

استخراج برچسب برای آگهی‌های وبسایت مبتنی بر درون‌سازی واژه

علی محمد زارع چاهوکی^{۱*}، محمدمهدی صدقیان^۲

^۱دانشیار دانشکده کامپیوتر، دانشگاه یزد، یزد، ایران

^۲دانشکده کامپیوتر، دانشگاه یزد، یزد، ایران

چکیده

در جامعه امروزی اخبار و آگهی‌ها، جایگاه به‌خصوصی در رشد و ترقی جامعه دارند. با مشخص کردن واژگان اصلی آگهی، می‌توان به مفهوم کلی آن پی برد. آماده‌سازی این واژگان به روش سنتی نیازمند صرف زمان و دانش تخصصی راجع به موضوع متن است. سایت ایده‌کاو، سامانه‌ای هست که به جمع‌آوری پیام‌ها و آگهی‌های تلگرام می‌پردازد. نیازمندی سامانه ایده‌کاو، استخراج کلمات کلیدی از آگهی‌های منتشر شده در تلگرام بوده است. کیفیت کلمات کلیدی استخراج شده، نقش بسزایی در بهبود سئو و آمار بازدید آگهی‌ها دارد. با استفاده از الگوریتم‌های درون‌سازی، می‌توان صحبت‌های محاوره‌ای و ساختار معنایی متن را استخراج کرد، از این رو در تشخیص کلمات کلیدی در آگهی‌های تلگرام که اغلب به صورت عامیانه منتشر می‌شوند، مفید واقع می‌شود. در این پژوهش با استفاده از داده‌های سامانه ایده‌کاو مدلی از روش‌های درون‌سازی پیاده‌سازی شده است. نوآوری استفاده شده در این پژوهش از ترکیب کردن روش‌های درون‌سازی واژه، بسامد کلمات و جایگاه کلمات ایجاد شده است. مدل درون‌سازی از کلمات دو کلمه‌ای ایجاد شده است. ایجاد مدل از کلمات دو کلمه‌ای، به این دلیل است که اغلب کلمات کلیدی از دو کلمه به بالا تشکیل شده‌اند. جهت نمایش بهتر ارزیابی‌ها، مدل آی‌کی (مدل پیشنهادی) با روش‌های آماری و روش‌های مبتنی بر گراف مقایسه شده است که نتایج به‌دست آمده نشان می‌دهد ترکیب مدل آی‌کی دو-گرم عملکرد بهتری در استخراج کلمات کلیدی نسبت به سایر روش‌ها به وجود آورده است.

مقاله پژوهشی

تاریخ دریافت:

۱۴۰۲/۰۸/۲۵

تاریخ پذیرش:

۱۴۰۳/۰۱/۲۱

کلیدواژه‌ها:

استخراج برچسب، بهینه‌سازی برای موتور جست‌وجو (سئو)، یادگیری عمیق، درون‌سازی واژه.

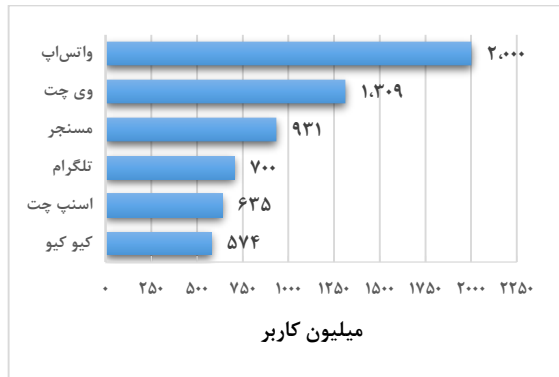
نویسنده مسئول:

chahooki@yazd.ac.ir

doi : 10.22034/ABMIR.2024.20879.1041

۱- مقدمه

بر اندروید، آی‌اواس، ویندوز، لینوکس و وب منتشر شده است. تلگرام با امکانات زیاد و کاربردی خود توانسته از سایر پیام‌رسان‌ها متمایز باشد و کاربران زیادی را به سوی خود جذب کند [۵]. در شکل ۱ به مقایسه محبوبیت برنامه‌های مختلف در حوزه شبکه‌های اجتماعی می‌پردازیم [۶].



