

A Consensus-based Auto-scaling Approach For Serverless Environments

Mobina Kashaniyan, Mehrdad Ashtiani, Amirhossein Ghassemi

Iran University of Science and Technology, School of Computer Engineering
mobina_kashanian@comp.iust.ac.ir , m_ashtiani@iust.ac.ir , a_ghassemi@comp.iust.ac.ir

Abstract

Efficient management of computing resources has always been a significant concern for users. An automatic scaling system can help in managing hardware resources by adapting to the system's performance history. It can increase or decrease resources automatically, without human intervention, based on predefined criteria. This ensures smooth program execution without any disruption caused by changes in the operating environment. This study focuses on serverless environments, which rely on functions. We model these functions using graph theory, analyze their dependencies, and identify the most critical bottlenecks in the graph. We then use two approaches, supervised and unsupervised, to predict the scalability of bottleneck resources. To be more sure of the scaling decision, the consensus mechanism compares the predictions of the models, and the best model's result is considered the final scaling decision, which creates consistency between the results obtained from the methods. Results show that supervised approaches perform better than unsupervised approaches in the automatic scaling problem. The models implemented in this research can determine the scaling result with 98% accuracy, which is a 2.5% improvement compared to previous works.

Keywords: Autoscaling, Cloud Environment, Machine Learning, Workload Prediction, Ensemble Learning.

ارائه یک رویکرد مقیاس پذیری خودکار منابع در محیط های بدون سرویس دهنده

مبینا کاشانیان، مهرداد آشتیانی، امیرحسین قاسمی

دانشگاه علم و صنعت ایران، دانشکده مهندسی کامپیوتر

Mobina_kashanian@comp.iust.ac.ir , m_ashtiani@iust.ac.ir , a_ghassemi@comp.iust.ac.ir

چکیده

مدیریت صحیح منابع همواره یکی از دغدغه های کاربران هنگام استفاده از محیط های محاسباتی بوده است. یک سامانه مقیاس پذیر خودکار می تواند با توجه به تاریخچه ی عملکرد سامانه، منابع سخت افزاری را مدیریت کند و در زمان مناسب برای جلوگیری از اختلال در روند اجرای برنامه با وفق دادن سامانه با محیط عملیاتی جدید منابع را به صورت خودکار، بدون دخالت انسان و بر اساس معیارهای تنظیم شده افزایش یا کاهش دهد. در این پژوهش، با در نظر گرفتن ویژگی منحصر به فرد محیط های بدون سرویس دهنده که فعالیتی بر اساس توابع دارند، توابع را با نظریه گراف مدل سازی کرده و با تحلیل وابستگی های آنها و پیدا کردن مهم ترین گلوگاه در گراف، دو رویکرد با نظارت و بدون نظارت را برای پیش بینی مقیاس پذیری منابع گلوگاه در نظر می گیریم و برای اطمینان بیشتر به تصمیم مقیاس پذیری، پیش بینی مدل ها توسط سازوکار اجماع نظر مقایسه می گردد و نتیجه ی بهترین مدل به عنوان تصمیم نهایی مقیاس پذیری در نظر گرفته می شود که به نوعی میان نتایج حاصل شده از روش ها سازگاری ایجاد کند. نتایج نشان می دهند که رویکردهای با نظارت در مقایسه با رویکردهای بدون نظارت در مسئله مقیاس پذیری خودکار بهتر عمل می کنند و مدل های پیاده سازی شده در این پژوهش، می توانند تا دقت ۹۸٪ نتیجه ی مقیاس پذیری را تعیین کنند که نسبت به کارهای پیشین انجام شده بهبود ۲۰۵ درصدی داشته است.

کلمات کلیدی

مقیاس پذیری خودکار، محیط بدون سرویس دهنده، یادگیری ماشین، رویکرد بدون نظارت، رویکرد با نظارت

می تواند منجر به افزایش هزینه ها شود، یا برعکس در اثر تخصیص منابع کمتر منجر به کاهش عملکرد سامانه شود که بسیار آسیب زنده است. روش دیگر استفاده از راه حل های مقیاس پذیر خودکار است که این روش راهکاری هوشمند برای تغییر پویا مقیاس و تأمین نیازهای منابع است، این روش باید بتواند با یادگیری شرایط محیطی، نیازهای منابع آینده را پیش بینی کند [۳]. طرح تخصیص منابع باید مقیاس پذیر، قابل انطباق و قابل اطمینان در برابر تغییرات در حجم کار باشد که نیازهای منابع آینده را به درستی پیش بینی کند و اقدامات مقیاس پذیری لازم را انجام دهد نگرانی های روش قبل را برطرف کند. اما اثربخشی چنین راه حل هایی به شدت به نوع مدل پیش بینی و همچنین کیفیت و کمیت داده های آموزشی بستگی دارد و هنگامی که منابع مورد نیاز به طور مناسب تخصیص داده شوند، عملکرد و کارایی سامانه به دلیل کمبود منابع کاهش نمی یابد. همچنین، سامانه های نرم افزاری بزرگ نیاز به ارائه درجه بالایی از اطمینان در کیفیت خدمات مانند زمان پاسخ گویی، توان عملیاتی و در

