

بهبود خوشه بندی اسناد وب با استفاده از الگوریتم تکاملی زنبورعسل

مهرنوش پاپری¹، هدی فهندژسعدی²

¹ کارشناسی ارشد مهندسی کامپیوتر، دانشگاه آزاد اسلامی واحد شیراز، ایران mehrnoosh.papary@yahoo.com

² کارشناسی ارشد مهندسی کامپیوتر، دانشگاه آزاد اسلامی واحد شیراز، ایران h.fahandezh@yahoo.com

چکیده

امروزه حجم صفحات وب در حال افزایش است. با گسترش وب و افزایش چشمگیر حجم اطلاعات، نیاز به روش‌ها و ترفند‌هایی برای دستیابی به داده‌ها و استخراج اطلاعات از آن داریم. به منظور دستیابی به اطلاعات مورد نیاز در وب، لازم است روشی خودکار برای خوشه بندی صفحات وب ارائه شود. درمسأله خوشه بندی صفحات وب، هدف دسترسی سریع و دقیق به اطلاعات وب می‌باشد. در این مقاله با بکارگیری تلفیقی از روش‌های آماری تخصیص پنهان دیریکله و الگوریتم هوش مصنوعی زنبورعسل اقدام به خوشه بندی اسناد وب نموده ایم. از معیار ارزیابی *F-measure* به عنوان معیار ارزیابی استفاده می‌شود و نتایج به دست آمده، نشان می‌دهد که میزان کیفیت فراخوانی و دقت خوشه بندی اسناد وب با روش پیشنهادی در مقایسه با سایر روش‌ها مثل الگوریتم‌های *GA*، *K-Means* بهبود داشته است.

واژه‌های کلیدی

کلمات کلیدی: داده کاوی، خوشه بندی، اسناد وب، الگوریتم زنبورعسل.

1. مقدمه

داده کاوی که به عنوان استخراج دانش از پایگاه داده‌ها نیز شناخته می‌شود، روالی برای استخراج دانش ناشناخته از داده‌ها است. کاوش اسناد بر اساس روش‌های داده کاوی به استخراج اطلاعات و دانش از اسناد می‌پردازد. خوشه بندی اسناد یکی از مهمترین روش‌های کاوش اسناد است که دسته بندی بدون سرپرست اسناد به گروه‌های مختلف می‌باشد. سیستم‌های رایج بازیابی اطلاعات و خوشه بندی اسناد بر کلمات کلیدی استوار می‌باشند. با توجه به اینکه کلمات کلیدی مختلف می‌توانند برای توصیف یک مفهوم استفاده شوند، این سیستم‌ها می‌توانند نتایج نادرست و ناقصی را ایجاد نمایند. همچنین روابط معنایی ممکن است بین کلمات موجود باشد که شناسایی آنها نیاز به استخراج دانش دامنه مورد نظر دارد. مهمترین گام‌ها در خوشه بندی اسناد نحوه نمایش اسناد و معیار اندازه‌گیری شباهت بین آنها است. این تحقیق بر بهبود کارایی خوشه بندی اسناد تمرکز دارد. [1]

1.1. خوشه بندی

خوشه بندی یکی از تکنیک‌های بسیار قدرتمند داده کاوی برای کشف موضوع در صفحات متنی می‌باشد، در ابتدا خوشه بندی در سیستم‌های بازیابی اطلاعات و برای افزایش یادآوری ودقت و پیدا کردن صفحات مشابه بکار می‌رفت بعدها برای جستجوی مجموعه‌ای از اسناد بطور اتوماتیک بکارگیری شد. [2]

2.1. خوشه بندی اسناد

خوشه بندی اسناد عبارت است از تبدیل حجم عظیمی از اسناد متنی به گروه‌هایی از متن‌های مشابه؛ که به هر کدام از این گروه‌ها یک خوشه گفته می‌شود. مسئله خوشه بندی اسناد متنی را می‌توان به صورت ساده‌تر، مسئله پیدا کردن اسناد مشابه و قرار دادن آنها کنار هم تعریف کرد. خوشه بندی نقش حیاتی در روش‌های بازیابی اطلاعات دارد. الگوریتم‌های خوشه بندی خوشه‌ها را بر اساس ویژگی‌های داده‌ها و اندازه‌گیری شباهت‌ها و یا عدم شباهت‌ها محاسبه می‌کنند.

