



بهبود شبکه عصبی تابع پایه شعاعی مبتنی بر خوشه‌بندی‌های K-Means++ و DBSCAN برای مجموعه داده‌های بزرگ

فرشته حاج قاضی استرآبادی^۱، رضا قائمی^{۲*}، یعقوب آزاد^۳

*نویسنده مسئول، دریافت: ۱۴۰۱/۰۵/۱۸، بازنگری: ۱۴۰۲/۰۸/۱۵، پذیرش: ۱۴۰۳/۰۲/۱۶

^۱ دانشجوی دکتری مهندسی کامپیوتر، واحد نیشابور، دانشگاه آزاد اسلامی، نیشابور، ایران

^۲ استادیار گروه مهندسی کامپیوتر، واحد قوچان، دانشگاه آزاد اسلامی، قوچان، ایران

^۳ دانشجوی دکتری مهندسی کامپیوتر، واحد نیشابور، دانشگاه آزاد اسلامی، نیشابور، ایران

چکیده

شبکه‌های عصبی تابع پایه شعاعی در مسائل تقریبی مفید هستند، اما برای آموزش در داده‌های بزرگ، نیاز به زمان طولانی دارند و در صورتی که داده‌ها نامعتبر یا پرت باشند، آموزش آن با خطا همراه است. خوشه‌بندی یکی از فن‌های داده‌کاوی است که در شبکه عصبی تابع پایه شعاعی به کار می‌رود. استفاده از شبکه عصبی تابع پایه شعاعی در بسیاری از حوزه‌ها مانند پردازش تصویر، طبقه‌بندی متن، بیومتریک و ریزآرایه به دلیل اندازه بزرگ مجموعه داده‌ها، نیاز به زمان و حافظه زیادی برای پردازش دارد. با این حال، خوشه‌بندی می‌تواند این مشکل را بهبود بخشد و مراکز توابع گوسی را به درستی در بین داده‌ها تعیین کند. در این مقاله، روش خوشه‌بندی جدیدی با استفاده از الگوریتم‌های K-Means++ و DBSCAN ارائه شده که می‌تواند سرعت و دقت شبکه عصبی تابع پایه شعاعی را بهبود دهد. در روش پیشنهادی، ابتدا داده‌ها توسط الگوریتم K-Means++ گروه‌بندی شده و سپس، الگوریتم DBSCAN به‌طور جداگانه برای هر گروه اعمال می‌شود و داده‌ها خوشه‌بندی می‌شوند. ایده این مقاله، در بهبود سرعت اجرای الگوریتم DBSCAN از طریق کاهش محاسبات با تقسیم فضای کاری به تعدادی ناحیه جداگانه و استفاده از آن در شبکه عصبی تابع پایه شعاعی است. نتایج مدل پیشنهادی در مقایسه با شبکه عصبی تابع پایه شعاعی استاندارد و مدل بهبودیافته RBF-KMC بر روی مجموعه داده‌های Image segmentation، Pendigit و Letters، Shuttle control بر اساس زمان آموزش و دقت نتایج بهتری داشته است.

کلمات کلیدی: خوشه‌بندی، داده‌های حجیم، سرعت، دقت، شبکه عصبی تابع پایه شعاعی، K-Means++، DBSCAN.

۱- مقدمه

بالای پردازش داده‌ها می‌تواند بسیار زمان‌بر شود. دسته‌ای از روش‌های رایج برای تجزیه و تحلیل داده‌ها به‌عنوان داده‌کاوی شناخته می‌شود که به معنای شناسایی مفید، قابل اعتماد، الگوهای داده ساده و قابل درک که داده‌های خام را به داده یا اطلاعات مفید تبدیل می‌کند [۳]. خوشه‌بندی یکی از فن‌های داده‌کاوی است که در شبکه عصبی تابع پایه شعاعی نیز به کار گرفته می‌شود. خوشه‌بندی داده‌ها حوزه مهمی از یادگیری بدون نظارت در نظر گرفته می‌شود که در آن داده‌ها را می‌توان بر اساس شباهت‌هایشان از دیدگاهی آگاهانه در کل مجموعه داده به گروه‌های مختلف تقسیم نمود [۴]. این خوشه‌بندی به شبکه عصبی تابع پایه شعاعی این امکان را می‌دهد تا نورون‌های خود را در مرکز داده‌ها قرار داده و تخمین بهتری از الگو

عصر داده‌های بزرگ منجر به توسعه و کاربرد فناوری‌ها و روش‌هایی باهدف استفاده از مقادیر زیادی داده برای حمایت از تصمیم‌گیری و فعالیت‌های کشف دانش شده است [۱]. در بین فن‌های داده‌کاوی، شبکه‌های عصبی از جمله شبکه عصبی تابع پایه شعاعی^۱ (RBF) از روش‌های پرکاربرد در مواجهه با مجموعه داده‌های مختلف از جمله داده‌های بزرگ است. حجم زیادی از داده‌ها باعث شده است تا محققان و صنایع در راه حل‌های محاسباتی برای تجزیه و تحلیل داده‌های بزرگ تجدیدنظر کنند. به‌عنوان نمونه، تأکید زیادی بر طراحی الگوریتم‌های جدید که در محاسبات کارآمدتر هستند [۲]. با این حال، عملیات تجزیه و تحلیل و بازیابی به دلیل هزینه‌های محاسباتی

شعاعی باعث افزایش دقت در انتخاب مراکز و همچنین، سرعت آموزش تا حد زیادی شده است.

در ادامه مقاله، بخش دوم ادبیات تحقیق را مرور و بررسی می‌کند. در بخش سوم، الگوریتم پیشنهادی تشریح می‌شود. بخش چهارم، نتایج آزمایش‌های الگوریتم پیشنهادی را تحلیل و ارزیابی می‌کند. در نهایت، بخش پنجم نتیجه‌گیری و پیشنهادهایی را برای تحقیقات آینده ارائه می‌دهد.

۲- کارهای مرتبط

امروزه سازمان‌ها به دنبال ثبت و تجزیه و تحلیل داده‌ها از طریق روش‌های کارآمد هستند تا بتوانند تصمیم‌گیری مبتنی بر داده را تحقق بخشند. کلان داده‌ها را می‌توان به داده‌های ساختاری و منظم برای پردازش آسان‌تر تبدیل نمود. در این مقاله، به داده‌های ساختاریافته و بزرگ پرداخته می‌شود. اصطلاح داده‌های بزرگ به داده‌هایی اطلاق می‌شود که برای استخراج دانش نیاز به پردازش و تجزیه و تحلیل روش‌ها و معماری‌های جدید دارند [۱۴]. با توجه به اندازه بزرگ این داده‌ها، تجزیه و تحلیل کارآمد آن‌ها از طریق فن‌های مرسوم موجود بسیار دشوار است. با توجه به مطالعات بررسی شده در تحقیق [۱۵]، منظور از داده‌های بزرگ مجموعه داده‌هایی است که نمی‌توانند توسط ابزارهای نرم‌افزاری معمولی روی یک ماشین در زمان قابل قبولی مدیریت و پردازش شوند.

در سال‌های اخیر، روش خوشه‌بندی مختلفی برای کلان داده‌ها معرفی شده است. روش‌های فوق برای اجرای الگوریتم‌های خوشه‌بندی و کار با داده‌های بزرگ از طریق بهبود مقیاس‌پذیری و سرعت آن‌ها استفاده می‌شوند. اخیراً میر جلیلی [۱۶] نشان داد که ترکیبی از الگوریتم‌های تکاملی مانند بهینه‌سازی ازدحام ذرات با شبکه عصبی تابع پایه شعاعی عملکرد خوبی در مسائل طبقه‌بندی و مسائل تقریب نشان می‌دهد. استفاده از الگوریتم تکاملی به‌عنوان ابزاری برای انتخاب مراکز دقیق‌تر نیز در بسیاری از متون اخیر گزارش شده است [۱۷-۱۸]. لذا الگوریتم تکاملی برای آموزش شبکه عصبی تابع پایه شعاعی روش مناسبی است، اگر سرعت آموزش شبکه و هزینه محاسبات نگرانی اصلی نباشد. در کنار افزایش محبوبیت شبکه‌های عصبی تابع پایه شعاعی، کاربرد این شبکه عصبی در زمینه‌هایی مانند طبقه‌بندی [۱۹]، تشخیص الگو [۲۰] و پیش‌بینی [۲۱] نیز افزایش داشته و استفاده و قابلیت اطمینان شبکه عصبی تابع پایه شعاعی در بسیاری از زمینه‌ها را نشان می‌دهد.

باین حال، دقت شبکه عصبی تابع پایه شعاعی به‌طور عمده به مراکز اولیه انتخاب‌شده از مجموعه داده، قبل از شروع آموزش شبکه بستگی دارد [۲۲-۲۳]. علاوه بر این، اندازه مجموعه داده‌های آموزشی و داده‌های نامعتبر موجود در مجموعه داده‌ها نیز نقش مهمی در تعیین سرعت و دقت آموزش شبکه ایفا می‌کنند [۲۴-۲۵]. همچنین، گزارش شده است که الگوریتم یادگیری برای آموزش شبکه‌ها ممکن است با افزایش مجموعه داده‌ها بدتر عمل کند [۲۶]. از این‌رو برای حل این مشکلات، محققان استفاده از الگوریتم‌های خوشه‌بندی را در انتخاب مراکز [۲۷-۳۰] برای شبکه عصبی تابع پایه شعاعی برای دستیابی به دقت و سرعت آموزش بهتر و جلوگیری از گنجاندن مجموعه داده‌های نامعتبر احتمالی در آموزش شبکه‌ها پیشنهاد کردند.

پرکاربردترین الگوریتم خوشه‌بندی در انتخاب مراکز، الگوریتم K-Means است، زیرا سریع‌ترین و کم پیچیده‌ترین الگوریتم است و در مقایسه با سایر الگوریتم‌های خوشه‌بندی موجود، دقت خوبی برای انتخاب مراکز دارد. باین حال، الگوریتم‌های خوشه‌بندی مرسوم موجود، الگوریتم‌های بی‌عیب و نقصی نیستند که دقت بالایی را در انتخاب مراکز تضمین نمایند. به‌این ترتیب، محققان راه‌هایی برای بهبود ضعف الگوریتم‌های خوشه‌بندی موجود را بررسی می‌کنند.

در ادبیات اخیر، در تحقیق [۳۱]، یک الگوریتم K-Means بهبودیافته با رویکردهای چند خوشه‌ای تکراری برای انتخاب مراکز و درعین حال به حداقل

داده‌ها انجام دهد. خوشه‌بندی در طیف وسیعی از زمینه‌ها مانند شناسایی مجدد وسیله نقلیه [۵]، حذف نویز تصویر [۶]، پردازش سری زمانی [۷] و تجزیه و تحلیل شبکه‌های اجتماعی مبتنی بر وب [۸] استفاده می‌شود. امروزه، روش‌های فوقی مرسوم در برخورد با داده‌های بزرگ به‌اندازه کافی کارآمد نیستند و به دلیل پیچیدگی محاسباتی بالا قابل اجرا نمی‌باشند [۹].

