



طراحی سیستم معاملاتی الگوریتمیک بر پایه یادگیری تقویتی عمیق مورد مطالعاتی: بورس اوراق بهادار تهران

سعید کاظمیان حسین آبادی^۱

سید محمدرضا داودی^{۲*}

محمد مشهدی زاده^۳

پارسا جوزی^۴

تاریخ دریافت: ۱۴۰۱/۱۱/۰۴ تاریخ پذیرش: ۱۴۰۲/۰۳/۱۳

چکیده

امروزه معاملات الگوریتمی استفاده گسترده‌ای در مدیریت معاملات دارد. سبد گردانی الگوریتمی نوع جدیدی از این سامانه‌هاست که از طریق آن سبد گردان با استفاده از ابزارهای الگوریتمی به بالا بردن کیفیت سود و کاهش ریسک‌های سبد خود کمک می‌کند. هدف از پژوهش حاضر طراحی سیستم معاملاتی الگوریتمیک بر پایه یادگیری تقویتی عمیق به کمک شبکه عصبی است. در این رویکرد عامل یا معامله‌گر در فضای جستجو برای یافتن پاداش بیشتر که همان بازده بیشتر می‌باشد، به جستجو می‌پردازد. عامل معامله‌گر با سیگنال‌های تکنیکی شامل شاخص قدرت نسبی، نوسانگر تصادفی، نشانگر همگرایی-واگرایی و قیمت‌های کمینه، بیشینه، بسته شدن و باز شدن مواجه می‌شود. یادگیری تقویتی عمیق جدول تابع ارزش یا کیفیت Q را با یک شبکه عصبی جایگزین می‌کند. شبکه عصبی مذکور در نهایت با دریافت ورودی حالت، یکی از سه عمل فروش، خرید و نگهداری را پیشنهاد می‌کند. این پیشنهاد به صورت سه احتمال با مجموع یک می‌باشد و پیشنهاد با حداکثر احتمال مورد پیاده‌سازی قرار می‌گیرد. نتیجه پیاده‌سازی سیستم معاملاتی یادگیری تقویتی عمیق بر روی شاخص کل بورس اوراق بهادار تهران در بازه ۱۳۹۱ تا ۱۴۰۰ نشان می‌دهد که سیستم پژوهش در میانگین و شاخص همگرایی-واگرایی دارای تفاوت معناداری با سه سیستم دیگر داشت. همچنین نسبت شارپ سیستم پژوهش نسبت به سه مدل دیگر رشد حداقل ۱/۴ برابری را نشان داد.

کلمات کلیدی: شبکه‌های عصبی، یادگیری تقویتی، یادگیری تقویتی عمیق، نوسانگر تکنیکی.

^۱ کارشناسی ارشد مدیریت مالی، واحد دهقان، دانشگاه آزاد اسلامی، دهقان، ایران. saeedkazemian1365@yahoo.com

^۲ استادیار، گروه مدیریت، واحد دهقان، دانشگاه آزاد اسلامی، دهقان، ایران (نویسنده مسئول). smrdavoodi@ut.ac.ir

^۳ استادیار، گروه مدیریت، واحد مبارکه، دانشگاه آزاد اسلامی، مبارکه، ایران. m.mashhadi@mau.ac.ir

^۴ کارشناسی ارشد مدیریت مالی، واحد دهقان، دانشگاه آزاد اسلامی، دهقان، ایران. Parsajozipersonal@gmail.com

۱. مقدمه

سرمایه‌گذاری و انباشت سرمایه در تحول اقتصادی کشور نقش بسزایی داشته است. اهمیت این عامل و نقش مؤثر آن را می‌توان به‌وضوح در سیستم کشورهای با نظام سرمایه‌گذاری مشاهده کرد. بدون شک بورس یکی از مناسب‌ترین جایگاه‌ها جهت جذب سرمایه‌های کوچک و استفاده از آن‌ها جهت رشد یک شرکت، در سطح کلان و نیز رشد شخصی فرد سرمایه‌گذار است (شمس و عطایی، ۱۳۹۵). از آنجایی که هدف و تعریف سرمایه‌گذاری، به تعویق انداختن مصرف جهت مصرف بیشتر و بهتر در آینده است؛ افراد با سرمایه‌گذاری انتظار دستیابی به سود مورد انتظار خود را دارند. بنابراین مهم‌ترین امر در این زمینه، خرید یک سهم به قیمت پایین و فروش آن به قیمت بالاتر است که این موضوع؛ به معنی تعیین مناسب موقعیت معاملاتی است. معاملات الگوریتمی در بازارهای مالی، به معنای استفاده از برنامه‌های کامپیوتری برای انتخاب سهام از بین صدها موقعیت موجود، ورود سفارش‌های خرید و فروش به سیستم‌های معاملاتی، کمک به اندازه‌گیری ریسک و بهینه‌سازی پرتفوی و کمک به مدیریت دارایی با بهره‌وری بالاتر می‌باشد. در این روش یک یا چند الگوریتم در انتخاب و اعمال این سفارش‌ها از جنبه‌های مختلف مانند زمان‌بندی، قیمت یا حجم معاملات آن تصمیم می‌گیرند. معاملات الگوریتمی در بسیاری از مواقع، بدون دخالت انسان انجام می‌شود (تیبات و دامین^۱، ۲۰۲۰). امروزه معاملات الگوریتمی استفاده گسترده‌ای در مدیریت معاملات بانک‌های سرمایه‌گذاری (تامین سرمایه‌ها)، صندوق‌های بازنشستگی (و شرکت‌های سرمایه‌گذاری) و صندوق‌های سرمایه‌گذاری

مشترک دارد. سبب گردانی الگوریتمی نوع جدیدی از این سامانه‌هاست که از طریق آن سبب گردان با استفاده از ابزارهای الگوریتمی به بالا بردن کیفیت سود و کاهش ریسک‌های سبب خود کمک می‌کند. همچنین با استفاده از این ابزار، سبب گردان قادر می‌شود برای تعداد زیادی مشتری سبب گردانی واقعا به‌صورت اختصاصی انجام دهد. هوش مصنوعی و به‌ویژه شبکه‌های عصبی از ابزارهای نوین در پیاده‌سازی معاملات الگوریتمی محسوب می‌شود (ونگ و همکاران^۲، ۲۰۱۵). نتایج تحقیقات زیادی نشان داده‌اند که رفتار قیمتی سهام از نوع غیرخطی می‌باشد و سیستم‌های غیرخطی موفقیت بیشتری در استخراج ویژگی‌های بازار و در نهایت تعیین استراتژی‌های معاملاتی سودمند دارند. سیستم‌های غیرخطی دارای انعطاف بیشتری در مدل‌سازی داده‌ها و درک ساختار و روابط پیچیده بین متغیرها هستند. یکی از این مدل‌های پرکاربرد غیرخطی، شبکه‌های عصبی^۳ می‌باشد. شبکه‌های عصبی سامانه‌های مدل‌سازی هستند که به‌صورت غیرخطی بر روی داده‌های ورودی تغییراتی را ایجاد می‌کنند و آن را به خروجی تبدیل می‌کنند. این سیستم‌ها از بدن انسان الهام گرفته‌اند و برای مدل‌سازی از مفاهیمی همچون نرون و توابع برانگیختگی استفاده می‌کند. یک شبکه عصبی برای کاربرد باید مورد آموزش قرار گیرد تا پارامترهای مجهول ساختاری آن مقداردهی بهینه شوند (کراتاماددی و همکاران^۴، ۲۰۲۱). در مبحث یادگیری ماشین و شبکه‌های عصبی مفهوم مهمی بانام یادگیری عمیق قرارداد. یادگیری عمیق^۵ شاخه‌ای از بحث یادگیری ماشین و هوش مصنوعی و مجموعه‌ای از الگوریتم‌هایی است که تلاش می‌کنند مفاهیم انتزاعی سطح بالا را با استفاده از یادگیری در سطوح و لایه‌های

