



دانشگاه تهران  
پردیس دانشکده‌های فنی  
دانشکده برق و کامپیوتر



گزارش تمرین شماره چهارم  
درس یادگیری تعاملی  
پاییز ۱۴۰۱

نام و نام خانوادگی  
سیاوش رزمی  
شماره دانشجویی  
۸۱۰۱۰۰۳۵۲

فهرست

چکیده	۳
سوال ۱ -	۴
۱	۴
۲	۵
سوال ۲ -	۶
منابع	۱۱



## چکیده

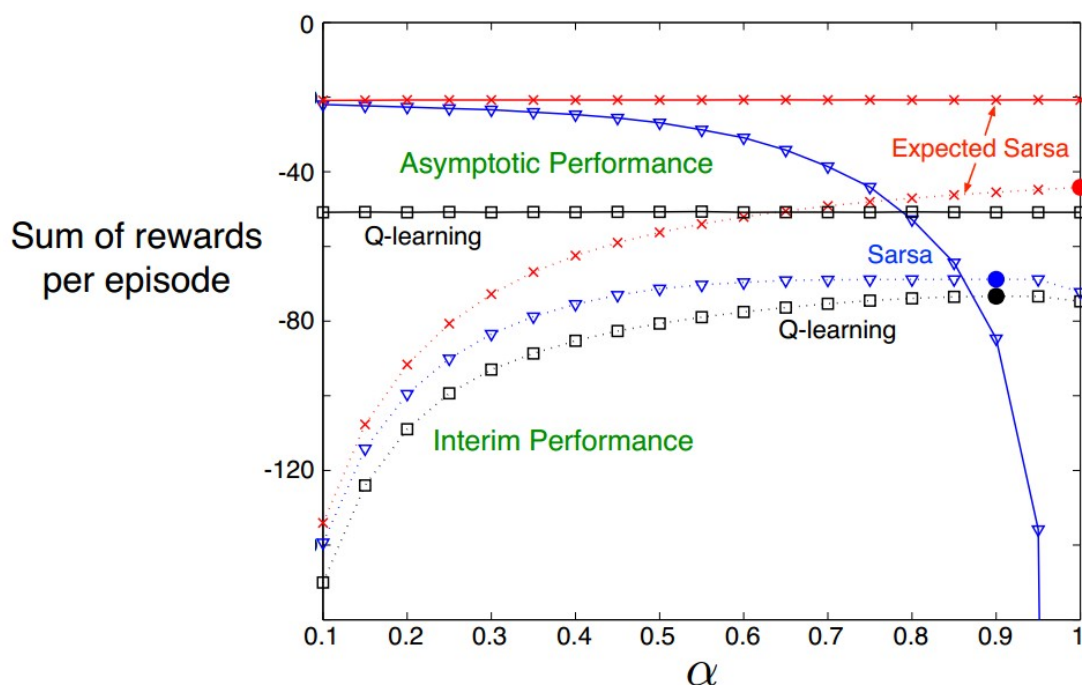
---

در این تمرین برخلاف تمرین‌های گذشته قصد داریم به حل مسأله با فرض ناشناخته بودن محیط یا به اصطلاح مسائل Model-free بپردازیم، در بخش اول به بررسی برخی از مسائل تحلیلی و شناخت الگوریتم‌های model-free می‌پردازیم و در بخش دوم با استفاده از این الگوریتم‌ها در محیط Taxi کتابخانه‌ی gym سعی در حل این مسأله می‌کنیم.

۱.۱. به طور کلی الگوریتم Expected Sarsa به دلیل متوسط گرفتن از تمامی مقادیر ارزش اکشن حالت‌ها نسبت به الگوریتم Sarsa دارای پیچیدگی محاسباتی بیشتری است اما انجام این کار باعث می‌شود که واریانس تخمین در این الگوریتم کمتر شده و همچنین Sampling efficiency بالاتر برود، بنابراین الگوریتم Expected Sarsa با تعداد sample های کمتری می‌تولند به سیاست بهینه برسد و دارای میزان Regret کمتری است اما میزان محاسبه و زمان مصرفی آن بیشتر است.

