

تاریخ دریافت مقاله: ۹۶/۱۰/۰۵
تاریخ پذیرش مقاله: ۹۷/۰۶/۰۷

بهبود الگوریتم فازی C-Means با الگوریتم ژنتیک برای انتخاب ویژگی‌ها در دسته‌بندی اسناد متنی

ندا محمودی جاریحان

مریی، گروه مهندسی کامپیوتر، واحد ارومیه، دانشگاه آزاد اسلامی، ارومیه، ایران
پست الکترونیکی: nmahmoodi510@gmail.com

فرهاد سلیمانیان قره چیق*

استادیار، گروه مهندسی کامپیوتر، واحد ارومیه، دانشگاه آزاد اسلامی، ارومیه، ایران
پست الکترونیکی: farhad@iaurmia.ac.ir

چکیده

الگوریتم فازی C-Means انجام می‌شود و این ویژگی‌ها به الگوریتم ژنتیک جهت بهبود در دسته‌بندی ارسال می‌گردند. روش پیشنهادی بر روی سه مجموعه داده مختلف Reuters21578, WEBKB, CADE 12 و بر اساس معیارهای ارزیابی مختلفی مورد آزمایش و ارزیابی قرار گرفته است. مقایسه نتایج روش پیشنهادی با سایر روش‌های مطرح در دسته‌بندی متون نشان می‌دهد که روش پیشنهادی عملکرد بهینه‌ای را در دسته‌بندی اسناد متنی دارد.

واژه‌های کلیدی: دسته‌بندی اسناد متنی، بهینه‌سازی، انتخاب ویژگی، الگوریتم فازی C-Means، الگوریتم ژنتیک.

۱- مقدمه

رشد چشمگیر منابع اطلاعاتی، فناوری‌های جدید و افزایش حجم پایگاه‌داده‌های الکترونیکی نیاز به دسته‌بندی متون را بیش از پیش آشکار می‌سازد. از آنجایی که حجم اطلاعات الکترونیکی و برخط روز به روز بیشتر می‌شود، دسترسی سریع و صحیح به منابع مهم، یکی از دغدغه‌های استفاده از این منبع اطلاعاتی بزرگ است. بسیاری از داده‌ها

افزایش روزافزون مستندات الکترونیکی در وب، لزوم دسته‌بندی آنان در دسته‌های مختلف را نشان می‌دهد. با توجه به حجم و دامنه وسیع اسناد متنی که به‌طور قابل توجهی از طریق محیط‌های برخط و سایر منابع قابل دسترسی می‌باشند، در صورت عدم دسته‌بندی مناسب، عمل بازیابی و پردازش اسناد متنی دسته‌بندی نشده با مشکلات زیادی مواجه می‌گردد. این نیاز منجر به ایجاد روش‌های نوین برای دسته‌بندی اطلاعات شده است. دسته‌بندی، تخصیص اسناد متنی یا ویژگی‌ها به یک یا چندین دسته است، به‌طوری که اسناد متنی با توجه به موضوعات یا میزان مشابهت ویژگی‌ها می‌توانند دسته‌بندی گردند. در ارائه روش‌های دسته‌بندی، استخراج و انتخاب ویژگی‌های کلیدی اسناد متنی از اهمیت بالایی برخوردار می‌باشد. در این مقاله روشی براساس بهبود الگوریتم فازی C-Means با الگوریتم ژنتیک برای انتخاب ویژگی‌ها در دسته‌بندی اسناد متنی ارائه شده است که در روش پیشنهادی انتخاب ویژگی‌های کلیدی متون از طریق

* نویسنده مسئول

در کاربردهای مختلف از جمله در اسناد متنی الکترونیکی، صفحات وب، ایمیل‌ها، اخبار و مقالات، نیاز به ابزار متن‌کاوی جهت دسته‌بندی خودکار متون را بیش از گذشته مبرم کرده است. بدیهی است دسته‌بندی اسناد موجب تسهیل در روند جستجو و بازیابی اطلاعات می‌شوند. از جمله زمینه‌های کاربردی دسته‌بندی می‌توان به متن‌کاوی، پیدا کردن مشابهت اسناد، ابهام زدایی، بازیابی اطلاعات، روش‌سازی و دسته‌بندی متن اشاره نمود. در تعریف کلی دسته‌بندی متون تخصیص اسناد متنی یا ویژگی‌ها (کلمات) به یک یا چندین دسته است، به طوری که اسناد متنی با توجه به موضوعات یا میزان مشابهت ویژگی‌ها می‌توانند دسته‌بندی گردند. دسته‌بندی متون تمام فعالیت‌هایی را که به نوعی به دنبال کسب دانش از متن هستند شامل می‌گردد و در سطوح بالا زیرمجموعه‌ای از متن‌کاوی و داده‌کاوی است. امروزه، روش‌های زیادی برای دسته‌بندی اسناد متنی خودکار خصوصاً روش‌های یادگیری ماشین وجود دارند.

دسته‌بندی متن به روش‌های یادگیری نظارت شده دلالت دارد، به طوری که این روش‌ها از نمونه‌های برچسب‌گذاری شده و از پیش تعریف شده در قسمت آموزش بهره برده و از اطلاعات به دست آمده در نمونه‌های آزمایشی در قسمت آزمایش سیستم برای دسته‌بندی استفاده می‌کنند. یادگیری نظارت شده یکی از انواع روش‌های یادگیری ماشین است که در آن ورودی و خروجی مشخص بوده و به عبارت بهتر ناظری در این روش یادگیری وجود دارد که اطلاعات را برای یادگیرنده فراهم می‌نماید و سیستم با استفاده از این اطلاعات سعی در پیدا کردن ساختاری مشخص در نمونه‌های داده دارد.

روش‌های دسته‌بندی متون مختلفی از جمله روش‌های آماری برای دسته‌بندی خودکار اسناد متنی وجود دارند. از جمله آن‌ها روش‌های k نزدیک‌ترین همسایه [۶]، الگوریتم ژنتیک [۲] و ماشین بردار پشتیبان [۵] را می‌توان نام برد. در اکثر روش‌هایی که برای دسته‌بندی متون ارائه شده‌اند

استخراج ویژگی‌های کلیدی از اهمیت بالایی برخوردار است و براساس این ویژگی‌ها دسته‌بندی متون انجام می‌گیرد. ثابت شده که تنها ۳۳ درصد کلمات در یک متن مفید هستند و می‌توان از آن‌ها برای استخراج اطلاعات استفاده نمود [۶]. اکثر کلمات موجود در متون برای رساندن هدف یک متن استفاده می‌شوند و گاهی اوقات تکراری هستند. به‌طور کلی هدف از انتخاب ویژگی‌های متون کم کردن حجم داده‌ها، تسریع در عملیات دسته‌بندی متون، کاهش زمان لازم برای آموزش، کاهش زمان محاسباتی و افزایش سرعت عملکرد روش پیشنهادی در دسته‌بندی متون می‌باشد. از این رو در این مقاله روشی براساس بهبود الگوریتم فازی C-Means با الگوریتم ژنتیک برای انتخاب ویژگی‌ها در دسته‌بندی اسناد متنی ارائه شده است. در روش پیشنهادی ویژگی‌های کلیدی و مرتبط با موضوع اصلی متون از طریق الگوریتم فازی C-Means انتخاب خواهد شد. سپس این ویژگی‌ها به‌عنوان ورودی‌های الگوریتم ژنتیک در نظر گرفته می‌شود. با توجه به عملکرد الگوریتم ژنتیک که یک الگوریتم بهینه‌سازی و مبتنی بر تکرار است و برای جستجو و یادگیری ماشین استفاده می‌شود [۷]، بهبود در عملکرد الگوریتم فازی C-Means برای دسته‌بندی متون ایجاد خواهد شد.

ساختار این مقاله به شرح زیر ساماندهی شده است: در بخش ۲ کارهای قبلی انجام شده در این حوزه بیان شده است. در بخش ۳ روش پیشنهادی به صورت کامل ارائه شده است. در بخش ۴ روش پیشنهادی به‌طور کامل بررسی و ارزیابی می‌شود و در نهایت در بخش ۵ به نتیجه‌گیری و کارهای آینده اشاره شده است.

۲- مروری بر کارهای گذشته

در سال‌های اخیر، حوزه دسته‌بندی متون به دلیل رشد زیاد پایگاه‌داده‌ها، اطلاعات در صفحات وب و حجم بالای داده‌های متنی در فضای اینترنت اهمیت ویژه‌ای پیدا کرده است و در این حوزه روش‌های مختلفی توسط پژوهشگران

ارائه گردیده است که در ذیل به تعدادی از آن‌ها اشاره شده است.

