

انتخاب ویژگی در متون فارسی با استفاده از ترکیب الگوریتم های فراشناختی

لیلا هاوشکی^۱، محبوبه شمس^۲، عبدالرضا رسولی کناری^۳

^۱ دانشجوی کارشناسی ارشد نرم افزار، دانشکده برق و کامپیوتر، دانشگاه صنعتی قم،
haveshtki.l@qut.ac.ir

^۲ استادیار، دانشکده برق و کامپیوتر، دانشگاه صنعتی قم،
shamsi@qut.ac.ir

^۳ استادیار، دانشکده برق و کامپیوتر، دانشگاه صنعتی قم،
rasouli@qut.ac.ir

چکیده

انتخاب ویژگی به طور گسترده در زمینه های متن کاوی برای ایجاد یک مدل با تعداد ویژگی های کمتر استفاده می شود. انتخاب ویژگی متن یک گام مهم در طبقه بندی متن است و به طور مستقیم بر عملکرد آن تأثیر می گذارد. در این مقاله یک روش بهبود انتخاب ویژگی برای طبقه بندی داده های بزرگ با استفاده از ترکیب الگوریتم های فراشناختی^۱ پیشنهاد می شود. در این روش از ترکیب الگوریتم بهینه سازی ازدحام ذره ها بر اساس لینک و الگوریتم جست و جوی گرانشی^۲ استفاده می شود. الگوریتم بهینه سازی ازدحام ذره ها بر اساس لینک، بهبودی از الگوریتم بهینه سازی ازدحام ذره ها^۳ است. در این پژوهش برای افزایش توان جست و جوی محلی این الگوریتم از الگوریتم جست و جوی گرانشی استفاده می شود. بر اساس ویژگی های انتخاب شده، مدل طبقه بندی کلاس بند نزدیک ترین همسایه^۴ ساخته می شود. در آخر نتایج بر اساس معیار ارزیابی مورد بررسی قرار می گیرد.

عملکرد الگوریتم پیشنهادی بر روی مجموعه داده همشهری مورد بررسی قرار گرفته است. این مجموعه داده با استفاده از کتابخانه هضم پیش پردازش شده است و دو مجموعه داده تصادفی ایجاد شده است. نتایج ما نشان می دهد که روش پیشنهادی با انتخاب تعداد ویژگی کمتر در دو مجموعه داده تولید شده به ترتیب به افزایش دقت ۳/۹۷٪ و ۱/۵۸٪ رسیده است.

کلمات کلیدی

انتخاب ویژگی، داده های بزرگ، الگوریتم جست و جوی گرانشی، الگوریتم بهینه سازی ازدحام ذره ها، متن کاوی، پیش پردازش متن.

۱- مقدمه

ماشین به طور موثر عمل کنند، اعمال پیش پردازش روی داده ها ضروری به نظر می رسد. یکی از تکنیک های پر کاربرد در پیش پردازش داده ها در یادگیری ماشین، انتخاب ویژگی است. انتخاب ویژگی فرآیندی جهت تشخیص ویژگی های مناسب و حذف ویژگی های نامناسب، تکراری یا داده های دارای اختلال است [۳]. چارچوب یک فرآیند انتخاب ویژگی معمولی شامل سه مرحله ای اساسی تولید کننده، تابع ارزیابی و شرط خاتمه است. در مرحله اول زیر مجموعه ای

در طی سال های اخیر حجم داده هایی با ابعاد بالا افزایش قابل توجهی پیدا کرده است. بنابراین روش های یادگیری ماشین در مواجهه با این حجم داده ها با تعداد ویژگی های زیاد دچار مشکل شده است. به همین جهت انتخاب ویژگی به یکی از زمینه های تحقیقی مهم و چالش برانگیز در یادگیری ماشین تبدیل شده است. برای اینکه روش های یادگیری