۱.۲. در مورد تفاوت سیاست بهینه و اپسیلون بهینه به نظر می‌رسد روند همگرایی در این دو الگوریتم در این دو نوع سیاست تقریباً یکسان است اما میزان Regret به دست آمده در الگوریتم بهینه بیشتر خواهد بود به این دلیل که میزان Reward بهینه به شکل واضح از الگوریتم اپسیلون بهینه بیشتر است.

۱.۳. در زمانی که محیط ما کاملاً Deterministic باشد الگوریتم Expected sarsa نسبت به تغییرات نرخ یادگیری حساسیت بسیار کمی دارد و با افزایش نرخ یادگیری در کوتاه مدت بازدهی آن افزایش پیدا کرده و در بلند مدت نیز تغییرات نرخ یادگیری تأثیر به سزایی در همگرایی آن ندارد، اما بازدهی الگوریتم Sarsa با نزدیک شدن نرخ یادگیری به ۱ کاهش پیدا میکند و در بلند مدت تنها با کاهش نرخ یادگیری به مقدار زیاد امکان همگرایی دارد، بنابراین میتوان گفت در الگوریتم Sarsa بایستی پس از مدتی نرخ یادگیری را کاهش داد تا الگوریتم امکان همگرایی داشته باشد، اما در الگوریتم Expected sarsa تغییرات نرخ یادگیری تأثیر به سزایی در عمل کرد الگوریتم ندارد و نرخ یادگیری در الگوریتم Sarsa بایستی سریع‌تر کاهش پیدا کند.



شکل ۱-۱: تأثیر تغییرات نرخ یادگیری در الگوریتم های Sarsa و Expected Sarsa

اما در مورد نرخ کاهش اپسیلون، به دلیل اینکه الگوریتم Sarsa Expected دارای Sample efficiency بالاتری نسبت به الگوریتم Sarsa است و تعداد سمپل کمتری نیاز دارد تا بتواند سیاست بهینه را پیدا کند به نظر می‌رسد در زمان کوتاه تری می‌تواند به سمت الگوریتم حریصانه تر حرکت کند و نرخ کاهش اپسیلون برای این الگوریتم بایستی بیشتر از الگوریتم Sarsa باشد.

۲

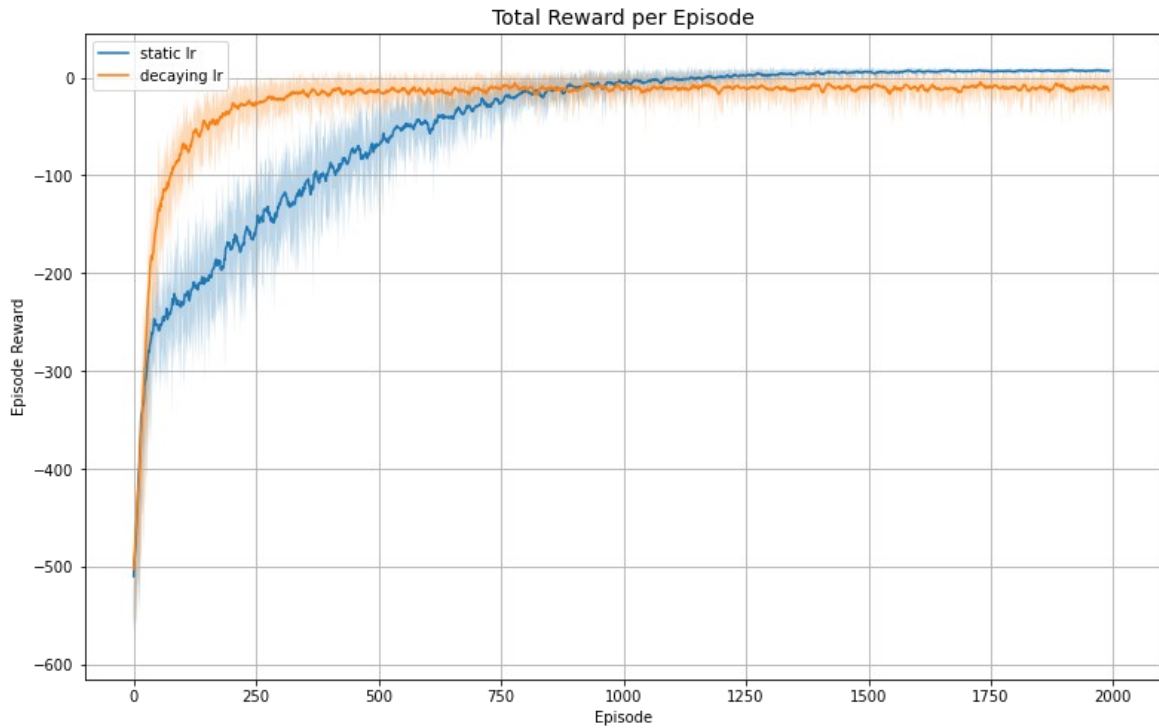
طبق خاصیت error reduction property در الگوریتم های n-step return تخمین آن‌ها از  $v_\pi$  در مرحله n در بدترین حالت از تخمین  $v_{t+n-1}$  بهتر است، درواقع:

$$\max_s |\mathbb{E}[G_{t:t+n}|S_t = s] - v_\pi(s)| \leq \gamma^n \max_s |V_{t+n-1}(s) - v_\pi(s)|$$

بدین معنا که در بدترین حالت هم تخمین این الگوریتم از  $\gamma^n$  برابر ماکزیمم تفاوت بین تخمین  $v_\pi$  و  $V_{t+n-1}$  کمتر مساوی است.

## سوال ۲ -

۱.

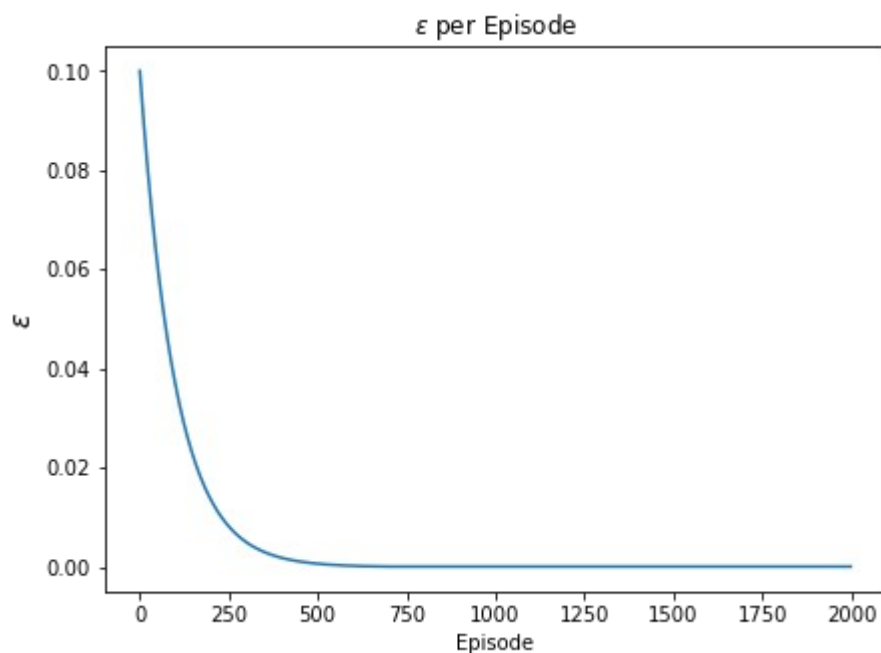


شکل ۲-۱: نتایج مقایسه الگوریتم Q learning با استفاده از نرخ یادگیری ثابت و کاهشی

همانطور که از شکل ۲-۱ قابل مشاهده است، الگوریتم با نرخ یادگیری کاهشی بسیار سریع تر به همگرایی می‌رسد و با توجه به بازه اطمینان ۹۵ درصد (هاله رنگی در اطراف مقدار میانگین) این اختلاف معنا دار است، اما میزان پاداش نهایی در الگوریتم با نرخ یادگیری ثابت در نهایت کمی بیشتر از حالت کاهشی است.

