

انتخاب ویژگی‌های تأثیرگذار در عملیات دسته‌بندی با استفاده از الگوریتم جستجوی فاخته (CSA)

محمد رضا فرهمند^۱، الهام سادات حجازی^۲

^۱ استادیار، هیئت علمی، دانشگاه آزاد اسلامی، واحد ابرکوه.

^۲ کارشناس ارشد، گروه کامپیوتر، دانشگاه آزاد اسلامی، واحد یزد.

نام نویسنده مسئول:

الهام سادات حجازی

چکیده

مسئله انتخاب زیرمجموعه ویژگی‌ها، به مفهوم شناسایی و انتخاب یک زیرمجموعه مفید از ویژگی‌ها، از میان مجموعه داده اولیه می‌باشد و همچنین مبحث مهمی در تحلیل میزان همبستگی در زمینه‌های دسته‌بندی می‌باشد که در کاهش ابعاد مجموعه ویژگی‌ها به کار می‌آید. این کار با حذف ویژگی‌هایی که ایجاد نویز می‌کنند و یا اینکه با دیگر ویژگی‌ها همبستگی کمی دارند، انجام می‌شود. در بسیاری از مجموعه داده‌ها برخی از ویژگی‌ها در تصمیم‌گیری نقشی ندارند و به نوعی می‌توان آن‌ها را اضافی تلقی نمود. پس انتخاب یک زیرمجموعه‌ی مناسب از ورودی‌ها می‌تواند هم در دقت طبقه‌بندی و هم در سرعت آن تأثیر داشته باشد در این مطالعه، با استفاده از الگوریتم جستجوی فاخته و دو معیار درجه تفکیک پذیری و اطلاعات متقابل برای محاسبه همبستگی در میان مشخصه‌ها، رویکرد جدیدی جهت انتخاب ویژگی‌های تأثیر گذار و بهینه، ارائه شده است. ویژگی‌های به دست آمده با سه دسته بندی کننده ماشین بردار پشتیبان، درخت تصمیم و K نزدیک‌ترین همسایه مورد ارزیابی قرار گرفته اند. سپس با مقایسه روش پیشنهاد شده با نتایج به دست آمده از کل مجموعه ویژگی، درجه F و روش‌های انتخاب ویژگی مبتنی بر همبستگی، نشان داده ایم که روش پیشنهاد شده به صورت کلی بسیار کارآمد می‌باشد. در انتها از آنجایی که ماشین بردار پشتیبان جهت دسته بندی، نتایج بهتری را ارائه داده است، با دیگر روش‌های انتخاب ویژگی به علاوه بردار پشتیبان نیز مقایسه انجام گرفته است که راهکار پیشنهادی نسبت به دیگر روش‌ها نتایج بهتری را حاصل نموده است. از آنجایی که روش پیشنهادی یک راهکار انتخاب ویژگی را معرفی می‌کند می‌توان از آن در سایر زمینه‌های تحقیقاتی نظیر تشخیص پزشکی، در سیستم‌های تعیین خسارت و یا سیستم‌های تولیدی استفاده نمود.

واژگان کلیدی: انتخاب ویژگی، الگوریتم جستجوی فاخته، درجه تفکیک پذیری، اطلاعات متقابل با همبستگی دسته‌بندی.

مقدمه

داده کاوی، به دلیل حجم انبوه داده های روزانه و ضرورت تبدیل این داده ها به اطلاعات مفید، سریع ترین زیرحوزه رو به رشد فناوری اطلاعات محسوب می شود [۱]. داده کاوی شامل پیش پردازش چندگانه (ادغام، فیلتر، تبدیل، کاهش، و غیره)، ارائه دانش و همچنین ارزیابی الگو می گردد [۲]. یکی از مراحل پیش پردازش اصلی، انتخاب ویژگی نام دارد که هدف آن پاک کردن ویژگی های غیر مرتبط و غیرفعال مجموعه داده های خاص است [۳]. وجود مقادیر کلان داده های موجود در مسائل واقعی در جهان و حضور داده های نامرتب و حشو روند تحلیل آنها را چالش برانگیز و اغلب نادرست می سازد [۴]. انتخاب ویژگی (FS) یک گام پیش پردازش است که ویژگی های نامرتب و زائد را حذف می کند و مجموعه نهایی از مفیدترین ویژگی ها را پیدا می کند که به نوبه خود به عملکرد بهتر در تکنیک های داده کاوی (به عنوان مثال طبقه بندی) منجر می شود [۵]. ویژگی های نامرتب و زائد، جستجوی بیشتری در بر می گیرند؛ چرا که آنها گوهایی با قابلیت تشخیص کم و قوانین لازم برای پیش بینی یا طبقه بندی نسبتاً آشکار و نیز ریسک بیش برآزش بالا می سازند. انتخاب زیر مجموعه های ویژگی نیاز به تعیین ویژگی مناسب برای به حداکثر رساندن دقت پیش بینی یا طبقه بندی دارد. در واقع، طبقه بندی سریع و دقیق، با استفاده از حداقل تعداد ویژگی ها انتخاب می شود. این ظاهراً از طریق انتخاب ویژگی به دست می آید [۶].

به طور گسترده طبقه بندی داده ها در وظایف مختلف و زمینه های چند رسانه ای مانند طبقه بندی تصویر [۷]، [۸]، طبقه بندی متن [۹]، طبقه بندی عواطف [۱۰]، تقسیم بندی [۱۱]، تشخیص شی [۱۲]، [۱۳]، [۱۴]، [۱۵]، تشخیص رویداد [۱۶]، تشخیص برجستگی [۱۷]، [۱۸]، ردیابی بصری [۱۹] و غیره کاربرد دارد. با این حال، با رشد انفجاری داده ها، و اینکه چنین داده هایی همیشه ساختار داده پیچیده دارند (به عنوان مثال چند کلاسی، چندین دیدگاه و بُعد بالا)، طبقه بندی داده ها هنوز هم یک کار بسیار پیچیده است [۲۰].

در این تحقیق روشی خودکار و سریع جهت انتخاب تأثیرگذارترین ویژگی های ارائه می گردد که در آن با استفاده از الگوهای آزمایشی و استفاده از الگوریتم های مختلف، توانایی هر ویژگی به تنهایی و در کنار سایر ویژگی ها ارزیابی و نهایتاً ویژگی هایی که بهترین میزان تابع ارزیابی عملکرد را از خود نشان دهند انتخاب می گردند.

انتخاب ویژگی

روش انتخاب ویژگی برای حذف اصطلاحات بی اهمیت و کاهش ابعاد کلان مجموعه ویژگی ها و بهینه سازی کارایی و اثربخشی طبقه بندی استفاده می شود [۲۱]. مسأله انتخاب ویژگی، یکی از مسائلی است که در مبحث یادگیری ماشین و همچنین شناسایی آماری الگو مطرح است. یافتن مجموعه ای بهینه از ویژگی ها چالش برانگیز و از نظر محاسباتی هزینه بر است. به تازگی، الگوریتم های فراابتکاری به عنوان ابزارهایی موثر و قابل اطمینان برای حل مسائل بهینه سازی (به عنوان مثال، یادگیری ماشین، مسأله داده کاوی، طراحی مهندسی و انتخاب ویژگی) معرفی شدند [۲۲]، [۳]. این مسئله در بسیاری از کاربردها (مانند طبقه بندی) اهمیت به سزایی دارد، زیرا در این کاربردها تعداد زیادی ویژگی وجود دارد، که بسیاری از آن ها یا بلااستفاده هستند و یا اینکه بار اطلاعاتی چندانی ندارند. حذف نکردن این ویژگی ها مشکلی از لحاظ اطلاعاتی ایجاد نمی کند ولی بار محاسباتی را برای کاربرد موردنظر بالا می برد و علاوه بر این باعث می شود که اطلاعات غیرمفید زیادی را به همراه داده های مفید ذخیره کنیم. برای مسئله انتخاب ویژگی، راه حل ها و الگوریتم های فراوانی ارائه شده است که بعضی از آن ها قدمت سی یا چهل ساله دارند. مشکل بعضی از الگوریتم ها در زمانی که ارائه شده بودند، بار محاسباتی زیاد آن ها بود، اگرچه امروزه با ظهور کامپیوترهای سریع و منابع ذخیره سازی بزرگ این مشکل، به چشم نمی آید ولی از طرف دیگر، مجموعه های داده ای بسیار بزرگ برای مسائل جدید باعث شده است که همچنان پیدا کردن یک الگوریتم سریع برای این کار مهم باشد. انتخاب ویژگی (حذف متغیر) در درک اطلاعات، کاهش نیاز به محاسبه، کاهش اثر مشقت بعدچندی و بهبود عملکرد پیش بینی کننده همواره موثر بوده است [۲۳].