جستجو میان حجم زیاد داده‌های ذخیره شده روی وب کار دشواری است، بنابراین برای دستیابی سریع و دقیق به داده‌های مورد نظر کاربر، باید خوشه بندی روی صفحات وب اعمال شود. میدانیم که خوشه بندی نقش مهمی در سیستم‌های جستجو دارد. زمانی که کاربر عمل جستجو را در بین صفحات وب انجام می‌دهد انتظار دارد بهترین و نزدیکترین نتیجه را در کمترین زمان از سیستم بگیرد. برای نیل به این هدف خوشه بندی اسناد وب کارآمدترین گزینه است. [5]

تاکنون روش‌های مختلفی برای خوشه بندی مطرح شده است که در این میان استفاده از روش‌های ابتکاری برای مجموعه‌ی داده‌های بزرگ مؤثرتر واقع شده است. الگوریتم‌های جستجوی ابتکاری الگوریتم‌هایی هستند که از فرآیندهای زیستی و بیولوژیکی موجود در طبیعت الهام گرفته شده‌اند و اکثریت آنها به صورت جمعیتی عمل می‌کنند. الگوریتم‌های بهینه‌سازی فرا ابتکاری نیز در واقع یک روش ابتکاری هستند که می‌توانند با کمی تغییر برای مسائل بهینه‌سازی مختلف استفاده شوند. امروزه با بزرگ شدن مسائل دنیای واقعی و اهمیت یافتن سرعت دسترسی به پاسخ، روش‌های کلاسیک پاسخگوی بسیاری از مسائل نیستند و بیشتر از الگوریتم‌های تصادفی استفاده می‌شود. از این رو در دهه‌های اخیر استفاده از الگوریتم‌های ابتکاری و فرا ابتکاری رشد چشمگیری داشته است. [3]. از جمله این روش‌ها می‌توان به الگوریتم بهینه‌سازی ازدحام ذرات، الگوریتم ژنتیک، الگوریتم کرم شب تاب و غیره اشاره کرد.

2. پیشینه پژوهش

مطالعات زیادی در زمینه خوشه بندی داده ها و اسناد وب انجام شده که در زیر به چند مورد اشاره شده است.

نجمه سخایی و همکاران، روشی جدید با استفاده از خوشه بندی و فازی سازی به منظور راه حلی بهینه برای کاوش در میان اسناد وب ارائه کرده اند. روش پیشنهادی علاوه بر کاوش اسناد وب، سازماندهی، دسترسی سریع و صحیح به اسناد و صفحات وب ساخت یافته و نیمه ساخت یافته انجام می دهد به طوری که علاوه بر ویژگی های ساختاری، ویژگی های محتوایی و معنایی اسناد نیز مورد کاوش قرار میگیرد. [7]

لاله فرج زاده و همکاران، یک روش خوشه بندی جدید مبتنی بر الگوریتم ژنتیک و رقابت استعماری ارائه داده است. یکی از مشکلات خوشه بندی اسناد متنی وب، کار با فضای داده با ابعاد بالا می باشد که این جا یک روش کاهش ابعاد فضای داده با استفاده از الگوریتم ژنتیک ارائه شده سپس با استفاده از الگوریتم رقابت استعماری به خوشه بندی اسناد پرداخته شده است. از معیار ارزیابی F -measure به عنوان معیار ارزیابی استفاده می شود و نتایج به دست آمده، نشان می دهد که در حالتی که از الگوریتم کاهش ابعاد استفاده شود دقت خوشه بندی بالا می رود و هم چنین روش خوشه بندی با استفاده از الگوریتم رقابت استعماری نسبت به روش های مورد مقایسه بهتر عمل می کند. [5]

