



طراحی یک سامانه گفتگوگر وظیفه‌گرا مشترک مبتنی بر یادگیری تقویتی عمیق

محمدجواد نصری‌لوشانی^۱، جواد سلیمی‌سرتختی^۲ و حسین ابراهیم‌پور کومله^۲

^۱ دانشجوی کارشناسی ارشد، گروه مهندسی کامپیوتر، دانشکده مهندسی برق و کامپیوتر، دانشگاه کاشان، کاشان، ایران
mohammad.j.nasri@gmail.com

^۲ استادیار، گروه مهندسی کامپیوتر، دانشکده مهندسی برق و کامپیوتر، دانشگاه کاشان، کاشان، ایران
{salimi, ebrahimpour}@kashanu.ac.ir

رسانی به مشتری است که یک عامل مجازی جایگزین انسان می‌شود و هزینه‌های مرکز تماس و خدمات پشتیبانی را کاهش می‌دهد.

سامانه‌های گفتگوگر را می‌توان به دو دسته سامانه‌های وظیفه‌گرا و سامانه‌های غیر وظیفه‌گرا، دسته‌بندی کرد. هدف سامانه‌های وظیفه‌گرا، کمک به کاربر برای کامل کردن وظایف معینی می‌باشد؛ برای مثال تهیه بلیط هواپیما و یا بکارگیری در رستوران. از سوی دیگر، سامانه‌های غیر وظیفه‌گرا روی گفتگو با انسان در حوزه‌های باز تمرکز دارند تا واکنش‌های صحیحی برای گفتگو با انسان ایجاد کنند؛ همانند یک ربات چت.

نحوه عملکرد سامانه‌های گفتگوگر وظیفه‌گرا به این صورت می‌باشد که این سامانه‌ها ابتدا پیام ارسالی انسان را درک می‌کنند (بخش فهم زبان طبیعی) و سپس آن را به عنوان یک حالت داخلی بیان می‌کنند و با توجه به سیاست و حالت گفتگو، کنشی اتخاذ می‌کنند (بخش مدیریت گفتگو) و در نهایت این کنش به شکل یک زبان طبیعی در آورده می‌شود (بخش تولید زبان طبیعی) و به انسان نمایش داده می‌شود.

گروهی از روش‌ها برای ساخت بخش‌های مختلف سامانه‌های گفتگوگر وظیفه‌گرا، روش‌های یادگیری ماشین هستند که به سه دسته یادگیری نظارت‌شده، یادگیری بدون نظارت و یادگیری تقویتی، تقسیم می‌شوند. سامانه‌های گفتگوگر نیاز به پیکره‌های بزرگ برای ساختن مدل‌های کارآمد دارند. از طرفی، روش‌های یادگیری تقویتی می‌توانند با استفاده از پیکره‌های کوچک شروع به یادگیری کنند و با گذشت زمان و تعامل با محیط، یادگیری انجام دهند و عملکرد خود را بهبود ببخشند. همچنین، روش‌های

چکیده - ساخت سامانه‌های گفتگوگر در سال‌های اخیر توجه زیادی به خود جلب کرده است. دسته‌ای از این سامانه‌ها، سامانه‌های گفتگوگر وظیفه‌گرا هستند که هدفشان رساندن انسان به مقصودش با انجام گفتگو در یک حوزه خاص می‌باشد؛ مثلاً رستوران. این سامانه‌ها از بخش‌های مختلفی تشکیل می‌شوند که اگر دو یا چند بخش، همزمان توسعه داده شوند، سامانه‌ی مشترک (Joint) نامیده می‌شود. یکی از روش‌هایی که برای توسعه این سامانه‌ها استفاده می‌شود، روش یادگیری تقویتی عمیق است. در یادگیری تقویتی عمیق، عامل که شبکه عصبی است با تعامل با محیط (کنش) در حالت‌های مختلف و دریافت پاداش از آن، آموزش می‌بیند. همچنین، در شروع یادگیری، عامل تعدادی کنش بصورت تصادفی انجام می‌دهد و به مرور زمان از دانشی که بدست آورده، استفاده می‌کند. در این مقاله، برای اینکه حالت مناسبی از محیط گفتگو ایجاد شود، از چسباندن نمایش جمله آخرین پیام ربات و انسان، استفاده شده است. همچنین، تابعی جدید برای کاهش احتمال انجام کنش تصادفی، بکارگرفته شده است. برای ارزیابی و مقایسه عملکرد روش ارائه شده با دو سامانه گفتگوگر دیگر، از شبیه‌ساز گفتگو در حوزه رستوران استفاده شده است. روش ارائه شده، پیشینه پاداش ۰/۲۹۹۳۷ را در ۲۷،۹۰۰ گام گفتگو بدست می‌آورد که نسبت به دو روش دیگر، با تعداد گفتگوهای کمتر، پاداش بیشتری بدست آورده است.

