

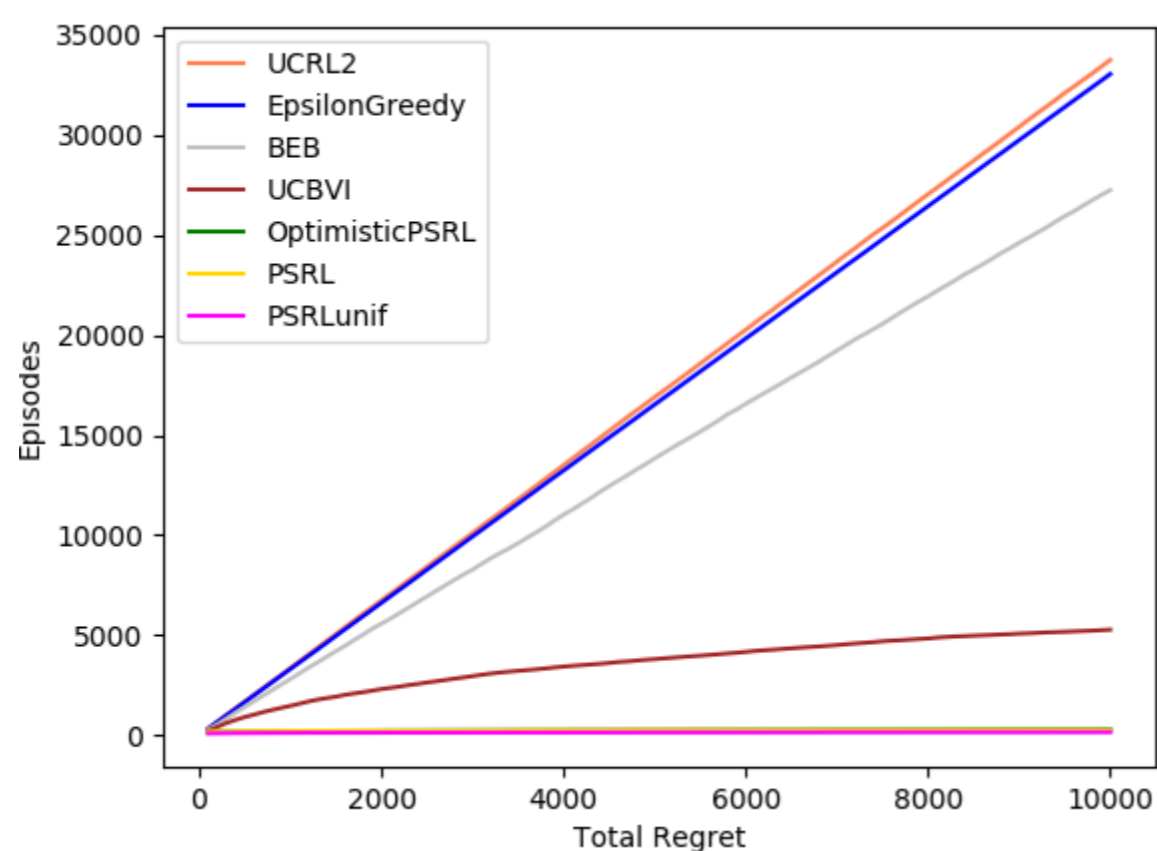
مدل سازی جدیدترین الگوریتم های یادگیری تقویتی در فرایندهای تصمیم گیری غیر متعارف مارکوف



