# Apprentissage par renforcement

### Sommaire



• Présentation générale



Objectif



• L'algorithme du Q-learning



• Application simple : Le labyrinthe



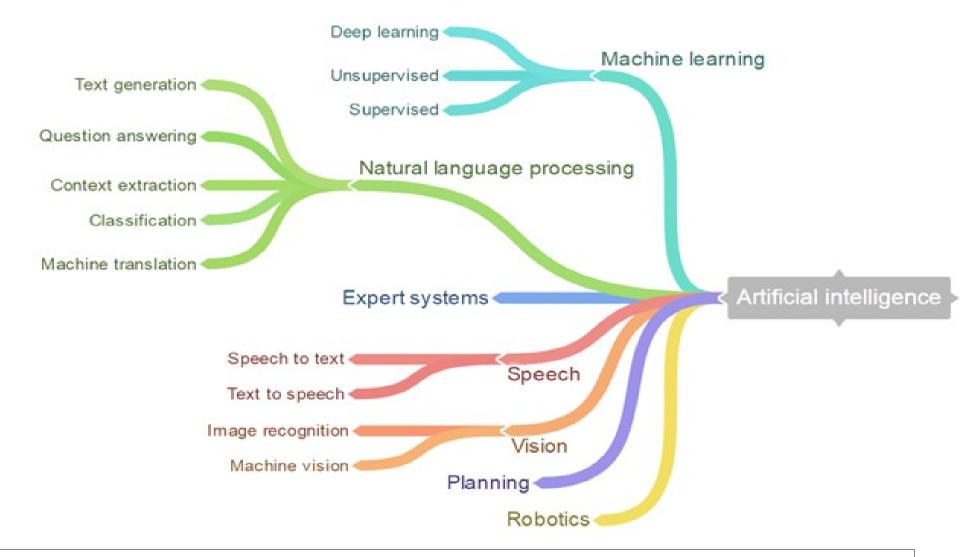
- Application : Pendule inversé
  - Modèle de l'environnement
  - Résultats
  - Limite et conclusion

六

# Présentation générale



### Vue d'ensemble



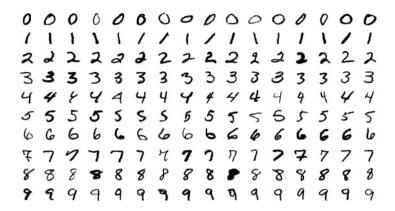
<u>Définition</u>: Intelligence artificielle

L'ensemble des théories et des techniques mises en œuvre en vue de réaliser des machines capables de simuler l'intelligence

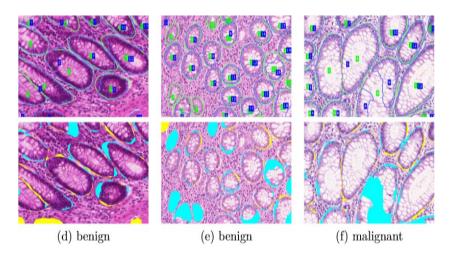


### Pourquoi?

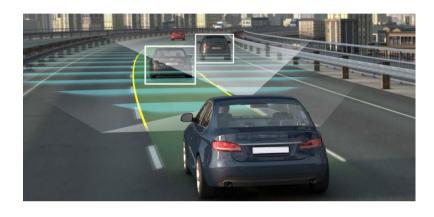
Résoudre des problèmes où la programmation classique est insuffisante



Reconnaissance de caractère



Détection de cancer



Voiture autonome

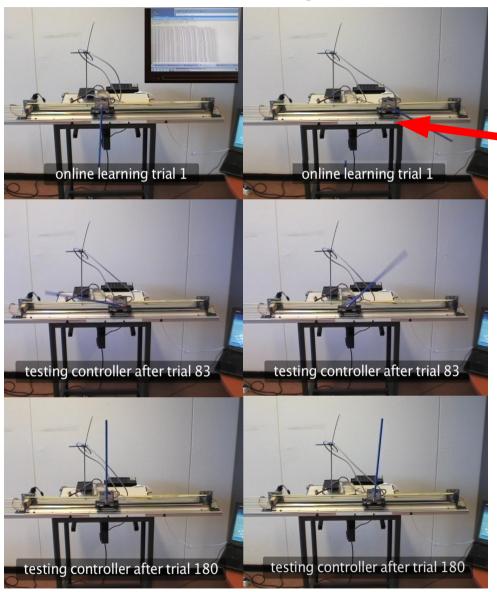


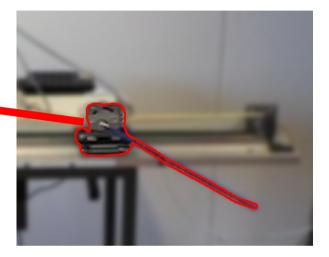
Optimisation procédé de production



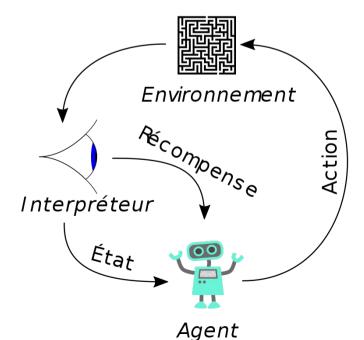
## L'apprentissage par renforcement

Plusieurs algorithmes basés sur un principe commun

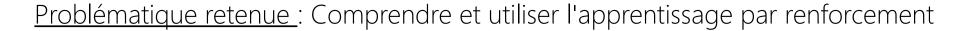




Pendule inversé











Objectif : Reproduire l'exemple du pendule inversé

### Sommaire



• Présentation générale



Objectif



• L'algorithme du Q-learning



• Application simple : Le labyrinthe

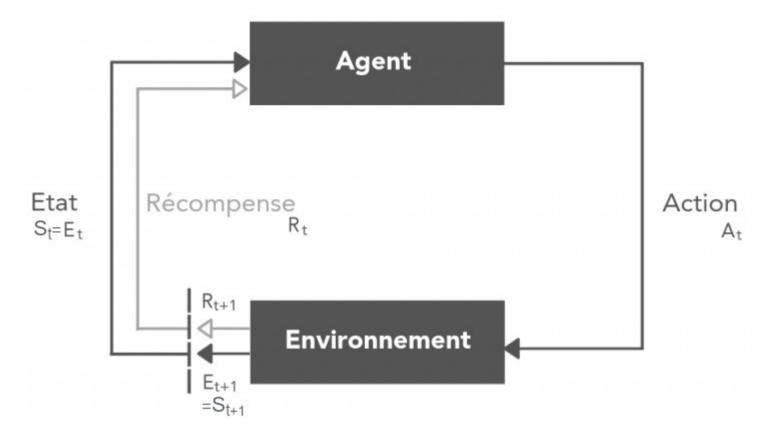


- Application : Pendule inversé
  - Modèle de l'environnement
  - Résultats
  - Limite et conclusion



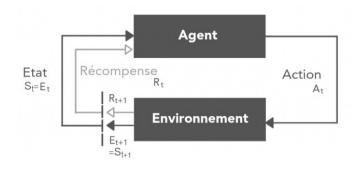
# L'algorithme du Q-learning





Q-Table		Actions					
		South (0) North (1)		East (2)	West (3)	Pickup (4)	Dropoff (5)
	0	0	0	0	0	0	0
				9			
		43					
			1.00				
States	328	-2.30108105	-1.97092096	-2.30357004	-2.20591839	-10.3607344	-8.5583017
			100				1.0
	499	9.96984239	4.02706992	12.96022777	29	3.32877873	3.38230603



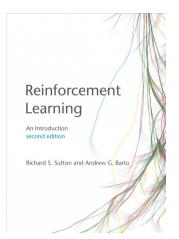


Initialized

Q-Table		Actions					
		0	0	0	0	0	0
States		0	0	0	0	0	0
			,				
		0	0	0	0	0	0



Q-Table		Actions						
		South (0)				Pickup (4)		
		0	0	0	0	0	0	
			100			20		
States		-2.30108105	-1.97092096	-2.30357004	-2.20591839	-10.3607344	-8.5583017	
		9.96984239	4.02706992	12.96022777	29	3.32877873	3.38230603	



[1] R.S. SUTTON, A. BARTO: Reinforcement Learning: An Introduction, 1998, MIT Press

$$Q(S_t, A_t) \leftarrow Q(S_t, A_t) + \alpha \left[ R_{t+1} + \gamma \max_{a} Q(S_{t+1}, a) - Q(S_t, A_t) \right]$$
  
Équation de Bellman

#### Q-learning (off-policy TD control) for estimating $\pi \approx \pi_*$

Algorithm parameters: step size  $\alpha \in (0,1]$ , small  $\varepsilon > 0$ Initialize Q(s,a), for all  $s \in \mathbb{S}^+$ ,  $a \in \mathcal{A}(s)$ , arbitrarily except that  $Q(terminal, \cdot) = 0$ 

