

ارائه یک روش ترکیبی از الگوریتم ژنتیک برای استفاده در خوش بندی داده ها

میثم رهنما فلاح^۱، مرضیه فریدی ماسوله^{۲*}، محمد رضا عسگری پور^۳

^۱دانشگاه آزاد اسلامی، واحد الکترونیکی، گروه کامپیوتر، تهران، ایران

* نویسنده مسئول مکاتبات

چکیده

داده های بزرگ یا عظیم داده ترجمه اصطلاح Big Data می باشد که معمولاً به مجموعه از داده ها اطلاق می شود که اندازه آنها فراتر از حدی است که با نرم افزارهای معمول بتوان آنها را در یک زمان معقول اخذ، دقیق سازی، مدیریت و پردازش کرد. مفهوم «اندازه» در داده های بزرگ بطور مستمر در حال تغییر است و به مرور بزرگتر می شود. از این رو با رشد روز افزون داده ها و نیاز به بهره برداری و تحلیل از این داده ها، بکارگیری زیرساخت های Big Data از اهمیت ویژه ای برخوردار شده است. ما به ارائه خلاصه ای از بررسی کلی برروی مسایل داده های حجمی می پردازیم که شامل فرصت ها و چالش های داده های حجمی، تکنیک ها و فناوری های فعلی می باشند. برای بهینه سازی از بزرگ الگوریتم های تکاملی و داده کاوی مطرح می شوند. الگوریتم های تکاملی (EA) معروف است که کشف کننده های بهتر در فضای جستجو در مقایسه با تکنیک های سنتی هستند؛ که روش های قطعی را در برمی گیرند و از مکانیزم ها و عملیات ابتدایی برای حل مسئله استفاده می کنند و طی یک سری از تکرارها به راه حل مناسب برای مسئله می رسانند. ما در اینجا از ترکیب الگوریتم K-میانگین و الگوریتم ژنتیک برای رسیدن به نتیجه مطلوب استفاده کردیم که الگوریتم ژنتیک یکی از مشهور ترین و متداول ترین الگوریتم ها در بین الگوریتم های تکاملی است. این الگوریتم یک روش سراسری کمینه سازی است که برای حل مسائل بهینه سازی کاربردهای بسیار دارد.

واژه های کلیدی: داده های بزرگ، الگوریتم های تکاملی، الگوریتم ژنتیک، داده کاوی، K-میانگین.

۱. مقدمه

داده های بزرگ یا عظیم داده ترجمه اصطلاح Big Data می باشد که معمولاً به مجموعه از داده ها اطلاق می شود که اندازه آنها فراتر از حدی است که با نرم افزارهای معمول بتوان آنها را در یک زمان معقول اخذ، دقیق سازی، مدیریت و پردازش کرد. مفهوم «اندازه» در داده های بزرگ بطور مستمر در حال تغییر است و به مرور بزرگتر می شود. داده های بزرگ (Big Data) مجموعه از تکنیک ها و تاکتیک هایی است که نیازمند شکل جدیدی از یکپارچگی هستند تا بتوانند ارزش های بزرگی را که در مجموعه های بزرگ، وسیع، پیچیده و متنوع داده پنهان شده اند، آشکار سازند. از این رو با رشد روز افزون داده ها و نیاز به بهره برداری و تحلیل از این داده ها، بکارگیری زیرساخت های Big Data از اهمیت ویژه ای برخوردار شده است.

عبارت Big Data مدت ها است که برای اشاره به حجم های عظیمی از داده ها که توسط سازمان های بزرگی مانند گوگل یا ناسا ذخیره و تحلیل می شوند مورد استفاده قرار می گیرد؛ اما به تازگی، این عبارت بیشتر برای اشاره به مجموعه های داده ای بزرگی استفاده می شود که به قدری بزرگ و حجمی هستند که با ابزارهای مدیریتی و پایگاه های داده سنتی و معمولی قابل مدیریت نیستند. مشکلات اصلی در کار با این نوع داده ها مربوط به برداشت و جمع آوری، ذخیره سازی، جستجو، اشتراک گذاری، تحلیل و نمایش آنها است. این مبحث، به این دلیل هر روز جذابیت و مقبولیت بیشتری پیدا می کند که با استفاده از تحلیل حجم های بیشتری از داده ها، می توان تحلیل های بهتر و پیشرفته تری را برای مقاصد مختلف، از جمله مقاصد تجاری، پژوهشی و امنیتی، انجام داد و نتایج مناسب تری را دریافت کرد. بیشتر تحلیل های مورد نیاز در پردازش داده های عظیم، توسط دانشمندان در علومی مانند هوشنگسازی، ژنتیک، شبیه سازی های پیچیده فیزیک، تحقیقات زیست شناسی و محیطی، جستجوی اینترنت، تحلیل های اقتصادی و مالی و تجاری مورد استفاده قرار می گیرد. حجم داده های ذخیره شده در مجموعه های داده ای Big Data، عموماً به خاطر تولید و جمع آوری داده ها از مجموعه بزرگی از تجهیزات و ابزارهای مختلف مانند گوشی های موبایل، حسگرهای محیطی، لاغ نرم افزارهای مختلف، دوربین ها، میکروفون ها، دستگاه های تشخیص RFID، شبکه های حسگر بی سیم و غیره با سرعت خیره کننده ای در حال افزایش است. در اینجا ضمن بررسی مفاهیم پایه ای در داده های حجمی، به بررسی راه حل های موجود برای مدیریت و بهره برداری از این نوع داده ها خواهیم پرداخت.

