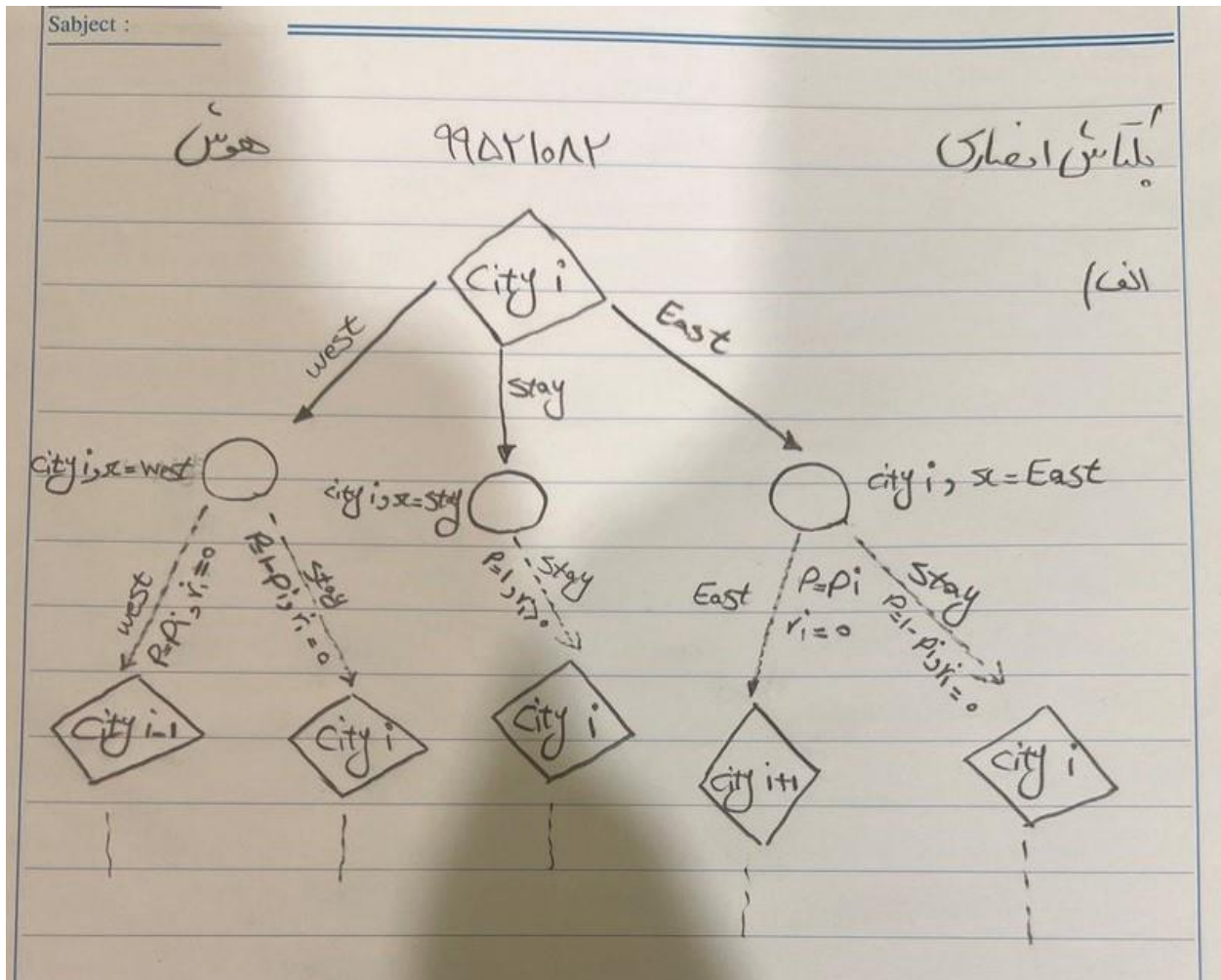


تمرین سری 12
بکتاش انصاری
99521082

بخش تئوری :



$$\gamma = \frac{1}{2} \text{ و } r_i = 1, \rho_i = 1 \quad (ب)$$

در ابتدا تمام مقادیر V را صفر می‌گذاریم؟
 چون $a = \text{stay}$ یا $s' = s$ ؟
 $V_0^{\text{stay}}(1) = 0$
 $V_1^{\text{stay}}(1) = \sum_s T(s, \text{stay}, s) [R(s, \text{stay}, s) + \frac{1}{2} V_0^{\text{stay}}(s)]$

$$\Rightarrow V_1^{\text{stay}}(1) = 1 \times [1 + \frac{1}{2} \times 0] = 1$$

$$V_2^{\text{stay}}(1) = 1 \times [1 + \frac{1}{2} \times 1] = \frac{3}{2}$$

2404

Date : / /

Subject :

$$V_{k+1}^{\text{stay}}(1) = 1 \times (1 + \frac{1}{2} V_k^{\text{stay}}(1))$$

← که در عبارت بالا در $k \rightarrow \infty$ مقدارش ۲ می‌شود
 $V^{\text{stay}}(1) = 2$