Loop for each episode:

Initialize S

Loop for each step of episode:

Choose A from S using policy derived from Q (e.g.,  $\varepsilon$ -greedy)

Take action A, observe R, S'

$$Q(S, A) \leftarrow Q(S, A) + \alpha [R + \gamma \max_{a} Q(S', a) - Q(S, A)]$$

 $S \leftarrow S'$ 

until S is terminal

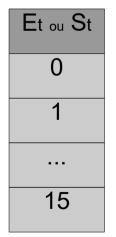


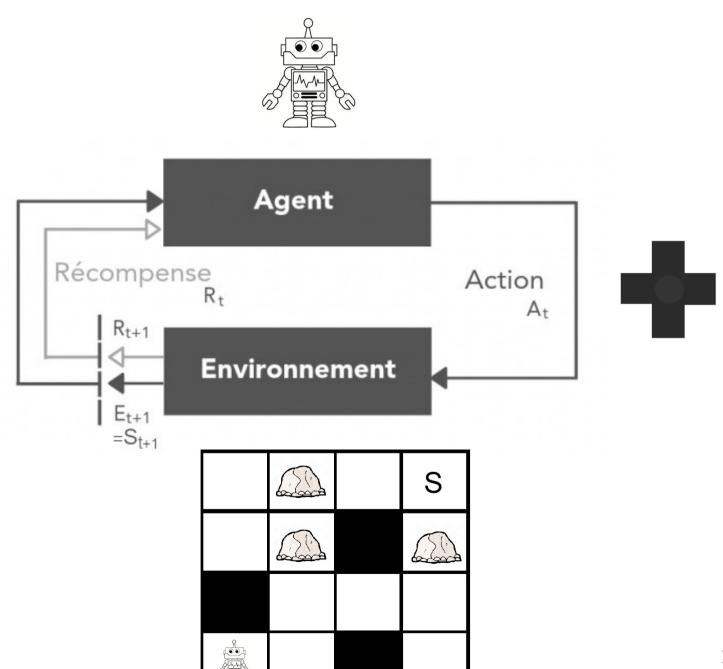
Application simple: Le labyrinthe



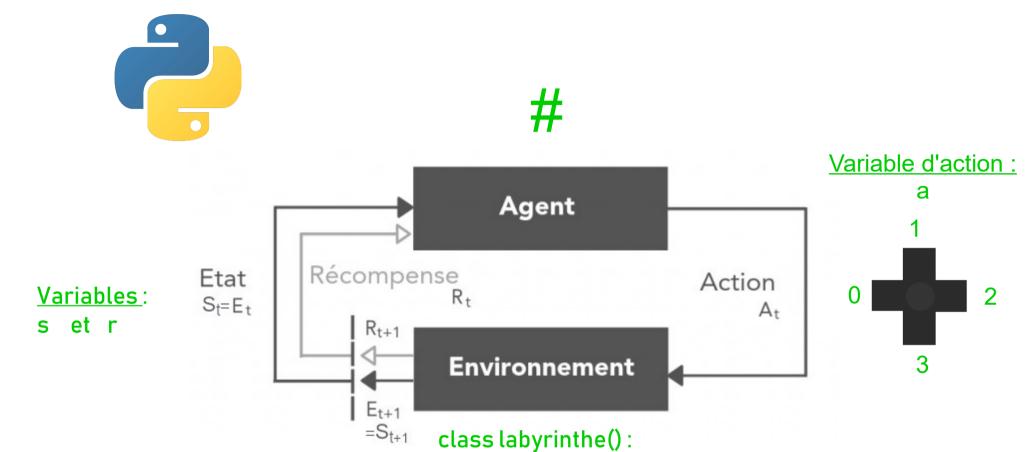
Élément	Rt	
Vide	-1	
Rocher	-10	
Sortie	+10	
Butée	-5	

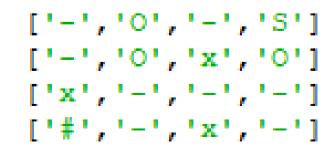


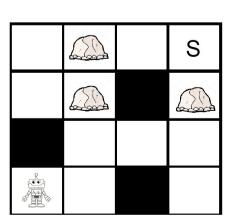














# Application simple: Le labyrinthe

L'environnement

L'agent/ Le Q-learning

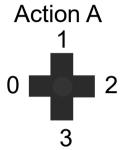


### Action A

# Environnement class labyrinthe()

État(=State) S

Récompense R



### État(=State) S

0	4	8	12
1	5	9	13
2	6	10	14
3	7	11	15

### Récompense R

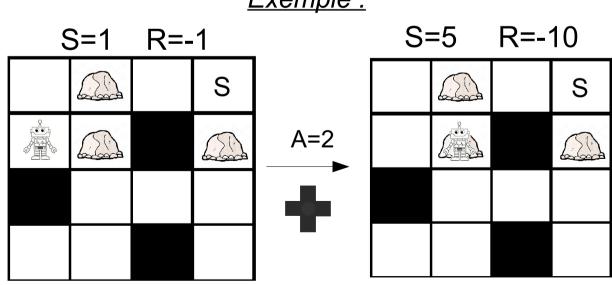
Élément	r
Vide	-1
Rocher	-10
Sortie	+10
Butée	-5

# self.position=[1,1] self.fin=False

self.obtenir\_id(x,y)
 #retoune 5 pour [1,1]
self.deplacer(3)



### Exemple:





# Application simple: Le labyrinthe

L'environnement

L'agent/Le Q-learning



### Agent

#### Q-learning (off-policy TD control) for estimating $\pi \approx \pi_*$

Algorithm parameters: step size  $\alpha \in (0,1]$ , small  $\varepsilon > 0$ 

Initialize Q(s, a), for all  $s \in S^+$ ,  $a \in A(s)$ , arbitrarily except that  $Q(terminal, \cdot) = 0$ 

Loop for each episode:

Initialize S

Loop for each step of episode:

Choose A from S using policy derived from Q (e.g.,  $\varepsilon$ -greedy)

Take action A, observe R, S'

$$Q(S, A) \leftarrow Q(S, A) + \alpha \left[ R + \gamma \max_{a} Q(S', a) - Q(S, A) \right]$$

$$S \leftarrow S'$$

until S is terminal

		Actions				
		0	1	2	3	
	0	0,846	0,657	0,214	0,981	
	1	0,624	0,321	0,247	0,458	
	2	0,148	0,854	0,312	0,174	
Etats	3	0,394	0,43	0,378	0,368	
	15	0,351	0,157	0,954	0,545	



#### Q-learning (off-policy TD control) for estimating $\pi \approx \pi_*$

Algorithm parameters: step size  $\alpha \in (0,1]$ , small  $\varepsilon > 0$ Initialize Q(s,a), for all  $s \in \mathbb{S}^+$ ,  $a \in \mathcal{A}(s)$ , arbitrarily except that  $Q(terminal, \cdot) = 0$ 

Loop for each episode:

Initialize S

Loop for each step of episode:

Choose A from S using policy derived from Q (e.g.,  $\varepsilon$ -greedy)

Take action A, observe R, S'

$$Q(S, A) \leftarrow Q(S, A) + \alpha \left[ R + \gamma \max_{a} Q(S', a) - Q(S, A) \right]$$

$$S \leftarrow S'$$

until S is terminal

		Actions				
		0	1	2	3	
	0	0,846	0,657	0,214	0,981	
	1	0,624	0,321	0,247	0,458	
	2	0,148	0,854	0,312	0,174	
Etats	3	0,394	0,43	0,378	0,368	
	 15	0,351	0,157	0,954	0,545	

0	4	8	12
1	5	9	13
2	6	10	14
3	7	11	15

S=3



#### Q-learning (off-policy TD control) for estimating $\pi \approx \pi_*$

Algorithm parameters: step size  $\alpha \in (0,1]$ , small  $\varepsilon > 0$ 

Initialize Q(s,a), for all  $s \in S^+$ ,  $a \in A(s)$ , arbitrarily except that  $Q(terminal, \cdot) = 0$ 

Loop for each episode:

Initialize S

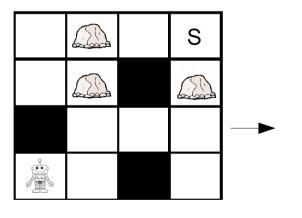
Loop for each step of episode:

Choose A from S using policy derived from Q (e.g.,  $\varepsilon$ -greedy) Take action A, observe R, S'

$$Q(S, A) \leftarrow Q(S, A) + \alpha \left[ R + \gamma \max_{a} Q(S', a) - Q(S, A) \right]$$

