



**HAL**  
open science

## Algorithmes de sortie du Piège de la Zone Ennuyeuse en apprentissage par renforcement

Landry Steve Noulawe Tchamanbe, Paulin Melatagia Yonta

► **To cite this version:**

Landry Steve Noulawe Tchamanbe, Paulin Melatagia Yonta. Algorithmes de sortie du Piège de la Zone Ennuyeuse en apprentissage par renforcement. *Revue Africaine de Recherche en Informatique et Mathématiques Appliquées*, 2021, Volume 34 - 2020 - Special Issue CARI 2020 (CARI 2020), 10.46298/arima.6748 . hal-02925738v3

**HAL Id: hal-02925738**

**<https://hal.science/hal-02925738v3>**

Submitted on 2 Jul 2021

**HAL** is a multi-disciplinary open access archive for the deposit and dissemination of scientific research documents, whether they are published or not. The documents may come from teaching and research institutions in France or abroad, or from public or private research centers.

L'archive ouverte pluridisciplinaire **HAL**, est destinée au dépôt et à la diffusion de documents scientifiques de niveau recherche, publiés ou non, émanant des établissements d'enseignement et de recherche français ou étrangers, des laboratoires publics ou privés.

## Algorithmes de sortie du Piège de la Zone Ennuyeuse en apprentissage par renforcement

Landry Steve Noulawe Tchamanbe<sup>1</sup> and Paulin Melatagia Yonta<sup>1,2</sup>

<sup>1</sup>Département d'informatique, Université de Yaoundé I, Cameroun

<sup>2</sup>Sorbonne Université, IRD, UMMISCO, F-93143, Bondy, France

\*E-mail : [stevenoulawe@gmail.com](mailto:stevenoulawe@gmail.com) — [paulinyonta@gmail.com](mailto:paulinyonta@gmail.com)

DOI : [10.46298/arima.6748](https://doi.org/10.46298/arima.6748)

Soumis le 01 Septembre 2020 - Publié le 24 Juin 2021

Volume : 34 - Année : 2021

Numéro spécial : CARI 2020

Éditeurs : Nabil Gmati, Mathieu Roche, Tri Nguyen-Huu, Laurent Debreu

---

### Résumé

Les algorithmes d'apprentissage par renforcement ont réussi au fil des années à obtenir des résultats impressionnants dans divers domaines. Cependant, ces algorithmes souffrent de certaines faiblesses mises en évidence par Refael Vivanti et al. qui peuvent expliquer la régression des agents même bien entraînés dans certains environnements : la différence de variance sur les récompenses entre les zones de l'environnement. Cette différence de variance conduit à deux problèmes : le Piège de la Zone Ennuyeuse (Boring Area Trap) et le Consultant Manipulateur. Nous observons que l'algorithme Adaptive Symmetric Reward Noising (ASRN) proposé par Refael Vivanti et al. présente des limites pour des environnements ayant les caractéristiques suivantes : longues durées de jeu et environnement à zones ennuyeuses multiples. Pour pallier ces problèmes, nous proposons trois algorithmes dérivés de l'algorithme ASRN nommés Rebooted Adaptive Symmetric Reward Noising (RASRN) : Continuous  $\varepsilon$  decay RASRN, Full RASRN et Stepwise  $\alpha$  decay RASRN. Grâce à deux séries d'expérimentations menées sur le problème du bandit à k bras, nous montrons que ces algorithmes permettent de mieux corriger le problème du piège de la zone ennuyeuse.

### Mots-Clés

Apprentissage par renforcement, Piège de la Zone Ennuyeuse, ASRN, bandit à k bras.

---

## I INTRODUCTION

L'apprentissage automatique [2] est un sous-domaine de l'intelligence artificielle qui se fonde sur des approches mathématiques et statistiques pour donner aux ordinateurs la capacité d'« apprendre » à partir de données, c'est-à-dire d'améliorer leurs performances à résoudre des tâches sans être explicitement programmés. L'apprentissage automatique à son tour se subdivise

en plusieurs sous-domaines dont les principaux sont : l'apprentissage supervisé, l'apprentissage non supervisé et l'apprentissage par renforcement. On parle d'apprentissage supervisé lorsque les données utilisées pour entraîner la machine sont associées à des *étiquettes* qui représentent les bonnes réponses à la tâche que la machine doit exécuter. Si par contre les données ne sont pas étiquetées et qu'il s'agit pour la machine de découvrir des motifs à l'intérieur de ces dernières, il s'agit de l'apprentissage non supervisé. Enfin lorsqu'il s'agit pour la machine d'apprendre à réaliser une séquence d'actions dans le but de maximiser une récompense totale, il s'agit d'apprentissage par renforcement.

Ce papier porte sur l'apprentissage par renforcement. Parmi les succès récents de l'apprentissage par renforcement nous pouvons citer : en 2015, Mnih et al. [6] qui ont montré que l'apprentissage par renforcement permet de créer un programme qui joue à des jeux Atari. En 2018, Hessel et al. [4] ont combiné plusieurs techniques de l'apprentissage par renforcement pour améliorer les performances de l'algorithme Deep Q-network sur le benchmark Atari 2600. AlphaGo Zero [9] est une nouvelle technique d'apprentissage par renforcement où l'agent apprend à jouer au jeu de Go en jouant contre lui-même. De nombreux algorithmes d'apprentissage par renforcement ont été proposés dans la littérature et malgré qu'ils ont en général de bonnes performances, ils présentent pour la plupart d'entre eux le problème du piège de la zone ennuyeuse qui lorsqu'il se produit dégrade considérablement leurs performances. Dans ce papier, nous proposons des algorithmes de correction de ce problème pour l'algorithme Q-learning.