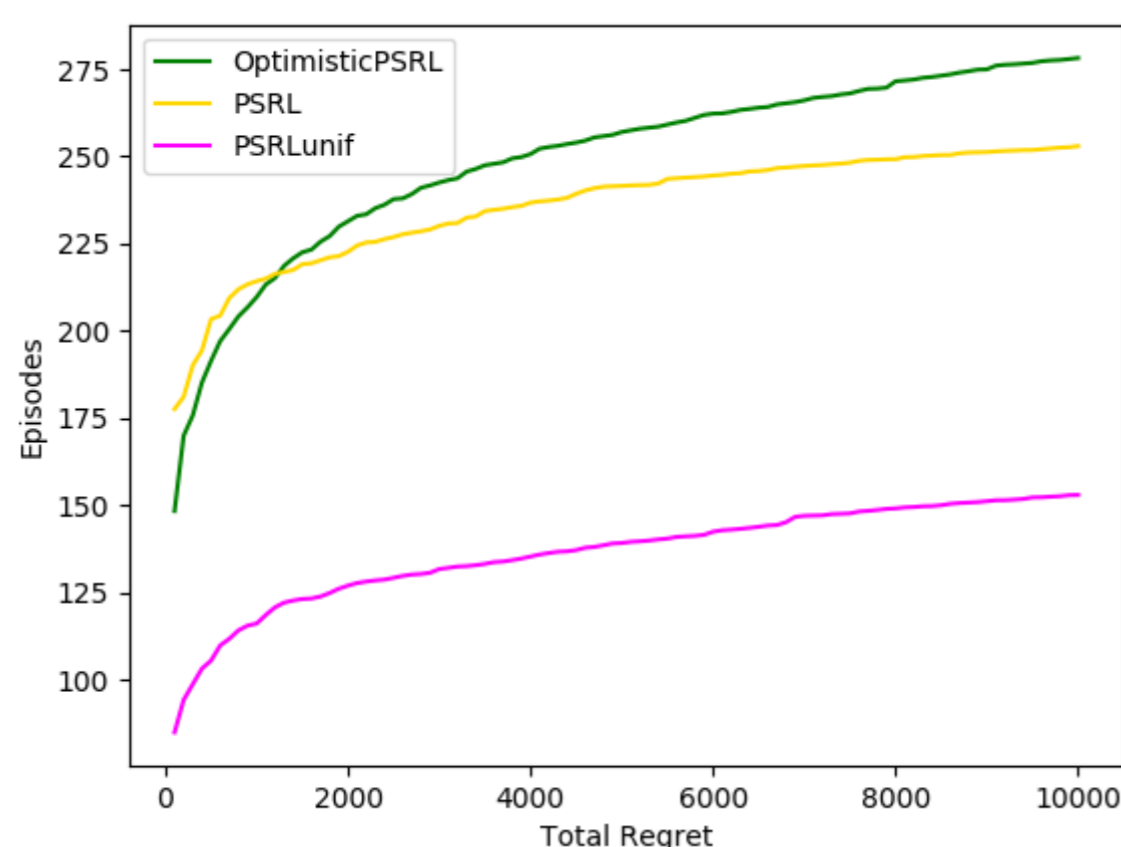
دانشجو: محمد بیات
استاد راهنما: دکتر خونساری
دانشکده مهندسی برق و کامپیوتر، دانشگاه تهران

نتایج

با اجرای الگوریتم های مختلف بر روی محیط River Swim به این نتیجه رسیدیم که الگوریتم های مبتنی بر روش های بیزی (خانواده ی PSRL) از سایر الگوریتم ها مجموع خطای کمتری دارند و به سرعت خطای آن ها کاهش داشته و فرم لگاریتمی به خود می گیرند.



در شکل زیر به صورت جزئی تر الگوریتم های PSRL را بررسی کرده ایم.



چکیده

امروزه شرکت ها و دولت های بسیاری سرمایه گذاری زیادی روی هوش مصنوعی انجام داده اند تا به صورت بهینه اهداف خود را محقق کنند. یکی از جنبه های مهم عامل های هوشمند توانایی یادگیری در محیط است تا بتواند با استفاده از تجربیات گذشته تصمیم بهتری در آینده بگیرند. به همین منظور، یادگیری تقویتی با تعامل با محیط بیرونی سعی در بهبود رفتار خود دارد و این رویکرد یادگیری تقویتی را یکی از گرایش های مهم هوش مصنوعی قرار داده است. در این پروژه، به پیاده سازی و مقایسه ی جدیدترین الگوریتم های یادگیری تقویتی می پردازیم، که محیط آن ها را به کمک فرایندهای تصمیم گیری مارکوف می توان مدل سازی کرد. در اینجا یادگیرنده فقط به یک رشته از مشاهدات دسترسی دارد. در این پروژه معیارهای کارایی مثل همگرایی، پیچیدگی نمونه و پشیمانی را اندازه گیری می کنیم. علاوه بر آن نگاهی بر تاثیر ویژگی های فرآیند تصمیم گیری غیرمتعارف مارکوف (مانند قطر) در مقایسه مان می اندازیم. در نهایت الگوریتم ها را از نظر مجموع خطا با یکدیگر مقایسه می کنیم.