 $S \leftarrow S'$ 

until S is terminal



		Actions				
		0	, 1	2	3	
	0	0,846	0,657	0,214	0,981	
	1	0,624	0,321	0,247	0,458	
	2	0,148	0,854	0,312	0,174	
Etats	3	0,394	0,43	0,378	0,368	
	15	0,351	0,157	0,954	0,545	

S=3



#### Q-learning (off-policy TD control) for estimating $\pi \approx \pi_*$

Algorithm parameters: step size  $\alpha \in (0, 1]$ , small  $\varepsilon > 0$ Initialize Q(s, a), for all  $s \in S^+$ ,  $a \in \mathcal{A}(s)$ , arbitrarily except that  $Q(terminal, \cdot) = 0$ 

Loop for each episode:

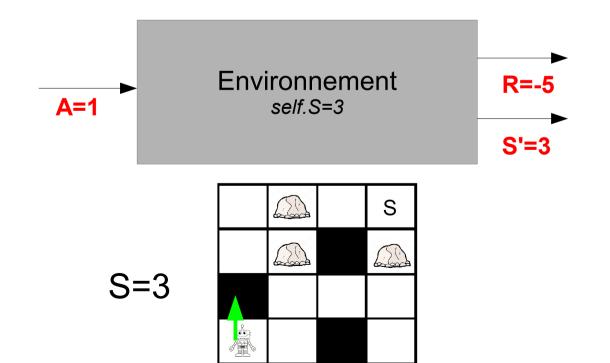
Initialize S

Loop for each step of episode:

Choose A from S using policy derived from Q (e.g.,  $\varepsilon$ -greedy)

Take action A, observe R, S'  $Q(S,A) \leftarrow Q(S,A) + \alpha \left[ R + \gamma \max_{a} Q(S',a) - Q(S,A) \right]$   $S \leftarrow S'$ 

until S is terminal





#### Q-learning (off-policy TD control) for estimating $\pi \approx \pi_*$

Algorithm parameters: step size  $\alpha \in (0,1]$ , small  $\varepsilon > 0$ Initialize Q(s,a), for all  $s \in S^+$ ,  $a \in \mathcal{A}(s)$ , arbitrarily except that  $Q(terminal, \cdot) = 0$ 

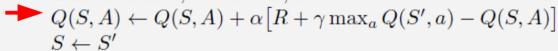
Loop for each episode:

Initialize S

Loop for each step of episode:

Choose A from S using policy derived from Q (e.g.,  $\varepsilon$ -greedy)

Take action A, observe R, S'



until S is terminal

#### Variables de l'environnement

$$A_{t=1}$$
  $R_{t+1}=-5$   $S_{t=3}$   $S_{t+1}=3$ 

		Actions				
	0	1	2	3		
0	0,846	0,657	0,214	0,981		
1	0,624	0,321	0,247	0,458		
2	0,148	0,854	0,312	0,174		
3	0,378	0.43	0,394	0,368		
15	0,351	0,157	0,954	0,545		
֡	3	1 0,624 2 0,148 3 0,378	0 1 0 0,846 0,657 1 0,624 0,321 2 0,148 0,854 3 0,378 0.43	1 0,624 0,321 0,247 2 0,148 0,854 0,312 3 0,378 0.43 0,394		

#### Paramètre de l'algorithme

$$\alpha$$
 (Learning\_rate) = 0.7  $\gamma$ = 0.99

$$Q(S_t, A_t) \leftarrow Q(S_t, A_t) + \alpha \left[ R_{t+1} + \gamma \max_{a} Q(S_{t+1}, a) - Q(S_t, A_t) \right]$$
  
 $Q(S=3,A=1) = 0.43 +0.7^*[-5 +0.99 * 0.43 - 0.43] =-9.85$ 

		Actions				
		0	1	2	3	
	0	0,846	0,657	0,214	0,981	
	1	0,624	0,321	0,247	0,458	
C4-4-	2	0,148	0,854	0,312	0,174	
Etats	3	0,378	-9,85	0,394	0,368	
	15	0,351	0,157	0,954	0,545	



#### Q-learning (off-policy TD control) for estimating $\pi \approx \pi_*$

Algorithm parameters: step size  $\alpha \in (0,1]$ , small  $\varepsilon > 0$ 

Initialize Q(s, a), for all  $s \in S^+$ ,  $a \in A(s)$ , arbitrarily except that  $Q(terminal, \cdot) = 0$ 

Loop for each episode:

Initialize S

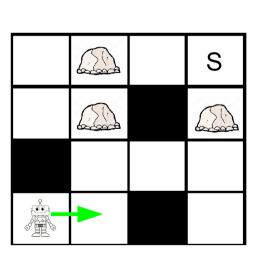
Loop for each step of episode:

 $\longrightarrow$  Choose A from S using policy derived from Q (e.g.,  $\varepsilon$ -greedy)

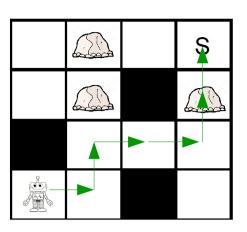
Take action A, observe R, S'  $Q(S,A) \leftarrow Q(S,A) + \alpha \left[ R + \gamma \max_{a} Q(S',a) - Q(S,A) \right]$   $S \leftarrow S'$ 

until S is terminal

		Actio <u>ns</u>				
		0	1	. 2	3	
Etats	0	0,846	0,657	0,214	0,981	_
	1	0,624	0,321	0,247	0,458	
	2	0,148	0,854	0,312	0,174	
	3	0,378	-9,85	0,394	0,368	5-3
	15	0,351	0,157	0,954	0,545	

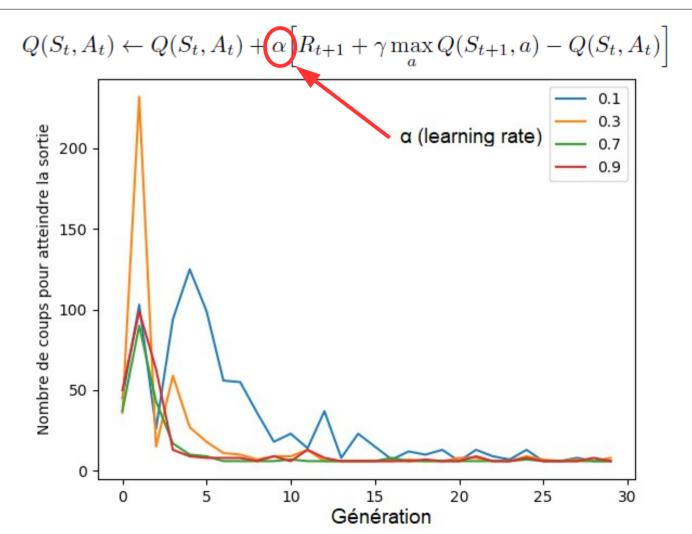






### **Conclusion**:

L'agent atteint la sortie en un minimum de coup (6 coups)



### Sommaire



• Présentation générale



Objectif



• L'algorithme du Q-learning



• Application simple : Le labyrinthe



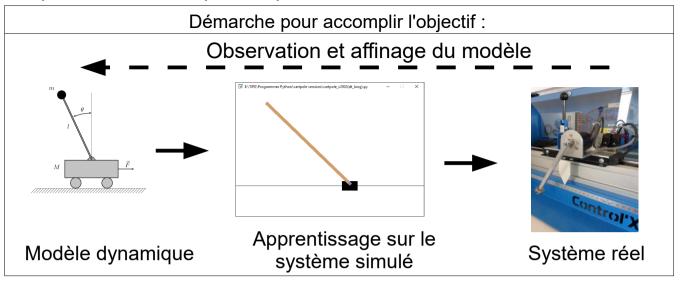
- Application : Pendule inversé
  - Modèle de l'environnement
  - Résultats
  - Limite et conclusion

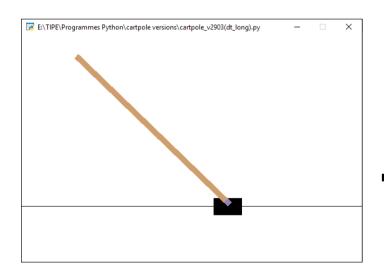


Application : Pendule inversé

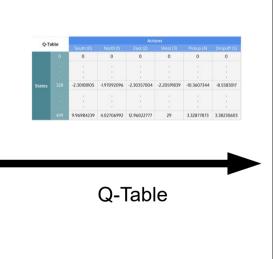
### Objectif: Reproduire l'exemple du pendule inversé