Dans la suite de ce papier, nous présenterons à la Section 2 quelques notions fondamentales de l'apprentissage par renforcement. Dans la Section 3, nous présenterons le problème du Piège de la Zone Ennuyeuse (Boring Area Trap) ainsi que celui du Consultant Manipulateur (Manipulative Consultant) et également les solutions proposées dans l'état de l'art au problème du Boring Area Trap. Dans la Section 4, nous décrivons la solution nommée Rebooted Adaptive symmetric Reward Noising (RASRN) que nous proposons à ce problème. Dans la Section 5, les résultats des expérimentations menées avec RASRN sur le problème du bandit à k bras seront présentés. La Section 6 conclura le papier.

## II NOTIONS FONDAMENTALES DE L'APPRENTISSAGE PAR RENFORCEMENT

Dans l'apprentissage par renforcement, on dispose d'un agent plongé dans un *environnement* caractérisé par un ensemble d'états ( $s$ ) et un ensemble d'actions ( $a$ ) qu'il peut effectuer. Chaque action donne lieu à une *récompense* ( $r$ ), bonus ou malus et permet d'accéder à un autre état (sauf pour des états dits terminaux). L'objectif de l'agent est de déterminer une *politique* (*stratégie*) *optimale* ; celle qui lui permettra de maximiser son gain à long terme [2]. Une politique notée  $\pi$  est une fonction qui à chaque état associe une action à exécuter. Dans certains cas, l'action exécutée par l'agent peut être le fruit du hasard, notamment au début de l'apprentissage où l'agent ne connaît pas encore son environnement : c'est ce que l'on appelle *l'exploration* qui est régulé grâce à un hyper-paramètre :  $\varepsilon$ , le *taux d'exploration*. Pour atteindre une politique optimale, l'agent procède de façon itérative en bouclant sur les deux phases suivantes [1, 2] :

1. **Phase d'évaluation** : pendant laquelle l'agent évalue la politique courante au travers de sa fonction d'utilité  $V$  (respectivement  $Q$ ).

2. **Phase d'amélioration** : pendant laquelle l'agent met à jour sa politique au travers de sa fonction d'utilité  $V$  (respectivement  $Q$ ).

On note  $V^\pi$  (respectivement  $Q^\pi$ ) la fonction d'utilité  $V$  (respectivement  $Q$ ) associée à la politique  $\pi$ . Les fonctions d'utilité sont des fonctions qui permettent d'estimer l'espérance de gain de l'agent étant donné un état  $s$  ou un couple état-action  $(s,a)$ . Ces fonctions sont représentées sous forme de tableau comme la Q-table pour la fonction  $Q$ . La politique optimale est caractérisée par les équations suivantes :

$$V^*(s) = \max_{\pi} V^\pi(s), \quad \forall s \in \mathcal{S} \quad (1)$$

$$Q^*(s, a) = \max_{\pi} Q^\pi(s, a), \quad \forall s \in \mathcal{S}, \quad \forall a \in \mathcal{A} \quad (2)$$

Avec  $\mathcal{S}$ , l'ensemble des d'états et  $\mathcal{A}$ , l'ensemble des actions. Pour atteindre cette politique optimale, il existe plusieurs familles de méthodes en fonction de la connaissance que l'on a de l'environnement [1, 2] :

- **Méthodes de la programmation dynamique** : dans le cas où on a une bonne connaissance de l'environnement et donc on peut utiliser les informations locales pour améliorer la politique.

$$V^\pi(s) = \mathbb{E}_{\pi} \{ r_t + \gamma V^\pi(s_{t+1}) | s_t = s \} \quad (3)$$

où  $s_t$  et  $s_{t+1}$  sont respectivement les états de l'agent aux instants  $t$  et  $t + 1$  et  $r_t$  est la récompense obtenue par l'agent après avoir exécuté une action dans l'état  $s_t$ .  $\gamma$  est le *taux d'actualisation* qui permet de définir l'importance des récompenses futures.

- **Méthodes de Monte Carlo** : lorsque l'on ne connaît pas l'environnement, on échantillonne une zone de ce dernier par exploration, puis l'on estime les valeurs des fonctions d'utilité, et l'on se sert de ces valeurs pour faire une mise à jour de la politique.

$$V^\pi(s_t) = V^\pi(s_t) + \alpha [R_t - V^\pi(s_t)] \quad (4)$$

où  $R_t$  est la moyenne des gains sur la zone échantillonnée et  $\alpha$  le *taux d'apprentissage*.

- **Méthodes des différences temporelles** : Il s'agit d'un compromis entre les méthodes de Monte Carlo et celles de la programmation dynamique. On commence par échantillonner une zone de l'environnement comme pour les méthodes de Monte Carlo, mais on prend également en compte la corrélation entre les états pour faire les mises à jour comme pour les méthodes de la programmation dynamique.

$$V_t^\pi(s_t) = V_t^\pi(s_t) + \alpha [r_t + \gamma V_t^\pi(s_{t+1}) - V_t^\pi(s_t)] \quad (5)$$

où  $V_t^\pi(s_t)$  est la fonction d'utilité  $V$  associée à la politique  $\pi$  à l'instant  $t$ .