انتخاب زیر مجموعه های ویژگی به تعیین ویژگی مناسب برای افزایش دقت پیش بینی یا طبقه بندی نیاز دارد. هدف اصلی تحقیقات موجود، تعیین زیر مجموعه ویژگی بهینه است. انتخاب ویژگی ها معمولاً بر اساس پارامترهای زمان محاسباتی و کیفیت راه حل های زیر مجموعه ویژگی های تولید شده انجام می شود. در واقع، طبقه بندی سریع و دقیق، با استفاده از حداقل تعداد ویژگی ها انتخاب می شود. بدیهی است که این امر از طریق انتخاب ویژگی ها به دست می آید [۲۴].

روش انتخاب ویژگی (FS)، شرایط را با بالاترین نمره بر اساس معیار پیش تعیین شده اهمیت شرایط حفظ می کند. به طور گسترده مشاهده شد که انتخاب ویژگی می تواند ابزاری قدرتمند برای ساده سازی و یا تسریع محاسبات باشد، و در صورت استفاده مناسب و صحیح، اتلاف و زیان کمتری در کیفیت طبقه بندی به همراه دارد. هدف اصلی روش انتخاب ویژگی، بهبود کارایی طبقه بندی و کارایی محاسباتی است. به رغم رویکردهای متعدد در ادبیات، انتخاب ویژگی هنوز هم در زمره موضوعات پژوهشی در حال پیشرفت قرار دارد. محققان همچنان

به دنبال تکنیک‌های جدید برای انتخاب ویژگی‌های متمایز هستند تا دقت طبقه‌بندی را بهبود ببخشند و زمان پردازش نیز کاهش یابد. تعداد زیادی از تکنیک‌های مبتنی بر فیلتر و تکنیک‌های بسته‌بند برای انتخاب ویژگی‌های متمایز در طبقه‌بندی متن وجود دارد. روش‌های انتخاب ویژگی خودکار (FS) عبارتند از حذف عبارات غیرآموزنده و کم‌اهمیت بر اساس آمارمجموع [۲۱]، [۲۵].

روش‌های انتخاب ویژگی

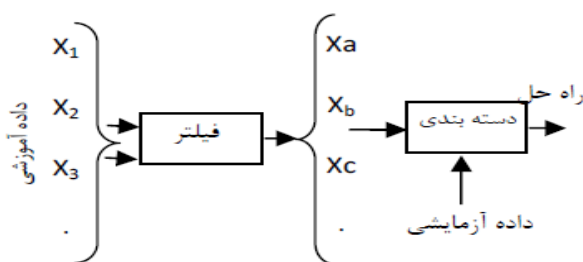
FS را می‌توان نوعی مساله بهینه‌سازی در نظر گرفت. با توجه به دشواری این مساله و داشتن تعداد زیادی از راه‌حل‌های محلی، الگوریتم‌های بهینه‌سازی تصادفی، روش‌هایی امیدوارکننده برای حل این مشکل هستند. هدف FS یافتن زیرمجموعه‌ای از ویژگی‌های M از مجموعه کامل با ویژگی‌های N است، این کار باعث بهبود عملکرد الگوریتم یادگیری (سرعت یادگیری و یا دقت طبقه‌بندی) می‌شود. FS می‌تواند به عنوان مساله بهینه‌سازی در نظر گرفته شود؛ چون زیرمجموعه (نزدیک) مطلوب را جستجو می‌کند [۵].

الگوریتم‌های انتخاب ویژگی به دو دسته اصلی یعنی الگوریتم‌های مبتنی بر بسته‌بند و الگوریتم‌های مبتنی بر فیلتر تقسیم می‌شوند [۲۶]. الگوریتم‌های مبتنی بر بسته‌بند به استفاده از الگوریتم‌های یادگیری ماشین برای ارزیابی آنها وابستگی دارد. از سوی دیگر، الگوریتم‌های مبتنی بر فیلتر از روش‌های آماری برای انتخاب زیرمجموعه ویژگی استفاده می‌کنند. با این حال، الگوریتم‌های مبتنی بر بسته‌بند نتایج بهتری کسب می‌کنند؛ و از لحاظ محاسباتی گران هستند. بنابراین، برای کاهش زمان محاسبات الگوریتم جستجوی هوشمند مورد نیاز است [۲۷] [۲۸].

روش‌های انتخاب ویژگی سنتی را می‌توان به روش فیلتر، بسته‌بند، و هیبرید تقسیم بندی کرد. روش هیبرید، همانطور که از نامش پیداست، ویژگی‌های هر دو روش فیلتر و بسته‌بند را ترکیب می‌کند. هدی و همکاران (۲۰۱۱) امتیاز رتبه‌بندی ویژگی فیلتر را در گام بسته‌بند گنجانند تا روند جستجو را تسریع بخشند. ژو و همکاران (۲۰۰۷) روش رتبه‌بندی فیلتر را در الگوریتم تکاملی ممتیک ادغام کردند، این روش جستجو و شناسایی زیرمجموعه ویژگی‌های اصلی را تسریع می‌کند [۲۹]. روش‌های حذف متغیر به طور کلی به روش‌های فیلتر (Filter) و بسته‌بندها (Wrapper) طبقه‌بندی شده‌اند. روش‌های فیلتر به عنوان پیش‌پردازش برای رتبه‌بندی ویژگی‌هایی عمل می‌کنند که در آن ویژگی‌های عالی رتبه انتخاب و به پیش‌بینی‌کننده اعمال می‌شوند. در روش بسته‌بندها، معیار انتخاب ویژگی عملکرد پیش‌بینی‌کننده است یعنی پیش‌بینی‌کننده بر روی الگوریتم جستجو پیچیده می‌شود تا زیرمجموعه‌ای را پیدا کند که بالاترین عملکرد پیش‌بینی‌کننده را نشان می‌دهد. روش‌های تعبیه شده شامل انتخاب متغیر به عنوان بخشی از روند آموزشی بدون تقسیم داده‌ها به مجموعه‌های آموزش و آزمون است [۲۳].

روش فیلتر

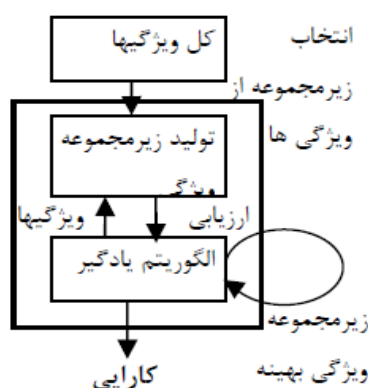
در این روش، از هیچ تابع دسته‌بندی استفاده نمی‌شود. به عبارت دیگر از هیچ فیدبکی از الگوریتم یادگیری اعمال شده استفاده نخواهد شد. این یک روش از پیش انتخاب‌شده‌ای است که مستقل از الگوریتم یادگیری ماشین اعمال شده می‌باشد. زیرمجموعه‌های ویژگی به وسیله مفاهیم دیگری ارزیابی می‌شوند. در الگوریتم Focus نوعی از روش‌های فیلتر (یک جستجوی جامعی برای آزمایش همه زیرمجموعه‌های ویژگی انجام می‌شود) و سپس این روش، زیرمجموعه‌ای با مینیمم تعداد ویژگی‌هایی را مشخص می‌کند که می‌توانند نمونه‌های مجموعه آموزشی را با دقت قابل قبولی دسته‌بندی کنند. Relief [۳۰] یک روش جستجوی تصادفی بر اساس مدل روش فیلتر می‌باشد که در این روش، به هر ویژگی بر اساس، مرتبط بودن با مفهوم هدف، یک وزن نسبت داده می‌شود و نمونه‌ها به‌طور تصادفی برای پیدا کردن ویژگی‌های مرتبط انتخاب می‌شوند. روش فیلتر به صورت زیر کار می‌کند: در مرحله اول، برای هر ویژگی یک رتبه اعتبار محاسبه می‌شود. سپس این رتبه‌ها مرتب‌شده و ویژگی‌هایی با کمترین رتبه حذف می‌شوند. (از یک آستانه برای رتبه‌های ویژگی استفاده می‌شود). سپس نتایج زیرمجموعه‌ای از ویژگی‌ها به عنوان ورودی به یک سیستم دسته‌بند داده می‌شود. شکل (۱) طرز کار روش فیلتر را نشان می‌دهد.



