

## یک روش وزن دهی مبتنی بر موقعیت واژه جهت مشابهت سنجی اسناد

مریم اسدی لنگرودی<sup>۱</sup>، سید ابوالقاسم میرروشندل<sup>۲</sup>

<sup>۱</sup>دانشجوی کارشناسی ارشد مهندسی فناوری اطلاعات، گروه کامپیوتر، پردیس دانشگاه گیلان، رشت  
maryam.as16@yahoo.com

<sup>۲</sup>استادیار دانشگاه، گروه کامپیوتر، دانشگاه گیلان، رشت  
mirroshandel@gmail.com

### چکیده

اندازه گیری میزان شباهت اسناد موجود در وب، از آن جهت دارای اهمیت است که در بسیاری از زمینه‌ها، مانند بازیابی اطلاعات، دسته بندی متون، خوشه بندی اسناد، سیستم های تشخیص تقلب و سرقت ادبی، خلاصه سازی متون، و سایر حوزه ها، نقش مهم و اساسی ایفا می کند. میزان درستی این تشخیص، می تواند کارایی، دقت و صحت هر یک از فرایندهای مذکور را بالا برد. در تمام روش های مشابهت سنجی اسناد، اصول کار، تشخیص شباهت سندها بر مبنای شناخت دقیق ویژگی های مشترک آنها است. از این رو باز نمود سند بر مبنای ویژگی های بارز آن بسیار با اهمیت است. روش های مهم مشابهت سنجی اسناد، شامل مدل های لغوی و معنایی مبتنی بر محتوا و مدل های مبتنی بر ساختار صفحه است. در روش های لغوی، ویژگی اصلی یک سند، وزن واژه های آن است. بنابراین وزن دهی صحیح واژه، می تواند یک گام موثر در نمایش دقیق ویژگی های بارز اسناد باشد.

در این مقاله، هدف اصلی، ارائه روشی بهبود یافته در وزن دهی واژه، برای تعیین میزان شباهت لغوی اسناد متنی است. روش وزن دهی پیشنهادی بر مبنای طرح TD - IDF و با تاکید بر اهمیت بیشتر واژه های خطوط ابتدایی اسناد، توانسته است، دقت و فراخوانی را در دسته بندی و میزان صحت را در خوشه بندی اسناد مجموعه TDT5 افزایش دهد.

### کلمات کلیدی

شباهت اسناد، معیار شباهت، وزن دهی واژه، خوشه بندی اسناد، دسته بندی اسناد، بازیابی اطلاعات، شباهت سنجی لغوی، TD - IDF

### ۱- مقدمه

• شیوه های شباهت سنجی اسناد می تواند مبتنی بر محتوا و مبتنی بر ساختار صفحه باشد [۸،۹]. در شیوه های مبتنی بر محتوا سه دسته بندی مهم شباهت سنجی لغوی، معنایی و ترکیبی وجود دارد. روش های لغوی به بررسی میزان هم پوشانی اسناد و تعداد واژه های مشترک میان آنها می پردازند. روش های معنایی اسناد را از نظر مفهوم و معنی واژه ها بررسی می کند شیوه های ساختاری نیز به بررسی مشابهت ساختاری صفحات، مانند ارجاعات و لینک های ورودی و خروجی یکسان در اسناد می پردازند. در تمام این روش ها، پس از پیش پردازش های لازم، استخراج ویژگی های اسناد صورت می گیرد و الگوریتم های سنجش

در بیشتر رهیافت های بازیابی اطلاعات توسط موتورهای جستجو، مبنای جستجو، تعیین میزان شباهت اسناد با پرس و جو و سپس رتبه بندی اسناد بر این اساس است [۱-۳]. در حوزه داده کاوی و دسته بندی و خوشه بندی اسناد نیز مبنای کار، روش هایی است که میزان شباهت سندها را با دقت بالا بررسی می کند [۴-۶]. مشابهت سنجی اسناد، در حوزه های دیگر نیز چون خلاصه سازی متن، سیستم های پرسش و پاسخ، ترجمه ماشینی، تشخیص تقلب و سرقت ادبی، نقش مهم ایفا می کند [۷].



طرح معرفی شد. در [۱۹] یک طرح بهبود یافته بر مبنای اصلاح عامل TF، ارائه شد. طرح پیشنهادی، اهمیت مدل توزیع یک واژه را در یک سند بررسی نمود و اهمیت واژگان را در توزیع یکنواخت و غیر متراکم در سند بیان کرد.

• وزن دهی دقیق واژه‌ها می‌تواند ویژگی‌های مهم اسناد را نمایش دهد. در این مقاله ابتدا مفاهیم اولیه روش‌های شباهت‌سنجی لغوی مبتنی بر محتوا و روش‌های وزن دهی، بیان می‌شود و سپس، یک روش جدید که با تغییر مدل وزن‌دهی واژه بر مبنای اهمیت بیشتر واژه‌های خطوط ابتدایی سند توانسته است، بهبود قابل توجهی در کارایی، دقت و فراخوانی دسته‌بندی و صحت خوشه‌بندی اسناد ایجاد نماید، پیشنهاد می‌گردد و حاصل اجرای این روش بر روی اسناد مجموعه متون TDT5 مورد مطالعه و ارزیابی قرار می‌گیرد.

