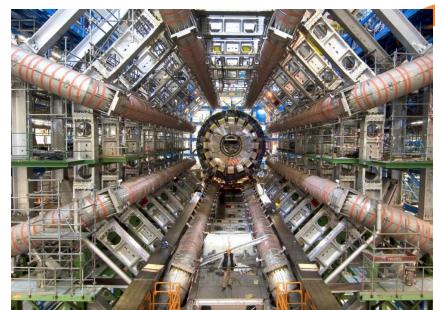
- What is a policy to you?
- How to treat a patient that is admitted to the hospital?
- How to treat an HIV patient?
- How to conduct a lab experiment?



| 41

► Policies, some intuition

- What is a policy to you?
- How to treat a patient that is admitted to the hospital?
- How to treat an HIV patient?
- How to conduct a lab experiment?
- How to collect data?



| 42

Value functions

- How good is it to be in a particular state



We define the value of being in state s , following a policy p , as:

$$V_p^p(s) = \mathbb{E}_{p,T} \left[\sum_{i=0}^T \gamma^i R(s_{t+i}, a_{t+i}, s_{t+i+1}) \mid s_t = s \right]$$

| 44

Value functions

- How good is it to be in a particular state



We define the value of being in state s , following a policy p , as:

$$V_p^p(s) = \mathbb{E}_{p,T} \left[\sum_{i=0}^T \gamma^i R(s_{t+i}, a_{t+i}, s_{t+i+1}) \mid s_t = s \right]$$

- How good is it to pick an action, when in a particular state



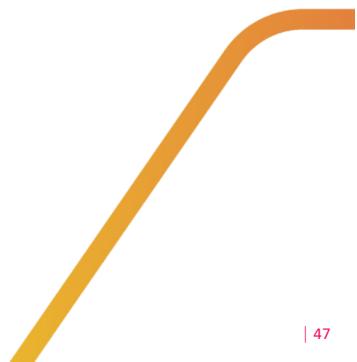
We define the value of being in state s , when taking action a , while following a policy p , as:

$$Q_p^p(s, a) = \mathbb{E}_{p,T} \left[\sum_{i=0}^T \gamma^i R(s_{t+i}, a_{t+i}, s_{t+i+1}) \mid s_t = s, a_t = a \right]$$

Value function

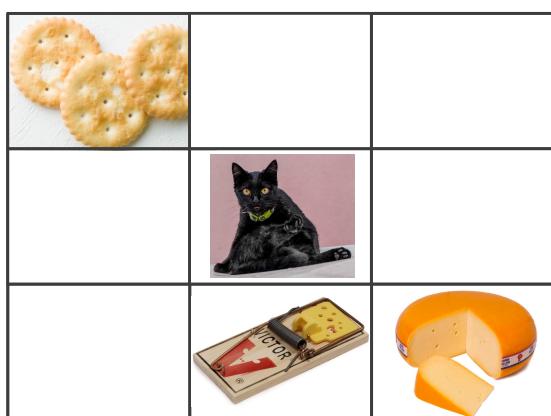
- Observation: Q*(s,a)
 - from an **optimal Q function**, we can derive an optimal policy
- Try to estimate Q*(s,a):
 - Q-learning

| 47

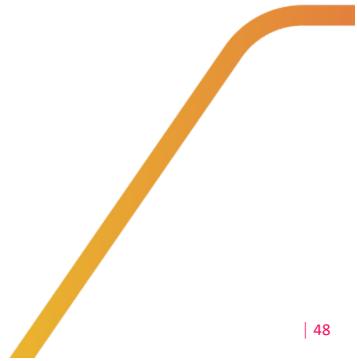


Q-learning

- Estimate Q*, by maintaining a Q-table:

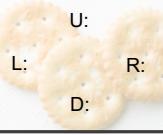
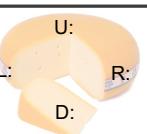


| 48

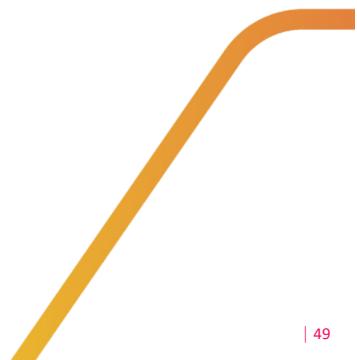


► Q-learning

- Estimate Q^* , by maintaining a Q-table:

	U: L: D:	U: L: D:	U: L: D:
U: L: D:		U: L: D:	U: L: D:
U: L: D:	U: L: D:		U: L: D:

| 49



► Q-learning

- Estimate Q^* , by maintaining a Q-table:

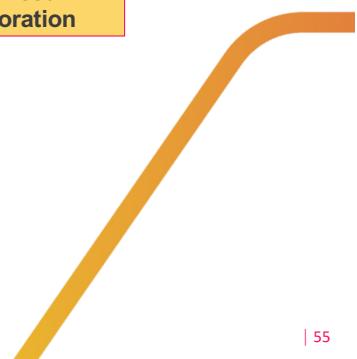
		
		
		

Exploit the best values

At the start,
all values are zero

We need:
exploration

| 55



■ Exploitation vs Exploration

- Crucial tradeoff to enable reinforcement learning
- Experiment:
 - What was the first beer you drank?



| 57

■ Exploitation vs Exploration

- Crucial tradeoff to enable reinforcement learning
- Experiment:
 - What was the first beer you drank?
 - Is your current favorite the same as your first beer?

| 58

Exploitation vs Exploration

- Crucial tradeoff to enable reinforcement learning
- Experiment:
 - What was the first beer you drank?
 - Is your current favorite the same as your first beer?