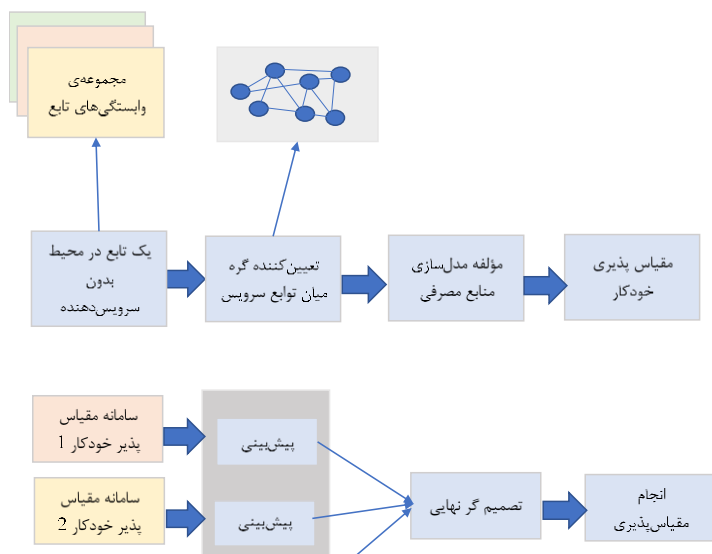
۱- مقدمه

در فضای محاسبات بدون سرویس دهنده سامانه هایی وجود دارند که می توانند پویا باشند و مقیاس و اندازه خود را در هر زمانی تغییر دهند. این سرویس ها نیاز دارند که آزادی عملکرد داشته و منابع خود را در هر لحظه ارتقا دهد تا سامانه بتواند بدون اختلال کار کند، با توجه به این امر امروزه نیاز به مدیریت منابع به صورت خودکار در این سامانه ها بسیار مشهود است [۱]. این در حالی است که سایر محیط های محاسباتی مانند خوشه ها و گریدها این امکان را به سادگی به کاربران نمی دهند و ارتقا منابع در این سامانه ها نیازمند روش های پیشرفته تری نسبت به فضاهای ابری و بدون سرویس دهنده است [۲]. یکی از روش های مقیاس پذیری خودکار بر اساس کنش و واکنش تصمیم به افزایش یا کاهش منابع می کند، این روش می تواند بر کارایی، عملکرد و هزینه سامانه های نرم افزاری از طریق تخصیص منابع اضافی اثر بگذارد و حتی

منابع بهینه عمل کرد [8]. آگاروال و همکاران یک عامل یادگیری تقویتی بازگشتی را برای مقیاس خودکار تابع بررسی می کنند و آن را با الگوریتم (PPO) ^۱ مقایسه می کنند و یک شبکه LSTM با پیشرفته ترین الگوریتم PPO ارائه کردند و نشان دادند که روابط بازگشتی برای مقیاس پذیری توابع بهتر عمل می کند [9].

۳- شرح الگوریتم پیشنهادی

پس از شناخت دقیق محیط بدون سرویس دهنده به یکی از مهم ترین ویژگی های این محیط که بر اساس توابع فعالیت دارند می رسمیم. در دنیای واقعی، توابع به یکدیگر وابستگی دارند و این وابستگی ها باعث کاهش عملکرد سامانه و افزایش پیچیدگی می شود و هرچقدر که این وابستگی ها بیشتر شوند یک جهنم وابستگی را ایجاد می کنند. در این پژوهش، توابع را با استفاده از نظریه گراف مدل سازی می کنیم و پس از آن مهم ترین گره یا گلوگاه گراف را پیدا می کنیم و منابع آن گره را تحت نظر می گیریم و سپس با استفاده از تعیین آستانه برای آن منابع مقیاس کاهش می دهیم، افزایشی یا بدون تغییر را در نظر می گیریم و بعد از آن پیش بینی مقیاس پذیری خودکار و منابع آینده را با سه مدل هوش مصنوعی پیشروی، پرسپترون، حافظه طولانی کوتاه مدت روی منابع آن گره انجام می دهیم و در نهایت میان نتایج به یک اجماع نظر می رسمیم که مصرف آینده گره چه خواهد شد و مقیاس پذیری آنچه می شود این کار را انجام می دهیم تا مهم ترین گره گراف که بیشترین وابستگی را دارد به کمبود منابع دچار نشود و یا در صورت هدر رفت منابع آن را کاهش دهیم تا بقیه گره ها به کمبود منابع دچار نشوند. روش ارائه شده از مؤلفه تعیین کننده گره میان توابع سرویس، مؤلفه مدل سازی منابع مصرفی و مؤلفه تصمیم گیرنده مقیاس پذیری تشکیل شده است. شکل (۱) معماری کلی ایده پژوهش را نشان می دهد.



شکل ۱: معماری کلی ایده پژوهش.