برای ایجاد یک دید مناسب در خصوص داده های حجمی و اهمیت آن، جامعه ای را تصور کنید که در آن جمعیت بطور نمایی در حال افزایش است، اما خدمات و زیرساخت های عمومی آن نتواند پاسخگوی رشد جمعیت باشد و از عهده مدیریت آن برآید. چنین شرایطی در حوزه داده در حال وقوع است؛ بنابراین نیازمند توسعه زیرساخت های فنی برای مدیریت داده و رشد آن در بخش هایی نظری جمع آوری، ذخیره سازی، جستجو، به اشتراک گذاری و تحلیل می باشیم. دستیاری به این توانمندی معادل است با شرایطی که مثلاً بتوانیم "هنگامی که با اطلاعات بیشتری در حوزه سلامت مواجه باشیم، با بازدهی بیشتری سلامت را ارتقا دهیم"، "در شرایطی که خطرات امنیتی افزایش پیدا می کند، سطح امنیت بیشتری را فراهم کنیم"، "وقتی که با رویدادهای بیشتری از نظر آب و هوایی مواجه باشیم، توان پیش بینی دقیقت و بهتری بدست آوریم"، "در دنیایی با خودروهای بیشتر، آمار تصادفات و حوادث را کاهش دهیم"، "تعداد تراکنش های بانکی، بیمه و مالی افزایش پیدا کند، ولی تقلب کمتری را شاهد باشیم"، "با منابع طبیعی کمتر، به انرژی بیشتر و ارزانتری دسترسی داشته باشیم" و بسیاری موارد دیگر از این قبیل که اهمیت پنهان داده های حجمی را نشان می دهد.

داده های حجمی اصطلاحی است که به مجموعه داده هایی اطلاق می شود که مدیریت، کنترل و پردازش آنها فراتر از توانایی ابزارهای نرم افزاری در یک زمان قابل تحمل و مورد انتظار است و چالش های بسیاری را منجر می شوند مانند دشواری ها و پیچیدگی هایی در گرفتن داده ها، ذخیره داده ها، تجزیه و تحلیل داده ها و غیره. نمونه هایی از بزرگ داده گزارش های وب، سامانه های بازشناسی با امواج رادیویی، شبکه های حسگر، شبکه های اجتماعی، متون و اسناد اینترنتی، جستجوهای اینترنتی، نجوم، مدارک پژوهشی، آرشیو عکس، آرشیو ویدیو، پژوهش های زمین شناسی و تجارت در مقیاس بزرگ هستند (چن و همکاران^۱، ۲۰۱۴).

¹ Chen et al.

داده‌های عظیم می‌توانند بر اساس مشخصات زیر تعریف شوند:

حجم: مقدار داده‌های تولید شده در این زمینه بسیار مهم است. اندازه داده‌ها ارزش و پتانسیل داده‌های مورد توجه به آن را تعیین می‌کند تا جایی که می‌توان تصمیم گرفت که داده عظیم محسوب می‌شود یا خیر.

تنوع: جنبه بعدی در داده‌های عظیم تنوع آن است. این بدان معنی است که دسته بندی داده‌های عظیم به ضرورت نیاز شناسایی شده توسط تحلیلگران دارد. این به افراد کمک می‌کند تا داده‌ها و ارتباطاتشان را دقیق‌تر تحلیل کنند تا از مزایا و رعایت اهمیت داده‌های عظیم به طور موثر استفاده کنند.

نرخ تولید: اصطلاح 'نرخ تولید' در این موضوع به سرعت تولید داده اشاره دارد و یا چگونگی سرعت تولید و پردازش داده‌ها برای پاسخگویی به خواسته و چالش‌های پیش رو در مسیر رشد و توسعه است.

ما به ارائه خلاصه ای از بررسی کلی برروی مسایل داده‌های حجمی می‌پردازیم که شامل فرصت‌ها و چالش‌های داده‌های حجمی، تکنیک‌ها و فناوری‌های فعلی می‌باشند. برای بهینه سازی عملکرد و تحلیل داده‌های بزرگ الگوریتم‌های تکاملی مطرح می‌شوند. الگوریتم‌های تکاملی (EA) معروف است که کشف کننده‌های بهتر در فضای جستجو در مقایسه با تکنیک‌های سنتی هستند. (باتاچاریا و همکاران^۱، ۲۰۱۶) که روش‌های قطعی را در برمی‌گیرند و از مکانیزم‌ها و عملیات ابتدایی برای حل مسئله استفاده می‌کنند و طی یک سری از تکرارها به راه حل مناسب برای مسئله می‌رسند. ما در اینجا الگوریتم ژنتیک را مطرح می‌کنیم که یکی از مشهورترین و متداول ترین الگوریتم‌های تکاملی است. این الگوریتم یک روش سراسری کمینه سازی است که برای حل مسائل بهینه سازی کاربردهای بسیار دارد. (وایتلی^۲، ۱۹۹۴)

۲. اهداف پژوهش

در این پژوهش قصد داریم پارامترهای الگوریتم ژنتیک را برای استفاده در کاربردهای داده‌های بزرگ با اتمات‌های یادگیر تنظیم کنیم و همچنین تکنیکی برای حل مشکل در ذخیره‌سازی اطلاعات در داده‌های بزرگ ارائه نمائیم.

۳. فرضیه‌های تحقیق

۱. می‌توان بر اساس مانشین‌های یادگیر پارامترهای الگوریتم ژنتیک را تنظیم کرد.
۲. می‌توان با استفاده از تکنیک‌هایی مشکلات داده‌های بزرگ را حل نمود.
۳. استفاده از الگوریتم بهینه ژنتیک در قیاس با دیگر الگوریتم‌های بهینه سازی در این حوزه بهتر جواب می‌دهد.

۴. مبانی نظری پژوهش

۴.۱ داده‌های حجمی

دانشمندان و محققان بر این باورند که داده‌های حجمی امروزه به اصطلاحی فراگیر و مبحثی مهم تبدیل شده در ادبیات فناوری اطلاعات است که مورد توجه دانشگاه‌ها، دولت‌ها و صنایع قرار گرفته شده است. همان‌طور که سازمان‌ها بزرگ می‌شوند، اطلاعات مرتبط با آنها نیز به صورت تصاعدی رشد می‌نماید و به تبع پیچیدگی‌های مرتبط با داده‌ها افزایش می‌یابد. سرعت داده‌های جدید در حال تولید سرسام آور است. بسیاری از سازمان‌های بزرگ در برنامه‌های کاربردی متفاوت‌شان داده‌های زیادی در فرمت‌های مختلفی را دارند. (کندی^۳، ۱۹۹۵) به همان میزان که داده‌ها گسترش می‌یابند، دسته بندی آنها با یک الگوریتم یا منطق مشخص بسیار دشوار می‌گردد. سازمان‌های بزرگ در واقع با این چالش مواجه هستند که تمامی اطلاعات را در یک پلتفرم نگهداری نمایند و یک دیدگاه ثابت به آن‌ها داشته باشند. این چالش منحصر به فرد برای درک تمامی داده‌هایی

¹ Bhattacharya et al.

