



Universidade do Minho
Escola de Engenharia
Departamento de Informática

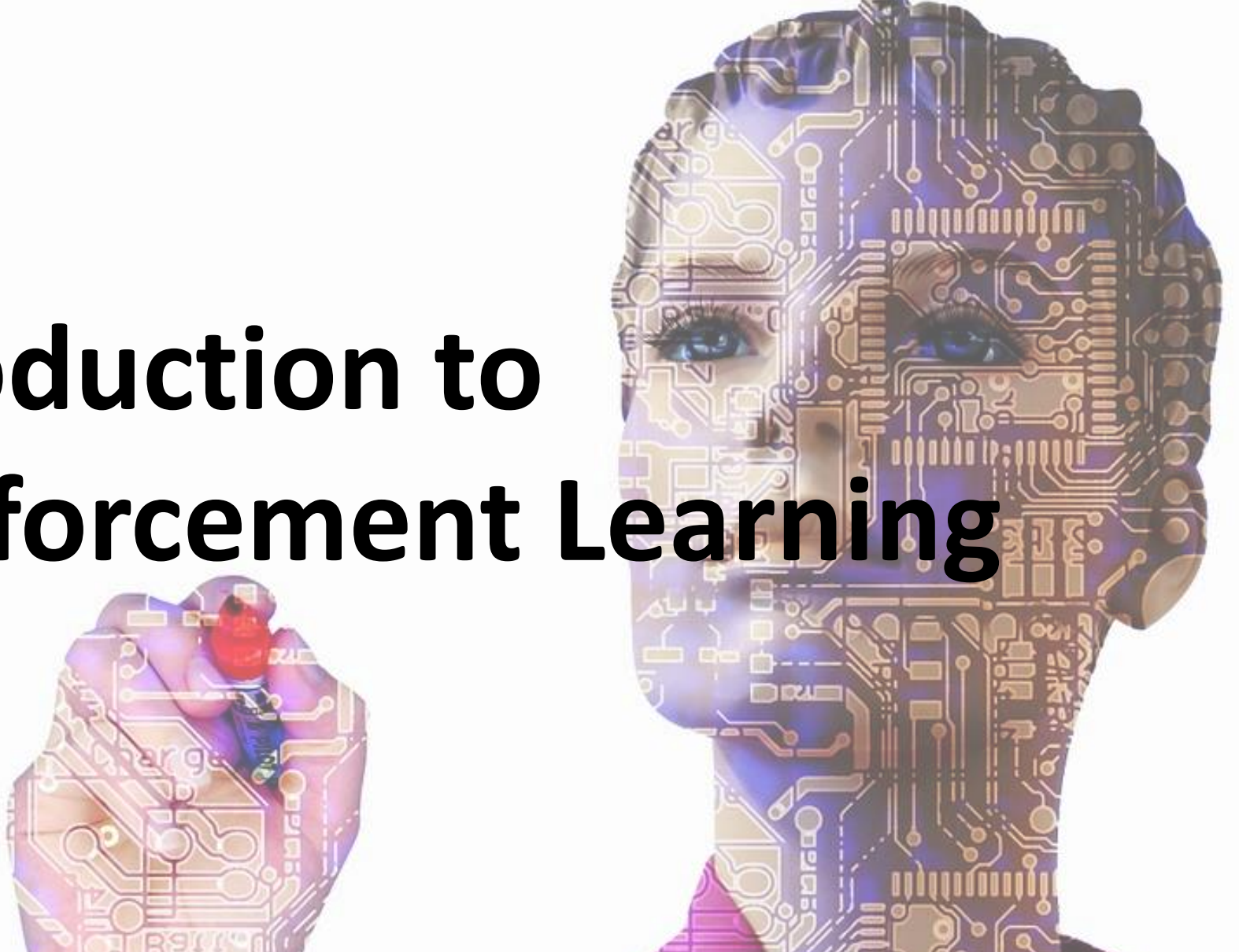
Mestrado Integrado em Engenharia Informática
Mestrado em Engenharia Informática
Computação Natural
2019/2020

Paulo Novais, Cesar Analide, Filipe Gonçalves

- Paulo Novais – pjon@di.uminho.pt
- Cesar Analide – cesar.analide@di.uminho.pt
- Filipe Gonçalves – fgoncalves@algoritmi.uminho.pt

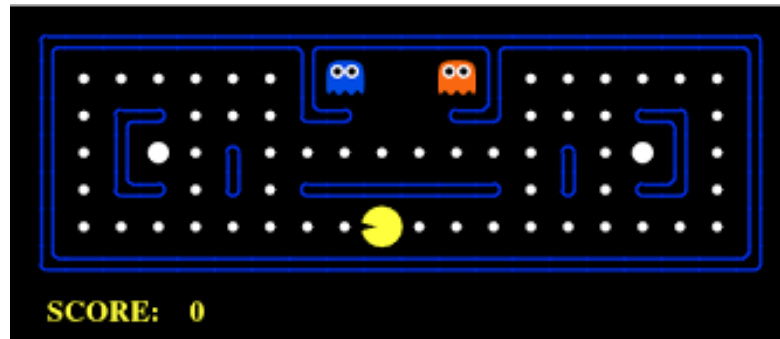
- Departamento de Informática
Escola de Engenharia
Universidade do Minho
- Grupo ISLab – (Synthetic Intelligence Lab)
- Centro ALGORITMI
Universidade do Minho

Introduction to Reinforcement Learning



- Reinforcement Learning supports automation by learning from the environment it is present in
 - So does Machine Learning and Deep Learning, using different strategies, with automation support
- Deep Learning and Machine Learning are learning processes, but which are most focused on finding patterns in the existing data
- Reinforcement Learning, on the other hand, learns by trial and error method, and eventually, gets to the right actions or the global optimum
 - Pros: it is not required to provide the whole training data as in Supervised Learning. Instead, a few chunks would suffice

- You have some sort of agent that “explores” some space
- As it goes, it learns the value of diferente state changes in diferente conditions
- Those values inform subsequeunte behaviour of the agent
- Examples:
 - Pac-Man
 - Cat & Mouse Game
 - Multi-armed Bandit problem
- Yields fast on-line performance once the space has been explored



- Applications where reinforcement systems are applied:
 - Self Driving Cars
 - Gaming
 - Robotics
 - Recommendation Systems
 - Advertising and Marketing

- The Multi-armed Bandit Problem (Exploration vs Exploitation Problem)



D1



D2



D3



D4



D5

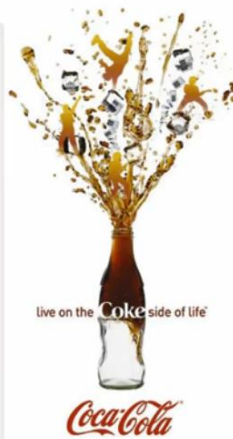
- The Multi-armed Bandit Problem (Exploration vs Exploitation Problem)



