

## Clustering Web Documents Using Ontology-Based Fuzzy Method

<sup>1</sup> Najmeh Sakhaee, <sup>2</sup>Fariba Salehi, <sup>3</sup>Majid Khalilian

<sup>1</sup> Ms Student of Software Engineer, Islamic Azad University Karaj Branch, Mechatronics College, Iran  
n.sakhaee.star@gmail.com

<sup>2</sup> Faculty member, Islamic Azad University Karaj Branch, Mechatronics College, Iran  
Fariba.salehi@kiauo.ac.ir

<sup>3</sup> Faculty member Islamic Azad University Karaj Branch, Mechatronics College, Iran  
Khalilian@kiauo.ac.ir

### Abstract

Web documents and web pages are expanding rapidly. Web search engines and web services use different methods to find web pages and documents in the massive amount of documents. However, organizing and analyzing a large amount of data is challenging. The problem with web page retrieval is that the information on the global web is in different formats and from different sources. The accuracy of data selection is essential and their compliance with user requests is a challenge in exploring the web. In order to provide an optimal solution for exploring web documents and organizing and providing quick and accurate access to structured and semi-structured Web documents and web pages, a new approach is proposed. The proposed method is based on the clustering and Web document fuzzation and the semantic and structure of web pages. In the proposed method for the reduction of dimension or features, the mapping of attributes to semantic domains is proposed. The results of the implementation of the proposed method in Python and MATLAB software show that the proposed method in categorizing and organizing web documents is appropriate for the quality of clusters and their density, and in the terms of the davies bouldin and silhouette index, they have suitable values.

**Keywords:** Clustering, Web Documents, Mining, Semantic Web.

## خوشه بندی اسناد وب با استفاده از روش فازی آنتولوژی محور

نجمه سخایی<sup>۱</sup>، فریبا صالحی<sup>۲</sup>، مجید خلیلیان<sup>۳</sup>

<sup>۱</sup> دانشجوی کارشناسی ارشد مهندسی نرم افزار، دانشکده مکترونیک، دانشگاه آزاد اسلامی واحد کرج، ایران  
n.sakhaee.star@gmail.com

<sup>۲</sup> عضو هیات علمی، دانشکده مکترونیک، دانشگاه آزاد اسلامی واحد کرج، ایران  
Fariba.salehi@kia.ac.ir

<sup>۳</sup> عضو هیات علمی، دانشکده مکترونیک، دانشگاه آزاد اسلامی واحد کرج، ایران  
Khalilian@kia.ac.ir

### چکیده

اسناد و صفحات وب در اینترنت به سرعت در حال گسترش هستند. موتورهای جستجو و خدمات رسان های وب برای یافتن صفحات وب و اسناد مورد نظر در میان حجم انبوهی از اسناد، از روش های مختلف استفاده می کنند. با این وجود سازمان دهی و تحلیل حجم وسیعی از داده ها چالش برانگیز است. مشکل مطرح در زمینه باز یابی صفحات وب، این است که اطلاعات موجود در وب وسیع جهانی در فرمت های مختلف و از منابع مختلف می باشند. صحت انتخاب داده ها ضروری بوده و تطابق آنها با درخواست کاربران به عنوان چالشی در کاوش وب می باشد. به منظور ارائه راه حلی بهینه برای کاوش در میان اسناد وب و سازمان دهی و دسترسی سریع و صحیح به اسناد و صفحات وب ساخت یافته و نیمه ساخت یافته در این تحقیق روشی جدید پیشنهاد شده است. روش پیشنهادی بر اساس خوشه بندی و فازی سازی اسناد وب و با توجه به معنا و ساختار صفحات وب می باشد. در روش پیشنهادی برای کاهش بعد یا ویژگی ها، نگاشت ویژگی ها به حوزه های معنایی پیشنهاد شده است. نتایج حاصل از پیاده سازی روش پیشنهادی در نرم افزار پایتون و متلب نشان می دهد روش پیشنهادی در دسته بندی و سازماندهی اسناد وب، از نظر کیفیت خوشه ها و تراکم آنها مناسب بوده و از نظر شاخص دیویس بولدین و سیلهوئت دارای مقادیر مناسبی می باشد.

### کلمات کلیدی

خوشه بندی، اسناد وب، کاوش، وب معنایی.

### ۱- مقدمه

لایه پنهان وب دور از دسترس بوده و قابل استفاده نمی باشد. مشکل مطرح در زمینه باز یابی صفحات وب، این است که اطلاعات موجود در وب وسیع جهانی در فرمت های مختلف و از منابع مختلف می باشند [2]. کنسرسیوم W3 اظهار داشته است که HTML، به علت محدودیت هایی همچون داده های نیمه ساخت یافته، حساسیت به حروف، تگ های از پیش تعریف شده و مسائلی از این دست، ساختار معنایی خوبی از محتویات صفحات وب توصیف نمی کند. بعدها برای غلبه بر این مشکلات تکنولوژی هایی همچون XML، Flash (با ویژگی های خوب طراحی) و مانند آن پدید آمد [3]. امروزه داده های نیمه ساخت یافته در وب افزایش یافته است و به دنبال آن استفاده از

امروزه اسناد و صفحات وب در اینترنت به سرعت در حال گسترش هستند. کاربران با استفاده از اینترنت و بسیاری از موتورهای جستجو می توانند اسناد و صفحات وب مورد نظرشان را پیدا کنند. سازمان دهی و تحلیل حجم وسیعی از داده ها چالش برانگیز بوده و گاهی غیرممکن می باشد [1]. موتورهای جستجو و خدمات رسان های وب، برای یافتن صفحات وب و اسناد مورد نظر در میان حجم انبوهی از اسناد، از روش های مختلف استفاده می کنند. با این وجود هنوز بخش عظیمی از وب توسط این خدمات رسان ها قابل کشف نبوده و به صورت



حاصله کافی برای مطالعه این چند میلیون پاسخ را نخواهند داشت و احتمالاً فقط پاسخ‌هایی را که در رتبه‌بندی‌های بالا قرار دارد مطالعه نمایند. در سال‌های اخیر توجه زیادی به سمت خوشه‌بندی معطوف شده است. خوشه‌بندی را می‌توان در رابطه با سازماندهی اسناد استفاده کرد تا بتوان آن‌ها را به‌طور خودکار در دسته‌های معنی‌دار تقسیم کرد. خوشه‌بندی یکی از تکنیک‌هایی است که کاربرد گسترده‌ای در تجزیه و تحلیل داده‌های کاوشی دارد چراکه به‌واسطه آن می‌توان ساختار طبیعی داده‌ها را معین نمود. در مقایسه با یادگیری نظارتی، خوشه‌بندی یک روش یادگیری نظارت‌نشده است از این‌رو هیچ نیازی به داده‌های برچسب دار در مدت پروسه خوشه‌بندی نیست؛ هدف خوشه‌بندی، گروه‌بندی اشیاء است، بطوریکه اشیاء در یک خوشه به همدیگر شبیه‌ترند تا با اشیاء دیگر خوشه‌ها [5].