## ۲- نمایش سند در شباهت سنجی لغوی

شباهت سنجی لغوی، اسناد را بر اساس میزان هم‌پوشانی آن‌ها بررسی می‌کند. در این روش، اندازه‌گیری میزان شباهت، بر اساس تعداد واژه‌های مشترک موجود در اسناد است و مفهوم و معنی واژه در زبان طبیعی حائز اهمیت نیست. ویژگی بارز سند، واژه‌های با وزن بالای آن است.

در این روش، برای نمایش سند از مدل فضای برداری سالتون استفاده می‌شود. هر سند به صورت کیسه‌ای از کلمات در یک فضای برداری، در نظر گرفته می‌شود و توالی و جایگاه واژه در سند، لحاظ نمی‌گردد. ابتدا پیش‌پردازش و عمل قطعه بندی اسناد صورت می‌گیرد و واژه‌های سند یا همان توکن‌ها استخراج می‌شوند. یک توکن می‌تواند یک واژه یا یک عبارت یا یک ngram از واژه‌های متوالی باشد و لزوماً یک کلمه در فرهنگ لغات یک زبان طبیعی نیست. در این مقاله یک واژه به عنوان یک توکن در نظر گرفته شده است. وزن هر واژه، بر اساس مدل‌های وزن‌دهی محاسبه می‌گردد و یک سند به صورت برداری از وزن واژه‌های آن در یک فضای برداری به ابعاد تعداد واژه‌های مجموعه نمایش داده می‌شود [۲۰]. در مجموعه‌ی اسناد  $D = \{t_1, t_2, \dots, t_m\}$  مجموعه‌ای از واژه‌های متمایز است، که در آن،  $t_i$  یک واژه است. مجموعه  $V$  معمولاً به عنوان مجموعه واژگان نامیده می‌شود، و  $m$  اندازه یا تعداد واژه‌های آن است.

وزن هر واژه  $t_i$  در سند  $d_j \in D$  برابر با  $w_{ij}$  است که اهمیت آن واژه را در آن سند، بیان می‌کند. برای واژه‌ای که در سند  $d_j$  نمایان نمی‌شود،  $w_{ij}=0$  خواهد بود. بنابراین هر سند  $d_j$  با یک بردار واژه به صورت فرمول (۱) نمایش داده می‌شود:

$$d_j = \{w_{1j}, w_{2j}, \dots, w_{mj}\} \quad (1)$$

که در آن  $d_j$  یک بردار سند در فضای  $m$  بعدی و  $w_{ij}$  وزن واژه است.

## ۲-۱- پیش پردازش اسناد برای نمایش

اسناد در وب، به صورت صفحات حاوی کدها و تگ‌های HTML است. برای تبدیل اسناد متنی به صورت برداری از واژه‌ها، ابتدا اعمال پیش‌پردازش روی آن‌ها صورت می‌پذیرد. ابتدا محتوای تگ‌ها از آن‌ها استخراج می‌شود و بدنه اصلی متن سند حاصل می‌گردد سپس تگ‌ها حذف می‌شوند. در بسیاری

• مبنای معیارهای تعیین شباهت، میزان شباهت ویژگی‌های اسناد را اندازه‌گیری می‌کنند [۱۲-۱۰].

• در تمام روش‌های مذکور، تشخیص دقیق ویژگی‌های با اهمیت اسناد، مدل نمایش سند، الگوریتم‌های شباهت‌سنجی و معیارهای تعریف شباهت حائز اهمیت است و بهبود در هر یک، می‌تواند منجر به تشخیص دقیق تر شباهت اسناد شود.

• در مدل شباهت سنجی لغوی اسناد به صورت برداری از وزن واژه‌های سند به عنوان ویژگی‌های اصلی نمایش داده می‌شوند. معروف ترین طرح وزن دهی در این روش، TF-IDF است. این طرح بر مبنای تعداد تکرار یک واژه در یک سند و کل اسناد مجموعه متون به هر واژه سند یک درجه اهمیت تخصیص می‌دهد. و جایگاه واژه در سند لحاظ نمی‌شود. این طرح دارای معایبی است. نخست این که یک طرح موردی است و زیر بنای مفاهیم ریاضی ندارد و دوم این که ابعاد وسیع و اندازه بزرگ مجموعه واژگان در یک مجموعه اسناد، حجم محاسبات را بسیار زیاد می‌نماید. این طرح، جایگاه واژه را در سند لحاظ نمی‌کند و تمام واژه‌ها در سند از فرمول وزن‌دهی یکسان استفاده می‌کنند. همچنین تنها نمایانگر اهمیت لغوی یک واژه است و معانی واژه‌ها را در نظر نمی‌گیرد. همچنین دارای معایبی چون عدم توجه به مدل توزیع واژه در سند، مدل توزیع بین کلاسی، اهمیت طول و محل قرارگیری واژه وسایر موردها است. مطالعات زیادی در این زمینه انجام گرفته است که به طور موردی تعدادی از آنها به شرح زیر بیان شده است.