معصومه رجب زاده شهرستانی و همکاران، در این مقاله چندین روش خوشه بندی صفحات وب مورد بررسی قرار گرفته است. این روش ها شامل خوشه بندی اسناد براساس ترکیب تگ های HTML، روش یادگیری ماشین K نزدیکترین همسایه، روش خوشه بندی بر اساس الگوریتم K-Means و خوشه بندی براساس پیوندهای خروجی است. سپس با استفاده از تمامی این الگوریتم ها دسته بندی صفحات وب مطابق با گروه بندی های از پیش تعیین شده انجام شده و پس از قرارگیری صفحات در گروه ها و دسته های مربوط به خود میزان درستی کار هر کدام از این الگوریتم ها را از طریق روش های ارزیابی خوشه بندی سنجیده شده و با سایر الگوریتم ها مقایسه شده و درصد موفقیت هریک را در گروه های مختلف محاسبه کرده و در نهایت به معرفی موفق ترین گروه در این آزمایش پرداخته شد. [6]

مشکل خوشه بندی متن یک مشکل بسیار جدی است و توجه بسیاری از محققان را به خود جلب کرده است. یک الگوریتم ترکیبی جدید برای حل مشکل خوشه بندی متن با توجه به ترکیبی از خوشه بندی k-means با الگوریتم MVO ابداع شده است. الگوریتم MVO (multi-verse optimizer algorithm) مبتنی بر جمعیت نشان دهنده موردی جدید از الگوریتم بهینه سازی است، که هدف آن حل بسیاری از مشکلات جدی بهینه سازی سراسری می باشد. [8]

گوسامی و همکاران، روشی جهت بهبود سرعت کا-میانه ارائه نمودند که از الگوریتم خوشه بندی رهبری در جهت بهبود سرعت کا-میانه استفاده می کند. یعنی به جای محاسبه مرکز خوشه با استفاده از تمام داده ها، سعی می کند مرکز خوشه را بر اساس مقادیر بدست آمده از خوشه بندی رهبری بدست آورد. [9]

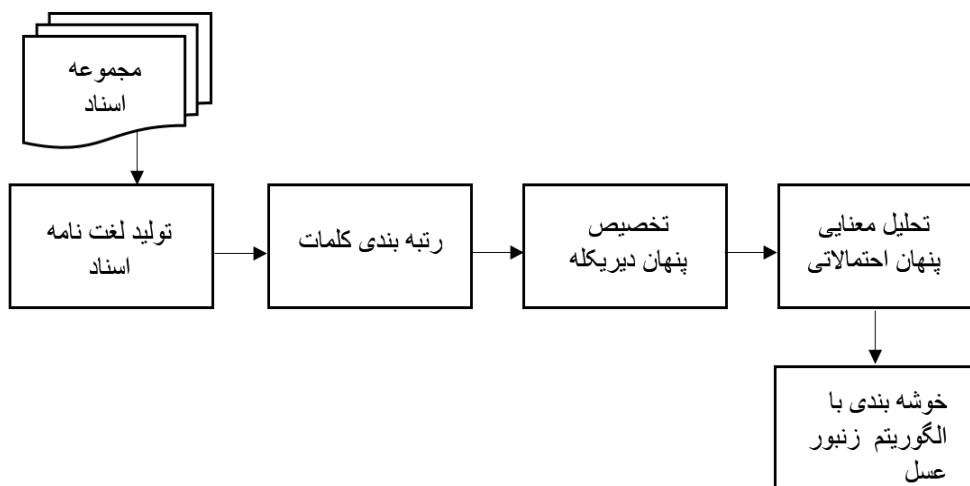
آشوکومار و دان، یک الگوریتم خوشه بندی مبتنی بر لینک را معرفی کرده اند که تنها از متن لینک برای خوشه بندی اسناد استفاده میکند. بسیاری از اسناد وب شامل تعداد زیادی لینک است که می تواند برای هدف خوشه بندی استفاده شود. [10]

3. مشخصات داده مورد استفاده

مجموعه داده DMOZ از 14 موضوع که برای هر موضوع برخی از صفحات وب انتخاب شده و در آن داده ها، را قرار داده اند و بر طبق مقاله رعنا فرصتی و همکاران در سال 2015 از سایت <https://raw.githubusercontent.com/gr33ndata/dmoz-urlclassifier/master/dmoz0409.csv> استخراج شده اند [11].