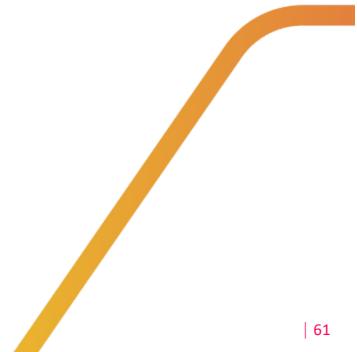


| 59

Q-learning

- Estimate Q^* , by maintaining a Q-table:

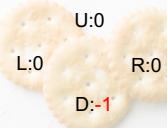
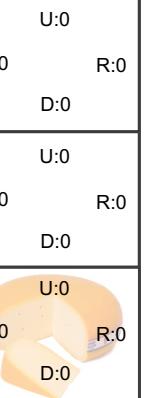
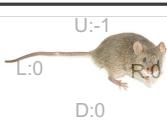
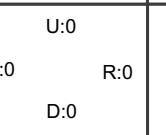
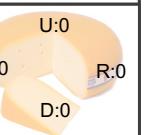
	U:0 L:0 D:0 R:0	U:0 L:0 D:0 R:0	U:0 L:0 D:0 R:0
U:0 L:0 D:0 R:0		U:0 L:0 D:0 R:0	U:0 L:0 D:0 R:0
	U:0 L:0 D:0 R:0	U:0 L:0 D:0 R:0	



| 61

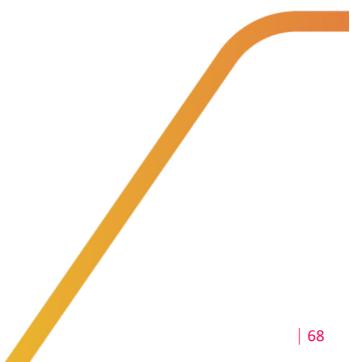
► Q-learning

- Estimate Q^* , by maintaining a Q-table:

 U:0 L:0 R:0 D:-1	 U:0 L:0 R:0 D:0	 U:0 L:0 R:0 D:0
 U:1 L:0 R:0 D:0	 U:0 L:0 R:0 D:0	 U:0 L:0 R:0 D:0
 U:-1 L:0 R:0 D:0	 U:0 L:0 R:0 D:0	 U:0 L:0 R:0 D:0

Balance exploration
and exploitation

| 68

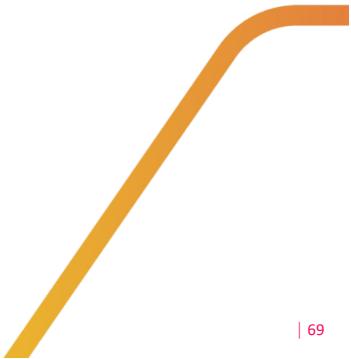


► Q-learning

- Estimate Q^* , by maintaining a Q-table:

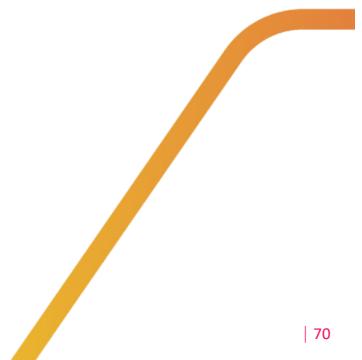
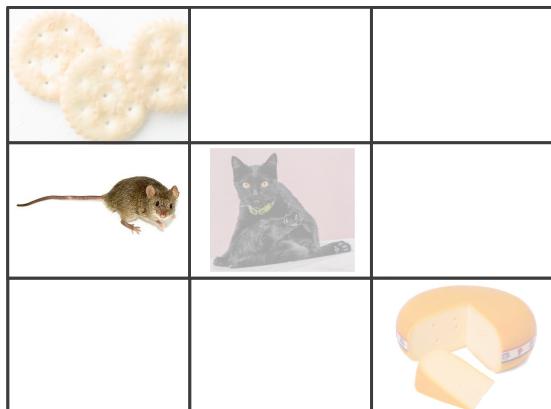
		
		
		

| 69



► Q-learning

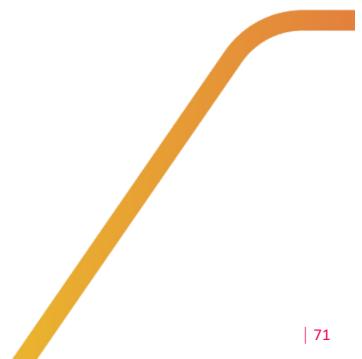
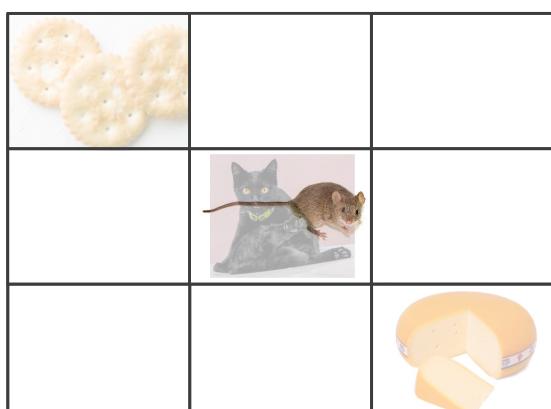
- Estimate Q^* , by maintaining a Q-table:



| 70

► Q-learning

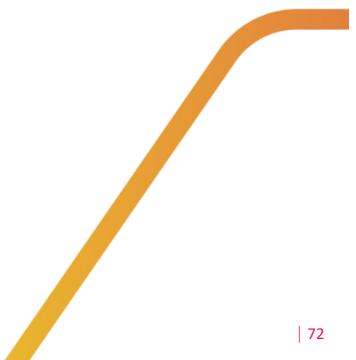
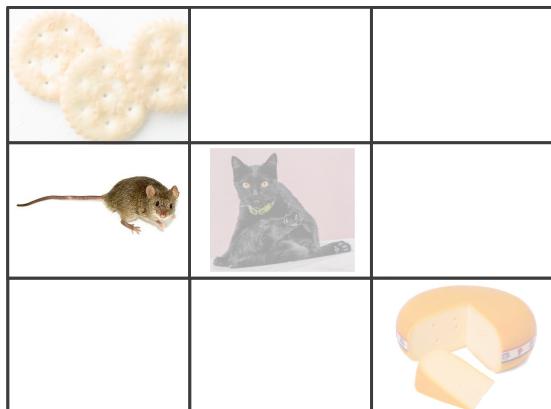
- Estimate Q^* , by maintaining a Q-table:



| 71

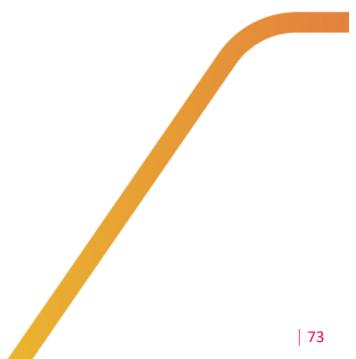
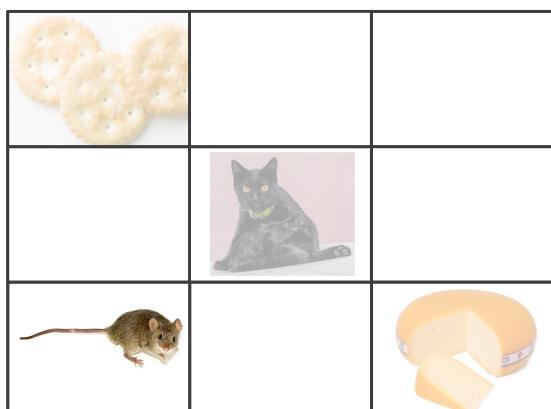
► Q-learning

- Estimate Q^* , by maintaining a Q-table:



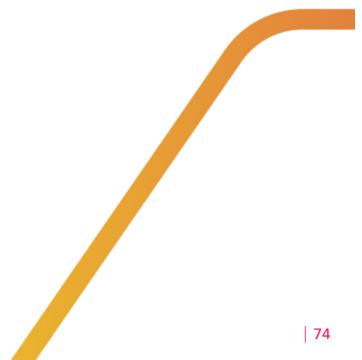
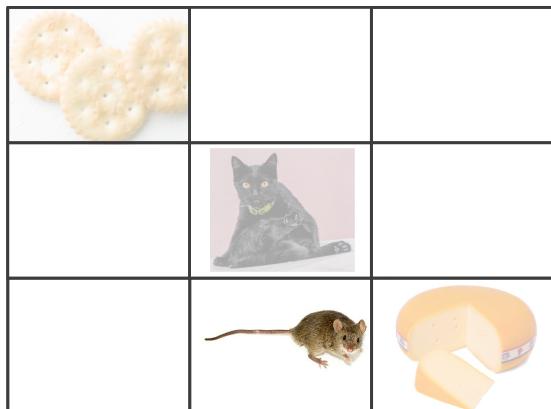
► Q-learning

- Estimate Q^* , by maintaining a Q-table:



► Q-learning

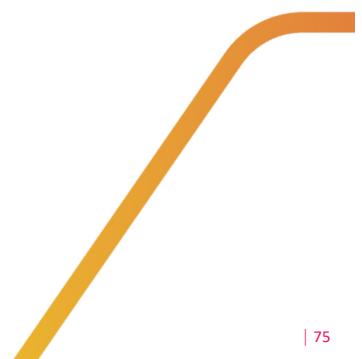
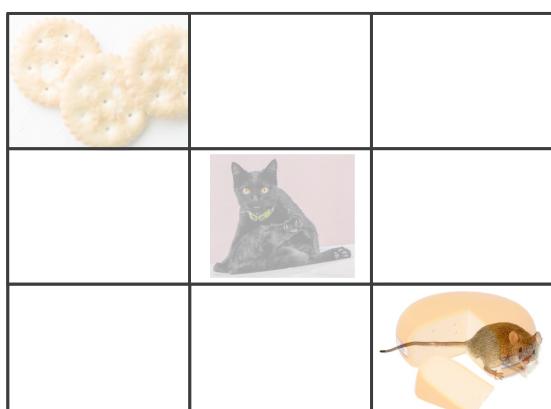
- Estimate Q^* , by maintaining a Q-table:



| 74

► Q-learning

- Estimate Q^* , by maintaining a Q-table:



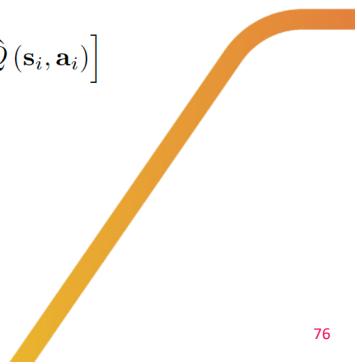
| 75

Q-learning

Watkins, 1999

```
Given: a learning rate  $\alpha$ 
 $\forall s \in \mathcal{S}, a \in \mathcal{A}$  : initialize  $\hat{Q}(s, a)$ 
for each episode  $e$  do
    Initialize  $s_0$ 
    for each step in  $e$ :  $i = 0, 1, \dots$  do
        Choose  $a_i$  from  $s_i$  using a policy derived from  $\hat{Q}(s, a)$ 
        Take action  $a_i$  and observe reward  $r_{i+1}$  and state  $s_{i+1}$ 
         $\hat{Q}(s_i, a_i) \leftarrow \hat{Q}(s_i, a_i) + \alpha [r_{i+1} + \gamma \max_{a'} \hat{Q}(s_{i+1}, a') - \hat{Q}(s_i, a_i)]$ 
         $s_i \leftarrow s_{i+1}$ 
    end
end
```

76



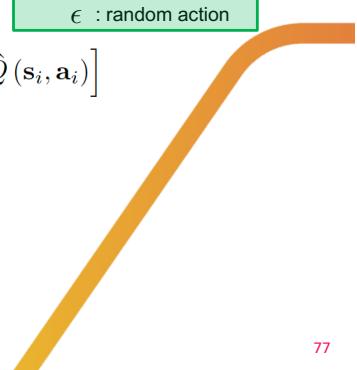
Q-learning

```
Given: a learning rate  $\alpha$ 
 $\forall s \in \mathcal{S}, a \in \mathcal{A}$  : initialize  $\hat{Q}(s, a)$ 
for each episode  $e$  do
    Initialize  $s_0$ 
    for each step in  $e$ :  $i = 0, 1, \dots$  do
        Choose  $a_i$  from  $s_i$  using a policy derived from  $\hat{Q}(s, a)$ 
        Take action  $a_i$  and observe reward  $r_{i+1}$  and state  $s_{i+1}$ 
         $\hat{Q}(s_i, a_i) \leftarrow \hat{Q}(s_i, a_i) + \alpha [r_{i+1} + \gamma \max_{a'} \hat{Q}(s_{i+1}, a') - \hat{Q}(s_i, a_i)]$ 
         $s_i \leftarrow s_{i+1}$ 
    end
end
```