بینا و رهگذر در [۳] دسته‌بندی خودکار متون فارسی با استفاده از الگوریتم یادگیری ماشینی، k نزدیک‌ترین همسایه را پیاده‌سازی کردند که نتایج دسته‌بندی خودکار متون فارسی، با استفاده از معیارهای شاخص‌گذاری همایند n تایی و حذف کلمه نهایی‌ها و برای دسته‌بندی متون از الگوریتم یادگیری ماشینی و k نزدیک‌ترین همسایه استفاده شده است. در نهایت به منظور ارزیابی و مقایسه نتایج، دو معیار دقت و یادآوری برای هر روش شاخص‌گذاری نیز محاسبه کرده‌اند. نتایج به دست آمده نشان داد که بهترین روش شاخص‌گذاری متون فارسی همایند n تایی می‌باشد و حذف کلمه نهایی‌ها نتایج را اندکی بهبود می‌بخشد.

موضوع برجسب زنی موضوعی متون فارسی، با استفاده از الگوریتم بردار فاصله اطلاعات را دسته‌بندی کرده و پیاده‌سازی شده است. نعمتی و بصیری در [۴] به ارائه روشی برای استخراج کلمات کلیدی با استفاده از معیار tf-idf پرداخته‌اند و می‌توان با در نظر گرفتن احتمال وقوع کلمات و نیز فاکتور، تمرکز کلمات مناسب برای دسته‌بندی متون را شناسایی نمود. همچنین الگوریتمی بر پایه فاصله بین بردارهای بسامد کلمات کلیدی ارائه دادند و مشاهده شد که هرس دستی کلمات غیر کلیدی که به اشتباه در زمره کلمات کلیدی قرار گرفته‌اند، سبب ایجاد تاثیر بسزایی در افزایش معیار دقت و بازخوانی در ارزیابی می‌گردد که کار هرس بر پایه دانش زبان شناختی و معنی شناختی انجام شد.

در [۸] از یادگیری با روش سنتی بی‌ساده با برجسب‌گذاری چندگانه برای دسته‌بندی متون استفاده شده است. در ابتدا، روش‌های استخراج ویژگی‌ها بر مبنای روش تحلیل مولفه‌های اصلی به حذف ویژگی‌های بی‌ربط و تکراری و زاید می‌پردازد. سپس در روش‌های انتخاب خصیصه برای افزایش کارایی این روش از الگوریتم ژنتیک استفاده شده تا زیرمجموعه‌ای از ویژگی‌هایی که بیشترین

مشابهت را دارند انتخاب شود. روش پیشنهادی از کارایی قابل قبولی بر روی داده‌های واقعی برخوردار است.

Feng و همکاران در [۹] روشی بر مبنای روش بی‌ساده ارائه کرده‌اند که از ویژگی‌های انتخاب شده برای دسته‌بندی بهینه‌تر استفاده می‌نماید. قابلیت‌های منحصربه‌فرد این روش با استفاده از روش شاخص‌گذاری ویژگی‌ها با روش شاخص‌گذاری عمومی است. نتایج این روش حاکی از آن است که با روش‌های شناخته شده دیگر دسته‌بندی مثل ماشین بردار پشتیبان قابل رقابت است.

Chen و همکاران در [۱۰] با استفاده از الگوریتم‌های ژنتیک و تکاملی تفاضلی به دسته‌بندی اسناد متنی پرداخته‌اند. آن‌ها از عملگرهای ادغام و جهش به منظور بهینه‌سازی تعداد دسته‌ها استفاده می‌کنند. نتایج آزمایش‌های آن‌ها نشان می‌دهد که الگوریتم ژنتیک کارا تر از الگوریتم تکاملی تفاضلی می‌باشد و همچنین الگوریتم ترکیبی نسبت به دو الگوریتم ژنتیک و تکاملی تفاضلی کارا تر است.

استفاده از روش k نزدیک‌ترین همسایه به منظور مشابه بودن داده‌ها در دسته‌ها در [۱۱] استفاده شده است. نتایج آزمایش‌ها نشان می‌دهد که k نزدیک‌ترین همسایه کارایی خوبی در مقایسه با K -Means دارد و اسناد متنی را بهتر دسته‌بندی می‌کند. در دسته‌بندی اسناد متنی هدف بهینه کردن دسته‌ها در حداقل یا حداکثر تعداد خوشه‌ها می‌باشد. برای دسته‌بندی اسناد متنی از الگوریتم K -Means استفاده کرده‌اند. همچنین Luo و همکاران برای افزایش کارایی K -Means از الگوریتم ژنتیک به منظور انتخاب کلمات مرتبط به هم برای دسته‌بندی مشابه استفاده کردند. الگوریتم ژنتیک این توانایی را دارد که تکرار و ایجاد نسل‌های مختلف دقت دسته‌بندی را بهبود دهد و هزینه محاسباتی را کاهش دهد [۱۲]. نتایج آزمایش‌ها نشان می‌دهد که الگوریتم ژنتیک دقت الگوریتم K -Means را تا حد زیادی افزایش داده است. از جمله زمینه‌های کاربردی دسته‌بندی می‌توان به

متن کاوی، پیدا کردن مشابهت اسناد، ابهام زدایی، بازیابی اطلاعات، روش‌سازی و دسته‌بندی متن اشاره نمود. در [۱۳] Revathi و همکاران از شبکه عصبی مصنوعی و الگوریتم ژنتیک برای دسته‌بندی پویای متون بر مبنای محتوای آنان در رده‌های تعریف شده استفاده کرده‌اند. الگوریتم ژنتیک به منظور بهینه‌سازی دسته‌بندی و انتخاب ویژگی‌های مناسب برای دسته‌بندی به‌کار رفته است. مجموعه داده رویترز برای یادگیری و آزمایش استفاده شده است. نتایج این روش پویا که حاصل از ترکیب شبکه عصبی مصنوعی و الگوریتم ژنتیک می‌باشد نشان دهنده آن است که این روش کارایی خوبی در دسته‌بندی متون دارد.

روش‌های یادگیری ماشین کاربرد گسترده‌ای در دسته‌بندی متون دارند، از جمله روش ماشین بردار پشتیبان در [۱۴] برای دسته‌بندی متون ارائه شده است. روش ارائه شده از روش انتخاب ویژگی چندرده‌ای ماشین بردار پشتیبان برای دسته‌بندی متون بهره می‌گیرد. نتایج حاکی از آن است که روش ارائه شده دقت خوبی در دسته‌بندی متون دارد. همچنین، Patra و Sing از شبکه‌های عصبی مصنوعی با بازخورد برای دسته‌بندی و یادگیری با نظارت برای دسته‌بندی خودکار اسناد متنی استفاده کرده‌اند. روش یادگیری در این روش از شبکه عصبی مصنوعی، الگوریتم پس‌انتشار می‌باشد و پنج مرحله برای پیش‌پردازش داده‌ها و متون ارائه کرده است. نتایج روش با الگوریتم پس‌انتشار به‌دست آمده و فاکتور مربوطه به‌عنوان روش وزن‌دار برای دسته‌بندی استفاده شده است.

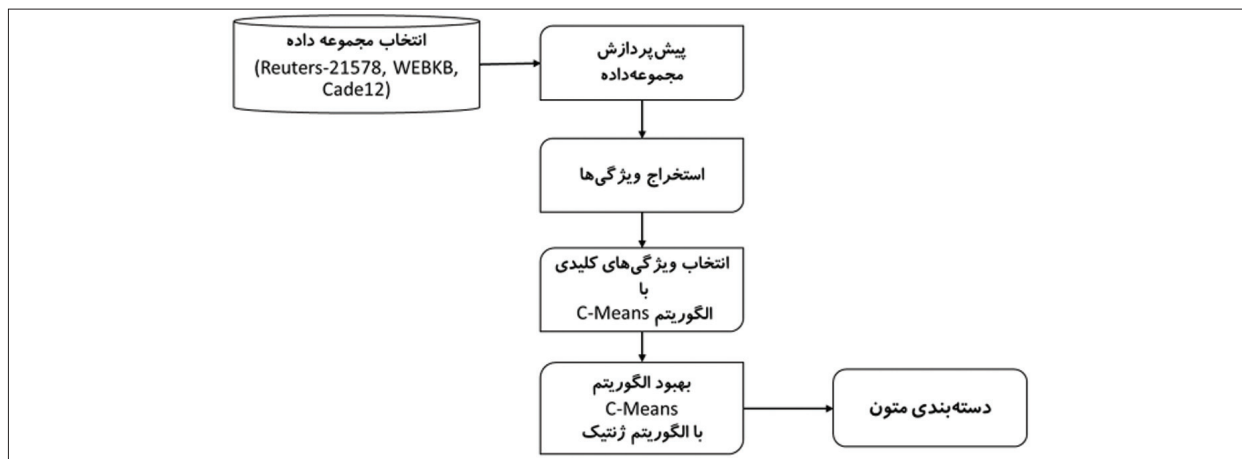
در [۱۵] Zhang و همکاران روشی نیمه نظارت شده برای دسته‌بندی متون ارائه کرده‌اند. این روش بر مبنای دسته‌بندی دسته‌های متون به‌صورت مجموعه‌ای از مولفه‌ها می‌باشد. در قسمت دسته‌بندی، از متون برچسب‌گذاری شده برای تعیین دسته‌های متنی و از متونی که برچسب‌گذاری نشده‌اند برای تعیین مراکز اقلیدسی دسته‌ها بهره می‌گیرد. زمانی که متن بدون برچسب‌گذاری برای دسته‌بندی در

