

Improving Text Mining with Featured Word Selection

¹ M.Amin Abolghasemi, ²Saeedeh Momtazi

¹ Master Student of Artificial Intelligence, Amirkabir University of Technology, Tehran, Iran
amin.a222@yahoo.com

² Assistant Professor of Artificial Intelligence, Amirkabir University of Technology, Tehran, Iran
momtazi@aut.ac.ir

Abstract

Text mining is one of the main tasks in web research that aims at classification or clustering available texts in the web for different applications, such as news analysis and social network analysis. Since a very large amount of textual data is available on the Web, reducing the dimension of data using feature extraction techniques plays an important role in improving the efficiency and effectiveness of the text mining algorithms. Various techniques have been proposed in machine learning tasks that can also be applied in the text mining domain. In this paper we study the available techniques and compare their impact on improving Persian text classification performance. Our experimental results on Hamshahri corpus shows that using an appropriate feature selection technique can improve the classification f-measure from 88.12% to 93.07%.

Keywords: Web Mining, Text Mining, Text Classification, Feature Selection

Archive of SID

بهبود متن کاوی با انتخاب کلمات ویژگی

محمدامین ابوالقاسمی^۱، سعیده ممتازی^۲

^۱ دانشجوی کارشناسی ارشد، گروه هوش مصنوعی دانشکده مهندسی کامپیوتر و فناوری اطلاعات، دانشگاه صنعتی امیرکبیر، تهران، ایران
amin.a222@yahoo.com

^۲ استادیار، گروه هوش مصنوعی دانشکده مهندسی کامپیوتر و فناوری اطلاعات، دانشگاه صنعتی امیرکبیر، تهران، ایران
momtazi@aut.ac.ir

چکیده

متن کاوی یکی از فعالیتهای اصلی در حوزه وب پژوهی محسوب می گردد که هدف آن دسته بندی یا خوشه بندی متون موجود در وب برای کاربردهای مختلف از جمله تحلیل خبر، تحلیل شبکه های اجتماعی و ... می باشد. با توجه به بالا بودن حجم دادگان موجود در وب برای پردازش های متن، کاهش ابعاد دادگان با کمک روش های استخراج ویژگی نقش مهمی را در بهبود کیفیت متن کاوی و همین طور بهینه سازی زمان اجرا ایفا می نماید. روش های متنوعی برای استخراج ویژگی در الگوریتم های یادگیری ماشین ارائه شده است که قابلیت کاربردی سازی در حوزه متن کاوی را دارند. در مقاله حاضر به بررسی الگوریتم های موجود در این زمینه پرداخته می شود و نتایج حاصل از این الگوریتم ها در استخراج کلمات ویژگی متون فارسی مقایسه می گردد. همچنین تاثیر به کارگیری انتخاب ویژگی در دسته بندی متون فارسی مورد تحلیل قرار می گیرد. نتایج به دست آمده در آزمایش ها بر روی پیکره همشهری فارسی نشان می دهد با کمک روش مناسب انتخاب ویژگی می توان نتایج دسته بندی متون فارسی را از ۸۸.۱۲٪ در معیار اف به ۹۳.۰۷٪ افزایش داد.

کلمات کلیدی

وب کاوی، متن کاوی، دسته بندی متون، انتخاب ویژگی

محاسباتی این فعالیت ها می گردد در حالیکه در اکثریت کاربردهای آن ها در وب کاوی سرعت بالای پردازش نقش کلیدی در میزان رضایت کاربر دارد. برای رفع این مشکل ایده استخراج ویژگی از مستندات پیشنهاد شده است که باعث کاهش حجم فضای برداری ایجاد شده برای هر سند می گردد. مهمترین مزیت انتخاب ویژگی، کاهش ابعاد داده و متعاقباً افزایش سرعت اجرای الگوریتم ها می باشد. به علاوه این کار باعث می شود تا میزان استعداد بیش برآزش الگوریتم های مورد استفاده در متن کاوی کاهش یابد. در این راستا در فعالیتهای مختلف حوزه یادگیری ماشین به منظور استخراج ویژگی های مناسب و بهینه معیارهای گوناگونی ارائه شده است [3]. اگرچه تمام این روش ها می توانند نقش موثر در استخراج ویژگی های کلیدی داده داشته باشند اما بررسی عمیق تری مورد نیاز است تا بتوان از میان روش های متنوع موجود بهترین گزینه را برای کاربردهای مختلف انتخاب نمود.