To balance exploration
and exploitation

$1 - \epsilon$: $\max_a \hat{Q}(s, a)$
 ϵ : random action

77



► Q-learning

Given: a learning rate α
 $\forall s \in \mathcal{S}, a \in \mathcal{A} : \text{initialize } \hat{Q}(s, a)$

```

for each episode  $e$  do
    Initialize  $s_0$ 
    for each step in  $e: i = 0, 1, \dots$  do
        Choose  $a_i$  from  $s_i$  using a policy derived from  $\hat{Q}(s, a)$ 
        Take action  $a_i$  and observe reward  $r_{i+1}$  and state  $s_{i+1}$ 
         $\hat{Q}(s_i, a_i) \leftarrow \hat{Q}(s_i, a_i) + \alpha [r_{i+1} + \gamma \max_{a'} \hat{Q}(s_{i+1}, a') - \hat{Q}(s_i, a_i)]$ 
         $s_i \leftarrow s_{i+1}$ 
    end
end

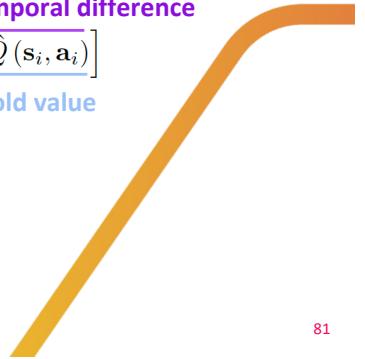
```

TD-learning

$$\hat{Q}(s_i, a_i) \leftarrow \hat{Q}(s_i, a_i) + \alpha [r_{i+1} + \gamma \max_{a'} \hat{Q}(s_{i+1}, a') - \hat{Q}(s_i, a_i)]$$

new value old value

81



► Q-learning

Given: a learning rate α
 $\forall s \in \mathcal{S}, a \in \mathcal{A} : \text{initialize } \hat{Q}(s, a)$

```

for each episode  $e$  do
    Initialize  $s_0$ 
    for each step in  $e: i = 0, 1, \dots$  do
        Choose  $a_i$  from  $s_i$  using a policy derived from  $\hat{Q}(s, a)$ 
        Take action  $a_i$  and observe reward  $r_{i+1}$  and state  $s_{i+1}$ 
         $\hat{Q}(s_i, a_i) \leftarrow \hat{Q}(s_i, a_i) + \alpha [r_{i+1} + \gamma \max_{a'} \hat{Q}(s_{i+1}, a') - \hat{Q}(s_i, a_i)]$ 
         $s_i \leftarrow s_{i+1}$ 
    end
end

```

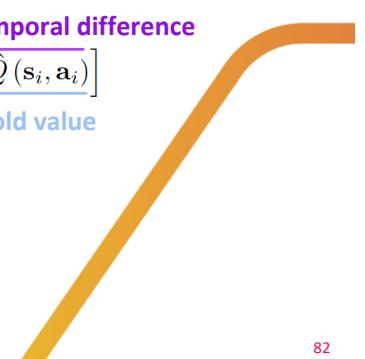
Convergence guarantees

TD-learning

$$\hat{Q}(s_i, a_i) \leftarrow \hat{Q}(s_i, a_i) + \alpha [r_{i+1} + \gamma \max_{a'} \hat{Q}(s_{i+1}, a') - \hat{Q}(s_i, a_i)]$$

new value old value

82



■ Q-learning

- Questions?
- One question for you:
 - “What is the disadvantage of using a Q-table?”

| 85

■ Reinforcement learning, some milestones

- Backgammon RL agent, Tesauro (1992)

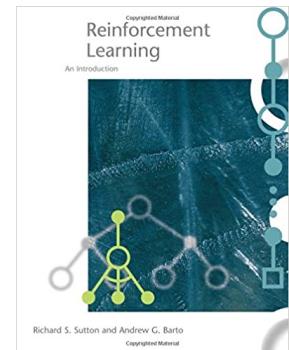


| 87

► Reinforcement learning, some milestones

- Backgammon RL agent, Tesauro (1992)
- Reinforcement Learning book, Sutton and Barto (1998)

RL bible



| 88

► Reinforcement learning, some milestones

- Backgammon RL agent, Tesauro (1992)
- Reinforcement Learning book, Sutton and Barto (1998)
- Atari agent, Deepmind (2015)



| 89

■ Reinforcement learning, some milestones

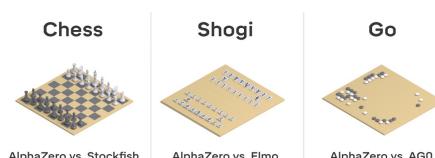
- Backgammon RL agent, Tesauro (1992)
- Reinforcement Learning book, Sutton and Barto (1998)
- Atari agent, *Deepmind* (2015)
- AlphaGo, *Deepmind* (2016)



| 90

■ Reinforcement learning, some milestones

- Backgammon RL agent, Tesauro (1992)
- Reinforcement Learning book, Sutton and Barto (1998)
- Atari agent, *Deepmind* (2015)
- AlphaGo, *Deepmind* (2016)
- AlphaGo Zero, *Deepmind* (2017)



| 91

Reinforcement learning, some milestones

- Backgammon RL agent, Tesauro (1992)
- Reinforcement Learning book, Sutton and Barto (1998)
- Atari agent, *Deepmind* (2015)
- AlphaGo, *Deepmind* (2016)
- AlphaGo Zero, *Deepmind* (2017)
- OpenAI Five, *OpenAI* (2018)



| 92

Reinforcement learning, some milestones

- Backgammon RL agent, Tesauro (1992)
- Reinforcement Learning book, Sutton and Barto (1998)
- Atari agent, *Deepmind* (2015)
- AlphaGo, *Deepmind* (2016)
- AlphaGo Zero, *Deepmind* (2017)
- OpenAI Five, *OpenAI* (2018)
- “Reward is enough”, Silver et al. (2021)