برای مرور مؤثر و جستجو در مجموعه عظیمی از اسناد وب و سازمان‌دهی آن‌ها نیازمند استفاده از روش‌های داده‌کاوی هستیم. روش‌هایی که در پاسخ به درخواست‌های کاربران اسناد و صفحاتی را ارائه نمایند که بیشترین ارتباط را با درخواست کاربر داشته باشد. در این تحقیق سعی بر ارائه راه‌حلی بهینه برای کاوش در میان اسناد وب و سازمان‌دهی و دسترسی سریع و صحیح به اسناد و صفحات وب با استفاده از خوشه‌بندی و فازی سازی شده است. در ادامه پس از مرور کارهای پیشین در بخش ۳ روش پیشنهادی توضیح داده می‌شود و در بخش ۴ پیاده‌سازی و ارزیابی نتایج شرح داده می‌شود و در پایان در بخش ۵ به جمع‌بندی مطالب بیان شده و در بخش ۶ پیشنهاد کارهای آتی می‌پردازیم.

## ۲- مرور کارهای پیشین

یادگیری نیمه نظارت شده، یادگیری با داده‌های برچسب دار و بدون برچسب، اخیراً توسط محققان زیادی مورد مطالعه واقع شده است. انواع الگوریتم‌های نیمه نظارت شده ای پیشنهاد شده اند، شامل آموزش همزمان [6]، بیز ساده ی نیمه نظارت شده [7]، ماشین‌های بردار پشتیبانی [8] TSVM، مدل خوشه بندی فازی [9] و روش‌های مبتنی بر گراف می‌باشند.

لیو و همکاران [10] یک الگوریتم جدید قطعه بندی صفحات وب بر اساس یافته‌های درخت گوموری هیو در گراف برنامه ریزی ارائه نمودند. الگوریتم ابتدا اطلاعات دیداری و ساختاری را از یک صفحه وب جدا کرده و در یک گراف وزندار بدون جهت قرار می‌داد به‌طوری‌که رؤس این گراف گره‌های برگ درخت DOM بوده و لبه‌ها یک رابطه وضعیت قابل رویت بین رؤس بودند. در نهایت گراف با الگوریتم خوشه بندی بر اساس درخت گرومی هیو پارتیشن بندی می‌شد.

در سال ۲۰۱۴ وانگ و همکاران [6] از نوعی خوشه بندی آنلاین برای خوشه بندی اطلاعات مربوط به ترافیک شبکه استفاده نمودند. لین و همکاران [11] نوعی خوشه بندی بازگشتی K-mean با استفاده از ماشین شاخص گذاری ناخالصی برای جاگذینی خوشه بندی K-mean معرفی کردند که در واقع ترکیبی از مرکز خوشه و نزدیکترین همسایگی بود و k نزدیکترین همسایگی را به صورت محلی انجام می‌داد.

در سال ۲۰۱۶ سو و همکاران [12] روش جدیدی از خوشه بندی نظارت نشده با عنوان خوشه بندی حلقه‌ای، ارائه دادند. در روش پیشنهادی آنها برای هر خوشه یک سرخوشه تعریف می‌شود که در آن سرخوشه اطلاعات مربوط به اعضای خوشه و فرمت خوشه ذکر می‌گردد.

XML به‌عنوان استاندارد برای داده‌های نیمه ساخت‌یافته، گسترش یافته است. محققان این حوزه برای بهبود صفحات وب، شروع به مهاجرت به صفحات وب XML ای نمودند و از تکنولوژی ادغام صفحات وب با معانی و محتوای صفحه برای فراهم‌سازی توصیف بهتر ساختار معنایی و توصیف محتوا استفاده کردند.

از جمله مشکلاتی که در زمینه طبقه بندی اسناد وب وجود دارد کشف کردن دانش مفید از متن نیمه ساخت‌یافته یا غیرساخت‌یافته است که توجه زیادی را به خود جلب کرده است. روش‌های داده‌کاوی سنتی فرض می‌کنند که اطلاعات به فرم پایگاه داده‌های رابطه‌ای هستند به همین دلیل برای بسیاری از کاربردها مانند اطلاعات الکترونیکی قابل دسترس به فرم نیمه ساخت‌یافته یا غیرساخت‌یافته مفید نیستند. داده‌های وب، داده‌های نیمه ساخت‌یافته هستند چون نه به‌طور کامل غیرساخت‌یافته هستند و نه به‌طور کامل ساخت‌یافته هستند.

وب معنایی خدمات وبی هوشمندتری را برای هماهنگی و مرتب نمودن داده‌های روی وب عرضه نموده است. در سال‌های اخیر صحت انتخاب داده‌های ضروری مطابق با درخواست کاربر و صدور آن‌ها در خروجی به‌عنوان چالشی در کاوش وب مطرح شده است. تحلیل و کاوش داده‌های XML از طرف انجمن کشف دانش و داده‌کاوی در سال‌های اخیر به محبوبیت رسیده است. مرور مؤثر صفحات وب و جستجوی مجموعه عظیمی از این اسناد نیازمند سازماندهی و استفاده از روش‌های داده‌کاوی می‌باشد. رشد داده‌های نیمه ساخت‌یافته و گستره استفاده از XML به‌عنوان استاندارد برای داده‌های نیمه ساخت‌یافته دلیل اصلی این نیاز می‌باشد [4]. وب معنایی می‌تواند به‌صورت کارا و مؤثر، منابع مختلف را بر اساس ارتباطات معنایی سازمان‌دهی، اشتراک‌گذاری، دسته‌بندی، ترکیب و مدیریت کند. ارتباطات معنایی، مطالعه مفاهیم در محیط‌های ارتباطی برای پشتیبانی از کاربردهای هوشمند در ارتباطات منابع در محیط‌های معنایی است، به‌طوری‌که ماشین‌ها و انسان‌ها توانایی درک یکدیگر را داشته باشند. زبان‌های معنایی توصیف دقیق نهادهای تعاملی مثل سرویس‌ها، منابع و کاربران در سطوح بالای تجرید را ممکن ساخته‌اند و نتیجه‌گیری خودکار در مورد این نوع ارائه را فراهم می‌سازند، بنابراین به‌انگیزش و درک دوطرفه میان نهادهای کم اشتراک یا دانش بدون اولویت در مورد همدیگر توجه می‌کنند. در وب معنایی از آنتولوژی استفاده می‌شود. یک آنتولوژی، توصیفی صریح و فرمال از مفاهیم یک دامنه خاص است (کلاس‌ها که گاهی اوقات مفاهیم نیز نامیده می‌شود)، ویژگی‌های هر یک از این مفاهیم و صفات خاص و خصایص مختلف این مفهوم را توصیف می‌کنند. آنتولوژی به‌منزله یک توصیف فرمال و صریح از واژه‌های یک دامنه خاص می‌باشد و ارتباطات میان آنتولوژی‌ها در وب یک دامنه وسیع را پوشش می‌دهند. هدف این ساختارها تسهیل تعامل عامل‌ها در وب است.