De ce qui précède, il va sans dire que les algorithmes les plus utilisés en apprentissage par renforcement sont de type différences temporelles. Nous pouvons citer entre autres [1, 2] : Q-learning (algorithme hors-politique, qui fait sa mise à jour sur la base de l'action qui maximise le gain de l'agent à l'état suivant), SARSA (algorithme sur-politique, qui fait sa mise à jour sur la base de la Q-table) et des algorithmes plus génériques type Deep Reinforcement Learning (algorithmes où la Q-table est approximée grâce à un réseau de neurones). Tous ces algorithmes, bien qu'ayant fait leurs preuves au fil du temps, présentent certaines faiblesses qui peuvent expliquer une chute d'efficacité des agents même bien entraînés dans certains environnements. Dans [10] Refael Vivanti et al. mettent en évidence deux problèmes des algorithmes d'apprentissage par renforcement : **Boring Area Trap** et **Manipulative Consultant**. Nous décrivons ces deux problèmes dans la section suivante.

### III PIÈGE DE LA ZONE ENNUYEUSE ET CONSULTANT MANIPULATEUR

Dans le contexte de l'apprentissage par renforcement, une zone ennuyeuse est une zone avec une faible variance et une faible espérance de gain sur les récompenses en comparaison aux autres zones. Une zone est considérée ici comme une action ou un couple état-action. Le **Boring Area Trap**, en français, **Piège de la Zone Ennuyeuse** est donc une conséquence de cette *différence de variance* entre les zones de l'environnement. Il s'agit plus précisément d'une situation dans laquelle, un agent se retrouve piégé, coincé dans une zone ennuyeuse ; une fois dans cette zone, due à sa faible variance, la politique (devenue ennuyeuse elle aussi) ne suggère alors que des actions qui vont conduire l'agent à y rester, d'où le terme *piège*.

Dans le cas des algorithmes de type Deep Reinforcement Learning, l'agent utilise un réseau de neurones pour prédire quelle est la meilleure action à exécuter étant donné un état. Ce réseau de neurones est donc considéré comme un consultant auprès de l'agent. On parle de **Manipulative Consultant**, en français, **Consultant Manipulateur** lorsqu'il y a conflit d'intérêt entre l'agent et son consultant. En effet, le réseau de neurones suggère des actions qui vont conduire l'agent à se retrouver piégé dans une zone ennuyeuse, mais qui le conduiront lui (le réseau de neurones) à améliorer la précision de ses prédictions.

Dans le cadre de ce papier, nous nous intéressons au problème du **Boring Area Trap**. Pour palier ce problème, bon nombre de solutions ont été proposées dans la littérature, parmi lesquelles :

- **Couper l'influence des récompenses sur la politique** [7, 8] : Cette solution a pour effet de diminuer la variance des zones intéressantes ce qui diminue la différence de variance entre les zones. Cependant, elle peut également réduire la variance des zones ennuyeuses, et dans ce cas le problème persiste.
- **L'utilisation d'un taux d'apprentissage faible** [10] : Supposons que nous soyons dans un environnement avec deux zones, une zone ennuyeuse ( $b$ ) et une zone intéressante ( $i$ ) et que l'algorithme d'apprentissage utilisé est Q-learning. Ainsi, la Q-table aura deux entrées  $Q_b$  et  $Q_i$  ; rappelons que l'agent se retrouve dans la zone ennuyeuse quand nous avons  $Q_b > Q_i$ . L'utilisation d'un taux d'apprentissage faible peut permettre de corriger le problème du Boring Area Trap comme évoqué dans [10], ceci s'explique par le

fait que la valeur de  $Q_b$  ne change pas drastiquement entre deux épisodes de même la comparaison entre  $Q_b$  et  $Q_i$  ne change pas drastiquement. Si l'exécution de l'algorithme commence avec cette inégalité fautive, pour qu'elle devienne vraie avec un  $\alpha$  très petit, il faut beaucoup de temps ; et pendant ce temps, il est possible que les récompenses échantillonnées permettent d'éviter la zone ennuyeuse. Cependant, cette solution conduit à une convergence extrêmement lente.

- **L'utilisation de l'exploration [10]** : Lorsque l'agent se retrouve piégé dans une zone ennuyeuse, la politique ne lui propose que des actions qui vont le conduire à y rester. Dans ce cas, l'exploration peut permettre à l'agent de sortir de cette zone car l'exploration est basée sur l'aléa et pas sur la politique devenue ennuyeuse elle aussi. Cependant, l'exploration change le comportement de l'agent et ne peut donc pas garantir que l'agent maximisera son gain sur le long terme ; il faut donc trouver une stratégie pour gérer le dilemme exploration/exploitation, de sorte que l'agent se comporte malgré tout de façon à maximiser son gain sur le long terme.

Refael Vivanti et al. ont proposé *Adaptive Symmetric Reward Noising (ASRN)* [10], qui consiste à ajouter du bruit aux récompenses de façon symétrique dans le temps et dans toutes les zones de l'environnement, et ce en fonction de leur variance. Ainsi, les zones de faibles variances verront leur variance augmenter mais les espérances de gain de ces zones resteront stables. Ce principe est mis en œuvre en deux phases comme il suit :

1. **Échantillonnage de l'environnement et détermination du bruit** : pendant les premières itérations de l'algorithme Q-learning :