شکل (۱): تعداد کاربران فعال ماهانه در ژانویه ۲۰۲۳

در این پژوهش از روش‌های مختلف آماری و درون‌سازی جهت استخراج کلمات کلیدی استفاده شده است. ارزیابی‌ها نشان می‌دهد که ترکیب کردن روش‌های آماری به همراه مدل شخصی وردتووک که از واژگان کلیدی دو کلمه‌ای ایجاد شده است، از کیفیت بالاتری برخوردار است. در این پژوهش نام استفاده شده برای روش پیشنهادی آی‌کی می‌باشد.

نوآوری استفاده شده در این پژوهش از ترکیب کردن روش‌های درون‌سازی واژه، بسامد کلمات و جایگاه کلمات ایجاد شده است. هدف ایجاد مدل از کلمات دو کلمه‌ای به این دلیل است که اغلب کلمات کلیدی از دو کلمه به بالا تشکیل شده‌اند. در آزمایش‌های انجام شده، مدل‌های ساخته شده از عبارت‌های دو کلمه‌ای، بهبود قابل توجهی نسبت به مدل‌های ساخته شده از عبارت‌های تک‌کلمه‌ای داشته‌اند.

تمرکز اصلی این مقاله بر ساخت مدل‌های درون‌سازی از عبارت‌های دو کلمه‌ای و مقایسه با روش‌های آماری برای استخراج کلمات کلیدی است. به همین جهت در بخش دوم نگاهی به پژوهش‌های انجام شده در زمینه استخراج کلمات کلیدی می‌اندازیم. در فصل سوم چگونگی پیاده‌سازی و روش پیشنهادی

در عصر جدید استفاده از روش‌های نوین، برای بازیابی اطلاعات امری مهم محسوب می‌شود. کلمات کلیدی، مجموعه‌ای از کلمات مهم در سند هستند که محتوای سند را توصیف می‌کنند. کلمات کلیدی اطلاعات نحوی مفیدی را برای پردازش متن فراهم می‌کنند. استخراج کلمات کلیدی، فرایندی خودکار برای تشخیص کلمات مهم به کاررفته در سند است. این کلمات، معنا و مفهوم سند را برای خواننده آشکار می‌کنند و می‌توانند در بسیاری از کارهای زبان طبیعی مانند تجزیه و تحلیل متون و خلاصه‌سازی استفاده شوند [۱].

کلمات کلیدی را می‌توان به‌عنوان برچسب برای اسناد و دسته‌بندی محتوا و محصولات سایت استفاده کرد. کلمات کلیدی، برچسب‌هایی هستند که سند یا محصول را به بهترین شیوه ممکن برای کاربران معرفی می‌کنند.

یکی از موارد مؤثر در سئو، کلمات کلیدی و برچسب‌ها هستند. از نظر سئو، کلمات کلیدی عبارت‌هایی هستند که کاربران در موتورهای جست‌وجو وارد می‌کنند. انتخاب کلمات کلیدی مناسب برای هر داده، موجب رضایت کاربران و افزایش سئو سایت می‌شود [۲-۴].

اکثر پژوهش‌های موجود در زمینه استخراج کلمات کلیدی، برای رسانه اجتماعی توئیتر، اسناد علمی و اسناد خبری صورت گرفته است و در حوزه آگهی‌های تلگرام، پژوهش‌های اندکی انجام شده است. تلگرام در ایران یکی از پرکاربردترین شبکه‌های اجتماعی است که روزانه حجم زیادی از اطلاعات را منتشر می‌کند. سایت ایده‌کاو سامانه‌ای هست که به جمع‌آوری پیام‌ها و آگهی‌های تلگرام می‌پردازد. با توجه به نیازمندی سامانه ایده‌کاو در استخراج کلمات کلیدی از آگهی‌های منتشر شده در تلگرام، در این پژوهش به دنبال ارائه روشی هستیم که بتواند کلمات کلیدی موجود در آگهی‌های تلگرام را شناسایی و استخراج کند. کیفیت کلمات کلیدی استخراج شده، نقش بسزایی در بهبود سئو و آمار بازدید آگهی‌ها دارد.

پیام‌رسان تلگرام، برنامه‌ای رایگان و مبتنی بر رایانش ابری است. تلگرام پیام‌رسانی است که امکان ارسال پیام به صورت متن، صوت، تصویر و ویدئو را امکان‌پذیر می‌کند. تلگرام برای دستگاه‌های مبتنی

- مرتب‌سازی عبارات بر اساس امتیازات
- انتخاب تعداد مشخصی از عبارات کاندیدا که دارای بیشترین یا کمترین وزن هستند

روند کلی رویکردهای آماری در شکل ۳ نشان داده شده است.