انباره‌های مختلفی وجود دارد که شامل حجم عظیمی از خدمات وب می‌باشد. یافتن روشی که خدمات موردنظر کاربر را از میان این انباره‌ها بیابد و خدمات یافت شده بیشترین ارتباط را به درخواست کاربر داشته باشد به مسئله‌ای چالش‌برانگیز تبدیل شده است. خدمت‌رسان‌های وب اغلب با حجم زیادی از خدمات وب روبرو می‌باشند که مشابه با درخواست کاربر می‌باشد. همان‌طور که اکثراً مشاهده شده است برای هر عمل جستجو در موتورهای جستجو معمولاً میلیون‌ها وب‌سایت و سند و تصویر مرتبط با درخواست کاربر در عرض کمتر از یک ثانیه به کاربر ارائه می‌شود درحالی‌که مسلماً کاربران فرصت و



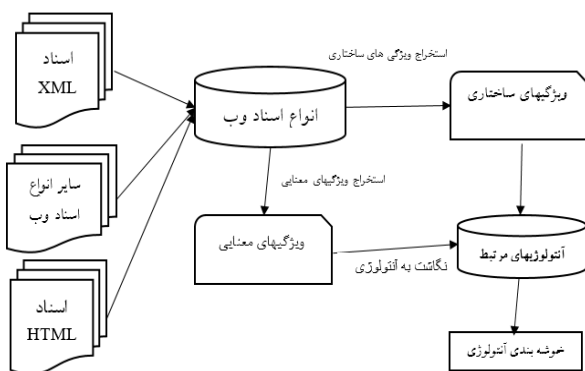
شکل(۱): مراحل روش پیشنهادی

در شکل ۲ معماری روش پیشنهادی نشان داده شده است. در ادامه هر یک از مراحل روش پیشنهادی شرح داده می شود.

جمع آوری اسناد و صفحات وب: جمع آوری اسناد و صفحات وب اولین گام برای کاوش صفحات وب می باشد. در این گام فایلها و اسناد وب به صورت اتوماتیک جمع آوری می شوند و در یک پایگاه داده ای ذخیره می شوند.

استخراج ویژگی های ساختاری: استفاده از تگها متداول ترین راه برای استخراج ویژگی های اسناد وب می باشد. در این روش متون بدون توجه به نام تگها به مفاهیم معنایی نگاشت می شوند. سپس مفاهیم استخراج شده از فضای آنتولوژی به نام تگها افزوده می شود. ویژگی های استخراج شده از نام تگها می توانند با بهره گیری از آنتولوژی لغوی به دانش برسند.

پاک سازی اسناد: پاک سازی اسناد وب شامل حذف تگها و کلماتی است که در لیست STOPWORD می باشد. یک لیست STOPWORD شامل کلماتی است که به صورت متناوب در زبان انگلیسی استفاده می شوند و معنی خاصی ندارند مانند "The".



شکل(۲): معماری روش پیشنهادی

استخراج ویژگی های محتوایی و کاهش بعد: داده های دارای بعد زیاد چالش های ریاضیاتی زیادی را در وظایف یادگیری ماشین پیش روی ما می گذارند. یکی از مسائل در رابطه با مجموعه داده های دارای بعد بالا این است که همه متغیرهای اندازه گیری شده برای درک پدیده مورد نظر حائز اهمیت هستند. در کاهش ابعاد سعی بر این است که بعد داده ها طبق یک تعداد از معیارها پیدا شود. که خود حوزه تحقیقی فعالی در یادگیری ماشین است. برای کاهش بعد پیشنهاد ما نگاشت آنتولوژی می باشد. از روی ویژگی های استخراج شده می توان آنتولوژی های مربوطه را مشخص نمود و

در سال ۲۰۱۶ ایوان و همکاران [13] روش خوشه بندی دسته ای داده های فازی با استفاده از رابطه غیر قابل تشخیص را ارائه دادند. آنها یک الگوریتم جایگزین با تصحیح روش K پارتیشن فازی برای خوشه بندی دسته ای داده های فازی ارائه دادند. الگوریتم پیشنهادی با استفاده از رابطه غیر قابل تشخیص میان داده ها خوشه بندی را انجام می داد، نتایج نشان می دهد که روش پیشنهادی در مقایسه با الگوریتم های قبلی با پایه K-means خوشه بندی خالص تری انجام می دهد و نیز زمان پاسخ کوتاهتری دارد.

در سال ۲۰۱۶ لیانگ ژانگ و همکاران در [14] از FCM برای دسته بندی داده های از دست رفته و نیمه کامل استفاده کردند. در این مقاله نویسندگان خوشه بندی جدیدی بر اساس FCM برای خوشه بندی داده های نیمه کامل با استفاده از مقادیر احتمالاتی هسته اطلاعات و معیار بیشترین شباهت ارائه کردند و در گام بعدی مدل خوشه بندی شان را با استفاده از بهینه ساز سه سطحی جایگزینی و ضریب لاگرانژ بهبود دادند. صحت نتایج به وسیله شاخص میانگین خطای اولیه، با ارجاع به دو شاخص به نام های طبقه بندی خطا و نسبت دهی خطا ارزیابی می شود. نتایج همگرایی پایدار، صحت بالا و قدرتمند در روش پیشنهادی را نشان می دهد.

با وجود تحقیق های فراوانی که برای سازمان دهی و تحلیل حجم وسیعی از داده ها انجام شده است، اما هنوز بخشهایی از وب به عنوان لایه پنهان دور از دسترس باقی است. مشکل مطرح در زمینه ارزیابی صفحات وب، این است که اطلاعات موجود در وب وسیع جهانی در فرمت های مختلف و از منابع مختلف می باشند. صحت انتخاب داده ها ضروری بوده و تطابق آنها با درخواست کاربران به عنوان چالشی در کاوش وب می باشد. به منظور ارائه راه حل بهینه برای کاوش در میان اسناد وب و سازمان دهی و دسترسی سریع و صحیح به اسناد و صفحات وب ساخت یافته و نیمه ساخت یافته در این تحقیق روشی جدید پیشنهاد شده است، که در ادامه به شرح آن می پردازیم.

