

① バンディット問題

bandit : スロットマシン

レバーを引くセラングに絵柄が変わる

マワりの並べ方によってコインがもらえる (0~n) 枚

99 月範 バンディット問題

説定

1本レバー・スロットマシンが複数台ある

スロットマシン毎に絵柄の出方が異なる。

プレイヤーはスロットマシンの情報をも知れない

(例えば) 1000回プレイして得るコインの枚数を最大化したい

用語

環境 Environment

スロットマシン

エージェント Agent

プレイヤー

行動 Action

1台を選んでプレイする

報酬 Reward

スロットマシンから出るコイン

スロットマシンの良さとは?

スロットマシン毎にもらえるコインの確率分布が異なる。

SMA

coin	0	1	5	10
p	0.7	0.5	0.12	0.03

SMA

coin	0	1	5	10
p	0.5	0.4	0.09	0.01

常に期待値の高いマシンを選ぶようプレイすればいい

価値 Value

報酬の期待値

行動価値 Action Value

行動に対して得られる報酬の期待値

$R_t \in \{0, 1, 5, 10\}$

: t回目には得られる報酬 Reward

$A_t \in \{a, b\}$

: t回目の行動 Action, a, b はスロットの種類

$E[R]$	R の期待値
$E[R A]$	A という行動をした場合の R の期待値
$E[R A=a] = E[R a]$	
$Q(A) = E[R A]$	行動価値 (Quality) 行動 A の価値
$Q(A)$	真の値 理論値 ← エージェントは知らない
$\hat{Q}(A)$	推定値

アルゴリズム

エージェントはスロットマシンの価値(確率分布や期待値)を知らない

→ 各スロットマシンの価値を推定する必要ある。

推定方法

SM	1回目	SMの価値	SM	1回目	2回目	3回目	SMの価値
1	0	0	1	0	1	5	2
2	1	1	2	1	0	0	0.33

$$Q(A=1) = 2, \quad Q(A=2) = 0.33$$

この期待値 $Q(A)$ は推定値ではなく SM_n の方が価値が高い

標本平均

大数の法則: 無限回のサンプリングで標本平均は真の値に一致する

実装

R_1, R_2, \dots, R_n から Q_n を求める。

単に平均すると $n+1$ 回目の行動で R_{n+1} を得たとき $R_1 \sim R_{n+1}$ を再度必要としてしまう。

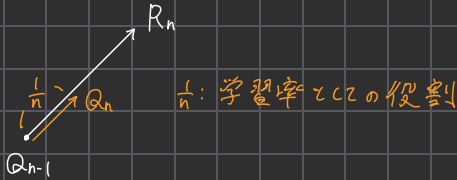
→ 漸化的に求める

$$Q_{n-1} = \frac{R_1 + \dots + R_{n-1}}{n-1} \quad \left| \quad \begin{aligned} Q_n &= \frac{1}{n} (R_1 + \dots + R_{n-1} + R_n) \\ &= \frac{1}{n} ((n-1)Q_{n-1} + R_n) \\ R_1 + \dots + R_{n-1} &= (n-1)Q_{n-1} \\ &= \left(1 - \frac{1}{n}\right) Q_{n-1} + \frac{1}{n} R_n \\ &= Q_{n-1} + \frac{1}{n} (R_n - Q_{n-1}) \end{aligned}$$

$$Q_n = Q_{n-1} + \frac{1}{n} (R_n - Q_{n-1})$$

α も γ も時間も小さい

$$Q_n = Q_{n-1} + \frac{1}{n} (R_n - Q_{n-1})$$



プレイヤーの戦略

- greedy (貪欲法) 各SMの価値の推定値が最大のものを常に選ぶ。
実験が少なく推定値の大小と真の値の大小が一致していると決めかねるのを「不確かさ」
- SMの価値を精度よく推定するため 様々なSMを試す。

活用Exploitation 経験から最適な行動をする (greedy)

← 真の最適を見逃しているかも

探索Exploration greedyでない行動を試みず。

トレードオフ

強化学習: 活用と探索のバランスをいかに取るか

- ϵ -greedy 法

$\epsilon = 0.1$ (例) の確率で 探索, $1-\epsilon$ 以外で活用
ランダムな行動を選ぶ

定常問題: 報酬の確率分布が定常

スロットマシンの勝率は定常だと。

非定常問題: 毎プレイで勝率が変化するなど

Q_n : 価値, R_n : 報酬

$$Q_n = \frac{R_1 + \dots + R_n}{n} = \frac{1}{n} R_1 + \frac{1}{n} R_2 + \dots + \frac{1}{n} R_n$$

← 重み

定常

$$Q_n = Q_{n+1} + \frac{1}{n} (R_n - Q_{n-1})$$

非定常

$$Q_n = Q_{n+1} + \alpha (R_n - Q_{n-1})$$

$$Q_n = \alpha R_n + \alpha(1-\alpha)R_{n-1} + \alpha(1-\alpha)^2 R_{n-2} + \dots + \alpha(1-\alpha)^{n-1} R_1 + (1-\alpha)^n Q_n$$

$R_1 \sim R_n$ の 指数 (加重) 移動平均