¹ Thibaut and Damien

² Wang et al

³ Neural network

⁴ Koratamaddi et al

⁵ Deep Learning

یادگیری ماشینی است که بر اساس پاداش دادن به رفتارهای دلخواه و/یا تنبیه رفتارهای نامطلوب است. به‌طور کلی، یک عامل یادگیری تقویتی قادر است محیط خود را درک و تفسیر کند، اقداماتی انجام دهد و از طریق آزمون و خطا یاد بگیرد. این اهداف بلندمدت به جلوگیری از توقف عامل در اهداف کمتر کمک می‌کند. با گذشت زمان، عامل یاد می‌گیرد که از منفی‌ها دوری کند و به دنبال مثبت باشد. این روش یادگیری در هوش مصنوعی به‌عنوان راهی برای هدایت یادگیری ماشینی بدون نظارت از طریق پاداش و جریمه پذیرفته شده است. معاملات الگوریتمی در بازار سرمایه به‌سرعت در حال گسترش می‌باشد. در دنیا درصد قابل توجهی از معاملات با این روش انجام می‌شود. امروز بخش عمده‌ای از فعالیت‌های بازار سرمایه از حالت دستی خارج شده و الکترونیکی انجام می‌شود و استفاده از فرصت‌های زودگذر و آربیتراژی نیازمند سرعت عمل بالا در شناسایی و اجرا است. طراحی سیستم‌های هوشمند جهت کسب چنین بازده‌هایی ضروری است. پژوهش حاضر نیز با استفاده از ابزارهای هوش مصنوعی نسبت به طراحی یک سیستم معاملات الگوریتمی اقدام می‌کند. تاکنون یادگیری تقویتی عمیق و سیستم معاملاتی مبتنی بر آن در بورس اوراق بهادار تهران مورد پیاده‌سازی و بررسی سودآوری قرار نگرفته است و لذا خلأ پژوهشی در تحقیق حاضر استفاده از این تکنیک به‌منظور پیاده‌سازی یک سیستم معاملاتی خودکار می‌باشد. بنابراین در پژوهش حاضر یک سیستم معاملاتی بر اساس رویکرد یادگیری تقویتی و به کمک شبکه‌های عصبی عمیق طراحی می‌شود.

مختلف مدل کنند. مطالعات بالینی نشان می‌دهند، که ساختار مغز پستانداران از معماری شبکه‌های عصبی عمیق بهره می‌برد که در آن، مفاهیم انتزاعی در لایه‌های مختلف، به ترتیب از مفاهیم و ویژگی‌های ساده تا مفاهیم سطح بالا، در نواحی مختلف قشر مغز پردازش می‌شوند. ایده یادگیری عمیق با الهام از ساختار طبیعی مغز انسان و به کمک امکانات و فن‌آوری‌های جدید، توانسته است در بسیاری از حوزه‌های مربوط به هوش مصنوعی و یادگیری ماشین، موفقیت‌های چشم‌گیری را کسب کند. از مزایای یادگیری عمیق می‌توان به یادگیری خودکار ویژگی‌ها، یادگیری چندلایه ویژگی‌ها، دقت بالا در نتایج، قدرت تعمیم بالا و شناسایی داده‌های جدید، پشتیبانی گسترده سخت‌افزاری و نرم‌افزاری، پتانسیل ایجاد قابلیت‌ها و کاربردهای بیشتر در آینده اشاره کرد. یادگیری تقویتی یکی از گرایش‌های یادگیری ماشینی است که از روانشناسی رفتارگرایی الهام می‌گیرد. این روش بر رفتارهایی تمرکز دارد که ماشین باید برای پیشینه کردن پاداشش انجام دهد. این مسئله، با توجه به گستردگی‌اش، در زمینه‌های گوناگونی مانند نظریه بازی‌ها، نظریه کنترل، تحقیق در عملیات، نظریه اطلاعات، سامانه چندعامله، هوش ازدحامی، آمار، الگوریتم ژنتیک، بهینه‌سازی بر مبنای شبیه‌سازی بررسی و استفاده می‌شود (ژانگ و لی^۱، ۲۰۲۲). یادگیری تقویتی با یادگیری با نظارت معمول دو تفاوت عمده دارد، نخست اینکه در آن زوج‌های صحیح ورودی و خروجی در کار نیست و رفتارهای ناکارآمد نیز از بیرون اصلاح نمی‌شوند، و دیگر آنکه تمرکز زیادی روی کارایی زنده وجود دارد که نیازمند پیدا کردن یک تعادل مناسب بین اکتشاف چیزهای جدید و بهره‌برداری از دانش اندوخته شده دارد. یادگیری تقویتی یک روش آموزش

¹ Zhang and Lei

به منظور واکاوی مفهوم سیستم معاملاتی الگوریتمیک بر پایه یادگیری تقویتی عمیق سوالات زیر مطرح می شود که در این پژوهش به آن پاسخ داده می شود:

۱- چگونه می توان به کمک یادگیری تقویتی عمیق یک سیستم معاملاتی طراحی کرد؟

۲- سودآوری سیستم معاملاتی طراحی شده چگونه است؟

در مبانی نظری به مطالعه اجمالی در مورد پیشینه پژوهش با تأکید بر سیستم معاملاتی الگوریتمیک بر پایه یادگیری تقویتی عمیق به کمک شبکه عصبی پرداخته می شود.

۲. مبانی نظری و مروری بر پیشینه پژوهش

پیشرفت های اخیر در حوزه فناوری های کامپیوتری، گسترش به کارگیری فناوری های معاملات الگوریتمی و خودکار، توسعه استراتژی های معاملاتی و به تبع آن سرعت بالای معاملات، منجر به افزایش پیچیدگی های سرمایه گذاری در بازار سهام شده اند (ژانگ و همکاران^۱، ۲۰۱۶). معاملات الگوریتمی نوعی از معاملات خودکار بوده که شامل برنامه های کامپیوتری برای ارسال سفارشات همراه با الگوریتم های تصمیم گیری هستند که این الگوریتم ها خود بر اساس پارامترهای منحصربه فرد سفارش مانند زمان، قیمت و یا مقدار سفارش می باشند. در بازارهای مالی الکترونیکی، معاملات الگوریتمی به معنای استفاده از برنامه های کامپیوتری برای ورود سفارش های معاملاتی است که سیگنال های معاملاتی آن توسط بخشی دیگر از سیستم تولید شده است (رستگار و صدقاتی پور، ۱۳۹۷).

یادگیری تقویتی می تواند در یک موقعیت عمل کند تا زمانی که بتوان پاداش روشنی را اعمال کرد. در مدیریت منابع سازمانی، الگوریتم های یادگیری تقویتی می توانند منابع محدودی را به وظایف مختلف اختصاص دهند تا زمانی که هدف کلی وجود داشته باشد. هدف در این شرایط صرفه جویی در زمان یا حفظ منابع خواهد بود. در رباتیک، یادگیری تقویتی به تست های محدود راه پیدا کرده است. این نوع یادگیری ماشینی می تواند به روبات ها توانایی یادگیری وظایفی را که معلم انسانی نمی تواند نشان دهد، انطباق مهارت های آموخته شده با یک کار جدید یا دستیابی به بهینه سازی علیرغم کمبود فرمول بندی تحلیلی در دسترس، ارائه دهد. یادگیری تقویتی همچنین در تحقیقات عملیات، نظریه اطلاعات، نظریه بازی، کنترل، بهینه سازی مبتنی بر شبیه سازی، سیستم های چندعاملی، هوش ازدحام، آمار و الگوریتم های ژنتیک استفاده می شود (یو و یان^۲، ۲۰۲۰).

صیادی نژاد و همکاران (۱۴۰۲) در پژوهش خود بیان کردند با افزایش محبوبیت و فراگیر شدن رمزارزها، ایجاد و توسعه روش های پیش بینی حرکت های قیمتی در این حوزه، توجهات زیادی را به خود جلب کرده است. در این بین مدل های یادگیری عمیق (DL) با ساختارهایی مانند حافظه طولانی کوتاه مدت (LSTM) و شبکه عصبی کانولوشنی (CNN) پیشرفت هایی در تحلیل این نوع از داده ها ایجاد کرده است. یکی دیگر از رویکردهایی که می تواند در تحلیل قیمتی بازار رمزارزها کارا باشد تجزیه سیگنال های از طریق الگوریتم هایی مانند تجزیه مد تجربی یکپارچه کامل (CEEMD) می باشد. با توجه به اهمیت مقوله

^۲ Yu and Yan

^۱ Zhang and Lei

پیش‌بینی در بازار رمز ارزها، در این تحقیق با ترکیب مدل‌های یادگیری عمیق و روش تجزیه مد تجربی یکپارچه کامل (CEEMD)، مدل هیبریدی $CEEMD-DL(LSTM)$ به‌منظور پیش‌بینی بازدهی قیمتی رمز ارز بیت‌کوین (به‌عنوان محبوب‌ترین رمز ارز) مورد استفاده قرار گرفته‌است. در این راستا از داده‌های روزانه قیمتی بیت‌کوین در دوره زمانی ۲۰۱۳/۰۱/۰۱ - ۲۰۲۲/۰۵/۲۸ استفاده گردید و نتایج به‌دست‌آمده با نتایج مدل‌های رقیب بر اساس معیارهای سنجش کارایی مقایسه شد. بر اساس نتایج به‌دست‌آمده، استفاده از مدل معرفی شده $(CEEMD-DL(LSTM))$ ، کارایی و دقت پیش‌بینی‌های بازدهی رمز ارز بیت‌کوین را افزایش داده‌است. بر همین اساس کاربرد این مدل به‌منظور پیش‌بینی در این حوزه پیشنهاد می‌گردد.