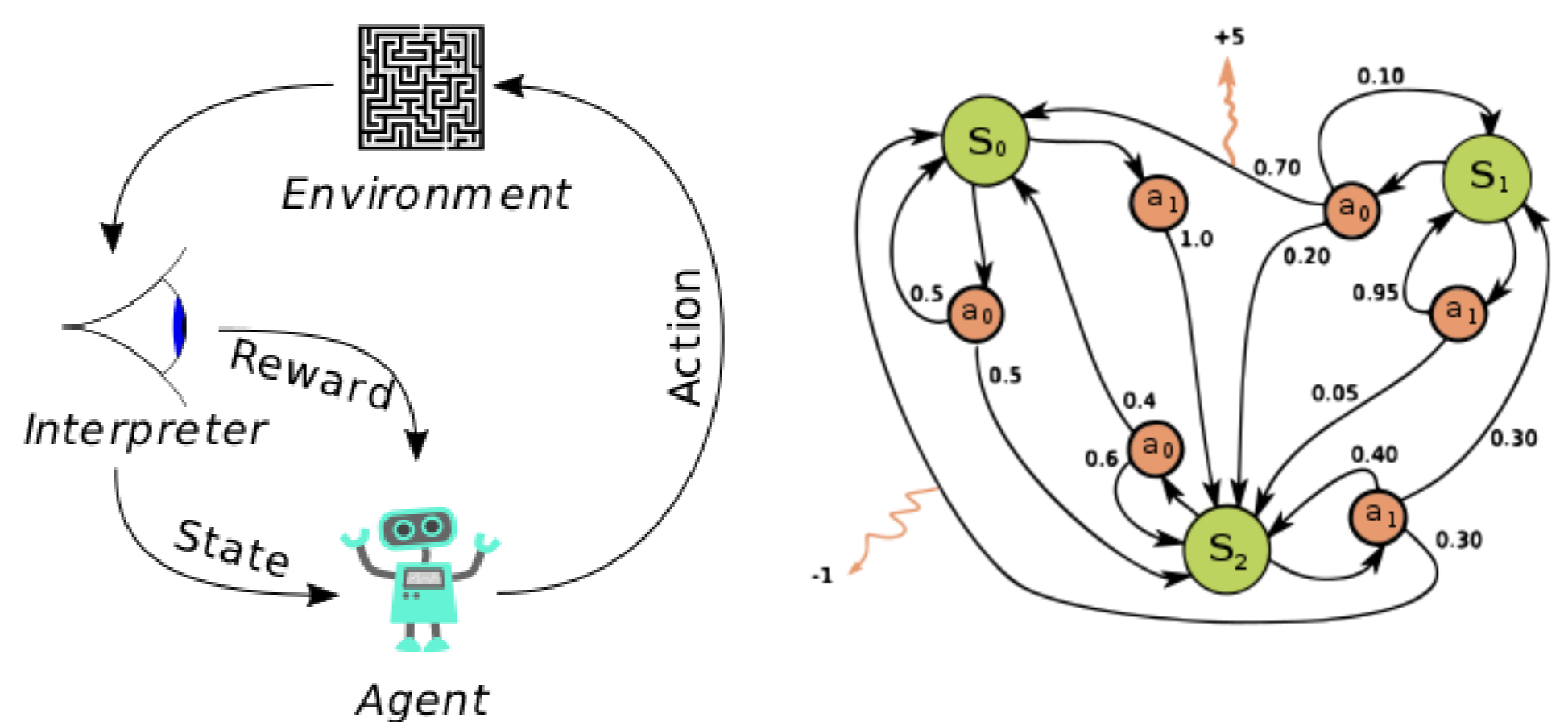
ساختار

مسائل یادگیری تقویتی معمولاً به کمک فرآیند تصمیم گیری مارکوف مدل سازی می شوند. هر مسئله ی فرآیند تصمیم گیری مارکوف به کمک حالات، عمل ها، ماتریس پاداش و ماتریس احتمالات مدل می شود. در این پروژه الگوریتم های مختلف همان به کمک کلاس عامل و خود فرآیند مارکوف به کمک یک کلاس به نام محیط در کد مدل شده اند. به کمک این مدل می توان محیط های مختلف را با الگوریتم های مختلف تست کرد و به نوعی زیرساختی برای توسعه و تست سریع الگوریتم ها در اختیار داریم.