۱- مقدمه

بازیابی اطلاعات و متن کاوی (اعم از دسته بندی و یا خوشه بندی متون) جز اساسی ترین فعالیت ها در حوزه وب کاوی می باشند. هدف از بازیابی اطلاعات یافتن مستندات مرتبط با کوئری ارائه شده توسط کاربر می باشد. وظیفه دسته بندی متون تخصیص هر سند به یک کلاس مشخص و از پیش تعیین شده می باشد که با کمک الگوریتم های یادگیری بانظارت انجام می پذیرد [1,6]. خوشه بندی متون نیز سعی در تقسیم بندی مستندات دارد با این تفاوت که در آن نوع خوشه ها از پیش تعیین شده نیستند و این تقسیم بندی تنها براساس شباهت میان مستندات و با کمک الگوریتم های یادگیری بدون نظارت صورت می گیرد.

در تمامی فعالیتهای فوق حجم بسیار بالایی از دادگان متنی مورد بازیابی، دسته بندی، و یا خوشه بندی قرار می گیرد که باعث بالا رفتن پیچیدگی

$$\begin{aligned}
 IG(w) &= - \sum_{i=1}^K P(c_i) \log P(c_i) \\
 &+ P(w) \sum_{i=1}^K P(c_i|w) \log P(c_i|w) \\
 &+ P(\bar{w}) \sum_{i=1}^K P(c_i|\bar{w}) \log P(c_i|\bar{w}) \quad (3)
 \end{aligned}$$

در رابطه (۳) w و c_i به ترتیب بیان گر کلمه و کلاس نام می باشند. همچنین احتمالات موجود در این رابطه به صورت زیر محاسبه می شوند.

$$\begin{aligned}
 P(c_i) &= \frac{N_i}{N} & P(w) &= \frac{N_w}{N} & P(\bar{w}) &= \frac{N_{\bar{w}}}{N} \\
 P(c_i|w) &= \frac{N_{iw}}{N_w} & P(c_i|\bar{w}) &= \frac{N_{i\bar{w}}}{N_{\bar{w}}}
 \end{aligned}$$

در روابط بالا N برابر با تعداد کل اسناد، N_i برابر با تعداد کل اسناد موجود در کلاس c_i ، N_w برابر با تعداد کل اسنادی که شامل کلمه w هستند، $N_{i\bar{w}}$ برابر با تعداد کل اسنادی که شامل کلمه w نیستند، N_{iw} برابر با تعداد کل اسنادی که در کلاس c_i بوده و شامل کلمه w نیز می باشند و در نهایت $N_{i\bar{w}}$ تعداد کل اسنادی که در کلاس c_i بوده ولی شامل کلمه w نمی باشند را بازنمایی می کند.

۲-۳- مربع کای

این معیار در مباحث آماری کاربرد بسیار متداولی دارد و هدف آن اندازه گیری میزان استقلال دو رویداد از هم می باشد. به طور خاص در حوزه متن کاوی این معیار کمک می کند تا بتوانیم اندازه بگیریم که به چه مقدار رخداد یک کلمه خاص و رخداد یک کلاس خاص از هم مستقل می باشند. مسئله قابل ذکر در مورد این معیار نادقیق بودن آن به دلیل وجود یک درجه آزادی است. این نادقیق بودن باعث می شود تا برخی ویژگی های نامناسب نیز در نظر گرفته شوند که البته آزمایش ها نشان داده است که این میزان ویژگی نویزی تاثیر چندانی روی دقت کلی دسته بند ندارد. رابطه (۴) روش محاسبه مربع کای را نمایش می دهد.

$$\chi^2(w, c_i) = \frac{N \cdot (N_{iw}N_{i\bar{w}} - N_{i\bar{w}}N_{iw})^2}{(N_{iw} + N_{i\bar{w}}) \cdot (N_{i\bar{w}} + N_{iw}) \cdot (N_{iw} + N_{i\bar{w}}) \cdot (N_{i\bar{w}} + N_{iw})} \quad (4)$$

پس از به دست آوردن مقدار بالا به ازای تمامی کلمات، آن ها را با استفاده از رابطه ی زیر رتبه بندی کرده و بهترین آن ها به عنوان ویژگی انتخاب می شوند.

$$\chi^2(w) = \sum_{i=1}^K P(c_i) \cdot \chi^2(w, c_i) \quad (5)$$

در مقاله حاضر سه نوع از این معیارها که در حوزه های مختلف یادگیری ماشین مورد استفاده قرار گرفته است برای متن کاوی و پردازش زبان طبیعی مورد بررسی قرار می گیرد. الگوریتم های مدنظر مقاله حاضر اطلاعات متقابل [8]، بهره اطلاعات [7] و مربع کای [4,5] می باشند. در این راستا ابتدا نحوه کاربردی کردن آن ها در حوزه متن کاوی معرفی می شود. سپس تفاوت خروجی این معیارها در استخراج ویژگی از متون فارسی بررسی می گردد و در نهایت تاثیر استفاده از این معیارها در نتایج دسته بندی متون مطالعه می شود.