اسفندیار و همکاران (۱۴۰۱) در پژوهش خود بیان کردند که این مقاله کاربرد معاملات الگوریتمی با تمرکز بر رویکرد یادگیری تقویتی برای بهینه‌سازی پرتفوی سهام‌های منتخب را بررسی می‌کند. این پژوهش از حیث هدف، کاربردی و از نظر نوع داده، کمی و از لحاظ روش، توصیفی - اکتشافی و از منظر طرح تحقیق، پس‌رویدادی است. جامعه آماری این پژوهش، ۶۷۲ شرکت بورس است که از این تعداد، داده‌های پنج شرکت (نمونه آماری) طی دوره زمانی ۱۳۹۶-۱۴۰۰ بررسی شده است. یافته‌های تحقیق در دوره‌های صعودی و نزولی بازار نشان داد که رویکرد یادگیری تقویتی در بازارهای صعودی و نزولی به‌صورت معناداری بر رویکرد خرید و نگهداری برتری دارد و عملکرد بهتری ارائه داده است و نتایج با عملکرد الگوریتم‌ها در بازارهای بورس سازگار است. نتایج آشکار کرد که از دیدگاه سودآوری، رویکرد یادگیری تقویتی نسبت به رهیافت خرید و نگهداری، عملکرد بهتر

و موثرتری داشته است؛ بنابراین، به‌کارگیری روش یادگیری تقویتی پیشنهاد می‌شود. دسارنج و همکاران (۱۳۹۹) در پژوهشی خود با هدف طراحی سیستم معاملات تکنیکی سهام با استفاده از مدل ترکیبی شبکه عصبی MLP و الگوریتم‌های تکاملی نشان دادند که مدل‌های ترکیبی MLP و الگوریتم‌های تکاملی عملکرد بهتر و معناداری نسبت به روش خرید و نگهداری و مدل $MLP-BP$ داشته است و مدل MLP_PSO بازدهی بیش‌تری نسبت به سایر مدل‌ها کسب کرده است. یافتیان و رستگار (۱۳۹۹) در پژوهش خود با هدف طراحی یک سیستم معاملاتی خودکار با استفاده از شبکه عصبی پیچشی نشان دادند که برای بررسی بازدهی و ریسک سیستم ارائه‌شده، یک روش برای خرید و فروش بر اساس نتایج مدل در زمان گذشته معرفی شده است که در ۸۰٪ موارد، این روش بازدهی بیشتری نسبت به استراتژی مرسوم خرید و نگهداری کسب کرده است. همچنین همواره از نظر معیارهای ریسک انحراف معیار و بیشترین افت بهتر عمل می‌کند. همچنین، نتایج نشان‌دهنده تأثیر زیاد کارمزد معاملات بورس اوراق بهادار تهران بر روی بازدهی مدل است. به‌گونه‌ای که مدل چند برابر سود کسب‌شده را برای پرداخت کارمزد از دست می‌دهد. خنجر پناه و همکاران (۱۳۹۷) در پژوهش خود با هدف کاربرد روش تکنیکال برای پیش‌بینی قیمت سهام: رویکرد مدل‌های احتمال غیرخطی و شبکه‌های عصبی مصنوعی نشان دادند که در آزمون نا پارامتری برابری نسبت‌ها، از لحاظ آماری مدل‌های ارائه‌شده تفاوت معناداری باهم نداشته‌اند، اما معیارهای سنجش خطا بیان می‌کند که مدل پروبیت، خطای کمتری در پیش‌بینی سهام در بازار بورس تهران دارد.

سهام استفاده کردند. آن‌ها هدف خود را افزایش دقت و سرعت پیش‌بینی نوسانات قیمت سهام با استفاده از مدل یادگیری تقویتی عمیق روش استراتژی گرادیان یا PG^4 قرار دادند. آزمایش‌های صورت گرفته نشان می‌دهند که دقت پیش‌بینی الگوریتم جدید و سرعت هم‌گرایی پاداش به‌طور قابل‌توجهی بالاتر از الگوریتم یادگیری تقویتی عمیق سنتی است. در نتیجه، الگوریتم جدید با شرایط نوسان بازار سازگارتر است. کراماتیکی و همکاران⁵ (۲۰۲۱) یک رویکرد یادگیری تقویتی عمیق جدید را برای آموزش مؤثر یک معامله‌گر خودکار هوشمند پیشنهاد می‌کنند، که نه تنها از داده‌های تاریخی قیمت سهام استفاده می‌کند، بلکه احساسات بازار را برای یک سبد سهام متشکل از شرکت‌های داو جونز درک می‌کند. نشان داده می‌شود که رویکرد پژوهش در مقایسه با خطوط پایه موجود در معیارهای استاندارد شده مانند نسبت شارپ و بازده سرمایه‌گذاری سالانه قوی‌تر است. تیبات و دامین⁶ (۲۰۲۰) در پژوهشی یک رویکرد نوآورانه مبتنی بر یادگیری تقویت عمیق (DRL) برای حل مشکل تعیین موقعیت بهینه (خرید، فروش و نگهداری) در یک معامله الگوریتمیک در طول فعالیت تجاری در بازارهای سهام را ارائه می‌دهند. استراتژی تجاری جدید DRL پیشنهاد شده به دنبال بیشینه‌سازی نسبت شارپ در طیف وسیعی از بازارهای سهام می‌باشد. یانگ و ژای⁷ (۲۰۲۰) برای پیش‌بینی قیمت سهام از یک روش ترکیبی در یادگیری عمیق استفاده کردند. ابتدا به کمک شبکه‌های عصبی پیچشی به استخراج ویژگی از ورودی‌ها شامل یک تنسور سه بعدی از ورودی‌های تکنیکی، قیمت و شاخص‌ها پرداختند و در ادامه به کمک

جین¹ (۲۰۲۳) در مقاله خود یک الگوریتم نوآورانه را برای حل مسئله سبد بهینه در فعالیتهای معاملاتی بازار سهام پیشنهاد می‌کند. استراتژی جدید معاملات پرتفوی پژوهش از سه ویژگی برای پیشی گرفتن از سایر استراتژی‌های معیار در یک محیط بازار واقعی استفاده می‌کند. ابتدا، یک مدل بهینه‌سازی پورتفولیو میانگین-واریانس را پیشنهاد می‌کند که راه‌حل آن بر اساس معماری بازیگر-منتقد است. برخلاف ادبیات موجود که انتظار بازده تجمعی را می‌آموزد، ماژول منتقد توزیع بازده تجمعی را با رگرسیون کمی می‌آموزد و ماژول بازیگر وزن بهینه پورتفولیو را با حداکثر کردن تابع هدف مدل بهینه‌سازی خروجی می‌دهد. ثانیه، از یک تابع تبدیل خطی برای تحقق فروش کوتاه‌مدت استفاده می‌کند، سوم، یک روش چند فرآیندی به نام $Ape-x$ برای تسریع سرعت آموزش یادگیری تقویتی عمیق استفاده شد. برای اعتبار سنجی رویکرد پیشنهادی خود، این پژوهش برای دو نمونه کارنامه نماینده یک تست انجام می‌دهد و مشاهده می‌کند که مدل پیشنهادی در این کار برتر از استراتژی‌های معیار است.