در پژوهشی یک سیستم برای طبقه بندی متون فارسی ارائه شده است. در این سیستم پس از پیش پردازش متن ها و استخراج ویژگی، برای کاهش ابعاد بردار ویژگی یک روش بهبود یافته انتخاب ویژگی مبتنی بر الگوریتم ازدحام ذره ها استفاده شده است. در نهایت روش های طبقه بندی بر بردار ویژگی کاهش داده شده، اعمال شده است. برای ارزیابی روش انتخاب ویژگی، طبقه بندی کننده های ماشین بردار پشتیبان و بیزین ساده به کار گرفته شده است [۱]. پژوهشی یک روش ترکیبی دو مرحله ای در قسمت انتخاب ویژگی با الگوریتم های بهره اطلاعات و همبستگی براساس زیرمجموعه انتخاب ویژگی^{۱۵} پیشنهاد داده است. در روش بهره اطلاعات با داشتن یک مجموعه آموزشی می توان بهره اطلاعات را برای هر واژه محاسبه نمود. حملاتی که بهره اطلاعاتی آن ها از یک حد آستانه کمتر است از فضای ویژگی حذف می شوند. سپس در روش همبستگی براساس زیرمجموعه انتخاب ویژگی، زیرمجموعه ویژگی ها براساس قابلیت پیش بینی به همراه درجه افزونگی بین آن ها ارزیابی می شود. در نتیجه ویژگی هایی که همبستگی بالایی با متغیر هدف داشته باشند و در عین حال همبستگی بین خود ویژگی ها پایین باشد ترجیح داده می شوند. در آخر از الگوریتم های ترکیبی ماشین برای بررسی نتایج استفاده شده است [۲]. در تحقیقی الگوریتمی پیشنهاد شده است که از یک استراتژی جست و جوی محلی استفاده می کند. این روش به گونه ای عمل می کند که ویژگی های غیر متمایز با احتمال بالا از ویژگی های مرتبط تر انتخاب می شوند. در این روش تعداد قابل توجهی از ویژگی های برجسته را با استفاده از یک طرح تعیین اندازه زیر مجموعه انتخاب می کند. این طرح بر روی یک منطقه محدود کار می کند و تلاش می کند تا یک زیرمجموعه با تعداد کم از ویژگی ها را ارائه دهد [۸]. دو مدل ترکیبی برای طراحی تکنیک های انتخاب ویژگی مختلف براساس الگوریتم بهینه سازی نهنگ^۷ پیشنهاد شده است. در مدل اول، الگوریتم تبرید شبیه سازی شده^۷ الگوریتم بهینه سازی نهنگ تعبیه شده است. برای بهبود بهترین راه حل بعد از هر تکرار، الگوریتم بهینه سازی نهنگ در مدل دوم استفاده می شود. در واقع با استفاده از الگوریتم تبرید شبیه سازی شده، افزایش بهره برداری از طریق جست و جوی مناطق امیدبخش در الگوریتم بهینه سازی نهنگ اتفاق می افتد [۷]. پژوهشی با استفاده از انتخاب ویژگی به تجزیه و تحلیل احساسات پرداخته است. در این پژوهش یک روش انتخاب ویژگی بر اساس فیلتر برای بخش حریصانه ای الگوریتم IG ساخته شده است. نتایج نشان می دهد که بهبود در انتخاب ویژگی باعث شده است که دقت طبقه بندی درصد قابل قبولی افزایش داشته باشد [۹]. پژوهشی الگوریتم انتخاب ویژگی به نام MA-HS را براساس دو الگوریتم Mayfly و Harmon Search ارائه کرده است. این دو الگوریتم های فراشناختی هستند که برای مسئله های بهینه سازی استفاده می شوند. و نتایج نشان می دهد ترکیب

کандید انتخاب می شود. در مرحله دوم این زیر مجموعه مورد ارزیابی قرار می گیرد. در مرحله سوم بهترین معیار قبلی مطابق با یک معیار ارزیابی خاص، مقایسه می شود. اگر زیر مجموعه تازه تاسیس شده از زیر مجموعه قبلی بهتر باشد، آنگاه آن آخرین زیر مجموعه خواهد بود. دو مرحله اول انتخاب ویژگی مبتنی بر جست و جو تا زمانی تکرار می شوند که یک معیار توقف معین ارضا شود. با توجه به معیار ارزیابی، الگوریتم های انتخاب ویژگی در مدل های فیلتر^{۱۰} لافاه^{۱۱} ترکیبی یا تعبیه شده^{۱۲} طبقه بندی می شوند [۴]. در روش فیلتر، ویژگی های مربوطه بدون استفاده از الگوریتم یادگیری تعریف می شود. از سوی دیگر، روش لافاه از روش یادگیری برای انتخاب ویژگی اطلاعاتی استفاده می کند. به طور کلی روش لافاه از روش فیلتر از نظر دقت طبقه بندی بهتر است اما روش فیلتر سریع تر و ساده تر است. روش ترکیبی از ادغام استفاده می کند. یعنی انتخاب ویژگی شامل ساختاری از مدل است [۵]. تکنیک های سنتی بسیاری برای انتخاب ویژگی وجود دارند. بهبود این تکنیک ها مشکل است زیرا آن ها فاقد مدل های ریاضی هستند [۶].

در دهه های گذشته الگوریتم های فراشناختی به هنگام حل مسائل بهینه سازی متنوع مانند طراحی مهندسی، یادگیری ماشین، برنامه ریزی استخراج داده ها، مشکلات تولید و ... بسیار قابل اطمینان بوده اند. انتخاب ویژگی یکی از حوزه هایی است که در آن الگوریتم های فراشناختی مورد استفاده قرار گرفته اند. الگوریتم های فراشناختی عملکرد برتری را در مقایسه با جست و جوی دقیق و کامل نشان می دهند زیرا جست و جوی کامل منجر به تولید همه راه حل های ممکن برای این مشکل می شود. به عنوان مثال، اگر مجموعه داده ها شامل N ویژگی باشد، آنگاه 2^N راه حل باید تولید و ارزیابی شود که به هزینه ای محاسبه بالایی نیاز دارد. جست و جوی تصادفی یک استراتژی دیگر برای انتخاب ویژگی است که گام بعدی را به طور تصادفی جست و جو می کند. الگوریتم های فراشناختی از جست و جوی تصادفی نیز بهتر در نظر گرفته می شوند، زیرا ممکن است به عنوان یک جست و جوی کامل در بدترین حالت عمل کنند. با وجود آنکه الگوریتم های فراشناختی، استراتژی پیدا کردن بهترین راه حل را در هر اجرا تضمین نمی کنند اما ممکن است راه حل قابل قبولی را در زمان معقول تعیین کنند [۷]. در سال های اخیر بسیاری از پژوهشگران در زمینه بهینه سازی انتخاب ویژگی کار کرده اند. با توجه به تحقیق های فراوان، ترکیب الگوریتم های فراشناختی کارایی بالاتری را در حل بسیاری از مسائل عملی نشان داده اند. تاکنون اکثر پژوهش های انجام شده بر روی سندهایی به زبان انگلیسی بوده است. در ادامه به کارهای مرتبط پیشین برای انتخاب ویژگی بر روی زبان های فارسی و انگلیسی پرداخته ایم.