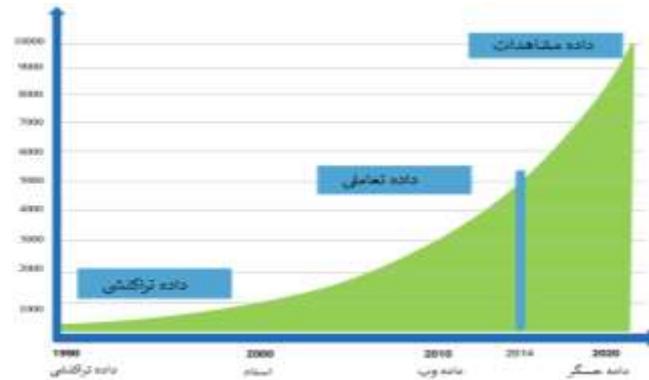
² Whitley

³ Kennedy

که از منابع متفاوت به دست می آیند و استخراج اطلاعات علمی مفید از آن ها، انقلاب داده های حجمی جهانی نامیده می شود. از سویی انجام کسب و کار برخلاف گذشته نیازمند تصمیم گیری هوشمند است. (آبینو و همکاران^۱، ۲۰۱۶)

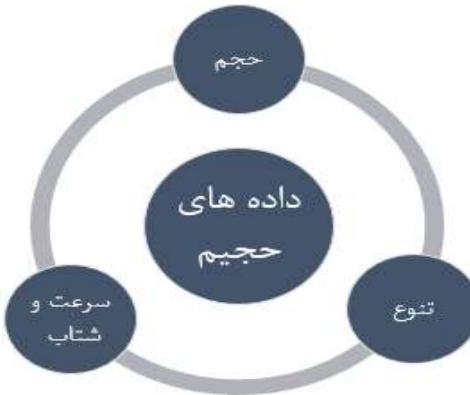
۴.۲ خصوصیات داده های حجمی

موسسه گارتنر (گروه متا) در سال ۲۰۰۱، سه بُعد از چالش ها و فرصت های پیش رو در حوزه رشد داده ها را مطرح کرد. این ویژگی ها به عنوان ویژگی های اصلی و معرف داده های بزرگ مطرح شدند. در واقع این سه واژه هستند که داده های حجمی را تعریف می کنند که در اصطلاح عامیانه به آن ها 3V گفته می شود:



شکل ۱: حجم در مقابل تنوع

سرعت و شتاب (Velocity) / افزایش سرعت تولید داده های ورودی و خروجی: اشاره به سرعت انتقال داده ها دارد. سرعت با افزایش تعداد منابع افزایش یافته است. داده ها از طریق برنامه های کاربردی و سنسورهای بسیار زیادی که در محیط وجود دارند با سرعت بسیار زیاد و به صورت بلاذرگ تولید می شوند که اغلب باید در لحظه پردازش و ذخیره شوند. رشد داده ها و انفجار رسانه های اجتماعی، نگاه ما را به داده ها تغییر داده است. امروزه مردم در شبکه های اجتماعی آخرين رخدادها را برای استفاده دیگران به روز رسانی می کنند. (چانگ و همکاران^۲، ۲۰۱۶)



شکل ۲: خصوصیات داده های حجمی

۴.۳ طبقه بندی داده های حجمی

داده های بزرگ برای درک بهتر ویژگی هایشان به دسته های مختلف طبقه بندی می شوند. شکل ۳ دسته بندی های متعدد از داده های بزرگ را نشان می دهد.

¹ Aibinu et al.

² Chang et al.



شکل ۳: طبقه بندی داده های حجمی

۴.۴ چالش ها در داده های بزرگ

خدمات داده های حجمی، چه در بستر سخت افزاری و چه نرم افزاری، دارای محدودیت هایی است. ما در اینجا برخی از این محدودیت های مهم که به طور مستقیم احساس می شود را لیست می کنیم: دستگاه های ذخیره سازی به یک محدودیت بزرگ تبدیل شده اند. هارد دیسک با تکنولوژی دسترسی تصادفی برای ذخیره سازی داده ها دارای محدودیت به ویژه برای انتقال ورودی/خروجی سریع هستند که برای پردازش داده های بزرگ نیاز است. از سایر محدودیت های ذخیره سازی، می توان محدودیت های طراحی الگوریتم از نظر تعریف ساختمان داده مناسب که جوابگوی دسترسی سریع برای مدیریت داده هاست نام برد. نیاز به طراحی بهینه و پیاده سازی نمایه سازی برای داده ها وجود دارد. ایده جدید ذخیره سازی کلید- مقدار و آرایش سیستم فایل پایگاه داده چالش هایی برای مدیریت داده های بزرگ هستند. (Ding¹, ۲۰۱۶)

۴.۵ تجزیه و تحلیل در داده های بزرگ

تجزیه و تحلیل در داده های بزرگ شامل شش حوزه می باشد که عبارتند از:

- ارتباط: حسگرها و شبکه ها
- ابر: محاسبات و داده ها برخط
- سایبر: مدل و حافظه
- محتوا و زمینه: معانی و ارتباطات
- جوامع: اشتراک و همکاری
- سفارشی سازی: شخصی سازی و ارزش تکنولوژی های داده های بزرگ

۴.۶ داده کاوی

امروزه دیگر فقدان اطلاعات یک مشکل نیست، بلکه عدم توانایی در استخراج اطلاعات مفید از داده ها یک مسئله مهم است. در دو دهه قبل توانایی های فنی بشر برای تولید و جمع آوری داده ها به سرعت افزایش یافته است. بطور کلی استفاده همگانی از وب و اینترنت به عنوان یک سیستم اطلاع رسانی جهانی ما را مواجه با حجم زیادی از داده و اطلاعات می کند. این رشد انفجاری در داده های ذخیره شده، نیاز مبرم وجود تکنولوژی های جدید و ابزارهای خودکاری را ایجاد کرده که به صورت

¹ Ding

هوشمند به انسان یاری رسانند تا این حجم زیاد داده را به اطلاعات و دانش تبدیل کند؛ داده کاوی به عنوان یک راه حل برای این مسائل مطرح می‌باشد. (لیل و همکاران^۱، ۲۰۱۷)

۴.۷ دسته بندی

دسته بندی از رایج‌ترین متدهای داده کاوی است که با عنوان تکنیک‌های یادگیری نظارتی^۲ شناخته می‌شود. دسته بندی مقادیر بعضی متغیرها را محاسبه می‌کند و بر مبنای نتایج حاصله، طبقه بندی می‌کند. الگوریتم‌های مورد استفاده در کلاس بندی عبارتند از: درخت تصمیم، شبکه عصبی و... .