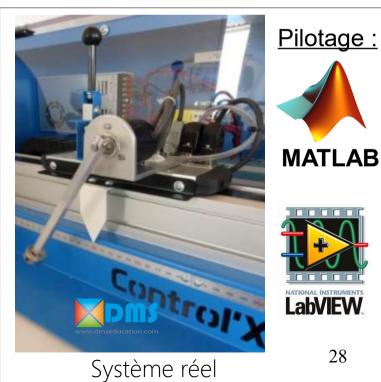


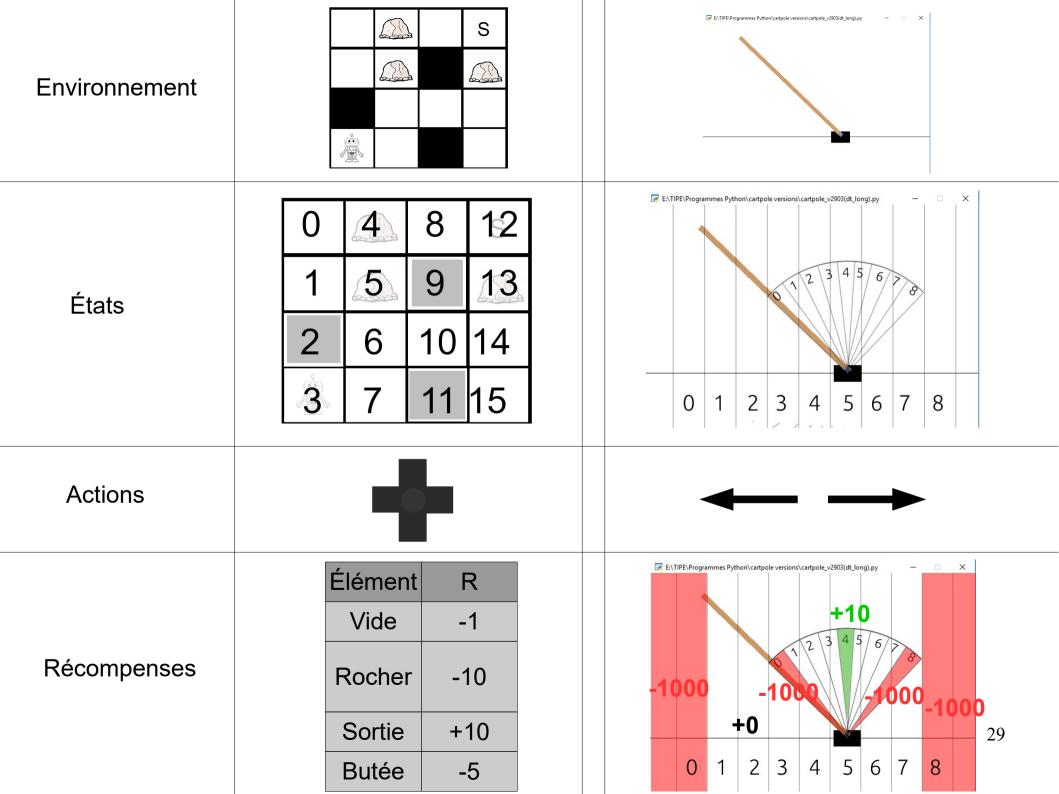




Système simulé









# Application : Pendule inversé

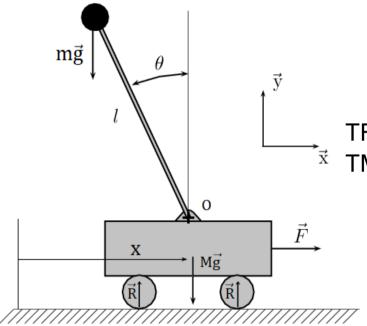
Modèle de l'environnement

Résultats

Améliorer les performances

# Enviro

### Environnement



<u>Hypothèse :</u> Liaisons parfaites

### Modèle mécanique

On isole {chariot+barre}

#### Bilan des Actions Mécaniques Extérieures :

Force de traction F Poids de la barre Poids du chariot Réaction du support

On applique le Principe Fondamental de la Dynamique en O

$$\overrightarrow{x} \ \ \text{TRD sur x} \begin{cases} (M+m)\ddot{x} = F + ml\ddot{\theta}cos(\theta) - ml\dot{\theta}^2sin(\theta) \\ ml^2\ddot{\theta} = mgsin(\theta)l + mlcos(\theta)\ddot{x} \end{cases}$$

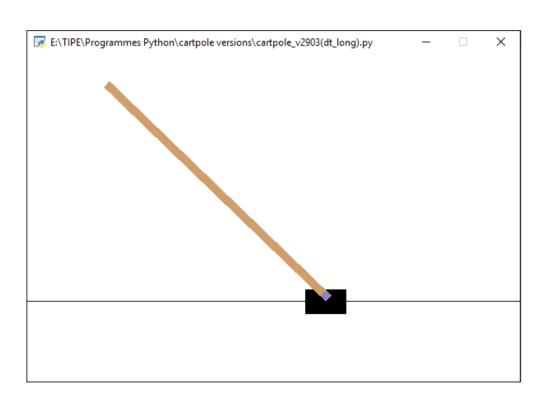
On découple les variables

$$temp = \frac{F - ml\dot{\theta}^2 sin(\theta)}{M + m}$$

$$\ddot{\theta} = \frac{\cos(\theta).temp + g\sin(\theta)}{l(1 - \frac{m\cos(\theta)^2)}{M+m})}$$
 
$$\ddot{x} = temp + \frac{ml\ddot{\theta}sin(\theta)}{M+m}$$



### Système simulé

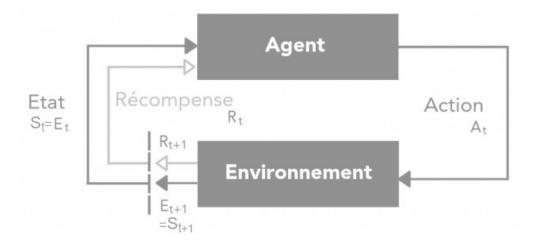


Fenêtre graphique de la bibliothèque pré-programmé

### Résolution numérique : Méthode d'Euler

```
x_dot = x_dot + self.tau * xacc
x = x + self.tau * x_dot
theta_dot = theta_dot + self.tau * thetaacc
theta = theta + self.tau * theta_dot
```

### Agent



### Identique à l'application simple du labyrinthe

```
Q-learning (off-policy TD control) for estimating \pi \approx \pi_*
Algorithm parameters: step size \alpha \in (0,1], small \varepsilon > 0
Initialize Q(s,a), for all s \in \mathbb{S}^+, a \in \mathcal{A}(s), arbitrarily except that Q(terminal, \cdot) = 0
Loop for each episode:
Initialize S
Loop for each step of episode:
Choose A from S using policy derived from Q (e.g., \varepsilon-greedy)
Take action A, observe R, S'
Q(S,A) \leftarrow Q(S,A) + \alpha \left[R + \gamma \max_a Q(S',a) - Q(S,A)\right]
S \leftarrow S'
until S is terminal
```