۲- معیارهای استخراج ویژگی

۱-۲- اطلاعات متقابل

یکی از متداول ترین معیارهای انتخاب ویژگی اطلاعات متقابل است. این معیار بررسی می کند که چه مقدار اطلاعات با وجود یا عدم وجود یک کلمه برای مدل یادگیری ما منجر به دسته بندی بهتر روی یک کلاس یا دسته می شود. رابطه (۱) روش محاسبه اطلاعات متقابل را نمایش می دهد.

$$MI(w, c_i) = \log \frac{P(w, c_i)}{P(w)P(c_i)} \quad (1)$$

در رابطه فوق w و c_i به ترتیب بیان گر کلمه و کلاس نام می باشند. $MI(w, c_i)$ اطلاعات متقابل آن ها را ارائه می نماید. احتمالات مورد نیاز به صورت زیر محاسبه می گردد.

$$P(c_i) = \frac{N_i}{N} \quad P(w) = \frac{N_w}{N} \quad P(w, c_i) = \frac{N_{iw}}{N}$$

که در آن N برابر با تعداد کل اسناد، N_i برابر با تعداد کل اسناد موجود در کلاس c_i ، N_w برابر با تعداد کل اسنادی که شامل کلمه w هستند و N_{iw} برابر تعداد اسناد کلاس c_i که شامل کلمه w می باشد است. پس از به دست آوردن مقدار بالا به ازای تمامی کلمات، آن ها را با استفاده از رابطه ی زیر رتبه بندی کرده و بهترین آن ها به عنوان ویژگی انتخاب می شوند.

$$MI(w) = \sum_{i=1}^K MI(w, c_i) \quad (2)$$

۲-۲- بهره اطلاعات

این معیار اختلاف آنتروپی را در دو وضعیت وجود یک ویژگی و عدم وجود یک ویژگی بیان می کند. در واقع بهره اطلاعات این مسئله که با وجود یک ویژگی مقدار آنتروپی چه مقدار کاهش می یابد را مورد ارزیابی قرار می دهد و سپس از بین ویژگی ها، آن هایی که بهره اطلاعات آن ها بالاتر است انتخاب می شوند. رابطه (۳) روش محاسبه بهره اطلاعات را نمایش می دهد.

شایان ذکر است با توجه به اینکه کلماتی که تعداد تکرار آن‌ها در کل پیکره بسیار کم است نقشی در کلمات کلیدی نخواهند داشت، می‌توان در یک مرحله پیش‌پردازش قبل از استخراج ویژگی‌ها کلماتی است که به ندرت در متن‌ها ظاهر شده اند را حذف نمود. این کار باعث می‌شود تا سرعت استخراج ویژگی در هنگامی که تعداد داده‌ها زیاد است افزایش

جدول (۱): ویژگی‌های بدست آمده با روش اطلاعات متقابل

domain-score	domain	score	Word
۰.۳۱۹۸۸۹	علمی فرهنگی	۱۳.۴۴۵۶۰۷	رایانه
۰.۳۲۳۷۷۷	اقتصادی	۱۳.۴۰۶۷۵۳	بورس
۰.۳۲۹۲۲۷	ورزش	۱۲.۸۰۰۷۵۰	قهرمانی
۰.۲۳۸۲۲۲	ورزش	۱۲.۶۹۸۶۲۹	گل
۰.۲۸۵۷۵۳	ورزش	۳.۱۴۴۸۲۴	جام
۰.۳۲۴۰۰۹	ورزش	۳.۱۴۰۷۴۹	فوتبال
۰.۲۹۸۲۴۳	علمی فرهنگی	۳.۱۱۷۴۳۹	نرم
۰.۳۲۶۳۹۴	ورزش	۱.۰۷۱۸۲۶	ورزشی
۰.۱۲۶۰۳۱	اقتصادی	۱.۰۵۳۱۶۵	تولید
۰.۳۱۱۶۵۴	علمی فرهنگی	۰.۹۶۹۱۹۷	اینترنت
۰.۲۲۲۴۳۱	اقتصادی	۰.۹۵۰۱۹۹	صنعتی
۰.۲۵۱۳۰۹	اقتصادی	۰.۹۴۳۸۱۳	صنایع
۰.۱۲۷۵۰۳	اقتصادی	۰.۹۲۷۵۳۰	عرضه
۰.۲۵۸۱۳۲	ورزش	۰.۹۱۷۸۲۳	مسابقه
۰.۲۹۲۸۱۷	ورزش	۰.۹۱۴۸۸۱	دیدار
۰.۳۱۱۹۵۸	اقتصادی	۰.۹۰۸۹۴۷	سهم
۰.۲۱۸۷۵۸	علمی فرهنگی	۰.۸۶۳۰۴۶	صفحه
۰.۲۰۹۱۸۵	اقتصادی	۰.۸۶۲۰۵۸	قیمت
۰.۱۱۴۰۷۷	اقتصادی	۰.۸۵۷۴۰۸	دولت
۰.۲۴۶۱۵۹	اقتصادی	۰.۸۲۴۲۸۷	میلیارد