شکل (۳): روند کلی رویکردهای آماری

وون و همکاران برای استخراج کلمات کلیدی از ویژگی‌های فراوانی کلمه، TF-IDF، جایگاه کلمه و طول کلمه استفاده کرده‌اند. در ویژگی طول کلمه، برای عبارات تک کلمه‌ای امتیاز یک و برای عبارات چند کلمه‌ای امتیاز دو را در نظر گرفته‌اند. امتیاز هر عبارت کاندیدا از حاصل ضرب چهار ویژگی ذکر شده به دست می‌آید. در مرحله آخر عبارت‌هایی که امتیاز بیشتری داشته باشد به عنوان کلمات کلیدی انتخاب می‌شوند. تعداد کلمات کلیدی انتخاب شده برای هر سند از رابطه (۱) محاسبه می‌شود که عبارت بیانگر تعداد کلمات سند است [۸].

$$N = 2.5 \times \log_{10}^{d_1} \quad (1)$$

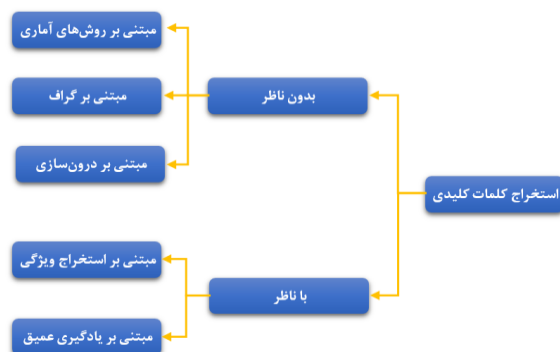
فلورسکو^۱ و همکاران، روشی مبتنی بر گراف و بدون ناظر، برای استخراج کلمات کلیدی پیشنهاد کرده‌اند. آنها گرافی از کلمات ایجاد کرده‌اند که گره‌های آن را کلمات و ارتباط بین کلمات را یال‌ها تشکیل داده‌اند. جهت امتیازدهی، از ویژگی‌های فراوانی کلمه، جایگاه کلمه و هم رخدادی کلمات استفاده کرده‌اند. با اعمال الگوریتم پیچ‌رنک وزن‌دار بر گراف، گره‌هایی که دارای بیشترین وزن بوده‌اند به عنوان کلمات کلیدی معرفی شده‌اند. آنها در مقایسه با چند روش بدون ناظر مانند تکست‌رنک، اکسپندرک، سینگل‌رنک و تاپیک‌رنک از عملکرد بهتری برخوردار بوده‌اند [۹].

مطرح می‌شود. در بخش چهارم آزمایش‌های تجربی بر روی دادگان مطرح می‌شود و در فصل پنجم به جمع‌بندی مطالب و ذکر پیشنهادها برای پژوهش‌های آینده می‌پردازیم.

۲- پیشینه پژوهش

روش‌های متعددی برای استخراج واژگان کلیدی وجود دارد. طبق دسته‌بندی انجام شده، رویکردهای استخراج عبارات کلیدی به دو دسته مبتنی بر یادگیری بدون ناظر و مبتنی بر یادگیری با ناظر تقسیم می‌شود [۷].

تقسیم‌بندی روش‌های استخراج واژگان کلیدی در شکل ۲ آمده است:



شکل (۲): دسته‌بندی روش‌های استخراج عبارات کلیدی

رویکردهای بدون ناظر اغلب مستقل از زبان و دامنه هستند و به مجموعه‌داده‌گان آموزشی برچسب‌دار نیازی ندارد. در مقابل، رویکردهای با ناظر دارای کیفیت بهتری نسبت به رویکرد بدون ناظر هستند؛ ولی به دادگان برچسب‌دار نیازمند می‌شوند.

۲-۱ رویکردهای آماری

در رویکردهای آماری، عبارات بر اساس اطلاعات آماری استخراج شده از متن به دست می‌آید. در این رویکرد کلمات کاندیدا با استفاده از شیوه‌های مختلف، رتبه‌بندی می‌شود. روند کلی استخراج عبارات کلیدی با استفاده از رویکردهای آماری، شامل موارد زیر می‌باشند:

- پیش‌پردازش متن
- استخراج ویژگی‌های آماری
- امتیازدهی عبارات

¹ Florescu

۲-۲ رویکردهای مبتنی بر درون‌سازی

بنجیو^۱ و همکارانش در سال ۲۰۰۳ عبارت درون‌سازی واژه را ابداع کرده‌اند. کلبرت^۲ و وستون اولین کسانی بودند که توانایی درون‌سازی واژه‌های از پیش آموزش‌دیده را در سال ۲۰۰۸ نشان دادند. آنها نه تنها درون‌سازی واژه را به‌عنوان یکی از ابزارهای سودمند برای پردازش زبان طبیعی، تثبیت کردند، بلکه معماری شبکه عصبی‌ای معرفی کردند که اساس بسیاری از روش‌های امروزی است. پیشرفت درون‌سازی واژه را می‌توان به کار میکولوف^۳ و همکارانش در سال ۲۰۱۳ که وُردتووک را ارائه کردند منتسب کرد. در سال ۲۰۱۴، پنینتگون^۴ و همکارانش الگوریتم بردار سراسری را ارائه کردند و با این کار نشان دادند درون‌سازی واژه وارد استفاده روزمره شده است [۱۰].

روش درون‌سازی واژه، یک روش شبکه عصبی پیش‌خور است که با استفاده از یک لایه پنهان، توانایی تولید بردار برای کلمات را دارد. این روش با هدف کاهش ابعاد و افزایش کارایی در پردازش زبان طبیعی طراحی شده است. روش درون‌سازی با شبکه عصبی بازگشتی متفاوت است، زیرا در شبکه عصبی بازگشتی، خروجی هر گره، به ورودی گره‌های بعدی منتقل می‌شود و حافظه‌ای از داده‌های قبلی را نگه می‌دارد. درحالی‌که در شبکه عصبی پیش‌خور، اطلاعات فقط در یک جهت، از لایه ورودی به لایه خروجی حرکت می‌کنند و هر گره فقط به گره‌های لایه بعد خود متصل می‌شود.

یکی از مهم‌ترین وظایف در پردازش زبان طبیعی، تبدیل متن به حالتی قابل فهم برای پردازش است. برای نمایش معنا و مفهوم متن، روش‌های درون‌سازی واژه معرفی شده‌اند. در الگوریتم‌های پیش آموزش‌دیده، نمایش برداری کلمات شامل روابط معنایی، نحوی و منطقی بین کلمات است [۱۱]. در درون‌سازی، کلمات مشابه دارای بردارهای مشابه هستند و در فضای برداری در فاصله نزدیک به هم قرار می‌گیرند. برای مثال کلمات ملک، خانه و منزل هر سه در فضایی نزدیک به هم قرار می‌گیرند. در درون‌سازی واژگان، معنای واژه با استفاده از کلمات اطرافش به دست می‌آید [۱۲].

درون‌سازی واژه از مدل‌های مختلفی جهت ایجاد بردار و پردازش کلمات استفاده می‌کند. روش‌های مبتنی بر شبکه‌های عصبی عموماً پیش‌بینی محور هستند که در آنها مدل‌ها برای پیش‌بینی کلمات محتوایی با استفاده از کلمات وسطی، یا برعکس پیش‌بینی کلمه وسط توسط مجموعه‌ای از کلمات محتوایی ساخته می‌شوند. در این الگوریتم، دو معماری مدل کیسه کلمات پیوسته و مدل اسکپ-گرام ارائه شده است:

- کیسه کلمات پیوسته: واژه فعلی با توجه به زمینه مورد بحث متن، پیش‌بینی می‌شود و برای کلمات پر تکرار، بردارهای بهتری ارائه می‌دهد.
- اسکپ-گرام: واژه فعلی بر اساس کلمات اطرافش پیش‌بینی می‌شود و با داده‌های کوچک بهتر کار می‌کند. این مدل برای کلمات کمیاب، نمایش برداری بهتری نسبت به مدل کیسه کلمات پیوسته ارائه می‌دهد.

هارمان^۵ و همکاران، برای استخراج کلمات کلیدی، روشی بر پایه امبد رنگ طراحی کرده‌اند. آنها علاوه بر محاسبه شباهت کسینوسی میان بردار سند درون‌سازی شده و بردار کلمه درون‌سازی شده، از عاملی به نام پاداش وزنی جهت امتیازدهی استفاده می‌کنند. برای محاسبه پاداش وزنی، برای کلماتی که در عنوان سند آمده‌اند امتیاز بیشتری نسبت به کلماتی که در محتوا سند آمده‌اند در نظر گرفته می‌شود [۱۳].