بنابراین یکی از بهترین ویژگی‌های شبکه‌های عصبی، توانایی آن در تعمیم و تقریب داده‌های نمونه بدون نیاز به تعیین معادله و ضرایب است، به‌ویژه زمانی که یک مدل ناشناخته، یک رابطه پیچیده ناشناخته را توصیف می‌کند و داده‌های آموزشی فراوانی دارد. به دلیل توانایی آن‌ها در تعمیم، شبکه‌های عصبی تابع پایه شعاعی معمولاً برای این منظور انتخاب می‌شوند [۱۰-۱۱]؛ اما در کلان داده‌ها، شبکه عصبی تابع پایه شعاعی در بسیاری از حوزه‌ها مانند پردازش تصویر، طبقه‌بندی متن، بیومتریک، ریزآرایه و غیره دارای اندازه مجموعه داده‌هایی به‌قدری بزرگ هستند که سیستم بلادرنگ نیاز به زمان و حافظه طولانی برای پردازش آن‌ها دارد. شبکه عصبی باید در کلیه نمونه‌های آموزشی موجود جستجو کند که به حافظه زیادی نیاز دارد و محاسبه فاصله تا مرکز را انجام دهد که در طول آموزش شبکه عصبی برای اهداف تقریبی، فرآیند کندی است. علاوه بر این، به دلیل ذخیره تمامی فواصل نمونه برای مجموعه داده‌های آموزشی در شبکه عصبی تابع پایه شعاعی، ممکن است فواصل نویز نیز ذخیره شوند که این می‌تواند باعث کاهش دقت تقریبی گردد.

الگوریتم DBSCAN [12]، یکی از محبوب‌ترین و کاربردی‌ترین الگوریتم‌های خوشه‌بندی است. باین حال، این الگوریتم در مواجهه با کلان داده‌ها به‌درستی عمل نمی‌کند. یکی از مهم‌ترین معایب آن، سرعت اجرای پایین است. با توجه به سرعت پیشرفت در دنیای امروز، زمان یکی از مهم‌ترین پارامترها است و باید مورد توجه ویژه قرار گیرد. برای غلبه بر مشکل پردازش کلان داده، روش‌های مختلفی از جمله بهینه‌سازی الگوریتم‌های مرسوم از نظر محاسبات و حتی مصرف حافظه، کاهش داده‌ها از طریق کاهش ابعاد یا نمونه‌های داده، پردازش جریان داده‌ها و استفاده از ابزارهای سنگین پیشنهاد شده است.

ایده این مقاله، در بهبود سرعت اجرای الگوریتم DBSCAN از طریق ایده کاهش محاسبات با تقسیم فضای کاری به تعدادی ناحیه جداگانه و استفاده از این خوشه‌بندی در شبکه عصبی تابع پایه شعاعی است. برای تقسیم داده‌ها به گروه‌های کوچک‌تر، از الگوریتم خوشه‌بندی معروف دیگری بنام ++K-Means [13] استفاده می‌گردد. این الگوریتم سرعت همگرایی مناسبی را نشان می‌دهد و می‌تواند تقسیم داده‌های خوبی را فراهم کند. لازم به ذکر است که اجرای K به معنی چند تکرار، این امکان را می‌دهد که بدون بار محاسباتی زیاد به نتایج دلخواه دست‌یابیم، چراکه ++K-Means نه برای به دست آوردن خوشه‌بندی نهایی داده‌ها، بلکه صرفاً برای تقسیم آن‌ها استفاده می‌شود.

به‌طور کلی، الگوریتم پیشنهادی به این صورت عمل می‌کند که ابتدا، داده‌ها توسط الگوریتم ++K-Means با تعداد کمی از تکرارها گروه‌بندی می‌شوند. سپس، الگوریتم DBSCAN به‌طور جداگانه برای هر قسمت اعمال می‌شود و داده‌ها خوشه‌بندی می‌شوند. در نهایت، بخش‌های مختلف یکپارچه‌شده‌اند. برای به دست آوردن نقاط اصلی در الگوریتم DBSCAN، اگر قطعات MinPts در شعاع ϵ وجود داشته باشد، باید برای هر قطعه داده تحلیل و آنالیز شود. از این‌رو، فاصله بین هر قطعه از داده‌ها و سایر داده‌های موجود در مجموعه داده باید محاسبه گردد؛ بنابراین، مقدار زیادی از محاسبات برای یک مجموعه داده بزرگ انجام می‌شود. باین حال در روش پیشنهادی، فضای داده‌ای که باید در نظر گرفته شود، کوچک‌تر است، زیرا DBSCAN به‌طور جداگانه در هر گروه اعمال می‌شود و در نتیجه، تعداد نقاط کمتری وجود دارد که فواصل آن‌ها باید محاسبه شود. این امر، زمان اجرا را به میزان قابل توجهی کاهش می‌دهد. استفاده از این خوشه‌بندی در شبکه عصبی تابع پایه

در تحقیق [۴۰]، یک الگوریتم خوشه‌بندی مبتنی بر چگالی سریع بر اساس الگوریتم DBSCAN و با استفاده از ایده کاهش محاسبات پیشنهاد شده است. الگوریتم ابتدا کلیه داده‌ها را بر اساس مختصات بعد مرتب می‌کند و سپس، با کاهش زمان اجرای هر پرس‌وجو منطقه یا فرکانس پرس‌وجوهای ناحیه، سرعت کل اجرای DBSCAN را افزایش می‌دهد. در تحقیق دیگری، یک الگوریتم خطی DBSCAN بر اساس هشینگ حساس به محل ارائه شده است تا سرعت یافتن نزدیک‌ترین گره‌های همسایه را افزایش دهد [۴۱]. مزیت استفاده از هشینگ حساس به محل این است که پیچیدگی زمانی و مقیاس داده را کاهش می‌دهد. الگوریتم دو بخش را در نظر گرفته است. شاخص هشینگ حساس به محل در بخش اول ساخته شده، در حالی که خوشه‌بندی در بخش دوم، توسط الگوریتم DBSCAN بر اساس شاخص بازیابی هشینگ حساس به محل^۴ انجام می‌شود که یک ضعف الگوریتم مشکل در تعیین مقادیر برای پارامترهای ورودی است.

یک الگوریتم خوشه‌بندی جدید در تحقیق [۴۲] ارائه شده است. این الگوریتم برای بهبود دقت خوشه‌بندی با ادغام K-Means++، با یک نقشه خودسازمان‌دهی^۵ پیشنهاد شده است. در واقع، K-Means++ وزن اولیه و مرکز اولیه را تعیین کرده و سپس، از نقشه خودسازمان‌دهی برای یافتن راه‌حل مناسب برای خوشه‌بندی نهایی داده‌ها استفاده شده است. الگوریتم نقشه خودسازمان‌دهی با هدف پیشنهاد خوشه‌بندی بهینه و با سرعت بالا، از نظر صرفه‌جویی در زمان آموزش و میزان خطا به درستی عمل نموده است. با این وجود، مشکل اصلی K-Means++ این است که ذاتاً ترتیبی است. در نتیجه، کارایی آن برای کلان داده محدود است.

دراهمکار پیشنهادی در این مقاله، نیازی به تکرار K-Means++ تا انتها نیست، زیرا تکرارها محدود هستند. در مقاله [۴۳]، روشی از نوع الگوریتم‌های مبتنی بر چگالی افزایشی برای مجموعه داده‌های بزرگ پیشنهاد شده است که هدف این روش، بهبود الگوریتم افزایشی DBSCAN با محدود کردن فضای جستجو به چند قسمت به جای استفاده از کل مجموعه داده بود. این امر فرآیند خوشه‌بندی را تسریع می‌کند. الگوریتم خوشه‌بندی سلسله مراتبی K-Means برای یافتن مراکز اولیه با کیفیت بالا برای مجموعه داده‌های بزرگ در تحقیق [۴۴] پیشنهاد شد. در این الگوریتم، داده‌ها در یک ساختار سلسله مراتبی کاهش می‌یابند؛ بنابراین، در گروه فن‌های کاهش داده‌ها قرار می‌گیرند. مرکزها از طریق ساختار K-Means سلسله مراتبی متشکل از چندین لایه به دست می‌آیند که اولین لایه شامل مجموعه داده اصلی است. لایه بعدی شامل مجموعه داده‌های کوچک‌تری است که از طریق مدل-هایی مشابه مجموعه داده اصلی ایجاد شده‌اند. به عبارت دیگر، مراکز هر خوشه نشان‌دهنده داده‌های آن خوشه است و به سطح بعدی منتقل می‌شود. از نظر محاسباتی فشرده، الگوریتم‌های سلسله مراتبی به‌طور معمول به‌عنوان خط مبنا یا برای بهینه‌سازی الگوریتم اصلی خوشه‌بندی استفاده می‌شوند و نمی‌توانند نتایج بسیار خوبی برای محاسبات در مقیاس بزرگ به دست آورند. روش دیگری در بهبود K-Means++ با رویکردی تکراری برای افزایش کیفیت و سرعت خوشه‌بندی حاصل از K-Means++ و K-Means از طریق حذف یک خوشه و تقسیم خوشه دیگر در تحقیق [۴۵] ارائه شده است. خوشه‌ها بر اساس مقایسه توابع سود و هزینه حذف و اضافه می‌شوند و مشکلی که در این الگوریتم وجود دارد، این است که محاسبات بیش‌ازحد دارد که باعث کاهش سرعت اجرا می‌شود.

در تحقیق دیگری، الگوریتمی برای افزایش سرعت DBSCAN در هنگام اجرای داده‌های بزرگ طراحی شده است [۴۶]. الگوریتم، ابتدا فضای داده را به یک ساختار شبکه تقسیم می‌کند و سپس، یک مقدار چگالی را به هر سلول شبکه اختصاص می‌دهد. سپس فضاهای شبکه را بررسی می‌کند تا فضاهای همسایه را پیدا کرده و متراکم‌ترین همسایه را در هر فضا مشخص کند. در مرحله بعدی، یعنی در مرحله تشکیل خوشه، الگوریتم زنجیره‌ای از متراکم‌ترین همسایگان را تشکیل می‌دهد و این کار را ادامه می‌دهد تا خوشه‌ها توسعه یابند. برای این منظور، بخش‌های شبکه

رساندن تابع هدف پیشنهاد شده است. این روش خوشه‌بندی سطح اول را انجام داده و مراکز نسل اول را ذخیره می‌کند و سپس، خوشه‌بندی سطح دوم را برای کوتاه کردن فواصل مراکز اجرا می‌کند. این روش دقت خوبی دارد، اما هزینه محاسباتی سنگینی است. در تحقیق [۳۲]، از یک K-Means سه سطحی برای تعیین مراکز استفاده شده است. این الگوریتم می‌تواند بر نویز غلبه کند و داده‌های نامعتبر را فیلتر کند، اما در بار محاسباتی پرهزینه است. کار در بهبود الگوریتم K-Means با استفاده از مفاهیم چگالی ادامه یافته تا جایی که در تحقیق [۲۷]، استفاده از چگالی را برای بهبود الگوریتم K-Means و به دست آوردن دقت مناسب معرفی می‌گردد. این روش همچنین نشان می‌دهد که الگوریتم K-Means نسبت به داده‌های نویز حساسیت کمتری دارد. این روش بر محاسبه تراکم هر شیء یا نقاط داده در مجموعه داده متمرکز است و تنها مراکز را انتخاب می‌کند که نقاط تراکم بالاتری در نزدیکی آن دارند. تحقیقات مشابهی از سال‌های اولیه در ادبیات [۳۳]، با استفاده از مفاهیم چگالی که بر محاسبه چگالی نقاط داده در اطراف مراکز تمرکز دارد، یافت شده است. در تحقیق [۳۴]، یک روش انتخاب مراکز اولیه برای الگوریتم K-Means با استفاده از شاخص‌های اتکینسون^۶ معرفی شده است. به جای انتخاب تصادفی مراکز اولیه برای ادامه الگوریتم K-Means برای تنظیم بهتر مراکز، نویسندگان برای استفاده از این رویکردها، خطای احتمالی داده‌های نامعتبر را از رویکردهای انتخاب تصادفی کاهش می‌دهند. شاخص‌های اتکینسون از مرز و نابرابری برای تعیین محدوده انتخاب مراکز برای الگوریتم K-Means استفاده می‌کنند. در تحقیق [۳۵]، از مفاهیم مشابه در الگوریتم K-Means پیشنهادی استفاده می‌کند. این الگوریتم از کران‌های بالا و پایین بین یک نقطه و مراکز اولیه تصادفی برای انتخاب مراکز مناسب استفاده می‌کند. این الگوریتم یک فرآیند پیچیده را برای انجام محاسبات اعمال کرده که برای محققان به‌راحتی قابل‌درک نیست. علاوه بر این، محققانی نیز از توابع هسته برای محاسبه فاصله در الگوریتم K-Means به‌منظور به دست آوردن نتایج طبقه‌بندی بهتر استفاده کردند [۳۶]. با استفاده از تابع هسته، مجموعه داده‌های پیچیده به ابعاد بالاتر نگاشت می‌شوند که مجموعه داده را به‌راحتی قابل‌تفکیک و طبقه‌بندی مجموعه داده‌ها تبدیل می‌کند.

