



یادگیری تقویتی در کنترل
تمرین کلاسی سوم: بررسی روش
Optimistic Initial Value

استاد: دکتر سعید شمعقدری

دانشجو: سیده ستاره خسروی

پائیز ۱۴۰۳

چکیده

در این تمرین به تحلیل و پاسخ تمرین کلاسی سری سوم پرداخته می شود.

واژه‌های کلیدی: یادگیری تقویتی، راهزن چند دست

فهرست مطالب

صفحه	عنوان
ب.....	فهرست مطالب
ج.....	فهرست تصاویر و نمودارها
۱.....	فصل ۱: شبیه‌سازی مسئله راهزن چنددست
۱.....	۱.۱ صورت سوال
۱.....	۱.۲ پاسخ تحلیلی و شبیه سازی
۴.....	۱.۳ منابع

فهرست تصاویر و نمودارها

صفحه

عنوان

- شکل ۱-۱: نمودار مقایسه‌ای روش‌های حل مسئله راهزن چنددست ۱
- شکل ۱-۲: خروجی شبیه‌سازی ۳

فصل ۱: شبیه‌سازی مسئله راهزن چنددست

۱.۱ صورت سوال

در گراف مقایسه روش های مختلف حل مسئله MAB در کتاب (شکل ۶-۲) بررسی کنید در روش Optimistic، چرا برای $Q_0=1$ ، عملکرد از سایر مقادیر بهتر است؟

۱.۲ پاسخ تحلیلی و شبیه سازی

ابتدا برای پاسخ به این سوال لازم است، به نمودار نگاهی بیندازیم.

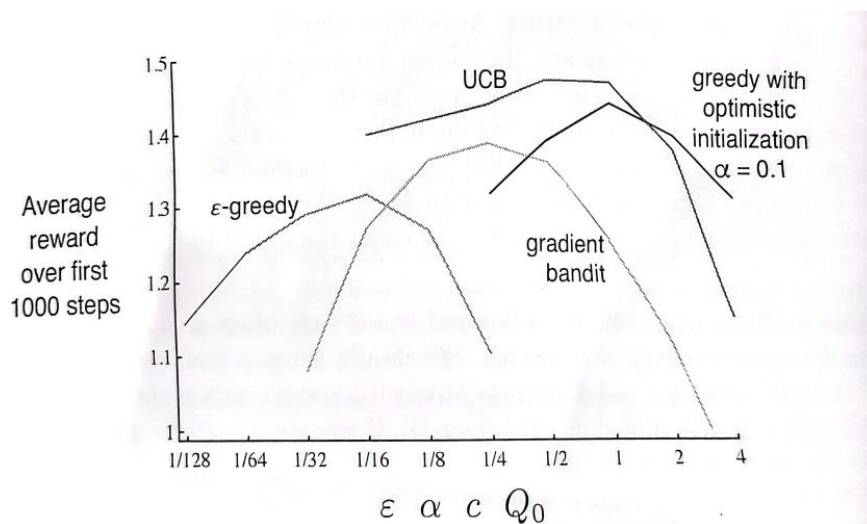


Figure 2.6: A parameter study of the various bandit algorithms presented in this chapter. Each point is the average reward obtained over 1000 steps with a particular algorithm and a particular setting of its parameter.

شکل ۱-۱: نمودار مقایسه ای روش های حل مسئله راهزن چنددست

گراف مربوط به روش شرایط اولیه خوشبینانه، بیانگر این است که در این روش، مسئله راهزن چنددست با روش حریصانه، گام ۰.۱ و شرایط اولیه خوشبینانه برابر با ۱ بیشترین پاداش میانگین را در ۱۰۰۰ گام اولیه به ارمغان می آورد. این در حالی است که بحث شد هر مقدار اولیه خوشبینانه بیشتر باشد بهتر است ولی در این نمودار چنین چیزی مشاهده نمی گردد.

باید به معادله‌ی زیر توجه کنیم.