کلمات کلیدی- نمایش جمله، حوزه رستوران، مدیریت گفتگو، DQL

۱. مقدمه

یک سامانه گفتگوگر^۱ یا دستیار مجازی، قادر به گفتگوی معنادار مبتنی بر متن (یا صوت) با کاربران انسانی است. یکی از اصلی‌ترین کاربردهای چنین سامانه‌هایی، محیط‌های خدمات

^۱ Dialogue System

۲-۲ یادگیری تقویتی و یادگیری Q

برخلاف روش‌های طبقه‌بندی که مطرح شد، در یادگیری تقویتی خبری از نمونه‌های آموزشی برای یادگیری وجود ندارد و عامل (یادگیر تقویتی) با استفاده از تعامل با محیط، یعنی انجام کنش a و دریافت پاداش r ، یادگیری انجام می‌دهد. حالت s ، موقعیت کنونی عامل در محیط است. کنش‌های مجاز، مجموعه کنش‌هایی است که عامل می‌تواند در حالت کنونی اتخاذ کند؛ و بازخورد محیط به کنش سیستم، در قالب پاداش مطرح می‌شود. سیاست نیز یعنی به ازای هر حالت، چه کنشی اتخاذ شود.

یکی از روش‌های معروف در یادگیری تقویتی، یادگیری Q است که نگاهی از زوج حالت و کنش را به مقادیری که ارزش Q نامیده می‌شوند، می‌برد. این مقادیر تحت عنوان یک جدول بنام جدول Q مطرح شود که می‌توان با استفاده از آن، سیاست مشخصی را برای انجام کنش‌های مطلوب در حالت‌های مختلف محیط پیدا کرد. مقادیر این جدول در فرآیند یادگیری با استفاده از رابطه (۱) بروز می‌شوند.

$$Q(s_t, a_t) \leftarrow Q(s_t, a_t) + \alpha(R_{t+1} + \gamma \max_{a'} Q(s_{t+1}, a') - Q(s_t, a_t)) \quad (1)$$

که t بیانگر لحظه a' کنش بدست آمده از بیشترین مقدار Q ، R پاداش فوری، α ضریب یادگیری و γ ضریب تخفیف می‌باشند.

۲-۳ سیاست ϵ -حریصانه^۲

در یادگیری تقویتی، عامل برای شروع تعامل با محیط، ابتدا در تعدادی گام بصورت تصادفی عمل می‌کند؛ چون از پیش دانشی ندارد تا بر اساس آن تصمیم بگیرد. در سیاست ϵ -حریصانه، عامل با احتمال از پیش تعیین شده ϵ ، کنش تصادفی انجام می‌دهد؛ به این صورت که اگر ϵ مقدار بیشتر از یک عدد تصادفی داشته باشد. در غیر این صورت، عامل بر اساس $Q(s, a)$ یاد گرفته شده، تصمیم می‌گیرد. مقدار ϵ بر اساس توابعی تنظیم می‌شود که مقدار آن در ابتدا بالا و با گذر گام‌های یادگیری، کاهش می‌یابد.

۲-۴ یادگیری عمیق

شبکه عصبی یکی از انواع روش‌های یادگیری ماشین است که از سه لایه ورودی، پنهان و خروجی تشکیل می‌شود که وقتی تعداد لایه‌های پنهان افزایش پیدا کند، آنرا یادگیری عمیق می‌نامند که از این روش‌ها بطور گسترده در حوزه متن استفاده می‌شوند.

یادگیری ماشین، نیاز به استخراج ویژگی از نمونه‌های آموزشی دارند اما روش‌های یادگیری عمیق می‌توانند خودشان با استفاده از ورودی که به آن‌ها داده می‌شود، استخراج ویژگی کنند.