- On estime la variance de chaque zone selon la formule suivante, inspirée de l'équation de mise à jour de l'algorithme Q-learning :

$$v_t = |Q_{new}(s_t, a_t) - Q_{old}(s_t, a_t)| \quad (6)$$

$$= |(1 - \alpha)Q(s_t, a_t) + \alpha(r_t + \gamma \max_a Q(s_{t+1}, a)) - Q(s_t, a_t)| \quad (7)$$

$$= |\alpha((r_t + \gamma \max_a Q(s_{t+1}, a)) - Q(s_t, a_t))| \quad (8)$$

- Pendant une période  $T$ , on échantillonne l'environnement. Il s'agit pour l'agent d'exécuter des actions de façon un peu aléatoire (car  $\varepsilon$  le taux d'exploration est élevé). Pendant cette période, on récupère chaque récompense  $r_t$  puis l'on calcule la variance  $v_t$  induite selon l'équation 8 et on stocke le couple  $(v_t, r_t)$  dans un tableau de taille  $T$ .
- Après cette période, on trie le tableau contenant les couples  $(v_t, r_t)$  selon l'ordre croissant des variances. On crée ensuite  $P$  paquets de taille  $N$  fixée chacun tels que  $p \in \{1, \dots, P\}$  avec :

$$T = \begin{cases} P \times N & \text{si } T \equiv 0 \pmod{P} \\ (P - 1) \times N + N' & \text{avec } N' < N \quad \text{Sinon.} \end{cases}$$

Sans nuire à la généralité, dans la suite on supposera que  $T \equiv 0 \pmod{P}$ . Dans chaque paquet, on range  $N$  récompenses consécutives (selon l'ordre croissant des variances associées). Parallèlement, on crée un tableau de taille  $P$  où on range la plus petite variance associée à chaque paquet de récompenses. À la fin de cette opération, on obtient  $P$  paquets de récompenses de taille de  $N$  chacun et pour chacun d'eux, on a une *unique* variance  $v_0$  qui lui est associée (la plus petite des variances associées au paquet).

- Une fois cette étape terminée, on calcule pour chaque paquet son écart-type selon la formule suivante :

$$S_p = STD_{1 \leq j \leq N}(r_j) \quad (9)$$

- On calcule ensuite l'amplitude du bruit à appliquer à chaque récompense en fonction de sa variance avec l'équation suivante :

$$N_p = \sqrt{S_{max}^2 - S_p^2} \quad \text{avec} \quad S_{max} = \max_p S_p \quad (10)$$

2. **Bruitage des récompenses** : tout au long des itérations suivantes de l'algorithme Q-learning, on détermine pour chaque récompense  $r_t$ , l'amplitude du bruit  $N_p$  à appliquer selon la variance  $v_t$  induite par la mise à jour. On échantillonne alors une récompense bruitée  $r'_t$  suivant la distribution :

$$r'_t \sim \mathcal{N}(r_t, N_p^2). \quad (11)$$

Il est à noter que plusieurs paquets peuvent faire référence à une même zone, il suffit pour cela que les variances  $v_0$  qui leurs sont associées soient proches ; ainsi les paquets de faibles variances feront référence à une zone ennuyeuse avec une intensité plus ou moins grande en fonction de leurs variances respectives et les paquets de fortes variances feront quant à eux référence à une zone d'intérêt également avec une intensité qui varie en fonction de leurs variances respectives.

À partir d'expérimentations menées sur le problème du bandit à 2 bras, Refael Vivanti et al. [10] montrent que l'algorithme ASRN résout le problème du Boring Area Trap. En reprenant les expérimentations menées par de Refael Vivanti et al. [10] sur de *longues durées de jeu*, on observe que le problème du Boring Area Trap persiste. Cette situation apparaît également dans le cas du bandit à *plusieurs bras (nombre de bras supérieur à 2)*, sur des durées de jeu moins longues que pour le cas du bandit à 2 bras. Ces deux caractéristiques à savoir : longues durées de jeu et environnement à zones ennuyeuses multiples, sont communes à plusieurs applications telles que : la conduite autonome [3, 11] et la peinture autonome [5].

Au vu des insuffisances que nous avons relevées dans la solution ASRN de Refael Vivanti et al. [10], nous proposons dans la section suivante des modifications de cet algorithme qui corrigent le problème du Boring Area Trap pour des applications ayant les caractéristiques sus-citées.

#### IV REDÉMARRAGE DE L'ALGORITHME ASRN

Le principe des algorithmes que nous proposons est de **redémarrer** à intervalle de temps régulier le processus d'échantillonnage de l'environnement et de détermination du bruit à appliquer

par zone. Dans le cas de l'algorithme ASRN ce processus est appliqué une seule fois, au début de l'algorithme. Au vu des contraintes posées par l'algorithme ASRN (à savoir que le taux d'exploration doit être élevé au début de l'algorithme), nous avons modifié les fonctions de variation des hyper-paramètres  $\varepsilon$  (taux d'exploration) et  $\alpha$  (taux d'apprentissage) de l'algorithme Q-learning. Nous proposons trois versions de l'algorithme modifié appelé **Rebooted ASRN (RASRN)**. Chacun de ces trois algorithmes RASRN est différent des autres par la manière dont varient les hyper-paramètres de l'algorithme Q-learning. Ces algorithmes sont :