• در [۱۳]، با استفاده از طرح TF-IDF و الگوریتم KNN، تعداد ۷۵۰۰ مقاله خبری کشور اندونزی با ۱۵ دسته موضوعی از پیش تعریف شده، دسته‌بندی شدند و صحت این دسته‌بندی حدود ۹۸٪ بیان شده است. در [۱۴] طرح وزن‌دهی TF-IDF, LSI, MULI, WORD بر روی اسناد به زبان‌های چینی و انگلیسی مورد مطالعه قرار گرفت که در دسته‌بندی اسناد انگلیسی، این طرح در مرتبه دوم صحت دسته‌بندی قرار گرفت. در [۱۵] لن و همکاران، عملکرد ۹ طرح سنتی از جمله Binary, TF, TF-IDF, TF-IDF, LOG TF-IDF, TF-IDF Prob, IDF و TF-Chi را در مجموعه داده‌های های بزرگ با الگوریتم دسته بندی SVM مقایسه کردند. آن‌ها نتیجه گیری نمودند که TF-Chi و TF طرح‌هایی با بهترین عملکرد است. همچنین یک طرح خاص وزن دهی به نام TF-RF را پیشنهاد دادند که از فراوانی مربوط یک واژه استفاده می‌کند و ادعا کردند که طرح جدید آن‌ها، دارای عملکرد بهتری نسبت به طرح‌های دیگر است. در [۱۶]، مشکل بارز برای طرح TF-IDF عدم توجه به روابط میان واژه‌های هم‌معنی و عدم شناخت روابط کلمات ارجاعی به ضمیر، مطرح شد و به این منظور از ترکیب پیکره‌متون و روش‌های معنایی با طرح فوق بهره گرفته شد. در [۱۷]، یک طرح بهبود یافته TF-IDF با تمرکز بر اهمیت مکانی واژه و همچنین طول یک واژه بیان شده است. این طرح دارای دقت و فراخوانی بالاتری نسبت به طرح اولیه است. در [۱۸] بیان شده است که طرح TF-IDF، در دسته‌بندی اسناد، توزیع بین کلاس‌ها و داخل کلاسی واژه یا عبارت  $t_i$  را در نظر نمی‌گیرد و مشکل اصلی را در قسمت IDF طرح بیان کرده است. سپس یک الگوریتم تکرار شونده جهت بهبود



مدل پیشنهادی، اهمیت موقعیت و جایگاه یک واژه را در یک سند در نظر گرفته و توانسته است مشابهت سنجی دقیق تری را حاصل کند.

### ۲-۳- معیار مشابهت سنجی

تعیین میزان شباهت یا فاصله بردارهای اسناد، نیازمند استفاده از معیارهایی برای اندازه گیری است. معیارهای زیادی برای این اندازه گیری وجود دارد بعضی از این معیارها، میزان شباهت و بعضی میزان فاصله را محاسبه می کنند که قابل تبدیل به یکدیگر است. از مهمترین آن ها در بررسی شباهت سنجی متون با در نظر گرفتن توکن های تک واژه ای، معیار شباهت کسینوسی، ضریب دایس، ضریب جاکارد، فاصله اقلیدسی، فاصله منهن و معیار کالباک لیبلر است.

وقتی اسناد، به صورت برداری از ویژگی ها نمایش داده می شوند، مقدار کسینوس زاویه بین بردارها می تواند معیار اندازه گیری میزان شباهت آن ها باشد. این اندازه در بازه [۰,۱] است. در صورتی که دو بردار منطبق بر هم و یا کاملاً شبیه به هم باشند، این عدد ۱ خواهد بود، شباهت کسینوسی از مشهورترین معیارهای تشخیص شباهت متون در حوزه های بازایی اطلاعات و خوشه بندی اسناد است [۲۴]. فاصله اقلیدسی نیز یک معیار متریک مطرح و موفق تعیین میزان فاصله است که قابلیت تبدیل به معیار شباهت را دارد. این معیار پیش فرض، در الگوریتم های خوشه بندی K - Means استفاده می شود.