رویکردهای زیادی برای ساخت سامانه‌های گفتگوگر وظیفه‌گرا ارائه شده است. در رویکردهای قدیمی‌تر، این سامانه‌ها در یک ساختار کلی به نحوی ساخته می‌شدند که هر بخش جداگانه طراحی، آموزش و ارزیابی می‌شد. طراحی جداگانه هر بخش، منجر به تکثیر خطا می‌شود؛ زیرا عملکرد هر بخش به خروجی بخش قبلی بستگی دارد [۱]. در سال‌های اخیر، بیشتر محققان روی رویکردهای قابل آموزش مشترک^۲ کار کرده‌اند که دو یا چند بخش، ادغام شده و بصورت مشترک چند وظیفه را انجام می‌دهند.

در این مقاله، یک سامانه گفتگوگر وظیفه‌گرا مشترک مبتنی بر یادگیری تقویتی عمیق، طراحی شده است که انسان و ربات بتوانند گفتگویی از نوع متن داشته باشند و پس از اتمام گفتگو، انسان به مقصود خود برسد.

ساختار مقاله در ادامه به این صورت تدوین شده است: در بخش دوم مفاهیم پایه توضیح داده می‌شود. در بخش سوم، کارهای مرتبط در سامانه‌های گفتگوگر مشترک بررسی می‌شوند. در بخش چهارم، روش پیشنهادی ما برای ساخت این سامانه‌ها ارائه می‌شود. در بخش پنجم نیز، عملکرد روش پیشنهادی ارزیابی می‌شود و در بخش ششم، نتیجه‌گیری بیان می‌شود.

۲. مفاهیم پایه

در این بخش مفاهیم پایه مطرح می‌شوند. ابتدا طبقه‌بندی و یادگیری تقویتی، بیان می‌شوند. در ادامه، یادگیری عمیق توضیح داده می‌شود. سپس روش‌های نمایش کلمه و نمایش جمله معرفی می‌شوند. بخش‌های مختلف سامانه‌های گفتگوگر وظیفه‌گرا نیز در همین بخش بیان می‌شوند.

۲-۱ طبقه‌بندی و الگوریتم بیز ساده^۳

یک دسته از روش‌های یادگیری ماشین، روش‌های طبقه‌بندی هستند که بر اساس برجسب نمونه‌های آموزشی، یاد می‌گیرند تا با استفاده از ویژگی‌های نمونه‌ها، آنها را طبقه‌بندی کنند. روش‌های گوناگونی برای طبقه‌بندی ارائه شده که یکی از این روش‌ها، الگوریتم بیز ساده است که کاربرد بسیاری در پردازش زبان‌های طبیعی دارد. در این روش، از قضیه بیز و فرض استقلال بین متغیرها برای طبقه‌بندی استفاده می‌شود.

² Jointly Trainable Models

³ Naïve Bayes

⁴ Epsilon Greedy Policy



۲-۵- یادگیری تقویتی عمیق

این روش همانطور که از نامش پیداست، دو روش یادگیری تقویتی و یادگیری عمیق را ترکیب می‌کند. در یادگیری تقویتی عمیق، بخش عامل یادگیری تقویتی با استفاده از یک شبکه عصبی، ساخته می‌شود. در یادگیری Q ، بجای جدول Q از یک شبکه عصبی برای یادگیری سیاست‌ها استفاده شود، که به این روش، یادگیری Q عمیق یا DQL گفته می‌شود. ورودی این شبکه عصبی حالت محیط و خروجی آن ارزش Q است.

۲-۶- نمایش کلمه *Glove*

روش‌های نمایش کلمه، کلمات را از یک لغت‌نامه به بردارهای عددی در یک فضای پیوسته، نگاشت می‌کنند. این روش‌ها برای کلمات مشابه، بردارهای مشابه‌ای تولید می‌کنند که از این بردارها برای نمایش کلمات استفاده می‌شود. روش‌های متفاوتی برای نمایش کلمه وجود دارد که یک از آنها روش *Glove* است. این روش یک روش یادگیری بدون نظارت می‌باشد که یک ماتریس از شمارش کلمات در اسناد ساخته و سپس با فاکتورگیری، ماتریسی با ابعاد کمتر می‌سازد، که هر سطر، بردار نمایش یک کلمه می‌باشد.

