

پیش گزارش آزمایش یازدهم (پاده سازی الگوریتم DQN)

امیرحسین احمدی آشتیانی ۹۹۲۳۵۰۱

محمدرضا امیری ۹۹۲۶۰۴۰

محمد مهدی نوروزی ۹۹۲۳۰۸۵

تفاوت یادگیری تقویتی با سایر روش ها:

تفاوتها

- نوع داده مورد استفاده:
 - یادگیری با نظارت: داده ها شامل ورودی ها و برچسب های خروجی مشخص هستند. مدل بر اساس ارتباط بین این ورودی ها و خروجی ها آموزش داده می شود.
 - یادگیری بدون نظارت: تنها ورودی ها در دسترس هستند، و مدل باید ساختار یا الگوهای موجود در داده را بدون برچسب کشف کند.
 - یادگیری تقویتی: از تعامل مدل با محیط استفاده می شود. مدل از طریق دریافت پاداش یا تنبیه آموزش می بیند تا استراتژی بهینه را برای دستیابی به هدف پیدا کند.
- تعامل با محیط:
 - یادگیری با نظارت و بدون نظارت: داده های آموزشی معمولاً ثابت و از پیش آماده شده هستند.
 - یادگیری تقویتی: مدل باید به صورت فعال با محیط تعامل کند و از بازخورد (پاداش ها یا تنبیهات) ناشی از تصمیمات خود بیاموزد.
- هدف آموزش:
 - یادگیری با نظارت: هدف، پیش بینی خروجی های دقیق برای داده های جدید است.
 - یادگیری بدون نظارت: کشف خوشه ها یا کاهش ابعاد داده است.
 - یادگیری تقویتی: یافتن استراتژی بهینه برای رسیدن به حداکثر پاداش جمعی در طول زمان.
- زمان بندی بازخورد:
 - یادگیری با نظارت: بازخورد (خطا) فوری و مشخص است.
 - یادگیری تقویتی: بازخورد معمولاً تأخیری است و وابسته به اقدامات زنجیره ای در محیط است.
- کاربردها:
 - یادگیری با نظارت: طبقه بندی و رگرسیون، مانند تشخیص تصویر یا پیش بینی قیمت.
 - یادگیری بدون نظارت: خوشه بندی داده ها یا سیستم های پیشنهاددهنده.
 - یادگیری تقویتی: بازی های کامپیوتری، رباتیک، کنترل ترافیک، و حل مسائل پیچیده مانند Go و شطرنج.

دلایل ادغام یادگیری عمیق با یادگیری تقویتی

- توانایی یادگیری ویژگی ها (Feature Learning): یادگیری عمیق قادر است ویژگی های پیچیده و غیرخطی را از داده ها استخراج کند، که برای یادگیری سیاست های پیچیده در یادگیری تقویتی مفید است. به ویژه در مواردی که ورودی داده ها فرم خام مانند تصاویر یا ویدیو است.
- مدیریت ابعاد بزرگ حالت (State Space): در یادگیری تقویتی سنتی، حل مسائل با ابعاد بزرگ و پیوسته بسیار سخت است. مدل های یادگیری عمیق مانند شبکه های کانولوشنی (CNNs) یا شبکه های بازگشتی (RNNs)، به کاهش این پیچیدگی کمک می کنند.
- توانایی عمومی سازی: شبکه های عمیق در یافتن الگوها و عمومی سازی بهتر داده ها توانمند هستند، که در محیط های پویا و ناشناخته مزیت زیادی دارد.

4. پیشرفت محاسباتی:

پیشرفت در سخت افزار (مانند GPU ها) و الگوریتم های بهینه سازی پیشرفته مانند Adam و RMSprop، امکان آموزش مدل های عمیق برای کاربردهای یادگیری تقویتی را فراهم کرده اند.

5. حل مسائل بسیار پیچیده:

یادگیری عمیق این امکان را فراهم می کند که یادگیری تقویتی در محیط هایی مانند بازی ها، رباتیک، و کنترل خودران ها موثرتر عمل کند.

دلایل اهمیت بافر Replay:

1. رفع مسئله همبستگی (Correlation) بین داده ها:

در یادگیری تقویتی، نمونه های متوالی حاصل از تعامل با محیط معمولاً به شدت به یکدیگر وابسته هستند. این وابستگی می تواند فرآیند یادگیری را مختل کند، زیرا فرض معمول در الگوریتم های بهینه سازی این است که نمونه های داده مستقل و به طور یکنواخت توزیع شده اند.

- Replay Buffer با ذخیره تجربیات و نمونه برداری تصادفی از آنها، همبستگی داده ها را کاهش داده و فرآیند یادگیری را پایدارتر می کند.

2. استفاده موثرتر از داده ها:

در روش های کلاسیک یادگیری تقویتی، هر تجربه تنها یک بار استفاده می شود و پس از آن دور ریخته می شود. Replay Buffer امکان استفاده چندباره از تجربیات گذشته را فراهم می کند، که منجر به کارایی بیشتر در استفاده از داده های تعامل با محیط می شود.

3. پایداری در فرآیند آموزش:

یادگیری مستقیم از داده های متوالی می تواند باعث نوسانات شدید در پارامترهای شبکه عصبی شود. Replay Buffer با نمونه گیری تصادفی از تجربیات متنوع، این نوسانات را کاهش داده و باعث پایداری فرآیند آموزش می شود.