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

$k=1$
 $a = \text{stay} \leadsto V_1^{\text{stay}}(1) = 1 \times [1 + \frac{1}{2} \times 0] = 1$

$a = R \leadsto V_1^R(1) = 1 \times [1 + \frac{1}{2} \times 0] = 1$

Right
 $V_1^*(1) = 1$, $V_2^*(1) = \frac{3}{2}$, $V_3^*(1) = \frac{5}{2}$... $V_k^*(1) = 2 - \frac{1}{2^k}$

$k \rightarrow \infty \Rightarrow V^*(1) = 2$

بخش ت)

چون مقدار گاما برابر 1 است پس ما مشکل کاهش ریوارد به ازای قدم های بلند نداریم پس میتوانیم از این مورد عبور کنیم.

حال از انجایی که ممکن است در یک قدم ما دارای ریوارد 0 یا منفی باشیم اما آن ریوارد صفر باعث یک ریوارد بهتر در قدم های بعدی ما باشد. پس میتوانیم مانند روش hill climbing در هر مرحله بهترین reward را انتخاب کنیم. که این روش مشکلی که دارد این است که ممکن است در local maximum گیر کنیم که برای آنکه این مشکل برطرف شود میتوانیم یک بازه انتخاب کنیم (مثلا تعداد n قدم از محل شروع) و تمام این n قدم را طی کنیم و از بین آن ها ماکسیمم مقدار را انتخاب کنیم.

Date : / /

Subject :

(S, a, r, S')	$Q_{1, \text{stay}}$	$Q_{1, \text{East}}$	$Q_{2, \text{West}}, Q_{2, \text{stay}}$	(\sum)
initial	0	0	0	0
$(1, \text{stay}, 4, 1)$ $\frac{1}{F} \times 0 + \frac{1}{F} \times (F + \frac{1}{F} \times 0) = 1$	$\frac{1}{F} \times 0 + \frac{1}{F} \times (F + \frac{1}{F} \times 0) = 1$	$\frac{1}{F} \times 0 + \frac{1}{F} \times (F + \frac{1}{F} \times 0) = 1$	0	0
$(1, \text{East}, 0, 2)$ $\frac{1}{F} \times F + \frac{1}{F} \times (0 + \frac{1}{F} \times 0) = 1$	$\frac{1}{F} \times F + \frac{1}{F} \times (0 + \frac{1}{F} \times 0) = 1$	$\frac{1}{F} \times F + \frac{1}{F} \times (0 + \frac{1}{F} \times 0) = 1$	0	0
$(2, \text{stay}, 6, 2)$	1	1	$\frac{1}{F} \times 0 + \frac{1}{F} \times (F + \frac{1}{F} \times 0) = 1$	$\frac{1}{F} \times 0 + \frac{1}{F} \times (F + \frac{1}{F} \times 0) = 1$
$(2, \text{West}, 0, 1)$	1	1	$\frac{1}{F} \times F + \frac{1}{F} \times (0 + \frac{1}{F} \times 0) = \frac{V}{F}$	$\frac{1}{F} \times F + \frac{1}{F} \times (0 + \frac{1}{F} \times 0) = \frac{V}{F}$
$(1, \text{stay}, 4, 1)$ $\frac{1}{F} \times 1 + \frac{1}{F} \times (F + \frac{1}{F} \times 1) = \frac{11}{F}$	$\frac{1}{F} \times 1 + \frac{1}{F} \times (F + \frac{1}{F} \times 1) = \frac{11}{F}$	$\frac{1}{F} \times 1 + \frac{1}{F} \times (F + \frac{1}{F} \times 1) = \frac{11}{F}$	$\frac{V}{F}$	$\frac{V}{F}$

$$\text{Sample} = R(S, a, S') + \gamma \max_{a'} Q(S', a')$$

$$Q(S, a) = (1 - \alpha) Q(S, a) + \alpha \text{ sample}$$

$$\alpha = \frac{1}{F}, \gamma = \frac{1}{F}$$

PADA

بخش عملی :

ابتدا توضیحی راجب پیاده سازی توابع می‌دهم :

تابع `getQvalue` :

از آنجایی که ما از Counter که یک دیکشنری خاص است استفاده می‌کنیم اگر key مورد نظر در دیکشنری وجود نداشته باشد به عنوان value به ما مقدار 0 را برمیگرداند و ValueError نخواهیم داشت. که اگر این مقدار برابر 0 باشد مقدار 0.0 ریترن میشود.