بنانی اسمیرس^۶ و همکاران، روشی به نام امبد رنگ پیشنهاد کرده‌اند. آنها برای شناسایی کلمات کلیدی، الگوریتمی با استفاده از برچسب‌زن ادات سخن طراحی کرده‌اند. الگوریتم استخراج کلمات کلیدی، از یک یا چند اسم که به همراه یک یا چند صفت آمده است، تشکیل می‌شود. آنها از دو الگوریتم درون‌سازی داک‌تووک^۷ و سنت‌تووک^۸ برای نمایش برداری کلمات کلیدی استفاده کرده‌اند. در مرحله آخر شباهت کسینوسی بین بردار سند درون‌سازی شده و بردار کلمات کلیدی نامزد، محاسبه شده و مشابه‌ترین کلمات، استخراج می‌شوند [۱۴].

⁵ Haarman

⁶ Bennani-Smires

⁷ Doc2Vec

⁸ Sent2vec

¹ Benjio

² Collobert

³ Mikulov

⁴ Jeffrey Pennington

۳- روش پیشنهادی

هدف اصلی پژوهش استخراج کلمات کلیدی از آگهی می‌باشد. این واژگان از ترکیب‌های دوکلمه یا بیشتر تشکیل شده‌اند که در این مقاله به تشخیص واژگان دوکلمه‌ای پرداخته شده است. نوآوری استفاده شده در این پژوهش، ساخت مدل اختصاصی وُردتووک و ترکیب آن با روش‌های آماری مانند جایگاه کلمه و بسامد کلمه در سند می‌باشد. مدل اختصاصی روی کلمات دوتایی از آگهی‌های سایت ایده‌کاو، آموزش دیده است. در این بخش ابتدا به معرفی دادگان و پیش‌پردازش می‌پردازیم. در بخش ۳-۳ مراحل استخراج کلمات کلیدی و در بخش ۳-۴ نحوه آموزش وُردتووک شرح داده می‌شود. در بخش‌های ۳-۵ الی ۳-۷ روش‌های ترکیب شده جهت امتیازدهی کلمات توضیح داده خواهد شد و در انتها شبه کدی از کل روش‌های پیاده‌سازی شده بیان می‌شود.

۳-۱ مجموعه داده

در این پژوهش از مجموعه آگهی‌های سامانه ایده‌کاو به تعداد ۲۰۸،۰۳۸ عدد که از گروه‌ها و کانال‌های تلگرامی جمع‌آوری شده‌اند، استفاده شده است. برای بالابردن کیفیت ارزیابی نتایج، فقط از آگهی‌های فارسی و حوزه املاک استفاده کرده‌ایم. جهت پالایش، فقط آگهی‌های حوزه املاک انتخاب شده‌اند. پس از پالایش، تعداد آگهی‌های حوزه املاک به تعداد ۱۰،۵۳۴ عدد رسیده است.

۳-۲ پیش‌پردازش

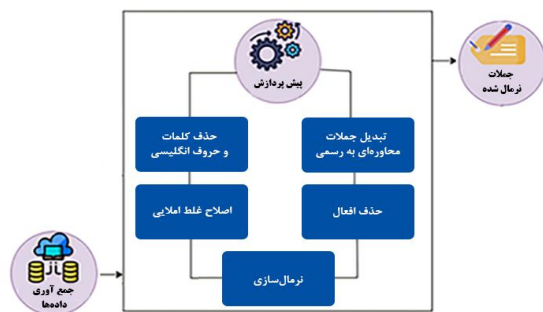
در پیام‌رسان تلگرام، افراد عموماً به سبک محاوره‌ای و مکالمه‌ای تبادل اطلاعات می‌کنند، از این رو پیش‌پردازش یکی از عوامل مؤثر در افزایش کیفیت داده‌ها محسوب می‌شود [۱۵]. پردازش زبان فارسی تا حدودی با پردازش زبان انگلیسی متفاوت است. در زبان انگلیسی هر حرف و هر کلمه، جدا و به شکلی خاص نوشته می‌شود، در حالی که در فارسی برخی از حروف به هم وصل و برخی جدا نوشته می‌شوند. اگر حروف، علائم نگارشی و کلمات فارسی به شکل‌های متفاوتی نوشته شوند، سیستم پردازشی قادر به تشخیص یکسان بودن کلمات نخواهد بود. برای مثال در صورت عدم پیش‌پردازش، سیستم دو کلمه «خانه» و «خونه» را

به‌عنوان دو کلمه مجزا در نظر می‌گیرد که به طور قابل توجهی نتایج را تحت تأثیر قرار می‌دهد [۱۶].