مرحله آزمایش وارد سیستم پیشنهادی می‌گردد، میزان مشابهت آنان با دسته‌های متنی سنجیده شده و بر مبنای این‌که نزدیک‌ترین دسته کدام است، برچسب آن دسته به آن متن جدید تعلق می‌گیرد. آزمایش سیستم پیشنهادی بر روی مجموعه داده‌های Reuters-21578 و TanCrop V1.0 انجام گرفته است. نتایج حاکی از آن است که رویکرد نیمه نظارت شده پیشنهادی از روش ماشین بردار پشتیبان و روش شبکه عصبی مصنوعی با الگوریتم پس‌انتشار بهتر عمل می‌کند. همچنین این روش کارایی قابل مقایسه‌ای با روش بیز ساده با حداکثر بازده داشته و پیچیدگی محاسباتی کمتری دارد. بروز رسانی مراکز دسته‌ها، در دسته‌بندی اسناد متنی کاری مهم می‌باشد زیرا با بهینه‌سازی مراکز دسته‌ها، روند دسته‌بندی ویژگی‌ها در اسناد متنی موجب کیفیت بهتر آن خواهد شد.

Karimov و همکاران در [۱۶]، به ایجاد روشی بهینه برای بهبود روند دسته‌بندی در الگوریتم K-Means پرداخته‌اند. آن‌ها روش تکاملی ترکیبی با استفاده از الگوریتم K-Means را برای دسته‌بندی اسناد متنی مطرح کردند. از روش فراابتکاری در این روش برای شناسایی نامزدهای مناسب برای تعیین مراکز اقلیدسی انتخابی در الگوریتم K-Means استفاده شده است. بهبود نتایج دسته‌بندی ۳۰٪ بیشتر از الگوریتم K-Means استاندارد است. از جمله روش‌های فراابتکاری دیگر استفاده از شبیه‌سازی تبریدی در دسته‌بندی است. کار انجام یافته در دسته‌بندی اسناد متنی با روش شبیه‌سازی تبریدی [۱۷] به دسته‌بندی اسناد متنی در زبان چینی با روش همایند n تایی می‌پردازد.

۳- روش پیشنهادی

عامل اصلی در دسته‌بندی اسناد متنی، استخراج و انتخاب ویژگی‌های کلیدی متون می‌باشد و انتخاب ویژگی قبل از دسته‌بندی متن امری اساسی و ضروری است. در واقع انتخاب ویژگی‌های کلیدی متون از میان بیش‌شمار ویژگی‌های استخراج شده، دشوار بوده و از طرفی انتخاب



شکل ۱: نحوه عملکرد روش پیشنهادی

ویژگی‌ها، ویژگی‌های موثری در دسته‌بندی باشند. در روش پیشنهادی از الگوریتم فازی C-Means برای یافتن اطلاعات و ویژگی‌های مهم متون استفاده شده است.

۳-۱ الگوریتم فازی C-Means

الگوریتم فازی C-Means به‌طور کلی نمونه‌ها را به C دسته تقسیم می‌شوند و تعداد C از قبل مشخص شده است. الگوریتم C میانگین دقیقاً شبیه به الگوریتم K میانگین است با فرق این‌که در آنجا K دسته داشتیم و در C میانگین تعداد دسته‌ها برابر با C است. در این بخش الگوریتم پایه فازی C-Means به‌طور خلاصه معرفی خواهد شد. تابع هدف این الگوریتم به صورت زیر تعریف می‌شود:

$$J = \sum_{i=1}^N \sum_{j=1}^C u_{ij}^m d^2(x_i - v_j) \quad (2)$$

در فرمول (۲) u_{ij} میزان تعلق نمونه x_i به مرکز دسته v_j و m درجه فازی بودن را تعیین می‌کند. در اینجا $d^2(x_i - v_j)$ همان فاصله اقلیدسی و برابر با $(x_i - v_j)$ است، به‌طوری که x_i نمونه x_i و v_j مرکز دسته v_j است. بر اساس تابع هدف معرفی شده در فرمول، معادلات بروزرسانی مراکز و توابع تعلق بدین صورت خواهد بود:

$$u_{ij} = \frac{1}{\sum_{i=1}^C \left(\frac{d^2(x_i - v_j)}{d^2(x_i - v_i)} \right)^{\frac{1}{m-1}}} \quad (3)$$

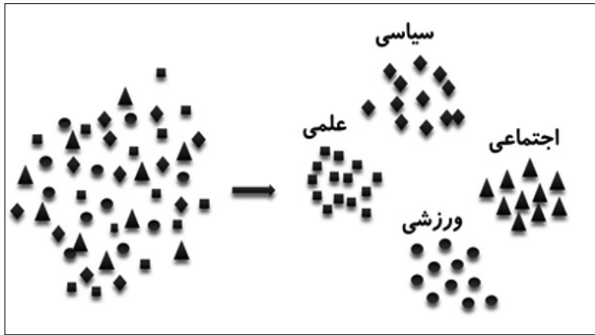
این ویژگی‌ها باعث افزایش کارایی در دسته‌بندی متون می‌گردد. لذا در این مقاله از الگوریتم فازی C-Means برای انتخاب ویژگی‌های کلیدی، به صورت جستجو در فضای کل ویژگی‌ها با حذف کلمات تکراری، زائد و غیر مرتبط استفاده شده است. در شکل (۱) نمای کلی از عملکرد روش پیشنهادی نشان داده شده است.

به‌طور کلی برای انتخاب ویژگی‌های کلیدی متون با الگوریتم فازی C-Means، تمام ویژگی‌های متون ابتدا باید استخراج شوند که از میان ویژگی‌های استخراج شده ویژگی‌های کلیدی و مرتبط با موضوع اصلی متن انتخاب گردد. برای استخراج ویژگی‌های متون از فرمول (۱) استفاده کردیم.

$$l_{tc} - \text{Weighting: } a_{ik} = \frac{\log(f_{ik} + 1.0) * \log(N/n_i)}{\sqrt{\sum_{j=1}^M [\log(f_{jk} + 1.0) * \log(N/n_j)]^2}} \quad (1)$$

در فرمول (۱)، f_{ik} فرکانس کلمه i در متن N ، k تعداد متون در مجموعه، M تعداد کلمات مجموعه پس از انجام عملیات کاهش و حذف کلمات اضافی و n_i : مجموع تعداد دفعاتی که کلمه i در هر مجموعه اتفاق افتاده است می‌باشد.

پس از استخراج ویژگی‌ها و لغات یک متن، این ویژگی‌ها تبدیل به برداری از کلمات می‌شوند و تمام لغات کلیدی، اسامی، افعال و ... مشخص می‌گردند. بعد از این مرحله باید ویژگی‌های بی‌فایده و نامرتب حذف گردند و ویژگی‌های کلیدی متون استخراج شود به‌طوری که این



شکل ۲: نمایی از الگوریتم فازی C-Means

تصادفی که جمعیت اولیه نامیده می‌شود شروع می‌شود و هر کروموزوم به‌عنوان فردی از جمعیت اولیه در نظر گرفته می‌شود. برای تولید نسل‌ها تمامی کروموزوم‌ها با یک تابع برازش متناسب با هر مسئله مورد ارزیابی قرار می‌گیرند و بهترین کروموزوم‌ها برای تولید نسل بعدی انتخاب می‌شوند. این الگوریتم، الگوریتم مبتنی بر تکرار است و برای بهینه‌سازی، جستجو و یادگیری ماشین مورد استفاده قرار می‌گیرد [۷].