■ Marketing - Ads Selection (Exploration vs Exploitation Problem)



D1



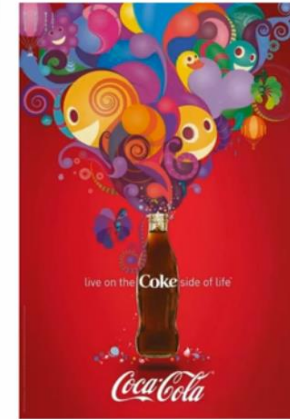
D2



D3



D4



D5

■ **Marketing - Ads Selection (Exploration vs Exploitation Problem)**

- We have d arms. For example, arms are ads that we display to users each time they connect to a web page
- Each time a user connects to this web page, that makes a round
- At each round n , we choose one ad to display to the user
- At each round n , ad i gives reward:
 $r_i(n) \in \{0, 1\}$: $r_i(n) = 1$ if the user clicked on the ad i , 0 if the user didn't
- Goal: maximize the total reward we get over many rounds

■ **Marketing - Ads Selection**

- Index: Person
- Column: Ads
 - 1: Person clicked on Ad
 - 0: Person ignored Ad

Index	Ad 1	Ad 2	Ad 3	Ad 4	Ad 5	Ad 6	Ad 7	Ad 8	Ad 9	Ad 10
0	1	0	0	0	1	0	0	0	1	0
1	0	0	0	0	0	0	0	0	1	0
2	0	0	0	0	0	0	0	0	0	0
3	0	1	0	0	0	0	0	1	0	0
4	0	0	0	0	0	0	0	0	0	0
5	1	1	0	0	0	0	0	0	0	0
6	0	0	0	1	0	0	0	0	0	0
7	1	1	0	0	1	0	0	0	0	0
8	0	0	0	0	0	0	0	0	0	0
9	0	0	1	0	0	0	0	0	0	0
10	0	0	0	0	0	0	0	0	0	0
11	0	0	0	0	0	0	0	0	0	0
12	0	0	0	1	0	0	0	0	0	0
13	0	0	0	0	0	0	0	0	1	0
14	0	0	0	0	0	0	0	1	0	0
15	0	0	0	0	1	0	0	1	0	0
16	0	0	0	0	0	0	0	0	0	0
17	0	0	0	0	0	0	0	0	0	0
18	0	0	0	0	0	0	0	1	0	0
19	0	0	0	0	0	0	0	0	1	0
20	0	1	0	0	0	0	0	1	0	0
21	0	0	0	0	1	0	0	0	0	1

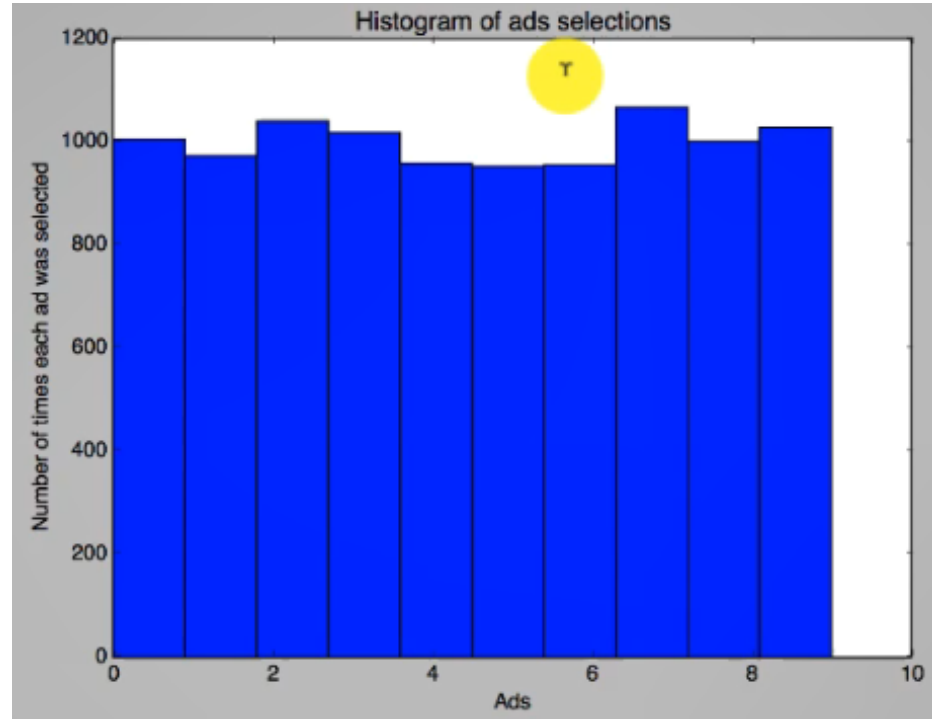
- Marketing - Ads Selection
 - Random Selection

```
1 # Random Selection
2
3 # Importing the libraries
4 import numpy as np
5 import matplotlib.pyplot as plt
6 import pandas as pd
7
8 # Importing the dataset
9 dataset = pd.read_csv('Ads_CTR_Optimisation.csv')
10
11 # Implementing Random Selection
12 import random
13 N = 10000
14 d = 10
15 ads_selected = []
16 total_reward = 0
17 for n in range(0, N):
18     ad = random.randrange(d)
19     ads_selected.append(ad)
20     reward = dataset.values[n, ad]
21     total_reward = total_reward + reward
22
23 # Visualising the results - Histogram
24 plt.hist(ads_selected)
25 plt.title('Histogram of ads selections')
26 plt.xlabel('Ads')
27 plt.ylabel('Number of times each ad was selected')
28 plt.show()
```

■ Marketing - Ads Selection

○ Random Selection

j ▲	Type	Size	Value
0	int	1	4
1	int	1	6
2	int	1	1
3	int	1	0
4	int	1	4
5	int	1	5
6	int	1	0
7	int	1	8
8	int	1	5
9	int	1	3
10	int	1	8
11	int	1	5
12	int	1	8
13	int	1	4
14	int	1	0