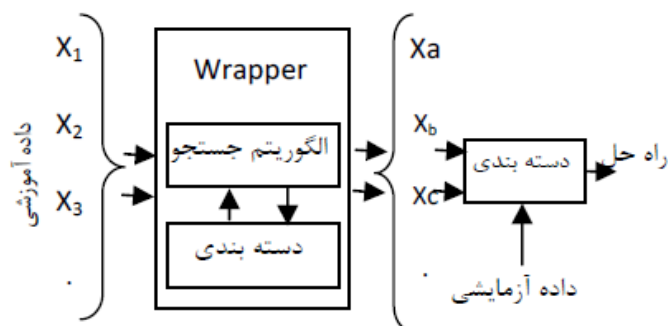
شکل ۱: طرز کار روش فیلتر

روش Wrapper

این روش به جعبه سیاه معروف است. در این روش از یک تابع دسته‌بندی برای ارزیابی شایستگی زیرمجموعه‌های ویژگی استفاده می‌شود. این روش از فیدبکی از الگوریتم یادگیری اعمال شده استفاده می‌کند. از یک الگوریتم ژنتیک هم برای جستجوی ویژگی‌های معتبر استفاده شده است. دلیل استفاده از الگوریتم ژنتیک این است که این الگوریتم می‌تواند یک جستجوی تصادفی را انجام دهد و مستعد گیر افتادن در مینیمم محلی نمی‌باشد. به هر حال عملگر باز ترکیبی استفاده شده در این الگوریتم، مفهوم راه‌حل‌های ترکیبی را دارد، در حالی که انتخاب‌های موفقیت‌آمیز ویژگی قبلی را حفظ می‌کند. به عبارت دیگر این روش، یک روش باز خورد می‌باشد که از الگوریتم یادگیری ماشین در فرآیند انتخاب ویژگی استفاده می‌کند. ارزیابی به وسیله اجرای الگوریتم استقرایی در طول فازهای یادگیری و آزمون در هر انتخاب ویژگی، انجام می‌شود. در این مقاله از شبکه عصبی برای این کار استفاده شده است. شکل (۲) انتخاب ویژگی را با استفاده از فرآیند Wrapper و شکل (۳) طرز کار این فرآیند را نشان می‌دهد.



شکل ۲: انتخاب ویژگی با استفاده از فرآیند Wrapper



شکل ۳: طرز کار روش Wrapper

الگوریتم جستجوی فاخته

روش جستجوی فاخته (CS) یک روش بهینه‌سازی فرا اکتشافی است که رویکردی تکاملی در جستجوی راه‌حل بهینه دارد و در سال ۲۰۰۹ توسط Yang و Deb پیشنهاد شده است. این روش از رفتار جالب توجه گونه‌هایی از پرنده‌ی فاخته در پرورش تخم الهام گرفته است و آن را با پرواز لووی که نوعی گشت تصادفی است ترکیب می‌کند. برخی از گونه‌های فاخته به جای ساختن لانه، تخم‌های خود را در لانه‌ی پرنده‌ای از گونه‌های دیگر می‌گذارند و آن‌ها را با تقلید از شکل تخم‌ها و جوجه‌های پرنده‌ی میزبان وادار به مشارکت در بقای نسل خود می‌کنند.

جستجوی فاخته یک الگوریتم فرا ابتکاری که تقلیدی از استراتژی تولید مثل هوشمندانه کوکو است. این پرنده به دقت لانه میزبان را از دیگر پرندگان انتخاب می‌کند و تخم خود را در میان تخم‌های موجود می‌گذارد. پرنده میزبان اشتباهی از تخم کوکو مراقبت می‌کند. با این حال برخی ممکن است تخم را تشخیص دهند و آن را از بین ببرند یا به خارج از لانه بفرستند. از رفتار هوشمندانه این پرنده برای توسعه یک الگوریتم بهینه‌سازی جدید تقلید شده است. این الگوریتم شامل مجموعه‌ای از لانه‌های میزبان با یک تخم است (راه‌حل) که هر

کدام یک نسل را ایجاد می‌کند. بهترین راه‌حل به دنبال تولید راه‌حلی جدید بر اساس یکی از راه‌حل‌های موجود و تغییر مشخصه‌های خاصی است.

همانند سایر الگوریتم‌های تکاملی، این روش هم با یک جمعیت اولیه کار خود را شروع می‌کند (جمعیتی متشکل از فاخته‌ها). این جمعیت از فاخته‌ها تعدادی تخم دارند که آن‌ها را در لانه تعدادی پرنده‌ی میزبان خواهند گذاشت. تعدادی از این تخم‌ها که شباهت بیشتری به تخم‌های پرنده میزبان دارند شانس بیشتری برای رشد و تبدیل شدن به فاخته بالغ خواهند داشت. سایر تخم‌ها توسط پرنده میزبان شناسایی شده و از بین می‌روند. میزان تخم‌های رشد کرده، مناسب بودن لانه‌های آن منطقه را نشان می‌دهند. هرچه تخم‌های بیشتری در یک ناحیه قادر به زیست باشند و نجات یابند به همان اندازه سود بیشتری به آن منطقه اختصاص می‌یابد. بنابراین موقعیتی که در آن بیشترین تعداد تخم‌ها نجات یابند پارامتری خواهد بود که الگوریتم فاخته قصد بهینه‌سازی آن را دارد.

روش پیشنهادی

کاهش بُعد می‌تواند با دو رویکرد به دست آید: استخراج ویژگی و انتخاب ویژگی. استخراج ویژگی بیانگر تبدیل خطی یا غیرخطی از فضای اصلی ویژگی به یک فضای جدید با ابعاد کمتر می‌باشد. از سوی دیگر، انتخاب ویژگی، زیرمجموعه‌ای از ویژگی را به واسطه انتخاب ویژگی‌های مهم از نمونه‌های اصلی بدون هیچ‌گونه تغییری، تولید می‌کند. به صورت کلی، یک زیرمجموعه ویژگی خوب می‌بایست دارای خصوصیات زیر باشد. اولاً، ویژگی‌های انتخاب شده به خوبی می‌تواند ارائه‌دهنده الگوهای ورودی باشد به طوری که ویژگی‌های نامربوط انتخاب نشده، تنها فضای جستجو را بزرگ‌تر کنند. دوماً، ویژگی‌های انتخاب شده دربرگیرنده همه اطلاعات استفاده شده برای تمایز گذاری الگوها با چندین برجسب دسته داده می‌باشد. چنین زیرمجموعه ویژگی‌ای می‌تواند به عملکرد رضایت‌بخش و مورد انتظار دسته‌بندی کننده فارغ از اینکه از چه الگوریتم آموزشی استفاده شده است، دست یابد. سوماً، حذف ویژگی‌های زائد می‌تواند اندازه مجموعه داده‌ها را کاهش دهد. در صورتی که ویژگی‌های نامربوط و زاید به صورتی نامناسب انتخاب شوند، فرآیند یادگیری را تحت‌الشعاع قرار داده و باعث می‌شود تا فرآیند آموزش ناکارآمد گردد.

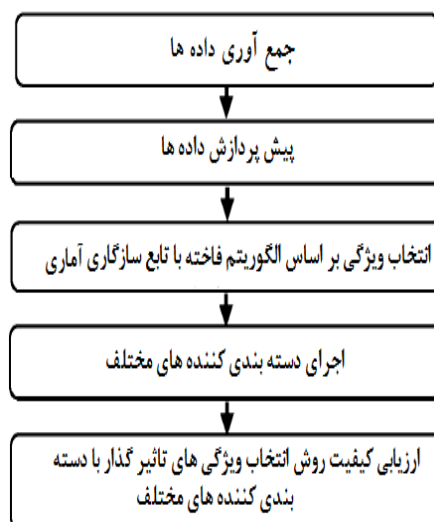
پیدا کردن زیرمجموعه ویژگی بهینه F_s از مجموعه بزرگ ویژگی F که منجر به رشد نمایی^۱ (تعریفی) فضای جستجو (2^d)، که d در اینجا دلالت بر اندازه F دارد) می‌شود که از نظر محاسباتی جذاب نمی‌باشد. بسیاری از روش‌های جستجوی غیرمستدل در پی پیدا کردن راه‌حل‌های نزدیک به بهینه بوده‌اند. از آنجایی که الگوریتم فاخته یکی از ابزارهای بهینه‌سازی می‌باشد که برای پیدا کردن راه‌حل‌ها در فضاهای جستجوی پیچیده و غیرخطی به گستردگی مورد استفاده قرار گرفته، به صورت طبیعی برای حل مسائل انتخاب ویژگی نیز بکار گرفته شده است.

الگوریتم بهینه‌سازی فاخته (COA) یک الگوریتم جدید جمعیت محور است که از روش زندگی گونه‌ای پرنده به نام فاخته الهام گرفته است. در این تحقیق ما رویکردی جدید را بر اساس COA برای انتخاب زیرمجموعه ویژگی معرفی می‌کنیم. برای بررسی کارایی الگوریتم، آزمایشات بر روی چند مجموعه داده انجام شده است و نتایج نشان می‌دهد که روش پیشنهادی می‌تواند راه‌حلی بهینه را برای مساله انتخاب زیرمجموعه ویژگی فراهم آورد.