4. روش پیشنهادی

شمای کلی مدل پیشنهادی در قالب چارت 1 آورده شده است. در ادامه به بررسی آن می‌پردازیم.



1.4. تولید لغت نامه از اسناد

با دریافت یک سند، عملیاتی به عنوان فاز پیش پردازش بر روی بخش‌های مختلف آن انجام می‌پذیرد. در مرحله نخست، یک بردار از کلمات سند توجه به نقش آن‌ها در جمله استخراج می‌گردد. در این پژوهش از روشی تحت عنوان Unigram استفاده کردیم. این الگوریتم فاصله میان کلمات را به عنوان کاراکتر جداکننده در نظر گرفته و یک بردار $I \times V$ را بر می‌گرداند. نقش کلمات چندان تأثیری در خوشه بندی یک سند نخواهد داشت به همین دلیل نیز از کلماتی همانند علائم، صفات، قید، ضمائر و ... چشم‌پوشی می‌کنیم. روند حذف نیز بدین صورت است که یک بردار از کلمات زائد تولید نموده و در صورت مشاهده آن‌ها را از بین می‌بریم.

2.4. رتبه بندی کلمات

منظور از رتبه بندی، تخصیص امتیاز به هر کلمه در بردار می‌باشد. امتیازدهی حالات مختلفی دارد که بررسی تعداد تکرار کلمه در هر سند به صورت محلی و در کورپس به صورت سراسری یکی از آن‌ها به حساب می‌آید. روش مورد استفاده محقق الگوریتم Tf-Idf است. در این جا هر کلمه به صورت جداگانه مورد بررسی قرار گرفته و تکرار آن در سند محاسبه می‌گردد. اگر تعداد D سند در یک کورپس داشته باشیم، هدفارائه کلماتی همانند w_d است که بیشترین تکرار را در این اسناد به خود اختصاص داده باشند.

$$w_d = f_{w,d} \times \log \frac{|D|}{f_{w,D}} \quad (1)$$

3.4. تخصیص پنهان دیریکله (LDA)

مدل سازی موضوعی یکی از قوی‌ترین تکنیک‌ها در متن‌کاوی است، روش‌های مختلفی برای مدل سازی موضوع وجود دارد. Latent Dirichlet Allocation (LDA) یکی از محبوب‌ترین‌ها در این زمینه است. [12]

این روش فرض میکند که اسناد موضوعات متعددی را نمایش می‌دهند. یعنی از کلماتی تشکیل شده است که هر یک متعلق به یک موضوع است و نسبت موضوعات داخل یک متن با هم تفاوت است. با توجه به این نسبتها میتوانیم آن متن را در یک موضوع خاص دسته بندی کنیم. الگوریتم تخصیص پنهان دیریکله را میتوان با فرمول 2 که شامل متغیرهای آشکار و پنهان فوق میباشد بیان کرد:

$$P(\beta, \theta, z, w) = \prod_{i=1}^k p(\beta_i | \eta) \prod_{d=1}^D P(\theta_d | \alpha) \quad (2)$$

$$(w_{d,n} | \beta_{1:k}, z_{d,n}) p(z_{d,n} | \theta_d) \prod_{n=1}^N p$$

از فرمول بالا میتوان به مجموعه‌های از وابستگی‌ها پی برد. مثلاً موضوع تخصیص داده شده به هر کلمه در یک سند $(z_{(d,n)})$ بستگی به نسبت موضوعات برای سند (d, θ) دارد و یا کلمه‌ی مشاهده شده در هر سند $(w_{(d,n)})$ بستگی به موضوع اختصاص داده شده به آن کلمه $(z_{(d,n)})$ و همه‌ی موضوعات $k:1$ (β) دارد. یعنی $(z_{(d,n)})$ موضوع مربوط به آن کلمه را مشخص میکند و سپس احتمال حضور آن کلمه در آن موضوع یعنی $(w_{(d,n)})$ مشخص میشود.