باتوجه به این‌که الگوریتم ژنتیک در هر تکرار چندین کروموزوم از فضای جستجو را در نظر می‌گیرد بنابراین شانس این‌که به یک ماکزیمم محلی همگرا شود کاهش می‌یابد. این الگوریتم جمعیت‌های کاملی از کروموزوم‌ها را تولید می‌کند. سپس هر کروموزوم را به‌صورت انفرادی امتحان می‌کند و با ترکیب محتویات آن‌ها یک جمعیت جدید را که شامل نسل بهبود یافته است تشکیل می‌دهد. صرف نظر از انجام یک جستجو و ملاحظه هم‌زمان، الگوریتم ژنتیک تعدادی از کروموزوم‌ها را با ماشین‌های موازی تطبیق می‌دهد زیرا در اینجا تکامل هر کروموزوم یک فرآیند مستقل است. لذا الگوریتم ژنتیک فقط نیاز به اطلاعاتی در مورد کیفیت رده حل‌های ایجاد شده به وسیله هر مجموعه از متغیرها دارد، در صورتی که بعضی از روش‌های بهینه‌سازی نیاز به اطلاعات یا حتی نیاز به شناخت کامل از ساختمان مسئله و متغیرها دارند. چون الگوریتم ژنتیک نیاز به چنین اطلاعاتی مشخصی از مسئله ندارد بنابراین قابل انعطاف‌تر از بیشتر روش‌های جستجو است.

$$v_j = \frac{\sum_{i=1}^N u_{ij}^m x_i}{\sum_{i=1}^N u_{ij}^m} \quad (4)$$

الگوریتم‌های سلسله مراتبی از یک معیار شباهت استفاده کرده و در هر مرحله داده‌ها را به دودسته تقسیم و در نهایت یک ساختار درختی از این تقسیم بندی تحت عنوان دندروگرام ایجاد می‌کنند [۱۸]. اما الگوریتم‌های افرازکننده، داده‌ها را به‌طور مستقیم درون چند دسته قرار می‌دهند. این الگوریتم‌ها خود به دو دسته سخت (انحصاری) و نرم (فازی) تقسیم می‌شوند.

در روش پیشنهادی الگوریتم فازی C-Means برای انتخاب ویژگی‌های موثر، ابتدا سردسته‌های (علمی، ورزشی، سیاسی و...) از دسته‌ها را ایجاد می‌کند و براساس فاصله اقلیدسی تمام ویژگی‌های کلیدی متون رده‌بندی می‌شوند. با این فرایند ویژگی‌های با فاصله خیلی دور از سردسته حذف شده و ویژگی‌های کلیدی انتخاب می‌گردند. در شکل (۲) نمایی از دسته‌بندی نشان داده شده است.

در الگوریتم فازی C-Means برای مشخص شدن تعداد مناسب ویژگی‌ها این الگوریتم در هر تکرار از میان بیش‌مار ویژگی‌های استخراج شده از متون تعداد ۱/۲ یا ۱/۳ یا ۱/۴ و ... ویژگی‌ها را انتخاب می‌کند و در نهایت فقط ۱/۳ این ویژگی‌ها را که جزء مهم‌ترین و باارزش‌ترین ویژگی‌ها در تشخیص متنی از متن دیگر می‌باشد انتخاب می‌نماید. پس از این‌که ویژگی‌های مهم و باارزش انتخاب شدند تمام این ویژگی‌ها به الگوریتم ژنتیک ارسال می‌گردند.

۳-۲ الگوریتم ژنتیک

الگوریتم ژنتیک روش بهینه‌سازی الهام گرفته از طبیعت جاندار (موجودات زنده) است که بر اساس قانون تکامل داروین (بقا بهترین) که می‌گوید: موجودات ضعیف‌تر از بین می‌روند و موجودات قوی‌تر باقی می‌مانند ارائه گردیده است. به‌طور کلی این الگوریتم با تولید راه‌حل‌های

جدول ۱: لیست دسته‌ها و موضوعات مجموعه داده Reuters 21578

Reuters 21578				
Topics #	train # docs	test # docs	other #	Total # docs
0	1828	280	8103	10211
1	6552	2581	361	9494
2	890	309	135	1334
3	191	64	55	310
4	62	32	10	104
5	39	14	8	61
6	21	6	3	30
7	7	4	0	11
8	4	2	0	6
9	4	2	0	6
10	3	1	0	4
11	0	1	1	2
12	1	1	0	2
13	0	0	0	0
14	0	2	0	2
15	0	0	0	0
16	1	0	0	1

ابتدا اجرا می‌شود تا تعداد دسته‌بندی نادرست را به حداقل برساند.

۴- بررسی و ارزیابی

در این بخش به منظور مشاهده نتایج، الگوریتم ژنتیک و روش پیشنهادی بر روی مجموعه داده WebKB [۲۰]، Reuters 21578 [۱۹] و Cade12 [۲۱] اجرا شده است. در جدول (۱) تا (۳) انواع دسته‌ها و موضوعات در مجموعه داده در دو بخش آزمایش و آموزش نشان داده شده است. برای محاسبه و ارزیابی عملکرد روش پیشنهادی از پنج معیار مطرح در طبقه‌بندی متون که عبارتند از: تعداد داده‌های درست طبقه‌بندی شده، تعداد داده‌های نادرست طبقه‌بندی شده دقت، بازخوانی و F Measure استفاده شده است.

همچنین بعضی از پارامترهای الگوریتم‌های فازی C-Means و ژنتیک که بر روی روند روش پیشنهادی تاثیر مهمی دارند در جدول (۴) نشان داده



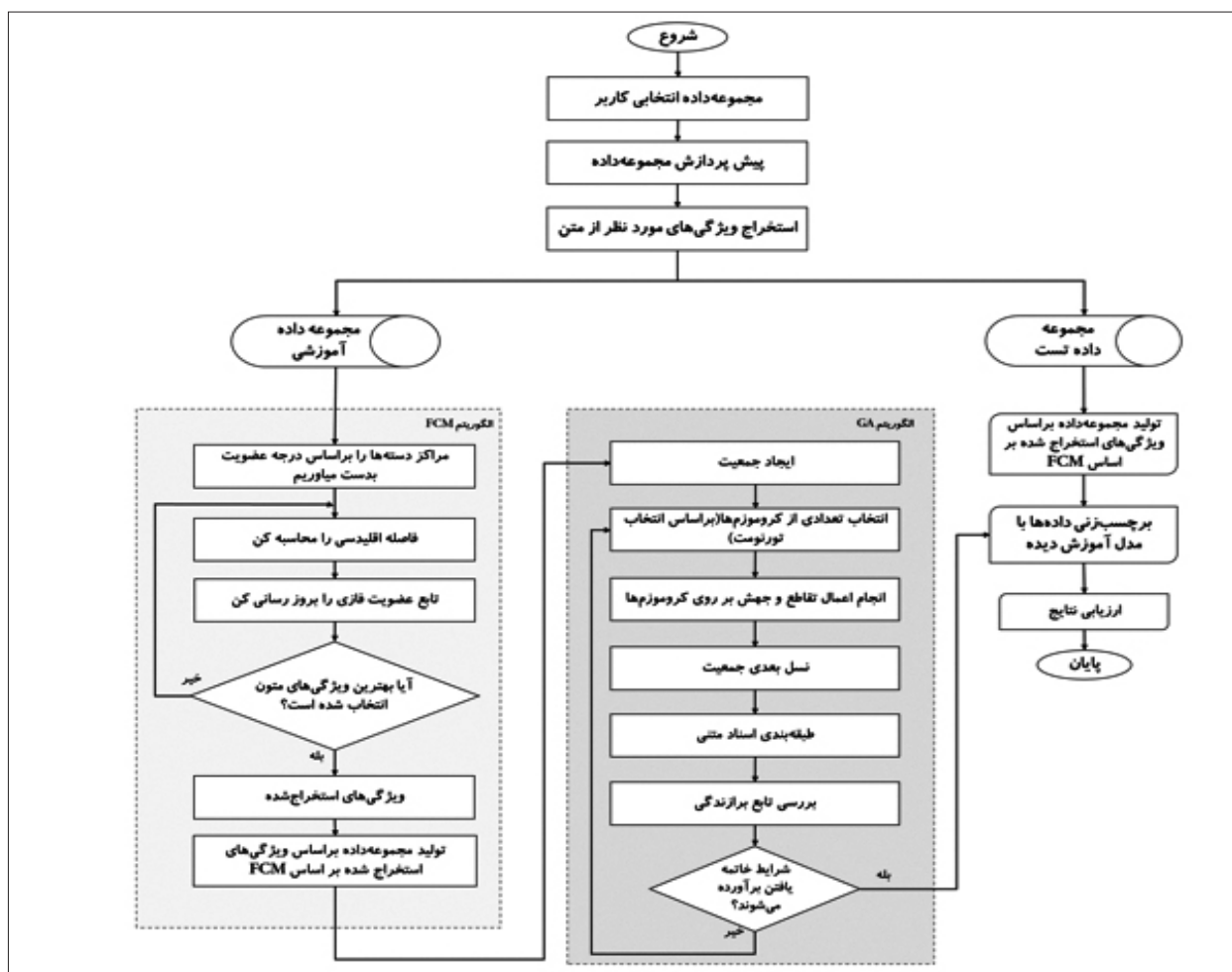
شکل ۳: هم‌ارزی عناصر ژنتیک در روش پیشنهادی