لی² و همکاران (۲۰۲۲) یک مدل یادگیری تقویتی جدید را برای پیاده‌سازی معاملات سهام پیشنهاد کردند که بازار سهام را از طریق داده‌های سهام، شاخص‌های فنی و نمودارهای شمعی تحلیل می‌کند و استراتژی‌های معاملاتی پویا را می‌آموزد. نتایج نشان داد که استراتژی معاملاتی می‌تواند در مقایسه با سایر استراتژی‌های تجاری سود بیشتری به دست آورد. ژانگ و لی³ (۲۰۲۲) از یادگیری تقویتی عمیق به‌منظور پیش‌بینی قیمت

⁵ Koratamaddi et al

⁶ Thibaut و Damien

⁷ Yang and Zhai

¹ Jin

² Li et al

³ Zhang and Lei

⁴ policy gradient strategy

می‌باشد. روش تحقیق در این پژوهش توصیفی-تحلیلی می‌باشد. این تحقیق بر اساس هدف یک تحقیق کاربردی است. با توجه به هدف، پژوهش حاضر در دسته کلی طراحی سیستم‌های معاملاتی مبتنی بر هوش مصنوعی جای می‌گیرد. ورودی سیستم معاملاتی شامل قیمت‌های باز و بسته شدن، بیشینه و کمینه قیمت به همراه شاخص‌های تکنیکی قدرت نسبی، شاخص تصادفی و همگرایی-واگرایی می‌باشد. خروجی سیستم معاملاتی مبتنی بر یادگیری تقویتی عمیق نیز شامل اتخاذ یکی از سه موقعیت خرید، فروش و نگهداری می‌باشد. این سیستم به‌جای استفاده از جدول q از یک شبکه عصبی عمیق برای رکورد نتایج حاصل از اقدامات استفاده می‌کند. پس از آموزش سیستم معاملاتی بر روی داده‌های آموزشی، دقت مدل پیش‌بینی بر روی داده‌های تست مورد ارزیابی قرار می‌گیرد. برای تجزیه و تحلیل اطلاعات از کد نویسی با زبان پایتون و پکیج تنسورفلو و همچنین زبان برنامه‌نویسی متلب استفاده خواهد شد. مدل پژوهش برگرفته از تحقیق لی و همکاران (۲۰۲۲) می‌باشد. جامعه آماری پژوهش شاخص کل بورس اوراق بهادار تهران به‌عنوان نماینده بازار می‌باشد و نمونه آماری به‌صورت زمانی شامل شاخص کل بورس اوراق بهادار تهران در بازه شروع سال ۱۳۹۱ تا پایان تیر ۱۴۰۱ می‌باشد.

۱.۳. تبیین مدل پژوهش

در این بخش به طراحی سیستم معاملاتی بر پایه یادگیری تقویتی عمیق پرداخته می‌شود. اما در ابتدا به معرفی سه نوسان‌نمای تکنیکی می‌پردازیم که بخشی از ورودی‌های سیستم معاملاتی می‌باشد. یکی از این نوسان‌ماهای معروف و پرکاربرد، قدرت نسبی یا RSI

شبکه‌های حافظه طولانی کوتاه‌مدت به پیش‌بینی پرداختند. نتیجه پژوهش بر روی ۵ شاخص از جمله $S\&P$ نشان می‌دهد که دقت روش ترکیبی از روش‌های دیگر چون بردار پشتیبان و ترکیب شبکه‌های عصبی با تحلیل مولفه اصلی بیشتر است. نبی پور و همکاران (۲۰۲۰) برای پیش‌بینی قیمت سهام از الگوریتم‌های مختلفی در حوزه هوش مصنوعی استفاده کردند. این الگوریتم‌ها شامل شبکه‌های عصبی، درخت‌های تصمیم، روش‌های جمعی، جنگل تصادفی و یادگیری حافظه طولانی کوتاه‌مدت می‌باشد. نتیجه پژوهش بر روی داده‌های تاریخی ۱۰ ساله از بورس اوراق بهادار تهران بر روی ۵ شاخص با ورودی‌های تکنیکی نشان می‌دهد که شبکه عصبی بازگشتی دارای بالاترین دقت می‌باشد. ژو و کسلج^۱ (۲۰۱۹) برای پیش‌بینی شاخص بازار از الگوریتم یادگیری عمیق شبکه‌های حافظه طولانی کوتاه‌مدت استفاده کردند. ورودی‌های سیستم پیش‌بینی شامل داده‌های تاریخی، نما گرهای تکنیکی و احساسات سهامداران می‌باشد. برای سنجش احساسات سهامداران از توییت منتشر شده در توییت استفاده شده است. نتیجه پژوهش نشان می‌دهد که سیستم پیش‌بینی با در نظر گرفتن احساسات سرمایه‌گذاران دارای دقت بالاتری نسبت به سیستم عمیقی می‌باشد که این ورودی را ندارد. همچنین نتیجه پژوهش نشان می‌دهد که توییت افراد دارای تعداد دنباله کنندگان بالا نقش مهمی در تبیین خروجی دارد.

۳. روش‌شناسی پژوهش

هدف پژوهش حاضر طراحی و ارزیابی عملکرد یک سیستم معاملاتی مبتنی بر یادگیری تقویتی عمیق

¹ Xu and Keselj

به طور معمول این دوره ۱۴ روزه می باشد. مفهوم این نوسانگر به این صورت است که همواره بزرگی قیمت نسبت به فاصله بیشترین و کمترین قیمت در دوره زمانی مورد بررسی سنجیده می شود. بنابراین هر چه $K\%$ بزرگتر شود و نزدیکی این عدد به ۱۰۰ نشان از پتانسیل بیشتر سهم برای ریزش قیمت دارد. از طرفی هر چه $K\%$ کمتر شود و نزدیکی این عدد به ۰ نشان از پتانسیل بیشتر سهم برای ریزش قیمت دارد. غالباً برای این نوسانگر نیز دو حد آستانه خرید افراطی و فروش افراطی در نظر گرفته می شود که معمولاً این دو خط به ترتیب خطوط ۸۰ و ۲۰ می باشند. میانگین متحرکها گروهی از عملگرهای هموار کننده می باشند که به جای سیگنال اصلی قیمت در هر لحظه، میانگین قیمتها را در یک فاصله زمانی مشخص محاسبه می کنند و به جای آن قرار می دهند. میانگین متحرکها با هموارسازی نمودار قیمت، روند حرکت آن را به صورت عینی تری به نمایش می گذارند. در صورتی که $\{p(t)\}_{t=1}^T$ سری قیمت یک سهام باشد، سری میانگین متحرک ساده به طول n برابر است با رابطه (۳):

$$MA_n(t) = \frac{p(t) + p(t-1) + p(t-n+1)}{n}$$

در صورتی که n عدد بزرگی باشد، میانگین متحرک روند تاریخی را بهتر نشان می دهد و هر چه n کوچکتر باشد حساسیت آن به داده های جدیدتر بیشتر می شود. با استفاده از ترکیب میانگین متحرکها نوسان نماهای جدیدی به نام نوسان نماهای واگرایی -

باشند که $n > m$ ، آنگاه در صورتی که $MA_n(t) < MA_m(t)$ سیگنالی برای خرید دریافت

می باشد. در تنظیمات این اندیکاتور عموماً از تنظیم ۱۴ روزه استفاده می شود. فرمول محاسبه این نوسان نما به صورت زیر می باشد:

(۱)

$$RSI = (100 - \frac{100}{1 + RS})$$

$$RS = \text{Average Gain} / \text{Average Loss}$$

صورت کسر RS برابر میانگین تمام افزایش قیمتها و مخرج برابر مجموع قدر مطلق تمام کاهش قیمتها در ۱۴ روز گذشته می باشد. از آنجا که RS همواره عددی مثبت است به راحتی می توان دید که RSI عددی مابین صفر و یک خواهد شد. غالباً و به صورت تجربی، دو سطح کلیدی ۳۰ و ۷۰ را به ترتیب به عنوان سطوح اشباع فروش (فزونی پیروزی بر شکست) و اشباع خرید (فزونی شکست بر پیروزی) می نامند.

نوسان نما استوکستیک یا SO نیز دارای نوسانی در بازه ۰ تا ۱۰۰ می باشد. در محاسبه این نوسانگر دو مقدار $D\%$ و $K\%$ به صورت رابطه (۲) محاسبه می شوند.

(۲)

$$K\% = \frac{(\text{Current Close} - \text{Lowest Low})}{(\text{Highest High} - \text{Lowest Low})} \times 100$$

که $D\%$ میانگین متحرک ساده سه روزه $K\%$ می باشد. در محاسبات بالا کمترین و بیشترین قیمت همواره در یک دوره زمانی مشخص محاسبه می شوند که همگرایی ساخته می شود که از آن برای تولید سیگنال خرید و فروش استفاده می شود. بنابراین در صورتی که $\{MA_n(t)\}, \{MA_m(t)\}$ دو سری میانگین متحرک

$$s_t = \begin{bmatrix} p_t \\ o_t \\ c_t \\ h_t \\ l_t \\ rsi_t \\ so_t \\ macd_t \end{bmatrix}$$

(۴)

می‌شود و اگر $MA_n(t) > MA_m(t)$ سیگنالی برای فروش دریافت می‌شود.