- **Continuous  $\varepsilon$  decay RASRN** : Il s'agit d'une version de RASRN où l'algorithme ASRN est redémarré avec la valeur *courante* de  $\varepsilon$ . Contrairement à l'algorithme ASRN qui démarre avec une valeur de  $\varepsilon$  à 1 qui décroît dans le temps (l'ensemble des valeurs de  $\varepsilon$  forme une suite géométrique décroissante de premier terme 1 et de raison  $\varepsilon$ -decay), ici, seule la première exécution de la phase d'échantillonnage de l'environnement se déroule avec une valeur de  $\varepsilon$  *maximale*, les autres se déroulent avec la valeur *courante*. L'autre particularité de cette version est que l'hyper-paramètre  $\varepsilon$  ne décroît pas jusqu'à 0 mais jusqu'à une valeur minimale (0.01 dans le cas de nos expérimentations). Cet algorithme n'est pas fidèle à l'algorithme ASRN car pour bien fonctionner, ce dernier a besoin d'une valeur de  $\varepsilon$  maximale lors de la phase d'échantillonnage de l'environnement, cependant, il permet à l'agent d'exploiter sa Q-table 99% du temps environ, Q-table qui est le résultat des exécutions précédentes.
- **Full RASRN** : Contrairement à l'algorithme précédent, celui-ci consiste à faire remonter la valeur de l'hyper-paramètre  $\varepsilon$  à 1 au début de chaque itération du démarrage de l'algorithme ; ceci pour rester fidèle à l'algorithme ASRN, mais au risque de faire fi de la Q-table construite lors des exécutions précédentes, et ce pendant les premiers instants qui suivent le redémarrage.
- **Stepwise  $\alpha$  decay RASRN** : Dans cette version, il s'agit de trouver un compromis entre les deux précédentes, rester fidèle à l'algorithme ASRN en faisant remonter la valeur de  $\varepsilon$  à 1 après chaque redémarrage tout en ne faisant pas fi des Q-tables des itérations passées, ceci en faisant décroître la valeur de l'hyper-paramètre  $\alpha$  au moment de chaque redémarrage. Pour baisser cette valeur, nous avons utilisé un autre hyper-paramètre  $\alpha$ -decay.

Dans la section suivante, nous présentons les expérimentations réalisées avec ces algorithmes qui montrent bien qu'ils améliorent l'algorithme ASRN.

## V EXPÉRIMENTATIONS

Nos expérimentations ont porté sur le problème du bandit à k bras (aussi appelé problème du bandit manchot). C'est un problème mathématique qui se formule de manière imagée de la façon suivante : un utilisateur (un agent), face à des machines à sous, doit décider quelles machines jouer. Chaque machine donne une récompense moyenne que l'utilisateur ne connaît pas a priori. L'objectif est de maximiser le gain total (cumul des récompenses obtenues) de l'utilisateur [2]. Ce problème peut également être considéré comme étant celui d'un agent devant une machine unique ayant plusieurs bras et chaque bras ayant une espérance de gain inconnue par le joueur.



Nous avons conduit deux séries d'expérimentations. Le but de la première série a été de mettre en évidence le fait que malgré l'algorithme ASRN, le Boring Area Trap persiste lorsque la durée du jeu augmente, mais également de montrer comment notre solution RASRN permet de pallier ce problème. Dans la seconde série d'expérimentations, le but était de vérifier que lorsque le nombre de bras du bandit augmente, l'algorithme ASRN ne permet plus de corriger efficacement le problème du Boring Area Trap. Nous avons également montré que notre solution RASRN se comporte mieux que l'algorithme ASRN dans ces cas.

## 5.1 Protocole expérimental

Dans le cas de la première série d'expérimentations, nous avons opté comme Refael Vivanti et al. [10], pour la configuration suivante : le nombre de bras du bandit est  $k=2$  ; l'agent n'a que deux possibilités de jeu : le *bras gauche* ou le *bras droit*. Du point de vue de l'algorithme Q-learning, il s'agit d'un environnement à un seul état et à deux actions *gauche* et *droit*. L'équation de mise-à-jour devient donc :

$$Q_{new}(a_t) \leftarrow (1 - \alpha)Q_{old}(a_t) + \alpha(r_t + \gamma \max_a Q(a)) \quad (12)$$

La Q-table a ainsi deux entrées  $Q_l$  et  $Q_r$  qui sont respectivement les espérances des gains cumulés associées aux bras gauche et droit. Sachant que le gain cumulé est  $R_t = \sum_{k=0}^{\infty} \gamma^k r_{t+k+1}$ , on va initialiser ces entrées avec les valeurs optimales :

$$Q_l = \sum_{t=0}^{\infty} \gamma^t \mu_l = \frac{\mu_l}{1 - \gamma} \quad \text{et} \quad Q_r = \sum_{t=0}^{\infty} \gamma^t \mu_r = \frac{\mu_r}{1 - \gamma} \quad (13)$$

où  $\mu_l, \mu_r$  sont les moyennes des récompenses associées respectivement au bras gauche et droit. On initialise les bras avec des valeurs optimales pour simuler un agent qui aurait parfaitement appris une politique optimale, mais qui, due à la différence de variance entre les récompenses des deux bras va régresser alors qu'il suit scrupuleusement la politique apprise.

Pour mener à bien nos expérimentations, nous avons initialisé chaque agent avec les valeurs suivantes :  $\mu_l = 0, \sigma_l = 0.5, \mu_r = 1$  et  $\sigma_r = 7$ . Les hyper-paramètres de l'algorithme Q-learning sont :  $\alpha = 0.1, \gamma = 0.9, \varepsilon = 1, \varepsilon\text{-min} = 0.01$  (dans le cas de l'algorithme *Continuous  $\varepsilon$  decay RASRN*),  $\varepsilon\text{-decay} = 0.999, \alpha\text{-decay} = 0.5$  (dans le cas de l'algorithme *Stepwise  $\alpha$  decay RASRN*). Comme on peut le remarquer en regardant les valeurs des écarts-types des bras gauche et droit, il y a une différence de variance de valeur 49, ce qui va forcément induire un problème de Boring Area Trap, raison pour laquelle nous avons choisi cet exemple extrême ; il permet de mettre en évidence le problème.