۴.۸ خوشبندی

خوشبندی از متدهای یادگیری بدون نظارت^۳ است که مجموعه‌ای از الگوهای را در گروه‌ها (یا خوشبندی) بخش بندی می‌کند. تحلیل‌های خوشبندی از گروه بندی مجموعه اشیاه داده ای به خوشبندی بر می‌گردد. به طور خاص، هیچ کلاس از پیش تعریف شده ای اختصاص داده نشده است. (جین و همکاران^۴، ۲۰۱۵)

۴.۹ کاربردهای خوشبندی

از آنجا که خوشبندی یک روش یادگیری بدون نظارت محسوب می‌گردد، در موارد بسیاری می‌تواند کاربرد داشته باشد. در بازاریابی (Marketing): دسته‌بندی مشتری‌ها به دسته‌هایی بر حسب رفتارها و نیازهای آنها از طریق مجموعه زیادی از ویژگی‌ها و آخرین خریدهای آنها.

زیست‌شناسی (Biology): دسته‌بندی حیوانات و گیاهان از روی ویژگی‌های آنها

کتابداری: دسته‌بندی کتاب‌ها

نقشه‌برداری شهری (Planning-City): دسته‌بندی خانه‌ها بر اساس نوع و موقعیت جغرافیایی آنها.

۴.۱۰ خوشبندی داده‌های حجمی

خوشبندی به عنوان یکی از روش‌های داده کاوی می‌تواند به تشخیص الگوهای پنهان در داده‌ها و تحلیل حجم عظیم داده‌های شبکه بپردازد. کاربرد خوشبندی داده‌ها تقریباً به وسعت همه حوزه‌های زندگی انسان است. خوشبندی داده یک تکنیک شناخته شده در زمینه‌های مختلف از علوم کامپیوتر و حوزه مرتبط است. می‌توان گفت هر کجا داده ای ذخیره و استفاده می‌شود، پتانسیل فراوانی برای ورود داده کاوی و خوشبندی داده در آن مشاهده می‌شود. (مائولیک^۵، ۲۰۰۰) بنابراین، روش‌های خوشبندی داده با طیف وسیعی از انواع داده و دانش نهفته در داده‌ها روبرو است. داده‌های بیشتر به تحلیل‌های دقیق‌تر می‌انجامد؛ تحلیل‌های دقیق‌تر نیز منجر به تصمیم گیری‌های مطمئن‌تر شده و در پایان تصمیمات بهتر، می‌تواند به معنای کارایی بیشتر عملیات و کاهش هزینه‌ها و ریسک‌ها باشد. خوشبندی ابزاری قوی جهت پردازش داده‌های تولیدشده توسط برنامه‌های مختلف می‌باشد. این تکنیک به عنوان یکی از روش‌های بدون نظارت تشخیص الگوهای پنهان شناخته می‌شود. (برتینو و همکاران، ۲۰۱۱)

¹ Leal et al

² Supervised

³ Unsupervised

⁴ Jin

⁵ Maulik

۴.۱۱ روش‌های خوشبندی داده‌های حجیم

به طور کلی، تکنیک‌های خوشبندی داده‌های بزرگ را می‌توان به دو دسته عمده طبقه بندی کرد: تکنیک‌های خوشبندی تک ماشین و تکنیک‌های خوشبندی چند ماشینه. اخیراً تکنیک‌های خوشبندی چند ماشینه به خاطر انعطاف‌پذیر بودن در مقیاس پذیری و ارائه زمان پاسخ سریع به کاربران، توجه بیشتری را بخود جذب کرده‌اند. (آنچالیا و همکاران^۱، ۲۰۱۳)

۴.۱۲ الگوریتم ژنتیک

در علوم کامپیوتر و ریاضیات الگوریتم جستجو الگوریتمی است که یک مساله را به عنوان ورودی می‌گیرد و پس از ارزیابی راه حل‌های ممکن یک راه حل برای آن مساله بر می‌گرداند. هنگامی که یک مساله را حل می‌کنیم بدنبال راه حل بهینه و یا به عبارتی بهترین پاسخ از بین پاسخ‌های ممکن هستیم. (آرورا^۲، ۲۰۱۴) در الگوریتم ژنتیک با الهام از طبیعت، هدف حل یک مساله جستجو و رسیدن به پاسخ بهینه است. الگوریتم ژنتیک به دلیل تقلید نمودن از طبیعت دارای چند اختلاف اساسی با روش‌های جستجوی مرسوم می‌باشد که در زیر به تعدادی از آنها اشاره می‌کنیم. (برمن^۳، ۲۰۱۳)

- الگوریتم ژنتیک با رشتهداری بیتی کار می‌کند که هر کدام از این رشتهدارها کلّ مجموعه متغیرها را نشان می‌دهد حال آنکه بیشتر روش‌ها به طور مستقل با متغیرهای ویژه برخورد می‌کنند.
- الگوریتم ژنتیک برای راهنمایی جهت جستجو، انتخاب تصادفی انجام می‌دهد که به این ترتیب به اطلاعات مشتق نیاز ندارد.