جدول (۲): ویژگی‌های بدست آمده با روش مربع کای

domain-score	domain	score	Word
۵۰۳.۴۸۱۰۹۷	ورزش	۶۳۳.۸۳۹۴۲	ورزشی
۴۲۷.۹۳۲۲۱۸	ورزش	۵۳۷.۲۶۲۵۷۶	تیم
۲۷۷.۰۸۵۵۵۶	ورزش	۳۴۷.۵۶۲۰۳۲	فوتبال
۲۲۷.۷۰۱۳۳۶	علمی فرهنگی	۲۸۸.۲۷۶۲۳۳	رایانه
۲۰۱.۱۸۱۸۰۲	ورزش	۲۷۹.۱۲۵۷۲۲	گروه
۱۹۹.۶۶۸۵۴۵	اقتصادی	۲۵۲.۱۳۴۹۷۵	بورس
۱۹۲.۳۷۸۵۰۷	ورزش	۲۴۰.۵۴۱۰۴۳	جام
۱۸۷.۵۱۰۸۶۳	علمی فرهنگی	۲۳۸.۹۰۹۶۷۱	اینترنت
۱۸۳.۵۷۱۶۳۴	ورزش	۲۳۰.۵۸۲۴۹۲	قهرمانی
۱۸۰.۸۷۸۲۳۲	ادب و هنر	۲۲۹.۰۷۶۶۱۹	هنری
۱۵۴.۴۸۹۷۶۲	اقتصادی	۲۰۷.۲۸۵۷۳۱	سرمایه
۱۵۲.۵۱۲۲۰۸	ادب و هنر	۱۹۷.۳۷۰۹۰۶	فیلم
۱۴۹.۲۴۷۵۱۸	ورزش	۱۸۷.۰۴۳۹۶۳	بازیکن
۱۳۸.۷۳۴۰۲۵	ورزش	۱۸۴.۴۶۸۲۰۰	بازی
۱۴۱.۸۱۱۰۱۱	اقتصادی	۱۸۲.۲۷۸۲۰۵	بانک
۱۳۹.۴۱۰۷۷۰	ورزش	۱۸۲.۰۱۳۳۴۹	دیدار
۱۴۲.۶۸۴۲۷۳	اقتصادی	۱۸۰.۹۶۱۲۷۷	سهم
۱۴۱.۷۵۵۴۵۳	علمی فرهنگی	۱۷۹.۴۷۹۷۱۰	سایت
۱۲۷.۹۰۰۴۳۳	اقتصادی	۱۷۰.۹۴۸۱۹۵	گذاری
۱۳۴.۹۹۸۸۳۸	علمی فرهنگی	۱۷۰.۴۶۸۷۹۴	نرم

۳- دسته‌بندی متون

برای دسته‌بندی متون از الگوریتم بیز ساده^۲ استفاده شده است. رویکرد بیز ساده با در نظر گرفتن استقلال ویژگی‌ها (مدل تک-نگاشتی^۳) عمل می‌کند. به طوری که در محاسبه احتمالات تنها به وجود یا عدم وجود یک ویژگی (کلمه) توجه خواهد داشت. در این روش نحوه تعیین کلاس یک سند به صورت زیر می‌باشد:

$$c = \operatorname{argmax}_{c_i} p(d|c_i)p(c_i) \quad (6)$$

$p(c_i)$ به عنوان احتمال اولیه^۴ و $p(d|c_i)$ به عنوان احتمال شباهت^۵ معرفی می‌شود. در مورد احتمال اولیه می‌توان گفت که این احتمال در مواقعی که داده‌ها به صورت یکنواخت در مجموعه داده بین کلاس‌های مختلف توزیع شده‌اند، تاثیر خاصی روی عمل دسته‌بندی ندارند، اما هنگامی که مجموعه داده به صورت غیر یکنواخت است، تاثیر بسیار زیادی خواهند داشت و هر اندازه که میزان آن به ازای یک دسته بیشتر باشد احتمال انتخاب آن دسته افزایش خواهد یافت. احتمال شباهت نیز براساس رابطه (۷) محاسبه می‌شود.