### ۳- مبنای روش پیشنهادی

در حوزه خلاصه سازی متن، امتیاز دهی به جملات، از اهمیت بالایی برخوردار است و یکی از مدل های مبنایی این حوزه به شمار می رود. در روش های خلاصه سازی متون، به دلیل ابعاد وسیع ویژگی های متن، مبنای کار، استخراج ویژگی های مهم تر و بارز سند است. مدل های مبتنی بر موضوع، مبتنی بر اهمیت جمله و مبتنی بر موقعیت واژه از جمله روش های با اهمیت در خلاصه سازی خودکار متون محسوب می شوند. در مدل های مبتنی بر موقعیت واژه، اهمیت واژه هایی که در موقعیت های مکانی خاص از متن قرار دارند، بالاتر است. به طور مثال در مقالات خبری و روزنامه های الکترونیکی، جمله اول متن، دارای رتبه بالاتری نسبت به سایر جملات است و یا جمله اول هر پاراگراف در مقالات خبری و جملات آخر در مقالات علمی از اهمیت بیشتری برخوردارند [۲۷-۲۵]. بر این مبنای روش پیشنهادی ارائه می گردد.

### ۴- روش وزن دهی پیشنهادی TF - IDF\_P

در روش پیشنهادی TF - IDF\_P، فرض بر این است که خطاهای جملات اولیه یک سند، می تواند حاوی چکیده مطالب کل سند و دارای بیشترین اهمیت محتوایی آن سند باشد. به طور مثال، در مقالات علمی خطوط اول، حاوی چکیده مقاله و در مقالات خبری، حاوی عناوین اصلی خبر است. از این رو، هنگام وزن دهی توکن ها بر اساس این طرح به توکن هایی که جزء ۲۵۰ کاراکتر نخست یک سند باشند و در خطاهای اولیه سند قرار بگیرند، ضریب ۵ برابر وزن حاصل از طرح اولیه TF - IDF اختصاص می یابد و سایر واژه ها از ضریب ۱ این طرح برخوردار خواهند بود. فرمول (۴) الگوریتم اجرایی این روش را نشان می دهد:

از الگوریتم ها، به جهت اهمیت بعضی از تگ ها مانند تگ عنوان و یا سایر تگ های مورد اهمیت در مدل های شباهت ساختاری، شاخص های ویژه ای برای آن ها در نظر گرفته می شود. در مرحله بعد، عمل قطعه بندی یا توکن بندی سند صورت می پذیرد. و واژه های سند استخراج می گردد. با توجه به ابعاد بزرگ داده ها، از پردازش هایی جهت کاهش ابعاد واژه ها استفاده می شود، مانند حذف علائم و اعداد. یکی از پردازش های مهم در این زمینه، حذف کلمات توقف است. کلمات توقف در یک زبان طبیعی واژه هایی است پرتکرار که بار معنایی خاصی ندارد. به عنوان مثال، در زبان انگلیسی کلماتی چون "And" و "The" کلمات توقف به شمار می آیند. سپس هر توکن، ریشه یابی می شود و ریشه یا مصدر اصلی واژه، جایگزین آن می گردد. بسیاری از مدل ها، از روش های دیگری نیز برای کاهش ابعاد بهره می جویند.

پس از پیش پردازش اسناد، عملیات وزن دهی به توکن ها بر اساس انواع طرح های وزن دهی، انجام می شود و بردار سند واژه ایجاد می گردد.

### ۲-۲- روش های وزن دهی واژه

وزن دهی واژه، یکی از مهم ترین مراحل در نمایش یک سند بر اساس اهمیت ویژگی های آن است. این امر می تواند ویژگی های شاخص یک سند را نمایان کند. عمل شباهت سنجی اسناد بر اساس تطبیق وزن واژه ها به عنوان ویژگی ها، انجام می گیرد. در روش های مختلف شباهت سنجی اسناد از طرح های گوناگون وزن دهی استفاده شده است.

یک طرح مهم وزن دهی در مدل نمایش برداری اسناد، طرح مشهور TF - IDF است. این طرح، یک عدد آماری ارائه می دهد که نشان دهنده درجه اهمیت یک واژه در یک سند از مجموعه اسناد یا پیکره متنی است [۲۱، ۲۲].

TF، معیاری است که ارزش آن، تعداد تکرار یک واژه در یک سند واحد است. در این حالت، واژه ی پرتکرار در یک سند، جزء با اهمیت ترین واژه های آن سند، محسوب می شود. از سوی دیگر اگر آن واژه، در سایر اسناد مجموعه و یا پیکره متنی نیز به کار رفته باشد، یک واژه شاخص برای آن سند محسوب نمی گردد. معیار معکوس تکرار سند یا IDF می تواند این حالت را متعادل نماید. این معیار عدد متعادل شده تعداد اسناد حاوی یک واژه را نسبت به کل اسناد مجموعه نشان می دهد. فرمول (۲) محاسبه این معیار را نشان می دهد [۲۳]:

$$IDF_i = \log \left( \frac{N}{df_i} \right) \quad (2)$$

که در آن N، تعداد کل اسناد مجموعه اسناد یا پیکره متنی و df<sub>i</sub> تعداد سندهایی است که واژه t<sub>i</sub> در آن ها وجود دارد.

