

## به نام خدا

استاد: دکتر هراتی- دکتر غیاثی

تمرین اول

یادگیری تقویتی

در این تمرین قصد داریم مساله ای را که در کلاس مطرح شد در مورد 10 armed bandits را بررسی کنیم .

در این تمرین قصد داریم حالت های مختلف ان را به ازای  $\lambda$  های مختلف بررسی و شکل ان را رسم کنیم .

در ابتدا ساده ترین الگوریتم را در نظر می گیریم که به هر اکشن یک میانگین از پاداش های بدست آمده را نسبت دهیم و می خواهیم نتایج الگوریتم  $\epsilon - Greedy$  را بررسی کنیم . مساله 10 armed bandits را که

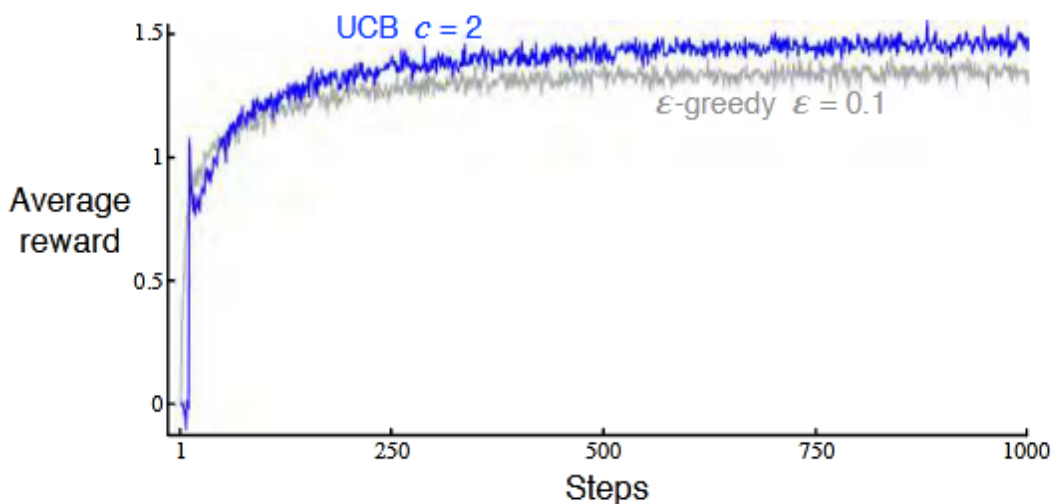
در کلاس مطرح شد را به ازای  $\epsilon$  های 0.1 و 0.5 و 0.7 و 1 بررسی کنید و نمودار میانگین امتیاز به ازای 1000 گام را رسم کنید .

1. با افزایش  $\epsilon$  نتایج چگونه خواهد شد ؟ چرا ؟

2. الگوریتم  $\epsilon - Greedy$  چه مزیتی دارد ؟

3. الگوریتم UCB چگونه می تواند  $Exploration$  را برای ما فراهم کند ؟

4. در مساله 10 armed bandits در حالتی که از الگوریتم UCB استفاده کنیم با  $c=2$  نتیجه به شکل زیر خواهد بود.



همانطور که مشاهده می کنید در گام 11 ما یک پیک داریم . توضیح دهید چه عاملی باعث به وجود آمدن پیک شده است.

برای انجام تمرین می توانید از کد های آماده لینک زیر استفاده کنید.

[https://github.com/ShangtongZhang/reinforcement-learning-an-introduction/blob/master/chapter02/ten\\_armed\\_testbed.py](https://github.com/ShangtongZhang/reinforcement-learning-an-introduction/blob/master/chapter02/ten_armed_testbed.py)

لطفا پاسخ تمرینات را در قالب یک فایل PDF در ویو اپلود کنید .

موفق باشید