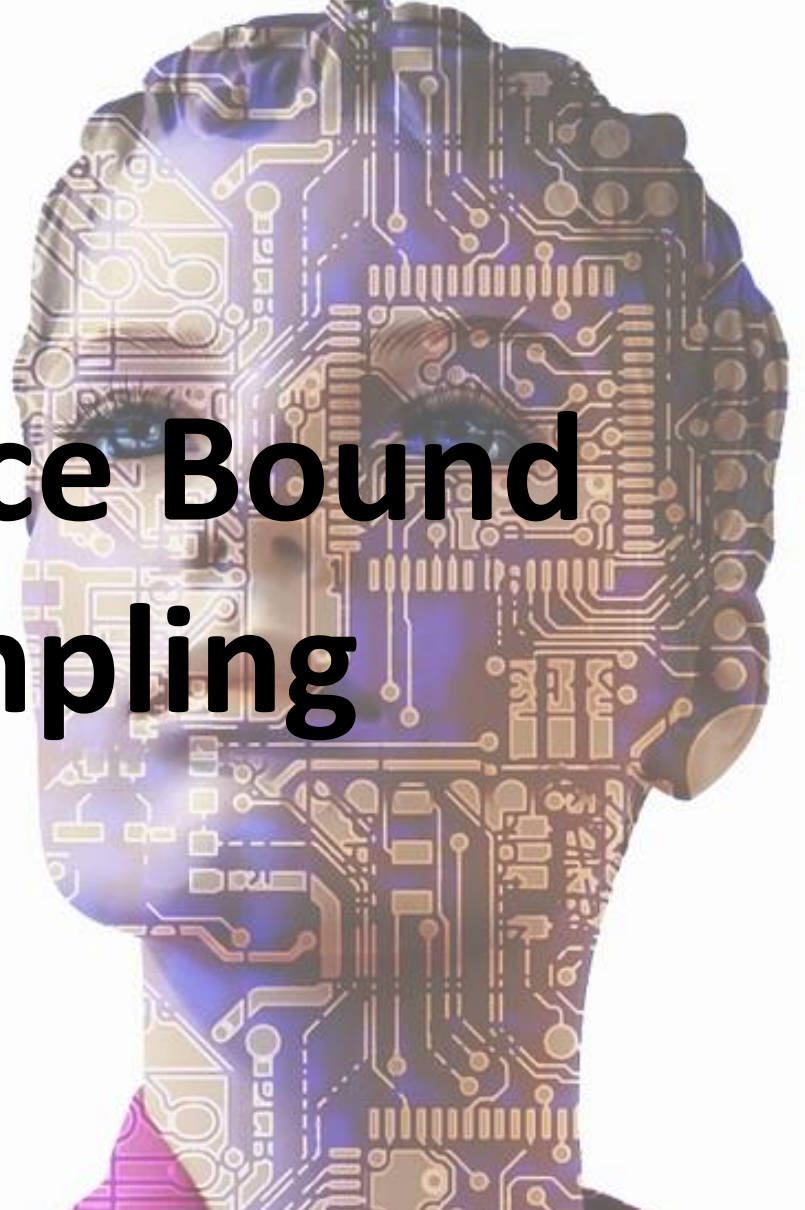
total_reward

int64

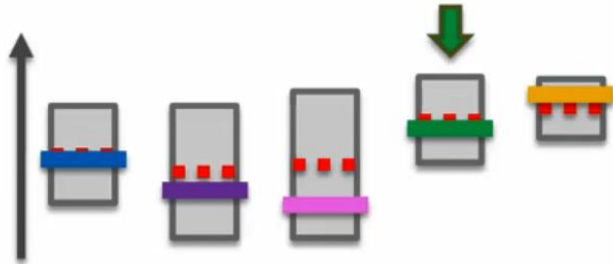
1

1193

Upper Confidence Bound vs Thomson Sampling

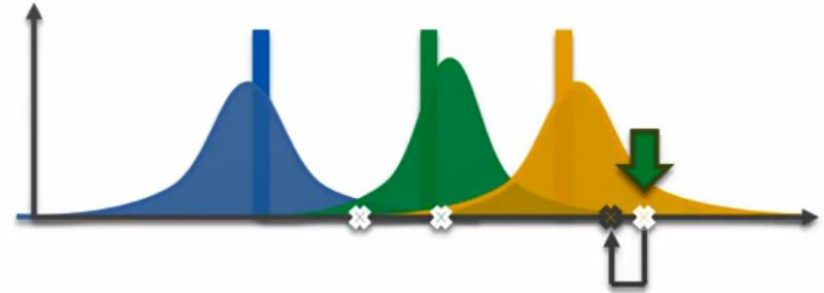


UCB



- Deterministic
- Requires update at every round

Thompson Sampling



- Probabilistic
- Can accommodate delayed feedback
- Better empirical evidence

■ **Upper Confidence Bound Algorithm**

Step 1. At each round n , we consider two numbers for each ad i :

- $N_i(n)$ - the number of times the ad i was selected up to round n ,
- $R_i(n)$ - the sum of rewards of the ad i up to round n .

Step 2. From these two numbers we compute:

- the average reward of ad i up to round n

$$\bar{r}_i(n) = \frac{R_i(n)}{N_i(n)}$$

- the confidence interval $[\bar{r}_i(n) - \Delta_i(n), \bar{r}_i(n) + \Delta_i(n)]$ at round n with

$$\Delta_i(n) = \sqrt{\frac{3 \log(n)}{2 N_i(n)}}$$

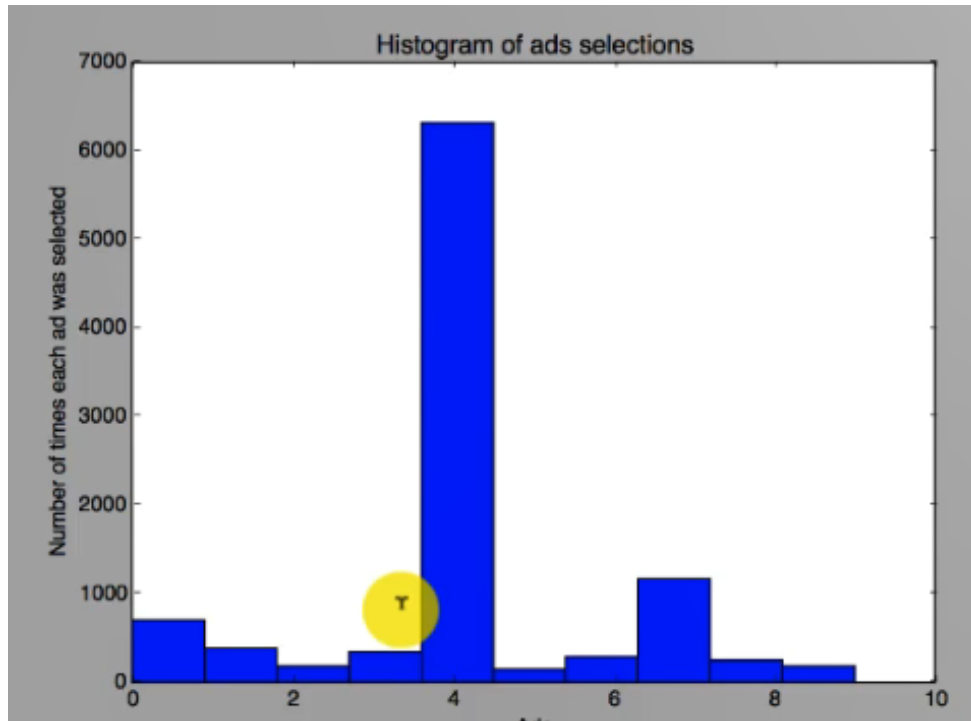
Step 3. We select the ad i that has the maximum UCB $\bar{r}_i(n) + \Delta_i(n)$.