۲-۷- نمایش جمله *MPNet*

برای نمایش جمله، می‌توان از میانگین گرفتن بردار نمایش کلمه کلمات آن جمله استفاده کرد. تولید نمایش جمله با استفاده از این روش، نتایج مطلوبی نمی‌دهد؛ چرا که این روش‌ها برای نمایش جمله ارائه نشده‌اند. از اینرو روش‌هایی پدید آمدند که جملات را به یک بردار عددی برای نمایش آنها، در فضای پیوسته نگاشت می‌کنند و نمایشی برای جملات می‌سازند. بر طبق آزمایش روی پیکره‌های بزرگ، یکی از برترین روش‌های نمایش جمله^۵ روش *MPNet* می‌باشد که هم از تکنیک *MLM* و هم از تکنیک *PLM* استفاده می‌کند [۲]. در *MLM*، وظیفه مدل، تولید جمله ورودی در خروجی است؛ به این صورت که تعدادی از کلمات ورودی به مدل، ماسک شده‌اند. در *PLM*، جمله ورودی در خروجی ساخته می‌شود؛ به این صورت که با استفاده از جایگشت، جای تعدادی از کلمات ورودی به مدل، عوض شده‌اند.

۲-۸- فهم زبان طبیعی

مدل‌های فهم زبان طبیعی، اطلاعات گفتگو را برای بخش مدیریت گفتگو فراهم می‌کنند. این اطلاعات شامل شناسایی دامنه، شناسایی مقصود و شناسایی اسلات‌ها یا تشخیص مفهوم می‌باشد. برای حل مسائل تشخیص دامنه و تشخیص مقصود، می‌توان از طبقه‌بندها استفاده کرد.

۲-۹- مدیریت گفتگو

بیشتر مدل‌های مدیریت گفتگو تلاش می‌کنند تا گفتگو را با حالت‌های مشاهده‌شده در مدل، مدیریت کنند. این روش‌ها به دو بخش ردیابی حالت و تولید کنش یا سیاست، تقسیم می‌شوند. ردیابی حالت، حالت گفتگوی کنونی را با توجه به تاریخچه گفتگو، بیان می‌کند. بخش تولید کنش، ارتباطی روی حالت‌ها و کنش‌های سیستم می‌سازد که یادگیری تقویتی برای این امر بطور گسترده استفاده می‌شود.

۲-۱۰- تولید زبان طبیعی

مدل‌های تولید زبان طبیعی، کنش انتخابی در بخش مدیریت گفتگو را تبدیل به یک زبان طبیعی می‌کنند. خروجی این مدل‌ها در قالب متن یا صوت، به انسان نمایش داده می‌شود.

۳. مرور کارهای مرتبط

در زمینه سامانه‌های گفتگوگر وظیفه‌گرا مشترک، روش‌های زیادی ارائه شده است. در بخش فهم زبان طبیعی، با استفاده از شبکه‌های عصبی بازگشتی، شناسایی مقصود و تشخیص مفهوم بصورت مشترک [۳] و هر سه وظیفه فهم زبان طبیعی بصورت مشترک [۴] را انجام می‌دهند. در بخش مدیریت گفتگو، با استفاده از یادگیری تقویتی، روشی ارائه شده است که هر دو وظیفه مدیریت گفتگو را بصورت مشترک انجام می‌دهد [۵]. در این روش، یک سامانه پیشنهاددهنده ساخته شده است که با استفاده از گفتگو، یک محصول به کاربر پیشنهاد می‌دهد.

در بخش‌های فهم زبان طبیعی و مدیریت گفتگو بصورت مشترک، روش *SimpleDS* طراحی شده است که برداری جهت نمایش پیام انسان و پیام ربات با استفاده از نمایش کلمه *Glove* می‌سازد و سپس این بردار را به عنوان ورودی به یک شبکه Q عمیق می‌دهد تا یادگیری انجام شود [۶]. در این روش از کاهش ϵ با تابع خطی در سیاست ϵ -حریصانه استفاده شده و همچنین از شبیه‌ساز کاربر بجای انسان استفاده شده و برای تعیین کنش‌های مجاز یادگیر تقویتی (یا ربات) از بیز ساده استفاده کرده‌اند. از طرفی، این روش قابلیت گفتگوی صوتی نیز داشت و در [۷]، این سامانه برای استفاده در چند دامنه بصورت همزمان توسعه داده شده است. روش *SCGSimpleDS* توسعه‌ای بر روش *SimpleDS* می‌باشد که با استفاده از شبکه‌های خود رمزگذار، بردار نمایش پیام انسان و پیام ربات را فشرده کرده و با استفاده از تشخیص احساسات گفتگوی انسان و خوشه‌بندی حالات محیط با استفاده از روش *K-Means*، بردار ورودی به شبکه تقویتی عمیق را، بازسازی کرده‌اند [۸]. دلیل تغییر حالت محیط در این کار، تولید حالت‌های بهتر از

⁵ www.sbert.net/docs/pretrained_models.html

از پیش تعیین‌شده، پر می‌کند و تولید گفتگو می‌کند. برای آموزش یادگیر تقویتی، نیاز به گام‌های زیاد و نمونه‌های آموزشی می‌باشد که استفاده از شبیه‌ساز، می‌تواند این مسئله را بخوبی پوشش دهد.