۴.۱۳ عملگرهای یک الگوریتم ژنتیک

به طور خلاصه الگوریتم ژنتیک از عملگرهای زیر تشکیل شده است:

• کدگزاری

این مرحله شاید مشکل‌ترین مرحله حل مسأله به روش الگوریتم ژنتیک به جای اینکه بر روی پارامترها یا متغیرهای مسأله کار کند، با شکل کد شده آنها سروکار دارد. یکی از روش‌های کد کردن، کد کردن دودویی می‌باشد که در آن هدف تبدیل جواب مسأله به رشتهداری از اعداد باینری (در مبنای ۲) است. (بهیه و همکاران^۴)

• ارزیابی

تابع برازندگی را از اعمال تبدیل مناسب بر روی تابع هدف یعنی تابعی که قرار است بهینه شود به دست می‌آورند. این تابع هر رشتهدار را با یک مقدار عددی ارزیابی می‌کند که کیفیت آن را مشخص می‌نماید. هر چه کیفیت رشتهدار جواب بالاتر باشد مقدار برازندگی جواب بیشتر است و احتمال مشارکت برای تولید نسل بعدی نیز افزایش خواهد یافت. (بوزا و همکاران^۵، ۲۰۱۴)

• ترکیب

مهتمترین عملگر در الگوریتم ژنتیک، عملگر ترکیب است. ترکیب فرآیندی است که در آن نسل قدیمی کروموزوم‌ها با یکدیگر مخلوط و ترکیب می‌شوند تا نسل تازه‌ای از کروموزوم‌ها بوجود بیاید. جفت‌هایی که در قسمت انتخاب به عنوان والد در نظر گرفته شدند در این قسمت ژن‌هایشان را با هم مبادله می‌کنند و اعضای جدید بوجود می‌آورند. ترکیب در الگوریتم ژنتیک باعث از بین رفتن پراکندگی یا تنوع ژنتیکی جمعیت می‌شود زیرا اجازه می‌دهد ژن‌های خوب یکدیگر را بیابند. (چانگ و همکاران، ۲۰۰۸)

¹ Anchalia et al.

² Arora

³ Berman

⁴ Bhih et al.

⁵ Buza et al.

• جهش

جهش نیز عملگر دیگری هست که جواب‌های ممکن دیگری را متولد می‌کند. در الگوریتم ژنتیک بعد از اینکه یک عضو در جمعیت جدید بوجود آمد هر ژن آن با احتمال جهش، جهش می‌یابد. در جهش ممکن است ژنی از مجموعه ژن‌های جمعیت حذف شود یا ژنی که تا به حال در جمعیت وجود نداشته است به آن اضافه شود. جهش یک ژن به معنای تغییر آن ژن است و وابسته به نوع کدگذاری روش‌های متفاوت جهش استفاده می‌شود. (چن و همکاران^۱، ۲۰۱۲)

• رمزگشایی

رمزگشایی، عکسِ عمل رمزگذاری است. در این مرحله بعد از اینکه الگوریتم بهترین جواب را برای مسئله ارائه کرد لازم است عکس عمل رمزگذاری روی جواب‌ها یا همان عمل رمزگشایی اعمال شود تا بتوانیم نسخه واقعی جواب را به وضوح در دست داشته باشیم. (سوکس و همکاران^۲، ۱۹۹۷)

۵. روش تحقیق

در روش خوشه بندی k -میانگین، تعداد خوشه‌ها باید توسط کاربر مشخص شود. تشخیص صحیح تعداد خوشه‌های یک مجموعه داده اغلب برای کاربر کار دشواری است. به همین دلیل روش‌های خوشه بندی که به طور خودکار قادر به تعیین تعداد خوشه‌های مجموعه داده باشند، مورد توجه قرار می‌گیرند. به منظور غلبه بر این نقاط ضعف، ایده استفاده از الگوریتم‌های بهینه سازی برای تعیین مراکز اولیه‌ی خوشه‌ها در الگوریتم k -میانگین مطرح شد. یکی از الگوریتم‌های بهینه سازی کارا در این زمینه، الگوریتم ژنتیک است که به میزان قابل ملاحظه‌ای نتایج خوشه بندی را بهبود می‌بخشد اما به دلیل زمان مصرفی بالای الگوریتم ژنتیک، استفاده از آن برای یافتن مراکز اولیه خوشه‌ها در k -میانگین باعث بالا رفتن زمان این الگوریتم به ویژه برای مجموعه داده‌های بزرگ می‌شود. همچنین، از آنجا که اغلب مجموعه‌های داده‌ای که روش‌های خوشه بندی به منظور تحلیل آنها ایجاد شده‌اند، شامل داده‌های مختلط (داده‌های شامل ویژگی‌های اسمی و عددی) می‌باشند، روش خوشه‌بندی که قابلیت کار بر روی داده‌های مختلط را دارد بسیار مورد توجه خواهد بود. در این بخش روشی برای خوشه بندی مجموعه داده‌های مختلط بزرگ با استفاده از k -میانگین و الگوریتم ژنتیک ارایه شده است. الگوریتم ژنتیک می‌تواند تعداد و مراکز اولیه خوشه‌ها را به طور خودکار پیدا کند. استفاده از الگوریتم ژنتیک برای یافتن مراکز خوشه‌های اولیه در k -میانگین تا حد زیادی باعث بهبود نتایج خوشه‌بندی می‌شود. اما با توجه به زمان مصرفی بالای الگوریتم ژنتیک، خوشه‌بندی مجموعه داده‌های بزرگ به این روش امکان پذیر نمی‌باشد. برای غلبه بر این مشکل و با توجه به ایده مطرح شده در اینجا ابتدا زیر مجموعه‌ای از مجموعه داده‌ی مورد نظر را انتخاب می‌کنیم. با اجرای الگوریتم ژنتیک روی این زیر مجموعه، بهترین مراکز اولیه‌ی خوشه‌ها را پیدا کرده و با ارسال آنها به k -میانگین و اجرای آن روی کل داده‌ها به خوشه‌بندی مجموعه‌ی داده مورد نظر می‌پردازیم. اجرای الگوریتم ژنتیک روی زیرمجموعه‌ای از داده‌ها، به جای کل آن‌ها، تا حد زیادی زمان مصرفی را کاهش می‌دهد. در ادامه به توضیح کامل مراحل کار می‌پردازیم.