- Marketing - Ads Selection
 - Upper Confidence Bound

```
1 # Upper Confidence Bound
2
3 # Importing the libraries
4 import numpy as np
5 import matplotlib.pyplot as plt
6 import pandas as pd
7
8 # Importing the dataset
9 dataset = pd.read_csv('Ads_CTR_Optimisation.csv')
10
11 # Implementing UCB
12 import math
13 N = 10000
14 d = 10
15 ads_selected = []
16 numbers_of_selections = [0] * d
17 sums_of_rewards = [0] * d
18 total_reward = 0
19 for n in range(0, N):
20     ad = 0
21     max_upper_bound = 0
22     for i in range(0, d):
23         if (numbers_of_selections[i] > 0):
24             average_reward = sums_of_rewards[i] / numbers_of_selections[i]
25             delta_i = math.sqrt(3/2 * math.log(n + 1) / numbers_of_selections[i])
26             upper_bound = average_reward + delta_i
27         else:
28             upper_bound = 1e400
29         if upper_bound > max_upper_bound:
30             max_upper_bound = upper_bound
31             ad = i
32     ads_selected.append(ad)
33     numbers_of_selections[ad] = numbers_of_selections[ad] + 1
34     reward = dataset.values[n, ad]
35     sums_of_rewards[ad] = sums_of_rewards[ad] + reward
36     total_reward = total_reward + reward
37
```

- Marketing - Ads Selection
 - Upper Confidence Bound

i	Type	Size	Value
9985	int	1	4
9986	int	1	4
9987	int	1	4
9988	int	1	4
9989	int	1	4
9990	int	1	4
9991	int	1	4
9992	int	1	4
9993	int	1	4
9994	int	1	4
9995	int	1	4
9996	int	1	4
9997	int	1	4
9998	int	1	4
9999	int	1	4



total_reward	int64	1	2178
upper_bound	float64	1	0.31017236647899182

■ **Thompson Sampling Algorithm**

Step 1. At each round n , we consider two numbers for each ad i :

- $N_i^1(n)$ - the number of times the ad i got reward 1 up to round n ,
- $N_i^0(n)$ - the number of times the ad i got reward 0 up to round n .

Step 2. For each ad i , we take a random draw from the distribution below:

$$\theta_i(n) = \beta(N_i^1(n) + 1, N_i^0(n) + 1)$$

Step 3. We select the ad that has the highest $\theta_i(n)$.

■ Thompson Sampling Algorithm

- Ad i gets rewards \mathbf{y} from Bernoulli distribution $p(\mathbf{y}|\theta_i) \sim \mathcal{B}(\theta_i)$.
- θ_i is unknown but we set its uncertainty by assuming it has a uniform distribution $p(\theta_i) \sim \mathcal{U}([0, 1])$, which is the prior distribution.
- Bayes Rule: we approach θ_i by the posterior distribution

$$\underbrace{p(\theta_i|\mathbf{y})}_{\text{posterior distribution}} = \frac{p(\mathbf{y}|\theta_i)p(\theta_i)}{\int p(\mathbf{y}|\theta_i)p(\theta_i)d\theta_i} \propto \underbrace{p(\mathbf{y}|\theta_i)}_{\text{likelihood function}} \times \underbrace{p(\theta_i)}_{\text{prior distribution}}$$

- We get $p(\theta_i|\mathbf{y}) \sim \beta(\text{number of successes} + 1, \text{number of failures} + 1)$
- At each round n we take a random draw $\theta_i(n)$ from this posterior distribution $p(\theta_i|\mathbf{y})$, for each ad i .
- At each round n we select the ad i that has the highest $\theta_i(n)$.

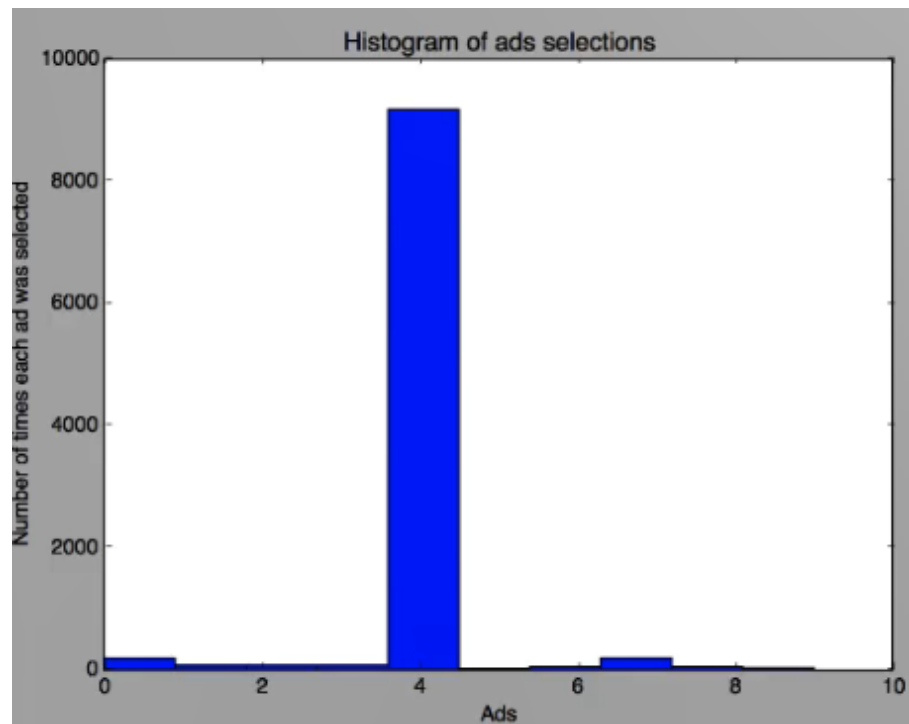
- Marketing - Ads Selection
 - Thompson Sampling

```
1 # Thompson Sampling
2
3 # Importing the libraries
4 import numpy as np
5 import matplotlib.pyplot as plt
6 import pandas as pd
7
8 # Importing the dataset
9 dataset = pd.read_csv('Ads_CTR_Optimisation.csv')
10
11 # Implementing Thompson Sampling
12 import random
13 N = 10000
14 d = 10
15 ads_selected = []
16 numbers_of_rewards_1 = [0] * d
17 numbers_of_rewards_0 = [0] * d
18 total_reward = 0
19 for n in range(0, N):
20     ad = 0
21     max_random = 0
22     for i in range(0, d):
23         random_beta = random.betavariate(numbers_of_rewards_1[i] + 1, numbers_of_rewards_0[i] + 1)
24         if random_beta > max_random:
25             max_random = random_beta
26             ad = i
27     ads_selected.append(ad)
28     reward = dataset.values[n, ad]
29     if reward == 1:
30         numbers_of_rewards_1[ad] = numbers_of_rewards_1[ad] + 1
31     else:
32         numbers_of_rewards_0[ad] = numbers_of_rewards_0[ad] + 1
33     total_reward = total_reward + reward
34
```