در تحقیق [۳۷]، روشی برای تعیین وزن به خوشه‌ها نسبت به واریانس معرفی شده است. این کار بینش جدیدی در مورد استفاده از رویکردهای آماری در الگوریتم یادگیری ماشین ایجاد می‌کند، با این حال در محاسبات بسیار پیچیده و پرهزینه است. در تحقیق [۳۸]، استفاده از فاصله ماهالانوبیس^۷ برای جایگزینی فاصله اقلیدسی در الگوریتم K-Means پیشنهاد شده است که البته بهبود قابل‌توجهی در دقت ایجاد نکرده است. به‌طور کلی، خوشه‌بندی را می‌توان به دو گروه اصلی [۳۹] طبقه‌بندی نمود، شامل روش‌های مبتنی بر یک ماشین واحد و روش‌های مبتنی بر ماشین‌های متعدد. الگوریتم‌های خوشه‌بندی بر اساس یک ماشین واحد روی یک ماشین اجرا می‌شوند و فقط می‌توانند از منابع آن ماشین استفاده کنند. الگوریتم‌های خوشه‌بندی مبتنی بر ماشین‌های متعدد بر روی دو یا چند ماشین اجرا می‌شوند و می‌توانند از منابع چندین ماشین به‌طور هم‌زمان استفاده کنند.

در روش‌های خوشه‌بندی بر اساس یک ماشین واحد، روش‌های پیشنهادی در تحقیق‌های قبلی برای حل مشکل پردازش داده‌های بزرگ اشاره می‌شود. روش‌های فوق را می‌توان به سه گروه تقسیم نمود. اولین فن، کاهش اندازه داده است که می‌تواند از طریق کاهش نمونه‌ها، ابعاد یا هر دو انجام شود. فن دوم، پردازش جریان داده است که پردازش حجم زیادی از داده‌ها را در یک ماشین ممکن می‌کند، اما از مشکلات خاصی رنج می‌برد. فن سوم، الگوریتم‌های بهینه‌سازی پردازش داده است که در آن، داده‌ها را می‌توان با بار محاسباتی کمتری پردازش کرد. ایده این مقاله، در گروه سوم روش‌ها بر اساس ماشین نهفته است.

ادغام خوشه‌های به‌دست‌آمده از فاز اول استفاده می‌شود. زمان اجرای این روش طولانی‌تر از روش‌های جدید ارائه‌شده در ادبیات تحقیق مطرح شده است. در واقع، تقریباً برابر با K-Means است که یک نقطه‌ضعف محسوب می‌شود.

در تحقیق [۵۳]، طرحی برای اجرای موازی DBSCAN در سال ۲۰۱۸ ارائه شد که سلول‌های کوچک داده را به‌طور تصادفی به موازی تقسیم می‌کند. در این تحقیق، یک الگوریتم موازی جدید DBSCAN معروف به RP-DBSCAN ارائه شد که در آن، از طرح تقسیم تصادفی داده‌ها و یک فرهنگ لغت سلولی دوسطحی استفاده شد که خلاصه‌ای فشرده از کل مجموعه داده بود. الگوریتم به‌طور هم‌زمان برای هر بخش از داده‌ها خوشه‌های محلی را پیدا کرده و سپس، آن‌ها را به‌منظور به دست آوردن خوشه‌بندی نهایی ادغام می‌کند. در این روش، هزینه‌های تقسیم فضایی به دلیل استفاده از توزیع داده‌های تصادفی کاهش می‌یابد. باین‌حال، زمان اجرا را طولانی می‌کند.

در آخرین روش‌های مطرح‌شده در سال ۲۰۲۰، بر اساس شباهت همسایگی، یک الگوریتم DBSCAN اصلاح‌شده ارائه‌شده است که با استفاده از روش درخت پوشش، همسایه‌های موازی هر نقطه را بازیابی می‌کند و بسیاری از محاسبات غیرضروری را برای اندازه‌گیری فاصله از طریق اصل نابرابری مثلثی بررسی می‌کند. مزیت قابل‌توجه این الگوریتم سرعت‌بالای آن است. در واقع، مناطق با تراکم بالا را به‌عنوان یک دسته طبقه‌بندی می‌کند و خوشه‌بندی فضایی موردنظر را به‌سرعت پیدا می‌کند. از آنجایی‌که این روش موازی است، اگر از فن‌هایی مانند کاهش نقشه استفاده شود، الگوریتم پیشنهادی بسیار سریع‌تر اجرا می‌شود [۵۴]. جدول ۱ خلاصه‌ای از روش‌های متفاوت مرور شده که مبتنی بر خوشه‌بندی DBSCAN است را نشان می‌دهد.

جدول ۱- روش‌های خوشه‌بندی مبتنی بر DBSCAN

مرجع	الگوریتم پیشنهادی	نقاط قوت	نقاط ضعف
[۴۰]	چگالی سریع در DBSCAN	دقت مناسب	سرعت پایین و زمان اجرای بالا
[۴۱]	هشینگ حساس به محل در DBSCAN	دقت و سرعت مناسب	دقت نامناسب
[۴۲]	K-Means++ با نقشه خودسازمان‌ده	سرعت مناسب	زمان اجرای بالا
[۴۳]	محدود کردن فضای جستجو در DBSCAN	زمان اجرای مناسب	دقت پایین
[۴۴]	K-Means سلسله مراتبی	سرعت مناسب	دقت پایین
[۴۵]	حذف و تقسیم خوشه در K-Means++	دقت و سرعت مناسب	دقت نامناسب
[۴۶]	تقسیم شبکه‌ای در DBSCAN	زمان اجرای مناسب	دقت پایین
[۴۷]	تقسیم ابرمکعب در DBSCAN	زمان اجرای مناسب	دقت نامناسب
[۴۸]	کاهش نمونه‌ها در DBSCAN	زمان اجرای مناسب	دقت پایین
[۴۹]	KNN-BLOCK-DBSCAN	زمان اجرای مناسب	دقت پایین
[۵۰]	Block K-Means++ DBSCAN	دقت و سرعت مناسب	دقت نامناسب
[۵۱]	K-Means++ MapReduce	سرعت مناسب	دقت پایین
[53]	اجرای موازی RP-DBSCAN	سرعت مناسب	دقت پایین
[54]	شباهت همسایگی در DBSCAN	زمان اجرای مناسب	دقت پایین

۳- الگوریتم پیشنهادی

به‌طور کلی، ساختار روش پیشنهادی از سه مرحله تشکیل شده است، شامل گروه‌بندی داده‌ها با استفاده از الگوریتم K-Means++، اجرای الگوریتم DBSCAN بر روی داده‌های هر گروه در مرحله قبل و ادغام خوشه‌های تولیدشده در گروه‌های مختلف. الگوریتم پیشنهادی K-Means++ نام‌گذاری شده است، چراکه داده‌ها در ابتدا با استفاده از الگوریتم K-Means++ به K گروه تقسیم می‌شوند و سپس،

را در نظر گرفته و متراکم‌ترین همسایه‌هایشان را به خوشه اضافه می‌کند. بخش بزرگی از زمان اجرای الگوریتم، صرف یافتن متراکم‌ترین همسایگان می‌شود که به‌عنوان یک ضعف الگوریتم در نظر گرفته می‌شود. در ادامه و در تحقیق [۴۷]، الگوریتم HCA-DBSCAN ارائه شد که در آن، ابرمکعب‌هایی را مانند شبکه‌هایی با قطر ϵ در فضای داده در نظر می‌گیرد و داده‌های هر ابرمکعب در یک خوشه قرار دارند. در هر ابرمکعب، چند نقطه که نزدیک به مرزها هستند، به‌عنوان عامل انتخاب می‌شوند. فقط عوامل هر جفت شبکه به‌جای کلیه نقاط موردبررسی و مقایسه قرار می‌گیرند. اگر فاصله بین دو عامل از دو شبکه کوچک‌تر از ϵ باشد، دو شبکه ادغام می‌شوند و داده‌های آن‌ها در یک خوشه قرار می‌گیرند. با استفاده از این عوامل، تعداد مقایسه‌هایی که برای محاسبه فواصل بین داده‌ها باید انجام شود، کاهش می‌یابد و درنهایت، الگوریتم با سرعت بالاتری اجرا می‌شود. به دلیل چندین تلاش برای جستجوی شاخص‌ها برای به دست آوردن شاخص‌های مجاور، این الگوریتم نیاز به‌صرف زمان زیادی دارد، موردی که به‌عنوان یک ضعف محسوب می‌شود.

ایده‌ای که در سال ۲۰۱۹ برای پیاده‌سازی DBSCAN بر روی داده‌های بزرگ ارائه شد، یک روش نمونه‌گیری برای کاهش اندازه داده‌ها است [۴۸]. این مقاله دو روش را برای نمونه‌گیری بهتر در DBSCAN پیشنهاد می‌کند. در نتیجه، زمان اجرای DBSCAN به دلیل کاهش تعداد نمونه‌ها کاهش می‌یابد. یکی از روش‌ها، نسخه اصلاح‌شده Rough DBSCAN است، درحالی‌که روش دیگر یک فن اکتشافی است. در تحقیق دیگری، یک روش تقریبی بنام KNN-BLOCK DBSCAN برای خوشه‌بندی کلان داده در سال ۲۰۱۹ پیشنهاد شد [۴۹]. برای این منظور، آن‌ها از یک الگوریتم نزدیک‌ترین همسایه برای شناسایی بلوک‌های اصلی (CB)، بلوک‌های غیرهسته‌ای (NCB) و بلوک‌های نویز (NOB) که به‌ترتیب شامل نقاط اصلی، نقاط مرزی و نقاط نویز هستند، استفاده شد. سپس CBهای قابل‌دسترسی با چگالی ادغام شدند، درحالی‌که NCBها به یک خوشه مناسب اختصاص داده شدند. این محققین در مقاله [۵۰]، یک روش خوشه‌بندی مبتنی بر شبکه بنام K-Means++ DBSCAN BLOCK را در دو نسخه برای داده‌های بزرگ با ابعاد بالا در سال ۲۰۲۰ پیشنهاد کردند. با استفاده از روش کاهش محاسباتی، این ایده دو فن را پیشنهاد می‌کند. فن اول، برای شناسایی بلوک‌های هسته داخلی که کلیه نقاط آن نقاط اصلی هستند، استفاده می‌شود. فن دوم، یک الگوریتم تقریبی با سرعت بالا است که پیش‌بینی می‌کند که آیا دو بلوک هسته داخلی قابل‌دسترسی به چگالی هستند یا خیر. علاوه بر این، یک درخت پوشش برای تسریع محاسبات چگالی برای نقاط بازدیدنشده استفاده می‌شود. مزیت این روش این است که می‌تواند داده‌هایی با ابعاد بالا را نسبتاً سریع خوشه‌بندی کند.