در این طرح وزن واژه t<sub>i</sub> در سند d<sub>j</sub> به صورت فرمول (۳)، خواهد بود که w<sub>ij</sub> وزن واژه است:

$$W_{ij} = TF_{ij} \cdot IDF_i \quad (3)$$

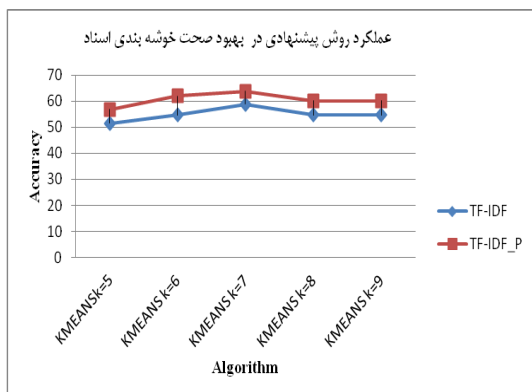
در این طرح، واژه های دارای بیشترین اهمیت و ویژگی شاخص یک سند است که در آن سند بیشترین تکرار و در سایر اسناد مجموعه کمترین رخداد را داشته باشد.

در این مقاله با تغییر در مدل وزن دهی واژه ها، بر اساس این طرح، روشی جدید با تاکید بر اهمیت واژه های خطوط اولیه اسناد مطرح شده است.

و TW، ماتریس سند واژه با وزن دهی پیشنهادی و با تعداد ۳,۰۳۸ سند و ۶۵ واژه حاصل گردید.

خوشه بندی اسناد، یک روش متن کاوی بدون نظارت است که اسناد متنی را به طور موثر بر اساس میزان شباهت آن‌ها، سازماندهی می‌کند. روش‌های خوشه بندی مبتنی بر چهار مفهوم اصلی، مدل نمایش داده‌ها، معیار شباهت، روش خوشه بندی، و الگوریتم خوشه بندی است. از آنجا که الگوریتم‌های خوشه بندی تفکیکی، مانند K - means نسبت به الگوریتم‌های خوشه بندی سلسله مراتبی، برای بررسی مجموعه داده‌های سندی بزرگ، مناسب‌تر است و نیاز به محاسبات کمتری نیاز دارد [۱۹]. برای پیاده سازی از این الگوریتم استفاده شده است [۲۸].

در پیاده سازی روش پیشنهادی، با استفاده از نرم افزار Weka از الگوریتم خوشه بندی تفکیکی K - means با معیار شباهت فاصله اقلیدسی در ۵ مرحله با در نظر گرفتن تعداد خوشه‌های ۹, ۸, ۷, ۶, ۵، استفاده شده است، که نمودار مقایسه‌ای آن در (شکل ۱) قابل مشاهده است:



شکل (۱): ارزیابی عملکرد TF-IDF\_P در صحت خوشه بندی

همان گونه که در (شکل ۱) مشخص است، روش وزن دهی پیشنهادی در تعداد خوشه‌های مختلف انتخابی دارای صحت بیشتر نسبت به روش وزن دهی بدون در نظر گرفتن تعداد کارکترهای اول بوده است و بهبود اجرای این روش، کاملاً در شکل نمایان است.

در مرحله بعد، از ارزیابی دیگر به این منظور استفاده شد. در مجموعه اسناد، می‌توان بر اساس تگ‌های مورد نظر یا با توجه به مشخصه‌های خاص سند، آن‌ها را در دسته‌هایی دارای برجسب مشخص، دسته بندی نمود. این مدل یک روش با نظارت در داده کاوی محسوب می‌شود. از مهم ترین الگوریتم‌های دسته بندی اسناد، روش K نزدیکترین همسایه است. این مدل بسیار به مدل‌های خوشه بندی تفکیکی، مشابه است و اشیاء را با توجه به مشابهت با اشیاء مجاور آن محاسبه می‌کند. این الگوریتم در دسته بندی متون، دقت و فراخوانی قابل قبولی را ارائه می‌دهد [۲۹]. روش‌های وزن دهی واژه تاثیر بسزایی در دسته بندی اسناد دارد.

در این مرحله، دو ماتریس طریق الگوریتم دسته بندی K نزدیکترین همسایه در چهار مدل تعداد همسایه‌های نزدیک ۵, ۳, ۱ و ۷ و موضوع خبر به عنوان کلاس نهایی، در ۱۰ کلاس موضوعی، دسته بندی شدند. در هر مرحله دقت و فراخوانی دسته بندی و در نهایت میانگین معیار F - Measure و

$$\text{IF POS}(t_{ij}) \leq 250 \text{ THEN} \\ W_{ij} = 5 * \text{TF}_{ij} \cdot \text{IDF}_i \quad (۴)$$

Else

$$W_{ij} = \text{TF}_{ij} \cdot \text{IDF}_i$$

که تابع POS موقعیت مکانی واژه  $t_{ij}$  را در سند، بیان می‌کند.

