

تمرین سری 7.2

امیرحسین باقری 9810
5621

سوال اول :

ابتدایین گفته شد که نیاز به TA مقدار $iteration$ های پای است.

و همچنین مقدار $iteration$ های $Value$ مقدار حسابات V^k

است و همچنین مقدار $iter$ های پای $Policy$ مقدار حسابات

در آید است که های π است.

در حالت کلی $Policy$ مقدار $iter$ های کمی که پای $converge$ نیاز دارد

اما هر $Policy$ $iter$ هزینه حساباتی بیش از $Value$ $iteration$ به $Value$ دارد.

بنابراین با مقدار $iteration$ های پای $Value$ شروع

می است.

ابتدا = 0

Value iteration:

for $i \leq \text{iter}$ & $(V_{k+1}^s \neq V_k^s \text{ for all } s)$ do

$$V_{k+1}(s) = \max_a \sum_{s'} T(s, a, s') [R(s, a, s') + \gamma V_k^*(s')]$$

$\underbrace{\hspace{10em}}_{O(s^2 A)}$

order of value iteration $\longrightarrow O(n s^2 A)$

Policy iteration:

① initialize random $\pi(s)$ first

②: iterate until values converge

$$V_{k+1}^{\pi_i}(s) \leftarrow \sum_{s'} T(s, \pi_i(s), s') [R(s, \pi_i(s), s') + \gamma V_k^{\pi_i}(s')]$$

③ extraction until convergence

$$\pi_{i+1}(s) = \arg\max_a \sum_{s'} T(s, a, s') [R(s, a, s') + \gamma V^{\pi_i}(s)]$$

stop, iteration order $\approx 10^4$

1

state ③ → $|s| \times |s| \times |a| = (s^2 a)$

state ② → $|s|^3$ order آن

است

Convergence of $V_{k+1}^{\pi_i}(s) \leftarrow \sum_{s'} P[R + \gamma V_k^{\pi_i}(s)]$

← اساس الگوریتمی Cse 473 (در علم) $O(s^3)$ لکه Zettlemyer

→ $O(|s|^3 + |s|^2 a)$ ← iterations است

نیاز این است که $|s| > a$ باشد که در اکثریت مواقع هست

در مقدار $iter$ $value$ از $Policy$ $iteration$ $iteration$ بیشتر است

$$Policy \rightarrow O(N (1S^3 + 1S^2A_1))$$

$$Value \leftarrow O(N 1S^2A_1)$$

$$Q(Policy) > O^{Value}(iter)$$

Value \Leftarrow در تعداد iter ها که یابی

سرع کی است که در صفحه قبل عنوان نموده ایم که اگر از A ایستای با A ایستای هم نباشد،
در صفحه قبل عنوان نموده ایم که اگر از A ایستای با A ایستای هم نباشد،
باز هم value سریع تر است.

سوال 2

جای حل سوال ۵ همین جا کنیم که تعدادی ۷ و ۹ را داریم.

Calculating policy from v^*

$$\pi^*(s) = \underset{a}{\operatorname{argmax}} Q^*(s, a) \rightarrow |S||A|$$

$$\hookrightarrow O(|S||A|)$$

بنابراین استفاده از Q ها بهتر است.

دقت کنید که غرضی که داده ایم که خطای Q و V را داریم.

چون جای extract policy باید الگوریتم converge را

در هر iteration پیدا کنیم. باید کار اضافی توان Q را این

پیدا کرد. بنابراین حساب Q ها و V ها

از نظر زمانی تفاوت زیادی ندارند.

$$O(Q) = O(|S||A|)$$

$$O(V) = O(|S|^2|A|)$$

بنابراین Q بهتر است.