۴-۲- حالت

برای ساخت حالت، از برداری استفاده می‌شود که از چسباندن آخرین بردار نمایش پیام ربات و پیام شبیه‌ساز که هر کدام از طریق نمایش جمله MPNet بدست آمده‌اند، ساخته می‌شود.

۴-۳- پرونده نمایش

از گفتگوهایی که در نمونه‌های آموزشی وجود دارد، یک پرونده نمایش ساخته خواهد شد که از آن برای آموزش طبقه‌بند بیز ساده استفاده می‌شود. در ساخت پرونده نمایش، از روی گفتگوها، بردار حالت آنها ساخته می‌شود و سپس کنش انتخابی ربات به آخر آن چسبانده می‌شود؛ از کنش ربات برای برچسب طبقه‌بند استفاده می‌شود.

۴-۴- پاداش

در این سامانه گفتگوگر، تابع پاداش از رابطه (۲) بدست می‌آید که از سه بخش تشکیل شده است. بخش Cr ، از اسلات‌های تاییدشده توسط کنش کنونی، تقسیم بر تعداد کل اسلات‌هایی که باید پر شود، بدست می‌آید. بخش Dr ، مقدار شباهت حالت کنونی به حالت‌های موجود در پرونده نمایش است؛ که توسط طبقه‌بند بیز ساده بدست می‌آید. بخش DL نیز مشوق گفتگوهای کارآمدتر است.

$$R = (Cr \times \alpha) + (Dr \times (1 - \alpha)) - DL \quad (2)$$

۴-۵- مقدار ϵ

در این سامانه گفتگوگر، دو روش پیاده‌سازی شده است که روش اول از تابع خطی برای کاهش مقدار ϵ استفاده می‌کند که SenSimpleDS نامیده شده و روش دوم از تابعی دیگر که در رابطه (۳) آمده است، برای کاهش مقدار ϵ استفاده می‌کند که SenSimpleDS+ نامیده شده است.

$$\epsilon_{new} = \max(\epsilon_{min}, (\epsilon \times \beta)) \quad (3)$$

که ϵ_{new} مقدار جدید و ϵ_{min} کمترین مقدار برای ϵ ؛ β ضریب کاهنده و ϵ مقدار فعلی ϵ هستند.

۵. آزمایش‌ها

در این بخش، عملکرد سامانه گفتگوگر ارائه‌شده، ارزیابی و با دو روش SimpleDS و SCGSimpleDS مقایسه می‌شود. برای این

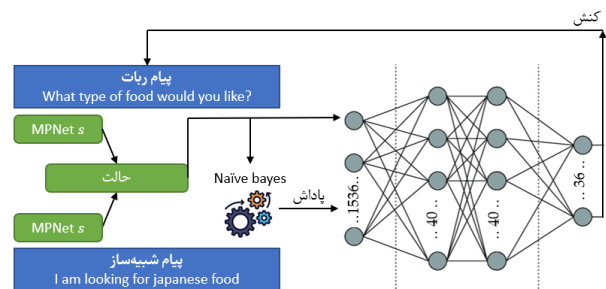
محیط می‌باشد. همچنین در این روش، با استفاده از شبکه‌های مولد تخصصی^۶، برای آموزش بیز ساده داده‌های بیشتری تولید کرده‌اند. در روشی دیگر که مبتنی بر یادگیری تقویتی عمیق برای فهم زبان طبیعی و مدیریت گفتگو می‌باشد [۹]؛ شبکه یادگیری تقویتی ابتدا بروی یک حوزه آموزش دیده و سپس روی حوزه دیگری با وزن‌های مدل قبلی آموزش می‌بیند. این روش از ابتدا، مجموعه کنش‌های دو حوزه را برای هر دو مدل در نظر می‌گیرد.

در بخش تولید زبان طبیعی بصورت مشترک، با استفاده از شبکه‌های عصبی بازگشتی، مدلی طراحی شده است که می‌تواند هر دو وظیفه تولید زبان را انجام دهد [۱۰]. در این روش، هم از یک روش رتبه‌دهنده و هم از یک روش تولید دنباله برای تولید جملات استفاده شده است.