$$P(d|c_i) = \prod_{w \in d} P(w|c_i) \quad (7)$$

۴- نتایج

۴-۱- انتخاب کلمات ویژگی از متون فارسی

برای بررسی تفاوت میان معیارهای مطرح شده دادگان پیکره همشهری مورد استفاده قرار گرفته‌است. در این پیکره در کنار هر سند موضوع آن سند نیز مشخص شده است [2].

دادگان مورد استفاده برای این پژوهش تعداد ۲۵۸۰ سند است که با یکی از موضوعات ورزش (۵۹۹ سند)، اقتصادی (۵۶۹ سند)، سیاسی (۱۱۷ سند)، ادب و هنر (۵۷۹ سند) و علمی فرهنگی (۵۷۳ سند) و یا محیط زیست (۱۴۳ سند) برچسب‌دهی شده است.

هدف از آزمایش‌های این بخش این است که بر اساس هر یک از سه الگوریتم مذکور یک بردار ویژگی بهینه با حداکثر امتیاز تولید نماییم. جداول ۱-۳ کلمات استخراجی براساس هر یک از معیارهای اطلاعات متقابل، بهره اطلاعات و مربع کای را نمایش می‌دهد. برای این کار ۲۰ کلمه اول این پیکره براساس هر یک از معیارها مرتب شده‌اند. علاوه بر نمایش امتیاز هر کلمه، دامنه‌ای که هر یک از کلمات بیشترین امتیاز را در آن کسب کرده مشخص شده‌است.

خروجی برای هر الگوریتم بصورت ۴ ستون زیر خواهد بود:

- $Word_i$ کلمه انتخاب شده برای بردار ویژگی است.
- $Score$ امتیاز کلی این کلمه در الگوریتم انتخاب ویژگی است.
- $domain$ اسم دامنه‌ای است که کلمه $Word_i$ بیشترین امتیاز را در آن بدست آورده است.
- $domain-score$ امتیاز دامنه‌ای است که کلمه $Word_i$ بیشترین امتیاز را در آن بدست آورده است.

$$\text{Recall}_i = \frac{\#True\ Positives}{\#True\ Positives + \#False\ Negatives} \quad (9)$$

$$F1\text{measure} = \frac{2 * \text{Recall} * \text{Precision}}{\text{Recall} + \text{Precision}} \quad (10)$$

دادگان آموزشی مورد استفاده در این بخش شامل یک پیکره‌ی ۲۵۸۰ سندی و دادگان آزمایشی شامل ۸۶۲ سند می‌باشد. بر این اساس ابتدا الگوریتم بیز ساده بدون استفاده از الگوریتم‌های انتخاب ویژگی روی پیکره‌ی آموزشی پیاده شده‌است. برای این منظور ۲۰۰۰ کلمه پرکاربرد متن را به‌عنوان ویژگی مدنظر قرار داده‌ایم. نتایج این آزمایش در قالب ماتریس آشفتگی^{۱۱} در جدول ۴ نمایش داده شده‌است. در این جدول عنوان هر سطر مشخص کننده برچسب صحیح و عنوان ستون مشخص کننده برچسب انتخاب شده توسط الگوریتم بیز ساده است. از طریق این جدول می‌توانیم تشخیص دهیم که به ازای هر زوج کلاس C_i و C_j، چه تعداد سند از کلاس C_j به طور اشتباه به عنوان C_i برچسب خورده‌اند.

