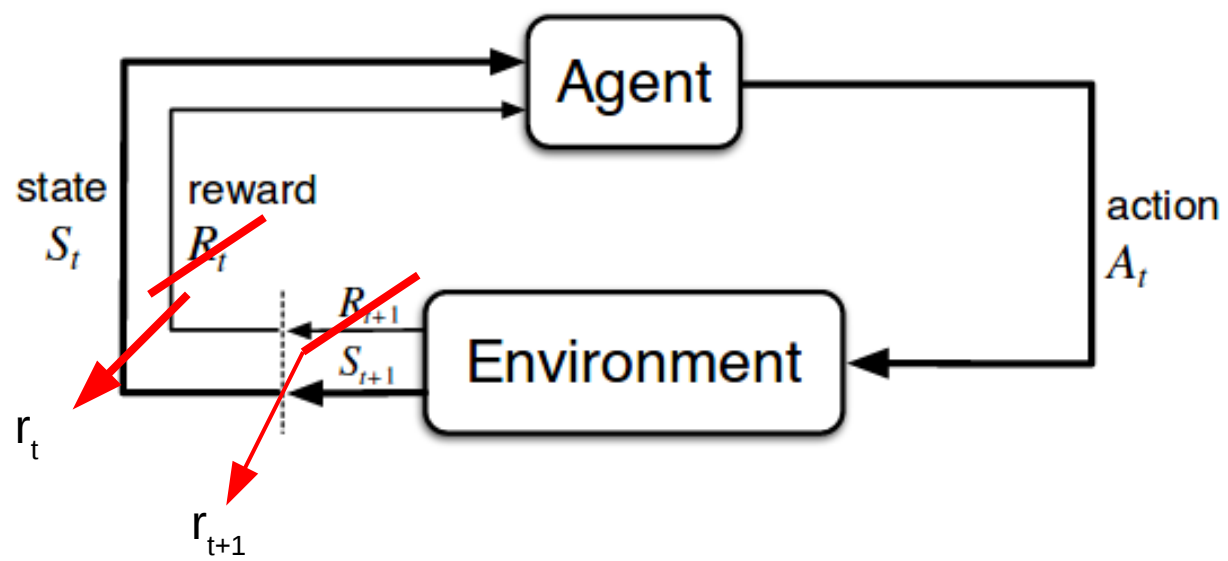# Base class

Classe mère : Agent : contient un code de base + commentaires + structure
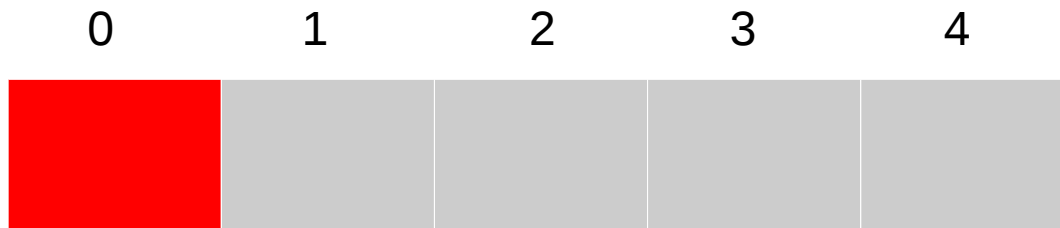
Classes filles : AgentRandom + Agents que vous allez implémenter

state
$S_t$

reward
$R_t$

$R_{t+1}$

$S_{t+1}$

action
$A_t$

Agent

Environment

$r_t$

$r_{t+1}$

# Q learning Monte Carlo

Step by step

# Initialization

| 0 | 1 | 2 | 3 | 4 |
|---|---|---|---|---|

**Reward function :**
-1 every step
10  when in 0

**Ɛ-greedy policy :**

$\pi(s) = random(a)$ with prob Ɛ
$\pi(s) = argmax_a Q(s,a)$ with prob 1- Ɛ

| s \ a | left | right |
|-------|------|-------|
| 0     |      |       |
| 1     | []   | []    |
| 2     | []   | []    |
| 3     | []   | []    |
| 4     | []   | []    |

Return table

| s \ a | left | right |
|-------|------|-------|
| 0     |      |       |
| 1     | 0    | 0     |
| 2     | 0    | 0     |
| 3     | 0    | 0     |
| 4     | 0    | 0     |