در الگوریتم ژنتیک تمام ویژگی‌های انتخاب شده یک متن توسط الگوریتم C-Means فازی به عنوان یک کروموزوم و ورودی‌های الگوریتم ژنتیک در نظر گرفته می‌شود. مجموع کروموزوم‌های الگوریتم ژنتیک، مجموع متون یک مجموعه داده را تشکیل می‌دهد. در شکل (۴)، شبه کد و در شکل ۵ روندنمای روش پیشنهادی نشان داده شده است. در شبه کد روش پیشنهادی مراحل مختلف روش پیشنهادی تشریح شده است.

در اجرای الگوریتم ژنتیک در هر مرحله، از ادغام دو طرفه و انتخاب تورنومنت تعدادی از متون انتخاب شده و به نسل بعدی انتقال داده می‌شود. بعد از اعمال ادغام به منظور اجتناب از همگرایی به بهینه محلی و ایجاد تنوع و گوناگونی در جمعیت با استفاده از عملگر جهش با نرخ ۰/۱ یک تعداد از کروموزوم‌های به دست آمده را تغییر می‌دهیم. بدین شکل هر متن براساس این‌که از چه ویژگی‌ها و کلمات کلیدی برخوردار است طبقه‌مربوط به آن‌ها، برایش پیش‌بینی می‌شود. در نهایت پس از تخصیص هر متن به یک دسته و پیش‌بینی دسته هر یک از متون، ستون هدف با رده‌بندی اولیه مقایسه می‌شود تا تعداد پیش‌بینی درست یا اشتباه روش پیشنهادی مشخص شود. در صورت تعداد زیاد پیش‌بینی روش پیشنهادی الگوریتم ژنتیک از

بهیود الگوریتم فازی C-Means با الگوریتم ژنتیک برای انتخاب ویژگی‌ها در دسته‌بندی اسناد متنی	
مجموعه داده انتخابی کاربر	ورودی‌ها
متون دسته‌بندی شده درست، نادرست و مقدار معیارهای ارزیابی دقت، دوباره خوانی، میانگین دقت و باز خوانی	خروجی‌ها
گام‌ها	گام
۱- بارگذاری داده‌های مجموعه داده‌ی انتخابی توسط کاربر ۲- تبدیل داده‌های موجود در مجموعه داده به شکل استاندارد و حذف داده‌های مبهم و تهی ۳- استخراج ویژگی‌های متون با استفاده از فرمول (۱) ۴- انتخاب ویژگی‌های کلیدی متون از میان ویژگی‌های استخراج شده متون	پیش پردازش مجموعه داده استخراج ویژگی‌های متون الگوریتم فازی C-Means
۵- کروموزوم‌ها براساس ویژگی‌های متون موجود در مجموعه داده جدید، برای ایجاد مجموعه جمعیت اولیه تولید می‌شوند ۶- تعدادی از کروموزوم‌ها به‌طور تصادفی انتخاب می‌شوند (از قانون انتخاب تورنومت استفاده می‌شود). ۷- عمل ترکیب بر روی دو کروموزوم دارای بالاترین تابع ارزیابی برای تولید نسل جدید اعمال می‌شود (تعداد پیش‌بینی درست و غلط به‌عنوان تابع ارزیابی الگوریتم ژنتیک در نظر گرفته شده است). ۸- عمل جهش بر روی کروموزوم اعمال شود. ۹- کروموزوم جدید به لیست کروموزوم‌های موجود اضافه شده و از گام ۲ ادامه می‌یابد ۱۰- شرط توقف دسته‌بندی تمام متون موجود در مجموعه داده است.	الگوریتم ژنتیک
۱۱- ذخیره‌سازی نتایج در روش پیشنهادی ۱۲- دسته‌بندی داده‌ها با روش آموزش دیده ۱۳- محاسبه معیارهای ارزیابی	دسته‌بندی متون موجود در مجموعه داده

شکل ۴: شبه کد روش پیشنهادی



شکل ۵: روندنمای مدل پیشنهادی

جدول ۲: لیست دسته‌ها و موضوعات مجموعه داده WebKB

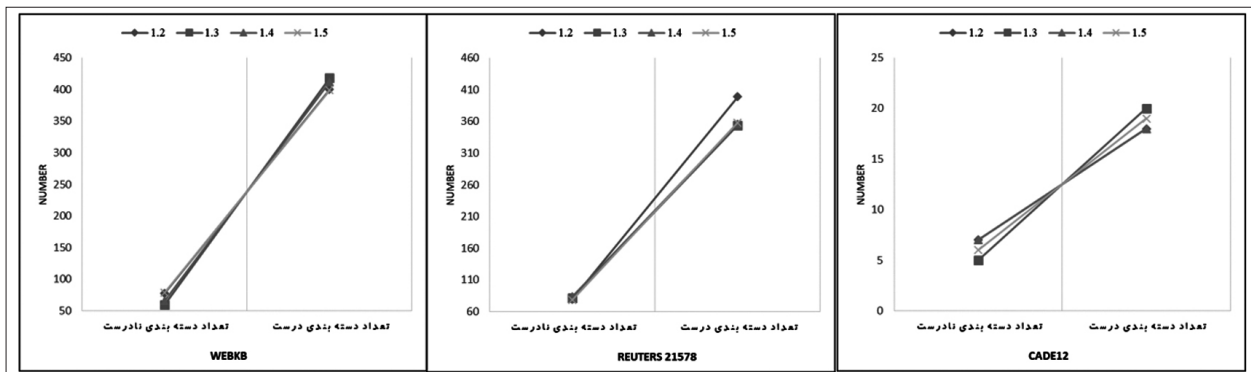
WebKB			
Class	train docs #	test docs #	Total # docs
project	336	168	504
course	620	310	930
faculty	750	374	1124
student	1097	544	1641
Total	2803	1396	4199

جدول ۳: لیست دسته‌ها و موضوعات مجموعه داده Cade 12

Cade12			
Class	# train docs	# test docs	Total # docs
01--servicos	5627	2846	8473
02--sociedade	4935	2428	7363
03--lazer	3698	1892	5590
04--informatica	2983	1536	4519
05--saude	2118	1053	3171
06--educacao	1912	944	2856
07--internet	1585	796	2381
08--cultura	1494	643	2137
09--esportes	1277	630	1907
10--noticias	701	381	1082
11--ciencias	569	310	879
12--compras-online	423	202	625
Total	27322	13661	40983

جدول ۴: پارامترهای مورد استفاده در روش پیشنهادی

پارامتر	مقدار
تعداد جمعیت اولیه	۱۰۰
تعداد تکرار	۵۰
نرخ ادغام	۰.۲۵
نرخ جهش	۰.۱
تابع ارزیابی	تعداد پیش‌بینی درست و غلط



شکل ۶: ارزیابی معیارهای تعداد دسته‌بندی درست و تعداد دسته‌بندی نادرست روش‌ها بر روی مجموعه داده‌های Reuters 21578, WEBKB و CADE12 با ۱/۲، ۱/۳، ۱/۴ و ۱/۵ ویژگی‌های انتخابی از کل ویژگی‌ها