چارچوب روش پیشنهادی

ما در این بخش یک سیستم هوشمند جدید را برای انتخاب ویژگی ارائه می‌دهیم. این سیستم بر روی انتخاب مبتنی بر الگوریتم جستجوی فاخته استوار است که برای انتخاب یک زیرمجموعه مناسب ویژگی‌ها، از یک معیار استقلال و تفکیک‌پذیری و یک معیار همبستگی مبتنی بر اطلاعات متقابل برای کاهش بعد داده‌ها و کسب دقت بالای دسته بندی کننده، تمرکز نموده است. از آنجایی که نسبت ماتریس پراکندگی بین دسته‌ای به پراکندگی دورن دسته‌ای زیرمجموعه‌ها، می‌تواند مشارکت آن‌ها در دسته‌بندی را ارائه دهد، زیرمجموعه ویژگی بهینه بر حسب نمره مستقل بودن دسته انتخاب می‌گردد. بعلاوه، متغیرهای زائد نیز در فرآیند انتخاب ویژگی مدنظر قرار گرفته‌اند. ماتریس پراکندگی میان دسته‌ای بر رابطه بین متغیرها بر حسب برجسب‌های دسته‌شان، دلالت می‌کند و هرچه که این معیار بزرگ‌تر باشد و در سوی دیگر افزایش همبستگی ویژگی‌ها با کلاس داده و کاهش افزونگی بین ویژگی‌های حاصل، نشان دهنده اینکه رو متغیر کمتر تکراری می‌باشند. شکل (۴) زیر مراحل روش پیشنهادی را به خوبی نشان می‌دهد.

¹ Exponential Growth



شکل ۴: فلوجارت کلی روش پیشنهادی

رویکرد پیشنهاد شده با داده‌های استاندارد تأیید شده و راه‌حل معرفی شده برای آموزش سه دسته‌بندی کننده مورد استفاده قرار گرفته‌اند، این دسته‌بندی کننده‌ها عبارتند از: ماشین بردار پشتیبانی (SVM)، درخت تصمیم (Dt) و K نزدیک‌ترین همسایه (KNN). نتایج ما همچنین به ترتیب با نتایج کسب شده توسط کل مجموعه ویژگی، نمره F^2 و انتخاب ویژگی مبتنی بر همبستگی^۲ (CFS)، مورد مقایسه قرار گرفته‌اند.

رویکرد پیشنهادی بر اساس الگوریتم جستجوی فاخته

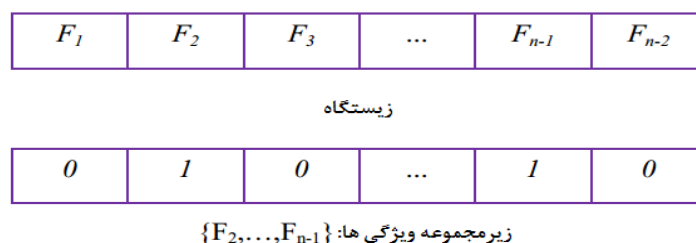
در این بخش، الگوریتم پیشنهادی برای انتخاب ویژگی ارائه شده است. گام‌های الگوریتم پیشنهادی با جزئیات در ادامه ارائه شده است.

ایجاد سکونت‌گاه اولیه فاخته

در الگوریتم ژنتیک و بهینه‌سازی ازدحام ذرات، هر راه‌حل از مسئله را به ترتیب کرموزوم و مکان ذره نامیدیم. اما در COA آن را سکونت‌گاه^۴ می‌نامیم. در یک مساله N بعدی، یک سکونت‌گاه آرایه‌ای $1 \times N$ است که نشان‌دهنده محل فعلی زندگی فاخته است. این آرایه به صورت زیر است.

$$habitat = (x_1, x_2, x_3, \dots, x_n)$$

در رویکرد پیشنهادی برای انتخاب ویژگی، هر سکونت‌گاه، یک رشته‌ای از اعداد دودویی است. هنگامی که مقدار متغیر ۱ است، پس این ویژگی انتخاب می‌شود و هرگاه ۰ است ویژگی متناظرش انتخاب نمی‌شود. شکل (۵) نمایش ویژگی را به عنوان یک سکونت‌گاه در رویکرد پیشنهادی را نشان می‌دهد.



شکل ۵: مثالی از بیان ویژگی‌ها در روش پیشنهادی

² F-score

³ Correlation-based feature selection

⁴ habitat

سود و یا شایستگی یک سکونت‌گاه تحت یک تابع شایستگی که در ادامه بیان می‌شود، تعریف می‌شود. بیشتر طبقه‌بندی‌کننده‌ها می‌توانند برای محاسبه بهره استفاده شوند. برای مثال K نزدیک‌ترین همسایه (KNN)، درخت تصمیم (DT) و ماشین‌های پشتیبان بردار (SVM) سه طبقه بند محبوب هستند

الگوریتم با N_{pop} سکونت‌گاه اولیه تصادفی در اندازه جمعیت آغاز می‌شود. سکونت‌گاه واقعی فاخته جایی است که آن‌ها با بیشترین فاصله از سکونت‌گاهشان تخم‌گذاری می‌کنند. این ماکزیمم محدوده را شعاع تخم‌گذاری^۵ (ELR) گویند که به صورت زیر تعریف می‌شود:

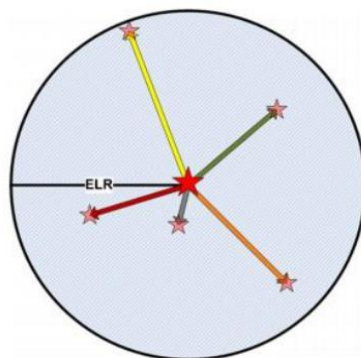
رابطه (۱)

$$ELR = \alpha \times \frac{\text{تعداد تخم‌های فاخته}}{\text{تعداد نهایی تخم‌ها}} \times (\text{Var}_{hi} - \text{Var}_{low})$$

در اینجا α یک عدد صحیح بوده که برای کنترل حداکثر اندازه ELR در نظر گرفته می‌شود. Var_{hi} و Var_{low} به ترتیب کران پایین و بالای تعداد تخم‌های هر فاخته می‌باشد. α عددی صحیح است، var_{hi} و var_{low} کران بالا و پایین به ترتیب هستند. مطابق با معادله بالا، ELR متناسب با تمام تخم‌ها، تعداد تخم‌های فعلی و محدودیت‌های متغیر است.

تخم‌گذاری فاخته

هر فاخته تخم‌گذاری را بطور تصادفی در دیگر لانه‌های پرندگان میزبان در محدوده ELR شان آغاز می‌کنند. شکل ۶ تصویری واضح از این مفهوم را ارائه می‌دهد. ستاره قرمز در وسط شکل، زیستگاه اولیه فاخته‌ای با ۵ تخم است. ستاره‌های دیگر محل هر یک از تخم‌های فاخته را نشان می‌دهد. بعد از تخم‌گذاری، تخم‌ها با کمترین مقدار بهره، شناسایی و از بین برده می‌شوند. دیگر تخم‌ها در لانه‌های میزبان رشد می‌کند و توسط پرندگان میزبان تغذیه می‌شوند. قابل توجه است که فقط یک تخم در هر لانه شانس رشد دارد زیرا فاخته بیشترین مقدار غذا آورده شده توسط پرنده میزبان را می‌خورد.



شکل ۶: نمایی از تخم‌گذاری تصادفی در الگوریتم فاخته

مهاجرت فاخته‌ها در روش پیشنهادی

وقتی فاخته‌ها رشد کردند و بالغ شدند در جمعیت خود زندگی می‌کنند. اما در زمان تخم‌گذاری، آن‌ها به جمعیت جدید و بهتری با مشابهت بیشتر تخم‌ها با پرنده میزبان مهاجرت می‌کنند. بعد از اینکه فاخته‌ها در نواحی مختلفی شکل گرفتند، جمعیت با بهترین میزان تابع هدف به عنوان نقطه هدف برای مهاجرت دیگر فاخته‌ها انتخاب می‌شود. بعد از استقرار، تمیز دادن اینکه فاخته متعلق به کدام گروه است سخت می‌باشد. برای حل این مساله، خوشه‌بندی انجام می‌گیرد. بعد از گروه‌بندی فاخته، ماکزیمم میانگین بهره هر گروه هدف را تعیین می‌کند. همان‌طور که قبلاً ذکر شد، فاخته‌ها سکونت‌گاه خود را برای تخم‌گذاری با انتقال همه فاخته‌ها به نقطه هدف بهبود می‌بخشند. نسخه اصلی COA بر روی مسائل پیوسته عمل می‌کند. چون انتخاب ویژگی مساله ای گسسته است، در این تحقیق یک روش مهاجرت جدید که برای مسائل گسسته مناسب است ارائه شده است. این عملگر به صورت فرمان‌های زیر است.