پس از معرفی سه نوسانگر به بررسی مدل یادگیری تقویتی عمیق پرداخته می‌شود. یکی از اجزای تشکیل‌دهنده یادگیری تقویتی مجموعه حالات می‌باشد. در سیستم معاملاتی پژوهش با توجه به پیوستگی قیمت و شاخص‌های تکنیکی، متغیرهای حالت پیوسته می‌باشند و لذا یادگیری Q بر اساس جدول قابلیت کاربرد ندارد و به جای آن از یادگیری تقویتی عمیق استفاده می‌شود. بردار حالت عبارت است از

که در آن

جدول ۱: معرفی اجزای بردار حالت در یادگیری تقویتی عمیق

متغیر	تعریف
p_t	قیمت روزانه
o_t	قیمت باز شدن
c_t	قیمت بسته شدن
h_t	بیشترین قیمت روزانه
l_t	کمترین قیمت روزانه
rsi_t	شاخص قدرت نسبی
so_t	شاخص تصادفی
$macd_t$	شاخص همگرایی-واگرایی

ماهیت سیستم معاملاتی در بورس اوراق بهادار، سه عمل یا اقدام تعریف می‌گردد.

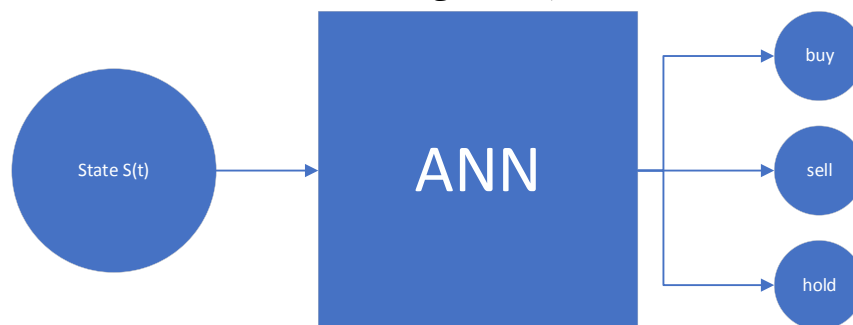
پس از معرفی مجموعه حالات در یادگیری تقویتی، به معرفی مجموعه اقدامات پرداخته می‌شود. با توجه به

جدول ۲: معرفی اجزای بردار اقدام

عمل	توصیف
<i>sell</i>	فروش به شرط اینکه دارایی قبلاً خریداری شده باشد(فاقد فروش استقراضی).
<i>buy</i>	خرید
<i>hold</i>	نگهداری

پیشنهاد به صورت سه احتمال با مجموع یک می باشد و پیشنهاد با حداکثر احتمال مورد پیاده سازی قرار می گیرد.

یادگیری تقویتی عمیق جدول تابع ارزش یا کیفیت Q را با یک شبکه عصبی جایگزین می کند. شبکه عصبی مذکور در نهایت با دریافت ورودی حالت، یکی از سه عمل فروش، خرید و نگهداری را پیشنهاد می کند. این



نمودار ۱: شبکه عصبی عمیق مورد استفاده در یادگیری تقویتی

t ، عامل یک عمل را انتخاب می کند a_t پاداشی را مشاهده می کند r_t وارد حالت جدیدی می شود s_{t+1} که ممکن است به حالت قبلی s_t و عملکرد انتخاب شده بستگی داشته باشد، به روز می شود. هسته اصلی الگوریتم معادله بلمن به عنوان یک به روزرسانی تکرار مقدار ساده است که از میانگین وزنی مقدار قدیمی و اطلاعات جدید استفاده می کند:

پاداش یک عمل در سیستم معاملاتی پژوهش، بازده ناشی از معامله می باشد. این پاداش در زمان فروش دارایی محاسبه خواهد شد و برای اقدامات نگهداری و خرید، پاداشی در نظر گرفته نمی شود. برای آموزش شبکه عصبی از معادله بلمن^۱ استفاده می شود. $Q: S \times A \rightarrow \mathbb{R}$ قبل از شروع یادگیری Q به یک مقدار ثابت احتمالاً دلخواه (انتخاب شده توسط برنامه نویس) مقداردهی اولیه می شود. سپس، در هر بار

¹ Bellman equation

$$Q^{new}(s_t, a_t) \leftarrow \underbrace{Q(s_t, a_t)}_{\text{old value}} + \underbrace{\alpha}_{\text{learning rate}} \cdot \underbrace{\left(\underbrace{r_t}_{\text{reward}} + \underbrace{\gamma}_{\text{discount factor}} \cdot \underbrace{\max_a Q(s_{t+1}, a)}_{\text{estimate of optimal future value}} - \underbrace{Q(s_t, a_t)}_{\text{old value}} \right)}_{\text{temporal difference}} \quad (5)$$

new value (temporal difference target)

تنظیم می‌شود. در بیشتر موارد، $Q(s_f, a)$ را می‌توان برابر با صفر در نظر گرفت.

ضریب یادگیری نشان می‌دهد که تا چه حد اطلاعات تازه به‌دست‌آمده بر اطلاعات قدیمی غلبه می‌کند. ضریب باعث می‌شود عامل چیزی یاد نگیرد (به‌طور انحصاری از دانش قبلی استفاده می‌کند)، درحالی‌که ضریب ۱ باعث می‌شود عامل فقط جدیدترین اطلاعات را در نظر بگیرد (نادیده گرفتن دانش قبلی برای کشف احتمالات).

در محیط‌های کاملاً قطعی، نرخ یادگیری $\alpha_t = 1$ بهینه است. هنگامی که مسئله تصادفی است، الگوریتم تحت برخی شرایط فنی روی نرخ یادگیری همگرا می‌شود که نیاز به کاهش آن به صفر دارد. در عمل، اغلب از یک نرخ یادگیری ثابت مانند $\alpha_t = 0.1$ استفاده می‌شود.

ضریب تخفیف γ اهمیت پاداش‌های آینده را تعیین می‌کند. ضریب ۰ عامل را تنها با در نظر گرفتن

پاداش‌های فعلی، یعنی r_t (در قانون به‌روزرسانی بالا) "کوتاه بین" (یا کوتاه بین) می‌کند، درحالی‌که ضریب نزدیک به ۱ باعث می‌شود آن را وادار کنید تا برای یک پاداش بلندمدت تلاش کند. اگر ضریب تخفیف برابر یا بیشتر از ۱ باشد، مقادیر عمل ممکن است متفاوت باشد.

برای $\gamma = 1$ ، بدون حالت پایانی، یا اگر عامل هرگز به یک حالت نرسد، تمام تاریخچه‌های محیطی بی‌نهایت طولانی می‌شوند و ابزارهای کاربردی با پاداش‌های اضافی

جایی که r_t پاداش دریافتی هنگام انتقال از حالت s_t به حالت s_{t+1} و α نرخ یادگیری است که $0 < \alpha \leq 1$. توجه داشته باشید که $Q^{new}(s_t, a_t)$ مجموع سه عامل است:

$(1-\alpha)Q(s_t, a_t)$: مقدار فعلی وزن‌دار شده توسط نرخ یادگیری. مقادیر نرخ یادگیری نزدیک به ۱ باعث می‌شود تغییرات در Q سریع‌تر شود.

αr_t : پاداش $r_t = r(s_t, a_t)$ برای دریافت اینکه آیا عمل a_t در حالت s_t (وزن دهی شده بر اساس میزان یادگیری) انجام می‌شود یا خیر.

$\alpha \gamma \max_a Q(s_{t+1}, a)$: حداکثر پاداشی که می‌توان از حالت s_{t+1} به دست آورد (وزن‌بندی شده بر اساس نرخ یادگیری و ضریب تخفیف)

یک اجرا از الگوریتم زمانی به پایان می‌رسد که حالت s_{t+1} حالت نهایی یا پایانی باشد. اگر ضریب تخفیف کمتر از ۱ باشد، مقادیر کنش محدود هستند حتی اگر مشکل بتواند حلقه‌های بی‌نهایت داشته باشد. برای همه حالت‌های نهایی s_f ، هرگز به‌روزرسانی $Q(s_f, a)$ نمی‌شود، اما روی مقدار پاداش r برای مشاهده s_f ،

عصبی اول را می‌گیرد و شبکه عصبی دوم از نو آموزش می‌بیند.

۴. یافته‌های پژوهش

مدل تحقیق از دو شبکه عصبی عمیق برای پیاده‌سازی جدول Q استفاده می‌کند. شبکه عصبی اول آموزش نمی‌بیند و فقط برای مقصد محاسبه احتمالات استفاده می‌شود. پس از محاسبه سمت راست معادله بلمن برای یک ورودی، از یک شبکه عصبی (که همان سیستم معاملاتی می‌باشد) برای آموزش استفاده می‌شود. بدین‌صورت که ورودی آن حالت و خروجی آن مقدار معادله بلمن می‌باشد که توسط شبکه عصبی اول آموزش دیده است. در فواصل زمانی منظم شبکه عصبی دوم، جای شبکه عصبی اول را می‌گیرد و شبکه عصبی دوم از نو آموزش می‌بیند. پارامترهای مورد استفاده در معادله بلمن در جدول ۳ ارائه شده است.