■ Marketing - Ads Selection

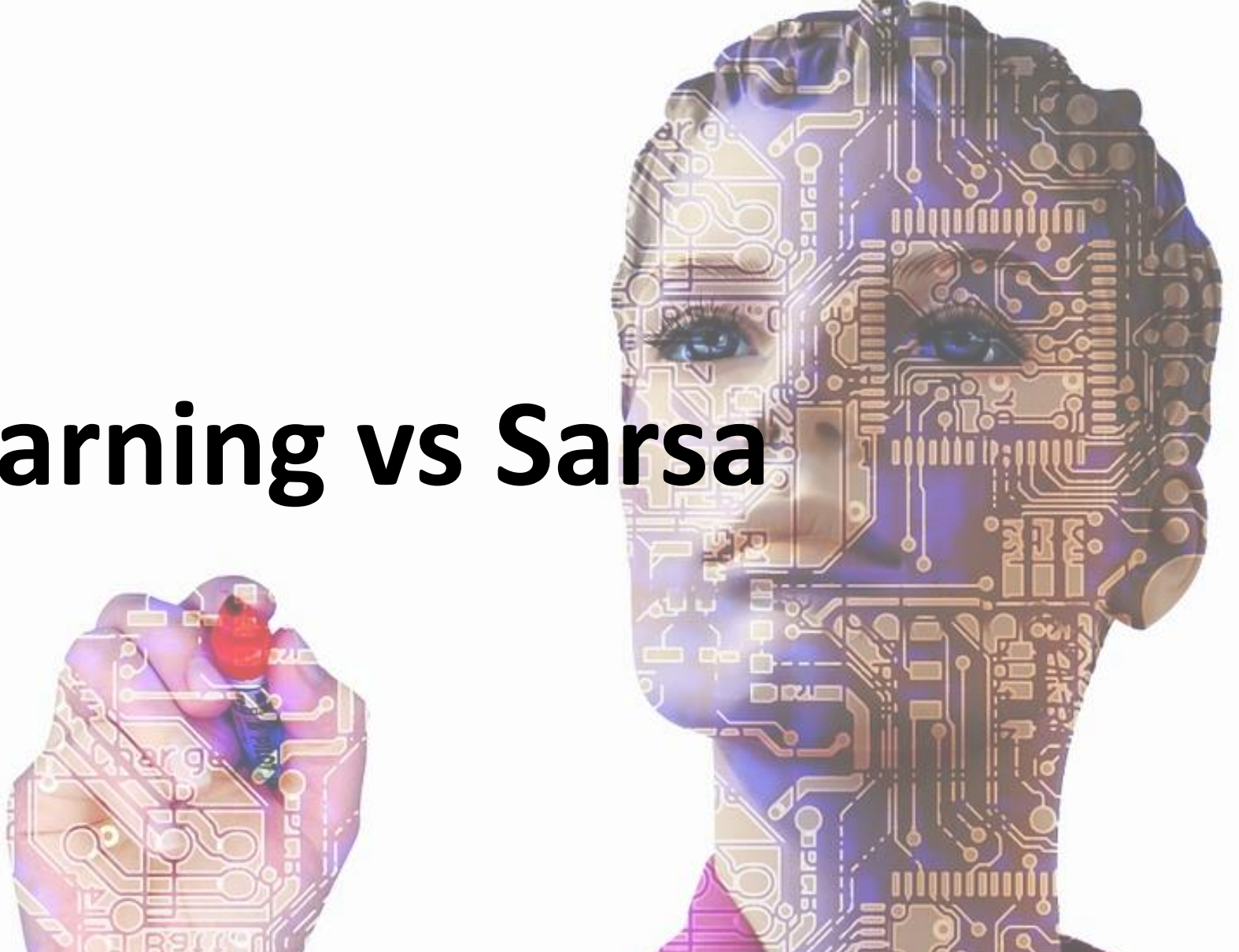
○ Thompson Sampling

i	Type	Size	Value
9985	int	1	4
9986	int	1	4
9987	int	1	4
9988	int	1	4
9989	int	1	4
9990	int	1	4
9991	int	1	4
9992	int	1	4
9993	int	1	4
9994	int	1	4
9995	int	1	4
9996	int	1	4
9997	int	1	4
9998	int	1	4
9999	int	1	4