- برای هر زیستگاه تکرار کن
- فاصله بلوک شهری را بین هر زیستگاه و نقطه هدف را محاسبه کن

- یک رشته باینری (S) اولیه با طول N به تعداد ویژگی‌ها با مقداردهی اولیه به صفر بساز
- تعدادی از سلول‌های باینری را به ۱ مقداردهی کن
- سلول‌های ۱ موجود در نقطه هدف را به همان سلول‌های زیستگاه‌ها در S کپی کن

از بین بردن فاخته در بدترین سکونت‌گاه‌ها

به دلیل تعادل جمعیت در پرندگان، پارامتری جدید تعریف شده است که ماکزیمم تعداد فاخته‌های زنده در جمعیت را محدود می‌کند. برای مدل‌سازی این محدودیت، N_{max} تعداد فاخته‌های زنده‌ای است که بهره بهتری دارند و دیگر فاخته‌ها می‌میرند.

تابع هدف روش پیشنهادی

در این قسمت برای انتخاب یک زیرمجموعه مناسب ویژگی‌ها، از یک معیار استقلال و تفکیک‌پذیری و یک معیار همبستگی مبتنی بر اطلاعات متقابل برای کاهش بعد داده‌ها و کسب دقت بالای دسته بندی کننده، تمرکز نموده شده است. معیار تفکیک پذیری با ماتریس پراکندگی میان دسته ای بر رابطه بین متغیرها بر حسب برچسب‌های دسته‌شان، دلالت می‌کند و هرچه که این معیار بزرگتر باشد و در سوی دیگر افزایش همبستگی ویژگی‌ها با کلاس داده و کاهش افزونگی بین ویژگی‌های حاصل، نشان دهنده این است که متغیرها کمتر تکراری می‌باشند. در بسیاری از رویکردهای انتخاب ویژگی قبلی تنها از یک تابع هدف بر اساس بررسی توزیع داده‌ها پشتیبانی می‌گردد که این امر منجر به حذف بسیاری از راه حل‌های کارآمدی می‌گردد که ممکن است بر اساس اهداف کارایی دیگری دارای قدرت باشند. در واقع ارزیابی این زیرمجموعه‌ها با یک روش انجام می‌شود که منجر به پس زدن زیرمجموعه‌هایی با خواص خوب با توجه به معیار دیگری می‌شود. در این زمینه، راه حل با افزونگی بسیار کم یا همبستگی بسیار بالا که دو هدف کارایی بسیار مهم می‌باشند، می‌تواند توسط فرآیند انتخاب رد شود. برای افزایش تنوع زیرمجموعه‌های انتخاب شده، در این پایان نامه فضای جستجوی بسیار بزرگ مسئله با یک روش چند منظوره مورد بررسی قرار می‌گیرد. برای این کار دو هدف کارایی بر اساس درجه همبستگی مبتنی بر پراکندگی داده‌ها و سپس همبستگی و افزونگی مبتنی بر اطلاعات متقابل تعریف می‌گردد. در اینجا به بررسی دو معیار کارایی مختلف و سپس تابع هدف روش پیشنهادی برای ارزیابی الگوریتم بهینه‌سازی جستجوی فاخته می‌پردازیم.

معیار اول: درجه تفکیک‌پذیری

از آنجایی که نسبت ماتریس پراکندگی بین دسته‌ای به پراکندگی دورن دسته‌ای زیرمجموعه‌ها، می‌تواند مشارکت آن‌ها در دسته‌بندی را ارائه دهد، زیرمجموعه ویژگی بهینه بر حسب میزان استقلال ویژگی‌ها در فرآیند دسته‌بندی انتخاب می‌گردد. هدف درجه تفکیک‌پذیری^۶، انتخاب ویژگی‌های بهینه برای دسته‌بندی می‌باشد. تصور کنید که $(x, y) \in (R_d \times Y)$ یک نمونه است که در اینجا R_d یک فضای ویژگی d بعدی و $Y = \{1, 2, \dots, c\}$ ، مجموعه برچسب دسته می‌باشد. نماد n_i بیانگر تعداد نمونه‌های متعلق به طبقه i ام بوده و N تعداد کلی نمونه‌ها را ارائه می‌دهد. فکر کنید که x_{ij} بر نمونه j ام در دسته i ام دلالت داشته، u میانگین نمونه همه دسته‌ها و u_i میانگین نمونه دسته i ام می‌باشد.

$$u_i = \frac{1}{n_i} \sum_{j=1}^{n_i} X_{ij} \quad \text{رابطه (۲)}$$

$$u = \frac{1}{N} \sum_{i=1}^c \sum_{j=1}^{n_i} X_{ij} \quad \text{رابطه (۳)}$$

پس ماتریس پراکندگی بین دسته‌ای S_b ^۷ و ماتریس پراکندگی میان دسته‌ای S_w ^۸، به صورت زیر تعریف می‌شوند:

$$S_b = c \sum_{i=1}^c n_i (u_i - u) \quad \text{رابطه (۴)}$$

⁶ The Separability Score

⁷ Between-Class Scatter Matrix

⁸ Within-Class Scatter Matrix

$$S_w = \sum_{i=1}^c \sum_{j=1}^{n_i} (u_i - x_{ij}) \quad \text{رابطه (۵)}$$

در اینجا، S_b فواصل میان بردار میانگین هر یک از دسته‌ها و میانگین کلی را اندازه‌گیری کرده در حالی که S_w میانگین پراکندگی دسته‌ها در حول بردارهای میانگین‌شان را اندازه‌گیری می‌کند. برای یک زیرمجموعه ویژگی مورد نظر، تفکیک‌پذیری دسته مبتنی بر ماتریس پراکندگی^۹، ارزیابی‌کننده نسبت ردیابی^{۱۰} یا عامل تعیین‌کننده ماتریس پراکندگی بین دسته‌ای با ماتریس پراکندگی درون دسته‌ای می‌باشد. تفکیک‌پذیری دسته Obj_1 به صورت زیر ارائه شده است:

$$Obj_1 = \frac{S_b}{S_w} \quad \text{رابطه (۶)}$$

یک زیرمجموعه با Obj_1 بزرگ، به عنوان یک زیرمجموعه خوب در نظر گرفته شده و به معنی پراکندگی درون دسته‌ای کوچک و پراکندگی بین دسته‌ای بزرگ می‌باشد. از این رو، یک Obj_1 بزرگ تضمین می‌کند که دسته‌ها توسط میانگین‌های پراکندگی‌شان به خوبی پراکنده شده‌اند. این معیاری ساده، قدرتمند و یکپارچه برای دسته‌بندی می‌باشد. ایده تفکیک‌پذیری دسته به منظور انتخاب زیرمجموعه بهینه برای دسته‌بندی مورد استفاده قرار گرفته، زیرا مشارکت‌های این زیرمجموعه‌ها برای دسته‌بندی را منعکس می‌کند.

معیار دوم: همبستگی مبتنی بر اطلاعات متقابل

برای افزایش تنوع زیرمجموعه‌های انتخاب شده، در این تحقیق فضای جستجوی بسیار بزرگ مسئله با یک روش چند منظوره مورد بررسی قرار می‌گیرد. ارزیابی کیفیت اهداف کارایی بیان شده براساس اطلاعات متقابل^{۱۱} صورت می‌پذیرد تا به طور جداگانه میزان همبستگی داده‌ها و افزونگی زیرمجموعه ویژگی‌های انتخابی را اندازه‌گیری کند. این دو معیار هر دو کیفیت اختصاصی ویژگی‌های انتخاب شده و کیفیت زیرمجموعه را اندازه‌گیری می‌کند.

در این بخش ما فرموله‌سازی معیار دیگری را برای انتخاب ویژگی به عنوان تابع هدف الگوریتم بهینه‌سازی جستجوی فاخته و نحوه محاسبه آن‌ها را ارائه می‌دهیم. اطلاعات متقابل یک شاخص خوب برای مطالعه وابستگی بین یک ویژگی و طبقه‌بندی و افزونگی بین ویژگی‌های تصادفی است. X و Y را به عنوان دو متغیر تصادفی با قوانین احتمال گسسته در نظر بگیرید. اطلاعات متقابل دو متغیر X و Y با $I(X, Y)$ نشان داده شده و از طریق $P(X)$ و $P(Y)$ و $P(X, Y)$ به صورت زیر تعریف می‌شود:

$$I(X; Y) = \sum_{y \in \Omega_Y} \sum_{x \in \Omega_X} P(x, y) \cdot \log \frac{P(x, y)}{P(x) \cdot P(y)} \quad \text{رابطه (۷)}$$

که در آن Ω_X و Ω_Y به ترتیب فضاهای نمونه X و Y به ترتیب هستند. هنگامی که دو متغیر X و Y به یکدیگر وابسته هستند، $I(X, Y)$ زیاد است و در نقطه مقابل هنگامی که X و Y به طور کامل از هم مستقل هستند، $I(X, Y)$ برابر با صفر است. بعد از محاسبه اطلاعات متقابل، ما دو هدف کارایی مختلف را برای انتخاب بهترین زیرمجموعه از بین راه‌حل‌های تولیدی توسط الگوریتم بهینه‌سازی جستجوی فاخته را دنبال می‌کنیم. اولین هدف کیفیت مجموعه ویژگی‌های انتخابی در راستای پیش‌بینی هدف داده‌ها که تحت عنوان کلاس بیان می‌گردد، را ارزیابی می‌کند و دومین هدف به بررسی افزونگی و تکراری بودن مجموعه ویژگی‌ها می‌پردازد. این دو هدف در ادامه بیان شده‌اند.