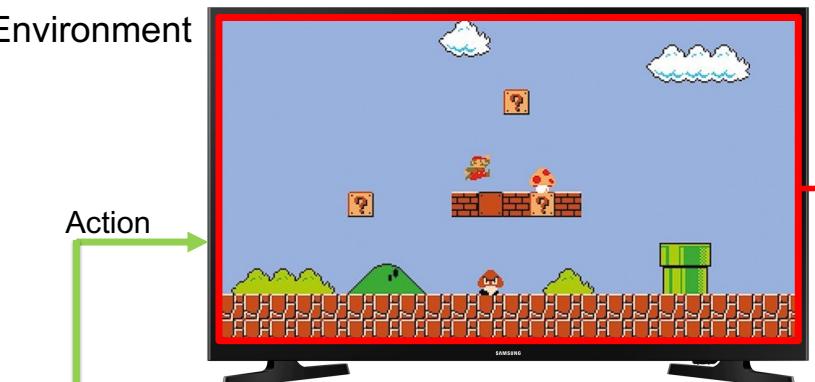


| 93

► Q-learning in a large state space

$$Q(s_t, a_t \mid \theta)$$

Environment



Function approximation

Action



Agent

State/Reward

| 94

► Deep Q-networks (DQN)

$$Q(s_t, a_t \mid \theta)$$

State
s

Action

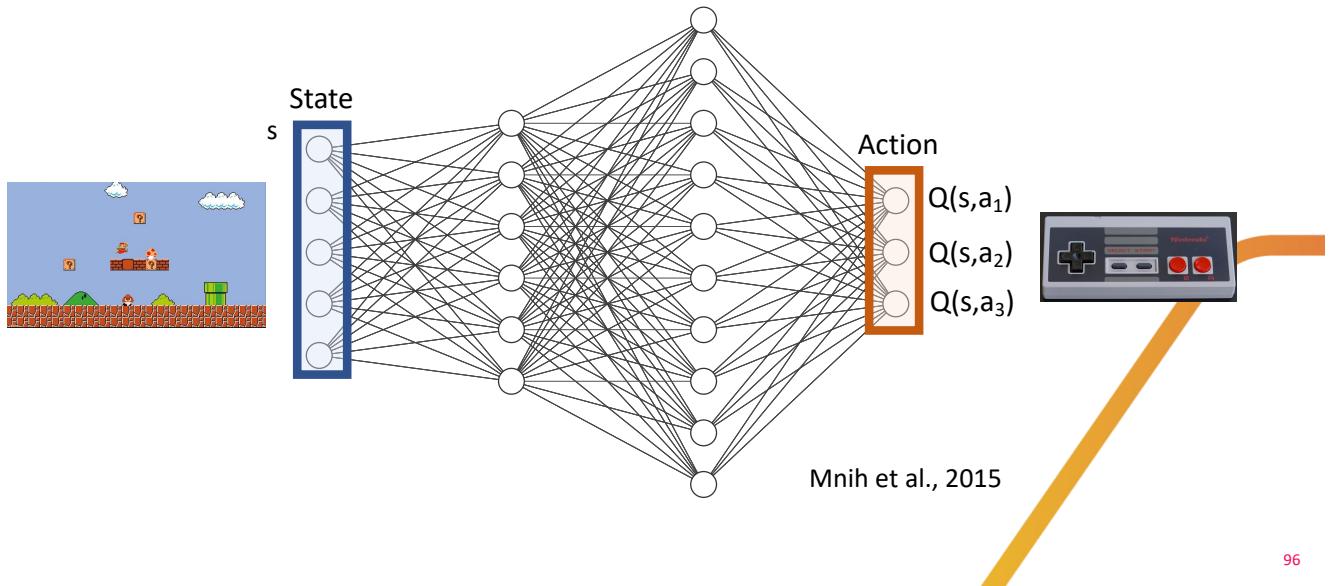
$$\begin{aligned} & Q(s, a_1) \\ & Q(s, a_2) \\ & Q(s, a_3) \end{aligned}$$

Mnih et al., 2015

| 95

Deep Q-networks (DQN)

$$Q(s_t, a_t | \theta)$$



96

Deep Q-networks

$$Q(s_t, a_t | \theta)$$

- Minimize temporal difference error, from mini-batch:

$$\{s_i, a_i, r_i, s_{i+1}\}_{i=0}^N$$

- We could use the mean squared error:

$$L(s, a) = \frac{1}{N} \sum_i (y_i - Q(s_i, a_i | \theta))^2 , \text{ where,}$$

$$y_i = r_i + \gamma \max_a Q(s_{i+1}, a | \theta)$$

Unstable

| 99

Deep Q-networks

$$Q(s_t, a_t \mid \theta)$$

- Minimize temporal difference error, from mini-batch:

$$\{s_i, a_i, r_i, s_{i+1}\}_{i=0}^N$$

- We could use the mean squared error:

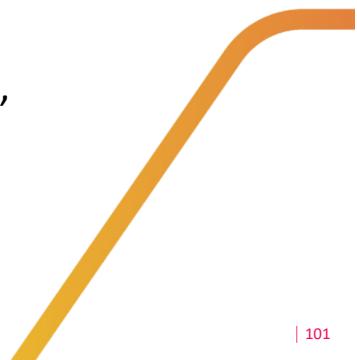
$$L(s, a) = \frac{1}{N} \sum_i (y_i - Q(s_i, a_i \mid \theta))^2, \text{ where,}$$