### ۳- روش پیشنهادی

به منظور ارائه راه حل بهینه برای کاوش در میان اسناد وب و سازمان دهی و دسترسی سریع و صحیح به اسناد و صفحات وب ساخت یافته و نیمه ساخت یافته در این تحقیق روشی جدید با استفاده از خوشه بندی و فازی سازی ارائه شده است که در آن علاوه بر ویژگی های ساختاری، ویژگی های محتوایی اسناد نیز مورد کاوش قرار می گیرد. فرایند اصلی روش پیشنهادی مطابق شکل ۱ از پنج گام به شرح زیر تشکیل شده است:

۱. جمع آوری اسناد و صفحات وب
۲. استخراج ویژگی های ساختاری
۳. پاک سازی اسناد
۴. استخراج ویژگی های محتوایی و کاهش بعد
۵. خوشه بندی و فازی سازی



means به صورت موضعی فاصله مربع میانگین بین اشیاء و مراکز خوشه‌های را به حداقل می‌رساند C-Means. فازی تابع هدف مشابهی دارد ولی K-means را گسترش داده و درجه اطلاعات عضویت را شامل می‌شود که حاکی از وجود اطمینان در تخصیصی به خوشه موردنظر است. این دو الگوریتم خوشه‌بندی را می‌توان با بهینه سازی توابع هزینه، حاصل نمود.

K-means را می‌توان با استفاده از الگوریتم EM روی ترکیبی از گاوسی‌های C تحت فرضیات معینی مدل کرد، که در آن C تعداد خوشه‌های K-means است. از آنجاییکه الگوریتم K-means تخصیص سختی از نقاط داده‌ای را در خوشه‌ها اجرا می‌کند که در آن هر نقطه داده‌ای بطور منحصر بفردی با یک خوشه مرتبط است الگوریتم EM تخصیص نرمی را بر اساس احتمالات پیشین صورت می‌دهد. مدل ترکیبی گاوسی را با اجزای C در نظر بگیرید که در آن میانگین این اجزا  $\mu_1, \dots, \mu_C$  هستند و ماتریس‌های کوواریانس معمول اجزای ترکیبی با  $\Sigma = \mathcal{E}I$  می‌شوند که در آن I ماتریس یکسانی را نشان می‌دهد و  $\mathcal{E}$  پارامتر مشترکی با همه اجزا است. با  $\mathcal{E} \rightarrow 0$  احتمال داده‌های کامل را می‌توان به شکل رابطه ۳ نوشت.

$$E_{Z|X, \theta^{old}} [\ln P(X, Z | \theta)] \quad (3)$$

$$= \frac{1}{2} \sum_{i=1}^N \sum_{c=1}^C P(Z_c | X_i; \theta) \|X_i - \mu_c\|^2 + const$$

از اینرو با به ماکزیمم رسانی احتمال داده‌های موردانتظار بالا می‌توان به نتیجه رسید. در واقع هر داده به یک خوشه تخصیص داده می‌شود همانطوریه در رابطه ۴ می‌بینیم:

$$P(Z_c | X_i; \theta) \quad (4)$$

$$= \begin{cases} 1 & \text{if } c = \operatorname{argmin} \|X_i - \mu_c\|^2 \\ 0 & \text{otherwise} \end{cases}$$

اساساً تابع احتمال پیشین به تخصیص سخت محدود نمی‌شود همانطوریکه در رابطه ۴ می‌بینیم دیگر تخصیص‌های سخت را نیز می‌توان به خوبی استفاده کرد. ارتباط بین K-means و EM توصیفی در بالا منجر می‌شود که یک الگوریتم K-means نیمه نظارت شده پیشنهاد دهیم. تابع احتمال پیشین با یک تابع عضویت فازی دارای محدودیت جایگزین می‌شود. سپس می‌توان عضویت فازی را داخل مرحله E الگوریتم EM جاگذاری کرد. در این مطالعه از ماتریس U استفاده شده و اطلاعات عضویت فازی نگه داشته می‌شود که در آن Uic درجه عضویت سند i در خوشه c را نشان می‌دهد. عبارت دیگر Uic برای نشان دادن احتمال پیشین  $P(Z_c | X_i)$  استفاده می‌شود. در مرحله M، Uic استفاده شده و پارامترهای مدل جدید محاسبه می‌شود که مراکز خوشه‌های  $\mu_1, \dots, \mu_C$  از ماکزیمم سازی تابع احتمال لاگ هستند.

علاوه بر این تعداد کمی مثال‌های برچسب دار موجود است. در نیمه-Kmeans فازی از مثال‌های برچسب دار برای دو هدف اصلی استفاده می‌شود. اول اینکه این مثال‌های برچسب دار می‌توانند موجب حدس اولیه مراکز شوند. عملاً الگوریتم‌های خوشه‌بندی نظیر K-means و C-Means فازی به حدس اولیه حساس هستند. از اینرو بزرها می‌توانند حدس بهتری برای الگوریتم حاصل کرده و اسناد را خوشه‌بندی کنیم. دوم اینکه

در ادامه کار به جای استفاده از ویژگی‌های زیاد از حوزه مربوطه استفاده می‌کنیم. برای نگاشت آنتولوژی از الگوریتم پیشنهادی زیر استفاده می‌شود: در الگوریتم نگاشت آنتولوژی بر اساس کلمات کلیدی، کلمات کلیدی به عنوان ورودیهای الگوریتم می‌باشد و پس از اجرای الگوریتم، کلمات کلیدی به حوزه‌های معنایی مربوطه نگاشت می‌یابند و خروجی الگوریتم آنتولوژیها می‌باشد. مثلاً اگر در سندی کلمات کلیدی شامل توپ، دروازه، شوط، گل، و غیره باشد، حوزه معنایی ورزش به همراه زیر خوشه‌های فوتبال، والیبال، بسکتبال و موارد مشابه به عنوان آنتولوژی نگاشت خواهد شد.

ورودی: کلمات کلیدی

خروجی: آنتولوژی

۱. شروع
۲. براساس وزن کلمات کلیدی آنها را مرتب کن.
۳. با توجه به وزن کلمات حوزه‌های مربوطه را مشخص کن
۴. برای هر حوزه وزن هر سند را محاسبه کن
۵. برای هر زیر مجموعه از حوزه وزن سند را محاسبه کن
۶. حوزه‌های مربوطه را استخراج کن
۷. پایان

خوشه‌بندی و فازی سازی: برای خوشه‌بندی از الگوریتم نیمه K-means فازی استفاده می‌شود [5]. این الگوریتم نوعی از الگوریتم K-means است ولی در آن از تکنیک EM برای به کارگیری یک تابع عضویت فازی استفاده شده و این امکان حاصل می‌شود که اشیاء به بیش از یک خوشه تعلق داشته باشند.