- در مرحله ی سوم بر اساس ویژگی های انتخاب شده مدل طبقه بندی کلاس بند نزدیک ترین همسایه ساخته خواهد شد و مورد ارزیابی قرار خواهد گرفت.

۲-۱-۱- پردازش متن

برای استخراج ویژگی از متن که داده ای غیرساخت یافته است ابتدا باید عملیات پیش پردازش متن انجام شود. عملیات شامل مراحل نرمال سازی^۱، جداسازی و حذف کلمه های توقف^۲ است. در نرمال سازی همه ی نویسه های متن با معادل استاندارد آن جایگزین می شوند. در این مرحله مطابق با یک سری قاعده دقیق و مشخص، فاصله ها و نیم فاصله های موجود در متن برای علامت هایی نظیر "ها" و "ی" غیرچسبان در انتهای کلمه ها و همچنین پیشوندها و پسوندهای فعل ساز نظیر "می"، "ام"، "ایم"، "اید" و موارد مشابه اصلاح می شوند. در جداسازی، تمام متن به صورت کلمه های جداگانه در می آیند. در مرحله ی حذف کلمه های توقف، کلمه هایی مانند: از، با، در، به، برای، اگر و ... که با وجود تکرار بسیار و حضور در اغلب سندها فاقد اطلاعات معنایی هستند حذف می شوند [۱۳].

۲-۱-۲- الگوریتم ازدحام ذره ها بر اساس لینک

الگوریتم ازدحام ذره ها در ابتدا برای مسائل بهینه سازی پیوسته طراحی شد اما پس از گذشت زمان با استفاده از تابع سیگموئید آن را به خوبی برای حل کردن مسائل بهینه سازی گسسته به الگوریتم دودویی ازدحام ذره ها تبدیل کردند. همانند دیگر الگوریتم های فراشناختی، الگوریتم ازدحام ذره ها از بعضی مشکل های ذاتی مانند همگرایی قبل از موقع و گیر انداختن خود در حداقل های محلی رنج می برد. در چنین شرایطی جست و جوی همسایه های شخصی برای پیدا کردن بهترین راه حل مفید است. در این الگوریتم رتبه هر ذره از طریق ماتریس ارتباط و همسایه هر ذره بر اساس رتبه ها محاسبه می شود. درواقع به جای یادگیری از بهترین مکان سراسری یافت شده از بهترین مکان همسایه یافت شده استفاده می شود. سرعت و موقعیت الگوریتم بهینه سازی ازدحام ذرات بر اساس لینک به ترتیب با توجه به معادله های (۱) و (۲) به روز می شود [۵].

$$v_i^j(t+1) = w * v_i^j(t) + rand_1 * c_1 * (pbest_i^j(t) - x_i^j(t)) + rand_2 * c_2 * (nbest_i^j(t) - x_i^j(t)) \quad (1)$$

$$x_i^j(t+1) = x_i^j(t) + v_i^j(t+1) \quad (2)$$

$v(t)$: سرعت فعلی ذره i

$x(i)$: مکان فعلی ذره i

الگوریتم های فراشناختی باعث بهبود انتخاب ویژگی می شوند [۱۰]. پژوهشی یک روش انتخاب جدید ویژگی برای خوشه بندی متن، به نام بهینه سازی ازدحام ذره ها بر اساس لینک را پیشنهاد داده است. این روش یک استراتژی انتخاب همسایه جدید در الگوریتم بهینه سازی ازدحام ذره ها را برای انتخاب ویژگی های برجسته معرفی می کند. عملکرد بهینه سازی ذره ها در الگوریتم ازدحام ذره ها با استفاده از بهترین مکانیسم به روز رسانی سری برای انتخاب ویژگی محدود شده است. در روش پیشنهادی به جای استفاده از بهترین سراسری، ذره ها بر اساس بهترین موقعیت همسایگی برای افزایش قابلیت بهره برداری و اکتشاف به روز می شوند. این ویژگی های برجسته سپس با استفاده از الگوریتم خوشه بندی کی - میانگین برای بهبود عملکرد و کاهش هزینه زمان محاسبه الگوریتم پیشنهاد شده، مورد آزمایش قرار گرفته است [۵].