$$y_i = r_i + \gamma \max_a Q(s_{i+1}, a \mid \phi)$$

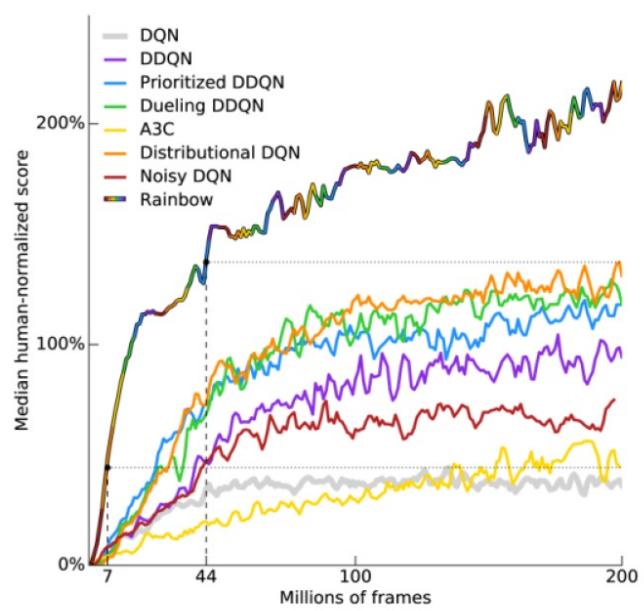
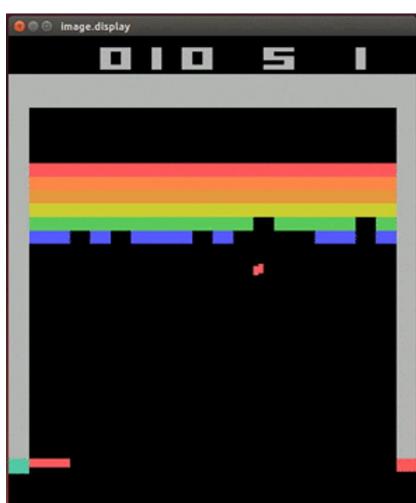
Use a separate target network

Unstable

| 101



Deep Q-networks



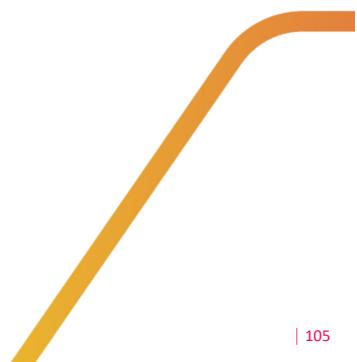
Hessel, 2017

| 103

► Deep Q-networks

- Questions?
- What about large/continuous action spaces? What kind of problem do you see?

| 105



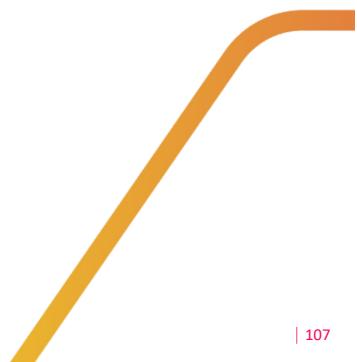
► Policy gradient

- Parameterize a policy directly p_θ
- Optimize the expectation of the policy gradient:

$$\mathbb{E}_{\theta, T} \left[\sum_{t=0}^{\infty} \frac{\nabla_{\theta} \log p_{\theta}(\mathbf{a}_t | \mathbf{s}_t)}{\text{log prob. of the policy's network output}} \hat{A}_t \right]$$

Gradient of the policy
log prob.
of the
policy's
network
output
Advantage

| 107

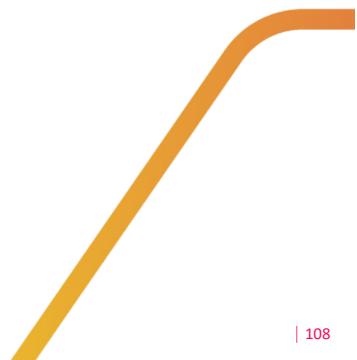


Policy gradient

- Parameterize a policy directly p_{θ}
- Optimize the expectation of the policy gradient:

$$\mathbb{E}_{p_{\theta}, T} \left[\sum_{t=0}^{\infty} \nabla_{\theta} \log p_{\theta}(\mathbf{a}_t | \mathbf{s}_t) \hat{A}_t \right]$$
$$\hat{A}_t = \underbrace{\left[\sum_{i=0}^{\infty} d^i R(\mathbf{s}_{t+i}, \mathbf{a}_{t+i}, \mathbf{s}_{t+i+1}) \right]}_{\text{Return}} - \underbrace{V^p(\mathbf{s}_t)}_{\text{Value function}}$$

| 108



Policy gradient

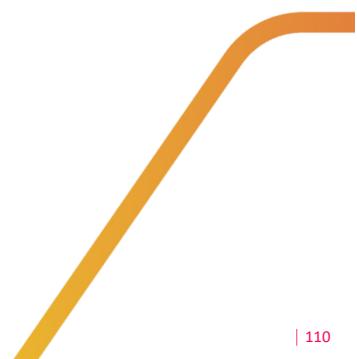
- Parameterize a policy directly p_{θ}
- Optimize the expectation of the policy gradient:

$$\mathbb{E}_{p_{\theta}, T} \left[\sum_{t=0}^{\infty} \nabla_{\theta} \log p_{\theta}(\mathbf{a}_t | \mathbf{s}_t) \hat{A}_t \right]$$
$$\hat{A}_t = \left[\sum_{i=0}^{\infty} d^i R(\mathbf{s}_{t+i}, \mathbf{a}_{t+i}, \mathbf{s}_{t+i+1}) \right] - V^p(\mathbf{s}_t)$$

- Estimate using:

$$\hat{g} = \frac{1}{|Y|} \sum_{y \in Y} \sum_t \nabla_{\theta} \log p_{\theta}(\mathbf{a}_t | \mathbf{s}_t) \hat{A}_t$$

| 110

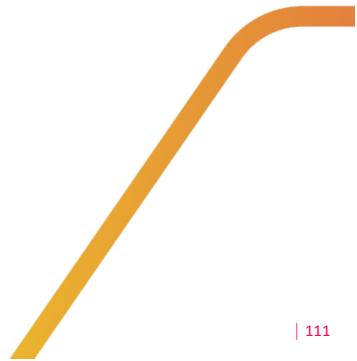


Policy gradient

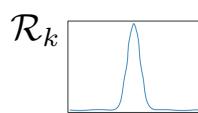
Williams, 1992

Given: a learning rate $\alpha(\cdot)$
Initialize policy parameter θ and value function V
for $i = 1, 2 \dots$ **do**
 Collect a set of trajectories Y using the current policy p_{θ_i}
 for each time step t and each trajectory $y \in Y$ **do**
 | Compute the advantage \hat{A}_t
 end
 Re-fit $V^p(s_t)$
 Compute the policy gradient estimate \hat{g}_i
 Update the policy
end