\* تصاویر عمل کرد عامل در ۲۰ مرحله اجرا برای تمامی مسائل در فولدر pictures به همراه گزارش ارسال شده است.

مقدار اپسیلون در این سؤال و تمامی سؤال‌های بعدی به شکل نمایی کاهش پیدا میکند، بدین صورت که یک تابع نمایی با نرخ کاهش ۰.۰۰۱ از ۱ شروع و تا ۰.۰۰۱ کاهش پیدا میکند سپس با ضرب این مقدار در مقدار اپسیلون مشخص شده در صورت سؤال (۰.۱) مقدار نهایی اپسیلون در هر اپیزود به دست می‌آید.

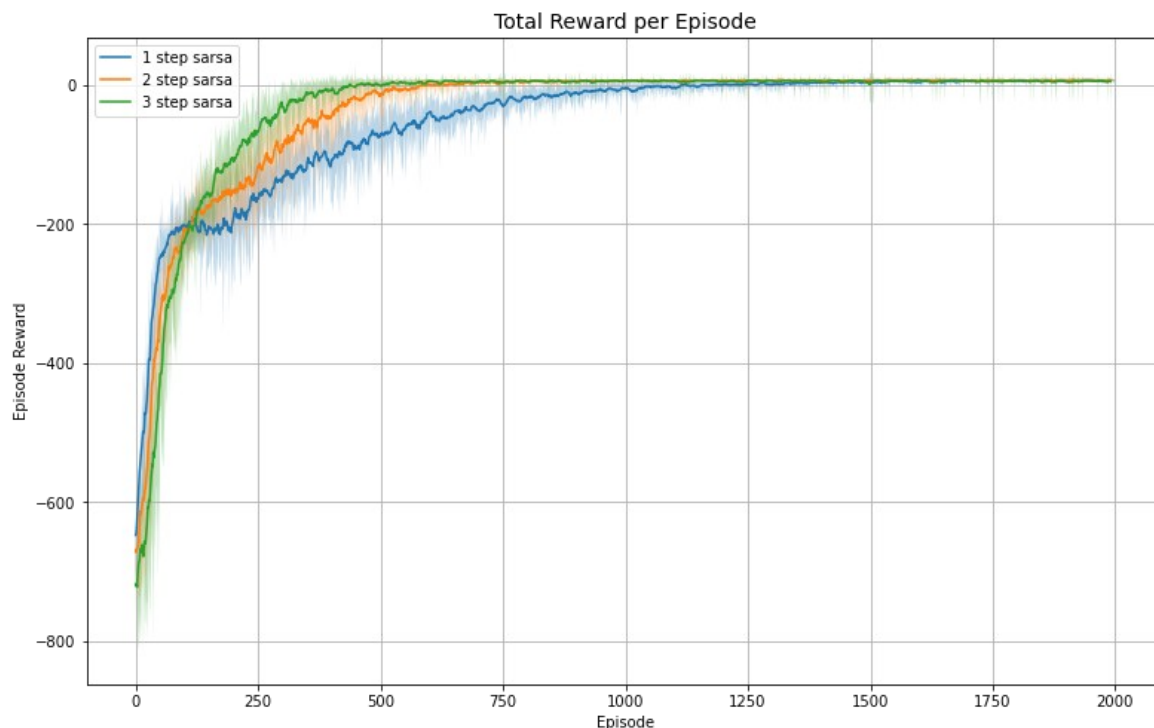


شکل ۲-۲: نحوه کاهش اپسیلون در هر اپیزود

۲.

در محیط معرفی شده تاکسی می‌تواند در ۲۵ خانه، مسافر در ۵ محل (چهار نقطه رنگی به همراه خود تاکسی) و ۴ مقصد است بنابراین به طور ۵۰۰ حالت مجزا در این مسأله تعریف می‌شود، حالت‌های غیر قابل دسترسی مربوط به مواقعی است که محل اولیه و مقصد مسافر یکسان باشد که برابر ۱۰۰ حالت می‌شود، اما در ۴ حالت از این ۱۰۰ حالت تاکسی هم در محل مقصد است که درواقع با یک اکشن به انتهای اپیزود می‌رسیم بنابراین در کل ۴۰۴ حالت در این محیط قابل دسترسی و ۹۶ حالت غیر قابل دسترسی می‌باشد. و برای شناسایی این حالت‌ها کافیهست که محل مسافر و مقصد و محل تاکسی را با شروط ذکر شده بررسی کنیم تا دسترسی پذیر یا ناپذیر بودن آن حالت را به دست آوریم.