برخی روش‌ها نیز ارائه شده‌اند که تمامی بخش‌های یک سامانه گفتگوگر را در یک مدل، به انجام رسانده باشند. در [۱] با استفاده از چند شبکه عصبی بازگشتی سلسه مراتبی، روشی برای این کار ارائه شده است که در آن برای مدیریت گفتگو از یادگیری تقویتی استفاده شده است.

۴. روش پیشنهادی

روشی که برای ساخت سامانه گفتگوگر وظیفه‌گرا مشترک ارائه شده، مبتنی بر یادگیری تقویتی عمیق می‌باشد که برای ساخت حالت از نمایش جمله MPNet در آن استفاده شده است که معماری آن در شکل ۱ قابل مشاهده است. برای آموزش یادگیری تقویتی عمیق، از شبیه‌ساز کاربر استفاده شده که نقش انسان را برای تولید گفتگو ایفا می‌کند. در ادامه، بخش‌های مختلف سامانه بیان خواهد شد.

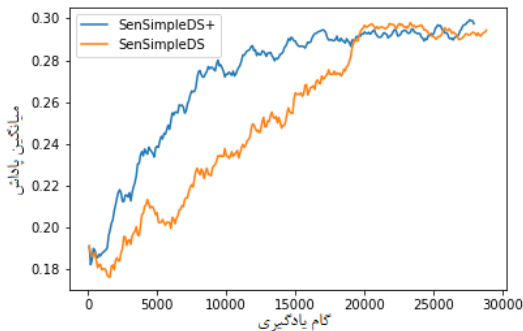


شکل ۱: معماری سامانه گفتگوگر وظیفه‌گرا مشترک ارائه‌شده

۴-۱- شبیه‌ساز کاربر

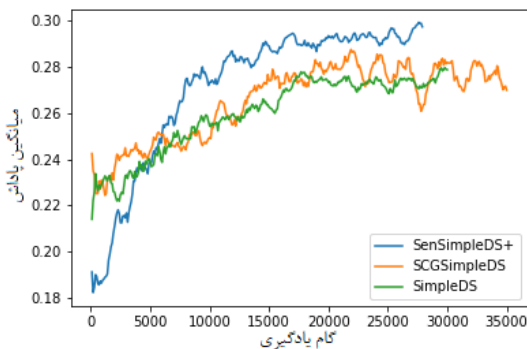
در این سامانه گفتگوگر، بجای انسان از شبیه‌ساز کاربر استفاده می‌شود. شبیه‌ساز، جملاتش را بصورت تصادفی با تعدادی اسلات

⁶ Generative Adversarial Network (GAN)



شکل ۲: پاداش جمع‌آوری شده در ۳۰۰۰ دور گفتگو توسط روش‌های ارائه شده؛ خط نارنجی SenSimpleDS و خط آبی SenSimpleDS+.

برای مقایسه عملکرد روش SenSimpleDS+ ارائه شده با SimpleDS و SCGSimpleDS، همه روش‌ها برای ۳۰۰۰ دور گفتگو با شبیه‌ساز، آموزش دیده‌اند و نمودار میانگین پاداش آنها برای ۱۰۰ گام یادگیری، رسم شده است که در شکل ۳ آورده شده است. دو روش SimpleDS و SCGSimpleDS، برترتیب تعداد ۳۰,۰۰۰ و ۳۲,۵۰۰ گام یادگیری را طی کرده‌اند که از روش SenSimpleDS+، بیشتر است. این یعنی روش SenSimpleDS+ در تعداد گفتگوهای کمتری انسان را به مقصودش می‌رساند. همچنین روش SenSimpleDS+ نسبت به دو روش دیگر، میانگین پاداش بیشتری را جمع‌آوری کرده است که در شکل ۳ نیز قابل مشاهده است. قبل از ۵,۰۰۰ گام، این روش میانگین پاداش کمتری دارد که یعنی این روش نسبت به دو روش دیگر دیرتر یادگیری انجام داده است، که دلیل بزرگتر بودن اندازه بردار حالت آن می‌باشد.



شکل ۳: نتایج بدست آمده توسط روش SenSimpleDS+ (خط آبی) و مقایسه آن با سایر روش‌ها؛ خط سبز SimpleDS و خط نارنجی SCGSimpleDS.

امر، سامانه در حوزه رستوران توسعه یافته و از پیکره ارائه شده توسط روش SimpleDS که قابل دسترس است^۷، استفاده می‌شود. در این بخش همچنین، عملکرد دو روش SenSimpleDS و SenSimpleDS+ در یادگیری نیز، مقایسه می‌شوند.