در روش پیشنهادی از مدل تخصیص پنهان دیریکله برای استخراج ویژگیهای معنایی متون استفاده شده است. روش استخراج ویژگی به این صورت است که ابتدا مجموعه داده متنی، به صورت یک ماتریس بیان میشود. سپس مدل تخصیص پنهان دیریکله بر روی آن اعمال و مقادیر احتمالی بدست آمده به عنوان بردار ویژگی در نظر گرفته میشود.

4.4 تجزیه و تحلیل معنایی پنهان احتمالاتی (PLSA)

تشابه معنایی مبتنی بر احتمال یا PLSA به بیان رابطه میان متغیرهای پنهان در یک سند با پارامترهای آشکار می‌پردازد. متغیرهای پنهان عبارتند از عناوینی که منجر به دسته بندی می‌گردند و متغیرهای آشکار عبارتند از بردار کلمات تشکیل دهنده اسناد در کورپس. در دیریکله هم، متغیرهای پنهان بر روی متغیرهای تصادفی مورد بررسی قرار می‌گیرند اما تفاوت این روش، توجه بیشتر به رخدادهای پیشین بر اساس تئوری بیز است. به بیان دیگر، هر عنوانی که سابقه تکرار کمتری داشته باشد انتخاب نمی‌شود. علاوه بر این، استفاده از بردارهای احتمال توزیع دیریکله و احتمال انتخاب عناوین در این روش وجود ندارند. در نتیجه ترکیب PLSA با خروجی‌های دیریکله باعث می‌شود تعداد عناوین کاهش چشمگیری داشته باشند. رابطه 3 و 4 به ترتیب کار بر روی عناوین استخراج شده بر اساس حداکثر تشابه کورپس را نشان می‌دهد. $n(d, w)$ تعداد تکرارهای عناوین در کورپس را نشان می‌دهد.

$$(3) \quad \iota = \sum_{d \in D} \sum_{w \in W} n(d, w) \log P(d, w)$$

$$(4) \quad P(d, w) = \sum_{z \in Z} P(z) p(w|z) p(d|z)$$

دلیل استفاده از \log در رابطه بالا افزایش حساسیت به انتخاب عناوین به دست آمده از دیریکله می‌باشد [13].

5.4. خوشه بندی به کمک الگوریتم کلونی زنبور عسل

بردارهای خروجی PLSA احتمال حضور عناوین در اسناد را به صورت جداگانه نشان می‌دهد. در این مرحله می‌بایستی مجموع احتمال به دست آمده از حضور عناوین در هر سند به همراه احتمال آن‌ها خوشه بندی شوند. این الگوریتم با بهره‌گیری از زندگی اجتماعی زنبور ها در کندو به عنوان عامل‌های هوشمند، سعی در حل مسائل بهینه‌سازی دارد. زنبورهای کلونی به سه دسته کارگر، جستجو کننده و مهاجم تقسیم بندی می‌شوند. در ابتدای کار، فضای مساله یا همان منابع غذایی در اختیار عامل‌ها قرار می‌گیرد. (تعداد زنبورهای کارگر برابر است با منابع غذایی موجود). زنبورهای کارگر توسط مهاجم‌ها رانده شده و منابع غذایی آن‌ها را به سرقت می‌برند. روال اجرای الگوریتم به صورت زیر می‌باشد [14].

- تولید n عامل هوشمند و قرار دادن آن‌ها در فضای مساله
- ترکیب زنبورهای کارگر، مهاجم و جستجو کننده در سرتاسر فضای مساله
- بررسی کیفیت پاسخ‌های ارائه شده توسط هر زنبور کارگر و به خاطر سپردن بهترین راه حل به واسطه مکان محلی.
- ذخیره کردن راه حل‌های ارائه شده تا لحظه فعلی توسط هر زنبور
- تکمیل روند حرکت زنبورها در فضای مساله و به دست آوردن راه حل‌های تازه تر
- بررسی شهد راه حل‌های ارائه شده توسط زنبورهای کارگر و انتخاب هر کدام به صورت تصادفی تابع انتخاب مسیره‌های بهینه بر اساس احتمال به صورت زیر می‌باشد