Q table

Q(a,s) = mean R(a,s)

# 1$^{st}$ iteration : performs episode

| 0 | 1 | 2 | 3 | 4 |
|---|---|---|---|---|

**Reward function :**
-1 every step
10  when in 0

**Ɛ-greedy policy :**

$\pi(s)$ = random(a) with prob Ɛ
$\pi(s)$ = argmax$_a$ Q(s,a) with prob 1- Ɛ

2 → 3 : r = -1
3 → 2 : r = -1
2 → 1 : r = -1
1 → 2 : r = -1
2 → 1 : r = -1
1 → 0 : r = 10

R = 5

Note :  Ɣ = 1

Q(a,s) = mean R(a,s)

| s \ a | left | right |
|-------|------|-------|
| 0 |  |  |
| 1 | [] | [] |
| 2 | [] | [] |
| 3 | [] | [] |
| 4 | [] | [] |

Return table

| s \ a | left | right |
|-------|------|-------|
| 0 |  |  |
| 1 | 0 | 0 |
| 2 | 0 | 0 |
| 3 | 0 | 0 |
| 4 | 0 | 0 |

Q table

# 1<sup>st</sup> iteration : update tables

| | 0 | 1 | 2 | 3 | 4 |



**Reward function :**
-1 every step
10  when in 0

**Ɛ-greedy policy :**

$\pi(s) = random(a)$ with prob $Ɛ$
$\pi(s) = argmax_a Q(s,a)$ with prob $1- Ɛ$

2 → 3 : r = -1
3 → 2 : r = -1
2 → 1 : r = -1
1 → 2 : r = -1
2 → 1 : r = -1
1 → 0 : r = 10

R = 5

Note :  ɣ = 1

$Q(a,s) = mean\ R(a,s)$

| s \ a | left | right |
|-------|------|-------|
| 0     |      |       |
| 1     | [5]  | [5]   |
| 2     | [5]  | [5]   |
| 3     | [5]  | []    |
| 4     | []   | []    |

Return table

| s \ a | left | right |
|-------|------|-------|
| 0     |      |       |
| 1     | 0    | 0     |
| 2     | 0    | 0     |
| 3     | 0    | 0     |
| 4     | 0    | 0     |

Q table

# 1<sup>st</sup> iteration : update tables

| 0 | 1 | 2 | 3 | 4 |

**Reward function :**
-1 every step
10 when in 0

**Ɛ-greedy policy :**

$\pi(s)$ = random(a) with prob Ɛ
$\pi(s)$ = argmax$_a$ Q(s,a) with prob 1- Ɛ

$2 \rightarrow 3 : r = -1$
$3 \rightarrow 2 : r = -1$
$2 \rightarrow 1 : r = -1$
$1 \rightarrow 2 : r = -1$
$2 \rightarrow 1 : r = -1$
$1 \rightarrow 0 : r = 10$

R = 5

Note : ɣ = 1

Q(a,s) = mean R(a,s)

| s \ a | left | right |
|-------|------|-------|
| 0 |   |   |
| 1 | [5] | [5] |
| 2 | [5] | [5] |
| 3 | [5] | [] |
| 4 | [] | [] |

Return table

| s \ a | left | right |
|-------|------|-------|
| 0 |   |   |
| 1 | 5 | 5 |
| 2 | 5 | 5 |
| 3 | 5 | 0 |
| 4 | 0 | 0 |

Q table

# 2$^{nd}$ iteration : perform episode

| 0 | 1 | 2 | 3 | 4 |
|---|---|---|---|---|

**Reward function :**
-1 every step
10  when in 0

**Ɛ-greedy policy :**

$\pi(s)$ = random(a) with prob Ɛ
$\pi(s)$ = argmax$_a$ Q(s,a) with prob 1- Ɛ

3 → 2 : r = -1
2 → 1 : r = -1
1 → 2 : r = -1
2 → 1 : r = -1
1 → 0 : r = 10

R = 6

| s \ a | left | right |
|-------|------|-------|
| 0 | | |
| 1 | [5] | [5] |
| 2 | [5] | [5] |
| 3 | [5] | [] |
| 4 | [] | [] |

Return table

| s \ a | left | right |
|-------|------|-------|
| 0 | | |
| 1 | 5 | 5 |
| 2 | 5 | 5 |
| 3 | 5 | 0 |
| 4 | 0 | 0 |

