

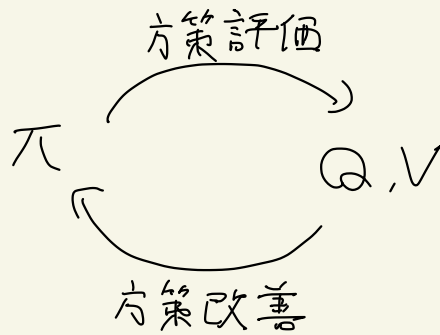
第2回 強化学習ゼミ (8/12)

ましろ



復習

強化学習とは..



の繰り返し.

今からやること.

- π をどのように定め,
 Q や V を どう改善 するのかを見ていく.
- ヘルマン方程式は 状態遷移確率 がわかっている
状況なら 解けるが; そうでないときは 解けない.
(あと計算量が爆発する)

- 状態遷移確率がわかっていない解法 Bellmanの
 → モデルベース ← 解ける

： わかっていない

→ モデルフリー

→ モデルフリーのアプローチを見ていく。

3.6 モンテカルロ法

- アイデア

たくさん言式行錯誤を繰り返す。

→ 状態の遷移は近似的に状態遷移確率に従う(大数の法則)

→ 表反面州の平均は基期待値に収束。

何度も繰返して得られた経路をもとに
 状態価値や状態行動価値を推定する方法を

モンテカルロ法 という。

1. やりここの大枠

1人下を M 回やり. ($m=1 \dots M$)

① 状態列 s_0, s_1, \dots, s_{T_m} を π によって生成する.

② 状態列に対する報酬列

r_0, r_1, \dots, r_{T_m} を観測する,

③ 報酬列に対して状態価値を推定

④ 状態価値を更新.

③をどうやるか?

→ {

- ・ 初回訪問モンテカルロ法.
- ・ 逐一訪問モンテカルロ法

初回言方問 モニテカルロ法

状態列 $s_1 \dots s_m$ について.

状態 s_t が初回の言方問なら、収益 $R(s_t)$ を

$$R(s_t) \leftarrow R(s_t) + \sum_{t'=t+1}^T r_{t'}$$

状態価値 $V(s_t)$ を

$$V(s_t) \leftarrow \frac{1}{m} R(s_t)$$

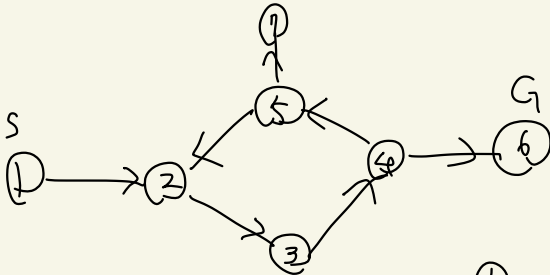
とす.

③の部分.

* 初回言方問という条件を入れているのは.

観測された状態列の中に同じ状態が複数回

あらわれて、各状態が不平等な言評価になるのをさけるため.



① → ② → ③ → ④ → ⑤ → ②

・ 3系 - モンテカルロ法.

初回 言方問 モンテカルロ法では状態価値を評価した,

モデルが不明なので, とるべき行動がわからない,

→ 状態行動価値を評価する

状態行動価値最大の行動を
すればよい

・ 問題点

状態行動価値 $Q(s, a)$ にて,

(s, a) が 観測されない $Q(s, a)$ が 評価できない,

例)

$a = \{0, 1, 2, 3, 4\}$ とする.

$\pi(a|s) = \frac{a}{10}$ のとき,

$a=0$ が 観測される確率は 0.

問題の解決へ

① すべての行動の確率が0でないと仮定する。

$$\left(\begin{array}{l} \cdot R(s_t) \leftarrow R(s_t) + \sum_{t'=t+1}^{T-1} r_{t'} \\ \cdot Q(s_t, a_t) \leftarrow \frac{1}{m} R(s_t) \end{array} \right.$$

とすればOK

「モンテカルロ - ES 法」 という。

(2) 方策オン型モンテカルロ制御

確率的な方策を用いることで、任意の状態 行動対ル (s, t)

について $\pi(s, a) > 0$ を保証する、

↖ 遺伝的アルゴリズムみたい。

1. ϵ -greedy 法.