پژوهشگران در اکثر کارهای پیشین به این موضوع توجه داشته اند که الگوریتم های فراشناختی را به صورتی ترکیب کنند تا دقت در انتخاب ویژگی را افزایش دهند. گرچه این الگوریتم ها در جست و جوی سراسری موفق عمل می کنند اما اغلب در جست و جوی محلی دچار مشکل می شوند. ما در این پژوهش با استفاده از الگوریتم گرانشی، جست و جوی محلی را در الگوریتم بهینه سازی ازدحام ذره ها بر اساس لینک بهبود بخشیده ایم.

۲- مطالب اصلی

یکی از مسائل مهم در متن کاوی، طبقه بندی متن ها است. طبقه بندی متن، انتساب سندهای متنی بر اساس محتوا به یک یا چند طبقه از قبل تعیین شده است. در این پژوهش نیز روشی بر اساس داده کاوی و الگوریتم های انتخاب ویژگی برای طبقه بندی متن ها ارائه شده است.

۲-۱- مفاهیم اولیه مورد نیاز

روش پیشنهادی در این پژوهش در سه مرحله انجام می شود:

- در مرحله ی اول پردازش متن صورت می گیرد. در این مرحله با استفاده از ابزارهای مناسب و بر اساس نیاز، متن جهت سایر پردازش ها و تعیین دسته آماده می شود.
- در مرحله ی دوم پس از پردازش متن و استخراج ویژگی از متن، با توجه به اینکه ویژگی های استخراج شده همگی جهت ساخت مدل مناسب نیستند، انتخاب ویژگی به کمک ترکیب دو الگوریتم بهینه سازی ازدحام ذره ها بر اساس لینک و جست و جوی گرانش صورت خواهد گرفت.

که در آن a ضریب نزولی، G ثابت گرانش اولیه، $iter$ تعداد تکرار فعلی و $max iter$ نشان دهنده ی حداکثر تعداد تکرار است.

$$F_{ij}^d(t) = G(t) \frac{(M_{pi}(t) \times M_{aj}(t))}{(R_{ij}(t) + \varepsilon)} \times (x_j^d(t) - x_i^d(t)) \quad (7)$$

که در آن M_{aj} جرم گرانشی مربوط به عامل j -ام، M_{pi} جرم گرانشی غیرفعال مربوط به عامل i -ام، ε یک ثابت کوچک و R_{ij} فاصله ی اقلیدسی بین دو عامل i و j است. جرم تمام عامل ها با استفاده از رابطه (۸) به روز رسانی می شود.

$$m_i(t) = \frac{fit_i(t) - worst(t)}{best(t) - worst(t)} \quad (8)$$

که در آن، $fit_i(t)$ ارزش شایستگی عامل i -ام در زمان t ، $best(t)$ قوی ترین عامل در زمان t و $worst(t)$ ضعیف ترین عامل در زمان t است. در نهایت الگوریتم جست و جوی گرانش با رسیدن به معیار نهایی متوقف می شود [۱۱].

۲-۱-۴- الگوریتم کلاس بند نزدیک ترین همسایه

الگوریتم کلاس بند نزدیک ترین همسایه یک الگوریتم دسته بندی مبتنی بر نمونه، با ناظر است. اساس کار این الگوریتم مقایسه ی میزان شباهت نمونه ی جدید با نمونه های موجود در داده آموزشی است. جهت اینکه الگوریتم، کلاس نمونه ی جدید را تشخیص دهد، K نزدیک ترین عضو مجموعه ی آموزشی نسبت به نمونه جدید را انتخاب می کند. کلاسی که دارای بیشترین عضو در این K عضو باشد به نمونه ی جدید تعلق می گیرد [۱۲].

۲-۲- طرح مسأله

یکی از مسأله های موجود در انتخاب ویژگی به وسیله ی الگوریتم های فراشناختی، همگرایی زودرس و گیر افتادن در نقاط محلی است. این مسأله باعث می شود که با انتخاب تعداد ویژگی های زیاد به خصوص در داده های بزرگ، زمان پردازش افزایش یابد. با بهبود عملکرد در نقاط محلی یک الگوریتم فراشناختی که در پژوهش های پیشین نتایج قابل قبولی داشته است می توان به این هدف رسید که با انتخاب تعداد ویژگی کمتر به دقت بالاتری دست یافت. همچنین از الگوریتم های ترکیبی فراشناختی برای انتخاب ویژگی در متون فارسی کم استفاده شده است. با پیش پردازش سندها در زبان فارسی و با استفاده از الگوریتم کلاس بند نزدیک ترین همسایه می توان انتخاب ویژگی در زبان فارسی را بدون به خطر انداختن دقت بهبود بخشید.

rand: عدد تصادفی بین ۰ و ۱

c: ضریب شتاب

pbest: بهترین مکان شخصی یافت شده توسط ذره i

nbest: بهترین مکان همسایه یافت شده توسط ذره i

۲-۱-۳- الگوریتم جست و جوی گرانش

هدف این الگوریتم پیدا کردن بهترین راه حل در فضای جست و جوی مسأله با استفاده از نظریه نیوتن است. این نظریه بیان می کند: «هر ذره در جهان، ذره دیگر را با یک نیرویی که به طور مستقیم با حاصل ضرب جرم آن دو ذره و به طور معکوس با مجذور فاصله ی بین آن ها رابطه دارد، جذب می کند.» الگوریتم جست و جوی گرانش را می توان به عنوان مجموعه ای از راه حل های نامزد در نظر گرفت که دارای جرم متناسب با مقدار خود در تابع شایستگی هستند. در طول اجرای الگوریتم تمام توده ها توسط نیروی جاذبه بین خودشان یکدیگر را جذب می کنند. هر چه جرم سنگین تر باشد، نیروی جاذبه بیشتر است. بنابراین، جرم های سنگین تر که به احتمال زیاد نزدیک ترین جرم ها به کمینه کلی هستند، جرم های دیگر را متناسب با فاصله آن ها جذب می کنند.