روش‌های خوشه‌بندی بر اساس چند ماشین نیز دارای تحقیقات متفاوتی است. در تحقیق [۵۱]، برای بهبود خوشه‌بندی کلان داده از طریق الگوریتم K-Means تلاش‌هایی انجام شده و برای این منظور از چارچوب MapReduce استفاده شده است. این مقاله بر روی بهینه‌سازی انتخاب مراکز خوشه‌بندی اولیه برای افزایش دقت و سرعت متمرکز شده است. به این منظور، تعدادی زیرمجموعه از کلان داده‌ها از طریق نمونه‌گیری انتخاب شدند. زیرمجموعه‌ها به‌صورت موازی به‌منظور به دست آوردن مراکزی که می‌توانند برای خوشه‌بندی مجموعه داده اصلی استفاده شوند، پردازش شدند. ادغام پس‌ازانجام نمونه‌برداری‌های موردنیاز صورت گرفت و در این روش، از K-Means++ برای انتخاب مراکز اولیه استفاده شد.

در پژوهش [۵۲]، مجموعه داده به بخش‌های کوچک‌تر تقسیم شده و به چندین گروه در یک خوشه از ماشین‌ها توزیع شده است. برای این منظور Apache Hadoop به‌عنوان یک چارچوب مقیاس‌پذیر و قدرتمند استفاده شد. K-Means در دو فاز بر روی ماشین‌ها پیاده‌سازی می‌شود که در مرحله اول، تعداد زیادی خوشه‌های اولیه ایجاد می‌شوند. برای این منظور، از خوشه‌بندی احتمالی برای انتخاب مراکز اولیه بهتر و کاهش تعداد دوره‌های همگرایی استفاده می‌شود. در فاز دوم، یک آستانه برای

فاصله تا نقاط دیگر کمک نماید. فواصل یک نقطه تا کلیه نقاط دیگر در مجموعه داده باید اندازه‌گیری شود تا مشخص شود که آیا یک هسته است یا خیر. گروه‌بندی ارائه‌شده توسط K-Means++ باعث می‌شود که محاسبات سنگین به‌ویژه در مجموعه‌های داده بزرگ، به‌جای این‌که در کل مجموعه داده انجام شود، به‌گروه‌ها محدود شود.

به‌عنوان نمونه، اگر مجموعه داده‌ای متشکل از یک میلیون قطعه داده را از طریق K-Means++ به ده گروه تقسیم نماییم، تنها صدهزار محاسبه به‌طور متوسط به‌جای محاسبه یک میلیون فاصله زمانی که DBSCAN اجرا می‌شود، انجام می‌شود. اکنون باید مشخص شود که آیا خوشه‌های تولیدشده توسط DBSCAN در گروه‌های K-Means++ مجاور وجود دارد که باید ادغام شوند. از آنجایی‌که الگوریتم DBSCAN هر گروه را جداگانه در نظر می‌گیرد و خوشه‌بندی را تنها با در نظر گرفتن داده‌های همان گروه انجام می‌دهد، خوشه‌های یک گروه باید همراه با گروه‌های مجاور بررسی شوند؛ بنابراین، اگر هیچ مرزی بین گروه‌ها وجود نداشته باشد، DBSCAN می‌تواند داده‌ها را در گروه‌های مجاور در هنگام انجام خوشه‌بندی مشاهده نماید.

۳-۲- ادغام خوشه‌های مرزی انتخاب‌شده

در این مرحله باید بررسی شود که آیا خوشه‌های به‌دست‌آمده از الگوریتم DBSCAN در صورتی‌که مرزهای گروهی بین داده‌ها وجود نداشته باشد، یا این‌که خوشه‌های مرزی به یکی ادغام می‌شوند. برای کاهش محاسبات، در این مرحله نیز دو دور هرس برای مواردی که نیاز به بررسی دارند، پیشنهاد می‌شود. دور اول، شامل تجزیه و تحلیل فواصل بین گروه‌های K-Means++ و اعمال هرس برای گروه‌های فواصل طولانی (بیش از ϵ) است که ادغام خوشه‌های داخلی را غیرممکن می‌کند. دور دوم، شامل هرسی می‌شود که وقتی خوشه‌های داخلی دو گروه برای ادغام و حذف خوشه‌هایی که ادغام آن‌ها غیرممکن است، بررسی می‌شود. پس از انجام این دو دور هرس، گروه‌های انتخاب‌شده توسط الگوریتم DBSCAN ادغام می‌شوند و دوباره بر روی داده‌های این خوشه‌ها اجرا می‌شود.

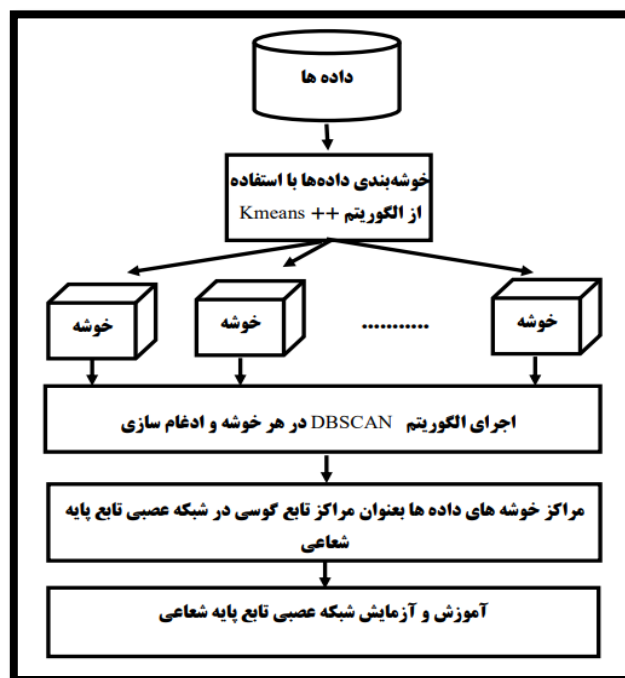
بنابراین، اگر فاصله بین دو مجموعه بیشتر از ϵ باشد، قطعاً هیچ نقطه‌ای از مجموعه دوم در همسایگی ϵ نقاط مجموعه اول قرار نخواهد گرفت و بالعکس (مقدار ϵ ثابت و برابر ۰.۰۱ در نظر گرفته شده است). بدیهی است که در صورت پیاده‌سازی DBSCAN روی داده‌های این دو مجموعه، ایجاد یک خوشه غیرممکن خواهد بود. سپس روش ادغام خوشه‌های مرزی تعیین‌شده در سه مرحله توضیح داده می‌شود. فواصل بین مجموعه‌ها در مرحله اول تعیین می‌شود و اگر فاصله آن‌ها از ϵ بیشتر باشد، مجموعه‌ها هرس می‌شوند. این اولین فرآیند هرس در این مرحله است. پس از آن، فواصل بین خوشه‌های موجود در مجموعه‌های باقی‌مانده از فرآیند هرس اول، دوباره محاسبه می‌شوند. اگر فاصله بین دو خوشه از ϵ بیشتر شود، فرآیند هرس دوم اتفاق می‌افتد. در نهایت، DBSCAN دوباره از داده‌های خوشه‌های تعیین‌شده در مرحله سوم اجرا می‌شود.

مرحله اول: بررسی فواصل بین گروه‌ها است و این هرس برای گروه‌های K-Means++ قبل از بررسی ادغام خوشه‌ها اعمال می‌شود. بر اساس ایده ارائه‌شده در این مقاله، امکان ادغام باید برای خوشه‌های واقع در گروه‌های مجاور بررسی شود. بدیهی است که ادغام خوشه‌های جفت گروه‌هایی که به‌جای همسایگی از یکدیگر دور هستند، غیرممکن است و این گروه‌ها را می‌توان از فرآیند محاسبات برای بررسی امکان ادغام خوشه‌های گروهی حذف نمود. لذا، گروه‌های به‌دست‌آمده از الگوریتم K-Means++ دوباره در نظر گرفته شده است و معادله (۱) برای آن‌ها محاسبه می‌شود.

$$dis_{K_{ij}} = dis_{KC_{ij}} - (dis_{r_i} + dis_{r_j}) \quad (1)$$

که در آن، $dis_{K_{ij}}$ نشان‌دهنده فاصله بین مراکز دو گروه در K-Means++ و dis_{r_i} و dis_{r_j} نشان‌دهنده فواصل بین مراکز دو گروه و دورترین نقاط روی آن‌ها

الگوریتم DBSCAN در هر گروه استفاده می‌شود. هر مرحله از الگوریتم به‌طور جداگانه در ادامه توضیح داده شده است. در شکل ۱، فلوچارت کلی از روش پیشنهادی نشان داده شده است.



شکل ۱- مراحل روش پیشنهادی

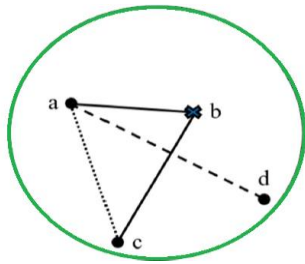
۳-۱- گروه‌بندی داده‌ها از طریق الگوریتم K-Means++

در این مرحله، الگوریتم K-Means++ روی کل مجموعه داده اجرا می‌شود تا داده‌ها به چند قسمت کوچک‌تر تقسیم شوند. از آنجایی‌که الگوریتم DBSCAN بر اساس محاسبه چگالی هر قطعه داده عمل می‌کند و هدف آن قراردادن داده‌های نزدیک به یکدیگر در یک خوشه است، مانند سایر روش‌های خوشه‌بندی، داده‌ها را نمی‌توان به‌صورت دلخواه یا تصادفی تقسیم کرد تا کارایی مناسب به دست آید. در روش پیشنهادی از K-Means++ که یک الگوریتم خوشه‌بندی است، استفاده می‌شود. نتایج نشان می‌دهد که اجرای K-Means++ در تعداد کمی از تکرارها نیز می‌تواند تقسیم مناسبی را بر روی داده‌ها ارائه دهد، یافته‌ای که تعادل قابل قبولی را نیز نشان می‌دهد. تنها تفاوت بین K-Means++ و K-Means کلاسیک این است که مراکز اولیه خوشه در K-Means++ هوشمندانه‌تر از K-Means انتخاب می‌شوند. از این رو می‌تواند تقسیم‌بندی مناسب‌تری را پیشنهاد کند.

همان‌طور که قبلاً ذکر شد، تنها هدف از به‌کارگیری K-Means++ تقسیم داده‌ها است، کاری که می‌تواند بدون بار محاسباتی زیاد بر روی سیستم انجام شود، حتی زمانی‌که الگوریتم در تعداد کمی از تکرارها اجرا شود؛ مانند سایر روش‌های خوشه‌بندی، الگوریتم K-Means++ داده‌ها را خوشه‌بندی می‌کند. با این حال در این مقاله، خوشه‌های K-Means++ همراه با عناوین گروه نام‌گذاری شده‌اند که با خوشه‌های به‌دست‌آمده از DBSCAN متفاوت باشند.