ϵ は十分に小さな正の値とする ($0 < \epsilon < 1$)

$$\pi(a_t | s_t) := \begin{cases} \frac{1-\epsilon}{|B(s_t)|} + \frac{\epsilon}{|A(s_t)|} & (a_t = \arg \max_{a \in A(s_t)} Q(s_t, a)) \\ \frac{\epsilon}{|A(s_t)|} & \text{otherwise.} \end{cases}$$

ただし、 $A(s_t)$: 時刻 s_t のときの行動空間

$$B(s_t) = \{ a \in A(s_t) \mid a = \arg \max_{b \in A(s_t)} Q(s_t, b) \}$$

※ 教科書だと $\sum \pi(a_t | s_t) = 1$ にならない、

($\epsilon = \frac{1}{2^n}$ とかすると良さそう?)

softmax 法

温度パラメータ $\tau \in \mathbb{R}_{>0}$ を用いて,

$$\pi(a_e | s_e) := \frac{\exp\left(\frac{Q(s_e, a)}{\tau}\right)}{\sum_{b \in A} \exp\left(\frac{Q(s_e, b)}{\tau}\right)}$$

τ は学習率をコントロールするパラメータで,

$$\lim_{\tau \rightarrow 0} \pi_{\text{greedy}}(a_e | s_e) = \lim_{\tau \rightarrow 0} \pi_{\text{softmax}}(a_e | s_e)$$

(証明)

文字列 w の i 番目の文字 b_i を用いて, $f_k: \mathbb{R}_{>0} \rightarrow \mathbb{R}_{>0}$ とし,

$$f_k(x) = \frac{\exp\left(\frac{b_k}{x}\right)}{\sum_{b_i \in A} \exp\left(\frac{b_i}{x}\right)}$$

が $x \rightarrow 0$ 付近での挙動を説明する.

$$f_k(x) = \frac{1}{\sum_{b_i \in A} \exp\left(\frac{b_i - b_k}{x}\right)} \left(\text{分母は } \exp\left(\frac{b_k}{x}\right) \text{ である} \right)$$

である

$$\sum_{b_i \in A} \exp\left(\frac{b_i - b_k}{x}\right) \quad (x > 0),$$

i) $b_i - b_k > 0$ のとき, ($b_k \neq \max(A)$ のとき)

$x > 0$ のとき,

$$\frac{b_i - b_k}{x} > 0.$$

$$\lim_{x \rightarrow 0} \frac{b_i - b_k}{x} = \infty$$

$$\lim_{x \rightarrow 0} \exp\left(\frac{b_i - b_k}{x}\right) = \infty$$

$$\lim_{x \rightarrow 0} f(x) = 0 \quad \left(\forall i, \exp\left(\frac{b_i - b_k}{x}\right) > 0 \text{ のとき} \right)$$

ii) $b_i - b_k < 0$ のとき,

$$\frac{b_i - b_k}{x} < 0.$$

$$\lim_{x \rightarrow 0} \frac{b_i - b_k}{x} = -\infty$$

$$\lim_{x \rightarrow 0} \exp\left(\frac{b_i - b_k}{x}\right) = 0.$$

iii) $b_i = b_k$ のとき,

$$\frac{b_i - b_k}{x} = 0.$$

$$\lim_{x \rightarrow 0} \exp\left(\frac{b_i - b_k}{x}\right) = 1.$$

i), ii), iii) 成立する。

- $b_k = \max(A)$ のとき

$$B = \{b \in A \mid b = \max(A)\}$$

より

$$\lim_{x \rightarrow 0} f(x) = \frac{1}{|B|}$$

- $b_k < \max(A)$ のとき

$$\lim_{x \rightarrow 0} f(x) = 0, \quad \square$$

③ 方策オプティミカル制御.

言平価・改善される方策 π (推定方策) と,

観察の行動する方策 π' (行動方策)

という2つの方策を用いる

→ 推定方策と行動方策では収益が異なる.

・ 改善した方だけ収益が良くなる.

→ 改善した際の方策ごとの差をなくすことを考える.

→ 重要度 (重み) を定義する.

推定方策 π において 状態行動対列 $(s_0, a_0) \dots (s_T, a_T)$ と 収益 R が観測 される 確率 ϵP ,

探索方策 π' において, 同じ状態行動対列 と 収益 R が観測 される 確率 $\epsilon P'$ となる,

この時, 確率的に 推定方策 では 探索方策 の $\frac{P}{P'}$ 倍 多く 観測 している,

- $\frac{P}{P'} = 1$ のとき 推定方策 と 探索方策 は 同一 (方策オン)
- $\frac{P}{P'} > 1$ 推定方策 では 多く 見ている
→ 強い意味をもつ (重要度が高い)
- $\frac{P}{P'} < 1$ 重要度 低め,
 $\frac{P}{P'}$ 倍 の 重み を 与える ($\frac{P}{P'}$ 回 対比 して みる)

Algo

$$R_{sum} \leftarrow R_{sum} + \frac{P}{P'} R$$
$$k \leftarrow k + \frac{P}{P'}$$
$$Q^\pi(s, a) = \frac{1}{k} R_{sum}$$

