

Projet Exploration/Exploitation

Ryan

Pierre Fumeron

Table des matières

1 Bandits-Manchots

2

Introduction

Dans ce projet, nous nous sommes intéressés au dilemme de l'exploration vs exploitation. Prenons l'exemple du bandit manchot pour illustrer ce problème : Imaginons plusieurs machines à sous, chaque machine a un gain et une probabilité de toucher ce gain qui lui sont propres. Un joueur, qui veut maximiser son profit, a donc plusieurs options : jouer constamment la même machine sans acquérir d'informations sur les autres et donc sans savoir s'il y a une meilleure solution. Ou bien, jouer toutes les machines pour en savoir plus sur leurs probabilités mais en gagnant moins d'argent que si on avait uniquement joué la "meilleure" machine (celle qui offre le meilleur profit). Ce problème trouve de nombreuses applications en intelligence artificielle mais également dans la vie quotidienne, par exemple : Vous êtes dans un restaurant, commandez-vous un plat que vous appréciez ou bien un plat que vous n'avez jamais goûté en espérant qu'il soit meilleur que celui que vous connaissez ?

Nous avons donc étudié des algorithmes du dilemme exploration vs exploitation avec 2 jeux différents. Premièrement, nous nous sommes intéressés à l'exemple du bandit-manchot décrit ci-dessus, avec comme objectif de maximiser le gain. Il faut pour cela trouver un bon compromis entre l'exploration des différentes machines puis l'exploitation de la machine la plus rentable.

Le deuxième jeu est celui du morpion où nous avons fait s'affronter différents algorithmes : aléatoire, UCT et Monte-Carlo).

1 Bandits-Manchots

Supposons une machine à sous à N leviers dénotés par l'ensemble $1, 2, \dots, N$. Chacun de ses leviers est une action possible parmi lesquelles le joueur doit choisir à chaque pas de temps : l'action choisie à l'instant t sera appelée a_t (un entier entre 1 et N). Nous supposons que la récompense associée à chaque levier i suit une distribution de Bernoulli de paramètre μ_i : avec une probabilité μ_i le joueur obtient une récompense de 1, avec une probabilité $1 - \mu_i$ il obtient 0. Les μ_i sont constants tout au long de la partie. Pour le joueur, le gain à la fin de T parties est la somme $G_T = \sum_{t=1}^T r_t$ (n.b. la récompense est aléatoire, le gain G_T est une variable aléatoire, tout comme r_t). Soit $i_t = \arg \max_i 1, \dots, N$ l'indice du levier choisi à l'instant t : $i_t = \arg \max_i \mu_i$. Si le joueur joue un autre levier que i_t , il obtient $r_t = 0$, avec μ_{i_t} la récompense aléatoire tirée de la distribution de Bernoulli de paramètre μ_{i_t} . On appelle regret $L_T = G_T - \max_i \sum_{t=1}^T r_t = G_T - \sum_{t=1}^T \mu_{i_t}$. L'objectif est donc de minimiser le regret.

Question 1