نکته پیاده سازی این تابع در این است که ما به عنوان کلید هر مقدار یک تاپل دوتایی می‌دهیم تا state و action هر دو مشخص شوند.

تابع `computeValueFromQvalues` :

در این تابع ما مقدار ماکسیمم Q ها را برای یک state مشخص و بین تمام اکشن های legal محاسبه می‌کنیم که اگر اکشن مجازی وجود نداشته باشد باید مقدار 0.0 ریترن شود.

تابع `computeActionFromQvalues` :

در این بخش مانند بخش قبل باید اکشن بهترین Q را خروجی بدهیم.

که در اینجا برای پیاده سازی من ابتدا مقدار تابع قبل را محاسبه کردم تا بهترین Q محاسبه شود و سپس کلید آن value را بدست آورده که تاپل دوتایی است و مقدار دوم تاپل را خروجی دادم که action مورد نظر خواهد بود.

تابع action :

در اینجا ما باید با احتمال epsilon محاسبه کنیم که چقدر agent ما از مقدار های محاسبه شده و از تجربه خودش استفاده کند و با چه احتمالی یک تجربه رندوم جدید را اعمال کند که برای این کار با استفاده از تابع flip_coin یک عدد رندوم بین صفر و یک ایجاد میکنیم که اگر آن عدد از epsilon کمتر بود یک مقدار رندوم در غیر اینصورت از مقادیر تجربه شده استفاده کند.

تابع update :

در این تابع نیز ما فرمول زیر را با استفاده از sample داده شده محاسبه میکنیم.

$$sample = R(s, a, s') + \gamma \max_{a'} Q(s', a')$$

no longer policy
evaluation!

- Incorporate the new estimate into a running average:

$$Q(s, a) \leftarrow (1 - \alpha)Q(s, a) + (\alpha) [sample]$$

حال به سوالات جواب میدهیم :

دلیل آنکه agent همواره مطابق جهتی که ما نشان میدهیم نمیروند این است که حرکت agent بصورت قطعی نیست و هر 4 جهت یک احتمالی دارند و قطعیتی برقرار نیست.

تصویر نهایی :



لاگ ترمینال در فایل زیپ موجود است.

برداشت من این است که هر چه تعداد iteration ها زیاد شود مسیر بهینه که به صورت بالا راست میباشد احتمال بیشتری پیدا کرده و ما از اینکه این مسیر بهینه است بیشتر مطمئن میشویم و احتمال جهت های این مسیر بیشتر میشود. و مسیر هایی که به ترمینال منفی ختم میشوند احتمال کمتری میگیرند.

: Crawler

مقادیر موجود :

: learning rate : Alpha

در اصل این مقیاس میزان یادگیری agent ما از روی سمپل جدید را مشخص میکند که در اصل باعث میشود که ما در Q جدید چ مقدار از مقدار قبلی Q و چه مقدار از sample جدید را استفاده کنیم.

این مقدار عددی بین صفر و یک است که اگر این مقدار برابر صفر باشد. Agent ما از sample یا action جدید چیزی یاد نمیگیرد و اگر مقدار آن برابر یک شود agent ما به طور کامل تمام دانش قبلی خود را ignore میکند و تنها از sample استفاده میکند. هر چه این مقدار زیاد باشد باعث میشود Q بصورت سریع مقدار عوض کند.

: Epsilon

این مقدار مربوط به انتخاب به روش epsilon-greedy است که ما در آن برای انتخاب action مورد نظر استفاده میکنیم که در کد آن را در تابع action پیاده کردیم. به این شکل که این مقیاس مشخص میکند که ما از یک action رندوم استفاده کنیم یا از اکشن های خوبی که تا الان داشتیم استفاده کنیم. اگر این مقدار صفر باشد ما هیچوقت explore نداریم و فقط exploit داریم و از دانش کنونی خود استفاده میکنیم و تجربه جدیدی نمیکنیم. اگر این مقدار برابر یک باشد ما هیچوقت از تجربه های گذشته خود استفاده نمیکنیم و فقط از مقادیر رندوم استفاده میکنیم. معمولاً این مقدار کوچک و نزدیک صفر است.

: Discount

مقیاسی است که در مقدار پیشبینی شده برای مقادیر پیش رو ضرب میشود تا مقدار Reward مقادیر دور با مقادیر نزدیک برابر نشود.

اگر مقدار آن را صفر در نظر بگیریم agent ما بصورت کامل reward های آینده را ignore میکند و تنها به reward کنونی اهمیت میدهد.