$$P_i = \frac{fit_i}{\sum_{i=1}^{SN} fit_n} \quad (5)$$

در این رابطه، fit_i برابر است با تابع تناسب به دست آمده در مکان i که با میزان شهد موجود در این مکان رابطه مستقیم دارد. SN نیز برابر است با تعداد پاسخ‌های مساله که این پارامتر با تعداد زنبورهای کارگر برابر است. می‌بایستی توجه داشت، در روند انتخاب مکان جدید منابع غذایی و به خاطر سپردن آن، زنبورها از رابطه زیر تبعیت می‌کنند.

$$v_{ij} = x_{ij} + \phi_{ij}(x_{ij} - x_{kj}) \quad (6)$$

در این جا، K و j اندیس‌های تصادفی و متفاوت از i هستند. همچنین ϕ_{ij} نیز یک تابع تصادفی نرمال بین $[-1, 1]$ می‌باشد که در اصل یک فاکتور وزن دهی جهت مقایسه منابع غذایی حاضر در یک مکان به حساب می‌آیند.

$$Fitness(i) = \frac{Intra Cluster Distance + Number Of Neighbors}{Intra Cluster Distance} \quad (7)$$

هدف از ارائه فرمول 7، افزایش قدرت در گره‌های سرخوشه‌ای است که بیشترین فاصله فراخوشه‌ای به همراه تعداد همسایگان مجاور را داشته و دارای کمترین فاصله درون خوشه‌ای باشند.

5. معیار ارزیابی

معیار ارزیابی در این جا پارامتر F-Measure است که یک فاصله هارمونیک میان بهترین پاسخ‌های ارائه شده و میانگین خطا می باشد. رابطه 8 الی 10 نحوه محاسبه این پارامتر را نشان می دهد. (C خروجی های به دست آمده از روش دسته بندی پیشنهادی و C* نتایج قابل پیش بینی است)

$$P(C, C^*) = \frac{|C \cap C^*|}{|C|} \quad (8)$$

$$R(C, C^*) = \frac{|C - C^*|}{|C|} \quad (9)$$

$$F - Measure(C, C^*) = \frac{2 \times P \times R}{P + R} \quad (10)$$

منظور از P و R در روابط بالا به ترتیب Precision و Recall می باشد.

6. نتایج شبیه سازی

در این بخش از تحقیق به بررسی نتایج به دست آمده از روش پیشنهادی می پردازیم. الگوریتم هایی که در این بخش مورد مقایسه قرار می گیرند عبارتند از : K-Means، ژنتیک، کلونی زنبور عسل و مدل پیشنهادی. تمامی این روش ها بر روی دیتاست DMOZ اجرا گردیده و نتایج هر کدام در جدول 1 قید شده است. (نتایج دیگر روش ها از مقاله رعنا فرصتی و همکاران در سال 2015 استخراج شده اند)

جدول 1: خروجی مدل پیشنهادی و دیگر روش های خوشه بندی اسناد

مقدار پارامتر F-Measure	نام الگوریتم
0.5423	K-Means
0.7194	GA
0.8574	BCO colony
0.91036	مدل پیشنهادی

7. نتیجه گیری

بررسی نتایج به دست آمده از مدل پیشنهادی نشان از افزایش کارایی نسبت به دیگر روش‌های خوشه‌بندی اسناد دارد. تولید بردار لغات با در نظر گرفتن یک مقدار حد آستانه و همچنین استخراج بردار ویژگی‌ها توسط الگوریتم تخصیص پنهان دیریکله منجر به دسته‌بندی معنایی اسناد گردیده است. البته می‌بایستی به این نکته اشاره داشت که آنالیز معنایی با استفاده از PLSA باعث می‌شود تعداد عناوین کاهش یافته و توزیع‌های با درصد احتمال بالا به عنوان خروجی برگردانده شوند. روش‌هایی که خوشه‌بندی را با مرکزیت و بر اساس میانگین‌گیری یا به طور کلی قطعه‌بندی انجام می‌دهند اغلب دارای مشکلاتی همانند افزایش وسعت خوشه، عدم حساسیت به نویز، عدم دسته‌بندی نمونه‌های مجاور و کارایی پایین در فضای مساله غیرخطی می‌باشند که این مشکلات با در نظر گرفتن پارامترهای چگالی همانند حداقل تعداد همسایه‌ها و فاصله شعاعی تا حد زیادی برطرف شده است. علاوه بر این، قابلیت الگوریتم کلونی زنبور عسل در جستجوی بهینه سراسری باعث می‌شود فضای مساله به صورت وسیعی مورد پوشش قرار گرفته و همگرایی در دوره‌های ابتدایی رخ دهد.