هدف اکثر الگوریتم‌های خوشه‌بندی، حداقل سازی تابع هزینه است که شامل سنجش تفاوتها بین اشیاء و نماینده خوشه می‌باشد. K-means به صورت موضعی فاصله مربع میانگین بین اشیاء و مراکز خوشه‌ای را به حداقل می‌رساند.

EM تکنیک متداول دیگری برای تجزیه و تحلیل الگوریتم‌های خوشه‌بندی است چراکه یک روش فرمول بندی شده آماری می‌باشد و اطلاعات مفصل تری راجع به نتیجه خوشه‌بندی می‌دهد. الگوریتم EM یک روش کلی برای یافتن راه حل‌های احتمال ماکزیمم در مدل‌های دارای متغیر نهان است. اگر مجموعه داده‌های مشاهده شده  $x = \{x_1, \dots, x_n\}$  باشد مجموعه پارامترهای مدل با  $\theta$  نشان داده می‌شود و مجموعه همه متغیرها نهان Z است، در مرحله E مقادیر پارامتر فعلی  $\theta^{old}$  را به کار برده و توزیع پیشین متغیرهای نهان را محاسبه می‌کنیم که با  $p(Z|X, \theta^{old})$  نشان داده می‌شود. سپس توزیع پیشین برای محاسبه مقدار امید ریاضی احتمال استفاده شده و مقدار پارامتری جدید  $\theta^{new}$  حاصل می‌شود. مقدار امید ریاضی احتمال داده‌های کامل در توزیع پیشین متغیرهای نهان با  $q(\theta, \theta^{old})$  نشان داده می‌شود که در رابطه ۱ می‌بینیم. در مرحله M پارامتر جدید  $\theta^{new}$  را با ماکزیمم سازی این تابع حاصل می‌کنیم که در رابطه ۲ مشاهده می‌شود.

$$q(\theta, \theta^{old}) \quad (1)$$

$$= E_{Z|X, \theta^{old}} [\ln P(X, Z | \theta)] \quad (2)$$

$$\theta^{new} = \operatorname{argmax}_{\theta} q(\theta, \theta^{old})$$

در بسیاری از الگوریتم‌های خوشه‌بندی هدف به حداقل رسانی تابع هزینه است که شامل سنجش اغتشاش بین اشیاء و نماینده خوشه می‌شود-K

بین  $H$  و  $\tilde{H}$  این است که هر سند در  $H$  با یک بردار ترنم نشان داده می‌شود؛ درحالی‌که هر ردیف در  $\tilde{H}$  بردار موضوعی نشان داده می‌شود. رده‌بندی اولیه اسناد برچسب دار هستند و این امکان را می‌دهند که از رده‌بندی‌های  $S_1, \dots, S_C$  برای محاسبه مراکز خوشه‌های اولیه تحت عناوین  $\mu_1, \dots, \mu_C$  استفاده کنیم. پروسه‌های بالا در خطوط ۲ الی ۴ الگوریتم نیمه K-means-فازی لیست می‌شوند.

وقتی که پارامترهای  $\mu_1, \dots, \mu_C$  حاصل شدند درجه عضویت  $U_{ic}$  را می‌توان با استفاده از تابع توزیع گاوسی با اندازه‌گیری فاصله محاسبه کرد. جدای از این هرگونه اطلاعات خوشه‌بندی رده‌ای شناخته می‌شود. سپس ماتریس درجه عضویت  $U$  نرمالیزه می‌شود. پروسه‌های بالا در خطوط ۶ الی ۱۷ الگوریتم نیمه K-means-فازی لیست می‌شوند.

وقتی که ماتریس درجه عضویت  $U$  تغییر می‌کند الگوریتم باید از احتمال پیشین جدید متغیر نهان برای محاسبه پارامترهای جدید استفاده کند که مراکز خوشه‌ای با اسناد برچسب دار و بدون برچسب هستند. در این مطالعه از بردار نرمالیزه سازی  $Z$  استفاده شده و فاکتور نرمالیزاسیون نشان داده می‌شود که در آن هر  $Z_c$  مجموع درجات عضویت اسناد در خوشه  $c$  را نشان می‌دهد. سپس الگوریتم می‌تواند مراکز خوشه‌ای جدا را محاسبه کند. پروسه بالا در خطوط ۱۸ الی ۲۳ الگوریتم نیمه K-means-فازی لیست می‌شوند. وقتی که الگوریتم همگرا است ماتریس  $U$  درجه عضویت را در خروجی می‌دهد.

علاوه بر تابع توزین گاوسی تشابه کسینوسی را می‌توان به خوبی استفاده کرد. در بازیابی اطلاعات (IR) یا پردازش زبان طبیعی (NLP)، مدل فضای برداری اغلب برای نشان دادن اسناد استفاده می‌شود که در آن هر سند به‌عنوان برداری نشان داده شده و هر بُعد متناظر با یک‌ترم متمایز است. تشابه کسینوسی اغلب برای اندازه‌گیری فاصله در IR یا NLP استفاده می‌شود. نتیجه تشابه کسینوسی عددی است مابین ۰ و ۱. جدای از این تابع عضویتی که در الگوریتم استفاده می‌شود باید نرمالیزه شود چراکه از توزیع احتمال پیشین حاصل می‌شود.

#### ۴- پیاده سازی و ارزیابی نتایج

به‌منظور ارائه راه‌حلی بهینه برای کاوش در میان اسناد وب و سازمان‌دهی و دسترسی سریع و صحیح به اسناد و صفحات وب ساخت‌یافته و نیمه ساخت‌یافته در این تحقیق روشی جدید با استفاده از خوشه‌بندی و فازی سازی ارائه شده است که در آن علاوه بر ویژگی‌های ساختاری، ویژگی‌های محتوایی و معنایی اسناد نیز مورد کاوش قرار می‌گیرد. داده‌های اولیه مورد استفاده در این تحقیق ۴۵ سند خبری از سایت پرس تی وی می‌باشد. ۱۵ خبر اقتصادی، ۱۵ خبر ورزشی و ۱۵ خبر سیاسی به صورت تصادفی از این سایت به عنوان داده‌های اولیه انتخاب شده‌اند. برای پیاده سازی روش پیشنهادی از نرم افزار پایتون و متلب استفاده شده است. برای استخراج ویژگی‌های ساختاری و پاکسازی اسناد از نرم افزار پایتون و از کتابخانه NLTK استفاده شده است. برای استخراج ویژگی‌های مفهومی و کاهش بعد از named entity recognition در NLTK استفاده شده است. ابتدا ۶ گروه، کلمه هدف با وزن یکسان تعریف کرده ایم سپس کلمات موجود به گروه مربوطه نگاشت شده است. جدول ۱ ویژگی‌ها و خروجی این مرحله را نشان می‌دهد. در ستون کمترین استفاده، تعداد

مثال‌های برچسب دار می‌توانند خوشه‌بندی را به سمت فضای جستجوی بهتری سوق دهند.