$$Q_{n+1} = Q_n + \alpha(R_n - Q_n) = Q_n + 0.1(R_n - Q_n) = 0.9Q_n + 0.1R_n$$

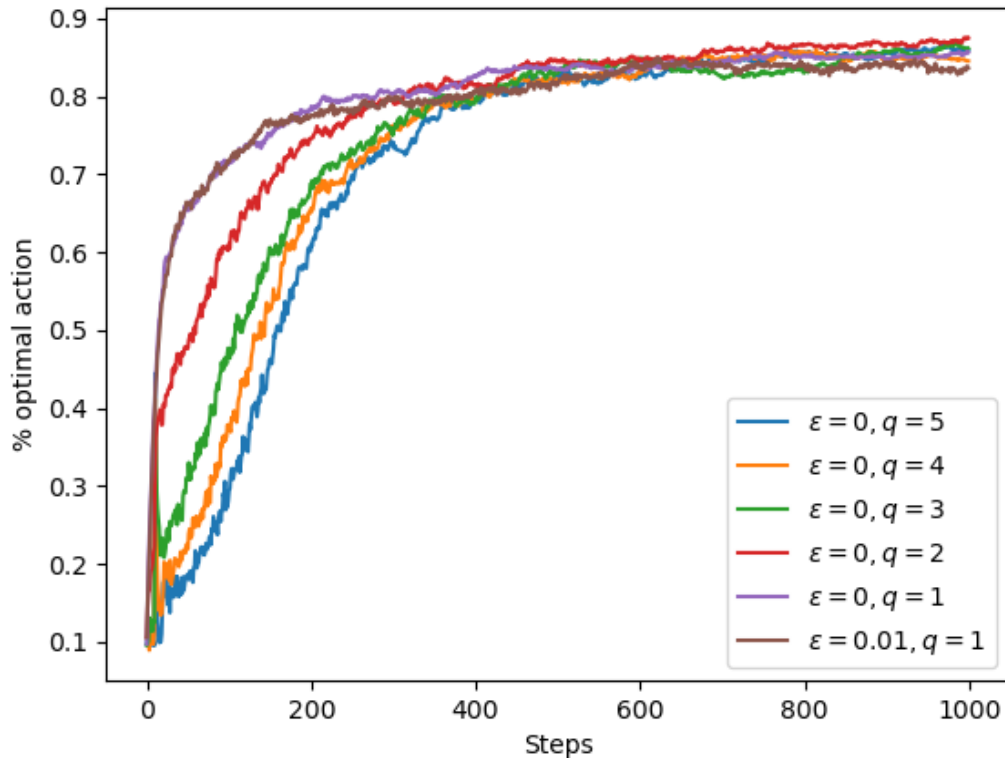
باید توجه نمود که هرچقدر مقدار Q_0 بیشتر باشد، اثر آن دیرتر از بین می‌رود و برای از بین رفتن اثر آن زمان بیشتری لازم است. با توجه به اینکه در اکتشاف عملاً یادگیری رخ نمی‌دهد، و انتخاب این مقدار اولیه با مقداری بیش از حد بالا، عامل ممکن است برای مدت طولانی‌تری به جست و جو ادامه دهد، و اصلاح این برآورد اولیه زمان بسیاری طول می‌کشد. این امر باعث تأخیر در رسیدن به بهترین عمل می‌گردد و ممکن است در مراحل اولیه عامل به عملکردی زیربهبوده برسد. ظاهراً تنظیم مقدار ۱ برای مقدار اولیه خوشبینانه، تعادل مناسبی میان اکتشاف و بهره‌برداری ایجاد می‌کند، به اندازه کافی برای انجام جست و جو خوشبینانه است، ولی انقدر هم زیاد نیست که عامل وقت زیادی را برای جست و جو صرف کند. برای مقادیر بالاتر از ۱، عامل ممکن است بیش از حد نیاز به جست و جو بپردازد.

در حالی که مقادیر اولیه بالا مثل ۵، می‌توانند برای جست و جو مفید باشند، تنظیم آنها با مقدار بیش از حد بالا می‌تواند باعث همگرایی کندتر و عملکرد ضعیف‌تر شود، زیرا عامل برای مدت طولانی به بررسی و جست و جو ادامه می‌دهد. مقدار اولیه ۱ در این مورد به نظر می‌رسد که یک تعادل مناسب بین جست و جو و بهره‌برداری را برقرار می‌کند.

برای بررسی بیشتر این موضوع، با استفاده از کدی که برای تمرین سری دوم قسمت UCB نوشته شده بود، شبیه سازی حریصانه با مقدار اولیه خوشبینانه را برای ۵ مقدار متفاوت انجام دادیم.

در شبیه سازی مقدار α برابر با ۰.۱ لحاظ کردیم، شرایط اولیه خوشبینانه نیز از ۱ تا ۵ لحاظ شدند، روش نیز greedy است و به تبع از آن مقدار ϵ برابر با صفر است. در یک حالت نیز با مقدار ϵ برابر با ۰.۰۱ نیز شبیه سازی را انجام دادیم که محل بحث نیست.

خروجی شبیه سازی در نمودار شکل ۲-۱ آورده شده است.



شکل ۱-۲: خروجی شبیه سازی

در نمودار فوق مشاهده می گردد که با افزایش مقدار اولیه خوشبینانه، مدت زمانی که عامل صرف جست و جو می کند و به خاطر آن عمل های غیربهبینه را انتخاب می کند و به تبع از آن پاداش کمتری دریافت می کند، طولانی تر می شود، در حالتی که مقدار اولیه خوشبینانه بیشتر از ۱ است، فرایند یادگیری و همگرایی کندتر می شود. مطابق توضیحات ارائه شده در صفحه قبل در اجرای با مقدار اولیه خوشبینانه ۱ تعادل بهتری میان اکتشاف و بهره برداری برقرار می گردد. به دلیل کند شدن روند همگرایی در مقادیر بالاتر از ۱ و اینکه نمودار ۱-۱ نیز مربوط به ۱۰۰۰ اجرای اولیه است، میانگین پاداش کسب شده با آن مقادیر کمتر از حالتی است که مقدار اولیه ۱ تنظیم می شود. در صورت ادامه یافتن فرایند شبیه سازی انتظار می رود که وضعیت حالتی که مقدار اولیه بیشتر از ۱ است بهتر شود. این موضوع از نمودار ۱-۲ نیز قابل استنتاج است، درواقع وقتی مقدار اولیه را افزایش می دهیم باید فرضت بیشتری نیز به عامل بدهیم تا اثر این مقدار اولیه نیز از بین برود.

۱.۳ منابع

کتاب Sutton and Barto

ریپازیتوری های:

<https://github.com/setarekhosravi/reinforcement-learning-an-introduction>