جدول(۴): ماتریس آشفتگی خروجی الگوریتم بیز ساده بدون انتخاب ویژگی

	Desired	Predicted						Recall
		علمی فرهنگی	اقتصادی	محیط زیست	ادب و هنر	سیاسی	ورزش	
	علمی فرهنگی	۱۵۹	۲	۵	۱۴	۵	۶	۰.۸۳
	اقتصادی	۵	۱۶۹	۳	۲	۹	۲	۰.۸۹
	محیط زیست	۰	۰	۴۲	۵	۱	۰	۰.۸۷
	ادب و هنر	۲	۲	۱	۱۸۶	۱	۲	۰.۹۵
	سیاسی	۰	۰	۰	۲	۳۷	۰	۰.۹۴
	ورزش	۳	۰	۱	۶	۳	۱۸۷	۰.۹۳
	precision	۰.۹۴	۰.۹۷	۰.۸۰	۰.۸۶	۰.۶۶	۰.۹۴	

مشاهده می‌شود که به طور کلی الگوریتم بیز با استفاده از ۲۰۰۰ کلمه پرکاربرد عملکرد خوبی داشته‌است، به‌جز دسته سیاسی و محیط زیست که به نسبت صحت پایین‌تری دارند. یکی از علت‌های این موضوع می‌تواند تفاوت توزیع دسته‌های مختلف و در واقع کم بودن تعداد داده‌های آموزشی این دو دسته باشد. مسئله قابل ذکر دیگر میزان پیش‌بینی‌های اشتباه به ازای یک دسته است. در ماتریس آشفتگی می‌بینیم که بیشترین میزان پیش‌بینی اشتباه و در واقع کمترین میزان فراخوانی مربوط به دسته علمی فرهنگی است. این موضوع می‌تواند به این دلیل باشد که کلمات این دسته تا حدی با سایر دسته‌ها همپوشانی دارد و در نتیجه میزان اطلاعاتی که این کلمات منتقل می‌کند در مقایسه با سایر دسته‌ها به میزان کمتری متمایز کننده می‌باشد. در بخش بعدی آزمایش‌ها سه روش مختلف انتخاب ویژگی مورد استفاده قرار گرفته‌اند و در هر یک از آزمایش‌ها ۲۰۰ کلمه اول استخراج شده توسط هر یک از روش‌های بهره اطلاعات، اطلاعات متقابل و مربع کای به‌عنوان ویژگی مورد استفاده قرار گرفته‌اند و سپس همانند بخش اول الگوریتم بیز ساده بر روی دادگان اجرا شده است.

نتایج بدست آمده در قالب ماتریس آشفتگی در جداول ۵-۷ آمده است. در این جداول نیز عنوان هر سطر مشخص کننده برچسب صحیح و عنوان ستون مشخص کننده برچسب انتخاب شده توسط الگوریتم بیز ساده است. همچنین

جدول (۳): ویژگی‌های بدست آمده با روش بهره اطلاعات

domain-score	domain	score	Word
۰.۲۶۱۲۰۸	ورزش	۰.۴۱۵۰۳۳	ورزشی
۰.۲۱۳۰۱۲	ورزش	۰.۳۳۶۶۹۱	تیم
۰.۱۱۶۸۳۴	ورزش	۰.۲۰۴۷۷۰	فوتبال
۰.۱۳۸۵۲۰	ورزش	۰.۲۰۱۷۳۷	گروه
۰.۰۹۶۷۴۹	علمی فرهنگی	۰.۱۷۰۶۵۲	رایانه
۰.۰۸۲۶۴۴	اقتصادی	۰.۱۵۱۰۵۹	بورس
۰.۰۸۶۱۰۰	ورزش	۰.۱۴۲۵۷۳	جام
۰.۰۷۹۵۹۹	علمی فرهنگی	۰.۱۴۰۱۸۲	اینترنت
۰.۰۷۱۸۴۱	ورزش	۰.۱۳۵۹۷۷	قهرمانی
۰.۰۷۶۵۰۹	ادب و هنر	۰.۱۳۱۲۱۶	هنری
۰.۱۰۲۵۷۹	محیط زیست	۰.۱۳۱۱۷۳	زیست
۰.۰۷۸۵۴۳	اقتصادی	۰.۱۲۶۴۶۹	سرمایه
۰.۰۶۶۷۴۸	ادب و هنر	۰.۱۱۹۱۵۸	فیلم
۰.۰۶۸۰۱۲	ورزش	۰.۱۱۴۹۷۵	بازی
۰.۰۶۳۷۳۱	اقتصادی	۰.۱۰۹۳۶۳	بانک
۰.۰۵۸۶۵۸	ورزش	۰.۱۰۹۲۱۴	دیدار
۰.۰۵۶۹۳۹	ورزش	۰.۱۰۸۱۱۴	بازیکن
۰.۰۵۵۸۳۴	علمی فرهنگی	۰.۱۰۵۸۷۷	سایت
۰.۰۵۹۰۴۵	اقتصادی	۰.۱۰۵۳۱۰	سهام
۰.۰۶۳۲۵۲	اقتصادی	۰.۱۰۴۷۶۹	گذاری

یابد و همچنین از لحاظ مصرف حافظه نیز بهبود خواهیم داشت. همانطور که در نتایج ارائه شده مشاهده می‌شود تمامی کلمات استخراج شده توسط سه روش کلمات مطرح و کلیدی در دامنه خود هستند که می‌توانند به خوبی موضوع آن دامنه را در ذهن خواننده متن تداعی نمایند. مقایسه سه روش نشان می‌دهد هر سه روش در استخراج کلمات موفق بودند و تفاوت زیادی در خروجی آن‌ها مشاهده نمی‌شود.