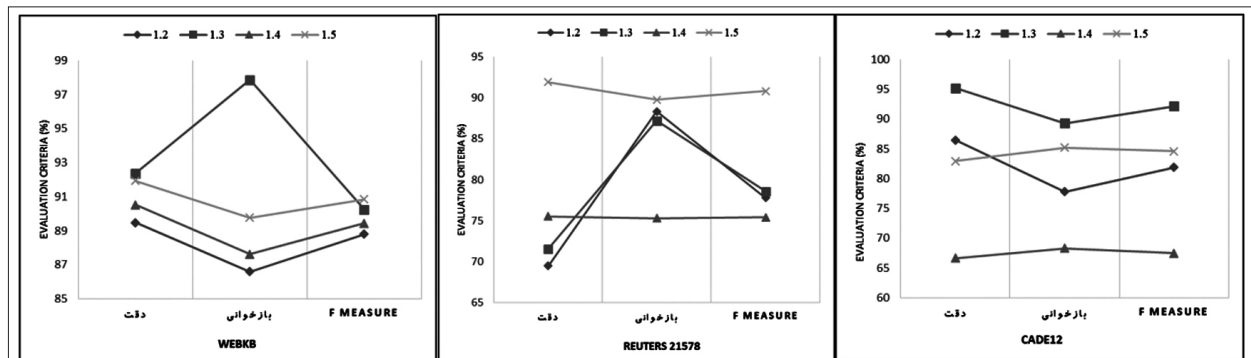
است. مقدار تمام این پارامترها بر مبنای آزمایش تعیین شده است. در جدول (۵) و شکل‌های (۶) و (۷) نشان داده شده است.

از آنجایی که انتخاب ویژگی‌های متون و انتخاب تعداد مختلف آن‌ها، اهمیت بالایی در دسته‌بندی درست اسناد متنی دارد. از این‌رو برای اطمینان از درستی انتخاب ۱/۳ تعداد ویژگی‌ها از میان بیش‌مار ویژگی‌های متون توسط الگوریتم فازی C-Means در روش پیشنهادی، ابتدا به بررسی تاثیر تعداد مختلف ویژگی‌های انتخاب شده از طریق الگوریتم فازی C-Means می‌پردازیم. نتایج حاصل از ارزیابی روش پیشنهادی با تعداد انتخاب ۱/۲، ۱/۳، ۱/۴ و ۱/۵ ویژگی‌های متون توسط الگوریتم فازی

با توجه به ارزیابی نتایج حاصل از انتخاب ۱/۲، ۱/۳، ۱/۴ و ۱/۵ ویژگی‌ها از میان بیش‌مار ویژگی توسط الگوریتم فازی C-Means انتخاب تعداد ۱/۳ ویژگی‌ها نتایج بهتری در مقایسه با تعداد ویژگی‌های مختلف انتخاب شده دارد. همچنین برای نشان دادن دقت به‌کارگیری همزمان دو الگوریتم ژنتیک و الگوریتم فازی C-Means نتایج پیاده‌سازی الگوریتم ژنتیک بدون دخالت الگوریتم فازی C-Means نیز در جدول (۶) نشان داده شده است.

جدول ۵: ارزیابی روش پیشنهادی با ویژگی‌های متفاوت انتخاب شده توسط الگوریتم فازی C-Means

مجموعه داده Cade 12				مجموعه داده Reuters 21578				مجموعه داده WebKB				
۱/۵	۱/۴	۱/۳	۱/۲	۱/۵	۱/۴	۱/۳	۱/۲	۱/۵	۱/۴	۱/۳	۱/۲	
۶	۷	۵	۷	۷۹	۸۴	۸۲	۸۰	۷۹	۶۶	۶۰	۷۸	دسته‌بندی نادرست
۱۹	۱۸	۲۰	۱۸	۳۹۹	۳۵۴	۳۵۶	۳۵۸	۳۹۹	۴۱۲	۴۱۸	۴۰۰	دسته‌بندی درست
۸۳	۶۷	۹۵	۸۶/۴۶	۹۲	۷۵/۵	۷۱/۵۵	۶۹.۵	۹۲	۹۰.۵	۹۲/۴	۸۹/۵	دقت
۸۵/۱۹	۶۸/۳	۸۹/۲۹	۷۷/۷۸	۸۹/۷۷	۷۵/۲۹	۸۷/۱۸	۸۸/۳۴	۸۹/۷۷	۸۷/۶	۹۷/۸۷	۸۶/۶	بازخوانی
۸۴/۵۹	۶۷/۵	۹۲/۲	۸۱/۸۸	۹۰/۸۳	۷۵/۳۹	۷۸/۵۹	۷۷/۷۹	۹۰/۸	۸۹/۴	۹۰/۳۳۸۰	۸۸/۸	F Measure



شکل (۷): ارزیابی معیارهای دقت، صحت و F Measure روش‌ها بر روی مجموعه داده‌های WEBKB، Reuters 21578 و CADE12 با ۱/۳، ۱/۴ و ۱/۵ ویژگی‌های انتخابی از کل ویژگی‌ها

جدول ۶: ارزیابی روش پیشنهادی بر روی مجموعه داده‌ها

مجموعه داده Cade 12		مجموعه داده Reuters 21578		مجموعه داده WebKB		
روش پیشنهادی	الگوریتم ژنتیک	روش پیشنهادی	الگوریتم ژنتیک	روش پیشنهادی	الگوریتم ژنتیک	
۵	۷	۸۲	۱۰۷	۶۰	۱۲۲	تعداد دسته‌بندی نادرست
۲۰	۱۸	۳۵۶	۳۳۱	۴۱۸	۳۵۶	تعداد دسته‌بندی درست
۹۵/۲۴	۷۵/۹۳	۷۱/۵۵	۶۸/۶۳	۹۲/۳۸	۸۷/۷۲	دقت
۸۹/۲۹	۷۹/۷۶	۸۷/۱۸	۷۵/۷۳	۹۷/۸۷	۸۵/۲۸	بازخوانی
۹۲/۲	۷۷/۸	۸۷/۵۹۴	۷۲/۲	۹۰/۲۴	۸۶/۵	F Measure

در جداول (۷) تا (۹) و شکل (۸) جمع‌آوری و نشان داده شده است.

۵. نتیجه‌گیری و کارهای آینده

نکته حائز اهمیت در ارائه یک روش دسته‌بندی دقیق، انتخاب ویژگی‌های مفید و حذف ویژگی‌های غیرمفید متون از میان تمام ویژگی‌های استخراج شده می‌باشد که براساس انتخاب ویژگی‌های مفید، دسته‌بندی متون با دقت و سرعت بالاتری امکان‌پذیر شود. برای تحقق این هدف در

براساس مقادیر جدول (۶) می‌توان نتیجه گرفت که روش پیشنهادی در تمام معیارها بر روی سه مجموعه داده Cade 12، WebKB و Reuters 21578، عملکرد بهتری از الگوریتم ژنتیک داشته است و به‌کارگیری همزمان دو الگوریتم C-Means فازی و ژنتیک در دستیابی به نتایج بهتر در دسته‌بندی اسناد متنی ضروری بوده است. در آخر نتایج معیارهای دقت، بازخوانی و F-Measure روش پیشنهادی با تعدادی از روش‌های ارائه شده مورد مقایسه قرار داده شده است. نتایج حاصل از مقایسه

جدول ۷: مقایسه روش پیشنهادی با سایر روش‌های ارائه‌شده در طبقه‌بندی اسناد متنی در مجموعه داده Reuters 21578