بر این اساس، دو سندی که در خطوط اولیه خود دارای واژه‌های مشترک باشند، امکان شباهت بیشتری با یکدیگر دارند. این مدل وزن دهی، توانسته است در بهبود عملکرد خوشه بندی و دسته بندی اسناد، تاثیر گذار باشد. به طور مثال اگر توکن "air"، در موقعیت محدوده ۲۵۰ کاراکتر اول یک سند، ۳ مرتبه تکرار شده باشد و در سایر مکان‌های سند، در مجموع ۶ مرتبه، رخ داده باشد، عدد تکرار این واژه در سند بر اساس روش پیشنهادی، عدد ۲۱ خواهد بود. در روش اولیه این عدد ۹ است.

#### ۴-۱- استخراج داده

در این روش، از اسناد حاوی متن به زبان انگلیسی مجموعه داده TDT5، استفاده شده است. لازم به ذکر است که فایل‌های این مجموعه، اسناد خبری است. از تعداد ۳,۳۲۳ فایل با ساختار HTML این مجموعه، تعداد ۱,۷۶۳,۶۹۹ تگ <doc> به عنوان سند مجزا، استخراج شده است. این اسناد توسط عملیات پیش پردازش، به قالب ساختارمند تبدیل گردید. از این تعداد، ۱,۲۵۸,۰۰۰ سند دارای تگ <keyword> که نمایانگر موضوع خبر است، انتخاب شد و سایر اسناد از مجموعه، خارج شدند. از تعداد مذکور سند، ۵,۷۷۶ موضوع خنجر منحصر بفرد حاصل شد. به منظور کاهش ابعاد داده، تعداد ۱۰، keyword یا همان موضوع خبری پرتکرار انتخاب گردید و از هر موضوع خبری واحد به صورت تصادفی تعداد ۱۰۰۰ سند برگزیده شد. در نهایت، تعداد ۱۰,۰۰۰ سند، به عنوان مجموعه داده‌های این آزمایش در نظر گرفته شد.

#### ۴-۲- پیاده سازی و ارزیابی روش

پس از اجرای پیش پردازش اولیه و حذف کلمات توقف، تعداد ۱,۰۵۴,۱۰۶ واژه از مجموعه اسناد به دست آمد. لازم به ذکر است، ۶۶۷ کلمه توقف زبان انگلیسی موجود در سایت رایگان ranks مورد استفاده قرار گرفت. سپس با استفاده از الگوریتم ریشه یابی پورتز، ریشه هر واژه جایگزین آن گردید. در مرحله اصلی، وزن دهی به دو روش TF-IDF یعنی بدون در نظر گرفتن اهمیت محل قرارگیری واژه و به روش پیشنهادی TF-IDF\_P، با در نظر گرفتن اهمیت ۵ برابری ۲۵۰ کاراکتر اول سند، وزن هر واژه محاسبه گردید. قبل از ساختن ماتریس سند واژه، به منظور کاهش و تعدیل زمان محاسباتی، بر اساس روش یادگیری، سه حد آستانه وزن، با مقادیر ۰.۰۶ و ۰.۰۷ و ۰.۰۸ برای انتخاب واژه‌ها در نظر گرفته شد که حد آستانه ۰.۰۷ قابل قبول ترین نتیجه را حاصل کرد.

دو ماتریس سند واژه از اسناد مجموعه و با وزن‌های بزرگتر یا مساوی کران ۰.۰۷، برای واژه‌ها ایجاد گردید. یک ماتریس با اعمال روش وزن دهی TF-IDF و ماتریس دیگر با روش وزن دهی TF-IDF\_P حاصل شد. TA، ماتریس سند واژه با وزن دهی عادی، با تعداد ۲,۸۰۰ سند و ۳۴ واژه

دقت و فراخوانی در دسته بندی با روش پیشنهادی بهبود یافته است و با تعداد ۵ نزدیکترین همسایه بهترین عملکرد را دارد. بالاتر است. در اجرای الگوریتم K-means، میزان صحت در تعداد خوشه‌های ۸ و ۹، با هم برابر است. در مقایسه با روش اولیه، بهترین صحت خوشه‌بندی با تعداد ۶ خوشه، با بهبود حدود ۸٪، حاصل شده است.

### ۵- نتیجه

روش‌های وزن‌دهی واژه در محاسبه دقیق اهمیت واژه‌های اسناد به عنوان ویژگی‌های بارز برای نمایش سند، بسیار حائز اهمیت است. این امر تاثیر مستقیم بر عملکرد مشابهت سنجی اسناد و در نتیجه، دسته بندی و خوشه بندی متون دارد. در این مقاله، به بررسی طرح TF-IDF و عملکرد آن اشاره شد و بیان گردید که این طرح، اهمیت موقعیت مکانی یک واژه را در متن در نظر نمی‌گیرد، سپس طرح TF-IDF\_P به عنوان یک روش بهبود یافته، پیشنهاد شد و عملکرد آن توانست بر صحت خوشه بندی و دقت و فراخوانی دسته بندی اسناد مجموعه متون TDT5 بهبودی قابل توجه حاصل کند.