همبستگی ویژگی‌ها با کلاس

برای هر زیرمجموعه از ویژگی‌ها، ما مفهوم ارتباط را که از طریق وابستگی بیان می‌گردد را تعریف می‌کنیم که در واقع از طریق محاسبه میانگین اطلاعات متقابل بین هر یک از ویژگی‌های موجود در مجموعه ویژگی‌های انتخابی با مجموعه هدف داده‌ها یعنی کلاس داده‌ها که از طریق متغیر C نشان داده می‌شود، به دست می‌آید. نحوه محاسبه این رابطه در زیر نشان داده شده است:

⁹ Scatter-Matrix-based class separability

¹⁰ The Ratio of The Trace

¹¹ Mutual Information (MI)

$$D_S = \frac{1}{|S|} \sum_{X_i \in S} I(X_i; c) \quad \text{رابطه (۸)}$$

که در آن $I(X_i, c)$ بیان‌کننده مقدار اطلاعات متقابل بین ویژگی X_i در زیرمجموعه ویژگی انتخابی کاندید و کلاس داده‌ها می‌باشد و بیان‌کننده این واقعیت است که این ویژگی با چه کیفیتی کلاس داده‌ها را بیان می‌کند.

کاهش افزونگی بین ویژگی‌های انتخابی

در یک راه‌حل انتخابی ممکن است، دو یا بیش از دو ویژگی در هدف اول یعنی ارتباط و همبستگی با کلاس بسیار خوب باشند، اما این ویژگی‌ها به نحوی دارای افزونگی بوده و بتوان با داشتن یکی از آن‌ها به مقادیر دیگر ویژگی‌ها دست یافت. در این شرایط می‌بایست ویژگی‌های انتخابی در راستای کاهش افزونگی و حذف ویژگی‌های تکراری بررسی گردند. در اینجا از اطلاعات متقابل برای ارزیابی افزونگی بین ویژگی‌ها استفاده می‌گردد. در اینجا دو متغیر i و j از یک مجموعه ویژگی‌های کاندید توسط الگوریتم بهینه‌سازی ازدحام ذرات به صورت $(X_i; X_j)_{i, j = 1, \dots, m_i \neq j}$ نشان داده شده و از طریق رابطه زیر محاسبه می‌گردد:

$$R_S = \frac{1}{|S|^2} \sum_{X_i, X_j \in S} I(X_i; X_j) \quad \text{رابطه (۹)}$$

هدف ما افزایش معیار اول یعنی همبستگی ویژگی‌ها با کلاس و کاهش هدف دوم یعنی افزونگی مابین ویژگی‌هاست. در اینجا از یک تابع هدف به صورت رابطه زیر استفاده می‌گردد:

$$Obj_2 = \frac{D_S}{R_S} \quad \text{رابطه (۱۰)}$$

در این تابع هدف همبستگی داده‌ها در صورت کسر و هدف دوم یعنی افزونگی در مخرج کسر قرار می‌گیرد. با بزرگ شدن صورت و کوچک شدن مخرج حاصل کسر بزرگ شده و کیفیت بالاتری را نشان می‌دهد. بنابراین یک تعادل و توازن ما بین این دو هدف کارایی برقرار می‌گردد.

تابع هدف نهایی الگوریتم بهینه‌سازی جستجوی فاخته

در بخش قبل دو هدف کارایی مختلف برای ارزیابی راه‌حل‌های تولید شده توسط هر فاخته در الگوریتم بهینه‌سازی جستجوی فاخته معرفی گردید. در این بخش یک تابع هدف کلی بر اساس این دو تابع هدف کارایی تعریف می‌گردد که از حاصل جمع دو تابع کارایی بیان شده حاصل می‌گردد. این تابع در رابطه زیر نشان داده شده است:

$$Fitness = Obj_1 + Obj_2 \quad \text{رابطه (۱۱)}$$

از آنجایی که در هر یک از این فاکتورها، بالاتر بودن بیانگر کیفیت بالاتر روش پیشنهادی می‌باشد، بنابراین بالابودن تابع هدف پیشنهادی برای هر راه‌حل تولید شده توسط الگوریتم بهینه‌سازی جستجوی فاخته بیانگر کارآمدی بالاتر آن می‌باشد.

ارزیابی روش پیشنهادی

روش پیشنهادی بر اساس پارامترهای متعددی مورد ارزیابی قرار گرفته است که در این فصل جزئیات شبیه‌سازی و پیاده‌سازی روش مورد نظر و همچنین نتایج حاصل شده مورد بررسی و ارزیابی قرار می‌گیرند.

ویژگی‌های محیط پیاده‌سازی

برای پیاده‌سازی روش پیشنهادی از نرم‌افزار MATLAB R2010a استفاده گردیده است. این شبیه‌سازی بر روی سیستمی کامپیوتری با پردازنده corei7 با میزان حافظه اصلی ۴ گیگابایت و بر روی سیستم عامل ویندوز ۷ اجرا شده است.

داده‌های مورد استفاده

داده‌های مربوط به داده‌های حجیم که در این تحقیق مورد استفاده قرار گرفته، مربوط به حمله‌های شبکه‌ها که در ۵ دسته قرار دارند و این مجموعه داده‌ها مربوط به مخزن داده UCI بوده و از آدرس اینترنتی^{۱۲} قابل دسترسی می‌باشد. جدول ۱ این داده‌ها را با مشخصات هر یک که بیانگر تعداد رکورد و ویژگی است، توصیف می‌کند.

جدول ۱: مشخصات مجموعه داده‌های مورد استفاده

تعداد طبقات داده	تعداد ویژگی‌ها	تعداد رکوردها	مجموعه داده
۵	۴۱	۴۹۴۰۲	KddCup99

در اینجا داده‌ها به صورت تصادفی بین دو مجموعه، یعنی ۷۰ درصد نمونه‌ها آموزشی و ۳۰ درصد نمونه‌های آزمایشی، تقسیم شده‌اند. ما الگوریتم فاخته را بر روی مجموعه آموزش، بکار گرفته‌ایم. در انتهای هر دور اجرا، بهترین زیرمجموعه‌ای که کشف گردیده را برای آموزش سه طبقه‌بندی کننده مختلف و ارزیابی عملکردشان مورد استفاده قرار می‌گیرند. پارامترهای بهینه الگوریتم فاخته، به صورت غیر مستدل توسط بکارگیری یک مجموعه از آزمایشات مقدماتی، که در جدول (۱) فهرست شده‌اند انتخاب گردیده و برای همه آزمایشات گزارش شده بکار گرفته شده‌اند. مدل‌های طبقه‌بندی کننده بر روی نمونه‌های آموزش ایجاد گردیده‌اند. عملکردهای الگوریتم‌های آموزشی بر روی نمونه‌های آموزشی مورد ارزیابی قرار گرفته و بر اساس تکرارهای متفاوت تقسیم داده‌ها، میانگین گیری شده‌اند.

ارزیابی کیفیت روش انتخاب ویژگی پیشنهادی

از چندین الگوریتم طبقه‌بندی کننده برای ارزیابی زیرمجموعه ویژگی به دست آمده، مورد استفاده قرار گرفته‌است. عملکرد الگوریتم پیشنهاد شده با طبقه‌بندی کننده‌های مختلف بر روی دو حالت کل ویژگی‌ها، ویژگی‌های انتخابی توسط الگوریتم جستجوی فاخته بر اساس پراکندگی و همبستگی داده‌ها در جدول (۱) نشان داده شده است. بهترین دقت‌های طبقه‌بندی با افزایش تعداد دفعات تکرار، نشان داده شده‌اند. به منظور تایید بیشتر کارآمدی و عمومیت روش‌مان، عملکردهای به دست آمده با عملکردهای به دست آمده از مجموعه کل ویژگی‌های مقایسه شده‌اند. دقت‌ها برای زمانی می‌باشند که همه ویژگی‌های و زیرمجموعه ویژگی بهینه فراهم شده‌اند. در این حالت از ۳ طبقه بندی کننده ماشین بردار پشتیبان، درخت تصمیم و K- نزدیک ترین همسایه استفاده می‌گردد.