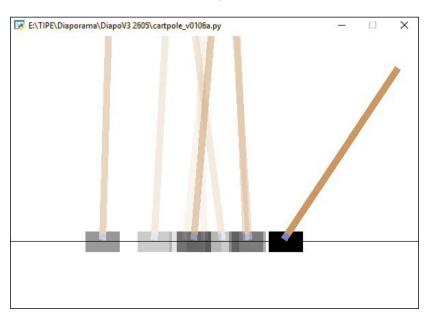


# Application : Pendule inversé

Modèle de l'environnement

Résultats

Améliorer les performances

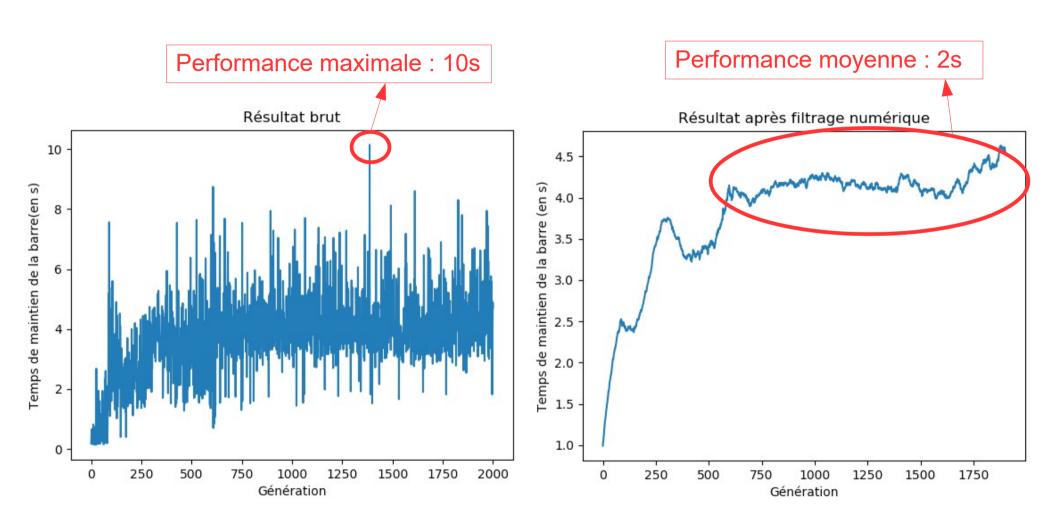


#### Paramètre de l'essai:

Discrétisation en 9 états

0 1 2 3 4 5 6 7 8

#### Apprentissage sur 2000 générations





# Application : Pendule inversé

Modèle de l'environnement Résultats

Améliorer les performances

#### Dilemme Exploration/Exploitation Fonction d'exploration

<u>Problème</u>: L'agent tend à reproduire un comportement passé : il ne peut pas découvrir de nouvelles manières de résoudre le problème

```
## START Choose A from S using
a=np.argmax(Q[s[0],s[1],s[2],s[3],:])
## END Choose A from S using

## START Choose A from S using

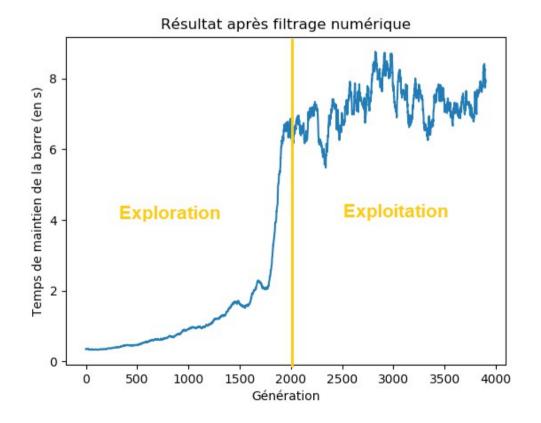
## Choisir action optimisée

## Choisir action optimisée

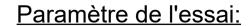
## END Choose A from S using

## END Choose A from S using

## END Choose A from S using
```



#### Problème : L'agent a un comportement trop grossier dû au découpage trop faible

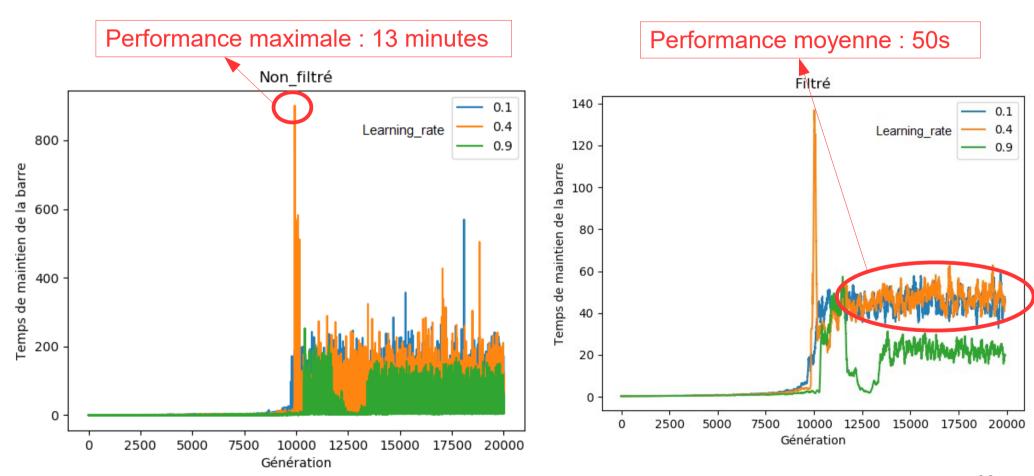


Discrétisation en 41 états

0 1 2 3 4 5 6 7 8

Apprentissage sur 20000 générations

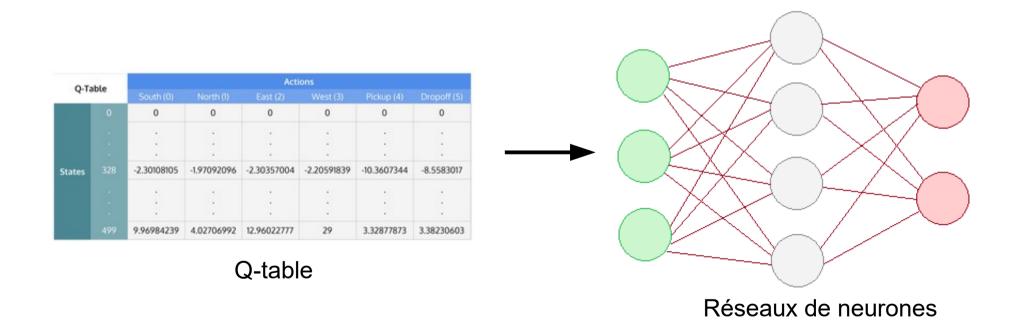
Observation secondaire: Influence du learning\_rate



## Limites

Augmenter le nombre d'états possibles => Amélioration des performances

Problème : Temps de calcul très important et limites des capacités machines atteintes



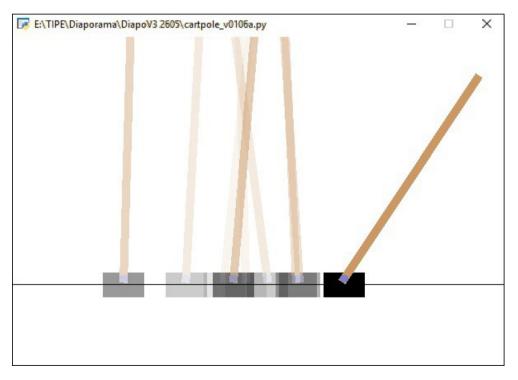
<u>Solution/ Ouverture</u>: Se libérer des contraintes de discrétisations pour traiter des états continus par les réseaux de neurones

# Conclusion



<u>Problématique retenue</u>: Comprendre et utiliser l'apprentissage par renforcement

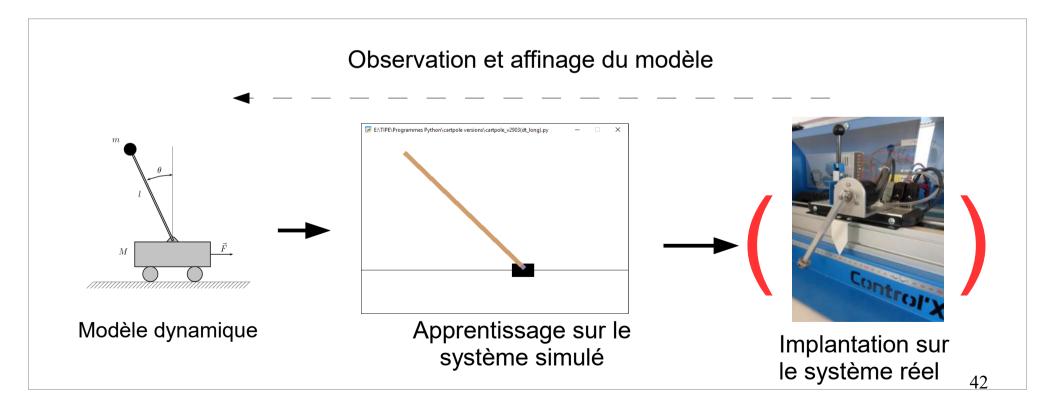




# Conclusion



<u>Problématique retenue</u>: Comprendre et utiliser l'apprentissage par renforcement