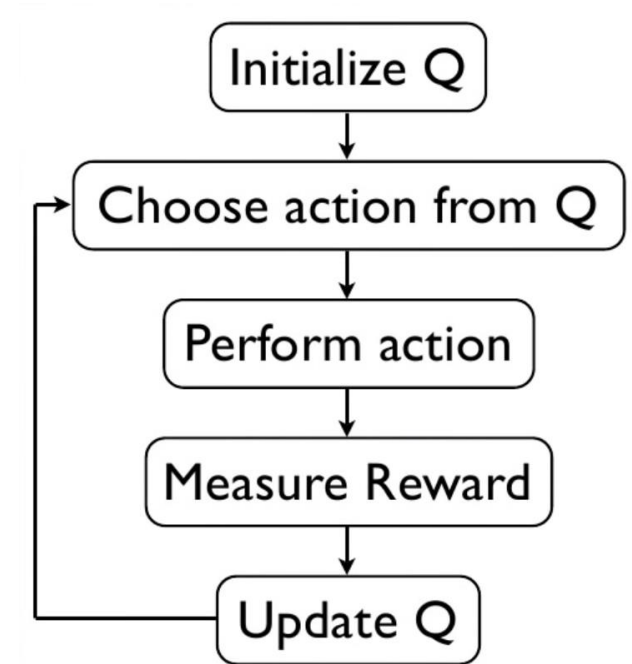


total_reward	int64	1	2613
--------------	-------	---	------

Q-Learning vs Sarsa



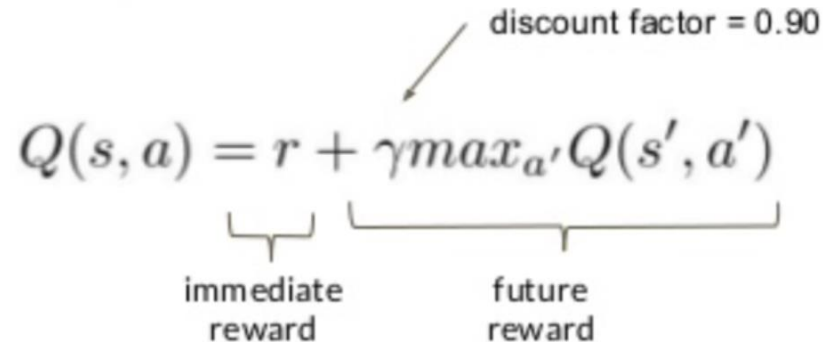
- A specific implementation of reinforcement learning:
- Determined by:
 - Set of environmental states S (called State)
 - Set of possible actions in those states A (called Actions)
 - Value of each state/action Q (called Q-value or Action-value)
- Start off with Q values of 0 / random-values
- Explore the space
- As bad things happen after a given state/actions, reduce its Q -value
- As rewards happen after a given state/action, increase its Q -value



- In RL we want to obtain a function $Q(S,A)$ that predicts the best action A in state S in order to maximize a cumulative reward
- This function can be estimated using Q-learning, which iteratively updates $Q(S,A)$ using the Bellman Equation

$$Q(s, a) = \underbrace{r}_{\text{immediate reward}} + \underbrace{\gamma \max_{a'} Q(s', a')}_{\text{future reward}}$$

discount factor = 0.90



The diagram shows the Bellman Equation $Q(s, a) = r + \gamma \max_{a'} Q(s', a')$. A bracket under the r is labeled "immediate reward". A bracket under the $\gamma \max_{a'} Q(s', a')$ term is labeled "future reward". An arrow points from the text "discount factor = 0.90" to the γ in the equation.

Pac-man game exemple (analyse the figure below):

- What are some state/actions here?
 - Pac-man has a wall to the West
 - Pac-man dies if he moves one step South
 - Pac-man just continues to live if going North or East
- You can “look ahead” more than one step by using a discount factor when computing Q (where S is previous state, S' is current state)



- Markov decision processes (MDPs) provide a mathematical framework for modeling decision making in situations where outcomes are partly random and partly under the control of a decision maker
 - MDP's are just a way to describe what was mentioned using mathematical notation
- States are still described as S and S'
- State transition functions are described as: $P_a(s, s')$
- “Q” values are described as a reward function: $R_a(s, s')$

- The Markov Decision Process can be designed as:
 - Set of states, S
 - Set of actions, A
 - Reward function, R
 - Policy, π
 - Value, V
- We take an action (A) to transition from our start state to our end state (S) -> getting rewards (R) for each action we take
- Our actions can lead to a positive reward or negative reward
- The set of actions we took define our policy (π) and the rewards we get in return defines our value (V)
 - i.e. **a policy is a mapping from state to action**

■ Take into account the following steps:

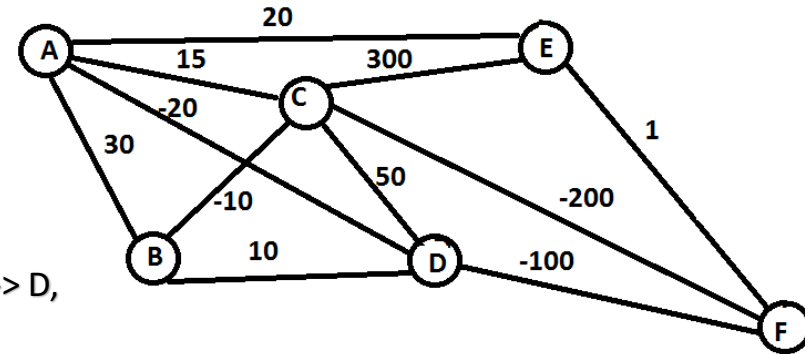
- Initialize Q Matrix (defines reward matrix, where lines define States and columns define Actions to transport to other States)
- Choose action from Q
- Perform action
- Measure Reward and Update Q
- Repeat

		Actions					
		0	1	2	3	4	5
States	0	0	0	0	0	0	0
	1	0	0	0	0	0	100
	2	0	0	0	0	0	0
	3	0	80	0	0	0	0
	4	0	0	0	0	0	0
	5	0	0	0	0	0	0

Q-values Matrix

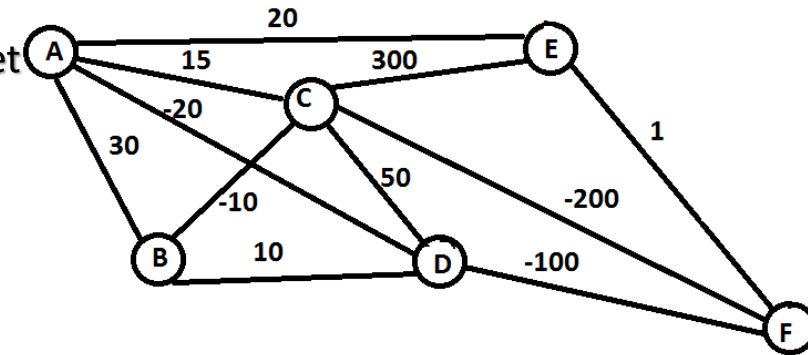
Example: Shortest path problem

- Task: Go from A to F, with as low cost as possible
- Numbers at each edge between two places represent the cost taken to traverse the distance
- Negative cost are actually some earnings on the way
- Value is the total cumulative reward when you do a policy:
 - The set of states are the nodes: {A, B, C, D, E, F}
 - The action to take is to go from one place to other: {A → B, C → D, etc}
 - The reward function is the value represented by edge, i.e. cost
 - The policy is the “way” to complete the task: {A → C → F}

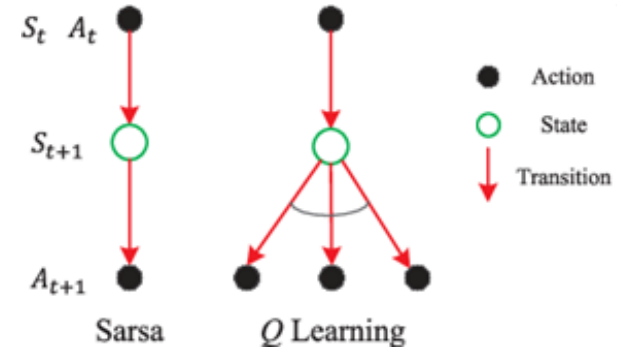


Example: Shortest path problem

- At place A, the only visible path is your next destination (a.k.a observable space)
- Algorithm can take a greedy approach and take the best possible next step, which is going from {A → D} from a subset of {A → (B, C, D, E)}
- At place D, it should go to place F, since it can choose from {D → (B, C, F)}. We see that {D → F} has the lowest cost and hence we take that path
- Our policy was to take {A → D → F} and our Value is -120
- The implemented algorithm is known as **epsilon greedy**
 - **Limitation:** it does not explore other alternatives (A → C → F)
 - **Solution:** add a probability of random exploration