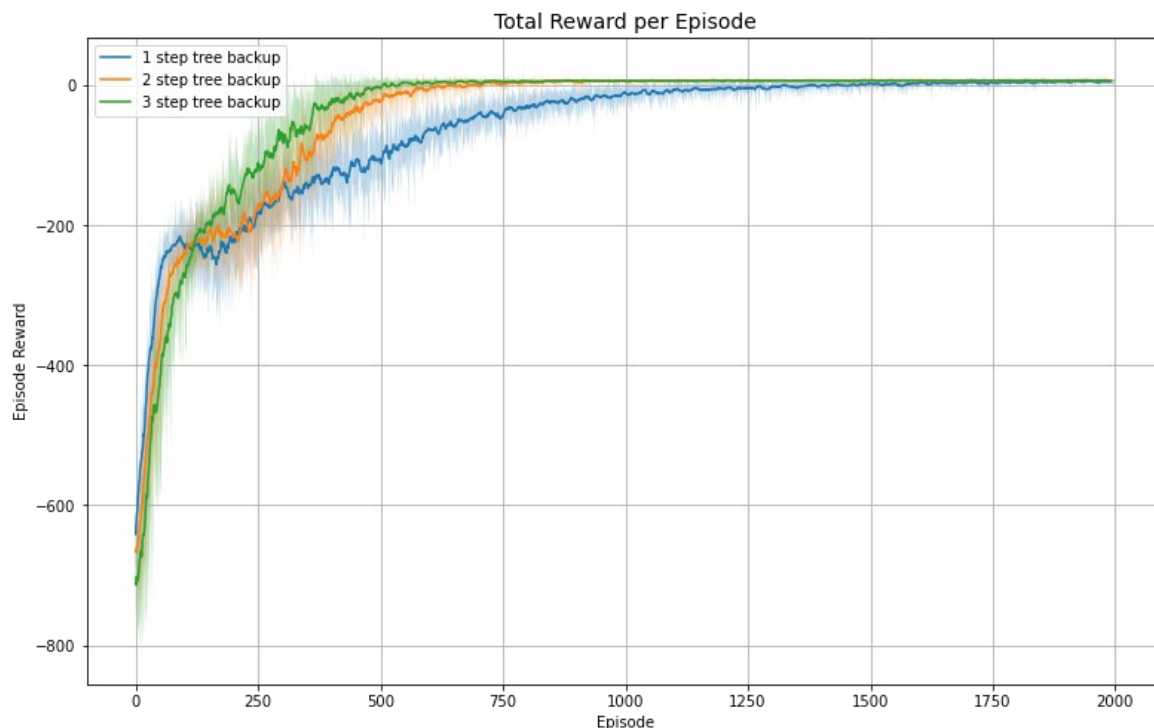
۳. هر دو الگوریتم  $n$  step sarsa و  $n$  step tree backup به ازای ۳ مقدار ۱ و ۲ و ۳ پیاده شده و نتایج زیر از اجرای آن‌ها حاصل شد:



شکل ۳-۲: نتیجه الگوریتم Sarsa برای سه مقدار مختلف

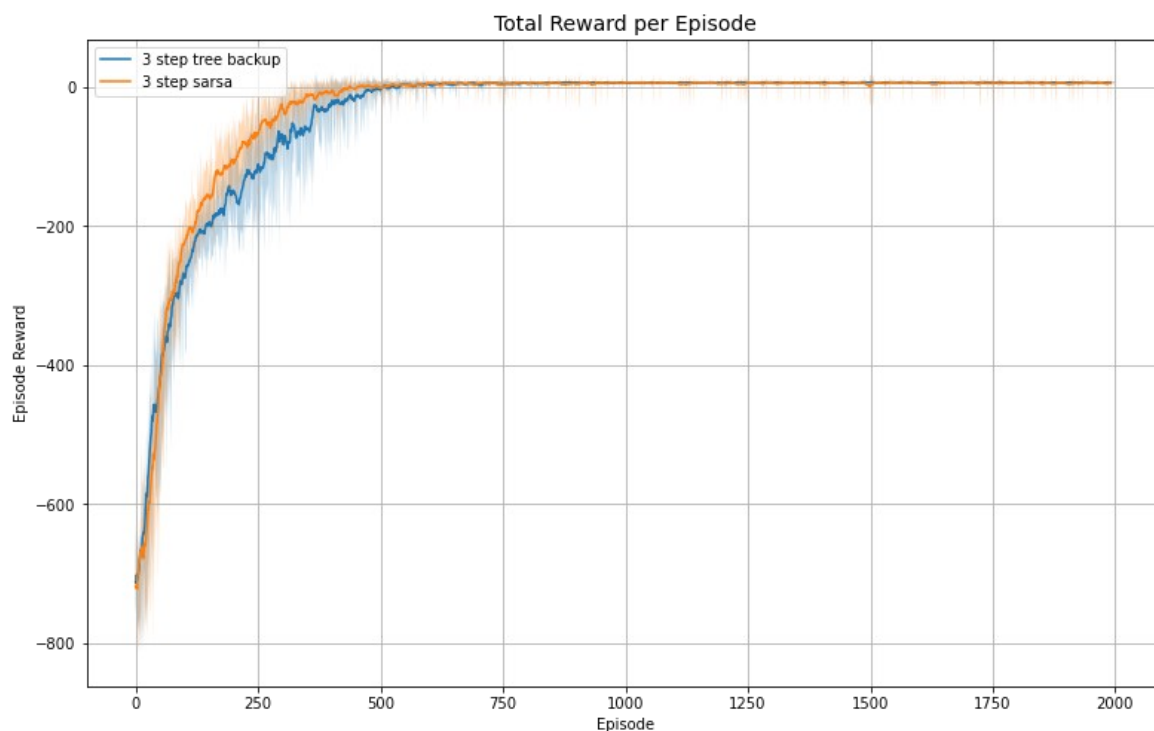
همانطور که از نتایج مشخص است مقدار  $n$  برابر با ۳ دارای کمترین میزان Regret در بین ۳ نوع است هر الگوریتم تا حدود ۵۰ اپیزود اول عمل کرد مشابهی دارند اما از آنجا به بعد با افزایش میزان  $n$  عمل کرد الگوریتم بهتر شده و در نهایت هر سه الگوریتم به یک مقدار همگرا شده‌اند با توجه به بازه اطمینان الگوریتم ها در  $n$  برابر با ۳ به شکل معنا داری از  $n$  برابر با ۱ حسرت کمتری دارد .





شکل ۴-۲: نتیجه الگوریتم tree backup به ازای n های مختلف

الگوریتم tree backup نیز عمل کردی حدوداً مشابه الگوریتم sarsa داشت و با افزایش n میزان حسرت آن کاهش پیدا میکند و در نهایت به یک مقدار یکسان همگرا می‌شوند، همچنین با مقایسه عمل کرد بهترین نسخه از هر کدام از الگوریتم‌ها به نظر می‌آید که نمیتوان گفت دو الگوریتم تفاوت معناداری در عمل کرد خود دارند زیرا بازه اطمینان هر دو تقریباً با همدیگر همپوشانی دارد.