جدول ۳: تنظیمات مورد استفاده در معادله بلمن

مقدار	پارامتر
۰/۱	ضریب یادگیری α
۰/۹۵	ضریب تحقیق یا تنزیل γ

تنظیمات مورد استفاده در تعادل بین جستجوی تصادفی و بهره‌گیری از اطلاعات موجود در جدول ۴ ارائه شده است.

مشخصات شبکه‌های عصبی عمیق مورد استفاده در یادگیری تقویتی عمیق در جدول ۵ ارائه شده است.

و تخفیف ناپذیر عموماً بی‌نهایت می‌شوند. حتی با یک ضریب تخفیف فقط کمی کمتر از ۱، یادگیری تابع Q منجر به انتشار خطاها و ناپایداری‌ها می‌شود، زمانی که تابع مقدار با یک شبکه عصبی مصنوعی تقریب زده شود. در این صورت، شروع با ضریب تخفیف کمتر و افزایش آن به سمت ارزش نهایی، یادگیری را تسریع می‌کند. برای طراحی سیستم معاملاتی بر پایه یادگیری تقویتی عمیق از دو شبکه عصبی استفاده می‌شود. شبکه اول برای محاسبه مقادیر Q در سمت راست معادله بلمن استفاده می‌شود. این شبکه عصبی آموزش نمی‌بیند و فقط برای مقصد محاسبه احتمالات استفاده می‌شود. پس از محاسبه سمت راست معادله بلمن برای یک ورودی، از یک شبکه عصبی (که همان سیستم معاملاتی می‌باشد) برای آموزش استفاده می‌شود. بدین‌صورت که ورودی آن حالت و خروجی آن مقدار معادله بلمن می‌باشد که توسط شبکه عصبی اول آموزش دیده است. در فواصل زمانی منظم شبکه عصبی دوم، جای شبکه

در فرآیند یادگیری تقویتی، در گام‌های نخست به جستجوی تصادفی اهمیت بیشتری داشته می‌شود و کم‌کم با جمع‌آوری اطلاعات از محیط و یادگیری بیشتر شبکه عصبی عمیق، به بهره‌گیری از اطلاعات کسب شده اهمیت داده می‌شود. برای این منظور پارامتر اپسیلون تعریف می‌شود که این ضریب را شبیه‌سازی می‌کند و به تدریج باگذشت زمان از مقدار آن کم می‌شود.

جدول ۴: تنظیمات جستجو

مقدار	پارامتر
۰/۹۵	ضریب توازن اولیه جستجوی تصادفی و بهره‌گیری ϵ

-۰/۰۰۱

مقدار کاهش ضریب توازن در هر اجرا

جدول ۵: مشخصات شبکه‌های عصبی عمیق مورد استفاده در یادگیری تقویتی عمیق

مقدار	پارامتر
شبکه عصبی عمیق پرسپترون	نوع
۳	تعداد لایه
Relu	تابع فعالساز لایه‌های میانی
سافت مکس	تابع فعالساز لایه آخر
۶	تعداد نرون لایه ورودی
۳	تعداد نرون لایه خروجی
آدام ^۱	الگوریتم بهینه ساز

تست مورد ارزیابی قرار گرفت. بدین صورت زمانی که خروجی شبکه عصبی برابر خرید می‌باشد، شاخص کل خریداری و یک موقعیت معاملاتی باز می‌شود و در ادامه زمانی که خروجی شبکه عصبی عمیق برابر فروش می‌باشد، شاخص خریداری شده، فروخته و موقعیت معاملاتی بسته می‌شود و بازده حاصل شده به‌عنوان یک نمونه از عملکرد سودآوری سیستم معاملاتی پژوهش ذخیره می‌گردد. مشخصات آماری موقعیت‌های معاملاتی حاصل شده در جدول ۶ ارائه شده است.

در هر ۱۰ تکرار، جای شبکه آموزش دیده جای شبکه پیش‌بینی را می‌گیرد. تعداد داده‌های بازده روزانه شاخص کل در بازه پژوهش ۲۶۵۴ داده می‌باشد که از ۱۶۵۴ داده برای آموزش مدل پژوهش استفاده گردید و ۱۰۰۰ داده نیز برای آزمون و بررسی عملکرد مدل پژوهش در داده‌های برون نمونه‌ای استفاده گردید. برای آموزش مدل پژوهش از ۱۰۰۰ تکرار استفاده گردید.

پس از آموزش سیستم معاملاتی پژوهش بر روی ۱۶۵۴ داده آموزشی روزانه، این سیستم بر روی ۱۰۰۰ داده

جدول ۶: مشخصات آماری معادل بازده روزانه سیستم معاملاتی پژوهش

مقدار	سیستم معاملاتی
۱۳	تعداد موقعیت معاملاتی

^۱ Adam

بیشینه بازده موقعیت	۰/۰۱۰۴
کمینه بازده موقعیت	-۰/۰۰۱۵
متوسط طول زمانی موقعیت	۲۱/۵۳۸۵
متوسط معادل بازده روزانه	۰/۰۰۱۸
انحراف معیار معادل بازده روزانه	۰/۰۰۳۷
نسبت شارپ	۰/۴۷۹۳

تعدیل شده نیز گفته می‌شود و بیشتر بودن این شاخص، مطلوب سرمایه‌گذاران می‌باشد.

در ادامه سودآوری سیستم‌های معاملاتی انفرادی شکل گرفته بر اساس شاخص‌های تکنیکی نیز مورد ارزیابی قرار می‌گیرد. در ابتدا سیستم معاملاتی شکل گرفته بر اساس شاخص قدرت نسبی با خرید در ۳۰ و فروش در ۷۰ بر روی داده‌های آزمون مورد محاسبه و شبیه‌سازی قرار گرفت که نتیجه در جدول ۷ ارائه شده است.

بر اساس اطلاعات حاصل شده، ۱۳ موقعیت معاملاتی توسط سیستم پژوهش شناسایی گردید که متوسط بازده روزانه‌ای برابر ۰/۰۰۱۸ را تولید می‌کند که این مقدار برابر متوسط بازده سالیانه‌ای برابر ۰,۹۱ می‌باشد. متوسط زمان هر معامله تقریباً برابر ۲۲ روز می‌باشد. نسبت شارپ سیستم پژوهش که از تقسیم بازده بر انحراف معیار محاسبه می‌شود برابر ۰/۴۷۹۳ می‌باشد و نشان می‌دهد که در هر واحد ریسک، بازده ای برابر ۰/۴۷۹۳ تولید می‌شود. شاخص نسبت شارپ بازده در مقیاس ریسک را محاسبه می‌کند و به آن بازده

جدول ۷: مشخصات آماری معادل بازده روزانه سیستم معاملاتی شاخص قدرت نسبی

مشخصات سیستم قدرت نسبی	مقدار
تعداد موقعیت معاملاتی	۱۴
بیشینه بازده موقعیت	۰/۰۰۷۸
کمینه بازده موقعیت	-۰/۰۰۳۳
متوسط طول زمانی موقعیت	۲۵/۷۷۲۷
متوسط معادل بازده روزانه	۰/۰۰۰۴
انحراف معیار معادل بازده روزانه	۰/۰۰۲۷
نسبت شارپ	۰/۱۴۹۳

روزانه‌ای برابر ۰/۰۰۰۴ را تولید می‌کند که این مقدار برابر متوسط بازده سالیانه‌ای برابر ۰/۱۶ می‌باشد. متوسط زمان

بر اساس اطلاعات حاصل شده، ۱۴ موقعیت معاملاتی توسط سیستم پژوهش شناسایی گردید که متوسط بازده

هر معامله تقریباً برابر ۲۵ روز می‌باشد. نسبت شارپ سیستم پژوهش که از تقسیم بازده بر انحراف معیار محاسبه می‌شود برابر ۰/۱۴۹۳ می‌باشد و نشان می‌دهد که در هر واحد ریسک، بازده ای برابر ۰/۱۴۹۳ تولید می‌شود.

در ادامه سیستم معاملاتی شکل گرفته بر اساس شاخص تصادفی با خرید در ۳۰ و فروش در ۷۰ بر روی داده‌های آزمون مورد محاسبه و شبیه‌سازی قرار گرفت که نتیجه در جدول ۸ ارائه شده است.