(s_t, a_t) は観測列 $(s_{t+1}, a_{t+1}), (s_{t+2}, a_{t+2}) \dots$
 は観測列の石置率 $p^\pi(s_t, a_t)$ の以下定義が成り立つ。

$$p^\pi(s_t, a_t) = P(s_{t+1} | s_t, a_t) \prod_{k=t+1}^{\infty} \pi(a_k | s_k) P(s_{k+1} | s_k, a_k)$$

重み $w(s_t, a_t)$ は

$$w(s_t, a_t) := \frac{p^\pi(s_t, a_t)}{p^{\pi'}(s_t, a_t)}$$

$$= \frac{\cancel{P(s_{t+1} | s_t, a_t)} \prod \pi(a_k, s_k) \cancel{P(s_{k+1} | s_k, a_k)}}{\cancel{P(s_{t+1} | s_t, a_t)} \prod \pi'(a_k, s_k) \cancel{P(s_{k+1} | s_k, a_k)}}$$

$$= \prod_{k=t+1}^{\infty} \frac{\pi(a_k | s_k)}{\pi'(a_k | s_k)} \quad (\text{状態序列の確率に依存しない})$$

推定が最も決定論的である場合

$$w(s_t, a_t) = \begin{cases} \prod_{k=t+1}^{\infty} \frac{1}{\pi'(s_k, a_k)} & a_k = \arg \max_{a \in A} Q(s_k, a) \\ 0 & \text{Other wise} \end{cases}$$

→ 最適行動以外は学習の意味がない。

3.7 TD学習

- モンテカルロ法は試行が終了なまでに更新できない
という問題点

→ この欠点を回避したのがTD学習

「1-step TD法」, 「k-step TD法」, 「TD(λ)法」
があって 今日は 1-step と k-step を見る。

1. 1-step TD法

名のとおし 1ステップごとに価値を改善する。

ベルマン方程式の解のことは。

$$\sum_{a \in A} \sum_{s_{t+1} \in S} \pi(a|s_t) p(s_{t+1}|s_t, a) = 1 \text{ より}$$

$$V(s_t) = r_t + \gamma V(s_{t+1}) \text{ である。}$$

$Q(s, a)$ も同様に (代入して確かめてみる) である。

$$Q(s, a) = r_t + \gamma Q(s_{t+1}, a_{t+1}) \text{ 解のことは。}$$

$V(S_t)$ と $Q(S_t, a_t)$ を最適解に近づけていくことを考える,

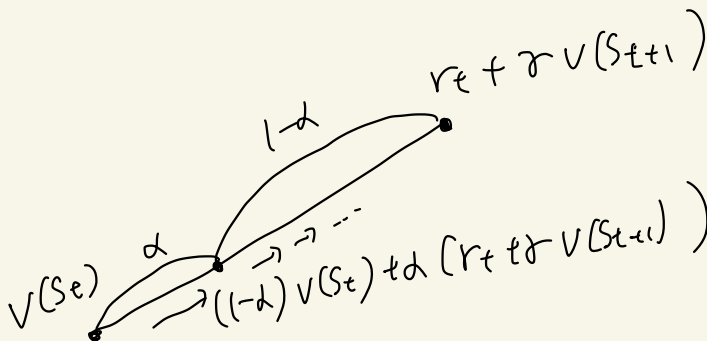
$$\begin{cases} V(S_t) \leftarrow V(S_t) + \alpha (r_t + \gamma V(S_{t+1}) - V(S_t)) \\ Q(S_t, a_t) \leftarrow Q(S_t, a_t) + \alpha (r_t + \gamma Q(S_{t+1}, a_{t+1}) - Q(S_t, a_t)) \end{cases}$$

～お気持です～

右辺を式変形すると、

$$(1-\alpha) V(S_t) + \alpha \underbrace{(r_t + \gamma V(S_{t+1}))}_{\text{解}}$$

図で書いてみる。



2. k-Step TD法

1-Step TD法だと局所解に123

(例は教科書を参考)

可能性があるので、1回先まで先読みすることと考える。

ベルマン方程式の解は

$$V(S_t) = r_t + \gamma V(S_{t+1})$$

であった。これをもう少し展開して、

$$V(S_t) = r_t + \gamma (r_{t+1} + \gamma V(S_{t+2}))$$

⋮

$$= \sum_{n=0}^{k-1} \left(\gamma^n r_{t+n} + \gamma^k V(S_{t+k}) \right)$$

(証明は数学的帰納法)

これをTD法の更新とする。

$$V(S_t) \leftarrow (1-\alpha) V(S_t) + \alpha \sum_{n=0}^{k-1} \left(\gamma^n r_{t+n} + \gamma^k V(S_{t+k}) \right)$$

・ Xソリット

1ヶ月先まで読むので、所解に陥りにくい。

・ ティXソリット

学習速度が遅い。

→ トレードオフ。

→ ∞ で モンテカルロ と一致する 書いてあげて
おいてか？ っ て 顔 になってます。

ここまで。