۶. یافته‌های پژوهش

۶.۱ مجموعه داده Iris

این مجموعه در واقع مجموعه‌ای از داده‌ها می‌باشد که شامل سه نمونه گل زنبق است که توسط فیشر^۳ در سال ۱۹۳۶ برای نشان دادن تکنیک‌های خطی تفکیک‌پذیر معرفی گردید. از این رو به نام مجموعه داده گل زنبق فیشر نیز خوانده می‌شود. از

¹ Chen et al.

² Cox et al.

³-Sir.Ronald Aylmer Fisher

طرف دیگر، به دلیل این که ادگار اندرسون^۱ نیز این مجموعه را به دلیل کیفیت تنوع جغرافیایی در شبکه جزیره گاسپه، گردآوری کرده است، به مجموعه داده زنبق اندرسون نیز مشهور می‌باشد.

جدول ۱: مراکز خوش به دست آمده با اجرای الگوریتم Kmeans-GA روی مجموعه داده Iris

مراکز خوش به داده Iris		
مرکز خوش سوم	مرکز خوش دوم	مرکز خوش اول
۵.۰۰۳۶	۶.۷۷۴۹	۵.۸۸۹۰
۳.۴۰۳۰	۳.۰۵۲۴	۲.۷۶۱۲
۱.۴۸۵۰	۵.۶۴۶۶	۴.۳۶۴۰
۰.۲۵۱۵	۲.۰۵۳۵	۱.۳۹۷۲

بنا به بهترین مرکز خوش به دست آمده روی مجموعه داده های Iris، مقدار حداقل تابع هزینه الگوریتم Kmeans-GA ۶۰.۵۷۵۹۵۹ گردیده است (جدول ۲). ادعا بر برتری روش پیشنهادی با توجه به نتیجه حاصله امکان‌پذیر می‌باشد.

۶.۲ مجموعه داده Wine

این مجموعه که از آزمایشگاه MCI گرفته شده، از فعل و انفعالات شیمیایی شراب به دست آمده است. این مجموعه داده، نتیجه‌ی جمع‌آوری ۳ نمونه شراب در مکان‌های مختلف ایتالیا می‌باشد. در این مجموعه پارامترها ($k=3$, $d=13$, $N=178$) می‌باشند. از این ۱۷۸ مورد، ۱۰۶ مورد مربوط به آموزش، ۳۶ مورد برای اعتبارسنجی و ۳۶ مورد باقی مانده جهت تست انتخاب گردیده‌اند. لازم به ذکر است که تمام ۱۳ ویژگی، پیوسته می‌باشند و هیچ‌کدام مقدار تهی ندارند. همچنین، توزیع داده‌ها در هر خوش، ۵۹، ۴۸، ۷۱، ۵۹ می‌باشد. ۱۳ ویژگی این مجموعه داده بدین شرح می‌باشد: الكل، اسید مالیک، خاکستر، خاصیت قلیایی خاکستر، منیزیم، اسید فنیک، فلاونوئید، فنول غیرفلاونوئید، پرانسوسیانینز، شدت رنگ، شکل، OD280/OD315 از شراب رقیق و پرولین.

با توجه به جدول ۲ برای مجموعه داده Wine، جواب بهینه الگوریتم پیشنهادی ۱۷۹۶۰۸۲.۷۵۹۵ می‌شود.

جدول ۲: پاسخ الگوریتم‌های موجود بر روی مجموعه داده Wine

تابع هزینه			الگوریتم‌های تکاملی
حداکثر هزینه	متوسط هزینه	حداقل هزینه	
۱۷۹۶۰۸۲.۷۵۹۵	۱۷۹۶۰۸۲.۷۵	۱۷۹۶۰۸۲.۷۵۹۵	Kmeans-GA
۱۷۹۶۴۲۳.۳۵۳۹	۱۷۹۶۳۱۴.۸۴	۱۷۹۶۲۰۶.۳۳۲۰	FCM
۱۸۱۷۳۰۶.۱۸۳۰	۱۸۰۷۱۵۳.۹۹	۱۷۹۷۰۰۱.۷۹۵۴	PSO
۱۸۶۴۴۲۳.۹۴۵۳	۱۸۳۱۹۹۴.۳۳	۱۷۹۹۵۶۴.۷۲۲	SA
۱۸۳۰۸۱۹.۲۰۹۷	۱۸۲۴۰۶۲.۷	۱۸۱۷۳۰۶.۱۸۳۰	TS
۱۸۱۱۳۲۴.۳۴۶۰	۱۸۱۰۶۳۸.۸۵	۱۸۰۹۹۵۳.۳۴۶۰	ACO
۱۸۱۸۷۴۶.۱۸۳۰	۱۸۱۱۹۸۰.۳۳	۱۸۰۵۲۱۴.۴۷۵۸	K-means

^۱-EdgarAnderson

۶.۳ مجموعه داده CMC

این مجموعه زیر مجموعه‌ای از بررسی‌های به انجام آمده جهت پیشگیری از شیوع بارداری خانم‌ها در اندونزی در سال ۱۹۸۷ می‌باشد. نمونه‌های مورد مطالعه خانم‌هایی بودند که در زمان مصاحبه یا باردار بودند و یا این که هنوز این موضوع برای آنها مشخص نشده بود. مشکلی که در این زمینه وجود داشت این بود که محققان نمی‌دانستند از کدام روش پیشگیری از بارداری خانم‌ها بر اساس خصوصیات دموگرافی و اجتماعی- اقتصادی او باید استفاده نمایند. با این فرض که از روش‌های کوتاه مدت و بلند مدت نیز نباید استفاده کنند. پارامترهای این مجموعه به این صورت ($N=۱۴۷۳$, $d=۱۰$, $k=۳$) معرفی گردیدند که سن، تحصیلات، تحصیلات همسر، تعداد فرزندان به دنیا آورده، مذهب، شغل، همسر، شاخص‌های زندگی استاندارد، دسترسی به رسانه‌های جمعی، استفاده از روش‌های جلوگیری از بارداری جزء شاخصه‌های این مجموعه می‌باشد. در جدول ۳ بهترین مرکز خوشة به دست آمده از اجرای برنامه بر روی مجموعه داده CMC آورده شده است.