در الگوریتم جست و جوی گرانش هر عامل یک راه حل نامزد تلقی می شود. در این الگوریتم تمام عوامل با مقادیر تصادفی مقداردهی اولیه می شوند. پس از مقداردهی اولیه، سرعت و موقعیت همه عوامل به ترتیب با استفاده از معادله های (۳) و (۴) محاسبه می شود.

$$v_i^d(t+1) = rand_i \times v_i^d(t) + a_i^d(t) \quad (3)$$

$$x_i^d(t+1) = x_i^d(t) + v_i^d(t+1) \quad (4)$$

که در معادله ها d نشان از بعد مسأله، t نشان از یک زمان خاص، $rand$ عدد تصادفی بین ۰ و ۱ است. a نیز نشان دهنده ی شتاب است. با توجه به قانون حرکت، شتاب یک عامل به طور مستقیم با نیروی برآیند و به طور معکوس با جرم آن رابطه دارد. شتاب با توجه به معادله (۵) محاسبه می شود.

$$a_i^d(t) = \frac{(F_i^d(t))}{M_{ii}(t)} \quad (5)$$

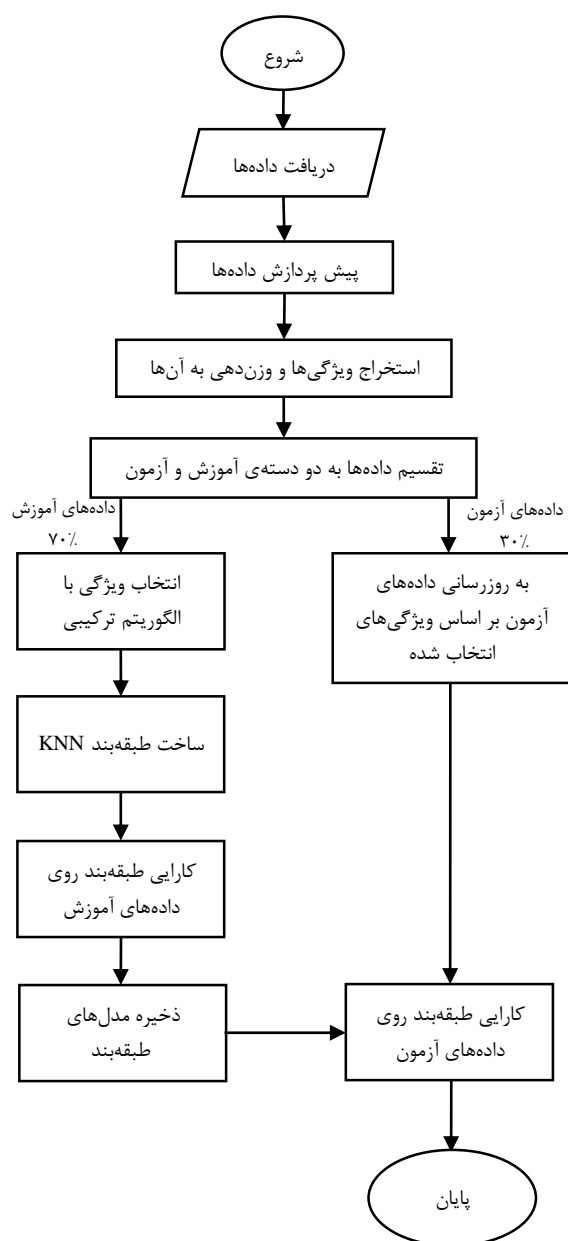
که در آن F_i^d نیروی گرانش و M_{ii} جرم جسم i -ام است. پارامترهای دیگر مانند ثابت گرانش و نیروی گرانش از روابط (۶) و (۷) محاسبه می شوند.

$$G(t) = G_0 \times \exp(-a \times iter / mat iter) \quad (6)$$

هر ذره استفاده خواهد شد. بر همین اساس رابطه (۱۱) نحوه محاسبه کارایی ذره ها در روش پیشنهادی را نشان می دهد.

$$Error = 1 - \frac{Correct\ Detect\ Sample}{Total\ Sample} \quad (11)$$

که $Correct\ Detect\ Sample$ بیان گر تعداد نمونه هایی است که به درستی تشخیص داده شده اند. $Total\ Sample$ تعداد کل نمونه ها است و $Error$ میزان خطا را نشان می دهد. هدف از انتخاب ویژگی، کاهش مقدار خطا در رابطه (۱۱) است. روش پیشنهادی در شکل شماره (۱) نشان داده شده است.