Q table

Q(a,s) = mean R(a,s)

# 2<sup>nd</sup> iteration : update tables

2$^{nd}$ iteration : update tables

| 0 | 1 | 2 | 3 | 4 |
|---|---|---|---|---|

**Reward function :**
-1 every step
10  when in 0

**Ɛ-greedy policy :**

$\pi(s)$ = random(a) with prob Ɛ
$\pi(s)$ = argmax$_a$ Q(s,a) with prob 1- Ɛ

3 → 2 : r = -1
2 → 1 : r = -1
1 → 2 : r = -1
2 → 1 : r = -1
1 → 0 : r = 10

R = 6

| s \ a | left | right |
|-------|------|-------|
| 0 |       |       |
| 1 | [5,6] | [5,6] |
| 2 | [5,6] | [5]   |
| 3 | [5,6] | []    |
| 4 | []    | []    |

Return table

| s \ a | left | right |
|-------|------|-------|
| 0 |   |   |
| 1 | 5 | 5 |
| 2 | 5 | 5 |
| 3 | 5 | 0 |
| 4 | 0 | 0 |

Q table

Q(a,s) = mean R(a,s)

# 2<sup>nd</sup> iteration : update episode

| 0 | 1 | 2 | 3 | 4 |

**Reward function :**
-1 every step
10  when in 0

**Ɛ-greedy policy :**

$\pi(s)$ = random(a) with prob Ɛ
$\pi(s)$ = argmax$_a$ Q(s,a) with prob 1- Ɛ

3 → 2 : r = -1
2 → 1 : r = -1
1 → 2 : r = -1
2 → 1 : r = -1
1 → 0 : r = 10

R = 6

| s \ a | left | right |
|-------|------|-------|
| 0     |      |       |
| 1     | [5,6] | [5,6] |
| 2     | [5,6] | [5]   |
| 3     | [5,6] | []    |
| 4     | []   | []    |

Return table

| s \ a | left | right |
|-------|------|-------|
| 0     |      |       |
| 1     | 5.5  | 5.5   |
| 2     | 5.5  | 5     |
| 3     | 5.5  | 0     |
| 4     | 0    | 0     |

Q table

Q(a,s) = mean R(a,s)

# 3$^{rd}$ iteration : perform episode

|   | 0 | 1 | 2 | 3 | 4 |
|---|---|---|---|---|---|

**Reward function :**
-1 every step
10  when in 0

**Ɛ-greedy policy :**

$\pi(s)$ = random(a) with prob Ɛ
$\pi(s)$ = argmax$_a$ Q(s,a) with prob 1- Ɛ

2 → 1 : r = -1
1 → 0 : r = 10

R = 9

| s \ a | left | right |
|-------|------|-------|
| 0 |  |  |
| 1 | [5,6] | [5,6] |
| 2 | [5,6] | [5] |
| 3 | [5,6] | [] |
| 4 | [] | [] |

Return table

| s \ a | left | right |
|-------|------|-------|
| 0 |  |  |
| 1 | 5.5 | 5.5 |
| 2 | 5.5 | 5 |
| 3 | 5.5 | 0 |
| 4 | 0 | 0 |

Q table

Q(a,s) = mean R(a,s)

# 3<sup>rd</sup> iteration : update tables

| 0 | 1 | 2 | 3 | 4 |

**Reward function :**
-1 every step
10  when in 0

**Ɛ-greedy policy :**

$\pi(s)$ = random(a) with prob $\mathcal{E}$
$\pi(s)$ = argmax$_a$ Q(s,a) with prob 1- $\mathcal{E}$

2 → 1 : r = -1
1 → 0 : r = 10

R = 9

| s \ a | left | right |
|-------|---------|-------|
| 0 | | |
| 1 | [5,6,9] | [5,6] |
| 2 | [5,6,9] | [5] |
| 3 | [5,6] | [] |
| 4 | [] | [] |

Return table

| s \ a | left | right |
|-------|------|-------|
| 0 | | |
| 1 | 5.5 | 5.5 |
| 2 | 5.5 | 5 |
| 3 | 5.5 | 0 |
| 4 | 0 | 0 |

Q table

Q(a,s) = mean R(a,s)