۳-۲- اجرای الگوریتم DBSCAN در هر گروه

در این مرحله، الگوریتم DBSCAN به‌طور جداگانه روی داده‌های هر یک از گروه‌های K-Means++ اجرا می‌شود. برای تعیین این‌که آیا یک نقطه یک هسته در DBSCAN است، باید فاصله آن تا کلیه نقاط دیگر مجموعه داده اندازه‌گیری شود. مقدار زیادی از محاسبات صرفاً برای به دست آوردن تعداد نقاط واقع در فاصله ϵ انجام می‌شود؛ بنابراین، تقسیم داده‌ها به گروه‌های کوچک‌تر و استفاده از DBSCAN به‌طور جداگانه در هر گروه می‌تواند به کاهش محاسبات موردنیاز برای اندازه‌گیری



شکل ۳: شماتیکی از نحوه محاسبه فاصله بین مرکز خوشه و دورترین نقطه (bc)

بر اساس فرآیند ذکر شده، کران بالای d را می توان با استفاده از DBSCAN با محاسباتی غیر قابل ملاحظه ای محاسبه کرد (با خط چین نشان داده شده است). برای محاسبه فاصله بین مرکز خوشه b و دورترین نقطه از آن که به صورت c نشان داده شده است، یک خط از a به c ترسیم کرده و مثلث abc تشکیل می شود. با توجه به مثلث فوق، قضیه نابرابری مثلث به صورت معادله (۲) نشان داده می شود [۵۵]:

$$bc < ab + ac \quad (2)$$

اگر مرکز خوشه b فرض کنیم، می توان فاصله بین نقطه a و مرکز b را محاسبه نمود. علاوه بر این، از آنجایی که d دورترین نقطه از a است، معادله (۳) صادق است [۵۵].

$$ac \leq ad \quad (3)$$

با توجه به مقدار نامعلوم ac و برای جلوگیری از محاسبات اضافی، معادله (۲) را می توان به صورت معادله (۴) و بر اساس معادله (۳) بازنویسی کرد [۵۵].

$$bc < ab + ad \quad (4)$$

از معادله (۴) مشخص است که مقدار bc قطعاً کمتر از مجموع ab و ad است. بنابراین، کران بالایی برای مقدار bc ، یعنی فاصله بین مرکز خوشه و دورترین نقطه داده در خوشه را می توان بر اساس معادله (۴) محاسبه کرد [۵۵]. با در نظر گرفتن دو خوشه از دو گروه، مانند شکل (۴)، از معادله (۵) برای به دست آوردن فاصله بین دو خوشه استفاده می شود [۵۵].

$$dis_{D_{ij}} = dis_{DC_{ij}} - (dis_{r_i} + dis_{r_j}) \quad (5)$$

در معادله (۵)، dis_{r_i} و dis_{r_j} نشان دهنده فواصل بین مراکز خوشه و دورترین نقاط از آن ها در همان خوشه ها، $dis_{D_{ij}}$ نشان دهنده فاصله بین دو خوشه DBSCAN و $dis_{DC_{ij}}$ نشان دهنده فاصله بین مراکز دو خوشه DBSCAN در دو گروه K-Means++ است. اگر مقدار $dis_{DC_{ij}}$ کمتر یا مساوی ϵ باشد، این دو خوشه را می توان ادغام نمود و در غیر این صورت هرس می شوند. فقط گروه هایی که در فواصل کمتر یا مساوی ϵ از یکدیگر قرار دارند مورد بررسی قرار می گیرند. این شرایط باعث می شود که بسیاری از محاسبات هرس شوند.

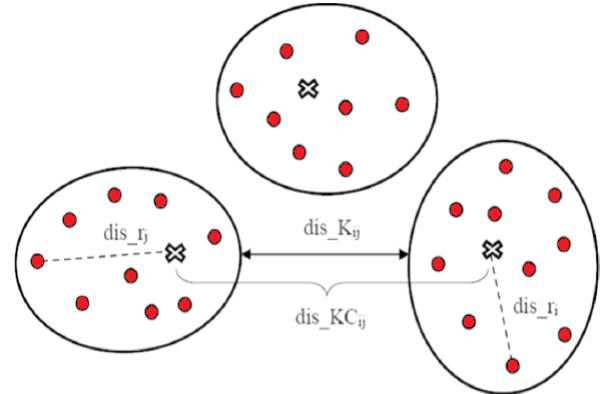
مرحله سوم: خوشه های مرزی با اجرای مجدد الگوریتم DBSCAN ادغام می شوند. پس از مشخص شدن امکان ادغام خوشه ها و انجام هرس پیشنهادی، داده های مربوط به آن هایی که می توانند ادغام شوند، با یکدیگر در نظر گرفته می شوند و الگوریتم DBSCAN بدون توجه به هر مرزی برای آن ها دوباره اعمال می شود. سپس، یا خوشه ها به همان ترتیب ادامه می دهند یا دو یا چند تا از آن ها در یک خوشه ادغام می شوند. بنابراین، کل فضای داده در نهایت توسط الگوریتم DBSCAN خوشه بندی می شود.

پیچیدگی زمانی روش پیشنهادی با در نظر گرفتن پیچیدگی هر سه مرحله آن محاسبه شده است. معادله (۶)، پیچیدگی زمانی رویکرد پیشنهادی را به صورت ترکیبی از پیچیدگی بیان می کند، با این تفاوت که محاسباتی که انجام می شوند بسیار پیچیده تر از محاسبات سنگین در DBSCAN استاندارد می باشند.

$$O((iter * k * n) + \left(\frac{n}{k}\right)^2 + \frac{k(k-1)}{2}m^2) \quad (6)$$

در این مراحل، از آنجایی که فواصل از k نقطه مرکزی تا n نقطه باید در الگوریتم K-Means++ به دست آید و محاسبات پارها تکرار می شوند، پیچیدگی زمانی

است. معادله فاصله بین دو گروه یعنی $dis_{K_{ij}}$ را با در نظر گرفتن بدترین حالت به دست می آورد. هنگامی که الگوریتم K-Means++ اعمال می شود، می توان این پارامترها را برای هر گروه بدون محاسبات قابل توجهی محاسبه و ذخیره نمود. فواصل به صورت شماتیک در شکل ۲ نشان داده شده اند. اگر مقدار $dis_{K_{ij}}$ بزرگ تر از ϵ باشد، قطعاً بعید است که خوشه های مرزی دو گروه ادغام شوند. بنابراین، نیازی به بررسی خوشه های داخلی کلیه گروه ها برای ادغام نیست.



شکل ۲: شماتیکی از نحوه محاسبه فاصله بین دو گروه [۵۵]

فقط گروه هایی که در فواصل کمتر یا مساوی ϵ از یکدیگر قرار دارند، مورد بررسی قرار می گیرند. این شرایط باعث می شود که بسیاری از محاسبات هرس شوند. به طور کلی، این هرس دومین موردی است که برای الگوریتم DBSCAN در الگوریتم پیشنهادی پس از تقسیم داده ها به گروه های کوچک تر اعمال می شود.

مرحله دوم: تجزیه و تحلیل احتمال ادغام خوشه های DBSCAN است. در این مرحله، امکان ادغام خوشه های داخلی برای گروه هایی که در مرحله قبل هرس نشده اند، تحلیل می شود. فواصل بین خوشه ها باید برای یافتن آن هایی که می توانند ادغام شوند، به دست آید. اگر فاصله کوچک تر یا مساوی ϵ باشد، ادغام می تواند اتفاق بیافتد، در غیر این صورت، احتمال صفر است و داده های مربوط به آن خوشه ها در مرحله ادغام در نظر گرفته نمی شوند. برای اعمال این هرس، باید یک مرکز برای هر خوشه DBSCAN در نظر گرفته شود که بر اساس داده های میانگین در خوشه محاسبه می شود. لازم به ذکر است که مرکز برای کلیه خوشه ها محاسبه نمی شود و در مقایسه با محاسباتی که بعداً هرس می شود، پیچیدگی محاسباتی بالایی را تحمیل نمی کند.

پارامتر دیگری که برای اعمال هرس مورد نیاز است، فاصله بین مرکز خوشه و دورترین نقطه داده در خوشه است. پارامتر را می توان با تغییر جزئی در الگوریتم DBSCAN و بدون بار محاسباتی قابل توجهی محاسبه نمود. همان طور که ذکر شد، DBSCAN از یک نقطه تصادفی شروع می شود و بررسی می کند که آیا آن نقطه یک هسته است یا خیر. اگر نقطه به عنوان یک هسته شناسایی شود، الگوریتم این روش را در نقاطی در فاصله ϵ از هسته تکرار می کند. به این معنا که ابتدا داده هایی که در فاصله حداکثر ϵ قرار دارند و سپس، کلیه داده هایی که در فاصله ϵ قرار دارند، بررسی می شوند. این روند تا زمانی ادامه می یابد که کل خوشه مشخص شود. با این حال، هنگام اجرای الگوریتم می توان فاصله داده هایی را که در آن محاسبه می شود، ذخیره نمود و بنابراین این می تواند حد بالایی برای فاصله بین نقطه شروع و دورترین فاصله از آن باشد.

شکل ۳ یک خوشه DBSCAN را نشان می دهد. فرض کنید که تشکیل خوشه به صورت تصادفی در نقطه a آغاز شده است. دورترین نقطه از a نام دارد که قطعاً یک نقطه مرزی نیز است و بنابراین، این می تواند حد بالایی برای فاصله بین نقطه شروع و دورترین نقطه از آن باشد.

عملکرد خطی هستند و خروجی‌های شبکه عصبی تابع پایه شعاعی برابر مجموع وزن دار شده خروجی‌های نرون‌های لایه پنهان متصل به لایه خروجی هستند. در واقع، خروجی y_k نرون k ام در لایه خروجی به صورت معادله (۸) محاسبه می‌شود:

$$y_k = \sum_{j=1}^m W_{kj} h_j \quad (8)$$

که در آن، W_{kj} وزن اتصال نرون j ام لایه پنهان و نرون k ام لایه خروجی و m تعداد نرون‌های لایه پنهان است. پارامترهای قابل تنظیم در فرآیند آموزش شبکه عصبی تابع پایه شعاعی متعدد هستند و کارایی شبکه در بعضی از موارد به شدت وابسته به این پارامترها و روش تطبیق آن‌ها است. این پارامترها مراکز توابع گوسی، میزان گستردگی آن‌ها و وزن‌های اتصالات طبقه خروجی می‌باشند. یک مسئله مهم نیز تعداد نرون‌های لایه پنهان است که باید به روش مناسبی آن را به دست آورد. نکته مهم این است که لایه‌های متفاوت یک شبکه RBF، وظایف متفاوتی انجام می‌دهند و منطقی است که روش‌های جداگانه‌ای برای تطبیق و آموزش هر یک در نظر گرفته شود.

بنابراین راه کار پیشنهادی این مقاله، متمرکز بر روی تعیین خوشه‌ها و مراکز خوشه با دقت و سرعت مناسب است تا شبکه عصبی تابع پایه شعاعی بتواند در مواجهه با داده‌های حجیم به خوبی عمل نماید.

۴- ارزیابی نتایج آزمایش

در این بخش، ابتدا تنظیمات و محیط سخت‌افزاری و شبیه‌سازی آزمایش و سپس، مجموعه داده‌های آزمون در آزمایش تشریح می‌شوند و در نهایت، نتایج حاصل شده در آزمایش ارزیابی می‌گردند.

۴-۱- تنظیمات و محیط سخت‌افزاری و شبیه‌سازی آزمایش

به منظور شبیه‌سازی آزمایش، محیط سخت‌افزار شامل رایانه‌ای با حافظه ۸ گیگابایت و پردازنده سه هسته‌ای اینتل و سیستم عامل ویندوز ۱۰ استفاده شده است. به علاوه، آزمایش‌ها در محیط نرم‌افزاری MATLAB-R2019b پیاده‌سازی شده است. با توجه به محاسبه زمان اجرا، کلیه روش‌ها بر روی این سیستم سخت‌افزاری و نرم‌افزاری پیاده‌سازی شده و نتایج آن ارزیابی شده است.

۴-۲- مجموعه داده‌های آزمون در آزمایش

در آزمایش‌ها، از چهار مجموعه داده آزمون استفاده شده است که شامل مجموعه داده Image segmentation دارای ۲۳۱۰ نمونه، Pendigit دارای ۱۰۹۹۲ نمونه، Letters دارای ۲۰۰۰۰ نمونه و Shuttle control دارای ۵۸۰۰۰ نمونه می‌باشند [۵۶]. این چهار مجموعه داده از پایگاه داده‌های UCI برداشته شده است [۵۶].

خصیصه‌های این چهار مجموعه داده آزمون در جدول ۲ نشان داده است.