شکل (۱): مراحل روش پیشنهادی

۲-۲- راه حل پیشنهادی

ایده اصلی ترکیب الگوریتم های ازدحام ذره ها و جست و جوی گرانشی در واقع توان جست و جوی سراسری ازدحام ذره ها و جست و جوی محلی الگوریتم گرانشی است. در همه ی الگوریتم های مبتنی بر جمعیت دو ویژگی توانایی الگوریتم برای جست و جوی تمام بخش ها و توانایی بهره برداری از بهترین راه حل در نظر گرفته می شود. یک الگوریتم مبتنی بر جمعیت باید این دو ویژگی را برای تضمین پیدا کردن بهترین راه حل داشته باشد. الگوریتم ازدحام ذره ها توانایی بالایی را در جست و جوی سراسری در تمام بخش ها دارد. در الگوریتم ازدحام ذره ها بر اساس لینک که پیش تر توضیح داده شد، با انتخاب بهترین همسایه، بهبودی در توانایی بهره برداری از بهترین راه حل ایجاد شد. اما الگوریتم همچنان برای توانایی بهره برداری از بهترین راه حل، دچار گیر افتادن در نقاط محلی می شود. توانایی بالای الگوریتم گرانشی در جست و جوی محلی و بهره برداری از بهترین راه حل، باعث بهبود الگوریتم ازدحام ذره ها بر اساس لینک شده است. در نتیجه، این الگوریتم ترکیبی به دقت بالایی در انتخاب ویژگی رسیده است.

۲-۲-۱- ترکیب الگوریتم بهینه سازی ازدحام ذره ها

بر اساس لینک و الگوریتم جست و جوی گرانشی

برای ترکیب این دو الگوریتم رابطه ی (۹) پیشنهاد شده است.

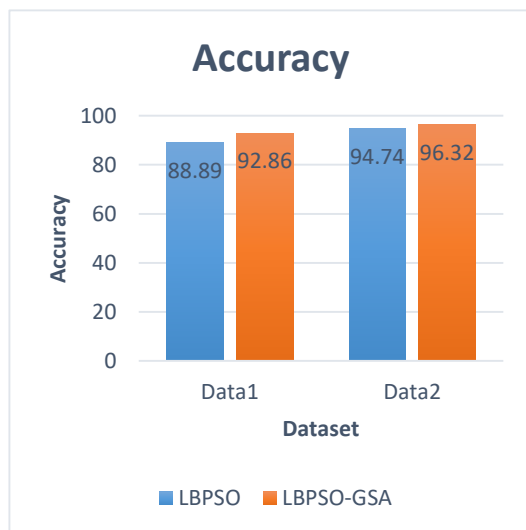
$$v_i(t+1) = w \times v_i(t) + c_1' \times rand \times ac_i(t) + c_2' \times rand \times (nbest - x_i(t)) \quad (9)$$

که در آن، c_1' ضریب شتاب، w تابع وزن و $ac_i(t)$ شتاب عامل i در تکرار t -ام است. موقعیت عوامل نیز در هر تکرار با توجه به رابطه (۱۰) به روز رسانی می شود.

$$x_i(t+1) = x_i(t) + v_i(t+1) \quad (10)$$

در این روش هر عامل به عنوان یک راه حل نامزد در نظر گرفته می شود. در ابتدا همه عوامل به صورت تصادفی مقداردهی اولیه می شوند. پس از مقداردهی اولیه، ثابت گرانش و نیروی گرانش به ترتیب با استفاده از رابطه های (۶) و (۷) محاسبه می شوند. سپس شتاب ذره ها به وسیله رابطه (۵) تعریف می شود. بهترین راه حل پیدا شده در هر تکرار باید به روز شود. پس از محاسبه شتاب و به روز رسانی بهترین راه حل، سرعت همه عامل ها با استفاده از رابطه (۹) محاسبه می شود. در نهایت، موقعیت عامل ها توسط رابطه (۱۰) به روز می شود. پس از انتخاب ویژگی، الگوریتم کلاس بند نزدیک ترین همسایه، برای تخمین دسته های داده ها استفاده شده است. انتخاب ویژگی در روش پیشنهادی روشی مبتنی بر لفافه است. در این روش از کارایی مدل طبقه بند برای محاسبه کارایی

همان طور که در جدول شماره (۱) آمده است تعداد ویژگی های هر مجموعه داده بسیار بالا است. Data1 دارای ۴۷۷۹ ویژگی است که با توجه به جدول شماره (۲) با انتخاب ۱۰۰۰ ویژگی توانسته ایم به دقت ۹۲/۸۶٪ برسیم اما همانطور که در شکل (۳) مشخص است روش ازدحام ذره ها بر اساس لینک (LBPSO)، تعداد ۲۳۰۰ ویژگی را انتخاب کرده است که با توجه به شکل (۲) به دقت ۸۸/۸۹٪ رسیده است. همین طور Data2 دارای ۴۵۱۲ ویژگی است که با توجه به جدول شماره (۲) ما در روش پیشنهادی (LBPSO-GSA) با انتخاب ۱۰۰ ویژگی توانسته ایم به دقت ۹۶/۳۲٪ برسیم. اما همانطور که در شکل (۳) مشخص است روش ازدحام ذره ها بر اساس لینک، تعداد ۵۰۰ ویژگی را انتخاب کرده است که با توجه به شکل (۲) به دقت ۹۴/۷۴٪ رسیده است. پس در این پژوهش ما با انتخاب ویژگی های کم تر و تعداد تکرار کم تر به دقت بالاتری رسیده ایم.