مقایسه با حالت عدم انتخاب ویژگی تاثیرگذار

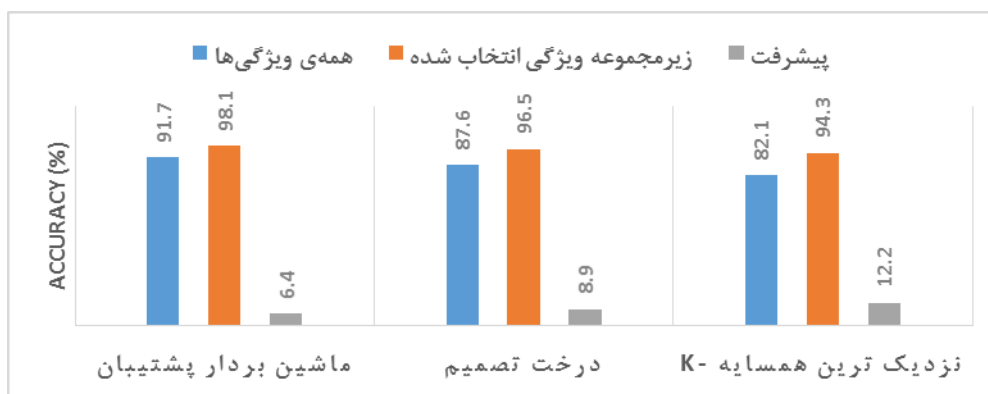
سه طبقه‌بندی کننده‌ای که در بالا مورد بحث و بررسی قرار گرفت برای ارزیابی زیرمجموعه ویژگی به دست آمده، مورد استفاده قرار گرفته‌اند. عملکرد الگوریتم پیشنهاد شده با سه طبقه‌بندی کننده در جدول (۲) نشان داده شده است.

جدول ۲: دقت طبقه‌بندی (درصد) به واسطه استفاده از کل ویژگی‌ها و ویژگی‌های بهینه انتخاب شده

داده‌ها/ طبقه بندی کننده	ماشین بردار پشتیبان	درخت تصمیم	K- نزدیک ترین همسایه
همه‌ی ویژگی‌ها	۹۱٫۷	۸۷٫۶	۸۲٫۱
زیرمجموعه ویژگی انتخاب شده	۹۸٫۱	۹۶٫۵	۹۴٫۳
پیشرفت	۶٫۴	۸٫۹	۱۲٫۲

ذکر این مطلب مناسب است که عملکرد به دست آمده با استفاده از زیرمجموعه ویژگی انتخاب شده در مقایسه با مجموعه کل ویژگی‌ها دارای برتری می‌باشد. به صورت ویژه، ماشین بردار پشتیبان، بالاترین دقت طبقه‌بندی را در مقایسه با درخت تصمیم و K نزدیک‌ترین همسایه دارد. از این رو، ماشین بردار پشتیبان به عنوان طبقه‌بندی کننده انتخاب شده است. ما به این نتیجه رسیدیم که K نزدیک‌ترین همسایه نسبتاً ضعیف تر از درخت تصمیم بر روی داده‌های مورد استفاده می‌باشد. با مقایسه نتایج کسب شده بر اساس مجموعه کل ویژگی‌ها، عملکرد طبقه‌بندی روش پیشنهادی پیشرفت کرده است. نمودار (۱) این مقایسه را به خوبی نشان می‌دهد.

¹² <https://archive.ics.uci.edu/ml/datasets/>



نمودار(۱): نمودار مقایسه دقت دسته‌بندی (درصد) به‌واسطه استفاده از کل ویژگی‌ها و ویژگی‌های بهینه انتخاب شده

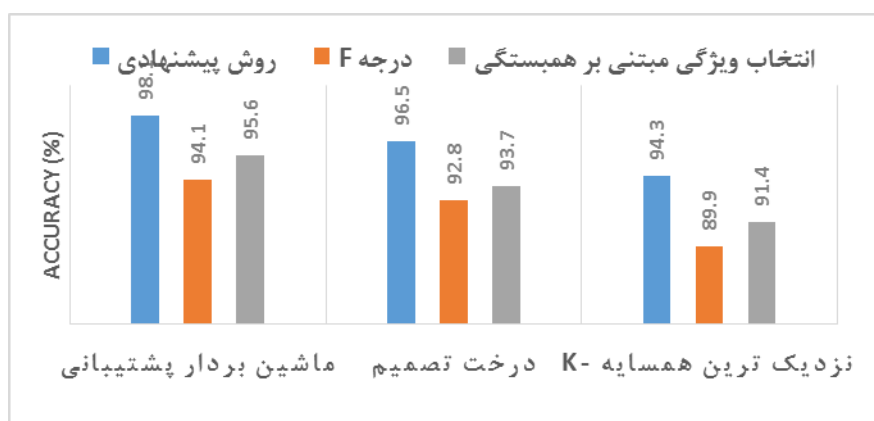
مقایسه با سایر روش‌های انتخاب ویژگی

به منظور تائید کارآمدی روش پیشنهاد شده انتخاب ویژگی مبتنی بر الگوریتم فاخته، نتایج با نتیجه‌های به دست آمده از دو روش انتخاب ویژگی دیگر که قبلاً نیز آن‌ها را در فصل دوم ذکر کرده‌ایم، مورد مقایسه قرار گرفته‌اند. نتایج طبقه‌بندی در جدول (۳) ارائه شده‌اند.

جدول ۳: مقایسه درصد دقت کسب شده توسط روش پیشنهادی و سایر روش‌ها

روش پیشنهادی	ماشین بردار پشتیبانی	درخت تصمیم	k- نزدیک ترین همسایه
روش پیشنهادی	۹۸،۱	۹۶،۵	۹۴،۳
درجه F	۹۴،۱	۹۲،۸	۸۹،۹
انتخاب ویژگی مبتنی بر همبستگی	۹۵،۶	۹۳،۷	۹۱،۴

همان‌طور که ملاحظه می‌کنید روش ما نتایج بهتری را در مقایسه با دیگر روش‌ها ارائه می‌دهد. روش ما با طبقه‌بندی کننده ماشین بردار پشتیبانی، بیشترین پیشرفت عملکرد طبقه‌بندی را در میان این سه روش انتخاب ویژگی دارد. به صورت کلی، روش پیشنهاد شده انتخاب ویژگی بر مبنای الگوریتم فاخته در امتداد ماشین بردار پشتیبانی بسیار کارآمد می‌باشد. نمودار (۲) این مقایسه را نشان می‌دهد.



نمودار(۲): نمودار مقایسه درصد دقت کسب شده توسط روش پیشنهادی و سایر روش‌های انتخاب ویژگی

همان‌طور که ملاحظه می‌کنید روش ما نتایج بهتری را در مقایسه با دیگر روش‌ها ارائه می‌دهد. روش مبتنی بر همبستگی دارای عملکرد نسبتاً بدتری در مقایسه با روش ما و عملکرد نسبتاً بهتری در مقایسه با روش انتخاب ویژگی مبتنی بر F-Score می‌باشد. روش ما با طبقه‌بندی کننده ترکیبی پیشنهادی، بیشترین پیشرفت عملکرد طبقه‌بندی را در میان این سه روش انتخاب ویژگی دارد. به صورت کلی، روش پیشنهادی برای انتخاب ویژگی بر مبنای الگوریتم جستجوی فاخته در امتداد طبقه‌بندی کننده‌های مختلف بسیار کارآمد می‌باشد.

مقایسه با روش های انتخاب ویژگی ارائه شده در تحقیقات قبلی

روش پیشنهادی را با سایر روش‌های که قبلاً در این زمینه کار کرده‌اند نیز مورد مقایسه قرار می‌گیرد. این مقایسه بر این دلالت دارد که روش ما دارای بالاترین دقت تشخیص می‌باشد. مقایسه‌ی موجود در جدول (۴) نشان می‌دهد که دقت روش ما مقداری بالاتر از دقت سایر روش‌ها می‌باشد، همچنین لازم به ذکر است که روش ما دارای پیچیدگی نسبتاً کمتری می‌باشد. بنابراین روش ما به صورت کلی بسیار کارآمد می‌باشد.