معیارهای ارزیابی					
F-Measure	بازخوانی	دقت			
۶۹,۷۷	۶۹,۳۲	۷۰,۲۳	TF		الگوریتم ادابوست [۱۸]
۶۸,۵۲	۶۵,۱۰	۷۲,۳۲	NORMTF		
۷۱,۱۸	۶۶,۸۷	۷۶,۰۹	LOGTF		
۷۲,۷۷	۶۸,۲۴	۷۷,۹۴	ITF		
۶۹,۹۴	۶۶,۰۷	۷۴,۲۹	SPARCK		
۶۸,۷۹	۶۵,۳۲	۷۲,۶۵	TF		[۱۸]
۶۶,۹۸	۶۴,۷۷	۶۹,۳۵	NORMTF		
۶۷,۰۵	۶۲,۰۰	۷۳,۰۰	LOGTF		
۷۲,۵۲	۶۹,۳۵	۷۶,۰۰	ITF		
۷۰,۰۷	۶۸,۱۱	۷۲,۱۴	SPARCK		
۷۲,۷۴	۶۹,۴۷	۷۶,۳۴	K=3	تعداد k مختلف	K نزدیک‌ترین همسایه [۲۲]
۷۰,۱۳	۶۷,۲۸	۷۳,۲۴	K=4		
۶۷,۵۱	۶۵,۱۲	۷۰,۰۸	K=5		
۸۶,۱۶	۷۹,۳۲	۹۴,۳۴	K=3	تعداد k مختلف	FS=80[22]
۸۳,۷۹	۷۶,۲۴	۹۳,۰۰	K=4		
۸۱,۸۵	۷۴,۲۰	۹۱,۲۷	K=5		
۸۸,۱۶	۸۲,۰۲	۹۵,۳۰	K=3	تعداد k مختلف	FS=120[22]
۸۷,۹۵	۸۱,۱۵	۹۶,۰۰	K=4		
۸۷,۰۵	۸۰,۲۰	۹۵,۱۸	K=5		
-	-	۹۱,۶۰	K=3	تعداد k مختلف	NB-K-Means[23]
-	-	۹۰,۱۰	K=4		
-	-	۸۳,۴۰	K=5		
-	-	۸۷,۳۲	K=3	تعداد k مختلف	KNN-K-Means[24]
-	-	۷۲,۳۷	K=4		
-	-	۶۷,۱۵	K=5		
۶۴/۶	۷۵/۹۶	۵۶/۱۶	۱/۲	تعداد ویژگی‌های متفاوت انتخاب شده	الگوریتم ژنتیک
۷۲/۲	۷۵/۷۳	۶۸/۶۳	۱/۳		
۷۹/۲	۸۲/۹۶	۷۵/۷۶	۱/۴		
۸۷/۲۴	۸۶/۱۷	۸۷/۸۹	۱/۵		
۷۷,۷۹۰	۸۸,۳۴	۶۹,۴۹	۱/۲	تعداد ویژگی‌های متفاوت انتخاب شده	روش پیشنهادی
۷۸,۵۹	۸۷,۱۸	۷۱,۵۵	۱/۳		
۷۵,۳۹۷	۷۵,۲۹	۷۵,۵۱	۱/۴		
۹۰,۸۳	۸۹,۷۷	۹۱,۹۲	۱/۵		

Reuters 21578

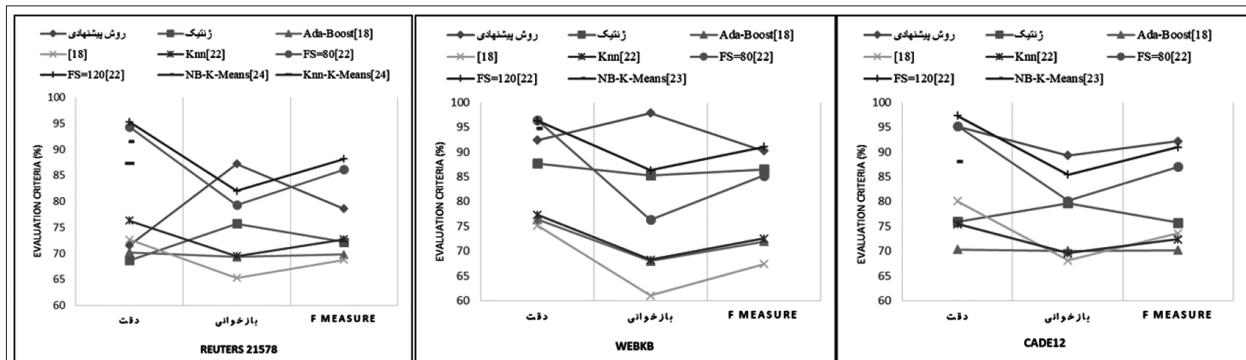
دیگر است، انتخاب شد و به الگوریتم ژنتیک برای بهبود در دسته‌بندی متون ارسال گردید. در روش پیشنهادی دسته‌بندی متون فقط براساس ۱/۳ ویژگی‌های انتخاب شده از طریق الگوریتم C-Means فازی انجام گرفت که این عمل باعث افزایش سرعت و دقت روش پیشنهادی در

این مقاله روشی براساس بهبود الگوریتم فازی C-Means با الگوریتم ژنتیک برای انتخاب ویژگی‌ها در دسته‌بندی اسناد متنی ارائه شده است که از طریق الگوریتم فازی C-Means از بین تمامی ویژگی‌های استخراج شده ۱/۳ ویژگی‌ها که مفیدترین ویژگی‌ها در جداسازی متنی از متن

جدول ۸: مقایسه روش پیشنهادی با سایر روش‌های ارائه شده در طبقه‌بندی اسناد متنی در مجموعه داده WEBKB

معیارهای ارزیابی				
F-Measure	بازخوانی	دقت		
۷۱٫۹۵	۶۸٫۰۰	۷۶٫۳۸	TF	
۷۱٫۰۷	۶۹٫۹۴	۷۲٫۲۴	NORMTF	
۷۰٫۵۲	۶۵٫۳۷	۷۶٫۵۴	LOGTF	
۶۹٫۹۹	۶۴٫۸۷	۷۵٫۹۸	ITF	
۷۱٫۴۸	۶۷٫۴۷	۷۶٫۰۰	SPARCK	
۶۷٫۳۸	۶۱٫۰۲	۷۵٫۲۳	TF	
۶۹٫۰۶	۶۶٫۵۵	۷۱٫۷۷	NORMTF	
۷۰٫۷۷	۶۸٫۰۸	۷۳٫۶۸	LOGTF	
۷۲٫۷۷	۶۹٫۷۷	۷۶٫۰۴	ITF	
۷۴٫۶۹	۷۰٫۳۶	۷۹٫۵۸	SPARCK	
۷۲٫۵۱	۶۸٫۲۴	۷۷٫۳۵	K=3	تعداد k مختلف
۶۹٫۱۹	۶۵٫۰۷	۷۳٫۸۷	K=4	
۶۵٫۴۷	۶۲٫۰۱	۶۹٫۳۴	K=5	
۸۵٫۲۳	۷۶٫۳۴	۹۶٫۴۷	K=3	تعداد k مختلف
۸۳٫۴۲	۷۴٫۲۱	۹۵٫۲۴	K=4	
۸۲٫۰۶	۷۳٫۴۱	۹۳٫۰۳	K=5	
۹۱٫۰۴	۸۶٫۳۰	۹۶٫۳۴	K=3	تعداد k مختلف
۸۷٫۷۹	۸۲٫۲۴	۹۴٫۱۵	K=4	
۸۵٫۸۱	۸۰٫۲۰	۹۲٫۲۷	K=5	
-	-	۹۴٫۸۰	K=3	تعداد k مختلف
-	-	۹۳٫۲۰	K=4	
-	-	۷۸٫۲۰	K=5	
۸۷/۵	۸۵/۷۷	۸۹/۳۷	۱/۲	تعداد ویژگی‌های متفاوت انتخاب شده
۸۶/۵	۸۵/۲۸	۸۷/۷۲	۱/۳	
۸۵/۷۱۸	۸۴/۳۵	۸۷/۱۳	۱/۴	
۸۷/۲۴	۸۶/۱۷	۸۷/۸۹	۱/۵	
۸۸/۸	۸۶/۵۹	۸۹/۴۷	۱/۲	تعداد ویژگی‌های متفاوت انتخاب شده
۹۰/۲۳۸۰	۹۷/۸۷	۹۲/۳۸	۱/۳	
۸۹/۴۳	۸۷/۶۲	۹۰/۵۲	۱/۴	
۹۰/۸۳۶	۸۹/۷۷	۹۱/۹۲	۱/۵	

WEBKB



شکل ۸: ارزیابی روش‌ها بر روی مجموعه داده‌های WEBKB و Reuters 21578 و CADE12 با ۱/۳ ویژگی‌های انتخابی از کل ویژگی‌ها

جدول (۹): مقایسه روش پیشنهادی با سایر روش‌های ارائه شده در طبقه‌بندی اسناد متنی در مجموعه داده CADE12