شکل (۲): مقایسه ی عملکرد دو الگوریتم LBPSO و LBPSO-GSA با استفاده از اندازه گیری دقت

شکل شماره (۳) تعداد ویژگی انتخاب شده برای هر داده را در دو روش نشان می دهد. تعداد ویژگی انتخاب شده در روش پیشنهادی به مراتب کمتر از روش LBPSO است. این موضوع نشان می دهد که روش پیشنهادی توانسته است ویژگی های برجسته تری را انتخاب کند.

برای شبیه سازی الگوریتم پیشنهادی از برنامه Matlab2020 استفاده شده است. برای ساخت مجموعه داده و استخراج اطلاعات از این سندها ابتدا عمل پیش پردازش بر روی سندها انجام شده است. عملیات پیش پردازش با استفاده از کتابخانه هضم [۱۴] در python و در محیط colab صورت گرفته است. پس از انجام پیش پردازش، خروجی در محیط متلب به صورت یک دیکشنری در آمده است. سپس عملیات TF-IDF روی آن انجام شده است. ما برای پیش پردازش متون فارسی از مجموعه داده همشهری استفاده کرده ایم. مجموعه داده همشهری مربوط به متن اخبار و گروه اخبارهای منتشر شده در روزنامه همشهری مربوط به سال های ۱۳۷۵-۱۳۸۷ است. مجموعه داده همشهری شامل ۱۰۰۰۰ سند در ده گروهی خبری به فرمت اکسل است. دو مجموعه داده از داده همشهری به صورت تصادفی با کلاس های متفاوت ایجاد کرده ایم. مجموعه داده اول با نام Data1، مجموعه ای با ۸۰ سند شامل کلاس های علم و فرهنگ، سیاسی و اتفاقات متفرقه است. مجموعه داده دوم با نام Data2، مجموعه ای با ۱۲۱ سند شامل کلاس های اقتصادی و سیاسی است. مشخصه های دو مجموعه داده در جدول شماره (۱) آورده شده است.

جدول شماره (۱): مشخصه های مجموعه داده

مجموعه داده	تعداد کلاس	تعداد سند	تعداد ویژگی ها
Data1	۳	۸۰	۴۷۷۹
Data2	۲	۱۲۱	۴۵۱۲

برای بررسی عملکرد الگوریتم پیشنهادی از معیار دقت برای طبقه بندی و میزان تعداد ویژگی انتخاب شده برای رسیدن به دقت بالا استفاده کرده ایم. دقت درواقع محاسبه ی درصد برچسب ها برای تعیین این که مدل تا چه اندازه خروجی را درست پیش بینی می کند. هرچقدر دقت بالاتر باشد ویژگی های مهم تری انتخاب شده اند. نتایج به دست آمده در جدول شماره (۲) آورده شده است.

جدول شماره (۲): مقدار پارامترها و نتایج

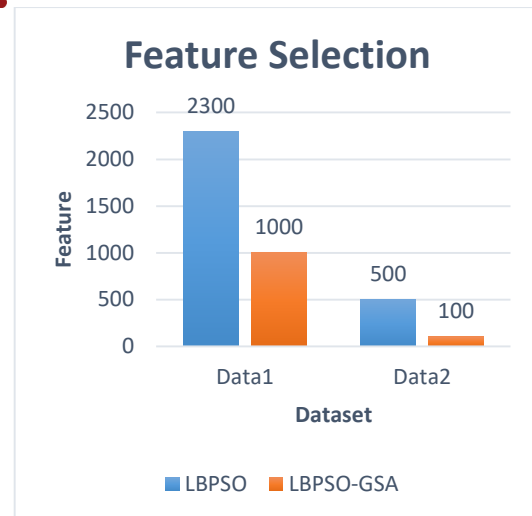
مجموعه داده	تعداد ویژگی انتخاب شده	تعداد تکرار	دقت
Data1	۱۰۰۰	۱۰	۹۲/۸۶
Data2	۱۰۰	۱۰	۹۶/۳۲

شکل شماره (۲) میزان دقت را در روش پیشنهادی و روش ازدحام ذره ها براساس لینک نشان می دهد. این نکته حائز اهمیت است که تعداد تکرار در روش پیشنهادی ۱۰ است و بسیار کمتر از تعداد تکرار ۱۰۰ در روش ازدحام ذره ها براساس لینک است. روش پیشنهادی در شکل های (۲) و (۳) با نام اختصاری LBPSO-GSA و روش ازدحام ذره ها بر اساس لینک با نام LBPSO نشان داده شده است.

