



ÉCOLE NATIONALE SUPÉRIEURE DE TECHNIQUES AVANCÉES

# ROB311 : APPRENTISSAGE POUR LA ROBOTIQUE

Parcours robotique

Année universitaire : 2020/2021

---

## TP3: Q-Learning et Pac-man

---

**Auteurs :**

M. AHMED YASSINE HAMMAMI

MME. HANIN HAMDİ

## Objectifs

L'objectif de ce TP est d'implémenter l'algorithme de Q-Learning en Python. En particulier, on va former un joueur IA du célèbre jeu d'arcade Pac-Man.

On va travailler sur un ensemble de fichiers et de bibliothèques Python afin de modifier certaines fonctions pour implémenter l'IA Q-Learning et Approximate Q-Learning pour le jeu Pac-Man.

## L'algorithme Q-Learning

1. Initialiser  $Q(s,a)$  arbitrairement, pour tout état  $s$  et pour toute action  $a$ .
2. Répéter pour chaque épisode :
  - (a) Trouver pour chaque action et pour chaque état de la matrice  $Q$  et calculer la fonction de valeur à partir de l'équation de Bellman :
$$V^*(S_t) = R(S_t) + \max_a \gamma \sum T(S_t, a, S_{t+1}) V^*(S_{t+1})$$
  - (b) Pour chaque état, choisir l'action optimale et mettre à jour  $Q$  selon :
$$Q(S_t) = \arg \max_a \sum T(S_t, a, S_{t+1}) V^*(S_{t+1})$$
  - (c) Répéter ces deux étapes pour chaque état jusqu'à convergence. Dans notre cas, on considère que l'algorithme converge si la matrice  $Q$  ne change pas.

## Implémentation des fonctions de la classe *QLearningAgent(ReinforcementAgent)* :

### *init()* :

Dans cette fonction, on va initialiser les Qvalues :

```
1 def __init__(self, **args):
2     ReinforcementAgent.__init__(self, **args)
3     self.QValues = util.Counter()
```

### *getvalues()* :

Cette fonction retourne la Q-valeur[state,action]

```
1 def getQValue(self, state, action):
2     return self.QValues[state, action]
```

### *computeValueFromQValues()* :

Cette fonction retourne l'action maximale parmi tous les actions possibles, s'il n'y a pas d'action possible, la fonction retourne 0.0.

```
1 def computeValueFromQValues(self, state):
2     all_values = []
3     for action in self.getLegalActions(state) :
```

```
1     vv=self.getQValue(state, action)
2     all_values.append(vv)
3
4     if (all_values):
5         return max(all_values)
6     else:
7         return 0.0
```

### ***computeActionFromQValues()*** :

Cette fonction retourne la meilleure action qu'on peut avoir dans un état. S'il n'y a pas d'action possible, elle retourne None.

```
1     def computeActionFromQValues(self, state):
2
3         #we start first of all to determine the legal actions
4         l_actions = self.getLegalActions(state)
5
6         value = self.getValue(state)
7         for action in l_actions:
8             if (value == self.getQValue(state, action)):
9                 return action
10        return None
```

### ***getAction()*** :

Elle calcule l'action dans l'état actuel avec la probabilité epsilon. On choisit une action aléatoire dans un état, sinon on choisit la "Policy" optimale actuelle.

```
1     def getAction(self, state):
2
3         legalActions = self.getLegalActions(state)
4         action = None
5
6         if (util.flipCoin(self.epsilon)):
7             action = random.choice(legalActions)
8         else:
9             action = self.getPolicy(state)
10
11        return action
```

### ***update()*** :

```
1     def update(self, state, action, nextState, reward):
2
3         #calculate the new Qvalue
4         new = (1 - self.alpha) * self.getQValue(state, action)
5         new += self.alpha * (reward + (self.discount * self.getValue(nextState)))
```

```
6 self.QValues[state, action] = new
```

## Implémentation des fonctions de la classe *ApproximateQAgent(PacmanQAgent)* :

Dans cette partie, on va implémenté un agent Q-learning approximatif. On commence tout d'abord par former l'agent puis le tester.

### *getQValue()* :

```
1 def getQValue(self, state, action):
2
3     features = self.featExtractor.getFeatures(state,action)
4     Qvalue = 0.0
5
6     for f in features:
7         QValue += self.weights[f] * features[f]
8
9     return QValue
```

L'idée ici est que QValue est représentée comme une somme linéaire pondérée des valeurs de caractéristiques pour un état, nous avons maintenant besoin d'un moyen pour obtenir la valeur de caractéristique de l'état afin de mettre à jour les poids.

Après avoir effectuer une action à partir d'un état, on passe à un nouvel état et on obtient un "reward", basée sur *features*, *reward* et *l'état suivant*.

### *update()* :

```
1 def update(self, state, action, nextState, reward):
2
3     QValue = 0
4     difference = reward + (self.discount * self.getValue(nextState) - self.getQValue(
5     ↪ state, action))
6     features = self.featExtractor.getFeatures(state, action)
7
8     for feature in features:
9         self.weights[feature] += self.alpha * features[feature] * difference
```

## Test et résultat

Pour tester l'ensemble des fonctions implémentées, on va lancer cette commande *python pacman.py -p PacmanQAgent -x 2000 -n 2010 -l smallGrid*.

Cette commande va apprendre à partir de 2000 épisodes d'entraînement, on testera ensuite l'IA résultante sur 10 matchs.

```
Training Done (turning off epsilon and alpha)
-----
Pacman emerges victorious! Score: 503
Pacman emerges victorious! Score: 503
Pacman emerges victorious! Score: 499
Pacman emerges victorious! Score: 503
Pacman emerges victorious! Score: 495
Pacman emerges victorious! Score: 503
Pacman emerges victorious! Score: 503
Pacman emerges victorious! Score: 499
Pacman emerges victorious! Score: 503
Pacman emerges victorious! Score: 503
Average Score: 501.4
Scores:      503.0, 503.0, 499.0, 503.0, 495.0, 503.0, 503.0, 499.0, 503.0, 503.0
Win Rate:    10/10 (1.00)
Record:      Win, Win, Win, Win, Win, Win, Win, Win, Win, Win
```

FIGURE 1 – Résultat du test

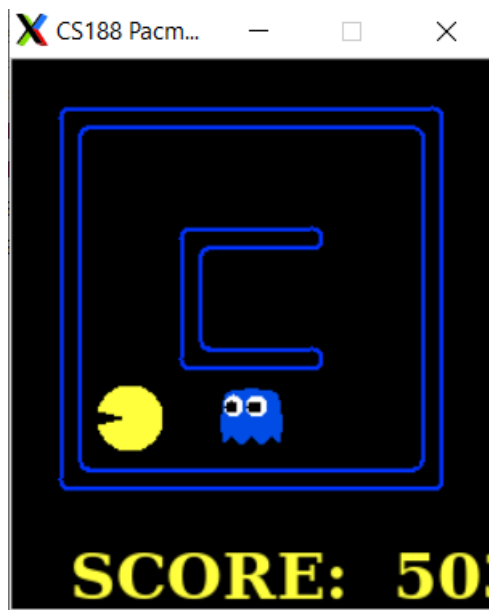


FIGURE 2 –

On remarque que le win Rate est égale à 10/10 ce qui affirme que notre algorithme est en train d'apprendre correctement dans la phase training.