

Présentation finale du PFE

Apprentissage par renforcement multi-agents

David Albert

INSA Rouen

21 février 2020

- 1 L'apprentissage
 - L'apprentissage automatique
 - L'apprentissage par renforcement
- 2 Résolution des jeux stochastiques
 - Jeux stochastiques
 - Théorie des jeux
 - Résolution des jeux stochastiques
- 3 L'API *marl*
 - Présentation
 - Les environnements Gym
 - Soccer
- 4 Expérimentations
 - Mono-agent
 - Multi-agent
- 5 Conclusion

- **Objectif** : Construire P_{model} proche de P_{reel}

- **Objectif** : Construire P_{model} proche de P_{reel}
- Supervisé:

$$P\left(\mathbf{y} = \text{cat} \mid \text{img}\right)$$



- **Objectif** : Construire P_{model} proche de P_{reel}
- Supervisé:

$$P\left(\mathbf{y} = \text{cat} \mid \text{img}\right)$$

- Non-supervisé:

$$\text{img} \sim P(\mathbf{x})$$

- **Objectif** : Construire P_{model} proche de P_{reel}
- Supervisé:

$$P\left(\mathbf{y} = \text{cat} \mid \text{img}\right)$$

- Non-supervisé:

$$\text{img} \sim P(\mathbf{x})$$

- Séquentiel:

$$P(X_{t+1} = \text{"pluie"} \mid X_t = \text{"soleil"})$$

- **Objectif** : Construire P_{model} proche de P_{reel}
- Supervisé:

$$P\left(\mathbf{y} = \text{cat} \mid \text{img}\right)$$

- Non-supervisé:


$$\sim P(\mathbf{x})$$

- Séquentiel:

$$P(X_{t+1} = \text{"pluie"} \mid X_t = \text{"soleil"})$$

- Renforcement:

$$P(X_{t+1} = \text{"Chomage"} \nearrow, R_{t+1} = -100 \mid X_t = \text{"Chomage"} \searrow, A_t = \text{"Reforme"})$$

$$\text{et } P(A_t = \text{"Reforme"} \mid X_t = \text{"Chomage"} \searrow)$$

- Etude d'un système composé de 2 éléments:
 - L'environnement (Le système économique français)
→ $p(s_t, r_t, o_t \mid s_1, r_1, o_1, a_1, s_2, \dots, s_{t-1}, r_{t-1}, o_{t-1}, a_{t-1})$
 - L'agent (M. Macron)
→ $p(a_t \mid r_1, o_1, a_1, \dots, a_t, r_t, o_t)$
- Planification vs Apprentissage
- Trois approches d'apprentissage:
 1. Estimation du modèle de transition \mathcal{T} + planification
 2. Estimation de la fonction de valeur q_π
 3. Optimisation direct de la politique d'action π

Problème: Prise de décision repose sur l'intégralité du passé.

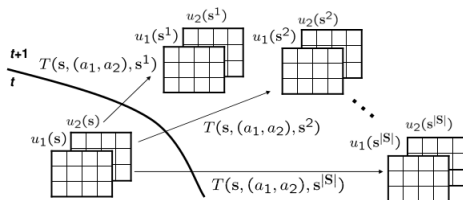
Solution: Utiliser un modèle simplifié d'évolution (MDP).

Definition

Les **Processus de décision de Markov** est un modèle d'apprentissage par renforcement qui suit la propriété de Markov ("le futur ne dépend que du présent").

Jeux stochastiques (SG) peuvent être vus comme :

1. extension des MDPs pour N agents (RL)
2. jeux matriciels à étapes (Théorie des Jeux)



Types de jeux :

- Coopératifs / **Non coopératifs**
- **Forme normale** / Forme extensive

Agents rationnels:

- exploite défaut des adversaires
- équilibre de Nash

Q-learning pour jeux à deux joueurs compétitifs.