### ضمائم

$$\text{Accuracy} = \frac{TP + TN}{TP + TN + FN + FP}$$

TP تعداد موارد دسته‌بندی مثبت صحیح  
TN تعداد موارد دسته‌بندی منفی صحیح  
FP تعداد موارد دسته‌بندی مثبت نادرست  
FN تعداد موارد دسته‌بندی منفی نادرست

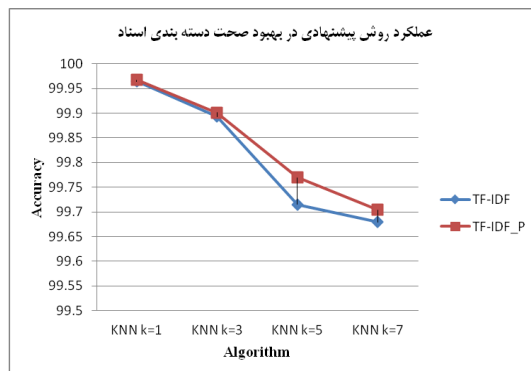
$$F - \text{measure} = \frac{2 \cdot \text{Precision} \cdot \text{Recall}}{\text{Precision} + \text{Recall}}$$

### مراجع

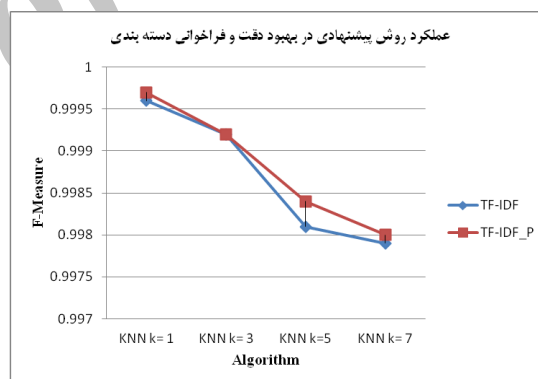
- [1] Singh, Jaskirat, and Mukesh Kumar. "A Meta Search Approach to Find Similarity between Web Pages Using Different Similarity Measures." *Advances in Computing, Communication and Control*. Springer Berlin Heidelberg, 150-160, 2011.
- [2] Liu, Bing. "Information retrieval and Web search." *Web Data Mining*. Springer Berlin Heidelberg, 211-268, 2011.
- [3] Roa-Valverde, Antonio J., and Miguel-Angel Sicilia. "A survey of approaches for ranking on the web of data." *Information Retrieval* 17, no. 4: 295-325, 2014.
- [4] Huang, Anna. "Similarity measures for text document clustering." *Proceedings of the sixth new zealand computer science research student conference (NZCSRSC2008)*, Christchurch, New Zealand. 2008.
- [5] Li, Yanjun, Congnan Luo, and Soon M. Chung. "A parallel text document clustering algorithm based on neighbors." *Cluster Computing* 18.2: 933-948, 2015.
- [6] Chim, Hung, and Xiaotie Deng. "Efficient phrase-based document similarity for clustering." *Knowledge and Data Engineering, IEEE Transactions on* 20.9: 1217-1229, 2008.

Accuracy یا صحت، جهت ارزیابی عملکرد، محاسبه گردید که مطابق (شکل ۲) و (شکل ۳)، بهبود حاصل از اجرای روش پیشنهادی مشهود است. سپس از طریق الگوریتم بیزین برای دسته‌بندی اسناد استفاده شد که در (شکل ۴) بهبود صحت مشخص شده است.

بهبود صحت دسته بندی بر حسب درصد، در تمام مراحل، به ویژه در شیوه اجرای الگوریتم مذکور با تعداد ۵ نزدیکترین همسایه از همه موارد بیشتر است. و این عدد از ۹۹.۷۱٪ به ۹۹.۷۷٪ افزایش یافته است. در تعداد همسایگان ۱ و ۳ کمترین تغییر، حاصل شده است.