جدول ۲: مشخصات مجموعه داده‌های آزمون

مجموعه داده	تعداد نمونه	تعداد کلاس	تعداد ویژگی‌ها	نوع ویژگی‌ها
Image segmentation	2310	7	19	عدد حقیقی
Pendigit	10992	44	16	عدد صحیح
Letters	20000	40	16	عدد صحیح
Shuttle control	58000	2	15	اسمی

۴-۳- معیارهای ارزیابی آزمایش

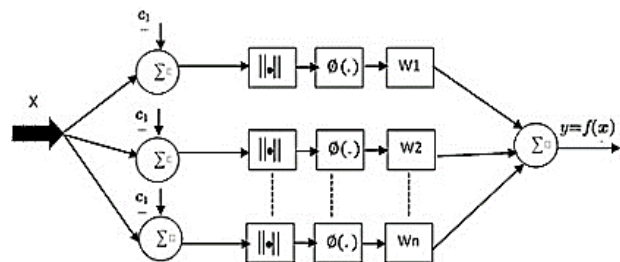
در کلیه آزمایش‌ها، از روش اعتبارسنجی K-fold با $K=10$ استفاده شده است و نتایج به صورت میانگین آمده است. برای ارزیابی سرعت، از ثانیه و برای دقت دسته‌بندی،

الگوریتم، $O((iter * k * n))$ است. تعداد تکرارها، تعداد مراکز اولیه و اندازه داده بر زمان اجرای الگوریتم تأثیر می‌گذارد. پیچیدگی زمانی الگوریتم DBSCAN نیز $O(n^2)$ است، زیرا باید فواصل بین نقاط در n جفت به دست آید. با توجه به مراحل ذکر شده و الگوریتم توسعه یافته برای روش پیشنهادی، پیچیدگی زمانی الگوریتم پیشنهادی، همان‌طور که در معادله (۶) بیان شده است، از سه بخش تشکیل شده است.

از آنجایی که الگوریتم K-Means++ برای اولین بار اجرا می‌شود و داده‌ها را گروه‌بندی می‌کند، اولین بخش پیچیدگی زمانی مربوط به K-Means++ است. بخش دوم مربوط به الگوریتم DBSCAN است که به طور جداگانه برای داده‌های موجود در k گروه اعمال می‌شود. بخش سوم پیچیدگی زمانی به مرحله نهایی در روش پیشنهادی می‌پردازد. در این مرحله از الگوریتم پیشنهادی، خوشه‌های موجود در گروه‌ها دوباره برای انتخاب خوشه‌های قابل ادغام جستجو می‌شوند. سپس، DBSCAN به داده‌های این گروه‌ها انتخاب شده، اعمال می‌شود. بنابراین، واضح است که پیچیدگی زمانی الگوریتم ارائه شده کمتر از DBSCAN است و زمانی که الگوریتم پیشنهادی روی مجموعه داده‌های مختلف اجرا می‌شود، به ویژه در داده‌های بزرگ، کاهش می‌یابد، چراکه ابتدا گروه‌بندی صورت پذیرفته و سپس، DBSCAN بر روی هر گروه انجام می‌شود.

۴-۳- مدل پیشنهادی

شبکه عصبی تابع پایه شعاعی به دلیل سادگی ساختار و همچنین روش کارآمد آموزش آن، مورد توجه بسیاری قرار گرفته است. شبکه عصبی RBF قابلیت تخمین همه منظوره را دارد، بنابراین شبکه عصبی RBF می‌تواند برای مسائل درون بایی استفاده شود. یک شبکه عصبی تابع پایه شعاعی گوسی، فرم غیر نرمالیزه شده تابع توزیع گوسی و بسیار غیرخطی است و ویژگی‌های مناسبی برای یادگیری افزایشی دارد. شبکه‌های عصبی گوسی برای یادگیری نگاشت‌های پیچیده تعریف شده‌اند، اما در شبکه عصبی تابع پایه شعاعی، انتخاب مراکز تابع گوسی اهمیت دارد و این کار با استفاده از الگوریتم K-Means انجام می‌شود. در مواجهه با داده‌های حجیم، الگوریتم خوشه‌بندی باید بتواند مراکز خوشه‌ها را با دقت و سرعت بالاتری به جای استفاده از الگوریتم K-Means به همراه داشته باشد. در مدل پیشنهادی بنام Means Scan++ در تعیین مراکز خوشه‌های گوسی در شبکه عصبی تابع پایه شعاعی استفاده می‌شود. یک شبکه RBF از دو لایه کاملاً متصل به یکدیگر بانام‌های لایه پنهان یا لایه RBF و لایه خروجی مطابق شکل ۴ تشکیل می‌شود.



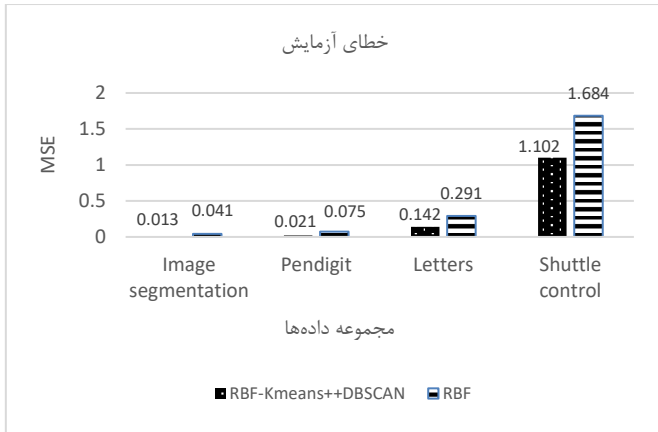
شکل ۴: معماری شبکه RBF [56]

مطابق شکل ۴، لایه ورودی شبکه عصبی RBF به طور مستقیم به لایه پنهان متصل می‌شود و خروجی نرون j ام لایه پنهان به صورت معادله (۷) به دست می‌آید:

$$h_j = \frac{\varphi(\|x - c_j\|)}{\sigma_j} \quad (7)$$

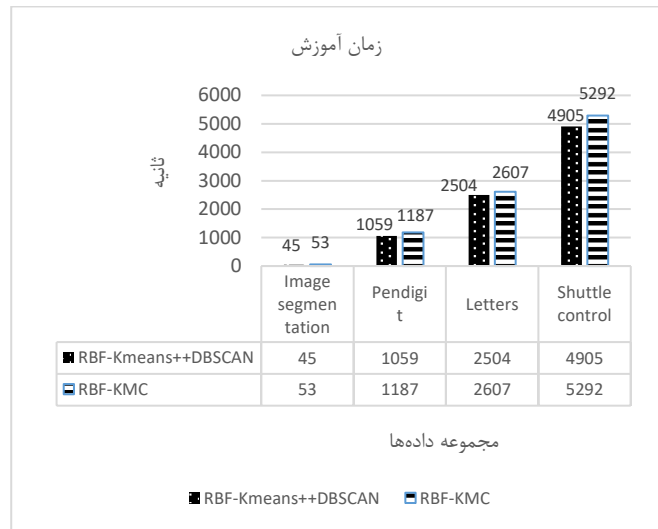
که در آن، h_j خروجی نرون j ام، φ تابع عملکرد غیرخطی RBF، X بردار ورودی، c_j مرکز نرون و σ_j میزان گستردگی مرکز نرون است. در مدل پیشنهادی، C_j توسط روش خوشه‌بندی پیشنهادی انجام می‌شود. عملکرد غیرخطی شبکه عصبی تابع پایه شعاعی به دلیل تابع عملکرد φ است. نرون‌ها در لایه خروجی، دارای تابع

Letters. برای مدل خوشه‌بندی پیشنهادی زمان لازم برای آموزش ۲۵۰۴ ثانیه بوده است، اما در شبکه عصبی استاندارد در حدود دو برابر و به میزان ۴۸۰۹ ثانیه است. میزان خطا در این مجموعه داده برای مدل خوشه‌بندی پیشنهادی در حدود ۰.۱۴۹٪ خطای کمتری از روش استاندارد است. در مجموعه داده Shuttle control، برای مدل خوشه‌بندی پیشنهادی زمان لازم برای آموزش ۴۹۰۵ ثانیه بوده است، اما در شبکه عصبی استاندارد ۷۰۵۴ ثانیه است. میزان خطا در این مجموعه داده برای مدل خوشه‌بندی پیشنهادی در حدود ۰.۵۸۲٪ خطای کمتری از مدل شبکه عصبی استاندارد بوده است.



شکل ۶: مقایسه خطای مدل پیشنهادی با روش RBF استاندارد

مقایسه مدل خوشه‌بندی پیشنهادی با مدل بهبودیافته شبکه عصبی تابع پایه شعاعی، با بهبود خوشه‌بندی الگوریتم K-Means سراسری تغییر یافته سریع [۵۶] بنام RBF-KMC انجام شده است. در شکل‌های ۷ و ۸ نشان می‌دهد که مدل خوشه‌بندی پیشنهادی در مجموعه داده‌های مشترک توانسته است که زمان آموزش کمتر و دقت بیشتری را در آزمایش داشته باشد.



شکل ۷: مقایسه زمان آموزش مدل پیشنهادی با الگوریتم RBF-KMC

در شکل ۷، میزان خطای آزمایش مدل خوشه‌بندی پیشنهادی با الگوریتم خوشه‌بندی K-Means سراسری تغییر یافته مقایسه شده است که نشان می‌دهد که خوشه‌بندی، میزان خطا در شبکه عصبی تابع پایه شعاعی را در هر دو روش کاهش داده است، زیرا مراکز تابع گوسی با دقت بهتری مشخص شده است، اما در مدل خوشه‌بندی پیشنهادی میزان خطای پایین‌تری به دست آمده است. همان‌طور که در شبیه‌سازی مدل پیشنهادی در شکل‌های ۷ و ۸ مشاهده می‌شود، می‌توان بیان نمود که در مجموعه داده Image segmentation، الگوریتم شبکه عصبی تابع پایه

از معیار خطای MSE استفاده شده است که در معادله (۹) آمده است و که در آن، پارامتر Y_i خروجی واقعی، \hat{Y}_i خروجی مدل است و n تعداد نمونه‌ها است.

$$MSE = \frac{1}{n} \sum_{i=1}^n (Y_i - \hat{Y}_i)^2 \quad (9)$$

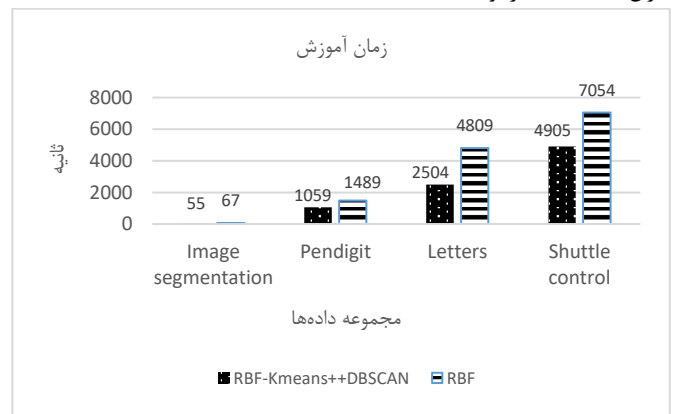
۴-۴- نتایج آزمایش

نتایج عددی برای بررسی عملکرد مدل خوشه‌بندی پیشنهادی با شبکه عصبی تابع پایه شعاعی استاندارد که با الگوریتم K-Means خوشه‌بندی می‌کند، مقایسه شده است. در جدول ۳، مقایسه زمان آموزش و میزان خطا در آزمایش مدل‌های شبکه عصبی بهبود یافته و مدل استاندارد آمده است.