1. Begin
2.  $H_i$  ( $\frac{H_i}{\|H_i\|}$  where  $H_i$  is the  $i$ th row of  $H$  and  $i=1 \dots N$ )
3.  $\hat{H}$ (PLSA\_Clustering( $H, K$ ))
4.  $\mu_C$  ( $\frac{1}{|S_c|} \sum_{d_i \in S_c} d_i$  where  $\mu_C$  is the center of  $c$ th cluster and  $c=1 \dots C$ )
5. repeat
6. for  $i=1$  to  $N$  do
7. for  $c=1$  to  $C$  do
8. If  $d_i$  is a document of  $S_c$  then
9.  $U_{ic} \leftarrow 1$
10.  $U_{ic'} \leftarrow 0$ . For  $c'=1 \dots C$  and  $c' \neq c$
11. Break
12. Else
13.  $U_{ic}(e^{-\|H_i - \mu_C\|^2 / 2\sigma^2})$
14. end
15. end
16. Normalize  $U_{ic}$  so that the sum of each row of  $U$  is 1.
17. end
18. For  $c=1$  to  $C$  do
19.  $\mu_C(\frac{1}{Z_c} \sum_{i=1}^N U_{ic} \times H_i)$
20. end
21. until convergence
22. return  $U$
23. end

#### شکل ۳ شبیه‌سازی الگوریتم نیمه K-means فازی

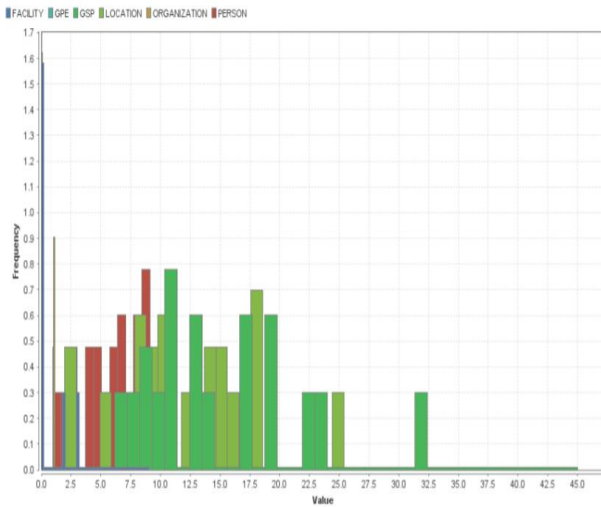
در روش پیشنهادی برای ورودی الگوریتم به جای ماتریس ترنم-سند از ماتریس آنتولوژی -سند استفاده می‌شود. ورودی‌های الگوریتم  $H$  ماتریس آنتولوژی-سند،  $K$  تعداد موضوعات،  $C$  تعداد خوشه‌ها و  $S_1, \dots, S_C$  هسته‌های خوشه‌های ۱ تا  $C$ . و خروجی  $U$  ماتریس عضویت سند-خوشه می‌باشد.

شکل ۳ شبیه‌سازی الگوریتم نیمه K-means فازی را نشان می‌دهد. تابع عضویت فازی مورد استفاده در این الگوریتم، تابع توزیع گاوسی هست که در معادله ۷ می‌بینیم که در آن  $\mu_C$  موقعیت مرکز و  $\sigma$  برای کنترل درجه عضویت  $X_i$  در خوشه  $c$  و  $|X_i - \mu_C|$  فاصله بین  $X_i$  و  $\mu_C$  است.

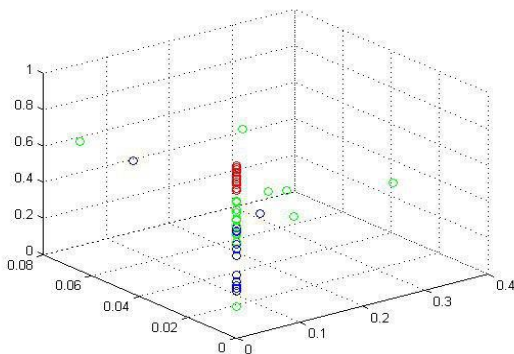
$$U_{ic} = e^{-\|X_i - \mu_C\|^2 / 2\sigma^2} \quad (5)$$

نقاط نزدیک به مرکز مهم خواهند بود و نقاط دور از مرکز نسبتاً کم‌اهمیت هستند. فاصله با استفاده از تعداد درجه فازی نشان داده می‌شود چراکه مقدار  $U_{ic}$  مقداری است مابین ۰ و ۱. جدای از این مراکز خوشه‌ای  $\mu_1, \dots, \mu_C$  پارامتر  $\Theta$  مدل هستند و به صورت مکرر آپدیت می‌شوند و ورودی‌های الگوریتم نیمه K-means فازی شامل ماتریس آنتولوژی -

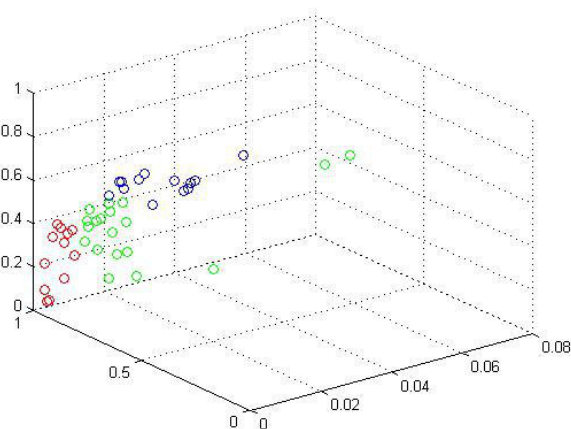
سند  $N \times M$ ، تعداد خوشه‌های  $C$ ، تعداد موضوعات  $K$  و رده‌بندی سندهای  $S_1, \dots, S_C$  می‌شوند. هر ردیف  $H$  سندی را نشان می‌دهد و هر ستون خصیصه ترمی سند را. در ابتدا هر ردیف از  $H$  باید نرمالیزه شود سپس ماتریس نرمال شده و تعداد موضوعات  $K$  به الگوریتم خوشه‌بندی PLSA داده شده و ماتریس موضوع-سند  $\tilde{H}$  حاصل می‌شود. تفاوت اصلی



شکل ۵ فراوانی ویژگیها در کل اسناد



شکل ۶ خوشه بندی بر اساس ۳ ویژگی Facility, Geopolitical Entity, Geo-Social-Political



شکل ۷ خوشه بندی اسناد بر اساس ۳ ویژگی Geopolitical Entity, Location و Geo-Social-Political

کمترین استفاده از این ویژگی در اسناد را نشان می دهد و ستون بیشترین استفاده، تعداد اسنادی که بیشترین استفاده از این گروه کلمات را داشته اند را نشان می دهد. ستون میانگین، میانگین استفاده از گروه کلمه در کل ۴۵ سند را نشان می دهد.