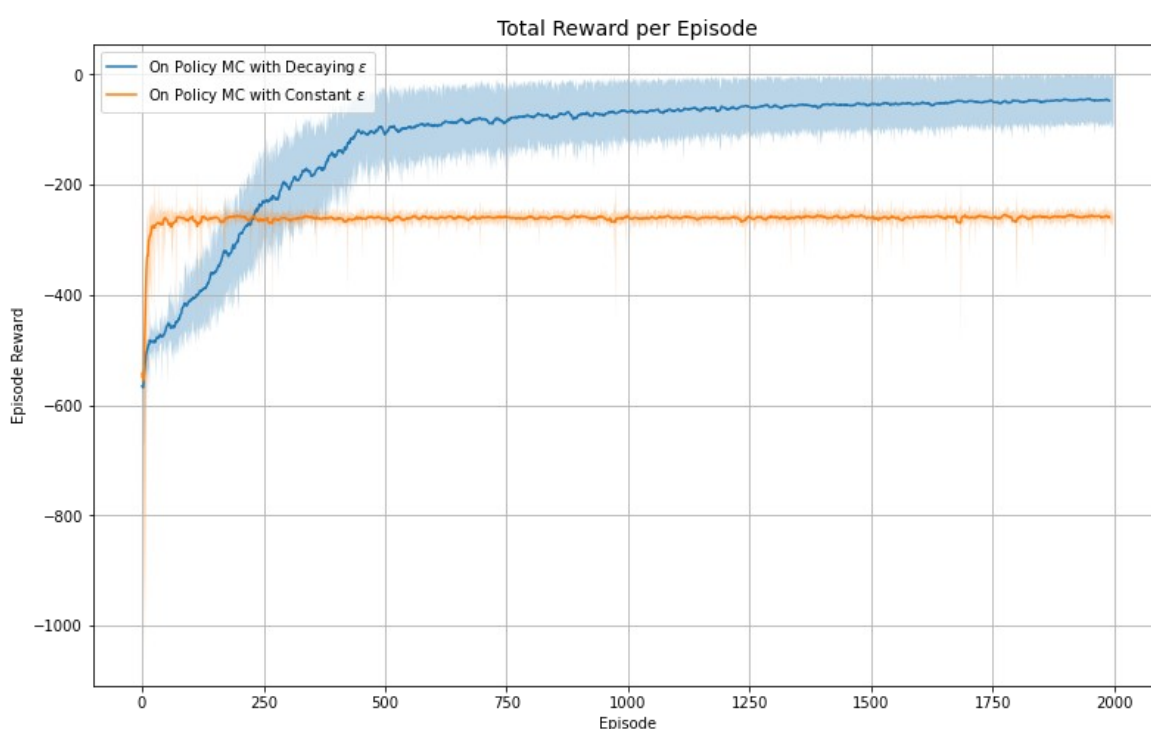


شکل ۵-۲: مقایسه عمل کرد بهترین نسخه Tree Backup و Sarsa

دلیل اینکه با افزایش میزان  $n$  سرعت همگرایی به سیاست بهینه در این الگوریتم ها بیشتر می شود این است که با زیادتر کردن مقدار  $n$  برای محاسبه ارزش هر اکشن حالت به جای استفاده از ارزش تخمینی از پاداش هایی که عامل به شکل واقعی از محیط کسب کرده استفاده میکنیم در نتیجه تخمین ما نسبت به حالت های با  $n$  کمتر دقیق تر خواهد بود و  $Q$  به دست آمده به  $Q$  بهینه سریعتر نزدیک می شود.

#### ۴.۱

الگوریتم On policy Monte carlo دو بار به ازای اپسیلون ثابت  $0.1$  و اپسیلون کاهشی  $0.5$  اجرا شد که طبق نتایج به دست آمده از شکل ۶-۲ در حالت اپسیلون ثابت الگوریتم بسیار سریعتر همگرا می شود اما به سیاست بهینه همگرا نشده و مقدار sub optimal دارد.



۶-۲: نتیجه الگوریتم On policy monte carlo به ازای مقادیر اپسیلون ثابت و کاهشی

#### ۵

سرعت یادگیری این الگوریتم نسبت به الگوریتم های قبلی کندتر است و دلیل آن نیز Sample inefficient بودن الگوریتم Monte carlo نسبت به الگوریتم های دیگر و همچنین نیاز به انجام کل اپیزود جهت به روز رسانی مقادیر  $Q$  است.

یکی از روش های افزایش سرعت الگوریتم پیاده سازی Batch sample learning است، به این معنا که از مشاهدات گذشته عامل در محیط چندین مرتبه جهت به روز رسانی استفاده کنیم.

R. S. Sutton and A. G. Barto, *An Reinforcement Learning: Introduction*. Mit Press, 2012

[١]