# 3<sup>rd</sup> iteration : update tables



```
      0       1       2       3       4
```

**Reward function :**
-1 every step
10  when in 0

**Ɛ-greedy policy :**

$\pi(s) = random(a)$ with prob $\epsilon$
$\pi(s) = \text{argmax}_a Q(s,a)$ with prob 1- $\epsilon$

2 → 1 : r = -1
1 → 0 : r = 10

R = 9

| s \ a | left | right |
|-------|--------|-------|
| 0 | | |
| 1 | [5,6,9] | [5,6] |
| 2 | [5,6,9] | [5] |
| 3 | [5,6] | [] |
| 4 | [] | [] |

Return table

| s \ a | left | right |
|-------|------|-------|
| 0 | | |
| 1 | 6.6 | 5.5 |
| 2 | 6.6 | 5 |
| 3 | 5.5 | 0 |
| 4 | 0 | 0 |

Q table

Q(a,s) = mean R(a,s)

# 4<sup>th</sup> iteration : perform episode



| 0 | 1 | 2 | 3 | 4 |
|---|---|---|---|---|

**Reward function :**
-1 every step
10  when in 0

**Ɛ-greedy policy :**

$\pi(s)$ = random(a) with prob Ɛ
$\pi(s)$ = argmax$_a$ Q(s,a) with prob 1- Ɛ

4 → 4 : r = -1
4 → 3 : r = -1
3 → 2 : r = -1
2 → 3 : r = -1
3 → 2 : r = -1
2 → 1 : r = -1
1 → 0 : r = 10

R = 4

| s \ a | left | right |
|-------|--------|--------|
| 0 | | |
| 1 | [5,6,9] | [5,6] |
| 2 | [5,6,9] | [5] |
| 3 | [5,6] | [] |
| 4 | [] | [] |

Return table

| s \ a | left | right |
|-------|------|-------|
| 0 | | |
| 1 | 6.6 | 5.5 |
| 2 | 6.6 | 5 |
| 3 | 5.5 | 0 |
| 4 | 0 | 0 |

Q table

Q(a,s) = mean R(a,s)

# 4<sup>th</sup> iteration : update tables

|   | 0 | 1 | 2 | 3 | 4 |
|---|---|---|---|---|---|



**Reward function :**
-1 every step
10  when in 0

**Ɛ-greedy policy :**

$\pi(s)$ = random(a) with prob Ɛ
$\pi(s)$ = argmax$_a$ Q(s,a) with prob 1- Ɛ

4 → 4 : r = -1
4 → 3 : r = -1
3 → 2 : r = -1
2 → 3 : r = -1
3 → 2 : r = -1
2 → 1 : r = -1
1 → 0 : r = 10

R = 4

| s \ a | left | right |
|-------|------|-------|
| 0 |  |  |
| 1 | [5,6,9,4] | [5,6] |
| 2 | [5,6,9,4] | [5,4] |
| 3 | [5,6,4] | [] |
| 4 | [4] | [4] |

Return table

| s \ a | left | right |
|-------|------|-------|
| 0 |  |  |
| 1 | 6.6 | 5.5 |
| 2 | 6.6 | 5 |
| 3 | 5.5 | 0 |
| 4 | 0 | 0 |

Q table

Q(a,s) = mean R(a,s)

# 4<sup>th</sup> iteration : update tables



**Reward function :**
-1 every step
10  when in 0

**Ɛ-greedy policy :**

$\pi(s)$ = random(a) with prob Ɛ
$\pi(s)$ = $\text{argmax}_a$ Q(s,a) with prob 1- Ɛ

4 → 4 : r = -1
4 → 3 : r = -1
3 → 2 : r = -1
2 → 3 : r = -1
3 → 2 : r = -1
2 → 1 : r = -1
1 → 0 : r = 10

R = 4

| s \ a | left | right |
|-------|------|-------|
| 0 | | |
| 1 | [5,6,9,4] | [5,6] |
| 2 | [5,6,9,4] | [5,4] |
| 3 | [5,6,4] | [] |
| 4 | [4] | [4] |

Return table

| s \ a | left | right |
|-------|------|-------|
| 0 | | |
| 1 | 6 | 5.5 |
| 2 | 6 | 4.5 |
| 3 | 5 | 0 |
| 4 | 4 | 4 |

Q table

Q(a,s) = mean R(a,s)