۴-۲- دسته‌بندی متون فارسی با انتخاب ویژگی

در این بخش به بررسی تاثیر انتخاب کلمات ویژگی در کیفیت دسته‌بندی متون فارسی پرداخته می‌شود.

برای این منظور از معیارهای صحت^{۱۲}، فراخوانی^{۱۳} و معیار اف^{۱۴} استفاده شده‌است که فرمول‌های آن‌ها در روابط (۸-۱۰) بیان شده‌اند. در مسئله دسته‌بندی مستندات، فراخوانی بیانگر درصد سندهای کلاس i است که درست دسته‌بندی شده‌اند و صحت بیانگر تعداد اسناد درست دست بندی شده کلاس i به نسبت کل سندهایی است که در دسته i قرار گرفته‌اند. از سوی دیگر به دلیل وجود یک مصالحه بین معیارهای صحت و فراخوانی از معیار اف که میانگین وزن دار همگن صحت و فراخوانی می‌باشد استفاده می‌شود.

$$\text{Precision}_i = \frac{\#True\ Positives}{\#True\ Positives + \#False\ Positives} \quad (8)$$

جدول (۷): ماتریس آشفتنی خروجی الگوریتم بیز ساده با انتخاب ویژگی به روش مربع کای

	Desired	predicted						recall
		علمی فرهنگی	اقتصادی	محیط زیست	ادب و هنر	سیاسی	ورزش	
علمی فرهنگی		۱۷۶	۳	۴	۴	۳	۱	۰.۹۲
اقتصادی		۴	۱۷۶	۵	۰	۵	۰	۰.۹۲
محیط زیست		۰	۱	۴۳	۰	۴	۰	۰.۸۹
ادب و هنر		۲	۰	۰	۱۸۵	۶	۱	۰.۹۵
سیاسی		۰	۰	۲	۳	۳۴	۰	۰.۸۷
ورزش		۰	۰	۲	۰	۰	۱۹۸	۰.۹۹
	precision	۰.۹۶	۰.۹۷	۰.۷۶	۰.۹۶	۰.۶۵	۰.۹۹	

جدول (۸): خلاصه نتایج رویکردهای مختلف انتخاب ویژگی

	بدون ویژگی	بهره اطلاعات	اطلاعات متقابل	مربع کای
Precision	۰.۸۶۱۶	۰.۹۰۹۳	۰.۸۷۵۶	۰.۸۸۶۶
Recall	۰.۹۰۱۶۷	۰.۹۵۳۱۳	۰.۹۲۱۹	۰.۹۲۶۴
F-measure	۰.۸۸۱۲	۰.۹۲۰۷	۰.۸۹۸۲	۰.۹۰۶۱

مراجع

- [1] Deepak Agnihotri, Kesari Verma, Priyanka Tripathi, *Variable Global Feature Selection Scheme for automatic classification of text documents*, Expert Systems with Applications, Volume 81, 2017.
- [2] Abolfazl AleAhmad, Hadi Amiri, Ehsan Darrudi, Masoud Rahgozar, and Farhad Oroumchian. *Hamshahri: A Standard Persian Text Collection*. Knowledge Based Systems, 2:382-387, 2009.
- [3] Huosong, X., & Jian, L. (2011). *The Research of Feature Selection of Text Classification Based on Integrated Learning Algorithm*. 10th International Symposium on Distributed Computing and Applications to Business, Engineering and Science, 20-22. 2011
- [4] Zhou Qian, Zhao Mingsheng, Hu Min. *Study of Feature Selection in Chinese Text Categorization*. Journal of Chinese Information Processing, 2004, 18(3):17-23.
- [5] Chen Tao, Xie Yangqun. *Literature Review of Feature Dimension Reduction in Text Categorization*. Journal of The China Society For Scientific and Technical Information, 24(6):690-695, 2005.
- [6] Alper Kursat Uysal, *An improved global feature selection scheme for text classification*, Expert Systems with Applications, Volume 43, 2016.
- [7] Shang Wenqian, Huang Houkuan, Zhu Haibin, et al. *A novel feature selection algorithm for text categorization*. Expert Systems with Application, 2007, 33:1-5.
- [8] Y Yang, J. Q Pedersen. *Feature selection in statistical learning of text categorization*. In the 14th Int Conf on Machine Learning, 1997:412-420.