این سامانه با زبان پایتون و در محیط گوگل کولب، پیاده‌سازی شده است؛ بطوری که عامل و محیط کاملاً از هم مستقل هستند و ارتباط بین آن‌ها متن پیام، حالت، پاداش و کنش می‌باشد. شبیه‌ساز کاربر دارای ۳۳ نوع گفتگو است که چهار اسلات را با مقادیر تصادفی پر می‌کند. بردار حالت با طول ۱۵۳۶، از چسباندن بردار نمایش MPNet ربات و شبیه‌ساز بدست می‌آید. تعداد شش فایل گفتگوی کامل در پیکره وجود دارد که پرونده نمایش از روی آنها ساخته می‌شود. در تابع پاداش، در بخش Cr تعداد کل اسلات‌هایی که باید پر شود برابر با سه؛ مقدار DL برابر با ۰/۱ و مقدار α برابر با ۰/۵، در نظر گرفته شده است. با این مقدار α ، سامانه بطور مساوی به اسلات‌های پر شده و شبیه‌بودن حالت فعلی به پرونده نمایش، اهمیت می‌دهد. برای روش SenSimpleDS+، مقدار β برابر با ۰/۹۹۹۸ در نظر گرفته شده است، تا مقدار ϵ جایی که می‌شود به آهستگی کاسته شود و در حدود گام یادگیری ۳۵,۰۰۰ به ϵ_{min} برسد. برای شبکه یادگیری تقویتی عمیق، لایه ورودی تعداد ۱۵۳۶ گره و دو لایه پنهان با ۴۰ گره، در نظر گرفته شده است. لایه خروجی نیز ۳۶ گره دارد که برابر با تمام کنش‌های ربات می‌باشد. همچنین، با استفاده از طبقه‌بند بیز ساده، کنش‌هایی که احتمال حداقل ۰/۰۰۱ داشته باشند به عنوان کنش‌های مجاز تعیین می‌شوند.

برای مقایسه دو روش SenSimpleDS و SenSimpleDS+، هر دو روش برای ۳۰۰۰ دور گفتگو با شبیه‌ساز، آموزش دیده‌اند. پاداش بدست‌آمده توسط آنها در ۱۰۰ گام یادگیری (یک گام گفتگو)، میانگین متحرک گرفته شده که در شکل ۲ قابل مشاهده است. از ابتدای آموزش تا قبل از گام ۲۰,۰۰۰، روش دوم پاداش بیشتری جمع‌آوری کرده است که دلیل آن استفاده از تابع دوم برای کاهش مقدار ϵ است که منجر به یادگیری سریع‌تر یادگیری تقویتی عمیق شده و از گام ۲۰,۰۰۰ به بعد، روش اول به میانگین پاداش بیشتری می‌رسد. روش اول و دوم، برترتیب تعداد ۲۸,۸۰۰ و ۲۷,۹۰۰ گام یادگیری را طی کرده‌اند که نشان می‌دهد روش دوم در گام‌های کمتری گفتگوها را به پایان رسانده که در شکل ۲ نیز مشخص است؛ دلیل آن نیز استفاده از تابع دوم، برای کاهش مقدار ϵ می‌باشد.

⁷ github.com/cuayahuitl/SimpleDS



یادگیری بهتر آن دارد. استفاده از تابع جدید برای کاهش ϵ نیز، بخوبی تناسبی میان استفاده از دانش کنونی و بدست آوردن دانش جدید (با کنش تصادفی)، ایجاد کرد، که منجر به یادگیری سریعتر عامل یادگیر تقویتی شد.