جدول ۸: مشخصات آماری معادل بازده روزانه سیستم معاملاتی شاخص تصادفی

مقدار	مشخصات سیستم معاملاتی شاخص تصادفی
۱۱	تعداد موقعیت معاملاتی
۰/۰۰۶۷	بیشینه بازده موقعیت
-۰/۰۰۲۳	کمینه بازده موقعیت
۲۰/۲۷۲۷	متوسط طول زمانی موقعیت
۰/۰۰۰۵۷	متوسط معادل بازده روزانه
۰/۰۰۳۰	انحراف معیار معادل بازده روزانه
۰/۱۹۱۲	نسبت شارپ

بر اساس اطلاعات حاصل شده، ۱۱ موقعیت معاملاتی توسط سیستم پژوهش شناسایی گردید که متوسط بازده روزانه‌ای برابر ۰/۰۰۰۵۷ را تولید می‌کند که این مقدار برابر متوسط بازده سالیانه‌ای برابر ۰/۲۳ می‌باشد. متوسط زمان هر معامله تقریباً برابر ۲۰ روز می‌باشد. نسبت شارپ سیستم پژوهش که از تقسیم بازده بر انحراف

معیار محاسبه می‌شود برابر ۰/۱۹۱۲ می‌باشد و نشان می‌دهد که در هر واحد ریسک، بازده ای برابر ۰/۱۹۱۲ تولید می‌شود. در ادامه به عملکرد سیستم معاملاتی شکل گرفته بر اساس شاخص همگرایی-واگرایی پرداخته می‌شود که نتیجه عملکرد آن در جدول ۹ ارائه شده است.

جدول ۹: مشخصات آماری معادل بازده روزانه سیستم معاملاتی شاخص همگرایی-واگرایی

مقدار	مشخصات سیستم معاملاتی شاخص همگرایی-واگرایی
۱۴	تعداد موقعیت معاملاتی
۰/۰۰۹۷۲۹	بیشینه بازده موقعیت
-۰/۰۱۶۵۳	کمینه بازده موقعیت
۳۲/۵۱۳۵۱	متوسط طول زمانی موقعیت
-۰/۰۰۰۴۴	متوسط معادل بازده روزانه
۰/۰۰۴۷۹۳	انحراف معیار معادل بازده روزانه
-۰/۰۹۱۰۳	نسبت شارپ

عملکرد چهار سیستم معاملاتی شامل یادگیری تقویتی عمیق، مبتنی بر شاخص قدرت نسبی، مبتنی بر شاخص تصادفی و مبتنی بر شاخص همگرایی و واگرایی مورد ارزیابی قرار گرفت. جدول ۱۰ نتایج چهار استراتژی معاملاتی را به صورت یکجا نشان می‌دهد.

بر اساس اطلاعات حاصل شده، ۱۴ موقعیت معاملاتی توسط سیستم پژوهش شناسایی گردید که متوسط بازده روزانه‌ای برابر $-0/00044$ را تولید می‌کند. از این رو سیستم شکل گرفته بر اساس شاخص همگرایی-واگرایی زیان ده می‌باشد.

جدول ۱۰: مشخصات آماری چهار استراتژی مورد استفاده در پژوهش

همگرایی- واگرایی	تصادفی	قدرت نسبی	یادگیری تقویتی عمیق	سیستم معاملاتی شاخص عملکرد
۱۴	۱۱	۱۴	۱۳	تعداد موقعیت معاملاتی
$0/009729$	$0/0067$	$0/0078$	$0/0104$	بیشینه بازده موقعیت
$-0/01653$	$-0/0023$	$-0/0033$	$-0/0015$	کمینه بازده موقعیت
$32/51351$	$20/2727$	$25/7727$	$21/5385$	متوسط طول زمانی موقعیت
$-0/00044$	$0/00057$	$0/0004$	$0/0018$	متوسط معادل بازده روزانه
$0/004793$	$0/0030$	$0/0027$	$0/0037$	انحراف معیار معادل بازده روزانه
$-0/09103$	$0/1912$	$0/1493$	$0/4793$	نسبت شارپ

جدول ۱۱: مقدار احتمال آزمون برابری میانگین برای چهار سیستم معاملاتی

شاخص همگرایی- واگرایی	شاخص تصادفی	شاخص قدرت نسبی	سیستم پژوهش	نوسانگر
$0/000$	$0/00912$	$0/0069$	سیستم پژوهش	
$0/000$	$0/873$		شاخص قدرت نسبی	
$0/000$			شاخص تصادفی	
			شاخص همگرایی-واگرایی	

روزانه کسب شده مورد به کمک آزمون مقایسات زوجی بررسی قرار گرفت.

در ادامه عملکرد سیستم‌های معاملاتی چهارگانه پژوهش از بابت تفاوت معناداری آماری بین میانگین‌های بازده

بر اساس اطلاعات جدول ۱۰ و ۱۱ سیستم پژوهش در میانگین دارای تفاوت معناداری با سه سیستم دیگر می‌باشد. سیستم شکل گرفته بر اساس شاخص قدرت نسبی و تصادفی دارای تفاوت معناداری نمی‌باشند و سیستم شکل گرفته بر اساس شاخص همگرایی-واگرایی دارای تفاوت معنادار با سه سیستم دیگر می‌باشد. همچنین نسبت شارپ سیستم پژوهش نسبت به سه مدل دیگر رشد حداقل ۱/۴ برابری را نشان می‌دهد.

۵. بحث و نتیجه‌گیری

در پژوهش حاضر یک سیستم معاملاتی بر اساس رویکرد یادگیری تقویتی و به کمک شبکه‌های عصبی عمیق طراحی گردید. در این رویکرد عامل یا معامله‌گر در فضای جستجو برای یافتن پاداش بیشتر که همان بازده بالاتر می‌باشد، به جستجو می‌پردازد. عامل معامله‌گر با سیگنال‌های تکنیکی شامل شاخص قدرت نسبی، نوسانگر تصادفی، شاخص همگرایی-واگرایی و قیمت‌های کمینه، بیشینه، بسته و باز شدن مواجه می‌شود. عامل این اطلاعات را در گام‌های زمانی متوالی دریافت می‌کند و بر اساس آن‌ها یک اقدام انجام می‌دهد که در اینجا شامل خرید، نگهداری و یا فروش دارایی می‌باشد. یادگیری تقویتی به دنبال طراحی یک استراتژی الگوریتمیک یا خودکار می‌باشد که در نتیجه آن میزان پاداش یعنی بازده خرید و فروش بیشینه شود که یکی از الگوریتم‌های مشهور در این زمینه Q -یادگیری می‌باشد. غالباً برای پیاده‌سازی یادگیری تقویتی از جدول Q استفاده می‌شود که در آن کیفیت اقدام‌های مختلف در وضعیت‌های مختلف در جدول ضبط و رکورد می‌شود. در مدل مورد بررسی پژوهش با توجه به ماهیت پیوسته متغیرهای حالت، جدول Q قابل استفاده نیست و به جای آن از یک شبکه عصبی عمیق استفاده می‌شود تا استراتژی معاملاتی یادگیری تقویتی مدنظر پژوهش را آموزش ببیند. در نهایت شبکه عصبی آموزش دیده در مرحله آزمایش می‌تواند با دریافت اطلاعات تکنیکی و

قیمتی مذکور، یکی از سه اقدام خرید، نگهداری و فروش را پیشنهاد دهد. در راستای سؤالات پژوهش به دو سؤال پاسخ داده شد.

۱- چگونه می‌توان به کمک یادگیری تقویتی عمیق یک سیستم معاملاتی طراحی کرد؟

برای طراحی سیستم معاملاتی بر پایه یادگیری تقویتی عمیق از دو شبکه عصبی استفاده می‌شود. شبکه اول برای محاسبه مقادیر Q در سمت راست معادله بلمن استفاده می‌شود. این شبکه عصبی آموزش نمی‌بیند و فقط برای مقصد محاسبه احتمالات استفاده می‌شود. پس از محاسبه سمت راست معادله بلمن برای یک ورودی، از یک شبکه عصبی (که همان سیستم معاملاتی می‌باشد) برای آموزش استفاده می‌شود. بدین صورت که ورودی آن حالت و خروجی آن مقدار معادله بلمن می‌باشد که توسط شبکه عصبی اول آموزش دیده است. در فواصل زمانی منظم شبکه عصبی دوم، جای شبکه عصبی اول را می‌گیرد و شبکه عصبی دوم از نو آموزش می‌بیند.