4. آموزش یکنواخت تر و بهتر:

بافر، تجربیات مختلف از زمان های مختلف را در خود ذخیره می کند، و این تنوع داده ها باعث می شود مدل به جای تمرکز بیش از حد بر تجربیات اخیر، سیاست بهینه تر و عمومی تری را بیاموزد.

نحوه عملکرد Replay Buffer:

- ذخیره تجربیات:

هر تعامل با محیط به صورت یک چهارگانه (s, a, r, s') شامل وضعیت کنونی s اقدام انجام شده a پاداش دریافت شده r و وضعیت بعدی s' در بافر ذخیره می شود.

- نمونه برداری تصادفی:

برای به روزرسانی شبکه Q به جای استفاده از داده های متوالی، تعدادی از تجربیات به صورت تصادفی از بافر انتخاب شده و استفاده می شوند.

مثال عملی:

در بازی هایی مانند **Atari**، تجربه هایی که بازی در ابتدا از محیط می گیرد، حاوی رفتارهای غیرتصادفی (و گاهی ناکارآمد) هستند. Replay Buffer این تجربیات را ذخیره کرده و به تدریج تجربیات بهینه تر را جایگزین تجربیات قدیمی تر می کند، که منجر به یادگیری بهتر می شود.

تعریف شبکه هدف:

در DQN، دو شبکه عصبی جداگانه استفاده می‌شود:

1. شبکه اصلی (Online Network):

شبکه‌ای که در هر مرحله از یادگیری به‌روزرسانی می‌شود و تصمیم‌گیری‌های مدل را بر اساس $Q(s, a; \theta)$ انجام می‌دهد.

2. شبکه هدف (Target Network):

نسخه‌ای از شبکه اصلی است که با وزن‌های ثابت (θ) برای مدت مشخصی ثابت نگه داشته می‌شود. این شبکه برای محاسبه مقدار هدف Q-value در به‌روزرسانی استفاده می‌شود.

اهمیت شبکه هدف در یادگیری:

1. کاهش نوسانات در یادگیری:

اگر شبکه اصلی هم برای محاسبه مقدار هدف و هم برای پیش‌بینی مقدار Q استفاده شود، به‌روزرسانی‌های وابسته به خود باعث ناپایداری در یادگیری می‌شود. شبکه هدف با ثابت نگه‌داشتن مقادیر هدف برای مدتی، این نوسانات را کاهش می‌دهد.

2. جلوگیری از انفجار مقادیر Q:

بدون شبکه هدف، مقادیر Q ممکن است به دلیل حلقه‌های بازخوردی (Feedback Loop) به شدت افزایش یافته و باعث یادگیری ناکارآمد شود.

3. پایداری در همگرایی:

شبکه هدف فرآیند همگرایی به سیاست بهینه را روان‌تر می‌کند، زیرا تغییرات تدریجی در مقادیر هدف باعث تنظیم بهتر وزن‌ها می‌شود.

4. تثبیت یادگیری در محیط‌های پویا:

در محیط‌هایی که پاداش‌ها یا حالات پیچیده‌اند، شبکه هدف تغییرات سریع را مدیریت کرده و به سیستم فرصت یادگیری بهتر می‌دهد.

توضیح کامل الگوریتم Deep Deterministic Polisy Gradient یا DDPG

الگوریتم DDPG یک روش یادگیری تقویتی (Reinforcement Learning) در دسته‌ی کنترل پیوسته است که بر مبنای دو شبکه عصبی (Actor و Critic) عمل می‌کند. این الگوریتم ترکیبی از DQN (Deep Q-Network) و Policy Gradient است که برای محیط‌های با فضای عمل پیوسته طراحی شده است.

مؤلفه‌های اصلی DDPG

1. شبکه Actor و Actor Target:

- شبکه Actor وظیفه تولید عمل (Action) را برای هر حالت (State) دارد.
- خروجی شبکه Actor، مقدار عمل در فضای پیوسته است که با استفاده از توابع فعال‌سازی مانند Tanh محدود می‌شود.

- شبکه Actor Target، نسخه کپی شده‌ای از Actor است و برای ثبات در به‌روزرسانی‌ها استفاده می‌شود.

2. شبکه Critic و Critic Target:

- شبکه Critic مقدار Q-value را برای جفت (State, Action) محاسبه می‌کند.
- هدف Critic ارزیابی کیفیت عمل پیشنهادی توسط Actor است.
- شبکه Critic Target مشابه Actor Target، نسخه‌ای پایدار از Critic است.

3. Replay Buffer:

- حافظه‌ای برای ذخیره تجربیات (Transitionها) شامل (State, Action, Reward, Next State, Done).
- با نمونه‌برداری تصادفی از این حافظه، مشکلاتی نظیر همبستگی داده‌ها و تغییرات نامنظم در داده‌های ورودی به شبکه کاهش می‌یابد.

4. نوین Ornstein-Uhlenbeck:

- برای ایجاد تصادفی‌بودن در انتخاب عمل و کمک به اکتشاف استفاده می‌شود.
- این نوین به‌طور خاص برای محیط‌های با فضای عمل پیوسته طراحی شده است.