زیر نویس ها

¹ Mutual Information

خلاصه‌ای از نتایج کلی هر ۴ حالت آزمایش شده در جدول ۸ نشان داده شده است.

همانطور که در نتایج مشاهده می‌شود هر سه روش پیشنهادی برای استخراج ویژگی باعث بهبود دقت سیستم دسته‌بندی متون شده‌اند. مقایسه سه روش با یکدیگر نشان می‌دهد روش بهره اطلاعات بیشترین بهبود را در نتایج به دست آورده و توانسته در مجموع ۵٪ معیار اف را افزایش دهد. این در حالی است که در اکثر روش‌های متن‌کاوی روش اطلاعات متقابل به عنوان روش استخراج ویژگی مورد توجه قرار گرفته است.

بررسی معیار فراخوانی دسته علمی فرهنگی در هریک از نتایج فوق نشان می‌دهد که نقص این مورد با روش‌های بهره اطلاعات و مربع کای از بین رفته است و شاهد بالا آمدن فراخوانی این دسته هستیم. این امر نشان می‌دهد که این روش‌ها موجب شده‌اند تا ویژگی‌های مناسب‌تری به ازای دسته علمی فرهنگی انتخاب شوند.

۵- نتیجه

در این مقاله به بررسی سه روش مختلف انتخاب ویژگی به نام‌های اطلاعات متقابل، بهره اطلاعات و مربع کای در متن‌کاوی پرداخته شد. نتایج این بررسی نشان داد اگرچه مقایسه کلمات انتخاب شده توسط هر سه روش کیفیت یکسانی را نشان می‌دهد و تفاوت مشهودی در خروجی سه روش مشاهده نمی‌شود، استفاده از هریک از این روش‌ها در دسته‌بندی متون فارسی به نتایج متفاوتی دست‌یافته‌است. در نهایت با استفاده از روش مربع کای برای انتخاب ویژگی دقت دسته بندی ۵٪ در معیار اف افزایش داده شد.

جدول (۵): ماتریس آشفتنی خروجی الگوریتم بیز ساده با انتخاب ویژگی به روش بهره اطلاعات

	Desired	predicted						recall
		علمی فرهنگی	اقتصادی	محیط زیست	ادب و هنر	سیاسی	ورزش	
علمی فرهنگی		۱۷۸	۵	۳	۲	۲	۱	۰.۹۳
اقتصادی		۶	۱۷۷	۳	۰	۴	۰	۰.۹۳
محیط زیست		۰	۰	۴۶	۲	۰	۰	۰.۹۵
ادب و هنر		۳	۰	۰	۱۸۲	۸	۱	۰.۹۳
سیاسی		۰	۰	۰	۱	۳۸	۰	۰.۹۷
ورزش		۰	۰	۲	۰	۱	۱۹۷	۰.۹۸
	precision	۰.۹۵	۰.۹	۰.۸	۰.۹	۰.۷	۰.۹	
		۱	۷	۵	۷	۱	۸	

جدول (۶): ماتریس آشفتنی خروجی الگوریتم بیز ساده با انتخاب ویژگی به روش اطلاعات متقابل

	Desired	predicted						recall
		علمی فرهنگی	اقتصادی	محیط زیست	ادب و هنر	سیاسی	ورزش	
علمی فرهنگی		۱۵۵	۷	۴	۱۲	۱۱	۲	۰.۸۱
اقتصادی		۳	۱۸۱	۳	۰	۳	۰	۰.۹۵
محیط زیست		۱	۱	۴۲	۳	۱	۰	۰.۸۷
ادب و هنر		۵	۰	۱	۱۸۳	۴	۱	۰.۹۴
سیاسی		۰	۰	۰	۱	۳۸	۰	۰.۹۷
ورزش		۰	۰	۲	۱	۲	۱۹۵	۰.۹۷
	Precision	۰.۹۴۵	۰.۹۵	۰.۸۰	۰.۹۱	۰.۶۴	۰.۹۸	

-
- 2 Information Gain
 - 3 Chi-Square
 - 4 Naïve Bayes
 - 5 Unigram
 - 6 Prior probability
 - 7 Likelihood probability
 - 8 Precision
 - 9 Recall
 - 10 F-measure
 - 11 Confusion Matrix

Archive of SID