مراجع

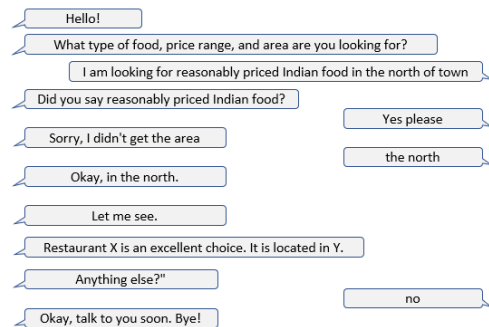
- [۱] B. Liu, G. Tur, D. Hakkani-Tur, P. Shah, and L. Heck, "Dialogue learning with human teaching and feedback in end-to-end trainable task-oriented dialogue systems," *arXiv preprint arXiv:1804.06512*, 2018.
- [۲] K. Song, X. Tan, T. Qin, J. Lu, and T.-Y. Liu, "Mpnnet: Masked and permuted pre-training for language understanding," *Advances in Neural Information Processing Systems*, vol. 33, pp. 16857-16867, 2020.
- [۳] X. Zhang and H. Wang, "A joint model of intent determination and slot filling for spoken language understanding," in *IJCAI*, 2016, vol. 16, no. 2016, pp. 2993-2999.
- [۴] Y. Shi, K. Yao, H. Chen, Y.-C. Pan, M.-Y. Hwang, and B. Peng, "Contextual spoken language understanding using recurrent neural networks," in *2015 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*, 2015: IEEE, pp. 5271-5275.
- [۵] Z. Yan, N. Duan, P. Chen, M. Zhou, J. Zhou, and Z. Li, "Building task-oriented dialogue systems for online shopping," in *Thirty-first AAAI conference on artificial intelligence*, 2017.
- [۶] H. Cuayáhuil, "Simpleds: A simple deep reinforcement learning dialogue system," in *Dialogues with social robots*: Springer, 2017, pp. 109-118.
- [۷] H. Cuayáhuil, S. Yu, A. Williamson, and J. Carse, "Scaling up deep reinforcement learning for multi-domain dialogue systems," in *2017 International Joint Conference on Neural Networks (IJCNN)*, 2017: IEEE, pp. 3339-3346.
- [۸] D. Zahra and S. S. Javad, "Design And Development Of A Dialogue System Using Deep Reinforcement Learning," in *12th International Conference on Information Technology, Computer and Telecommunication*, 2021.
- [۹] V. Ilievski, C. Musat, A. Hossmann, and M. Baeriswyl, "Goal-oriented chatbot dialog management bootstrapping with transfer learning," *arXiv preprint arXiv:1802.00500*, 2018.
- [۱۰] O. Dušek and F. Jurčiček, "Sequence-to-sequence generation for spoken dialogue via deep syntax trees and strings," *arXiv preprint arXiv:1606.05491*, 2016.

بیشینه پاداش جمع‌آوری شده روش‌های ارائه‌شده و دو روش SimpleDS و SCGSimpleDS در جدول ۱ آورده شده‌است که بیشترین مقدار آن مربوط به روش SenSimpleDS+ می‌باشد.

جدول ۱: بیشینه پاداش بدست‌آمده توسط روش‌های مختلف در ۳۰۰۰ دور گفتگو با شبیه‌ساز.

نام روش	بیشینه پاداش
SimpleDS [6]	۰/۲۷۹۷
SCGSimpleDS [8]	۰/۲۸۷۶
SenSimpleDS	۰/۲۹۷۹۵
SenSimpleDS+	۰/۲۹۹۳۷

در ادامه در شکل ۴، یک نمونه از گفتگوی ربات و شبیه‌ساز کاربر، توسط روش SenSimpleDS+ آورده شده‌است.



شکل ۴: نمونه از گفتگوی روش ارائه شده در زبان انگلیسی؛ سمت چپ ربات و سمت راست شبیه‌ساز.

۶. نتیجه‌گیری

در سال‌های اخیر، توسعه سامانه‌های گفتگوگر وظیفه‌گرا مشترک، بسیار مورد توجه قرار گرفته است. یکی از روش‌های توسعه این سامانه‌ها، استفاده از یادگیری تقویتی عمیق می‌باشد. در این مقاله، با استفاده از یادگیری تقویتی عمیق و شبکه MPNet که برای نمایش جمله بکار می‌رود، روشی بنام SenSimpleDS برای توسعه سامانه گفتگوگر وظیفه‌گرا مشترک ارائه شد. در این روش، برای تعیین کنش‌های مجاز یادگیری تقویتی عمیق، از الگوریتم بیز ساده استفاده شد. با اضافه کردن تابع نمایشی برای کاهش ϵ ، روش SenSimpleDS+ روش دیگری بود که ارائه شد. عملکرد روش‌های ارائه‌شده در حوزه رستوران، ارزیابی شد و نتایج آن با روش‌های دیگر با استفاده از گفتگو با شبیه‌ساز مقایسه شد، که نسبت به آنها، بهبود ایجاد کرده‌اند. استفاده از نمایش جمله برای ساخت بردار حالت، توانست نمایش بهتری از محیط بسازد و در نتیجه عامل یادگیری تقویتی عمیق، پاداش بیشتری دریافت کرد؛ که نشان از