شکل ۴ میزان استفاده از هر ویژگی در خوشه های تعیین شده را نشان می دهد. در شکل ۴، تعداد سطرها معرف تعداد خوشه ها می باشد که با توجه به اینکه نوع داده های اولیه، شامل سه نوع خبر ورزشی، اقتصادی و سیاسی می باشد پس ۳ خوشه داریم که در جدول ۱ با ۳ سطر نشان داده شده است. هر ستون این جدول بیانگر یک ویژگی می باشد. ستون اول نمایانگر ویژگی اول و میانگین نمونه های موجود در هر خوشه را نشان می دهد. به عنوان مثال ستون اول نمایانگر ویژگی facility بوده و در خوشه اول به طور میانگین ۰.۰۰۵۷ بار از این ویژگی استفاده شده است.

نتایج حاصل از خوشه بندی در شکل ۵ نشان داده شده است. در این پیاده سازی پس از مرحله کاهش بعد به ۶ ویژگی رسیدیم و فراوانی این ۶ ویژگی در کل ۴۵ سند طبق نمودار شکل ۵ می باشد.

نمودار شکل ۶ نمایش سه بعدی نتایج خوشه بندی در نرم افزار متلب بر اساس ۳ ویژگی Facility, Geopolitical Entity, Geo-Social-Political را نشان می دهد.

نمودار شکل ۷ خوشه بندی اسناد بر اساس ۳ ویژگی Geopolitical Entity, Geo-Social-Political, Location را نشان می دهد.

شکل ۸ خوشه بندی اسناد بر اساس ۳ ویژگی Geopolitical, Location و Organization را نشان می دهد.

شکل ۹ خوشه بندی اسناد در نرم افزار متلب بر اساس ۳ ویژگی Location, Person, Organization را نشان می دهد.

جدول ۱ استخراج ویژگیهای مفهومی

کلمه هدف	میانگین استفاده	بیشترین استفاده	کمترین استفاده
Facility	0.533	۹	۰
Geopolitical Entity	0.067	۱	۰
Geo-Social-Political	15.689	۴۵	۳
Location	14	۴۰	۱
Organization	0.400	5	۰
Person	9.022	28	1

	1	2	3	4	5	6
1	0.0057	0.0012	0.1979	0.1669	0.0044	0.1128
2	0.0056	0.0011	0.1935	0.1632	0.0043	0.1103
3	0.0117	0.0024	0.4043	0.3410	0.0091	0.2304

شکل ۴ میزان استفاده از ویژگیها در هر خوشه



نتایج حاصل از خوشه‌بندی با معیار سیلهوت در روش پیشنهادی محاسبه شده است و به شرح زیر می‌باشد:

سیلهوت = ۰.۵۲

## ۴-۲- ارزیابی نتایج خوشه‌بندی با معیار

### دیویس بولدین<sup>۲</sup>

شاخص دیویس بولدین، معیاری است که از شباهت بین دو خوشه (Rij) استفاده می‌کند که بر اساس پراکندگی یک خوشه (si) و عدم شباهت بین دو خوشه (dij) تعریف می‌شود. معمولاً شباهت بین دو خوشه طبق رابطه ۷ تعریف می‌شود.

$$R_{ij} = \frac{s_i + s_j}{d_{ij}} \quad (7)$$

که در آن dij و si با روابط ۸ و ۹ محاسبه می‌شوند.

$$d_{ij} = d(v_i + v_j) \quad (8)$$

$$s_i = \frac{1}{||c_i||} \sum_x d(x, v_i) \quad (9)$$

با توجه به مطالب بیان شده و تعریف شباهت بین دو خوشه شاخص دیویس بولدین به صورت رابطه ۱۰ تعریف می‌شود.

$$DB = \frac{1}{|n_c|} \sum_{i=1}^{x_c} R_i \quad (10)$$

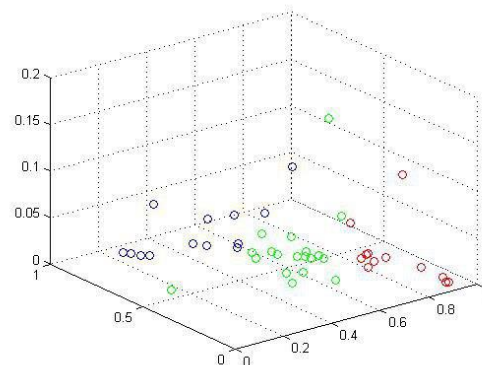
این شاخص در واقع میانگین شباهت بین هر خوشه با شبیه‌ترین خوشه به آن را محاسبه می‌کند. می‌توان دریافت که هر چه مقدار این شاخص بیشتر باشد، خوشه‌های بهتری تولید شده است.

نتایج حاصل از خوشه‌بندی با معیار دیویس بولدین برای روش پیشنهادی به قرار زیر می‌باشد:

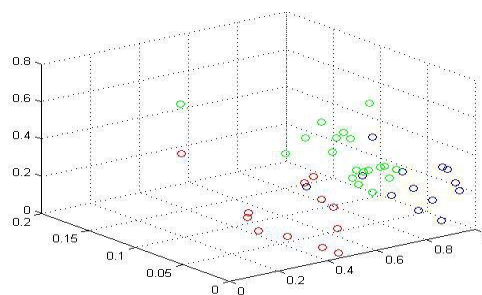
دیویس بولدین = ۰.۹۵

## ۵- جمع بندی

در این تحقیق روشی جدید با استفاده از خوشه‌بندی و فازی سازی، به منظور ارائه راه‌حلی بهینه برای کاوش در میان اسناد وب ارائه شده است. روش پیشنهادی علاوه بر کاوش اسناد وب، سازمان‌دهی، دسترسی سریع و صحیح به اسناد و صفحات وب ساخت یافته و نیمه ساخت یافته را انجام می‌دهد. طوری که علاوه بر ویژگی‌های ساختاری، ویژگی‌های محتوایی و معنایی اسناد نیز مورد کاوش قرار می‌گیرد. روش پیشنهادی با ۴۵ سند خبری از سایت پرس تی وی به عنوان داده‌های اولیه مورد ارزیابی و پیاده سازی می‌باشد. برای پیاده سازی روش پیشنهادی از نرم افزارهای متلب و پایتون استفاده شد و نتایج پیاده سازی با توجه به شاخص دیویس بولدین و سیلهوت نشان می‌دهد که روش پیشنهادی از نظر کیفیت خوشه‌ها و تراکم آنها مناسب می‌باشد.