دسترس بودن خدمات دارند، بدون چنین اطمینانی، ارائه دهندگان خدمات ممکن است کاربران خود را از دست بدهند و با کاهش درآمد مواجه شوند [4]. در این پژوهش، هدف ارائه یک سازوکار مقیاس پذیری خودکار در محیط های بدون سرویس دهنده است به نحوی که با کمک یادگیری ماشین بتوان مدلی ارائه داد که با ایجاد اجماع نظر میان نتایج، عملکرد پیش بینی تخصیص منابع را بهبود ببخشد. به طور خلاصه می توان به نوآوری های اصلی این کار به شرح زیر اشاره کرد:

۱. در نظر گرفتن ویژگی اصلی یک محیط بدون سرویس دهنده و مدل سازی فراخوانی توابع با استفاده از تئوری گراف.
 ۲. یافتن مهم ترین گره در گراف برای نظارت بر منابع آن ها.
 ۳. استفاده از سه مدل یادگیری ماشین برای پیش بینی نیازهای منابع آینده گرهی مهم.
 ۴. استفاده از الگوریتم اجماع برای یافتن بهترین مدل از بین مقادیر پیش بینی شده برای اطمینان در تصمیم گیری.
- ساختار این مقاله به شرح زیر است: بخش ۲ پیشینه تحقیق را مورد بحث قرار می دهد. رویکرد پیشنهادی در بخش ۳ ارائه شده است و در بخش ۴ به صورت جامع مراحل الگوریتم ارائه شده است. ارزیابی روش های پیشنهادی در بخش ۵ ارائه و مورد بحث قرار گرفته است. در نهایت، بخش ۶ مقاله را با ارائه نتایج به پایان می رساند.

۲- پیشینه تحقیق

این بخش به بررسی پژوهش های مرتبط و روش های مقیاس پذیری خودکار می پردازد. روش های مقیاس پذیری خودکار اهدافی مانند بهبود دقت پیش بینی، انطباق با SLA یا هزینه منابع را مورد هدف قرار می دهند که بر اساس روش آن پژوهش ها طبقه بندی می شوند. قوانین مبتنی بر آستانه به طور گسترده برای بهینه سازی استفاده از منابع استفاده می شود و اجرای موفقیت آمیز این قوانین مستلزم توجه دقیق به جزئیات و درک جامع مسئله است. ناتوانی در تنظیم دقیق آستانه ها می تواند منجر به تخصیص منابع غیر بهینه شود و باعث کاهش عملکرد یا ایجاد هزینه های غیر ضروری شود. رومی و همکاران نشان می دهند که آستانه ها باید به دقت تنظیم شوند تا از نوسانات در سامانه جلوگیری شود [5]. یادگیری تقویتی بدون هیچ دانش قبلی با توجه به حجم کاری ورودی قادر به تعیین بهترین اقدام مقیاس پذیری برای هر برنامه هستند. سیستم الستیک داکر برای مقیاس پذیری عمودی از یادگیری تقویتی برای بهینه سازی منابع استفاده می کند [6]. تئوری صف اغلب برای مدل سازی برنامه های کاربردی اینترنتی استفاده می شود و برای تخمین معیارهای کارایی مانند طول صف یا میانگین زمان انتظار مفید است. نظریه کنترل روشی برای مدیریت و مقیاس پذیری خودکار سامانه ها، نگه داشتن بار پردازنده و سایر متغیرهای کنترل شده در سطوح دلخواه از طریق تنظیمات ورودی و محاسبات ریاضی است. تجزیه و تحلیل سری زمانی روش دیگری است که شامل بررسی منظم منابع در طول زمان است و برای تشخیص الگو و پیش بینی آینده استفاده می شود. مارتینز و همکاران این ایده را مطرح کرده است که سامانه، الگوهای بارکاری دوره ای یا فصلی که تکرار می شود را یاد بگیرد و در پیش بینی ها بکار گیرد. [7]. آناستاسیوس و همکاران با استفاده از چندین روش یادگیری تقویتی، به مدیریت خودکار بارهای کاری پویا با تضمین کیفیت خدمات (QoS) پرداختند و در نهایت روشی را انتخاب کردند که در استفاده از

۴- مراحل روش پیشنهادی

۴-۱- مؤلفه تعیین کننده گره میان توابع سرویس

محیط بدون سرویس دهنده محیطی است که متکی به توابع است و در مباحث نرم افزاری توابع همیشه پر استفاده ترین عضو برنامه نویسی است. وظیفه این مؤلفه، تعیین گره مهم میان توابع در حال اجرا در محیط بدون سرویس دهنده است. هنگامی که کاربر درخواستی را برای اجرا بر یکی از توابع مدنظر ارسال می کند، روابط میان توابع در محیط های بدون سرویس دهنده مورد بررسی قرار می گیرد و ساختمان داده گرافی متناظر با روابط توابع در نظر گرفته می شود. می خواهیم در گراف وابستگی ها تابعی که در معرض استفاده بیشتر است را پیدا کنیم و هرگاه منابع این تابع را بر حسب نیاز تأمین کنیم می توانیم اجرای پیوسته و دائم و بدون هدر رفت منابع داشته باشیم.

۴-۱-۱- درجه مرکزیت Degree Centrality

درجه مرکزیت در گراف ها یک معیار است که میزان مهم بودن یک گره در گراف را اندازه گیری می کند. این معیار میزان تعداد اتصال هایی که به یک گره متصل هستند را اندازه گیری می کند و گره هایی که درجه مرکزیت بالاتری دارند، به عنوان گلوگاه در گراف تلقی می شوند.

$$D_i = \frac{1}{n} \sum_{j=1}^n X_{ij} \quad (1)$$

در (۱) X_{ij} نشان دهنده میزان اتصال هر گره است.