- State-Action-Reward-State-Action (SARSA) very much resembles Q-learning
- Key difference: SARSA learns the Q-value based on the action performed by the current policy instead of the greedy policy



- SARSA (on-policy learner):

$$Q(s_t, a_t) \leftarrow Q(s_t, a_t) + \alpha[r_{t+1} + \gamma Q(s_{t+1}, a_{t+1}) - Q(s_t, a_t)]$$

a_{t+1} s_{t+1}

- Q-learning (off-policy/greedy learner):

$$Q(s_t, a_t) \leftarrow Q(s_t, a_t) + \alpha[r_{t+1} + \gamma \max_a Q(s_{t+1}, a) - Q(s_t, a_t)]$$

a

 SARSA(λ): Learn function $Q : \mathcal{X} \times \mathcal{A} \rightarrow \mathbb{R}$

Require:

 States $\mathcal{X} = \{1, \dots, n_x\}$

 Actions $\mathcal{A} = \{1, \dots, n_a\}, \quad A : \mathcal{X} \Rightarrow \mathcal{A}$

 Reward function $R : \mathcal{X} \times \mathcal{A} \rightarrow \mathbb{R}$

 Black-box (probabilistic) transition function $T : \mathcal{X} \times \mathcal{A} \rightarrow \mathcal{X}$

 Learning rate $\alpha \in [0, 1]$, typically $\alpha = 0.1$

 Discounting factor $\gamma \in [0, 1]$
 $\lambda \in [0, 1]$: Trade-off between TD and MC

procedure QLEARNING($\mathcal{X}, A, R, T, \alpha, \gamma, \lambda$)

 Initialize $Q : \mathcal{X} \times \mathcal{A} \rightarrow \mathbb{R}$ arbitrarily

 Initialize $e : \mathcal{X} \times \mathcal{A} \rightarrow \mathbb{R}$ with 0.

 \triangleright eligibility trace

while Q is not converged **do**

 Select $(s, a) \in \mathcal{X} \times \mathcal{A}$ arbitrarily

while s is not terminal **do**
 $r \leftarrow R(s, a)$
 $s' \leftarrow T(s, a)$
 \triangleright Receive the new state

 Calculate π based on Q (e.g. epsilon-greedy)

 $a' \leftarrow \pi(s')$
 $e(s, a) \leftarrow e(s, a) + 1$
 $\delta \leftarrow r + \gamma \cdot Q(s', a') - Q(s, a)$
for $(\tilde{s}, \tilde{a}) \in \mathcal{X} \times \mathcal{A}$ **do**
 $Q(\tilde{s}, \tilde{a}) \leftarrow Q(\tilde{s}, \tilde{a}) + \alpha \cdot \delta \cdot e(\tilde{s}, \tilde{a})$
 $e(\tilde{s}, \tilde{a}) \leftarrow \gamma \cdot \lambda \cdot e(\tilde{s}, \tilde{a})$
 $s \leftarrow s'$
 $a \leftarrow a'$
return Q

 Q-learning: Learn function $Q : \mathcal{X} \times \mathcal{A} \rightarrow \mathbb{R}$

Require:

 States $\mathcal{X} = \{1, \dots, n_x\}$

 Actions $\mathcal{A} = \{1, \dots, n_a\}, \quad A : \mathcal{X} \Rightarrow \mathcal{A}$

 Reward function $R : \mathcal{X} \times \mathcal{A} \rightarrow \mathbb{R}$

 Black-box (probabilistic) transition function $T : \mathcal{X} \times \mathcal{A} \rightarrow \mathcal{X}$

 Learning rate $\alpha \in [0, 1]$, typically $\alpha = 0.1$

 Discounting factor $\gamma \in [0, 1]$
procedure QLEARNING($\mathcal{X}, A, R, T, \alpha, \gamma$)

 Initialize $Q : \mathcal{X} \times \mathcal{A} \rightarrow \mathbb{R}$ arbitrarily

while Q is not converged **do**

 Start in state $s \in \mathcal{X}$
while s is not terminal **do**

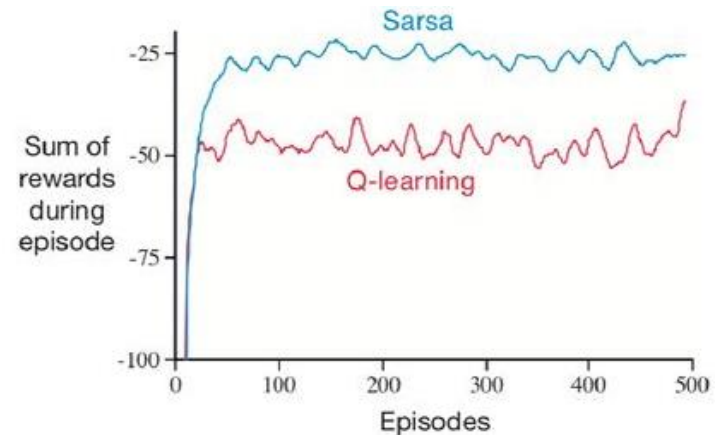
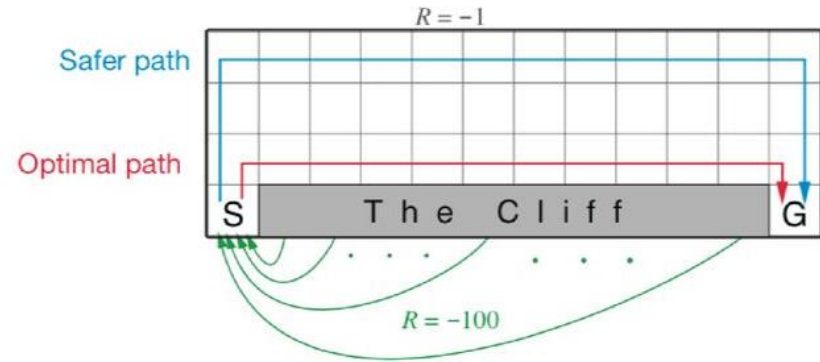
 Calculate π according to Q and exploration strategy (e.g. $\pi(x) \leftarrow \arg \max_a Q(x, a)$)

