



به نام او

## درس یادگیری تقویتی

پروژه اول، پوریا صفائی ۹۸۱۱۰۴۰۲

### سوال ۱

تابع `generate random samples` که در فایل کد قابل مشاهده است وظیفه تولید داده های تصادفی نرمال با امید ریاضی  $q_*(a)$  و واریانس ۱ را دارد. (تمام کدهای مربوط در فایل `project 1.py` در فایل زیپ قرار داده شده است).

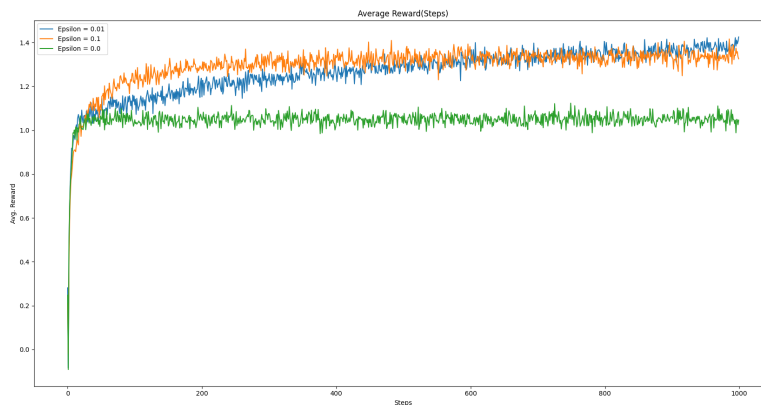
### سوال ۲

با توجه به شکل کشیده شده در کتاب، مقادیر تقریبی هر یک از  $q_*(a)$  ها را به صورت زیر مینویسیم:

$$q_*(1) = 0.25, q_*(2) = -0.6, q_*(3) = 1/5, q_*(4) = 0.5, q_*(5) = 1/25, \\ q_*(6) = -1/5, q_*(7) = -0.2, q_*(8) = -1/10, q_*(9) = 0.9, q_*(10) = -0.5$$

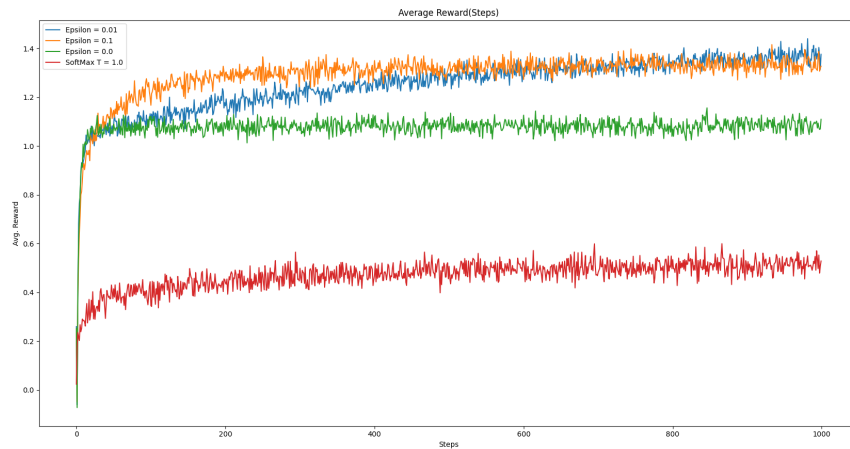
### سوال ۳

همانطور که در نمودار زیر میتوان مشاهده کرد، با فرض اجرای ۱۰۰۰ گام و در هر گام اجرای ۲۰۰۰ مرحله حاصل الگوریتم اپسیلون-حریصانه برای مقادیر مختلف اپسیلون قابل مشاهده است (مقدار اپسیلون صفر معادل همان الگوریتم حریصانه ساده است). همانطور که میتوان این تصویر را با عکس کتاب مشاهده کرد، شباهت کامل مشهود است. به ازای اپسیلون کوچک، پاداش متوسط با شیب کمی در حال افزایش است و به دلیل اینکه اکسپلوریشن کمتر از حالت اپسیلون بزرگتر است، نقطه بهینه خود را کمی دیرتر پیدا میکند. اما چون در این حالت اکسپویتهیشن بیشتر است، پس نمودار با یک شیب بزرگتر از حالت دیگر در حال افزایش میانگین پاداش است. حالت اپسیلون صفر نیز با پیدا کردن اولین اکشن که بزرگتر از صفر پاداش دارد، تا انتهای الگوریتم همان عمل را انجام میدهد.



## سوال ۴

در این بخش نیز میتوان نتیجه سافت مکس را برای مقدار دمای  $T = 1$  مشاهده کرد. همانطور که میبینیم، میزان متوسط پاداش در این حالت نسبت به دیگر حالت ها به طور متوسط کمتر به دست آمده، چرا که دلیل آن اسن است که سافت مکس به اندازه الگوریتم های قبلی حریصانه نیست و اجازه اکسپلوریشن را به ایجنت میدهد. از طرف دیگر چون متوسط توزیع های مختل نزدیک به هم است، پس توزیع احتمالاتی ای که سافت مکس برای حرکت بعدی در نظر میگیرد نیز یکنواخت تر از الگوریتم های حریصانه قبلی است. اما از طرف دیگر میتوان مشاهده کرد که نمودار الگوریتم سافت مکس میتواند صعودی باشد و اگر در مرور زمان اگر مقدار  $T$  را کاهش دهیم، این الگوریتم پس از مدتی اکتشاف در ابتدا نهایتا به بهترین جواب میل خواهد کرد. (در اینجا چون مقدار  $T$  در مرور زمان ثابت قرار داده شده است، سافت مکس ابتدا امتوسط ارزش همه عمل ها را پیدا میکند و سپس پس از آن با یک توزیع متناسب با پاداش هر عمل، عمل بعدی خود را انتخاب میکند.)



اما همانطور که در نمودار زیر میبینیم (نمودار بنفش)، اگر مقدار  $T$  را به مرور زمان کاهش دهیم (در اینجا ما دما را با ضریب  $0.05$  در هر مرحله کاهش میدهیم)، مشاهده میکنیم که الگوریتم سافت مکس نیز در اینجا به بهینه مقدار ممکن میل میکند.