- construit $Q_{s,a,o}$ pour chaque joueur
- $Q_{s,..}$ jeu matriciel



$Q_1(s)$ et $Q_2(s)$

$A_1 \setminus A_2$	Front	Back	Left	Right
Front	10 -10	10 -10	6 -6	10 -10
Back	8 -8	7 -7	5 -5	7 -7
Left	3 -3	5 -5	7 -7	5 -5
Right	8 -8	7 -7	5 -5	7 -7

Mise à jour:

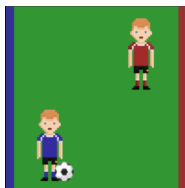
$$Q_i(s, a, o) = Q_i(s, a, o) + \alpha [r_t + \gamma \max_a \min_o Q_i(s', a, o) - Q_i(s, a, o)]$$

Choix de la politique:

$$\pi_i(s) = \arg \max_a \min_o Q_i(s, a, o)$$

Q-learning pour jeux à deux joueurs compétitifs.

- construit $Q_{s,a,o}$ pour chaque joueur
- $Q_{s,..}$ jeu matriciel



$Q_1(s)$ et $Q_2(s)$

$A_1 \setminus A_2$	Front	Back	Left	Right
Front	10 -10	10 -10	6 -6	10 -10
Back	8 -8	7 -7	5 -5	7 -7
Left	3 -3	5 -5	7 -7	5 -5
Right	8 -8	7 -7	5 -5	7 -7

Interprétation: Q-learning pessimiste : $Q_i(s, a) = \min_o Q_i(s, a, o)$

Problèmes:

- Doit visiter indéfiniment tous les triplet ($etat/a_1/a_2$) pour converger
- Politique déterministe
- Seulement 2 joueurs

Motivation: Améliorer directement la politique

Policy Hill Climbing (PHC):

- Idée similaire aux *Actor-Critic*
- *Critic* : Q-learning standard

$$Q_i(s, a_i) = Q_i(s, a_i) + \alpha[r_t + \gamma \max_{a'_i} Q_i(s', a'_i) - Q_i(s, a_i)]$$

- *Actor* : ↗ proba de choisir l'action qui maximise Q

$$\pi_i(s, a) = \pi_i(s, a) + \begin{cases} \delta & \text{si } a = \operatorname{argmax}_{a'} Q(s, a') \\ \frac{-\delta}{|A_i|-1} & \text{sinon} \end{cases}$$

Win or Learn Fast (WoLF)

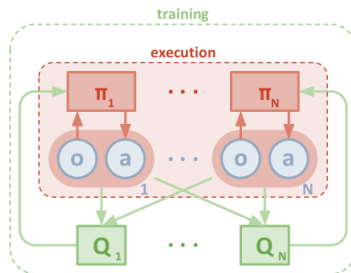
Principe d'apprentissage selon lequel un agent qui gagne mettra moins vite à jour son modèle (δ_{win} & δ_{lose}).

Motivation: Nouvelles méthodes de gradient de politique

Multi-agent Deep Deterministic Policy Gradient (MADDPG):

- Politiques déterministes/Actions continues
- *Critic*: Deep Q-learning
- *Actor*: Descente de gradient

$$\nabla_{\theta_i} J = \frac{1}{S} \sum_j \nabla_{\theta_i} \mu_i(o_i^j) \nabla_{a_i} Q_i(s^j, a_1^j, \dots, a_i^j, \dots, a_N^j) |_{a_i = \mu_i(o_i)}$$



- Vision généralisée de l'apprentissage multi-agent
- Système d'agents "**mixtes**" :
 - Agents entraînaables, Agents pré-entraînés, Bots, Humains, etc
- Algorithmes disponibles:

Single-agent algorithms

Q-learning	DQN	Actor-Critic	DDPG	TD3
✓	✓	✓	✓	✗

Multi-agent algorithms

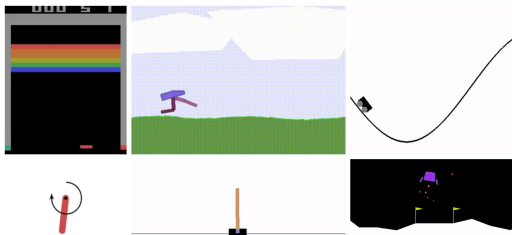
minimaxQ	PHC	JAL	MAAC	MADDPG
✓	✓	✗	✓	✓

- Documentation sphinx: <https://blavad.github.io/marl>
- Formalisme de jeux stochastiques
 - 1 récompense / agent
 - pas de communication explicite
- Compatibilité : environnements *Gym*