الگوریتم های پیاده سازی شده عبارتند از: UCRL2 ، Epsilon Greedy و انواع مختلف الگوریتم های PSRL

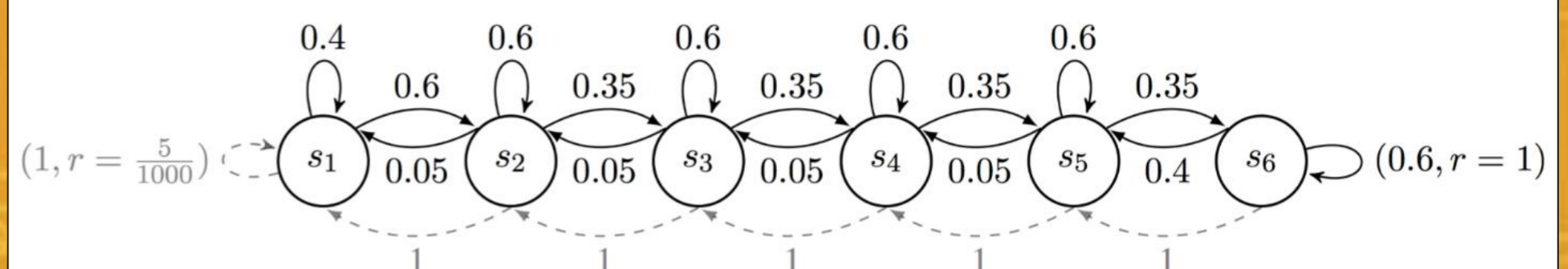
جمع بندی

با توجه به پیشرفت روزافزون یادگیری ماشین و نیاز به تست الگوریتم های مربوط به آن، نیاز به چارچوبی برای پیاده سازی، شبیه سازی الگوریتم ها در یک محیط و مقایسه ی آن ها در قالب نمودار احساس می شود. در این پروژه سعی شده که یک چارچوب برای برطرف کردن نیازهای بالا توسعه داده شود. این راهکار با جدا کردن محیط که در واقع مسئله ی فرآیند تصمیم گیری مارکوف را مدل سازی می کند و جدا کردن عامل ها که در واقع همان الگوریتم های مورد نظر ما هستند، صورت گرفته است. برای تست یک محیط و یک الگوریتم جدید کافی است در قدم اول یک محیط جدید با پارامترهای مارکوف آن مسئله بسازیم و در قدم دوم کلاس عامل را گسترش داده و توابع مربوط به بروزسانی سیاست را پیاده سازی کنیم. در نهایت به کمک توابع مربوط به محیط و عامل را می گیرند و شبیه سازی را انجام می دهند، نتایج را در غالب فایل CSV داریم. در مرحله ی بعد می توانیم خروجی این الگوریتم را با الگوریتم های دیگر در غالب نمودار با یکدیگر مقایسه کنیم.



نمونه ای از فرآیند تصمیم گیری مارکوف مدلی ساده از مسئله ی یادگیری تقویتی

محیط تست این پروژه نیز محیط استاندارد River Swim است که دارای ۶ حالت بوده و نسبتاً قطر طولانی ای دارد.



مراجع اصلی

1. R. Fruit, M. Pirotta, A. Lazaric, R. Ortner, "Efficient Bias-Span-Constrained Exploration-Exploitation in Reinforcement Learning", 2018. [Online]. Available: <https://arxiv.org/abs/1802.04020>
2. A. Tolver, "An Introduction to Markov Chain", 2016, University of Copenhagen, Department of Mathematical Sciences
3. P. L. Bartlett, A. Tewari, "REGAL: A Regularization based Algorithm for Reinforcement Learning in Weakly Communicating MDPs", 2012. [Online]. Available: <https://arxiv.org/abs/1205.2661>