۴-۱-۲- بردار مرکزیت ویژه

بردار مرکزیت ویژه محاسبه می کند که چقدر یک گره توسط گره های دیگر مهم می شود و گره هایی که ارتباطات بسیاری با گره مورد نظر دارند، می توانند به عنوان مهم ترین گره ها در شبکه شناخته شوند و با ماتریس مجاورت گراف و قوانین ضرب ماتریسی و دترمینان ماتریس محاسبه می شود.

$$C_v = \left(\frac{1}{\lambda} \right) \sum (A(u, v) * C(u)) \quad (2)$$

$$A * C_v = \lambda * C_v \quad (3)$$

در (۲) C_v بردار مرکزیت ویژه v را نشان می دهد. λ بزرگ ترین مقدار ویژه ماتریس مجاورت گراف، $A(u, v)$ عنصر ماتریس مجاورت میان دو گره و $C(u)$ مقدار بردار مرکزیت ویژه گره u است. در (۳) A ماتریس مجاورت گراف و λ مقدار ویژه است.

۴-۲- مؤلفه مدل سازی منابع مصرفی

ایده اصلی این مؤلفه این است که با داشتن میزان حافظه و پردازنده می توان مقیاس پذیری را بر اساس قوانین مبتنی بر آستانه تعیین کرد که آستانه ها بر اساس شناخت صحیح از سامانه تعیین می شوند و خبرگان سامانه می توانند در صورت نیاز آستانه ها را اضافه یا کم کنند. در این مؤلفه، ابتدا برای هر پارامتر تعیین کننده در مسئله، مقادیر آستانه تعریف می شوند و سپس بر اساس مقایسه ی بین پارامترها، تصمیم گیری مربوط به مقیاس پذیری انجام می شود.

۴-۳- مؤلفه تصمیم گیرنده مقیاس پذیری

هدف و تمرکز اصلی این پژوهش رسیدن به مقیاس پذیری است، برای مؤلفه مقیاس پذیری دو روش را در نظر می گیریم روش با نظارت و روش بدون نظارت. در روش بدون نظارت مدل یادگیری ماشین از روی شباهت داده ها به یک مقیاس می رسد و همه ی نمونه ها را بر اساس شباهت گروه بندی می کند اما در روش یادگیری با نظارت ما یک محرک یادگیرنده داریم که در مؤلفه مدل سازی منابع مصرفی بر اساس آستانه، آن محرک را ایجاد کردیم.

۴-۳-۱- روش یادگیری بدون نظارت

در روش بدون نظارت، روش های تجزیه و تحلیل مؤلفه های اصلی، خوشه بندی، نقشه های خود سازمان دهی در نظر گرفته شده اند. تجزیه و تحلیل مؤلفه های اصلی یک روش کاهش ابعاد بدون نظارت است که در تجزیه و تحلیل داده ها و استخراج ویژگی های مهم کمک می کند. خوشه بندی^۲ یک الگوریتم یادگیری ماشینی است که در آن داده ها بر اساس شباهت ها در گروه های متمایز قرار می گیرند. نقشه های خود سازمان دهی^۳ یک روش کلاسیک برای خوشه بندی داده ها و تجزیه و تحلیل الگوها استفاده می شود که مرکز خوشه به عنوان نماینده آن گروه عمل می کند.

۴-۳-۲- روش یادگیری با نظارت

در روش با نظارت، مدل های یادگیری ماشین شبکه پرسپترون چند لایه، شبکه پیچشی، شبکه حافظه بلند مدت کوتاه مدت در نظر گرفته شده اند و دلیل استفاده از این سه شبکه، فراوانی بیشتر و عملکرد بهتری است که در مسئله مقیاس پذیری دارند. شبکه عصبی پرسپترون چند لایه^۵ نوعی شبکه عصبی است که از چندین لایه از به هم پیوسته تشکیل شده است. حافظه طولانی کوتاه مدت^۶ یک نوع از شبکه های عصبی بازگشتی^۷ است که برای پردازش داده های دنباله ای و حفظ اطلاعات بلند مدت استفاده می شود. شبکه های عصبی پیچشی^۸ برای داده های ترکیبی استفاده می شوند. در روش یادگیری با نظارت تمامی پارامترهای مؤلفه تصمیم گیرنده مقیاس پذیری به عنوان یک معیار تعیین کننده محسوب شدند و با توجه به میزان تخطی از آستانه، مقدار مؤلفه مقیاس پذیری تعیین شد و شامل سه حالت انجام مقیاس پذیری، عدم انجام مقیاس پذیری و کاهش مقیاس پذیری است.

۴-۴- مؤلفه اجماع نظر میان پیش بینی ها

اگر چندین شبکه عصبی مختلف با ساختارها، پارامترها و بهینه سازیهای مختلف را استفاده کنیم، می توانیم با روش های اجماع نظر به بهترین خروجی برسیم. میانگین گیری مدل بیزی یک روش اجماع است که پیش بینی های چند مدل را ترکیب می کند و همزمان عدم قطعیت مرتبط با هر مدل را در نظر می گیرد. شکل (۲) به کد این اجماع را نشان می دهد.