L'environnement : la clé de l'apprentissage

- Environnements *Gym*

- 2 attributs : *observation_space*, *action_space*
- 3 méthodes : *reset()*, *step(action)*, *render()*



- Jeux stochastiques \Leftrightarrow à priori non-coopératif

- Formalisme coopératif : $r_i = r_j \forall i, j$

- Jeux stochastiques \Leftrightarrow pas de communication explicite

- Communication implicite de i vers j : $o_j^{(t+1)} = \text{concat}(c_j^{(t+1)}, a_i^{(t)})$

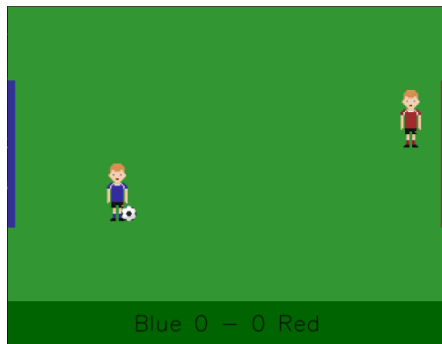
Deux types d'environnements:

- 1) Environnement discret
- 2) Environnement continue

Deux types d'environnements:

1) Environnement discret

- Grille $h \times w$ (ici 4×5)
- Action : $a \in \{none, front, back, left, right\}$
- Deux types d'observations
 - Etat exact : $s \in [0, \dots, n \times (hw)^n]$
 - Map : $s \in \{0, 1\}^h \times \{0, 1\}^w \times \{0, 1\}^3$

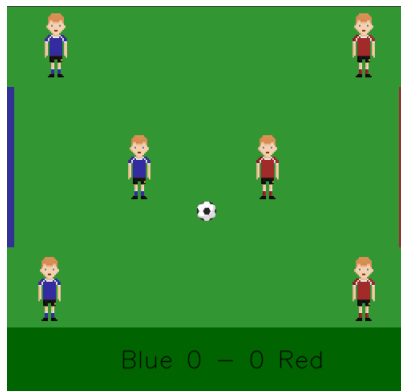


Deux types d'environnements:

2) Environnement continue

→ Action : $a \in \{none, front, back, left, right\}$

→ Observation : $o \in [-1, 1]^{5+2 \times n}$



Minimax-Q vs Random

- Minimax-Q : 93.6 %
- Random : 9.4 %
- Match nul : $\leq 2\%$

Minimax-Q vs Q-learning

- Minimax-Q : 100.0 %
- Q-learning : 0.0 %
- Match nul : $\geq 70\%$

Q-learning vs Random

- Q-learning : 95.9 %
- Random : 4.1 %
- Match nul : $\leq 1\%$

MinimaxQ vs PHC

- Minimax-Q : 58.2 %
- PHC : 41.8 %
- Match nul : 1 %

PHC vs Random

- PHC : 82.3 %
- Random : 17.7 %
- Match nul : $\leq 10\%$

Expérimentations → *DQN*

Soccer 5x5

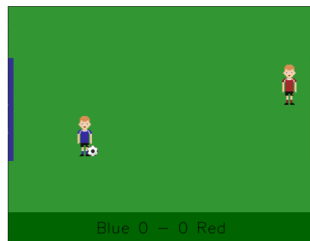
Problème: Etat de l'environnement devient vite trop grand.

Solution:

- DQN sur terrain 5x5
- CNN pour estimer la valeur Q



DQN



- Compétences développées:
 - Théorie sur l'apprentissage multi-agents
 - Programmation python et pytorch
 - Compréhension de méthodes de RL récentes (ex: DDPG)
 - Compréhension des problématiques liées au MARL
 - Développement d'un environnement Gym