محيط هاي ناشناخته

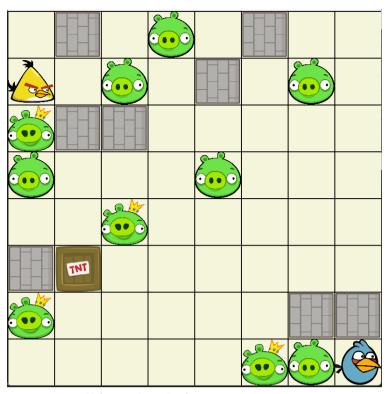
در این بخش باید الگوریتمی را پیاده سازی کنید تا بتواند در محیط ناشناخته فعالیت کند و یک سیاست برای حل مسئله استخراج کند. نام این محیط Unknown Angry Birds است.

1-2-معرفي محيط

یک محیط Grid به ابعاد 8x8 هستیم. در این محیط باید پرنده زرد که در ابتدای بازی در موقعیت (0,0) یعنی بالا-چپ قرار دارد به پرنده آبی برسد که در موقعیت (7,7) یعنی پایین-راست قرار دارد. پرنده زرد دارای چهار کنش بالا، پایین، چپ و راست میباشد. در صورت رسیدن به پرنده آبی *** امتیاز کسب خواهید کرد و بازی خاتمه میبابد. هر کنش پرنده زرد نیز یک امتیاز منفی در پی خواهد داشت.

- در این محیط ۸ **مانع** قرار دارند که عامل نمی تواند از آنها عبور کند. این موانع در ابتدای ایجاد محیط بصورت تصادفی در خانه ها قرار می گیرند.
- در این محیط \wedge خوک وجود دارد که در صورت برخورد پرنده زرد با آنها عامل \wedge ۲۵ امتیاز مثبت دریافت کرده و آنها را از بین میبرد. موقعیت قرارگیری این \wedge خوک در خانهها در ابتدای ایجاد محیط بطور تصادفی مشخص می شود.
- در این محیط ۲ ملکه خوکی قرار دارند که در صورت برخورد عامل با هر ملکه، ۱۹۰۰ امتیاز منفی دریافت میکند. موقعیت این ملکه ها نیز در ابتدای ایجاد محیط بطور تصادفی تعیین می شود.
- در این محیط یک TNT قرار دارد که در صورت برخورد با آن عامل • ۲ امتیاز منفی دریافت کرده و بازی خاتمه می یابد. موقعیت TNT بطور تصادفی در ابتدای تعیین می شود.
- عامل فقط **۱۵۰ کنش** فرصت دارد تا به پرنده آبی برسد. در صورتیکه تعداد کنشهای عامل بیشتر از این تعداد شود، • • • • ۱ امتیاز منفی دریافت کرده و بازی خاتمه می یابد.
- مانند بخش قبلی عامل دارای کنشهای تصادفی است. کنش در نظر گرفته شده برای عامل با احتمال خاصی انجام خواهد شد که به این کنش، کنش اصلی گفته می شود. کنش اصلی دارای دو کنش همسایه است که آنها نیز احتمال خاصی دارند و ممکن است یکی از آنها به جای کنش اصلی انجام شود. توجه کنید که توزیع احتمالاتی کنشها در یک بازی ثابت ولی در بازی مختلف می تواند فرق داشته باشد. کنش اصلی و کنشهای همسایه آن به شرح زیر هستند:

كنش همسايه	كنش همسايه	کنش اصلی
راست	چپ	بالا
راست	چپ	پایین
پایین	بالا	چپ
پایین	بالا	راست