۵- ارزیابی روش پیشنهادی

برای پیاده سازی و ارزیابی روش ارائه شده در این گزارش، از زبان پایتون و برای پیاده سازی مدل های یادگیری ماشین، از چارچوب کراس بر روی بستر گوگل کولب استفاده کردیم. داده های این پژوهش واقعی هستند و همگی از منابع مختلف مانند کگل^۹ و گیت هاب^{۱۰} جمع آوری گردیده اند و متناسب با نیاز پژوهش

برای آزمایش مؤلفه مدل‌سازی منابع مصرفی، مقدار تخطی از میزان مشخص مصرف پردازنده در چهار هسته پردازنده، حافظه، زمان اجرای تابع و همچنین مقدار بی‌استفاده ماندن پردازنده، حافظه را ملاک قرار داده‌ایم. برای تصمیم‌گیری در مورد مقیاس‌پذیری دو راه پیش‌روی داشتیم که توابع را برحسب شباهت در یک گروه‌بندی قرار دهیم (رویکردی بدون نظارت) و روش دقیق‌تر، روش یادگیری با نظارت است که برای این روش نیازمند داده‌های آموزشی که در آن مقیاس‌پذیری دخیل باشد وجود داشت، به همین منظور دو روش را پیاده‌سازی کردیم و بعد از آن پیش‌بینی مقیاس‌پذیری و منابع را برای گام زمانی آینده توسط روش با نظارت انجام دادیم و در بین خروجی‌های پیش‌بینی‌شده با روش میانگین‌گیری بیزی اجماع نظر گرفتیم.

۶- تحلیل نتایج

در این بخش نتایج روش پیشنهادی همراه با نمودارهای آن ارائه و برای بررسی کارایی و دقت مؤلفه پیش‌بینی‌کننده منابع، روش پیشنهادی با روش دیگری در همین حوزه مقایسه می‌شود. روش‌های بدون نظارت در مسئله پژوهش ما که مسئله‌ای حساس است نامناسب است چراکه گروه‌بندی به‌درستی انجام نشده‌است و خطای زیادی دارد. در پژوهش ما هدف بهبود پیش‌بینی مقیاس‌پذیری است، بنابراین روش‌های با نظارت بسیار بهتر از روش‌های بدون نظارت عمل می‌کنند. نتایج به تفکیک روش‌ها در شکل (۵) تا (۷) قابل مشاهده است.

برای ارزیابی مدل‌های یادگیری با نظارت، ۸۰ درصد مجموعه داده به‌عنوان داده تمرین و ۱۰ درصد را به‌عنوان داده آزمون و ۱۰ درصد را نیز برای ارزیابی استفاده کرده‌ایم و سپس مقادیر تابع‌های هزینه مانند میانگین خطای مربعات و میانگین خطای مطلق و دقت را روی مجموعه داده آزمون به دست آوردیم. شبکه حافظه طولانی کوتاه‌مدت توانسته است نسبت به سایر شبکه‌های عصبی در این مسئله طبقه‌بندی که برحسب سه برچسب مقیاس‌پذیری انجام می‌شد بهتر عمل کند. همچنین برای اطمینان از خروجی شبکه‌ها از ارزیابی K-Fold Cross Validation استفاده کردیم و داده‌ها را به‌صورت تصادفی در ۵ مرحله آزمایش کردیم و خروجی هر مرحله نیز در جدول (۱) مشخص است.

شکل (۸) مقایسه‌ای از نتایج دقت و خطای پیش‌بینی یادگیری مقیاس‌پذیری برای مدل‌های با نظارت را نشان می‌دهد.

جدول (۲) مقایسه‌ای از نتایج پیش‌بینی منابع در گام کنونی و آینده را برای مدل‌های با نظارت نشان می‌دهد.

۶-۱- بررسی با پژوهش مرتبط

در نهایت نیز پژوهش خود را با روشی که در پژوهش [۱۰] ارائه شده بود مقایسه کردیم. این مقایسه از این جهت انجام شد که مدل‌هایی که استفاده شده بود (پیش‌بینی و حافظه کوتاه مدت طولانی) در هر دو پژوهش یکسان بوده است. پژوهش مشابه هر دو مدل را ادغام کرده و از مدل ترکیبی پیش‌بینی-حافظه کوتاه‌مدت طولانی استفاده کرده و معیارهای خطا را برای مصرف CPU گزارش کرده است. در این مقاله ما سه مدل هوش مصنوعی را که به طور خاص پیش‌بینی، حافظه کوتاه مدت طولانی و شبکه پرسپترون چندلایه هستند را به صورت جداگانه آزمایش کرده و سپس برای مقایسه عادلانه فقط تأثیر دو مدل

با یکدیگر ادغام شدند. برای مؤلفه تعیین‌کننده گره میان توابع سرویس با استفاده از داده‌ها، وابستگی‌های بین توابع را به گراف تبدیل کردیم و هرچقدر که گره‌ای بار محاسباتی سنگینی را به دوش بکشد منابع بیشتری را اشغال می‌کند بنابراین با شناخت این تابع می‌توانیم آن را به بهترین شکل مدیریت کنیم و بیشتر از هر گره دیگری منابع آن را ارتقا دهیم. روش اولی که برای در نظر گرفتن مهم‌ترین گره پیاده‌سازی شد روش درجه مرکزیت بود و روش دوم مرکزیت بردار ویژه است تا به نتیجه‌ای که در روش قبل گرفتیم مطمئن‌تر شویم. با توجه به شکل (۳) و شکل (۴) در هر دو روش توانستیم به مهم‌ترین گره در گراف برسیم که می‌توانیم با اطمینان بالایی نتیجه بگیریم که مهم‌ترین گره در گراف است.