 $a \leftarrow \pi(s)$
 $r \leftarrow R(s, a)$
 \triangleright Receive the reward

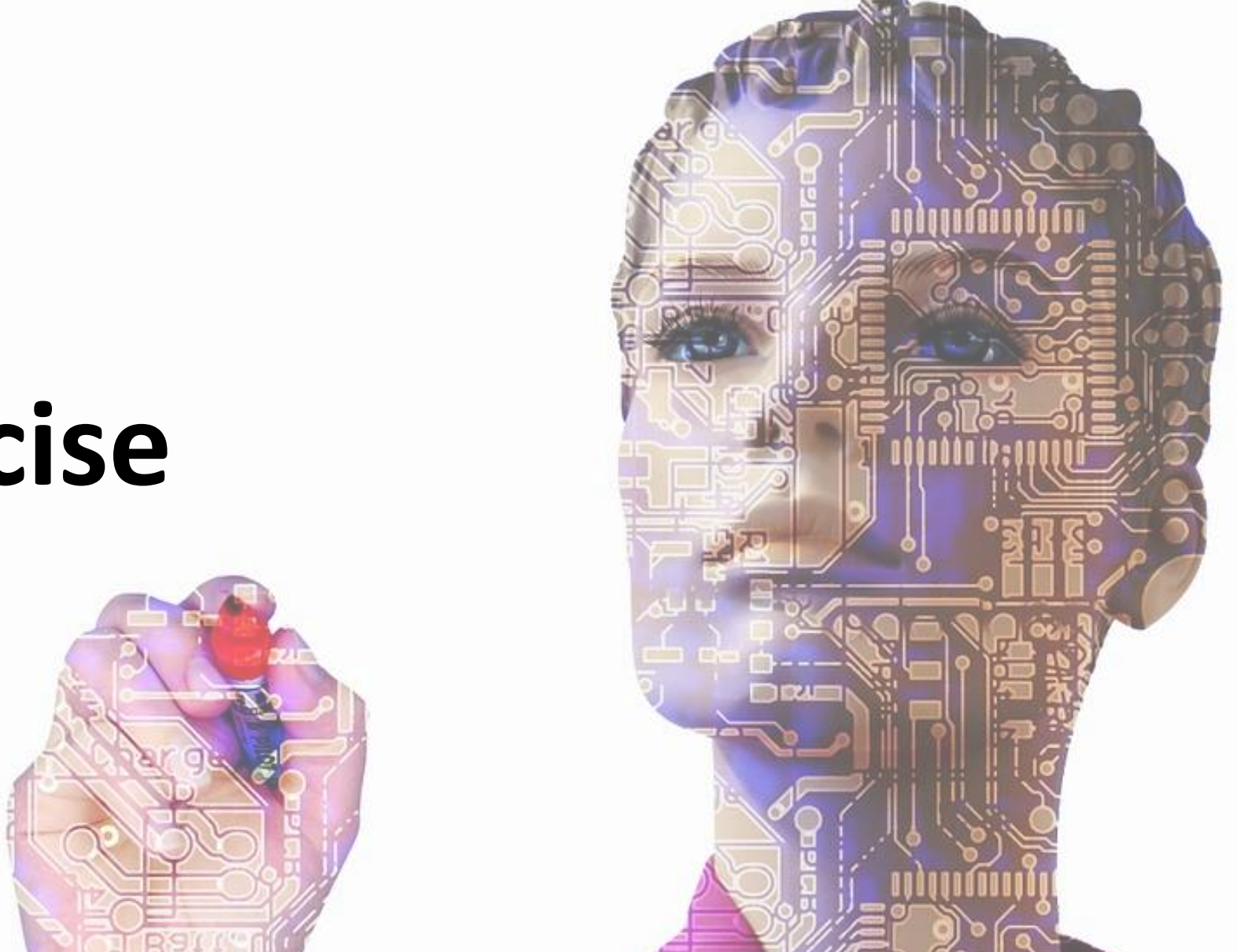
 $s' \leftarrow T(s, a)$
 \triangleright Receive the new state

 $Q(s', a) \leftarrow (1 - \alpha) \cdot Q(s, a) + \alpha \cdot (r + \gamma \cdot \max_{a'} Q(s', a'))$
 $s \leftarrow s'$
return Q

- The Cliff Walking Exercise
 - ϵ -greedy action selection
 - $\epsilon=0.1$ (fixed)
- Sarsa (on-policy learner)
 - Learns longer but safer
- Q-learning (off-policy learner)
 - Learns the optimal path
- If ϵ were reduced:
 - Both converge to the optimal policy



Exercise



▪ Game to be solved:

- Possible actions: Up, Down, Left, Right
- (1,1) -> Wall, can't go here
- (2,0) -> start position
- (0,3) -> terminal (+1 reward)
- (1,3) -> terminal (-1 reward)

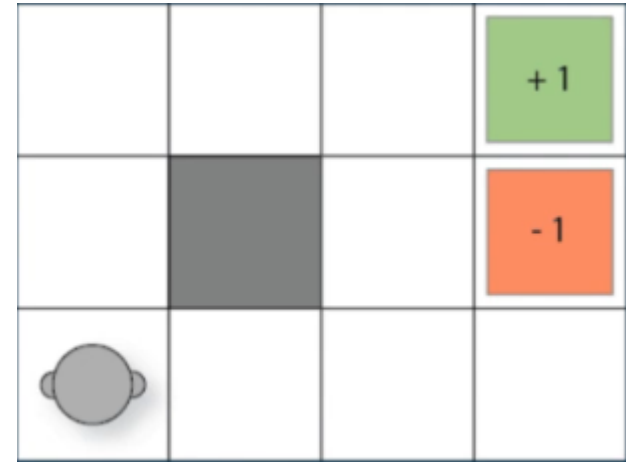
▪ Formulation:

- 12 positions
- 11 states (where the robot is)
- 4 actions
- Small game!
- But many concepts to be learned

▪ Compare the performance of the SARSA and Q-Learning algorithms

▪ Development Guideline:

- <https://towardsdatascience.com/reinforcement-learning-implement-grid-world-from-scratch-c5963765ebff>



Gridworld

Implementing Reinforcement Learning:

- Python Markov Decision Process Toolbox:
 - <http://pymdptoolbox.readthedocs.org/en/latest/api/mdp.html>
- Cat & Mouse Example:
 - <https://github.com/studywolf/blog/tree/master/RL/Cat%20vs%20Mouse%20exploration>
- Pac-Man Example:
 - <https://inst.eecs.berkeley.edu/~cs188/sp12/projects/reinforcement/reinforcement.html>



Universidade do Minho

Escola de Engenharia

Departamento de Informática

Mestrado Integrado em Engenharia Informática
Mestrado em Engenharia Informática
Computação Natural
2019/2020

Paulo Novais, Cesar Analide, Filipe Gonçalves