Nous avons monitoré la Q-table pour déterminer à quel moment l'agent entre dans la zone ennuyeuse, le temps qu'il y met, s'il parvient à en sortir et quand cela se produit. Les algorithmes en compétition ici sont l'algorithme ASRN (celui de Refael Vivanti et al. [10]) et les versions de RASRN que nous proposons. Le nombre de redémarrages des algorithmes RASRN a été fixé à cinq. Nous avons pour chacun de ces algorithmes instancié 10 agents qui ont résolu le problème du bandit à 2 bras dans les conditions énoncées ci-dessus. Les durées de jeu que nous avons retenues pour nos expérimentations sont : 40 000, 50 000, 60 000, 400 000, 500 000, 600 000, 1 000 000. Les métriques utilisées dans le cadre de cette série d'expérimentations sont :

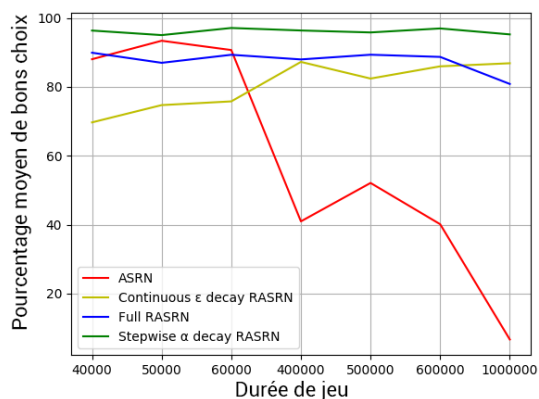


FIGURE 1 – Pourcentage moyen de bons choix :  $k=2$ .

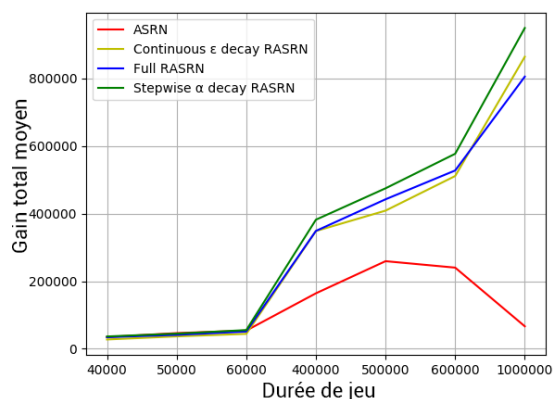


FIGURE 2 – Gain total moyen :  $k=2$ .

- **Pourcentage moyen de bons choix pour les 10 agents :** Au vu des conditions initiales de nos expérimentations, le bon choix correspond au choix du bras *droit*. Comme son nom l'indique, il s'agira de calculer le pourcentage de bons choix pour chacun des 10 agents et d'en faire la moyenne. Cette métrique a été retenue car elle nous renseigne sur la qualité du choix opéré par l'agent, et par conséquent s'il est victime d'un Boring Area Trap.
- **Gain total moyen pour les 10 agents :** Cette métrique nous a semblé importante car le but final de l'agent est de maximiser son gain total.

Dans le cadre de notre seconde série d'expérimentations, nous avons travaillé avec des environnements qui correspondent aux versions du bandit à  $k$  bras où  $k$  vaut 2, 3, 4 et 5. L'objectif de cette série d'expérimentations a été de montrer que lorsque le nombre de zones, notamment de zones ennuyeuses croît la solution ASRN de Refael Vivanti et al. [10] voit ses résultats chuter, notamment en terme de pourcentage de bons choix. Nous avons également vérifié que grâce à nos différentes versions de RASRN, nous apportons une solution à ce problème. Les variantes du bandit à  $k$  bras utilisées ici comme environnement à chaque fois, ont été initialisées de la façon suivante : lorsqu'il s'agit d'une zone ennuyeuse ( $b$ ),  $\mu_b = 0$ ,  $\sigma_b = 0.5$ , et lorsqu'il s'agit de la zone intéressante ( $i$ ) (il y en a une à chaque fois, peu importe le nombre de bras),  $\mu_i = 1$ ,  $\sigma_i = 7$ . Dans ces conditions, le choix du bras intéressant est considéré comme étant le bon, et tout autre choix est considéré comme étant mauvais,  $b \in \{1, \dots, k\} \setminus \{i\}$ . La métrique utilisée dans ces expérimentations est le *pourcentage moyen de bons choix* effectué par les agents. Les hyper-paramètres de l'algorithme Q-learning sont les mêmes que dans la première série d'expérimentations. Les durées de jeu retenues pour cette série d'expérimentations sont : 500 000 et 600 000 épisodes.

## 5.2 Résultats et discussions

Dans le cadre de la première série d'expérimentations, la Figure 1, montre que lorsque la durée du jeu croît, le pourcentage moyen de bons choix effectué par les agents a une tendance décroissante lorsque ces derniers utilisent l'algorithme ASRN de Refael Vivanti et al. [10]. On

TABLE 1 – Pourcentage moyen de bons choix :  $k=2,3,4$  et 5.