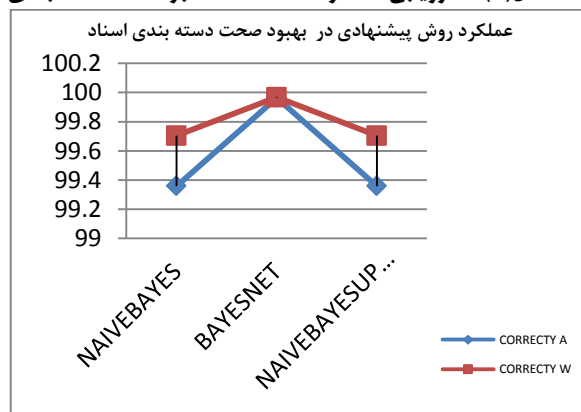


شکل (۲): ارزیابی عملکرد TF-IDF\_P در صحت دسته‌بندی



شکل (۳): ارزیابی عملکرد TF-IDF\_P در صحت دسته‌بندی

شکل (۴): ارزیابی عملکرد TF-IDF\_P بر صحت دسته‌بندی



- [23] Yoo, Jong-Yeol, and Dongmin Yang. "Classification Scheme of Unstructured Text Document using TF-IDF and Naive Bayes Classifier." , 2015.
- [24] Lin, Yung-Shen, Jung-Yi Jiang, and Shie-Jue Lee. "A similarity measure for text classification and clustering." Knowledge and Data Engineering, IEEE Transactions on 26.7: 1575-1590, 2014
- [25] Krishna, RVV Murali, and Ch Satyananda Reddy. "Extractive Text Summarization Using Lexical Association and Graph Based Text Analysis." Computational Intelligence in Data Mining—Volume 1. Springer India, 261-272, 2016.
- [26] Kallimani, Jagadish S., K. G. Srinivasa, and B. Esvara Reddy. "Summarizing News Paper Articles: Experiments with Ontology-Based, Customized, Extractive Text Summary and Word Scoring." Cybernetics and Information Technologies 12.2: 34-50, 2012
- [27] Osman, Ahmed Hamza, et al. "An improved plagiarism detection scheme based on semantic role labeling." Applied Soft Computing 12.5: 1493-1502, 2012.
- [28] Huang, Faliang, et al. "Clustering web documents using hierarchical representation with multi-granularity." World Wide Web 17.1: 105-126, 2014.
- [29] Kumar, Jayant, Peng Ye, and David Doermann. "Learning document structure for retrieval and classification." Pattern Recognition (ICPR), 2012 21st International Conference on. IEEE, 2012.
- [7] Gomaa, Wael H., and Aly A. Fahmy. "A survey of text similarity approaches." International Journal of Computer Applications 68.13: 13-18, 2013.
- [8] Gipp, Bela. "Citation-based Document Similarity." Citation-based Plagiarism Detection. Springer Fachmedien Wiesbaden, 43-55, 2014.
- [9] Liu, Hongyan, et al. "Measuring similarity based on link information: A comparative study." Knowledge and Data Engineering, IEEE Transactions on 25.12: 2823-2840, 2013.
- [10] Zesch, Torsten, and Iryna Gurevych. "Wisdom of crowds versus wisdom of linguists—measuring the semantic relatedness of words." Natural Language Engineering 16.: 25-59, 2010.
- [11] Chahal, Premjeet, Monika Singh, and Sudhakar Kumar. "Ranking of web documents using semantic similarity." Information Systems and Computer Networks (ISCON), 2013 International Conference on. IEEE, 2013.
- [12] Machnik, Łukasz. "Documents Clustering Techniques." Annales UMCS Sectio AI Informatica 2.1: 401-411, 2015.
- [13] Hakim, Ari Aulia, et al. "Automated document classification for news article in Bahasa Indonesia based on term frequency inverse document frequency (TF-IDF) approach." Information Technology and Electrical Engineering (ICITEE), 2014 6th International Conference on. IEEE, 2014.
- [14] Zhang, Wen, Taketoshi Yoshida, and Xijin Tang. "A comparative study of TF\* IDF, LSI and multi-words for text classification." Expert Systems with Applications 38.3: 2758-2765, 2011.
- [15] Lan, Man, et al. "A comprehensive comparative study on term weighting schemes for text categorization with support vector machines." Special interest tracks and posters of the 14th international conference on World Wide Web. ACM, 2005
- [16] Huynh, Minh Chau, Pham Duy Thanh Le, and Trong Hai Duong, (2015), "Improved Vector Space Model TF/IDF Using Lexical Relations.", International Journal of Advanced Computer Research 5.21 334.P.
- [17] Wang, Xingheng, et al., (2012), "Text clustering based on the improved TFIDF by the iterative algorithm.", Electrical & Electronics Engineering (EEESYM), 2012 IEEE Symposium on. IEEE..
- [18] Wang, Xingheng, et al., (2012), "Text clustering based on the improved TFIDF by the iterative algorithm.", Electrical & Electronics Engineering (EEESYM), 2012 IEEE Symposium on. IEEE.
- [19] Xia, Tian, et al, (2012), "An improved Global Weight Function of Terms based on Pearson's Chi-square statistics.", Information Science and Control Engineering 2012 (ICISCE 2012), IET International Conference on. IET.
- [20] Handojo, Andreas, Adi Wibowo, and Yovita Ria. "Document Searching Engine Using Term Similarity Vector Space Model on English and Indonesian Document." Intelligence in the Era of Big Data. Springer Berlin Heidelberg, 165-173, 2015.
- [21] Zhang, Wen, Taketoshi Yoshida, and Xijin Tang. "A comparative study of TF\* IDF, LSI and multi-words for text classification." Expert Systems with Applications 38.3: 2758-2765, 2011.
- [22] Hakim, Ari Aulia, et al. "Automated document classification for news article in Bahasa Indonesia based on term frequency inverse document frequency (TF-IDF) approach." Information Technology and Electrical Engineering (ICITEE), 2014 6th International Conference on. IEEE, 2014.