جدول ۳: مقایسه نتایج مدل پیشنهادی

مدل	روش پیشنهادی		روش استاندارد	
	RBF K-Means++	DBScan	RBF	
مجموعه داده	خطا (MSE) در آزمایش	خطا (MSE) در آموزش	خطا (MSE) در آزمایش	خطا (MSE) در آموزش
	0.013	0.041	0.013	0.041
	45	67	45	67
	0.021	0.075	0.021	0.075
Image segmentation	2504	4809	2504	4809
	4905	7054	4905	7054
Pendigit	1.102	1.684	1.102	1.684
	4905	7054	4905	7054
Letters	4905	7054	4905	7054
	4905	7054	4905	7054
Shuttle control	4905	7054	4905	7054
	4905	7054	4905	7054

در شکل ۵، مدت زمان آموزش برای مجموعه داده‌های آزمایش شده، آمده است که نشان می‌دهد با مدل خوشه‌بندی پیشنهادی، زمان آموزش کاهش یافته است و هرچه مجموعه داده حجیم‌تر باشد، کاهش زمان آموزش محسوس‌تر است. مدل خوشه‌بندی پیشنهادی توانسته است که با خوشه‌بندی سریع و دقیق به شبکه عصبی تابع پایه شعاعی کمک کند تا آموزشی در زمان کمتر و دقت دسته‌بندی بالاتری در آزمایش داشته باشد. می‌توان بیان نمود که در مجموعه داده Image segmentation، الگوریتم شبکه عصبی تابع پایه شعاعی پیشنهادی ۴۵ ثانیه زمان آموزش داشته است، در صورتی که نسبت به شبکه عصبی تابع پایه شعاعی استاندارد میزان ۱۲ ثانیه کمتر است. همچنین، خطای مدل پیشنهادی برای فرآیند آزمایش به اندازه ۰.۰۲۸٪ نسبت به شبکه عصبی تابع پایه شعاعی استاندارد کمتر بوده است. برای مجموعه داده Pendigit، مدل خوشه‌بندی پیشنهادی نیاز به ۱۰۵۹ ثانیه آموزش داشته است که به میزان ۴۳۰ ثانیه از مدل شبکه عصبی تابع پایه شعاعی استاندارد کمتر است. به علاوه، خطای مدل خوشه‌بندی پیشنهادی در آزمایش به میزان ۰.۵۴٪ کمتر بوده است.



شکل ۵: مقایسه زمان مدل پیشنهادی با الگوریتم RBF استاندارد

در شکل ۶، میزان خطای آزمایش برای مجموعه داده‌های آزمون آمده است که نشان می‌دهد با مدل خوشه‌بندی پیشنهادی، میزان خطا در شبکه عصبی تابع پایه شعاعی کاهش یافته است، زیرا مراکز تابع گوسی در این شبکه عصبی به درستی تنظیم شده است. در مجموعه داده‌های Letters و Shuttle control که تعداد نمونه‌های بیشتری را دارند، نتایج اختلاف بیشتری را نشان می‌دهد. در مجموعه داده

مجموعه داده‌های مشترک توانسته است که زمان آموزش کمتر و همچنین دقت بیشتری را به همراه داشته باشد.

روش استاندارد	روش پیشنهادی		مدل
	RBF	RBF K-Means++ DBScan	
خطا (MSE)	خطا (S)	خطا (MSE) در	زمان (S) در آموزش
در آموزش	در آموزش	آزمایش	در آموزش
0.016	63	0.013	45
0.026	1352	0.021	1059
0.149	3805	0.142	2504
1.114	6552	1.102	4905

جدول ۵: مقایسه نتایج مدل پیشنهادی با روش RBF HCA-DBSCAN

۵- نتیجه‌گیری و کارهای آینده

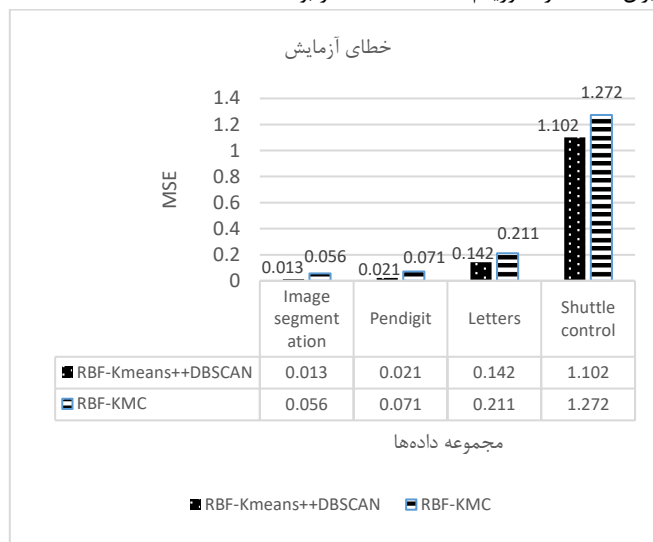
در این مقاله، یک مدل شبکه عصبی تابع پایه شعاعی جدید با بهبود خوشه‌بندی در مرحله بدون نظارت این شبکه ارائه شد. در خوشه‌بندی مطرح شده که مناسب برای داده‌هایی با تعداد بالا می‌باشند، K-Means++ در کنار DBSCAN قرار گرفت که توانسته است در مواجهه با داده‌های بزرگ، سرعت در خوشه‌بندی و دقت در دسته‌بندی را برای شبکه عصبی تابع پایه شعاعی به همراه داشته باشد. نتایج بر روی چهار مجموعه داده با ویژگی‌های متنوع و تعداد نمونه‌های مختلفی آزمایش شد و ارزیابی نتایج نشان داد که الگوریتم پیشنهادی توانسته است نسبت به روش استاندارد شبکه عصبی تابع پایه شعاعی، سرعت آموزش و دقت دسته‌بندی را بهبود دهد.

به‌عنوان کارهای آینده، می‌توان به دو مورد اشاره نمود. در بهبود مدل پیشنهادی، استفاده از خوشه‌بندی بدون نظارت مبتنی بر انبوه داده‌ها می‌تواند با روش‌های دیگری بررسی شود، زیرا در مدل پیشنهادی تأثیر زیادی بر عملکرد مدل داشته است. به‌علاوه، استفاده از مدل پیشنهادی در کاربردهای دنیای واقعی نیز می‌تواند آزمایش شود.

۶- مآخذ

- [1] V. Storey and I. Song, "Big data technologies and management: What conceptual modeling can do," *Data Knowledge Engineering*, no. 108, pp. 50-67, 2017.
- [2] M. Ianni, E. Masciari, G. Mazzeo, and C. Zaniolo, "Efficient big data clustering," *22nd International Database Engineering and Applications Symposium*, ACM, pp. 103-109, 2018.
- [3] P. Arora, D. Deepali, and S. Varshney, "Analysis of K-Means and K-Medoids algorithm for big data," *Proceeding Computer Science*, no. 78, pp. 507-512, 2016.
- [4] A. Jain, M. Murty, and P. Flynn, "Data clustering: a review," *ACM Computer Survey*, vol. 31, no. 3, pp. 264-323, 1999.
- [5] J. Zhu, M. Zeng, J. Huang, S. Liao, C. Cai, and L. Zheng, "Vehicle re-identification using quadruple directional deep learning features," *IEEE Transactions on Intelligent Transportation Systems*, vol. 21, no. 1, pp. 410-420, 2020.
- [6] S. Liu, M. Liu, P. Li, J. Zhao, Z. Zhu, and X. Wang, "SAR image denoising via sparse representation in Shearlet domain based on continuous cycle spinning," *IEEE Transactions on Geoscience and Remote Sensing*, vol. 55, no. 5, pp. 2985-2992, 2017.
- [7] S. Pei, T. Shen, X. Wang, C. Gu, Z. Ning, X. Ye, and N. Xiong, "3DACN: 3D augmented convolutional network for time series data," *Information Science*, no. 513, pp. 17-29, 2020.
- [8] S. Qiao, T. Li, H. Li, J. Peng, and H. Chen, "A new blockmodeling based hierarchical clustering algorithm for web social networks," *Engineering Applications of Artificial Intelligence*, vol. 25, no. 3, pp. 640-647, 2012.
- [9] D. Che, M. Safran, and Z. Peng, "From big data to big data mining: challenges, issues, and opportunities," *International Conference on Database Systems for Advanced Applications*, Springer, pp. 1-15, 2013.
- [10] S. Mirjalili, "Evolutionary Radial Basis Function Networks," *Studies in Computational Intelligence*, pp. 105-139, 2019.

شعاعی پیشنهادی ۴۵ ثانیه زمان آموزش داشته است، در صورتی که نسبت به الگوریتم RBF-KMC میزان ۸ ثانیه کمتر است. به‌علاوه، خطای مدل خوشه‌بندی پیشنهادی برای فرآیند آزمایش به‌اندازه ۰.۰۴۳٪ نسبت به الگوریتم RBF-KMC کمتر بوده است. برای مجموعه داده Pendigit، مدل خوشه‌بندی پیشنهادی نیاز به ۱۰۵۹ ثانیه آموزش داشته است که به میزان ۱۱۸ ثانیه از الگوریتم RBF-KMC کمتر است. همچنین، خطای مدل خوشه‌بندی پیشنهادی به میزان ۰.۰۵٪ از الگوریتم RBF-KMC کمتر بوده است.



شکل ۸: مقایسه خطای مدل پیشنهادی با الگوریتم RBF-KMC

در مجموعه داده Letters، برای مدل خوشه‌بندی پیشنهادی زمان لازم برای آموزش ۲۵۰۴ ثانیه بوده است، اما در الگوریتم RBF-KMC در حدود ۲۶۰۷ ثانیه است. میزان خطا در این مجموعه داده برای مدل خوشه‌بندی پیشنهادی در حدود ۰.۰۷۱٪ خطای کمتری نسبت به الگوریتم RBF-KMC است. در مجموعه داده Shuttle control، برای مدل خوشه‌بندی پیشنهادی زمان لازم برای آموزش ۴۹۰۵ ثانیه بوده است، اما در شبکه عصبی استاندارد ۵۲۹۲ ثانیه است. میزان خطا در این مجموعه داده برای مدل خوشه‌بندی پیشنهادی در حدود ۰.۱۵٪ خطای کمتری نسبت به الگوریتم RBF-KMC بوده است.

در جداول ۴ و ۵، الگوریتم‌های پیشنهاد شده در مرجع [۵۷] شامل K-DBSCAN و HCA-DBSCAN با مدل خوشه‌بندی پیشنهادی مقایسه شده است. همان‌طور که در جدول ۴ مشاهده می‌شود، مدل خوشه‌بندی K-DBSCAN زمان آموزش بیشتری نسبت به مدل خوشه‌بندی پیشنهادی دارد و دقت آن نیز اختلاف زیادی با نتایج مدل خوشه‌بندی پیشنهادی ندارد.