| 111



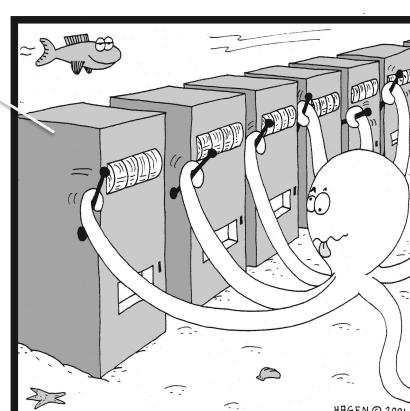
Multi-armed bandits



Set of arms: $\{A_k\}_k^K$

Stochastic reward: r_k

Arm's mean: $\mathbb{E}[r_k] = \mu_k$

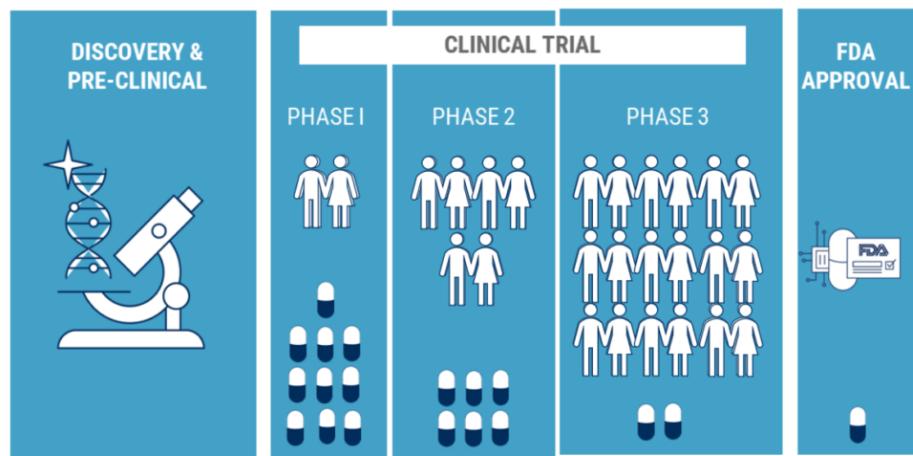


| 115



► Multi-armed bandit application

- Clinical trials



| 116

► Cumulative regret

- ϵ - greedy
 - Exploit with probability $1 - \epsilon$
- UCB
 - $A_t = \operatorname{argmax}_k \left[\hat{\mu}_k + c \sqrt{\frac{\ln(t)}{N_t(k)}} \right]$



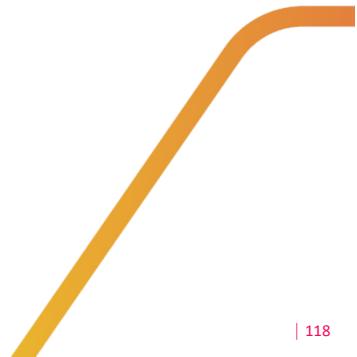
| 117

■ Use prior knowledge: Bayesian philosophy

- It's all about belief
 - Prior belief
 - Updating our belief, given data

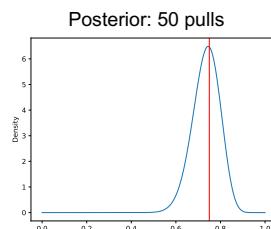
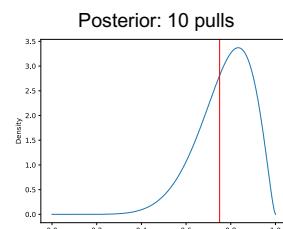
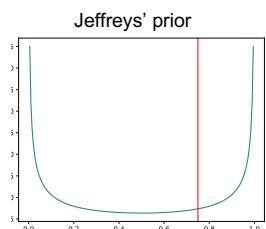


| 118



■ Bayesian bandits

- Prior belief over means: $\pi(\cdot)$
- History: $\mathcal{H}^{(t-1)} = \{a^{(i)}, r^{(i)}\}_{i=1}^{t-1}$
- Posterior: $\pi(\cdot \mid \mathcal{H}^{(t-1)})$



| 122



■ Cumulative regret: Thompson sampling

Given: $\pi(\cdot)$ and $\mathcal{H}^{(0)} = \emptyset$

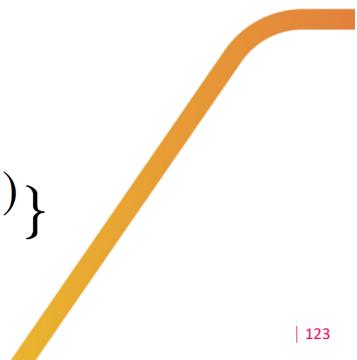
for $t = 1, \dots, +\infty$ **do**

$\theta^{(t)} \sim \pi(\cdot | \mathcal{H}_{t-1})$
 $a^{(t)} = \sigma_1(\theta^{(t)})$
 $r^{(t)} \leftarrow \text{Pull arm } a^{(t)}$
 $\mathcal{H}^{(t)} \leftarrow \mathcal{H}^{(t-1)} \cup \{a^{(t)}, r^{(t)}\}$

end

(Thompson, 1933; Chapelle, 2011)