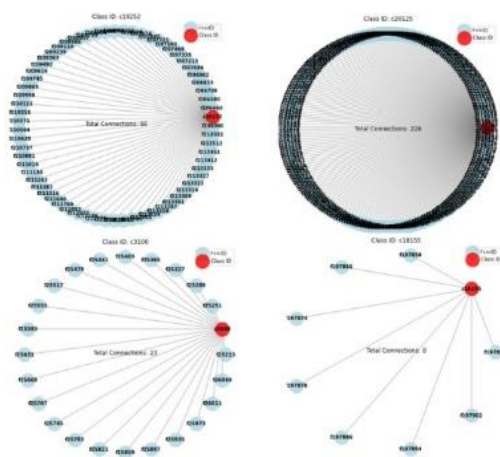
Algorithm 1 Consensus of predictions

INPUTS:
 $\hat{n}_{mlp_{t+1}}$ ← Prediction of mlp neural network $t+1$
 $\hat{n}_{cnn_{t+1}}$ ← Prediction of cnn neural network $t+1$
 $\hat{n}_{lstm_{t+1}}$ ← Prediction of lstm neural network $t+1$

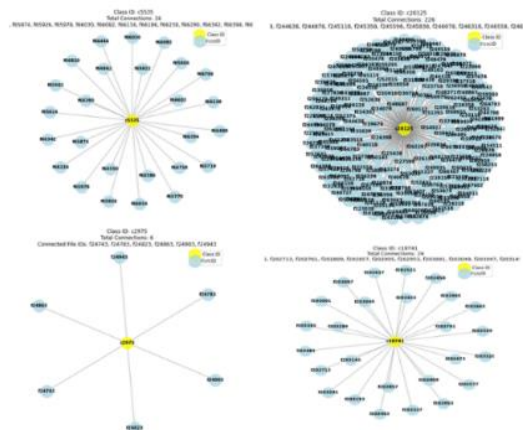
OUTPUTS:
 b_{avg}

Begin
 $\mu = 0$
 $\sigma = 1$
 $bayes_avg_preds = []$
for i **in** $range(len(test_data))$:
 $likelihood_prod = (1 / (\sigma * \sqrt{2 * \pi})) * \exp(-1/2 * ((\hat{n}_{mlp_{t+1}}[i] - \mu) / \sigma)^2) * (1 / (\sigma * \sqrt{2 * \pi})) * \exp(-1/2 * ((\hat{n}_{cnn_{t+1}}[i] - \mu) / \sigma)^2) * (1 / (\sigma * \sqrt{2 * \pi})) * \exp(-1/2 * ((\hat{n}_{lstm_{t+1}}[i] - \mu) / \sigma)^2)$
 $weights = \exp(-1/2 * ((\hat{n}_{mlp_{t+1}}[i] - \mu) / \sigma)^2) * \exp(-1/2 * ((\hat{n}_{cnn_{t+1}}[i] - \mu) / \sigma)^2) * \exp(-1/2 * ((\hat{n}_{lstm_{t+1}}[i] - \mu) / \sigma)^2)$
 $total_weight = sum(weights)$
 $weights = [weight / total_weight \text{ for } weight \text{ in } weights]$
 $posterior = sum(weight * likelihood_prod \text{ for } weight \text{ in } weights)$
 $bayes_avg_preds.append(posterior)$
 $b_{avg} = max(bayes_avg_preds)$
Return b_{avg}
End

شکل ۲: شبکه کد اجماع نظر پیش‌بینی‌ها



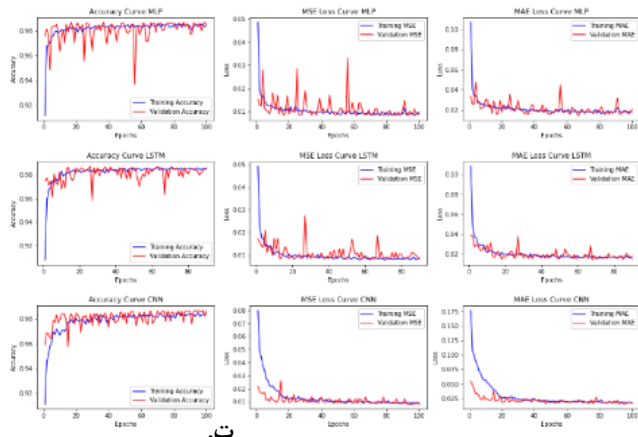
شکل ۳: گراف وابستگی‌ها با روش درجه مرکزیت.



شکل ۴: گراف وابستگی‌ها با روش مرکزیت بردار ویژه

۷- نتیجه گیری

در این پژوهش به بررسی ارائه یک رویکرد مقیاس پذیری خودکار در محیط بدون سرویس دهنده پرداخته شد و در ادامه ارزیابی جامعی بر روی الگوریتم پیشنهادی پیاده سازی شده با دو رویکرد با نظارت و بدون نظارت انجام شد.