Algorithmes	Durées de jeu	500000				600000			
		k=2	k=3	k=4	k=5	k=2	k=3	k=4	k=5
ASRN		52.09	40.95	26.43	14.29	40.12	13.08	14.46	14.50
Continuous $\varepsilon$ decay RASRN		82.39	76.47	67.35	63.61	85.93	74.56	71.24	56.64
Full RASRN		89.31	85.26	81.58	79.18	88.67	81.83	78.73	71.54
Stepwise $\alpha$ decay RASRN		95.77	96.00	89.83	90.70	96.93	93.17	90.70	81.57

observe également sur cette figure que lorsque les agents utilisent les algorithmes RASRN, malgré l'augmentation de la durée du jeu, le pourcentage moyen de bons choix reste assez élevé, avec une prédominance de l'algorithme Stepwise  $\alpha$  decay RASRN sur les autres. En effet, la courbe de l'algorithme Stepwise  $\alpha$  decay RASRN est au-dessus de toutes les autres pour les durées de jeu considérées sur la figure ; en particulier, très au-dessus de la courbe de l'algorithme ASRN lorsque la durée de jeu est supérieure ou égale à 400 000 épisodes. La Figure 2 montre que lorsque la durée du jeu augmente, le gain total moyen des agents a tendance à décroître au bout d'un certain temps (lorsque le jeu à une durée supérieure ou égale à 500 000 épisodes) pour le cas de l'algorithme ASRN ; par contre, sa tendance est toujours croissante avec les algorithmes RASRN, avec cette fois encore une prédominance de l'algorithme Stepwise  $\alpha$  decay RASRN sur les autres. On note également sur cette figure que pour les durées de jeu entre 40 000 et 60 000 épisodes, les courbes des quatre algorithmes sont quasiment au même niveau, mais lorsque la durée de jeu est supérieure ou égale à 500 000 épisodes, les courbes RASRN ont une tendance croissante *exponentielle*, tandis que la courbe de ASRN à une tendance décroissante. Les Figures 3, 4, 5 et 6 quant à elles montrent l'évolution des Q-values  $Q_l$  et  $Q_r$  d'un agent pour les exécutions du bandit à 2 bras sur un jeu d'une durée de 1 000 000 d'épisodes. On observe clairement sur la Figure 3 le problème du Boring Area Trap. En effet, sur cette figure on peut observer à partir de l'épisode 100 000 que la courbe rouge reste au-dessus de la courbe verte, et ce de façon permanente jusqu'à la fin du jeu ; ceci traduit le fait que l'agent choisit systématiquement le bras gauche (ennuyeux) au détriment du bras droit (intéressant) : il est piégé dans la zone ennuyeuse. Les Figures 4, 5 et 6 quant à elles montrent clairement une absence du phénomène du Boring Area Trap observé sur la Figure 3. En effet, sur ces figures on n'observe à aucun moment la courbe rouge au-dessus de la courbe verte de façon permanente ; on observe d'ailleurs le contraire à certains moments sur les Figures 5 et 6.

Dans le cadre de notre deuxième série d'expérimentations, le Tableau 1 présente les résultats obtenus. La ligne ASRN du tableau indique clairement que le pourcentage moyen de bons choix baisse avec l'augmentation du nombre de bras du bandit, mais également avec l'augmentation de la durée du jeu. Ce qui valide expérimentalement l'hypothèse que nous avons formulée selon laquelle l'algorithme ASRN a du mal à corriger efficacement le Boring Area Trap lorsque le nombre de zones ennuyeuses de l'environnement augmente. En observant les résultats obtenus par les autres algorithmes : Continuous  $\varepsilon$  decay RASRN, Full RASRN et Stepwise  $\alpha$  decay RASRN, le même phénomène de baisse est observé. Cependant, les algorithmes que nous proposons arrivent à maintenir à un seuil assez élevé ( $> 55\%$  pour les cas illustrés dans le tableau) le pourcentage moyen de bons choix. On note également que l'impact de nos algorithmes RASRN dans la résolution du Boring Area Trap est d'autant plus important que la durée du jeu est grande et que le nombre de bras est grand ; comme indiqué dans le Tableau 1.

## VI CONCLUSION

Le Boring Area Trap ou Piège de la Zone Ennuyeuse est un problème dont souffrent les algorithmes d'apprentissage par renforcement dans des environnements où il existe une différence de variances sur les récompenses entre les zones de ces environnements ; cette différence de variance conduit malheureusement les agents à préférer les zones de faibles variances, mais également de faibles espérances de gain. Pour pallier ce problème, Refael Vivanti et al. [10] ont proposé l'algorithme Adaptive Symmetric Reward Noising (ASRN) qui consiste à ajouter du bruit aux récompenses de façon symétrique dans le temps et dans toutes les zones de l'environnement, et ce en fonction de leurs variances. Dans ce papier, nous avons dans un premier temps montré expérimentalement que malgré cet algorithme, le problème persiste dans certaines situations (longues durées de jeu et environnement à zones ennuyeuses multiples). Dans un second temps, nous avons proposé trois variantes de l'algorithme ASRN qui améliorent le pourcentage moyen de bons choix des agents, ainsi que leur gain total moyen. Le principe de ces algorithmes nommées RASRN est le suivant : redémarrer à intervalle de temps régulier le processus de détermination du bruit à appliquer par zone. Ils sont différents les uns des autres par la manière dont sont mis-à-jour les hyper-paramètres de l'algorithme Q-learning. Il s'agit de Continuous  $\epsilon$  decay RASRN, Full RASRN et Stepwise  $\alpha$  decay RASRN. Grâce aux expérimentations menées sur le problème du bandit à k bras, nous avons pu mettre en évidence le fait que ces algorithmes améliorent l'algorithme ASRN. Les différences de performances entre les algorithmes RASRN et l'algorithme ASRN croissent avec l'augmentation des durées des jeux et aussi du nombre de zones ennuyeuses. En plus d'apporter une justification théorique aux résultats obtenus, nous nous proposons dans la suite de ce travail d'étendre nos algorithmes RASRN aux algorithmes SARSA et Deep Q-network, ainsi qu'à d'autres applications.