جدول ۳: مراکز خوشة به دست آمده با اجرای الگوریتم پیشنهادی روی مجموعه داده CMC

مراکز خوشه مجموعه داده CMC		
مراکز خوشه سوم	مراکز خوشه دوم	مراکز خوشه اول
۳۳.۵۳۸۷	۴۳.۹۹۱۵	۲۴.۰۳۳۴
۳.۰۸۹۸	۲.۸۵۰۷	۲.۹۸۳۱
۳.۵۱۱۶	۳.۳۵۲۱	۳.۴۶۲۷
۳.۶۳۰۲	۴.۸۲۳۴	۱.۷۶۴۲
۰.۷۹۰۱	۰.۸۱۴۴	۰.۹۲۸۲
۰.۶۹۷۱	۰.۷۶۹۴	۰.۷۹۴۷
۲.۰۹۹۱	۱.۸۸۳۲	۲.۳۱۸۴
۳.۲۶۸۱	۳.۳۴۱۷	۲.۹۱۲۷
۰.۰۶۷۳	۰.۱۱۲۳	۰.۰۴۵۴
۲.۰۷۳۰	۱.۶۲۴۰	۱.۹۸۴۹

۶.۴ مجموعه داده Vowel

این مجموعه داده‌ها شامل ۸۷۱ نمونه می‌باشند که برخلاف سه مجموعه‌ی قبل، مجموعه داده‌ها که صداهای آواهار تلوگو هندی می‌باشند، به ۶ خوشه تقسیم می‌گردند و ویژگی‌های مورد بررسی در این مجموعه ۳ شاخصه می‌باشد. به عبارت دیگر، این مجموعه داده نمونه‌هایی که دارای بعد کم، متوسط و زیاد می‌باشند؛ را پوشش می‌دهد. در واقع، این ۳ معیار شامل متکلم، آوا و ورودی می‌شود. محدوده‌ی تغییرات هر نطق مقادیر صحیح را به خود اختصاص می‌دهد که برای متکلم‌ها، بین ۰ تا ۰.۸۹، برای آواها بین ۰ تا ۱۰ و برای مقادیر ورودی بین ۰ تا ۹ تخمین زده شده است؛ به عبارت دیگر، پارامترهای این مجموعه را می‌توان به صورت ($N=۸۷۱$, $d=۳$, $k=۶$) نشان داد.

جدول ۴: مراکز خوش به دست آمده با اجرای الگوریتم پیشنهادی روی مجموعه داده Vowel

مراکز خوش مجموعه داده Vowel						
مراکز خوش ششم	مراکز خوش پنجم	مراکز خوش چهارم	مراکز خوش سوم	مراکز خوش دوم	مراکز خوش اول	
۶۴۴.۴	۴۴۲.۶	۳۶۱.۲	۴۰۹.۹	۴۱۵.۴	۵۱۰.۱	
۱۲۹۰.۷	۹۹۷.۲	۲۲۹۳.۹	۲۰۹۴	۱۰۲۷.۵	۱۷۷۳	
۲۲۹۸.۱	۲۶۶۵.۴	۲۹۷۰.۲	۲۶۵۳.۵	۲۳۴۵.۹	۲۵۲۰.۸	

۷. نتیجه گیری

امروزه کاربرد داده کاوی در اکثر علوم به طور چشم‌گیر مشاهده می‌گردد. واضح است در صورتی که بستر مناسبی جهت استفاده از این علم مهیا نگردد، از تکنولوژی روز و بهره‌گیری از پیشرفت‌های به دست آمده دور خواهیم ماند. خوشبندی یکی از ابزار داده کاوی محسوب می‌گردد. از این رو، سهم به سزایی از تحقیقات اخیر معطوف به این روش می‌باشد. الگوریتم Kmeans یک روش یادگیری بدون نظرارت است که در آن تعداد خوش‌ها از قبل تعیین نگردیده و علاوه بر آن یک داده می‌تواند به صورت همزمان به چندین خوش تعلق داشته باشد. متأسفانه این الگوریتم علی‌رغم پیاده‌سازی آسان، با مشکلاتی چون وابستگی به شرایط اولیه، همگرایی زودرس و گیرکردن در بهینه محلی روبه‌رو است.

در صورتی که مقدار اولیه مناسب برای الگوریتم Kmeans انتخاب گردد، امکان همگرایی به نقاط بهینه وجود دارد. لذا، با بهره‌گیری از الگوریتم ژنتیک به عنوان روش بهینه‌سازی نوین، این محدودیت مرتفع گردید. قابل توجه است که نسخه‌ی اصل الگوریتم ژنتیک به تنها‌ی در رفع این مشکل کارساز نبود. به همین دلیل، از روش تولید عدد متضاد جهت انتخاب جمعیت اولیه و روش استراتژی خود تطبیق، برای بهبود عملکرد این الگوریتم استفاده شد. این استراتژی از دو قانون جهش استفاده می‌کند. قانون اول از همگرایی زودرس و قانون دوم از سکون و گیر کردن در بهینه محلی جلوگیری می‌کند.

به این ترتیب، نتایج به دست آمده از ترکیب این دو الگوریتم، علاوه بر نوآوری، محدودیتها را نیز پوشش داد و با توجه به مقایسه‌ی نتایج، مقاوم بودن و کارایی این روش تضمین گردید.

با توجه به شبیه‌سازی‌های انجام گرفته در اغلب موارد، روش ترکیبی جدید بهترین جواب را با حداکثر دقیقت نتیجه می‌دهد. در این روش ترکیبی، علاوه بر این که محدودیت وابستگی به شرایط اولیه تا حدود بسیار زیادی خنثی گردیده، مشکل همگرایی دائمی به پاسخ بهینه محلی نیز تا حد قابل قبولی از بین رفته است.

شبیه‌سازی بر روی مجموعه داده‌ای متنوع اجرا گردیده است. پاسخ‌های نهایی به دست آمده گویای این مطلب می‌باشد که وابستگی به تعداد نمونه‌ها، تعداد خوش‌ها و تعداد پارامترها وجود ندارد و الگوریتم ترکیبی عملکرد مناسبی دارد؛ بنابراین، این الگوریتم را می‌توان بر مجموعه داده‌ها مشابه اعمال نمود و نتیجه مناسب را به دست آورد.