جدول ۴: مقایسه نتایج مدل پیشنهادی با الگوریتم RBF K-DBSCAN

روش استاندارد	روش پیشنهادی		مدل
	RBF K-DBSCAN	RBF K-Means++ DBScan	
خطا (MSE)	خطا (S)	خطا (MSE) در	زمان (S) در آموزش
در آموزش	در آموزش	آزمایش	در آموزش
0.013	52	0.013	45
0.022	1163	0.021	1059
0.145	2564	0.142	2504
1.112	4926	1.102	4905

در جدول ۵ نیز مدل خوشه‌بندی HCA-DBSCAN دارای دقت خوبی است که نزدیک به مدل خوشه‌بندی پیشنهادی عمل کرده است، اما زمان آموزش بالاتری داشته است. میزان بهبود نتایج در مدل خوشه‌بندی پیشنهادی نسبت به مدل بهبود یافته شبکه عصبی تابع پایه شعاعی با بهبود خوشه‌بندی الگوریتم K-Means سراسری تغییر یافته سریع [۵۶]، نشان می‌دهد که مدل خوشه‌بندی پیشنهادی در

- [35] Y. Ding, Y. Zhao, X. Shen, M. Musuvathi, and T. Mytkowicz, "Yinyang K-Means: A drop-in replacement of the classic K-Means with consistent speedup," *Icml-2015*, 2015.
- [36] R. Jothi, S. K. Mohanty, and A. Ojha, "DK-Means: a deterministic K-Means clustering algorithm for gene expression analysis," *Pattern Analysis and Applications*, 2017.
- [37] G. Tzortzis and A. Likas, "The MinMax k-Means clustering algorithm," *Pattern Recognition*, 2014.
- [38] I. Melnykov and V. Melnykov, "K-Means algorithm with the use of mahalanobis distances," *Statistics and Probability Letters*, 2014.
- [39] A. Shirkhorshidi, S. Aghabozorgi, T. Wah, and T. Herawan, "Big data clustering: a review," *International Conference on Computational Science and its Applications*, Springer, pp. 707–720, 2014.
- [40] B. LIU, "A fast density-based clustering algorithm for large databases," *International Conference on Machine Learning and Cybernetics*, IEEE, pp. 996–1000, 2006.
- [41] Y. Wu, J. Guo, and X. Zhang, "A linear DBSCAN algorithm based on Location Sensitive Hashing," *International Conference on Machine Learning and Cybernetics*, IEEE, vol. 5, pp. 2608–2014.
- [42] Y. Dogan, D. Birant, and A. Kut, "SOM++: Integration of self-organizing map and K-Means++ algorithms," *International Conference on Machine Learning and Data Mining in Pattern Recognition*, Springer, pp. 246–259, 2013.
- [43] A. Bakr, N. Ghanem, and M. Ismail, "Efficient incremental density-based algorithm for clustering large datasets," *Alex Engineering Journal*, vol. 54, no. 4, pp. 1147–1154, 2015.
- [44] T. Xu, H. Chiang, G. Liu, and C. Tan, "Hierarchical K-Means method for clustering large-scale advanced metering infrastructure data," *IEEE Transaction Power Delivery*, vol. 32, no. 2, pp. 609–616, 2015.
- [45] H. Ismkhan, "K-Means++: An iterative clustering algorithm based on an enhanced version of the k-means," *PattRecogn*, no. 79, pp. 402–413, 2018.
- [46] D. Brown, A. Japa, and Y. Shi, "A fast density-grid based clustering method Daniel Brown," *IEEE 9th Annual Computing and Communication Workshop and Conference*, IEEE, pp. 48–54, 2019.
- [47] V. Mathur, J. Mehta, and S. Singh, "HCA-DBSCAN: HyperCube accelerated density based spatial clustering for applications with noise," *33rd Conference on Neural Information Processing Systems (arXiv preprint)*, 2019.
- [48] D. Luchi, A. Rodrigues, and F. Varejao, "Sampling approaches for applying DBSCAN to large datasets," *Pattern Recognition Letters*, no. 117, pp. 90–96, 2019.
- [49] Y. Chen, L. Zhou, S. Pei, Z. Yu, Y. Chen, X. Liu, J. Du, and N. Xiong, "KNN-BLOCK DBSCAN Fast Clustering for Large-Scale Data," *IEEE Transactions on Systems, Man, and Cybernetics: Systems*, pp. 1–15, 2019.
- [50] Y. Chen, L. Zhou, N. Bouguila, C. Wang, Y. Chen, J. Du, "BLOCK-DBSCAN Fast clustering for large scale data," *Pattern Recognition*, vol. 109, no. 107627, 2020.
- [51] X. Cui, P. Zhu, X. Yang, K. Li, and C. Ji, "Optimized big data K-Means clustering using MapReduce," *J Supercompu*, vol. 70, no. 3, pp. 1249–1259, 2014.
- [52] A. Sinha, and P. Jana, "A novel K-Means based clustering algorithm for big data," *Conference on Advances in Computing, Communications and Informatics*, IEEE, pp. 1875–1879, 2016.
- [53] H. Song, and J. Lee, "RP-DBSCAN: A superfast parallel DBSCAN algorithm based on random partitioning," *International Conference on Management of Data*, pp. 1173–1187, 2018.
- [54] S. Li, "An Improved DBSCAN Algorithm Based on the Neighbor Similarity and Fast Nearest Neighbor Query," *IEEE Access*, no. 8, pp. 47468–47476, 2020.
- [55] <https://archive.ics.uci.edu/ml/index.php>
- [56] B. Hajji, A. Mellit, G. Marco Tina, A. Rabhi, J. Launay, and S. E. Naimi, *Proceedings of the 2nd International Conference on Electronic Engineering and Renewable Energy Systems*. Lecture Notes in Electrical Engineering, 2021.
- [57] N. Gholizadeh, H. Saadatfar, and N. Hanafi, "K-DBSCAN: An improved DBSCAN algorithm for big data," *The Journal of Supercomputing*, no. 77, pp. 6214–6235, 2021.
- [11] A. Osmanović, S. Halilović, L. A. Ilah, A. Fojnica, and Z. Gromilić, "Machine learning techniques for classification of breast cancer," *IFMBE Proceedings*, 2019.
- [12] M. Ester, H. Kriegel, J. Sander, X. Xu, "A density-based algorithm for discovering clusters in large spatial databases with noise," *2nd International Conference on Knowledge Discovery and Data Mining*, pp. 226–231, 1996.
- [13] A. David, and V. Sergej, "K-Means++: the advantages of careful seeding," *Proceedings of the eighteenth annual ACM-SIAM symposium on Discrete algorithms*, ACM, pp. 1027–1035, 2007.
- [14] A. Katal, M. Wazid, and R. Goudar, "Big data: Issues, challenges, tools and good practices," *6th International Conference on Contemporary Computing*, IEEE, pp. 404–409, 2013.
- [15] S. Shahrivari, "Beyond batch processing: Towards real-time and streaming big data," *Computers*, vol. 3, no. 4, pp. 117–129, 2014.
- [16] S. Mirjalili, "Evolutionary Radial Basis Function Networks," *Studies in Computational Intelligence*, pp. 105–139, 2019.
- [17] C. C. Liao, "Genetic K-Means algorithm-based RBF network for photovoltaic MPP prediction," *Energy*, 2010.
- [18] M. Diez, S. Volpi, A. Serani, F. Stern, and E. F. Campana, "Simulation-Based Design Optimization by Sequential Multi-criterion Adaptive Sampling and Dynamic Radial Basis Functions," *Computational Methods in Applied Sciences*, 2019.
- [19] M. Smolik, V. Skala, and Z. Majdisova, "Advances in Engineering Software Vector field radial basis function approximation," *Advances in Engineering Software*, vol. 123, no. 17, pp. 117–129, 2018.
- [20] M. Diez, S. Volpi, A. Serani, F. Stern, and E. F. Campana, "Simulation-Based Design Optimization by Sequential Multi-criterion Adaptive Sampling and Dynamic Radial Basis Functions," *Computational Methods in Applied Sciences*, 2019.
- [21] T. Wangchamhan, S. Chiewchanwattana, and K. Sunat, "Efficient algorithms based on the k-means and Chaotic League Championship Algorithm for numeric, categorical, and mixed-type data clustering," *Expert Systems with Applications*, 2017.
- [22] Y. Li, X. Wang, S. Sun, X. Ma, and G. Lu, "Forecasting short-term subway passenger flow under special events scenarios using multiscale radial basis function networks," *Transportation Research Part C: Emerging Technologies*, 2017.
- [23] V. K. Chauhan, A. Sharma, and K. Dahiya, "Faster learning by reduction of data access time," *Application Intelligence*, 2018.
- [24] G. Afendras and M. Markatou, "Optimality of training/test size and resampling effectiveness in cross-validation," *Journal of Statistical Planning and Inference*, vol. 16, pp. 1–16, 2018.
- [25] V. Chouvatut, W. Jindaluang, and E. Boonchieng, "Training set size reduction in large dataset problems," *Intelligent Computation Science Engineering Conference*, pp. 1–5, 2015.
- [26] W. A. Yousef and S. Kundu, "Learning algorithms may perform worse with increasing training set size: Algorithmdata incompatibility," *Computational Statistics and Data Analysis*, vol. 74, pp. 181–197, 2014.
- [27] G. Zhang, C. Zhang, and H. Zhang, "Improved K-Means algorithm based on density Canopy," *Knowledge-Based System*, 2018.
- [28] W. L. Zhao, C. H. Deng, and C. W. Ngo, "K-Means: A revisit," *Neurocomputing*, 2018.
- [29] T. Wangchamhan, S. Chiewchanwattana, and K. Sunat, "Efficient algorithms based on the K-Means and Chaotic League Championship Algorithm for numeric, categorical, and mixed-type data clustering," *Expert Systems with Applications*, 2017.
- [30] S. Maldonado, E. Carrizosa, and R. Weber, "Kernel Penalized K-Means: A feature selection method based on Kernel K-Means," *Information Science (Ny)*, 2015.
- [31] H. Ismkhan, "K-Means++: An iterative clustering algorithm based on an enhanced version of the K-Means," *Pattern Recognition*, 2018.
- [32] S. S. Yu, S. W. Chu, C. M. Wang, Y. K. Chan, and T. C. Chang, "Two improved K-Means algorithms," *Applied Soft Computing*, 2018.
- [33] J. Jędrzejowicz and P. Jędrzejowicz, "Distance-based online classifiers," *Expert Systems with Applications*, vol. 60, pp. 249–257, 2016.
- [34] S. Kant and I. A. Ansari, "An improved K means clustering with Atkinson index to classify liver patient dataset," *International Journal of System Assurance Engineering Management*, 2016.

3. Mahala Nobis

4. Location Sensitive Hashing (LSH)

1. Radial Basis Function (RBF)

2. Atkinson

شماره تماس: ۰۹۱۱۲۷۸۱۷۱۲

- 5. Self-Organizing Map (SOM)
- 6. K-Nearest Neighbors (KNN)

سرگذشت: فرشته حاج قاضی دارای مدرک کارشناسی ارشد مهندسی کامپیوتر گرایش نرم افزار از موسسه غیرانتفاعی میرداماد گرگان در سال ۱۳۹۹ و در حال حاضر، دانشجوی دکتری در دانشگاه آزاد اسلامی واحد نیشابور است. زمینه‌های تحقیقاتی او محاسبات گرید و شبکه عصبی است.

معرفی نویسندگان:



نام و نام خانوادگی: رضا قائمی

شماره تماس: ۰۹۱۵۳۱۱۵۷۷۲

سرگذشت: رضا قائمی دارای مدرک دکتری مهندسی کامپیوتر گرایش هوش مصنوعی از دانشگاه ملی یوپی ام مالزی در سال ۲۰۱۱ و در حال حاضر، استادیار گروه مهندسی کامپیوتر و کارشناسی ارشد هوش مصنوعی در دانشگاه آزاد اسلامی واحد قوچان است. زمینه‌های تحقیقاتی او یادگیری ماشین، داده‌کاوی، محاسبات نرم، کلان داده‌ها و اینترنت اشیا است.



نام و نام خانوادگی: یعقوب آراد

شماره تماس: ۰۹۱۵۱۸۸۸۴۹۹

سرگذشت: یعقوب آراد دارای مدرک کارشناسی ارشد مهندسی کامپیوتر گرایش نرم افزار از دانشگاه آزاد اسلامی واحد شیروان در سال ۱۴۰۰ و در حال حاضر، دانشجوی دکتری در دانشگاه آزاد اسلامی واحد نیشابور است. زمینه‌های تحقیقاتی او شبکه‌های حسگر بی‌سیم است.



نام و نام خانوادگی: فرشته حاج قاضی استرآبادی