مراجع

- [۱] رشدی، اکرم، اکبریور، شاهین، "روش ترکیبی انتخاب ویژگی برای متن کاوی فارسی مبتنی بر الگوریتم های تکاملی"، همایش ملی مهندسی رایانه و مدیریت فناوری اطلاعات، CSITM01، شرکت علم و طلوع فرزین، ۱۳۹۳.
- [۲] هاشمی، سید محسن، "بهبود دسته بندی متون فارسی با ترکیب روش دو مرحله ای انتخاب ویژگی و الگوریتم های یادگیری ماشین"، کنفرانس بین المللی یافته های نوین پژوهشی در مهندسی برق و کامپیوتر، COMCONF01، ۹، موسسه آموزش عالی نیکان، ۱۳۹۴.
- [3] Kumar, V. and S. Minz, *Feature selection: a literature review*. SmartCR, 2014. 4(3): p. 211-229.
- [4] Li, Y., T. Li, and H. Liu, *Recent advances in feature selection and its applications*. Knowledge and Information Systems, 2017. 53(3): p. 551-577.
- [5] Kushwaha, N. and M. Pant, *Link based BPSO for feature selection in big data text clustering*. Future Generation Computer Systems, 2018. 82: p. 190-199.
- [6] Abualigah, L.M. and A.T. Khader, *Unsupervised text feature selection technique based on hybrid particle swarm optimization algorithm with genetic operators for the text clustering*. The Journal of Supercomputing, 2017. 73(11): p. 4773-4795.
- [7] Mafarja, M.M. and S. Mirjalili, *Hybrid whale optimization algorithm with simulated annealing for feature selection*. Neurocomputing, 2017. 260: p. 302-312.
- [8] Moradi, P. and M. Gholampour, *A hybrid particle swarm optimization for feature subset selection by integrating a novel local search strategy*. Applied Soft Computing, 2016. 43: p. 117-130.
- [9] Lim, H. and D.-W. Kim, *Generalized Term Similarity for Feature Selection in Text Classification Using Quadratic Programming*. Entropy, 2020. 22(4): p. 395.
- [10] Bhattacharyya, T., et al., *Mayfly in harmony: A new hybrid meta-heuristic feature selection algorithm*. IEEE Access, 2020. 8: p. 195929-195945.
- [11] Mirjalili, S. and S.Z.M. Hashim, *A new hybrid PSOGSA algorithm for function optimization*. in 2010 international conference on computer and information application. 2010.
- [12] Kumar, R.R., P. Viswanath, and C.S. Bindu, *Nearest neighbor classifiers: a review*. Int. J. Comput. Intell. Res, 2017. 13(2): p. 303-311.
- [13] <https://blog.text-mining.ir/text-preprocessing-o2wzgzkagzyj/>
- [14] <https://www.sobhe.ir/hazm/>

زیر نویس ها



شکل (۳): تعداد ویژگی انتخاب شده برای دو مجموعه داده با روش های LBPSO و LBPSO-GSA

با توجه به نتایج و تحلیل داده ها می توان متوجه این موضوع شد که داده هایی با تعداد کلاس بالاتر، در اغلب موارد نیاز به انتخاب ویژگی بیشتری دارند تا بتوانند به دقت قابل قبولی برسند. زیرا به علت وجود موضوع های متنوع، ویژگی های برجسته بیشتری دارند. مجموعه داده اول با وجود آنکه تعداد سند کمتری دارد، دارای ویژگی بیشتری است. به همین علت برای رسیدن به دقت مطلوب تعداد ویژگی بالاتری برای انتخاب شده است.

۳- نتیجه

این مقاله یک روش انتخاب ویژگی جدید با ترکیب الگوریتم های ازدحام ذره ها بر اساس لینک و جست و جوی گرانش ارائه می کند. انتخاب ویژگی نقشی مهم در کاهش هزینه محاسبه، ساده کردن یادگیری مدل و ارتقا توانایی طبقه بندی کننده در مسأله های طبقه بندی دارد. در این پژوهش روش انتخاب ویژگی روی داده های فارسی انجام شده است. برای ایجاد داده ها دو مجموعه داده تصادفی با کلاس های متفاوت از مجموعه داده همشهری ایجاد شده است. پیش از انجام انتخاب ویژگی، داده ها پیش پردازش شده اند و توسط روش TF-IDF وزن دهی شده اند. روش پیشنهادی جست و جوی محلی را توسط الگوریتم جست و جوی گرانش افزایش می دهد. این بهبود موجب شده است که با انتخاب تعداد ویژگی کمتر و انجام تعداد تکرار کمتر دقت بالاتری به دست آید.

۳ Big Data
۴ meta-heuristic
۵ Link based BPSO (LBPSO)
۶ Gravitational Search Algorithm (GSA)
۷ Particle swarm optimization (PSO)
۸ K-nearest neighbors (KNN)
۹ Filter
۱۰ Wrapper
۱۱ Hybrid or embedded
۱۲ Support vector machines (SVM)
۱۳ Naive Bayes
۱۴ Information gain

۱۵ Correlation based on Feature Selection (CFS)
۱۶ Whale Optimization Algorithm
۱۷ Simulated Annealing Algorithm
۱۸ iterated greedy
۱۹ Normalize
۲۰ Tokenize
۲۱ Stop Words
۲۲ Sigmoid function
۲۳ BPSO
۲۴ Term Frequency- Inverse Document Frequency
۲۵ label