معیارهای ارزیابی						
F-Measure	بازخوانی	دقت				
۷۰,۱۸	۷۰,۰۵	۷۰,۳۲	TF		[18]Ada-Boost	CADE12
۷۲,۳۹	۶۸,۳۴	۷۶,۹۴	NORMTF			
۶۹,۲۵	۶۵,۰۸	۷۴,۰۰	LOGTF			
۷۳,۸۱	۷۱,۲۴	۷۶,۵۷	ITF			
۶۹,۴۷	۶۵,۴۶	۷۴,۰۱	SPARCK			
۷۳,۶۰	۶۸,۰۸	۸۰,۱۰	TF		[18]	
۷۱,۴۶	۶۹,۱۱	۷۳,۹۸	NORMTF			
۷۰,۳۲	۶۵,۰۰	۷۶,۵۸	LOGTF			
۸۱,۳۴	۸۰,۲۵	۸۲,۳۶	ITF			
۷۳,۶۸	۷۱,۶۸	۷۵,۰۹	SPARCK			
۷۲,۴۱	۶۹,۵۸	۷۵,۴۸	K=3	تعداد k مختلف	Knn[22]	
۶۹,۵۸	۶۷,۳۲	۷۲,۰۰	K=4			
۶۷,۴۲	۶۶,۴۱	۶۸,۴۶	K=5			
۸۷,۰۷	۸۰,۱۴	۹۵,۳۰	K=3	تعداد k مختلف	FS=80[22]	
۸۵,۶۵	۷۹,۳۳	۹۳,۰۷	K=4			
۸۳,۵۶	۷۶,۸۴	۹۱,۵۶	K=5			
۹۱,۰۲	۸۵,۴۷	۹۷,۳۴	K=3	تعداد k مختلف	FS=120[22]	
۸۹,۲۲	۸۴,۱۹	۹۴,۹۰	K=4			
۸۶,۸۳	۸۱,۲۲	۹۳,۲۷	K=5			
-	-	۸۸,۱۰	K=3	تعداد k مختلف	NB-K-Means[23]	
-	-	۸۵,۸۰	K=4			
-	-	۷۰,۳۰	K=5			
۶۴/۵۷۷	۶۷/۷۸	۶۱/۶۷	۱/۲	تعداد ویژگی‌های متفاوت انتخاب شده	الگوریتم ژنتیک	
۷۷/۷۹۶	۷۹/۷۶	۷۵/۹۳	۱/۳			
۷۵/۷۷۱	۷۹/۶۹	۷۲/۲۲	۱/۴			
۷۶/۷۳۸	۷۰/۳۷	۸۴/۳۸	۱/۵			
۸۱/۸۸۸	۷۷/۷۸	۸۶/۴۶	۱/۲	تعداد ویژگی‌های متفاوت انتخاب شده	روش پیشنهادی	
۹۲/۱۶۵	۸۹/۲۹	۹۵/۲۴	۱/۳			
۶۷/۴۷۵	۶۸/۳	۶۶/۶۷	۱/۴			
۸۴/۵۹	۸۵/۱۹	۸۲/۹۶	۱/۵			

دسته‌بندی اسناد متنی می‌شود.

در این مقاله از پنج معیار ارزیابی مطرح در طبقه‌بندی متون که عبارتند از: تعداد داده‌های درست طبقه‌بندی شده، تعداد داده‌های نادرست طبقه‌بندی شده، دقت، بازخوانی و F Measure برای محاسبه و ارزیابی عملکرد روش پیشنهادی استفاده شده است. براساس ارزیابی روش

پیشنهادی در تمام معیارها بر روی سه مجموعه داده Cade12 و WebKB, Reuters 21578 نتایج حاصله بیانگر عملکرد بهتر روش پیشنهادی در مقایسه با الگوریتم ژنتیک بود و به‌کارگیری همزمان دو الگوریتم C-Means فازی و ژنتیک در دستیابی به نتایج بهتر در دسته‌بندی اسناد متنی ضروری بوده است.

منابع:

57, pp. 1124-1130, 2015.

[15] W. Zhang, X. Tang, T. Yoshida, TESC: An Approach to Text Classification Using Semi-Supervised Clustering, Knowledge-Based Systems, Vol. 75, pp. 152-160, 2015

[16] J. Karimov, M. Ozbayoglu, Clustering Quality Improvement of K-Means Using a Hybrid Evolutionary Model, Procedia Computer science, Vol. 61, pp. 38-45, 2015.

[17] C. H. Chang, Simulated Annealing Clustering of Chinese Words for Contextual Text Recognition, Pattern Recognition Letters, Vol. 17, No. 1, pp. 57-66, 1996.

[18] Majidpour, Hiwa, and Farhad Soleimanian Gharehchopogh. "An Improved Flower Pollination Algorithm with AdaBoost Algorithm for Feature Selection in Text Documents Classification". Journal of Advances in Computer Research, pp: 1-11(2008).

[19] <http://archive.ics.uci.edu/ml/datasets/Reuters+21778+Text+Categorization+Collection> [Last Access: 08.08.2017]

[20] Craven, M., McCallum, A., PiPasquo, D., Mitchell, T., & Freitag, D. (1998). Learning to extract symbolic knowledge from the world wide web. In: DTIC Document.

[21] <http://ana.cachopo.org/datasets-for-single-label-text-categorization> [Last Access: 08.08.2017]

[22] Allahverdipour, Ali, and Farhad Soleimanian Gharehchopogh. "An Improved K-Nearest Neighbor with Crow Search Algorithm for Feature Selection in Text Documents Classification." Journal of Advances in Computer Research 9.2, pp: 1-11(2018).

[23] A. Allahverdipour, F.S. Gharehchopogh, A New Hybrid Model of K-Means and Naïve Bayes Algorithms for Features Selection in Text Documents Classification, Journal of Advances in Computer Research, Vol. 8, No. 4, pp: 73-86, 2017.

[24] R. Habibpour, K. Khalilpour, A New Hybrid K-means and K-Nearest-Neighbor Algorithms for Text Document Clustering, International Journal of Academic Research, Vol. 6 Issue 3, pp. 7984, 2014.

[۱] نعمتی شهلا، بصیری محمد احسان،، دسته بندی اسناد با استفاده از الگوریتم KNN، دهمین کنفرانس مهندسی برق ایران، ۱۳۸۶.

[۲] جلیلی سعید، گرانی شیما، دسته بندی متون با رویکرد الگوریتم ژنتیک، پانزدهمین کنفرانس مهندسی برق ایران، ۶۷-۷۳، ۱۳۸۶.

[۳] بینا بهاره و رهگذر مسعود و ده موبد آذین،، دسته بندی خودکار متون فارسی، انجمن کامپیوتر ایران، ۴-۱۳، ۱۳۸۶.

[۴] نعمتی شهلا، بصیری محمد احسان، دسته بندی اسناد فارسی با استفاده از الگوریتم knn. ۸-۱، ۱۳۹۲.

[5] Shao, Z., Yang, S., Gao, F., Zhou, K., & Lin, P. A New Electricity Price Prediction Strategy Using Mutual Information-Based SVM-RFE Classification. Renewable and Sustainable Energy Reviews, 70, 330-341, (2017).

[6] Liu, L., Kang, J., Yu, J., & Wang, Z. (2005). "A Comparative Study on Unsupervised Feature Selection Methods for Text Clustering", Proc. of the Conference on Natural Language Processing and Knowledge Engineering, IEEE NLP-KE'05. Proceedings of 2005 IEEE International Conference on (pp. 597-601). IEEE.

[7] Webber, W. E. 2010. Measurement in Information Retrieval Evaluation. Ph.D Thesis. The university of Melbourne <http://www.umiacs.umd.edu/~wew/wew-thesis-PhD.pdf> [Last Access: 08.08.2017]

[8] M. L. Zhang, J. M. Pena, V. Robles, Feature Selection for Multi-Label Naive Bayes Classification, Information Sciences, Vol. 179, Iss. 19, pp. 3218-3229, 2009

[9] G. Feng, J. Guo, B. Y. Jing, T. Sun, Feature Subset Selection Using Naive Bayes for Text Classification, Pattern Recognition Letters, Vol. 65, pp. 109-115, 2015.

[10] Chen Y., Qin B., Liu T., Liu Y., Li Sheng, The Comparison of SOM and K-Means for Text Clustering", Vol. 3, No. 2, pp. 268-274, 2010.

[11] C. Luo, Y. Li, S.M. Chung, Text Document Clustering Based on Neighbors, Data & Knowledge Engineering, Vol. 68, pp. 1271-1288, 2009.

[12] Verma H., Kandpal E., Pandey B., Dhar J., A Novel Document Clustering Algorithm Using Squared Distance Optimization Through Genetic Algorithms", International Journal on Computer Science and Engineering, Vol. 02, No. 05, pp. 1875-1879, 2010.

[13] N. Revathi, A. Peter, S.J. Kumar, Web Text Classification Using Genetic Algorithm and a dynamic neural Network Model, International Journal of Advanced Research in Computer Engineering and Technology (IJARCET), Vol. 2, Iss. 2, February 2013.

[14] B. Ramesh, J. G. R. Sathiaselvan, An advanced Multi Class instance selection based Support Vector Machine for Text Classification, Procedia Computer Science, Vol.