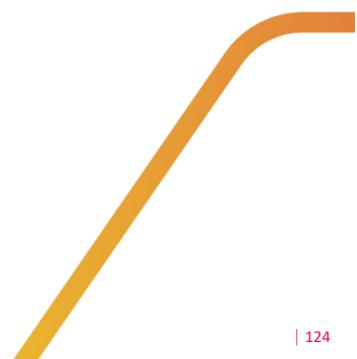
| 123



■ Pure exploration

- Decision making
 - e.g., Simulation of prevention strategies in a compute-intensive model
- Best-arm identification
 - fixed budget

| 124

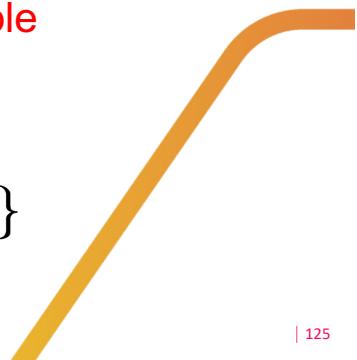


► Top-two Thompson sampling

```
Given:  $\pi(\cdot)$  and  $\mathcal{H}^{(0)} = \emptyset$ 
for  $t = 1, \dots, +\infty$  do
     $\theta^{(t)} \sim \pi(\cdot | \mathcal{H}_{t-1})$ 
     $a^{(t)} = \sigma_1(\theta^{(t)})$  Resample
     $r^{(t)} \leftarrow$  Pull arm  $a^{(t)}$ 
     $\mathcal{H}^{(t)} \leftarrow \mathcal{H}^{(t-1)} \cup \{a^{(t)}, r^{(t)}\}$ 
end
```

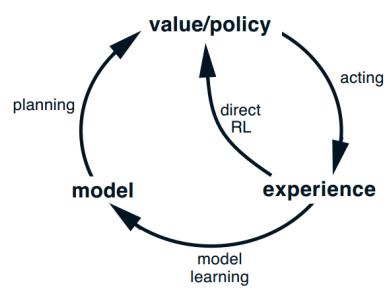
(Russo, 2016)

| 125

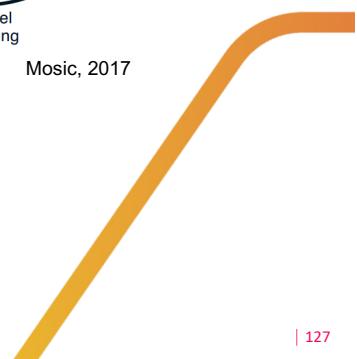


► Things we did *not* talk about

- Mixing planning and learning



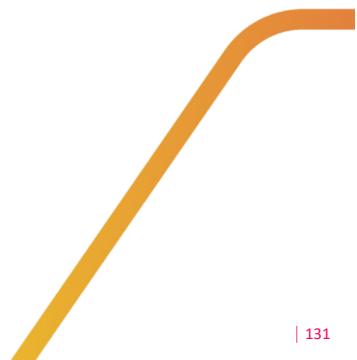
| 127



► Things we did *not* talk about

- Mixing planning and learning
- Multi-objective RL
- Multi-agent RL
- Safety
- Partially observable MDPs

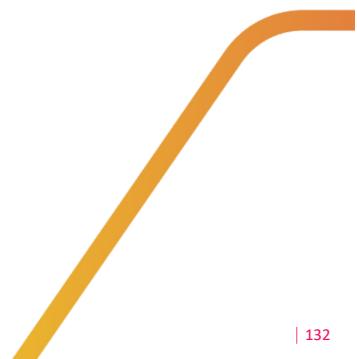
| 131



► Reading material

- RL book by Sutton and Barto
 - Mixing planning and learning
- Connecting RL with biomed:
 - Guiding the mitigation of epidemics with RL, Libin, 2020
 - <https://ai.vub.ac.be/education/research/phd-theses>
- Background and code
 - <https://spinningup.openai.com/>

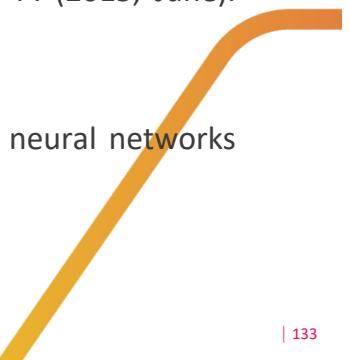
| 132



■ Reading material

- DQN paper
 - Mnih, Volodymyr, et al. "Human-level control through deep reinforcement learning." *Nature* 518.7540 (2015): 529-533.
- Policy gradient
 - Schulman, J., Levine, S., Abbeel, P., Jordan, M., & Moritz, P. (2015, June). Trust region policy optimization. In *ICML* (pp. 1889-1897).
- AlphaGo paper
 - Silver, David, et al. "Mastering the game of Go with deep neural networks and tree search." *Nature* 529.7587 (2016): 484-489.

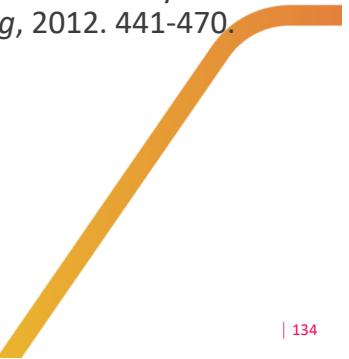
| 133



■ Reading material

- Multi-objective RL:
 - Ruijters, Diederik M., et al. "A survey of multi-objective sequential decision-making." *Journal of Artificial Intelligence Research* 48 (2013): 67-113.
- Multi-agent RL:
 - Nowé, Ann, Peter Vrancx, and Yann-Michaël De Hauwere. "Game theory and multi-agent reinforcement learning." *Reinforcement Learning*, 2012. 441-470.

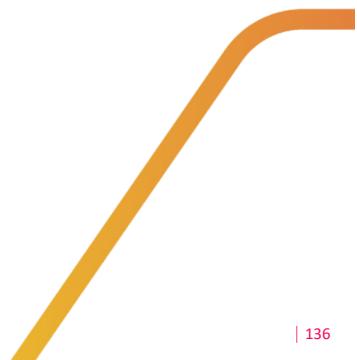
| 134



► Graphic credits

- Cats/dogs/crackers/drugs/ICU/bike/robot/LHC/beer/westvleteren/Las Vegas: unsplash
- Mouse/mouse trap/Atari/Bayes: Wikipedia
- Cheese: Kaasdok.nl
- Go board game: <https://www.go-jigs.eu/what-is-go/>
- Lee Sedol: New Yorker
- Clinical trials: <https://www.cbinsights.com/>
- Telecom application: Ericson website
- ICU covid: sciensano

| 136



► Thank you! Any more questions?



 Pieter.Libin@vub.be
 @PieterLibin

| 137