جدول ۴: مقایسه نتایج

روش کار	دقت طبقه بندی (Accuracy)
تبرید شبیه سازی شده + ماشین بردار پشتیبان [۳۱]	۹۴,۹
بهینه سازی ازدحام ذرات + ماشین بردار پشتیبان [۳۲]	۹۶,۱
الگوریتم ژنتیک + ماشین بردار پشتیبان [۳۳]	۹۶,۷
روش پیشنهادی	۹۸,۱

نتیجه گیری

انتخاب ویژگی به منظور انتخاب زیرمجموعه مناسب از بین ویژگی‌های استخراج شده، انجام می‌شود. انتخاب ویژگی یک مسئله‌ی جستجو در فضای بزرگی از راه‌حل‌ها شامل ترکیبات متفاوتی از ویژگی‌ها می‌باشد و منجر به بهبود محاسبه کار آیی دسته‌بندی، ایجاد دسته‌بندی‌های سریع و کم‌هزینه می‌گردد. در این تحقیق، یک رویکرد برای کاهش ابعاد ویژگی‌های مجموعه‌های داده‌ای با استفاده از الگوریتم جستجوی فاخته ارائه شده است. از آنجاکه تعداد ویژگی‌ها زیاد است لازم است برای کار آیی بالاتر طبقه‌بندی کننده ویژگی‌های مؤثر انتخاب و بقیه حذف شوند. از الگوریتم فاخته به این منظور استفاده شده است. از نکات مهم، انتخاب تابع هدف، ماکزیمم نمودن تابع هدف و نحوه کدگذاری داده‌ها است. نتایج نشان داد که روش انتخاب ویژگی بر اساس الگوریتم فاخته بر مبنای تابع سازگاری نسبت به سایر روش انتخاب ویژگی‌ها برتری قابل توجهی دارد. همچنین دسته‌بندی استفاده شده بر مبنای داده‌های استخراجی توسط الگوریتم فاخته، نسبت به سایر روش‌های معرفی شده برای طبقه‌بندی مناسب دارای دقت تشخیص بالاتری می‌باشد که بیانگر کیفیت روش معرفی شده از جنبه‌های مختلف می‌باشد.

منابع و مراجع

- [1] J. Han, J. Pei, and M. Kamber, *Data mining: concepts and techniques*. Elsevier, 2011.
- [2] H. Liu and H. Motoda, *Feature selection for knowledge discovery and data mining*, vol. 454. Springer Science & Business Media, 2012.
- [3] Al-Tashi, Q., Abdulkadir, S. J., Rais, H. M., Mirjalili, S., & Alhussian, H. (2019). Binary Optimization Using Hybrid Grey Wolf Optimization for Feature Selection. *IEEE Access*, 1–1. doi:10.1109/access.2019.2906757
- [4] Li, Y. , Li, T. , & Liu, H. (2017). Recent advances in feature selection and its applications. *Knowledge and Information Systems*, 53 (3), 551–577 .
- [5] M. Mafarja et al. (2019), Binary grasshopper optimisation algorithm approaches for feature selection problems, *Expert Systems With Applications* 117 (2019) 267–286, <https://doi.org/10.1016/j.eswa.2018.09.015>
- [6] Ahmed KHARRAT and Mahmoud NEJI, (2019), Feature selection based on hybrid optimization for magnetic resonance imaging brain tumor classification and segmentation, *Applied Medical Informatics Original Research* Vol. 41, No. 1 /2019, pp:9-23.
- [7] G. Cheng, C. Yang, X. Yao, L. Guo, and J. Han, “When deep learning meets metric learning: remote sensing image scene classification via learning discriminative cnns,” *TGRS*, vol. 56, no. 5, pp. 2811–2821, 2018.
- [8] G. Cheng, P. Zhou, and J. Han, “Duplex metric learning for image set classification,” *TIP*, vol. 27, no. 1, pp. 281–292, 2018.
- [9] B. Tang, S. Kay, and H. He, “Toward optimal feature selection in naive bayes for text categorization,” *TKDE*, vol. 28, no. 9, pp. 2508–2521, 2016.
- [10] F.Wu, Z. Yuan, and Y. Huang, “Collaboratively training sentiment classifiers for multiple domains,” *TKDE*, vol. 29, no. 7, pp. 1370–1383 ,2017.
- [11] J. Han, R. Quan, D. Zhang, and F. Nie, “Robust object cosegmentation using background prior,” *TIP*, vol. 27, no. 4, pp.1639–1651, 2018
- [12] D. Zhang, D. Meng, L. Zhao, and J. Han, “Bridging saliency detection to weakly supervised object detection based on selfpaced curriculum learning,” in *IJCAI*, 2016, pp. 3538–3544.
- [13] J. Han, D. Zhang, G. Cheng, N. Liu, and D. Xu, “Advanced deep-learning techniques for salient and category-specific object detection: a survey,” *SPM*, vol. 35, no. 1, pp. 84–100, 2018.
- [14] D. Zhang, J. Han, L. Zhao, and D. Meng, “Leveraging priorknowledge for weakly supervised object detection under a collaborative self-paced curriculum learning framework,” *IJCV*, p. DOI.<https://doi.org/10.1007/s1126>, 2018.
- [15] G. Cheng, J. Han, P. Zhou, and D. Xu, “Learning rotation-invariant and fisher discriminative convolutional neural networks for object detection,” *TIP*, vol. 28, no. 1, pp. 265–278, 2019.
- [16] X. Chang, Y.-L. Yu, Y. Yang, and E. P. Xing, “Semantic pooling for complex event analysis in untrimmed videos,” *TPAMI*, vol. 39, no. 8, pp. 1617–1632, 2017.
- [17] D. Zhang, D. Meng, and J. Han, “Co-saliency detection via a self-paced multiple-instance learning framework,” *TPAMI*, vol. 39, no. 5, pp. 865–878, 2016.
- [18] J. Han, G. Cheng, Z. Li, and D. Zhang, “A unified metric learningbased framework for co-saliency detection,” *TCSVT*, vol. 28, no. 10, pp. 2473–2483, 2018.
- [19] S. Zhang, X. Yu, Y. Sui, S. Zhao, and L. Zhang, “Object tracking with multi-view support vector machines,” *TMM*, vol. 17, no. 3, pp. 265–278, 2015.
- [20] Xu, J., Han, J., Nie, F., & Li, X. (2019). Multi-view Scaling Support Vector Machines for Classification and Feature Selection. *IEEE Transactions on Knowledge and Data Engineering*, 1–1. doi:10.1109/tkde.2019.2904256

- [21] B. S. Harish, M. B. Revanasiddappa, (2017), A Comprehensive Survey on various Feature Selection Methods to Categorize Text Documents, *International Journal of Computer Applications* (0975 – 8887) Volume 164 - No.8, April 2017.
- [22] Al-Tashi et al., (2019), Binary Optimization Using Hybrid Grey Wolf Optimization for Feature Selection, *IEEE*, DOI 10.1109/ACCESS.2019.2906757, *IEEE Access*
- [23] G. Chandrashekar, F. Sahin, (2014), A survey on feature selection methods *Computers and Electrical Engineering* 40 (2014) 16–28.
- [24] Kharrat, A., & Neji, M. (2018). A Hybrid Feature Selection for MRI Brain Tumor Classification. *Innovations in Bio-Inspired Computing and Applications*, 329–338. doi:10.1007/978-3-319-76354-5_30
- [25] Forman George. An extensive empirical study of feature selection metrics for text classification. *Journal of machine learning research*, 3(Mar):1289–1305, 2003.
- [26] Kohavi R, John G (1997) Wrappers for feature subset selection. *Artif Intell* 97(1):273–324.
- [27] Guyon I, Elisseeff A (2003) An introduction to variable and attribute selection. *Machine Learning Research* 3:1157–1182.
- [28] Sayed, G. I., Hassanien, A. E., & Azar, A. T. (2017). Feature selection via a novel chaotic crow search algorithm. *Neural Computing and Applications*. doi:10.1007/s00521-017-2988-6
- [29] Melo, A., & Paulheim, H. (2017). Local and global feature selection for multilabel classification with binary relevance. *Artificial Intelligence Review*. doi:10.1007/s10462-017-9556-4
- [30] I. Mejía-Guevara and Ákuri-Morales, "Evolutionary feature and parameter selection in support vector regression," *MICAI 2007: Advances in Artificial Intelligence*, pp. 399-408, 2007.
- [31] S. Lin, Z. Lee, S. Chen, and T. Tseng, "Parameter determination of support vector machine and feature selection using simulated annealing approach," *Applied soft computing*, vol. 8, pp. 1505-1512, 2008.
- [32] S. Lin, K. Ying, S. Chen, and Z. Lee, "Particle swarm optimization for parameter determination and feature selection of support vector machines," *Expert systems with applications*, vol. 35, pp. 1817-1824, 2014.
- [33] C. Huang and C. Wang, "A GA-based feature selection and parameters optimization for support vector machines," *Expert systems with applications*, vol. 31, pp. 231-240, 2012.