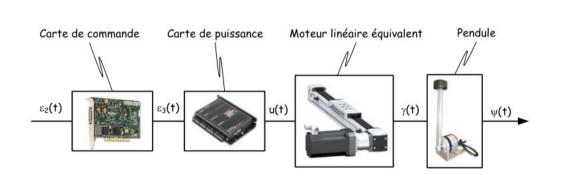
# Annexe

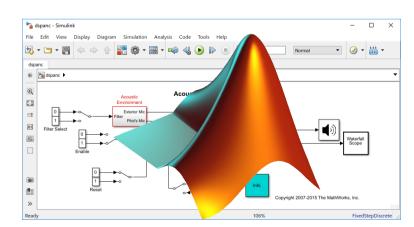
### Système réel







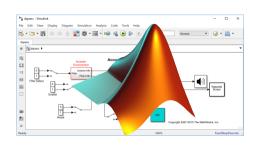




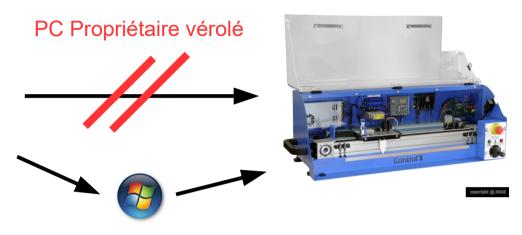


### Système réel

#### Problème 1 : Problème de communication entre le logiciel et le système

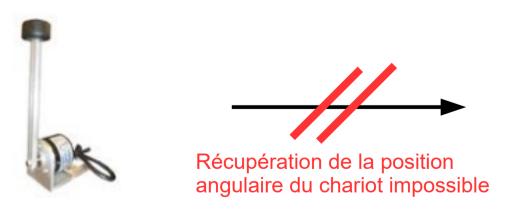


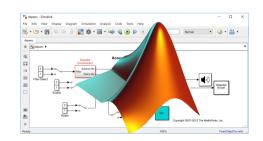
Pilotage par MATLAB



Solution : Réinstallation du système d'exploitation

#### Problème 2: Communication du module impossible



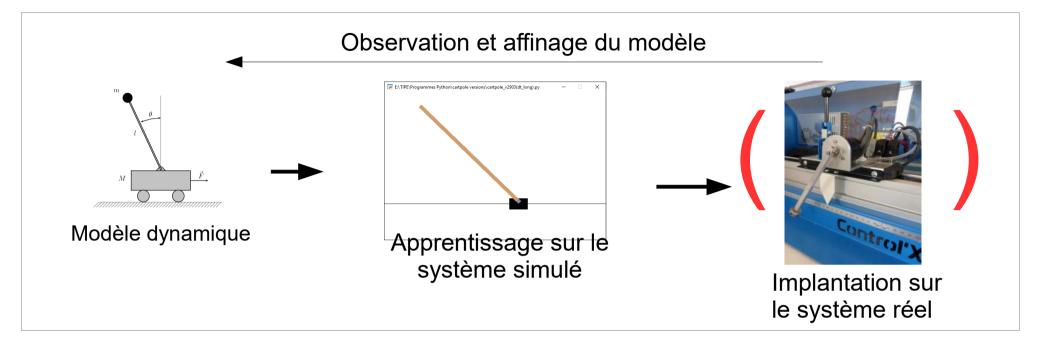


Pilotage par MATLAB

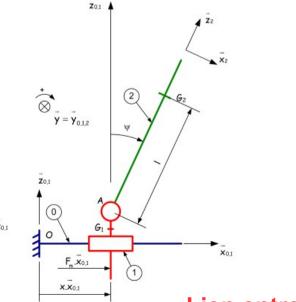
Solution: Non trouvée



#### Système réel



#### Modèle plus fin prenant en compte les frottements



$m_{tot}$ . $\ddot{x} + m_2 J. \ddot{\psi}$ . $\cos \psi - m_2 J. \dot{\psi}^2$ . $\sin \psi$	= $F_m - F_f$ .sign( $\dot{x}$ ) - $f_v . \dot{x}$
$I.\ddot{\psi} + m_2 l.\ddot{x}.\cos\psi = m_2.g.l.\sin\psi$	$-C_{f_2}$ .sign( $\dot{\psi}$ ) - $f_{\omega_2}$ . $\dot{\psi}$

Composant	Modèle linéaire
CNA de la carte de commande	$\varepsilon_3(\dagger) = \varepsilon_2(\dagger)$
Carte de puissance	$u(\dagger) = B.\epsilon_3(\dagger)$
Moteur linéaire équivalent	
Pendule	$I.\ddot{\psi}(\dagger) + m_2.l.\ddot{x}(\dagger) = m_2.g.l.\psi(\dagger) - f_{\omega_2}.\dot{\psi}(\dagger)$

46

nouvelle position[0]+=1



#### Code Python du labyrinthe

L'environnement

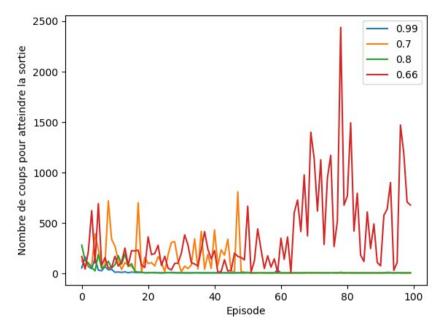
```
##Verifier si bordures
       if nouvelle position[0]<0 or nouvelle position[0]>=len(self.
           recompense=self.Recompenses[3]
           None
       ##Verifier si mur
       elif nouvelle position in self.Mur:
           recompense=self.Recompenses[3]
           None
       ##Verifier si rocher
       elif nouvelle position in self.Rocher:
           recompense=self.Recompenses[1]
           self.position=nouvelle position
       ##Verifier si sortie
       elif nouvelle position == self.sortie:
           recompense=self.Recompenses[2]
           self.fin=True
           self.position=nouvelle position
       ##Aucun problème donc on autorise le déplacement
           self.position=nouvelle position
           recompense=self.Recompenses[0]
       return recompense
   def afficher(self):
       ##Affiche l'état du labvrinthe
       for i in range(len(self.Grille)):
           for j in range(len(self.Grille[0])):
               if [i,j] in self.Mur:
                   self.Grille[i][j]='x'
               elif [i,j] in self.Rocher:
                   self.Grille[i][i]='0'
               elif [i,j] == self.position:
                   self.Grille[i][j]='#'
               elif [i,j]==self.sortie:
                   self.Grille[i][j]='S'
               else:
                   self.Grille[i][j]='-'
       print(self.Grille)
                                                           2/2
# FIN Définition de l'environnement
```

## FIN Algorithme d'apprentissage

##START Algorithm parameters: step size...

```
Coups: 0
[['-' '0' '-' ['S']
 ['-' 'x' '0']
 ['x' '-' '-' '-']
 ['-' '#' 'x' '-']]
Coups: 1
[['-' '0' '-' 'S']
 ['-' 'v' 'x' 'O']
 ['x' '#' '-' '-']
 ['-' '-' 'x' '-']]
Coups: 2
[['-' '0' '-' 'S']
 ['-' 'v' 'x' 'O']
['x' '-' '#' '-']
 ['-' '-' 'x' '-']]
Coups: 3
[['-' '0' '-' 'S']
 ['-' '0' 'x' '0']
 ['x' '-' '-' '#']
 ['-' '-' 'x' '-']]
Coups: 4
[['-' '0' '-' 'S']
 ['-' 'v' 'x' '0']
 ['x' '-' '-' '-']
 ['-' '-' 'x' '-']]
Coups: 5
[['+' '-' '0' '-' |#']
 ['-' 'x' '0']
 ['x' '-' '-' '-']
 ['-' '-' 'x' '-']]
>>>
```

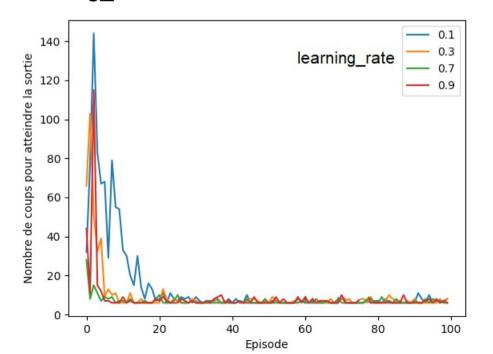




Mesure rendant compte de l'influence des deux paramètres de l'algorithme

$$Q(S_t, A_t) \leftarrow Q(S_t, A_t) + \alpha \left[ R_{t+1} + \gamma \max_{a} Q(S_{t+1}, a) - Q(S_t, A_t) \right]$$

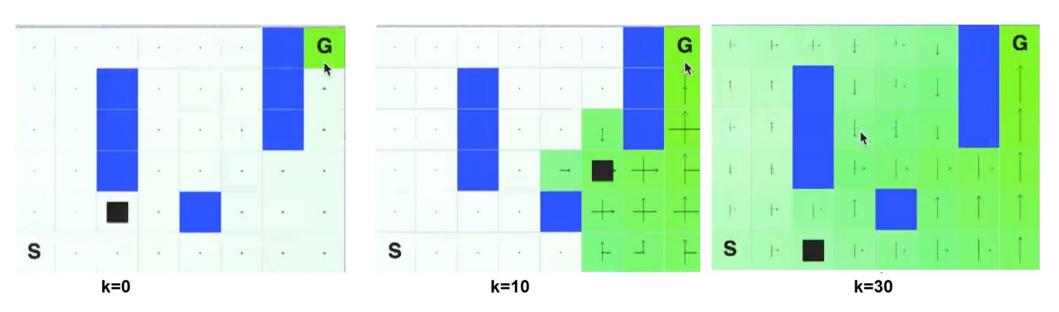
#### Gamma pour Learning\_rate fixé



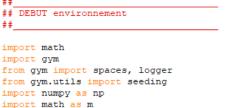
Learning\_rate pour Gamma fixé



Visualisation de la propagation de la récompense de la sortie vers les états précédents



Capture d'écran de la conférence de R.S Sutton [2]



class CartPoleEnv(gym.Env):