شکل ۸ خوشه بندی اسناد بر اساس ۳ ویژگی، Geo-Social-Political, Organization و Location



شکل ۹ خوشه بندی اسناد بر اساس ۳ ویژگی، Person و Organization و Location

## ۴-۱- ارزیابی نتایج خوشه بندی با معیار سیلهوت<sup>۱</sup>

کیفیت خوشه‌ها در روشهای یادگیری بدون نظارت از طریق روشهای ارزیابی درونی انجام می‌شود. این روشها ارزیابی می‌کنند که خوشه‌ها تا چه حد از هم جدا هستند و تا چه حد به هم فشرده‌اند. نمونه‌ای از این شاخص‌ها، شاخص سیلهوت می‌باشد که طبق رابطه ۶ تعریف می‌شود.

$$s(i) = \frac{b(i) - a(i)}{\max\{b(i), a(i)\}} \quad (6)$$

a(0) به صورت میانگین فاصله بین 0 و سایر اشیاء خوشه‌ای که 0 به آن تعلق دارد، تعریف می‌شود. به روش مشابهی b(0) کمترین میانگین فاصله بین 0 و همه خوشه‌هایی است که 0 به آنها تعلق ندارد.

مقدار این شاخص بین -۱ و ۱ است. مقدار a(0) فشرده‌گی خوشه‌ای را که 0 به آن تعلق دارد نشان می‌دهد. این مقدار هر چه کمتر باشد خوشه فشرده‌تر است. مقدار b(0) نشان می‌دهد که 0 چه قدر از سایر خوشه‌ها جدا است. هر چه b(0) بیشتر باشد 0 از سایر خوشه‌ها بیشتر جدا شده است. بنابراین در شاخص سیلهوت هر چه مقدار شاخص به یک نزدیک می‌شود خوشه فشرده‌تر است و از سایر خوشه‌ها دورتر است بنابراین حالت مطلوبی است. در شرایطی که شاخص سیلهوت منفی باشد، به این معناست که 0 به اشیاء خوشه دیگری غیر از خوشه‌ای که به آن تعلق دارد، نزدیکتر است. این حالت نامطلوب است و باید از بروز آن جلوگیری کرد.

<sup>2</sup> davies bouldin

<sup>1</sup> silhouette



features," C. Nedellec and C. Rouveirol, editors, *European Conf on Machine Learning (ECML)*, , 1998.

- [9] Y. Hamasuna, Y. E. , "Semi-supervised fuzzy C-means clustering using clusterwise tolerance based pairwise constraints," in *Proceeding of the 2010 IEEE International Conference on Granular Computing, GRC'10, IEEE Computer Society.*, 2010.
- [10] Xinyue Liu, Xianchao Zhang, Ye Tian., " Webpage Segmentation based on Gomory-Hu Tree Clustering," *Undirected Planar Graph.NSFC*, 2010.
- [11] K. S. T. C. Lin WC, " CANN: An intrusion detection system based on combining cluster center and nearest neighbors.," *Knowledge-based systems* , pp. 13-21., 2014.
- [12] S. P. S.-j. H. Soo-Hoon Moona, "Energy efficient data collection in sink-centric wireless sensor networks: A cluster-ring approach," *Computer Communications*, pp. 1-14, 2016.
- [13] Iwan Tri Riyadi Yanto, Younes Saadi, Dedy Hartama, Dewi Pramudi Ismi, Andri Pranolo, "A framework of fuzzy partition based on Artificial Bee Colony for categorical data clustering," *IEEE*, 2016.
- [14] Liyong Zhang, Wei Lu, Chongquan Zhong, "Fuzzy C-Means Clustering of Incomplete Data Based on Probabilistic Information Granules of Missing Values," *Knowledge-Based Systems*, 2016.

## ۶- پیشنهاد کارهای آتی

برای کارهای آتی پیشنهاد می گردد تا تمرکز بیشتری روی راه کارهای استخراج ویژگی های اسناد و کاهش ابعاد انجام گیرد. همچنین در این تحقیق از روش خوشه بندی و فازی سازی برای دسته بندی اسناد استفاده شد، برای تحقیقات آینده، تمرکز بر روی روش های جدیدی که دقت و سرعت و کارایی بالایی برای دسته بندی اسناد داشته باشند، پیشنهاد می گردد.

## مراجع

- [1] C. W. Choo, E. Austern., "Scanning the business environment: acquisition and use of information by managers," *M.E. Williams (ED), annual Review of information Science and Technology, Learned Information, Inc. Medford.*, pp. 250-256, 1996.
- [2] Maide Abedini Bagha, F. L., " Classifying web pages and documents based on expected cross entropy and weighted vote schema.," in *The international conference on new researches in engineering science*, Tehran, 2016.
- [3] Krishna Murthy. A, Suresha, " XML URL classification based on their semantic structure orientation for web mining applications," in *International conference on information and communication technologies (ICICT 2014)*, 2015.
- [4] Thasleena N.T., V. S. , "Enhanced Associative Classification of XML Documents Supported by Semantic Concepts," in *International Conference on Information and Communication Technologies (ICICT 2014)*. , 2014.
- [5] T.-H. H.-H. Chien-Liang Liu, "Clustering documents with labeled and unlabeled documents using fuzzy semi-Kmeans," *Fuzzy Sets and Systems* 221 , p. 48-64, 2013.
- [6] G. T. Q. R. C. M. M. F. Z. X. Wang W, "Autonomic Intrusion detection : Adaptively detecting anomalies over unlabeled data streams in computer networks," *Expert Systems with Applications* , p. 4062-4080., 2014.
- [7] A. M. S. T. a. T. M. K. Nigam, " Text classification from labeled and unlabeled documents using em.," *Machine Learning*, p. 103-134, 2000.
- [8] T. Joachims., " Text categorization with support vector machines: Learning with many relevant