شكل ٢-١ محيط Unknown Angry Birds

۲-۲-مراحل پیادهسازی

- در ابتدا باید با استفاده از الگوریتمهای موجود در محیطهای ناشناخته، یک سیاست در محیط استخراج کنید. انتخاب الگوریتم برعهده شما میباشد. توصیه میشود الگوریتمهای مختلف را پیادهسازی کرده و عملکرد هر کدام را بررسی کنند.
- رو فعالیت به بخش نسبت به بخش قبلی در این است که عامل معمولاً نمی تواند تنها با یک دور (episode) از فعالیت در محیط سیاست بهینه را استخراج کند. به همین دلیل لازم عامل دورهای متعددی از فعالیت در محیط را سپری کند و دانش استخراج شده از آنها را با هم ترکیب کند تا یادگیری تحقق پیدا کند. به عنوان مثال اگر دانش استخراج شده از محیط در قالب یک جدول Q باشد، عامل باید در ابتدای هر دور جدول ذخیره شده قبلی را بازیابی کرده، در حین فعالیت در آن دور آن را بروزرسانی کرده و دوباره در پایان دور آن را ذخیره کند، تا در نهایت به یک سیاست خاص همگرا شود.
- مال باید در حین اجرای الگوریتم اطلاعاتی را استخراج کنید تا بتوانیم همگرایی الگوریتم شما را بررسی کنیم. بعد از هر episode که جدول *Q بروزرسانی شد، معیار زیر را محاسبه کنید. نام آن Value Difference می باشد.

Value Diffrence =
$$\sum_{s \in S} \sum_{a \in A} |Q^{(k+1)}(s,a) - Q^{(k)}(s,a)|$$

episode کنید که این فرمول بر روی جدول Q^* قبل از بروزرسانی (اشاره به اندیس Q^*) و پس از بروزرسانی در یک Q^* (اشاره به اندیس Q^*) را در نظر می گیرد. شما در هر episode از اجرای الگوریتم باید این مقدار را محاسبه کرده و ذخیره کنید و در نهایت نمودار تغییرات آنرا رسم کنید. توجه کنید در صورت همگرایی الگوریتم شما، این معیار باید به مرور کاهش یابد. زمانیکه Value Difference کمتر از یک مقدار مثل Q^* بشود به معنای آن است که می توانیم اجرای الگوریتم را متوقف کنیم. . مقدار Q^* عدد بسیار کوچک (به عنوان مثال Q^* 0.001 است.

- در راستای ارزیابی سیاست استخراج شده نیاز است که روی یک نقشه، کنشی که سیاست برای آن حالت پیشنهاد داده است را نشان دهید. یک نقشه به ابعاد 8x8 بسازید و روی هر حالت کنش پیشنهادی سیاست را نمایش دهید.
 - پس از همگرایی الگوریتم، با توجه به سیاست استخراج شده در ماتریس * عامل باید به فعالیت در محیط بپردازد.

۲-۳-ارزیابی

در ابتدا الگوریتم شما مورد بررسی قرار گرفته و صحت و شیوه پیادهسازی آن ارزیابی خواهد شد. سپس همگرایی الگوریتم شما با توجه به نموداری که در بخش ۲-۲ معرفی شد بررسی خواهد شد.

پس از اجرای الگوریتم و بررسی همگرایی آن، عملکرد عامل برای فعالیت در محیط بر اساس سیاست استخراج شده مورد ارزیابی قرار می گیرد. مانند بخش قبلی عامل در تعدادی محیط و در هر محیط ۵ دور فعالیت خواهد کرد. هر دور از فعالیت عامل با رسیدن به یکی از شرایط پایان بازی متوقف شود و که میانگین امتیازاتش در همه دورها محاسبه می شود. در نهایت نمره شما براساس بر آیند امتیازاتی که عامل در محیطهای مختلف کسب کرده است، مطابق با جدول زیر مشخص خواهد شد.

نمره	میانگین امتیاز دریافتی	
100 %	mean score > 1100	
70 %	mean score > 850	
50 %	mean score > 600	
0 %	برخورد با TNT (در صورتیکه سیاست استخراج شده عامل را به سمت TNT هدایت کند)	

۲-4-بخش امتیازی

پیادهسازی یک شبکه عصبی مبتنی بر Deep QLearning با استفاده از کتابخانههای Pytorch یا TensorFlow. ورودی شبکه عصبی باید دارای ۴ خروجی باشد(به تعداد کنشها). عامل برای تصمیم گیری در انتخاب کنش در هر حالت باید از این شبکه استفاده کند.

موفق باشيد