#### Code Python du pendule inversé L'environnement

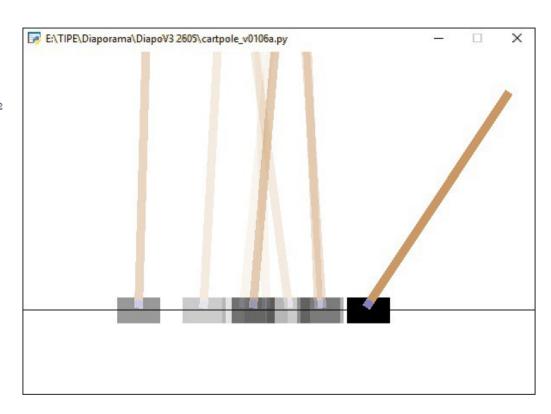
```
metadata = {
    'render.modes': ['human', 'rgb array'],
    'video.frames per second' : 50
def init (self):
    #Paramètres de la modélisation:
    self.gravitv = 9.8
    self.masscart = 1.0
    self.masspole = 0.1
    self.total mass = (self.masspole + self.masscart)
    self.length = .5 # actually half the pole's length
    self.polemass length = (self.masspole * self.length)
    self.force mag = 10.0
   self.tau = 0.02 # dt pour euler
    self.dt system= 0.02  # temps de réaction du système
    #Paramètres de l'instance
    self.Recompenses=[-10,0,1,5] #Récompenses associées à [Zône rouge, Neutre, Zô
    #Limite de fin d'instance et d'intervalle de dicrétisation
    self.x limit= 1 #Zône rouge pour le chariot
   self.theta limit= 30 * 2 * math.pi / 360 #Cône rouge pour la barre
    self.x dot limit=15
    self.theta dot limit=15
    #Paramètre de discrétisation
    self.nb discretisation x=12
    self.nb discretisation x dot=self.nb discretisation theta=self.nb discretisat
    #Intervalles de dicrétisation
    self.X=np.linspace(-self.x limit, self.x limit, self.nb discretisation x)
    self.X dot=np.linspace(-self.x dot limit, self.x dot limit, self.nb discretis
    self.THETA=np.linspace(-self.theta limit, self.theta limit, self.nb discretis
    self.THETA dot=np.linspace(-self.theta dot limit, self.theta dot limit, self.
    high = np.array([
       self.x limit,
       np.finfo(np.float32).max,
       self.theta limit,
       np.finfo(np.float32).max])
    #Autres paramètres
   self.action_space = spaces.Discrete(2)
    self.observation space = spaces.Box(-high, high, dtype=np.float32)
    done=False #Valeur déclarant la fin de l'instance
    self.seed()
    self.viewer = None
    self.state = None
    self.steps beyond done = None
def seed(self, seed=None):
    self.np_random, seed = seeding.np_random(seed)
    return [seed]
```

```
def step(self, action):
    #Modélisation du pendule
    for i in range(int(self.dt system/self.tau)):
        state = self.state
        x, x dot, theta, theta dot = state
        force = self.force mag if action == 1 else -self.force mag
        costheta = math.cos(theta)
        sintheta = math.sin(theta)
        temp = (force - self.polemass length * theta dot * theta dot * sintheta) / self.1
        thetaacc = (self.gravity * sintheta - costheta* temp) / (self.length * (1 - self
        xacc = temp - self.polemass length * thetaacc * costheta / self.total mass
        x dot = x dot + self.tau * xacc
        x = x + self.tau * x dot
        theta dot = theta dot + self.tau * thetaacc
        theta = theta + self.tau * theta dot
        self.state = (x, x dot, theta, theta dot)
        #Associer etat(state) continue à ceux discrets
        self.state discrete=self.continu to discrete(self.state)
    pos cart, vit cart, pos pole, vit pole=self.state discrete
    #Condition de fin de partie, si valeur sors des limites autorisés
    last index state=len(self.X)
    mid index state=len(self.X)//2
    if pos pole==0 or pos pole==last index state or vit pole==0 or vit pole==last index :
        done=True
    else:
        done=False
    #Attribution des récompenses
    negative reward, neutral reward, positive reward, outstanding reward=self.Recompenses
    if (pos pole==5 and (pos cart in [mid index state-1,mid index state,mid index state+
        reward=positive reward
    elif pos pole==mid index state:
        reward=outstanding reward
    elif pos pole==0 or pos pole==last index state or vit pole==0 or vit pole==last index
        reward=negative reward
        reward=neutral reward
    return self.state discrete, reward, done
def reset(self):
    self.state = self.np random.uniform(low=-0.05, high=0.05, size=(4,))
    self.steps beyond done = None
    self.state discrete=self.continu to discrete(self.state)
    return np.array(self.state discrete)
def continu to discrete(self, state):
    state discrete=[0,0,0,0]
    x,x dot,theta,theta dot=state
    for i in range(len(self.X)-1):
        if self.X[i] <= x < self.X[i+1]:</pre>
            state discrete[0]=i
        if self.X dot[i] <= x dot < self.X dot[i+1]:</pre>
            state discrete[1]=i
        if self.THETA[i] <= theta < self.THETA[i+1]:</pre>
            state discrete[2]=i
        if self.THETA dot[i] <= theta dot < self.THETA dot[i+1]:</pre>
            state discrete[3]=i
    return state_discrete
```

```
def render(self, mode='human'):
    screen width = 600
    screen height = 400
    world width = self.x limit*2
    scale = screen width/world width
    carty = 100 # TOP OF CART
    polewidth = 10.0
    polelen = scale * (2 * self.length)
    cartwidth = 50.0
    cartheight = 30.0
    if self.viewer is None:
        from gym.envs.classic control import rendering
        self.viewer = rendering.Viewer(screen width, screen height)
        1,r,t,b = -cartwidth/2, cartwidth/2, cartheight/2, -cartheight/2
        axleoffset =cartheight/4.0
        cart = rendering.FilledPolygon(((1,b), (1,t), (r,t), (r,b)))
        self.carttrans = rendering.Transform()
        cart.add attr(self.carttrans)
        self.viewer.add geom(cart)
        1,r,t,b = -polewidth/2,polewidth/2,polelen-polewidth/2,-polewidth/2
        pole = rendering.FilledPolygon([(1,b), (1,t), (r,t), (r,b)])
        pole.set color(.8,.6,.4)
        self.poletrans = rendering.Transform(translation=(0, axleoffset))
        pole.add attr(self.poletrans)
        pole.add attr(self.carttrans)
        self.viewer.add geom(pole)
        self.axle = rendering.make circle(polewidth/2)
        self.axle.add attr(self.poletrans)
        self.axle.add attr(self.carttrans)
        self.axle.set color(.5,.5,.8)
        self.viewer.add geom(self.axle)
        self.track = rendering.Line((0,carty), (screen width,carty))
        self.track.set color(0,0,0)
        self.viewer.add geom(self.track)
        self. pole geom = pole
    if self.state is None: return None
    # Edit the pole polygon vertex
    pole = self. pole geom
    1,r,t,b = -polewidth/2,polewidth/2,poleen-polewidth/2,-polewidth/2
    pole.v = [(1,b), (1,t), (r,t), (r,b)]
    x = self.state
    cartx = x[0]*scale+screen width/2.0 # MIDDLE OF CART
    self.carttrans.set_translation(cartx, carty)
    self.poletrans.set rotation(-x[2])
    return self.viewer.render(return rgb array = mode=='rgb array')
def close(self):
    if self.viewer:
        self.viewer.close()
        self.viewer = None
```



### Code Python du pendule inversé *L'environnement*



## FIN environnement



### Code Python du pendule inversé *L'environnement*

```
def step(self, action):
    #Modélisation du pendule
    for i in range(int(self.dt_system/self.tau)):
       state = self.state
       x, x dot, theta, theta dot = state
       force = self.force mag if action == l else -self.force mag
        costheta = math.cos(theta)
        sintheta = math.sin(theta)
        temp = (force - self.polemass length * theta dot * theta dot * sintheta) / self.total mass
        thetaacc = (self.gravity * sintheta - costheta* temp) / (self.length * (1 - self.masspole * costheta * costheta / self.total mass))
        xacc = temp - self.polemass length * thetaacc * costheta / self.total mass
        x dot = x dot + self.tau * xacc
        x = x + self.tau * x dot
        theta_dot = theta_dot + self.tau * thetaacc
        theta = theta + self.tau * theta dot
        self.state = (x,x_dot,theta,theta_dot)
        #Associer etat(state) continue à ceux discrets
        self.state discrete=self.continu to discrete(self.state)
```