شکل ۸: میانگین خطای مطلق و مربعات و دقت مراحل تمرین و آزمون مامی شبکه ها در روش بانظارت

جدول (۱): نتایج دقت و خطای پیش بینی یادگیری مقیاس پذیری

MAE	MSE	F1-Score	Recall	Precision	Accuracy	
0.02	0.04	0.98	0.98	0.98	0.984	MLP-Fold 1
0.01	0.02	0.99	0.99	0.98	0.99	MLP-Fold 2
0.02	0.03	0.98	0.99	0.98	0.9875	MLP-Fold 3
0.01	0.02	0.99	0.99	0.98	0.989	MLP-Fold 4
0.02	0.02	0.98	0.99	0.98	0.987	MLP-Fold 5
0.01	0.01	0.99	0.99	0.99	0.992	LSTM-Fold 1
0.01	0.02	0.99	0.99	0.98	0.9895	LSTM-Fold 2
0.01	0.01	0.99	0.99	0.98	0.9905	LSTM-Fold 3
0.01	0.01	0.99	0.99	0.98	0.9905	LSTM-Fold 4
0.01	0.01	0.99	0.99	0.98	0.9905	LSTM-Fold 5
0.02	0.02	0.98	0.99	0.98	0.9875	CNN-Fold 1
0.02	0.03	0.98	0.99	0.98	0.9865	CNN-Fold 2
0.05	0.09	0.97	0.97	0.97	0.9715	CNN-Fold 3
0.01	0.02	0.99	0.99	0.98	0.99	CNN-Fold 4
0.01	0.02	0.99	0.99	0.98	0.9895	CNN-Fold 5

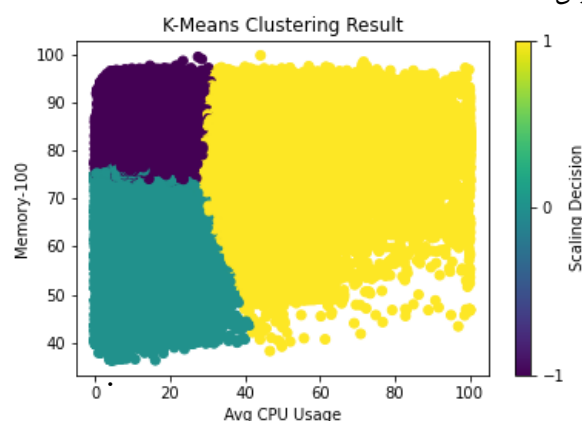
جدول (۲): مقادیر به دست آمده از معیارهای ارزیابی مدل پیشنهادی

مقادیر به دست آمده از معیارهای ارزیابی مدل پیشنهادی

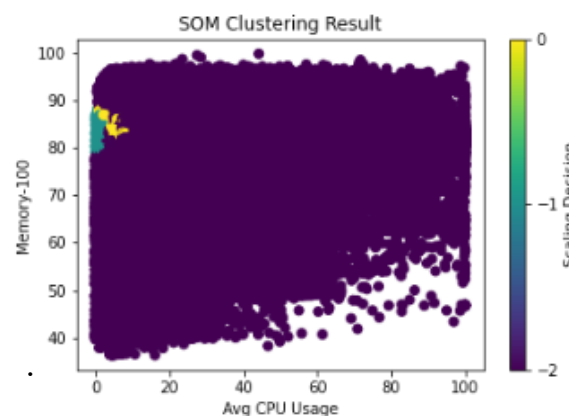
Actual Value for t	MLP Prediction for t	LSTM Prediction for t	CNN Prediction for t	Bayesian Avg Prediction for t
10.3520393	10.236974	10.364485	9.263559	9.934422
5.2734375	5.160522	5.405163	4.7756	5.103557
29.3650742	29.177689	29.725983	29.035881	29.304556
Actual Value for t+1	MLP Prediction for t+1	LSTM Prediction for t+1	CNN Prediction for t+1	Bayesian Avg Prediction for t+1
64.19699	63.02586	64.61939	46.664474	59.63367
64.19621	62.944954	64.57348	46.92247	59.659576
64.20816	63.03141	64.62794	47.328598	59.811844

پیشگی و حافظه کوتاه مدت طولانی را در نظر گرفتیم و نشان دادیم که با در نظر گرفتن شرایط مشابه مدل پیشنهادی ما بهتر از ادغام شبکه ها عمل میکند. این امر نشان دهنده کارا بودن ترکیب مکانیزم های اجماع نظر در مدل های یادگیری ماشین است. در جدول (۳) خلاصه ای از گزارش ها را مشاهده می کنید و منظور از CI بازه اطمینان است.

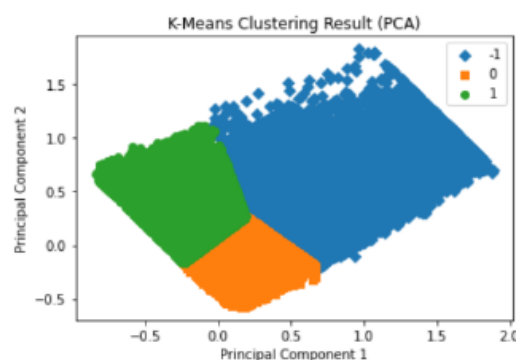
مقایسه نشان می دهد مدل پیشنهادی با مقدار RMSE برابر ۰.۶۹ برای آموزش شبکه پیشگی و ۰.۱۹ برای شبکه حافظه طولانی کوتاه مدت خطای پایین تری نسبت به مقاله مورد بحث با مقدار ۷.۱۴ دارد و مدل پیشنهادی با نرخ پایین تری از MAE و RMSE و MSE همراه است و بازه اطمینان نیز گزارش شده است.