منابع

- [1] مریم امیری و حسن ختن لو، 1392. خوشه‌بندی اسناد مبتنی بر آنتولوژی و رویکرد فازی. نشریه فناوری اطلاعات و ارتباطات ایران، سال پنجم شماره 17.
- [2] عبدالکریم الهی. 1389. بهبود خوشه‌بندی اسناد بر مبنای یک اندازه‌گیری مشابهت جدید. چهارمین کنفرانس داده کاوی.
- [3] میترا مقصودی، مهدیه قزوینی. 1398. کاهش ابعاد داده‌ها با استفاده از یک روش ترکیبی انتخاب ویژگی. سومین کنفرانس ملی کامپیوتر، فناوری اطلاعات و کاربردهای هوش مصنوعی.
- [4] پریسا چابکی بناب، علی رضا حاجی اسکندر، سیدعلی شریفی. 1398. خوشه‌بندی اسناد وب فارسی با استفاده از ترکیب روش‌های داده کاوی و یک الگوریتم تکاملی. ششمین کنفرانس بین‌المللی یافته‌های نوین علوم و تکنولوژی با محوریت علم در خدمت توسعه.
- [5] لاله فرج زاده، حسن ختن لو، منصور اسماعیل پور. 1394. خوشه‌بندی اسناد وب با استفاده از الگوریتم ژنتیک و رقابت استعماری. کنفرانس ملی فن آوری، انرژی و داده با رویکرد مهندسی برق و کامپیوتر.
- [6] معصومه رجب زاده شهرستانی، فاطمه احمدی آبکناری، غلامحسین اکباتانی فرد. 1395. طبقه‌بندی محتوایی صفحات وب با استفاده از راهکارهای خوشه‌بندی. نخستین کنفرانس ملی تحقیقات بین‌رشته‌ای در مهندسی کامپیوتر، برق، مکانیک و مکاترونیک.
- [7] Najmeh. S., Fariba S, Majid Kh, 2018. Clustering Web Documents Using Ontology-Based Fuzzy Method.
- [8] Abasi, A.K., Khader, A.T., Al-Betar, M.A., Naim, S., Alyasseri, Z.A.A. and Makhadmeh, S.N., 2020. A novel hybrid multi-verse optimizer with K-means for text documents clustering. *Neural Comput. Appl.*, 32(23), pp.17703-17729.
- [9] Gowtham, S., Goswami, M. . 2014. An approach for document Pre-Processing and K Means Algorithm Implementation. *Advances in Computing*.
- [10] Ashokkumar P, Don S ,2019. Link-Based Clustering Algorithm for Clustering Web Documents.
- [11] Forsati, R., Keikha, A., & Shamsfard, M. 2015. An improved bee colony optimization algorithm with an application to document clustering. *Neurocomputing*, 159, 9-26.
- [12] Hamed Jelodar · Yongli Wang · Chi Yuan · Xia Feng · Xiahui Jiang · Yanchao Li · Liang Zhao.2018 .Latent Dirichlet Allocation (LDA) and Topic modeling: models, applications, a survey
- [13] Hofmann, T. 2017. Probabilistic latent semantic indexing. In *ACM SIGIR Forum* (Vol. 51, No. 2, pp. 211-218). ACM.
- [14] Xue, Y., Jiang, J., Zhao, B., & Ma, T. 2017. A self-adaptive artificial bee colony algorithm based on global best for global optimization. *Soft Computing*, 1-18.