۲- سودآوری سیستم معاملاتی طراحی شده چگونه است؟

سیستم معاملاتی پژوهش پس از آموزش بر روی ۱۶۵۴ داده، بر روی ۱۰۰۰ داده تست مورد ارزیابی سودآوری قرار گرفت. نتایج نشان داد که سیستم پژوهش در میانگین دارای تفاوت معناداری با سه سیستم دیگر می‌باشد. سیستم شکل گرفته بر اساس شاخص قدرت نسبی و تصادفی دارای تفاوت معناداری نمی‌باشند و سیستم شکل گرفته بر اساس شاخص همگرایی-واگرایی دارای تفاوت معنادار با سه سیستم دیگر می‌باشد. همچنین نسبت شارپ سیستم پژوهش نسبت به سه مدل دیگر رشد حداقل ۱/۴ برابری را نشان می‌دهد. با توجه به نتایج حاصل شده از سیستم معاملاتی پژوهش پیشنهادهای زیر ارائه می‌گیرد.

۱- با توجه به نسبت شارپ بالای سیستم پژوهش نسبت به سیستم‌های انفرادی شکل گرفته بر اساس ابزارهای تکنیکی به سرمایه‌گذاران ریسک‌گریز پیشنهاد می‌شود تا از سیستم معاملاتی پژوهش برای ترکیب سیگنال‌های تکنیکی و اتخاذ موقعیت‌های معاملاتی استفاده کنند.

۲- توجه شود که عملکرد سودآوری سیستم معاملاتی پژوهش باید با توجه به سهام موردنظر مورد بر روی داده‌های آزمایشی مورد ارزیابی قرار گیرد. همچنین پیشنهاد می‌شود تا با تغییر ورودی‌های مدل و استفاده از طیف وسیعی از متغیرهای تکنیکی و بنیادی دقت سیستم معاملاتی ارزیابی و بهترین مدل انتخاب گردد.

۶. منابع

اسفندیار، مهدی، کرامتی، محمد علی، غلامی جمکرانی، رضا، کاشفی نیشابوری، محمدرضا (۱۴۰۱). بهینه‌سازی پرتفوی سهام در بورس اوراق بهادار تهران (کاربرد رهیافت یادگیری تقویتی). فصلنامه علمی مدل‌سازی اقتصادی، ۱۶(۵۸)، ۵۱-۶۶.

خنجرپناه، حسین، دوروش، داود، شوال پور، سعید، جبارزاده، آرمن. (۱۳۹۷). کاربرد روش تکنیکال برای پیش‌بینی قیمت سهام: رویکرد مدل‌های احتمال غیرخطی و شبکه‌های عصبی مصنوعی. راهبرد مدیریت مالی، ۶(۳)، ۵۹-۷۹.

رستگار، محمد علی، صداقتی‌پور، امین. (۱۳۹۷). ارائه سیستم معاملات الگوریتمی برای قرارداد آتی سکه طلا مبتنی بر داده‌های درون-روزی. دانش سرمایه‌گذاری، ۲۸(۲)، ۴۹-۶۸.

سارنج، علیرضا، قاسمی، احمدرضا، ارم، اصغر، تهرانی، رضا. (۱۳۹۹). طراحی سیستم معاملات تکنیکی سهام با استفاده از مدل ترکیبی شبکه عصبی MLP و الگوریتم‌های تکاملی. دانش مالی تحلیل اوراق بهادار، ۱۳(۴۵)، ۴۷-۶۴.

شمس، شهاب‌الدین؛ عطایی، بهروز. (۱۳۹۵). شناسایی دستکاری قیمت سهام از طریق مدل ترکیبی الگوریتم ژنتیک - شبکه عصبی مصنوعی و مدل SQDF. راهبرد مدیریت مالی، ۴(۳)، ۱۴۹-۱۷۱.

صیادی نژاد، اسماعیل زاده، & رستمی. (۲۰۲۳). ارائه مدل پیش‌بینی بازدهی بیت‌کوین با استفاده از روش هیبریدی یادگیری عمیق-الگوریتم تجزیه سیگنال (CEEMD-DL). اقتصاد مالی، ۱۷(۶۲)، ۲۱۷-۲۳۸.

صیادی نژاد، سکینه، اسماعیل زاده، علی، رستمی، محمدرضا (۱۴۰۲). ارائه مدل پیش‌بینی بازدهی بیت‌کوین با استفاده از روش هیبریدی یادگیری عمیق-الگوریتم تجزیه سیگنال (CEEMD-DL). اقتصاد مالی، ۱۷(۶۲)، ۲۱۷-۲۳۸.

یافتیان، امیرحسین، رستگار، محمدعلی. (۱۳۹۹). طراحی یک سیستم معاملاتی خودکار با استفاده از شبکه عصبی پیچشی. چشم انداز مدیریت مالی، ۱۰(۳۱)، ۱۸۴-۱۵۳.

Jin, B. (2023). *A Mean-VaR Based Deep Reinforcement Learning Framework for Practical Algorithmic Trading*. *IEEE Access*, 11, 28920-28933.

Koratamatti, P., Wadhvani, K., Gupta, M., & Sanjeevi, S. G. (2021). *Market sentiment-aware deep reinforcement learning approach for stock portfolio allocation*. *Engineering Science and Technology, an International Journal*, 24(4), 848-859.

Li, Y., Liu, P., & Wang, Z. (2022). *Stock Trading Strategies Based on Deep Reinforcement Learning*. *Scientific Programming*, 2022.

Nabipour, M., Nayyeri, P., Jabani, H., Mosavi, A and Salwana, E. and S., Shahab. (2020). *Deep Learning for Stock Market Prediction*. *Entropy* 2020, 22(8), 840.

Using Convolutional Neural Network and Long Short-Term Memory. Mathematical Problems in Engineering, 20,1-13.

Yu, P & Yan, X. (2020). Stock price prediction based on deep neural networks. Neural Computing and Applications. 32. 10.

Zhang, J., & Lei, Y. (2022). Deep Reinforcement Learning for Stock Prediction. Scientific Programming, 2022.

Thibaut, T., Damien, E (2020). An Application of Deep Reinforcement Learning to Algorithmic Trading, ArXiv2021.

Wang,L.,Wang.Z and Zhao,SH. (2015). Stock market trend prediction using dynamical Bayesian factor graph, Expert Systems with Applications,45,6263-6270.

Xiao-dan Zhang, Ang Li, Ran Pan, (2016),Stock trend prediction based on a new status box method and AdaBoost probabilistic support vector machine. Applied Soft Computing 49 , 385–398.

Xu,Y., and V. Keselj. (2019). Stock Prediction using Deep Learning and Sentiment Analysis, 2019 IEEE International Conference on Big Data (Big Data), Los Angeles, CA, USA, 2019, pp. 5573-5580.

Yang, C., Zhai, J., Tao, G. (2020). Deep Learning for Price Movement Prediction

Designing algorithmic trading strategy based on deep reinforcement learning

Case study: Tehran Stock Exchange

Saeed Kazemian hossinabadi¹
Sayyed Mohammad Reza Davoodi^{2}*
Mohammad Mashhadizadeh³
Parsa Jozzi⁴

Abstract

Today, algorithmic trading is widely used in trading management. Algorithmic portfolio management is a new type of these systems through which the portfolio manager helps to increase the quality of profit and reduce the risks of his portfolio using algorithmic tools. The purpose of this research is to design an algorithmic trading system based on deep reinforcement learning with the help of neural network. In this approach, the agent or trader searches the search space to find more reward, which is the same as more return. The trader is faced with technical signals including relative strength index, stochastic oscillator, convergence-divergence indicator and minimum, maximum, closing and opening prices. Deep reinforcement learning replaces the Q value or quality function table with a neural network. Finally, upon receiving the state word, the mentioned neural network suggests one of the three actions of selling, buying and holding. This proposal is in the form of three possibilities with a total of one, and the proposal with the maximum probability is implemented. The result of the implementation of the deep reinforcement learning trading system on the total index of the Tehran Stock Exchange in the period of 2013 to 2014 shows that the research system had a significant difference with the other three systems. The system formed based on the relative and random power index has no significant difference, and the system formed based on the convergence-divergence index has a significant difference with the other three systems. Also, the Sharpe ratio of the research system compared to the other three models showed a growth of at least 1.4 times.

Keywords: *neural networks, reinforcement learning, deep reinforcement learning, technical oscillator.*

¹Master of Financial Management, Dehaghan Branch, Islamic Azad University, Dehaghan, Iran. Saeedkazemian1365@yahoo.com

²Assistant Professor, Dehaghan Branch, Islamic Azad University, Dehaghan, Iran. (Corresponding author). Smrdavoodi@ut.ac.ir

³Assistant Professor, Mobarakeh Branch, Islamic Azad University, Mobarakeh, Iran. m.mashhadi@mau.ac.ir

⁴Master of Financial Management, Dehaghan Branch, Islamic Azad University, Dehaghan, Iran. Parsajozipersonal@gmail.com