#### Résolution par Euler

```
## DEBUT Algorithme d'apprentissage
def Q learning(env,learning rate,gamma,iteration):
    import random as r
    ##START Algorithm parameters: step size...
    1r=learning rate
    v=gamma
    iteration=iteration
    ##END Algorithm parameters: step size...
    ##START Initialize Q(s,a)...arbitrarly execpt that Q(terminal,.) =
    n states=env.nb discretisation x
    n actions=2
    Q=np.zeros([n states, n states, n states, n states, n actions])
    for i in range(len(Q)):
        for j in range(len(Q[0])):
                for k in range(len(Q[0][0])):
                    for 1 in range(len(Q[0][0][0])):
                            for m in range(len(Q[0][0][0][0])):
                                    Q[i,j,k,l,m]=r.random()
    ##END Initialize Q(s,a)...arbitrarly execpt that Q(terminal,.)=0
    ##START Loop for each episode
    for i in range (iteration):
        #Initialize S
        end=False
        s=env.reset()
        ##START Loop for each step of epsiode
        while end==False:
            ## START Choose A from S using
            a=np.argmax(Q[s[0],s[1],s[2],s[3],:])
            ## END Choose A from S using
            #Take action A
            observation=env.step(a)
            #Observe R,S'
            sl,r,end=observation
            \#Refresh Q: Q(S,A)=Q(S,A)+...
            Q[s[0], s[1], s[2], s[3], a] = Q[s[0], s[1], s[2], s[3], a] + 1r*(r + y * np.max(Q[s1]))
            #S becomes S'
            s=sl.copv()
        ##END Loop for each step of epsiode
    ##END Loop for each episode
    return Q
## FIN Algorithme d'apprentissage
env=CartPoleEnv()
learning rate=0.8
gamma=0.99
iteration=201
```

Q learning(env,learning rate,gamma,iteration)



Code Python du pendule inversé L'agent

```
self.viewer = None
Classic cart-pole system implemented by Rich Sutton et al.
                                                                                       self.state = None
Copied from http://incompleteideas.net/sutton/book/code/pole.c
permalink: https://perma.cc/C9ZM-652R
                                                                                       self.steps beyond done = None
import math
import gym
                                                                                  def seed(self, seed=None):
from gym import spaces, logger
                                                                                       self.np random, seed = seeding.np random(seed)
from gym.utils import seeding
import numpy as np
                                   Code Python du pendule inversé
                                                                                      return [seed]
class CartPoleEnv(gym.Env):
                                          L'environnement (alpha)
                                                                                  def step(self. action):
                                                                                       assert self.action space.contains(action), "%r (%s) invalid"%(action,
      A pole is attached by an un-actuated joint to a cart, which moves al
                                                                                       state = self.state
      This environment corresponds to the version of the cart-pole problem
                                                                                       x, x dot, theta, theta dot = state
   Observation:
                                                                                       force = self.force mag if action == 1 else -self.force mag
      Type: Box(4)
      Num Observation
                                   Min
                                                                                       costheta = math.cos(theta)
      0 Cart Position
                                  -4 8
                                               4 8
                                                                                       sintheta = math.sin(theta)
         Cart Velocity
                                -Inf
                                                Inf
         Pole Angle
                                  -24 deg
                                                24 dea
                                                                                       temp = (force + self.polemass length * theta dot * theta dot * sinthet;
            Pole Velocity At Tip -Inf
                                                                                       thetaacc = (self.gravity * sintheta - costheta* temp) / (self.length *
                                                                                       xacc = temp - self.polemass length * thetaacc * costheta / self.total
   Actions:
      Type: Discrete(2)
                                                                                       if self.kinematics integrator == 'euler':
      Num Action
                                                                                           x = x + self.tau * x dot
            Push cart to the left
            Push cart to the right
                                                                                           x dot = x dot + self.tau * xacc
                                                                                           theta = theta + self.tau * theta dot
      Note: The amount the velocity that is reduced or increased is not fi
                                                                                           theta dot = theta dot + self.tau * thetaacc
   Reward:
      Reward is 1 for every step taken, including the termination step
                                                                                       else: # semi-implicit euler
   Starting State:
                                                                                           x dot = x dot + self.tau * xacc
      All observations are assigned a uniform random value in [-0.05..0.05
   Episode Termination:
                                                                                           x = x + self.tau * x dot
      Pole Angle is more than 12 degrees
                                                                                           theta dot = theta dot + self.tau * thetaacc
      Cart Position is more than 2.4 (center of the cart reaches the edge
      Episode length is greater than 200
                                                                                           theta = theta + self.tau * theta dot
      Solved Requirements
                                                                                       self.state = (x, x dot, theta, theta dot)
      Considered solved when the average reward is greater than or equal t
                                                                                       done = x < -self.x threshold \
                                                                                                or x > self.x threshold \
   metadata = {
                                                                                                or theta < -self.theta threshold radians \
       'render.modes': ['human', 'rgb array'],
       'video.frames per second' : 50
                                                                                                or theta > self.theta threshold radians
                                                                                       done = bool(done)
   def init (self):
      self.gravity = 9.8
                                                                                       if not done:
      self.masscart = 1.0
                                                                                           reward = 1.0
      self.masspole = 0.1
      self.total mass = (self.masspole + self.masscart)
                                                                                       elif self.steps beyond done is None:
      self.length = 0.5 # actually half the pole's length
                                                                                           # Pole just fell!
      self.polemass length = (self.masspole * self.length)
      self.force mag = 10.0
                                                                                           self.steps beyond done = 0
      self.tau = 0.02 # seconds between state updates
                                                                                           reward = 1.0
      self.kinematics integrator = 'euler'
                                                                                       else:
      # Angle at which to fail the episode
                                                                                           if self.steps beyond done == 0:
      self.theta threshold radians = 12 * 2 * math.pi / 360
                                                                                                logger.warn("You are calling 'step()' even though this environ
      self.x threshold = 2.4
                                                                                           self.steps beyond done += 1
      \sharp Angle limit set to 2 * theta threshold radians so failing observat
                                                                                           reward = 0.0
      high = np.array([
          self.x threshold * 2,
          np.finfo(np.float32).max,
                                                                                       return np.array(self.state), reward, done, {}
          self.theta threshold radians * 2,
          np.finfo(np.float32).max])
                                                                                  def reset(self):
      self.action_space = spaces.Discrete(2)
                                                                                       self.state = self.np random.uniform(low=-0.05, high=0.05, size=(4,))
      self.observation space = spaces.Box(-high, high, dtype=np.float32)
                                                                                       self.steps beyond done = None
      self.seed()
                                                                                       return np.array(self.state)
      self.viewer = None
      self.state = None
```

```
def render(self, mode='human'):
    screen width = 600
    screen height = 400
    world width = self.x threshold*2
    scale = screen width/world width
    carty = 100 # TOP OF CART
    polewidth = 10.0
    polelen = scale * (2 * self.length)
    cartwidth = 50.0
    cartheight = 30.0
    if self.viewer is None:
        from gym.envs.classic control import rendering
        self.viewer = rendering.Viewer(screen width, screen height)
        1,r,t,b = -cartwidth/2, cartwidth/2, cartheight/2, -cartheight/2
        axleoffset =cartheight/4.0
        cart = rendering.FilledPolygon([(1,b), (1,t), (r,t), (r,b)])
        self.carttrans = rendering.Transform()
        cart.add attr(self.carttrans)
        self.viewer.add geom(cart)
        1, r, t, b = -polewidth/2, polewidth/2, polelen-polewidth/2, -polewidth/2
        pole = rendering.FilledPolygon([(1,b), (1,t), (r,t), (r,b)])
        pole.set color(.8,.6,.4)
        self.poletrans = rendering.Transform(translation=(0, axleoffset))
        pole.add attr(self.poletrans)
        pole.add attr(self.carttrans)
        self.viewer.add geom(pole)
        self.axle = rendering.make circle(polewidth/2)
        self.axle.add attr(self.poletrans)
        self.axle.add attr(self.carttrans)
        self.axle.set color(.5,.5,.8)
        self.viewer.add geom(self.axle)
        self.track = rendering.Line((0,carty), (screen width,carty))
        self.track.set color(0,0,0)
        self.viewer.add geom(self.track)
        self. pole geom = pole
    if self.state is None: return None
    # Edit the pole polygon vertex
    pole = self. pole geom
   1,r,t,b = -polewidth/2,polewidth/2,polelen-polewidth/2,-polewidth/2
   pole.v = [(l,b), (l,t), (r,t), (r,b)]
    x = self.state
    cartx = x[0]*scale+screen width/2.0 # MIDDLE OF CART
    self.carttrans.set translation(cartx, carty)
    self.poletrans.set rotation(-x[2])
    return self.viewer.render(return rgb array = mode=='rgb array')
def close(self):
    if self.viewer:
                                                        3/3
        self.viewer.close()
        self.viewer = None
```



Code Python du pendule inversé *L'environnement (alpha)* 



#### Code Python du pendule inversé *Mesures*