عدم وابستگی الگوریتم به پارامترهای موجود، f_{\max} , N , γ , α ، یکی دیگر از برتری‌های روش ارائه شده می‌باشد.

۸. پیشنهادها

با توجه به تلاش‌های انجام شده می‌توان، ترکیب دو الگوریتم فازی C-means و الگوریتم قربانی را به عنوان کار آینده پیشنهاد نمود.

الگوریتم خوش بندی C میانگین (Fuzzy C-Means): یکی از مهمترین و پرکاربردترین الگوریتم‌های خوش بندی، الگوریتم C میانگین می‌باشد. در این الگوریتم نمونه‌ها به C خوش تقسیم می‌شوند و تعداد C از قبل مشخص شده است. در نسخه فازی این الگوریتم نیز تعداد خوش‌ها (C) از قبل مشخص شده است. الگوریتم (FCM) یکی از معروف‌ترین تکنیک‌های استفاده شده برای خوش بندی است.

الگوریتم جهش ترکیبی قورباغه (SFLA): یک الگوریتم تکاملی و مبتنی بر جمعیت متاهیورستیک جدید است. این الگوریتم سریع است و قابلیت جستجوی سراسری بسیار خوبی دارد.

منابع

1. Aibinu, Abiodun Musa, et al. "A novel Clustering based Genetic Algorithm for route optimization." *Engineering Science and Technology, an International Journal* 19.4 (2016): 2022-2034.
2. Anchalia, Prajesh P., Anjan K. Koundinya, and N. K. Srinath. "MapReduce design of K-means clustering algorithm." *Information Science and Applications (ICISA)*, 2013 International Conference on. IEEE, 2013.
3. Arora, Saurabh, and Inderveer Chana. "A survey of clustering techniques for big data analysis." *Confluence The Next Generation Information Technology Summit (Confluence)*, 2014 5th International Conference-. IEEE, 2014.
4. Berman, Jules J. *Principles of big data: preparing, sharing, and analyzing complex information*. Newnes, 2013.
5. Bertino, Elisa, et al. "Challenges and Opportunities with Big Data." (2011).
6. Bhattacharya, Maumita, Rafiqul Islam, and Jemal Abawajy. "Evolutionary optimization: a big data perspective." *Journal of network and computer applications* 59 (2016): 416-426.
7. Bhih, Amhmed A., Princy Johnson, and Martin Randles. "EM Clustering Approach for Multi-Dimensional Analysis of Big Data Set."
8. Buza, Krisztian, Gábor I. Nagy, and Alexandros Nanopoulos. "Storage-optimizing clustering algorithms for high-dimensional tick data." *Expert Systems with Applications* 41.9 (2014): 4148-4157.
9. Chang, Dongxia et al. "A dynamic niching clustering algorithm based on individual-connectedness and its application to color image segmentation." *Pattern Recognition* 60 (2016): 334-347.
10. Chang, Fay, et al. "Bigtable: A distributed storage system for structured data." *ACM Transactions on Computer Systems (TOCS)* 26.2 (2008): 4.
11. Chen, CL Philip, and Chun-Yang Zhang. "Data-intensive applications, challenges, techniques and technologies: A survey on Big Data." *Information Sciences* 275 (2014): 314-347.
12. Chen, Yanpei, Sara Alspaugh, and Randy Katz. "Interactive analytical processing in big data systems: A cross-industry study of mapreduce workloads." *Proceedings of the VLDB Endowment* 5.12 (2012): 1802-1813.
13. Cox, Michael, and David Ellsworth. "Managing big data for scientific visualization." *ACM Siggraph*. Vol. 97. 1997.
14. D. Whitley, A genetic algorithm tutorial, *Statistics and computing*, 4 (1994) 65-85.
15. J. Kennedy, R. Eberhart, Particle swarm optimization, 1 (1995)
16. Jin, Xiaolong, et al. "Significance and challenges of big data research." *Big Data Research* 2.2 (2015): 59-64.
17. Leal, Emilcy J. Hernández, Néstor D. Duque Méndez, and Julián Moreno Cadavid. "Big Data: an exploration of research, technologies and application cases." *TecnoLógicas* 20.39 (2017).

18. R.J. Kuo, C.H. Mei, F.E. Zulvia, C.Y. Tsai, An application of a metaheuristic algorithm-based clustering ensemble method to APP customer segmentation, *Neurocomputing*, 205 (2016) 116
19. U. Maulik, S. Bandyopadhyay, Genetic algorithm-based clustering technique, *Pattern Recognition*, 33 (2000) 1455-1465.
20. Y. Ding, X. Fu, Kernel-based fuzzy c-means clustering algorithm based on genetic algorithm, *Neurocomputing*, 188 (2016) 233-238.

Providing a Hybrid Approach of Genetic algorithms for Use in Data Clustering

Meysam Rahnamay Fallah¹, Marzieh Faridi Masouleh² *, Mohammad Reza Askarpour³

^{1, 2, 3} Islamic Azad University, Electronic Unit, Computer Department, Tehran, Iran

* corresponding author

Abstract

Big Data is usually referred to as a set of data that exceeds the extent which can be obtained, managed and processed using standard software within a reasonable time. The concept of size in big data is constantly changing and becoming gradually bigger. Hence, with the increasing growth of the data and the need to exploit and analyze these data, it is particularly important to apply the Big Data infrastructures. We provide a summary of a comprehensive review of big data issues, including the opportunities and challenges of big data and the current techniques and technologies. Evolutionary algorithms (Eas) and data mining are used for optimizing the performance and analyzing the big data. EAs are known as better discoverers in search space than the traditional techniques. They include definitive methods and use basic mechanisms and operations to solve the problem and, provide a suitable solution to the problem through a series of repetitions. We have used a combination of the k-mean algorithm and the genetic algorithm in order to arrive at a desirable result, and the genetic algorithm is one of the most popular and commonly used algorithms among the evolutionary algorithms. This algorithm is a global minimization method with many uses for solving optimization problems.

Keywords: big data, evolutionary algorithms, genetic algorithm, data mining, K-mean.