## RÉFÉRENCES

- [1] R. BELLMAN, «A Markov decision process», *Journal of Mathematics and Mechanics*, Vol. 6, No. 5 (1957), pp. 679-684.
- [2] A. CORNUÉJOLS, L. MICLET, Y. KODRATOFF, «Apprentissage artificiel, Concepts et algorithmes», *Éditions Eyrolles*, 2003.
- [3] A. EL SALLAB, M. ABDU, E. PEROT, S. YOGAMANI, «End-to-end deep reinforcement learning for lane keeping assist», *arxiv.org*, arXiv preprint arXiv :1612.04340, 2016.
- [4] M. HESSEL, J. MODAYIL, H. VAN HASSELT, T. SCHAUL, G. OSTROVSKI, W. DABNEY, D. HORGAN, B. PIOT, M. AZAR, D. SILVER, «Rainbow : Combining Improvements in Deep Reinforcement Learning», *Thirty-Second AAAI Conference on Artificial Intelligence*, 2018.
- [5] B. JIA, C. FANG, J. BRANDT, B. KIM, D. MANOCHA., «Paintbot : A reinforcement learning approach for natural media painting», *arxiv.org*, arXiv preprint arXiv :1904.02201, 2019.
- [6] V. MNIH, K. KAVUKCUOGLU, D. SILVER, A. A. RUS, «Human-level control through deep reinforcement learning », *Nature*, Vol. 518, No. 7540, 2015.
- [7] J. SCHULMAN, S. LEVINE, P. ABBEEL, M. JORDAN, P. MORITZ, «Trust region policy optimization », *In International Conference on Machine Learning*, (pp. 1889-1897), 2015.

- [8] J. SCHULMAN, F. WOLSKI , P. DHARIWAL, P. RADFORD, O. KLIMOV, «Proximal policy optimization algorithms», *arxiv.org*, arXiv preprint arXiv :1707.063477, 2017.
- [9] D. SILVER ET AL., «Mastering the game of Go without human knowledge», *Nature*, Vol. 550, No. 24270 , 2017.
- [10] R. VIVANTI, T. D. SOHLBERG-BARIS, S. COHEN, O. COHEN, «Adaptive Symmetric Reward Noising for Reinforcement Learning», *arxiv.org*, arxiv :1901.10144v1, 2019.
- [11] C. YOUNG, JIANBO LU, D. FILEV, P. TSIOTRAS, «Advanced planning for autonomous vehicles using reinforcement learning and deep inverse reinforcement learning», *Robotics and Autonomous Systems*, Vol. 14, p. 1-18, April 2019.

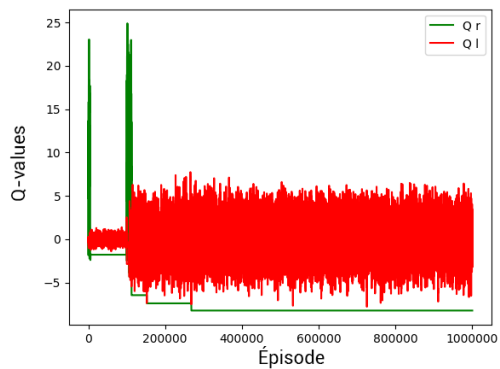


FIGURE 3 – Q-values : ASRN,  $k=2$ .

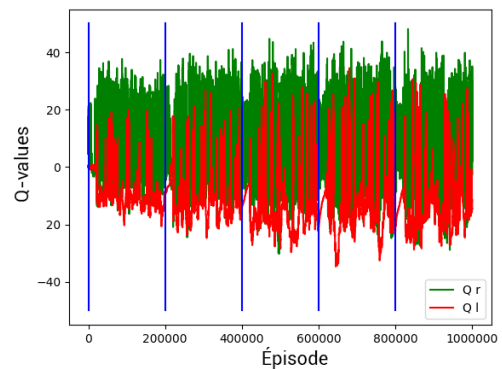


FIGURE 4 – Q-values : Continuous  $\epsilon$  decay RASRN,  $k=2$ .

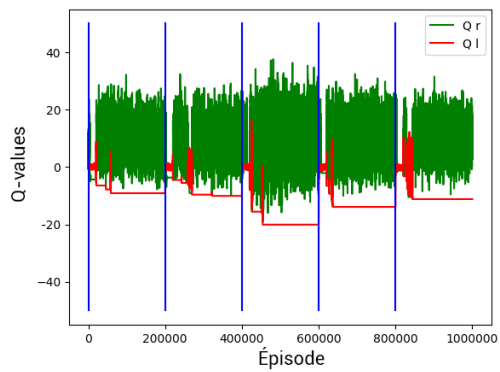


FIGURE 5 – Q-values : Full RASRN,  $k=2$ .

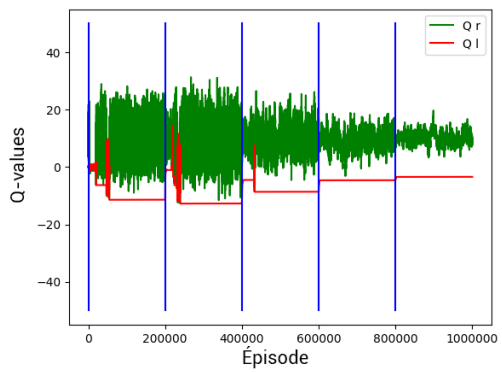


FIGURE 6 – Q-values : Stepwise  $\alpha$  decay RASRN,  $k=2$ .