شکل ۵: مقیاس پذیری با روش خوشه بندی



شکل ۶: مقیاس پذیری با روش نقشه های خود سازمان دهی



شکل ۷: مقیاس پذیری با روش تجزیه و تحلیل مؤلفه اصلی

- 7 RNN
- 8 CNN
- 9 Kaggle Code Metrics Dataset [SoftwareProjectStructure](#)
- 10 Github [Hardware Performance Archive](#)

جدول ۳: مقادیر به دست آمده از معیارهای ارزیابی مدل پیشنهادی

روش پیشنهادی این پژوهش					روش پیشنهادی پژوهش مرتبط					
	Model	MSE	MAE	RMSE		Model	MSE	MAE	RMSE	
Train	MLP	2.40	1.25	1.55		CNN-LSTM	26.90	2.64	5.18	
		0.46	0.58	0.68						
	CNN	0.07	0.15	0.26						
	LSTM									
Test	MLP	2.48	1.27	1.57		CNN-LSTM	51.09	3.72	7.14	
		0.47	0.59	0.69						
	CNN	0.03	0.15	0.19						
	LSTM									
CI	MLP	2.11-2.89	1.22-1.33	1.45-1.70		CNN-LSTM	188.11-208.49	8.56-9.15	13.71-14.47	
		0.44-0.52	0.57-0.61	0.66-0.72						
	CNN	0.03-0.04	0.14-0.16	0.17-0.22						
	LSTM									

نتایج نشان می‌دهد که الگوریتم پیشنهادی با استفاده از ترکیب مدل‌های شبکه عصبی پرسپترون و حافظه طولانی کوتاه مدت و پیش‌بینی می‌تواند خطاهای پیش‌بینی را به حداقل برساند و دقت بالاتری را ارائه دهد. مشاهدات به وضوح نشان می‌دهد که الگوریتم پیشنهادی با پایین‌ترین مقادیر MAE و RMSE و MSE، عملکرد بهتری نسبت به روش مقایسه شده دارد.

مراجع

- [1] Jonas, Eric, Schleier-Smith, Johann, Tsai, Chia-Che, "Cloud Programming Simplified: A Berkeley View on Serverless Computing," arXiv preprint arXiv: 1902.03383, Feb 9, 2019.
- [2] J. V. Bibal Benifa, D. Dejeu, "RLPAS: Reinforcement Learning-based Proactive Auto-Scaler for Resource Provisioning in Cloud Environment," Journal of Mobile Networks Application, vol. 24, no. 4, pp. 1348-1363, 2019.
- [3] Almeida, V., Arlitt, M., Rolia, J., "Analyzing a web-based system's performance measures at multiple time scales," Journal of ACM Sigmetrics Performance Evaluation Review, vol. 30, no. 2, pp. 3-9, Sep. 2002.
- [4] Golshani, E., Ashtiani, M., "Proactive auto-scaling for cloud environments using temporal convolutional neural networks," Journal of Parallel and Distributed Computing, vol. 154, no.4, pp. 119-141, 2021.
- [5] Al-Roomi, M., Al-Ebrahim S., Buqrais, S., Ahmad, I., "Cloud Computing Pricing Models: A Survey," Journal of Grid and Distributed Computing, vol. 6, no. 5, pp. 93-106, 2013.
- [6] Llorido-Bostrán, T., Miguel-Alonso, J., Lozano, J. A., "Comparison of Auto-scaling Techniques for Cloud Environments," Journal of Parallelism, Nov 2013, Paris, France, p.p 56-64.
- [7] Martínez, R. G., Li, Z., Lopes, A., Rodrigues, L., "Augure: Proactive reconfiguration of cloud applications using heterogeneous resources," in Proceedings of the 2017 IEEE 16th International Symposium on Network Computing and Applications (NCA), pp. 1-8, Cambridge, MA, USA, November, 2017.
- [8] Anastasios Z., Eleni F., Nikos F., Symeon P., "Reinforcement learning-assisted autoscaling mechanisms for serverless computing platforms", Simulation Modelling Practice and Theory, Vol. 116, 2022.
- [9] Siddharth A. and Rodriguez M.A. and Buyya R., "A Deep Recurrent-Reinforcement Learning Method for Intelligent AutoScaling of Serverless Functions", 2023.
- [10] Ouham, S., Hadi, Y., Ullah, A., "An efficient forecasting approach for resource utilization in a cloud data center using CNN-LSTM model," Journal of Neural Computing and Applications, vol. 33, no. 16, pp. 10043-10055, 2021.

زیر نویس

- ¹ model-free Proximal Policy Optimisation
- ² Principal Component Analysis
- ³ Clustering
- ⁴ Self-Organizing Maps
- ⁵ MLP
- ⁶ LSTM (Long Short-Term Memory)