جهت پیش‌پردازش متون، از کتابخانه متن‌کاوی استفاده شده است. در این پیش‌پردازش اقدامات زیر انجام شده‌اند:

- حذف کلمات و حروف انگلیسی، اعداد، آدرس ایمیل و شماره موبایل
- حذف علائم نگارشی و شکلک‌ها
- حذف فاصله‌های اضافی
- حذف لینک‌های تبلیغات
- تبدیل جمله‌های محاوره‌ای به رسمی
- اصلاح غلط‌های املائی
- حذف افعال

شکل ۴ نمای کلی از گردش کار در مرحله پیش‌پردازش را نشان می‌دهد.



شکل (۴): گردش کار استخراج واژگان کلیدی

۳-۳ مراحل استخراج کلمات کلیدی

در این پژوهش، روشی برای استخراج کلمات کلیدی بر پایه وُردتووک ارائه شده است. سامانه پیشنهادی شامل شش مرحله اصلی پیش‌پردازش، آموزش اختصاصی وُردتووک روی کلمات دوتایی، محاسبه احتمال رخ دادن کلمه در متن، محاسبه فاصله کلمات، محاسبه جایگاه کلمه در متن و سپس تعیین امتیاز برای هر کلمه هست که روال کلی آن در شکل ۵ نشان داده شده است.

اندازه بردار برای هر کلمه در این مدل برابر با ۴۰۰ می‌باشد. نمونه‌ای از مدل ساخته شده در شکل ۶ قابل مشاهده است.

	۰	۱	۲	۳	۴	۵	۶	۷	۸	۹	۱۰	۱۱	۱۲	۱۳	۱۴
آن روز	0.0942	0.1791	0.0753	0.1816	-0.1711	-0.0666	0.0326	0.5454	-0.0017	0.1600	-0.1053	0.3779	-0.1819	0.2179	0.1119
سایر جوانان	0.0197	0.2624	0.0291	0.2594	0.0071	0.4883	0.1706	0.0671	0.0714	0.1817	0.0219	0.0705	0.1076	0.2094	0.1119
مشاوران	0.0089	0.1942	0.1629	-0.0103	0.0282	-0.2856	0.2840	0.3981	0.1963	0.1401	-0.2887	-0.0175	-0.0670	0.0123	0.0187
فرمان‌بردار	-0.0085	0.3078	0.0279	-0.0118	0.0540	-0.1078	0.2740	0.3904	0.1488	-0.0407	-0.2041	0.1709	0.1784	0.1030	0.1119
رغم سوسوم	0.0746	0.2257	0.0210	0.2580	0.1474	0.0865	0.3852	0.2034	0.2767	0.1212	-0.2842	0.0828	-0.2467	0.1458	0.0819
فرمانت کلیت	-0.0748	0.3191	0.0425	-0.1605	-0.0272	-0.4785	0.0705	0.0799	0.0392	-0.0106	-0.1779	-0.2107	-0.2898	0.2299	0.2119
الفرانس	0.0084	0.1808	0.0711	0.0770	0.0259	0.2246	0.0203	0.0421	0.0182	0.0805	0.3684	0.1709	0.3078	0.1053	0.1119
مطهر بودن	0.2047	0.1801	0.0443	0.4960	0.2639	-0.2133	0.1716	0.2223	-0.1864	0.2669	-0.0714	-0.0172	-0.0159	0.1512	0.1119
مطهر بودن	0.0292	0.2792	0.2151	0.1783	0.0181	-0.0264	0.3783	0.1514	0.0718	0.0708	-0.2708	0.0037	0.0427	0.0445	0.0119
مدت تک	-0.0089	0.4848	0.0536	0.1030	0.1134	0.3480	0.1934	0.3461	0.0915	0.1709	0.0804	0.3669	0.1059	0.0423	0.1119
مدت تک برگ	0.0653	0.5047	0.0169	-0.0269	0.0146	-0.5345	0.1458	0.0409	-0.0871	-0.1139	-0.5336	0.3535	-0.1628	-0.2569	0.0119
تک برسان گفت	-0.0441	0.4839	0.0871	-0.0944	-0.1424	0.0652	0.0789	0.2793	0.0714	0.2200	0.0216	-0.2762	0.2240	0.1119	0.1119
مطهری املک	0.1091	0.1912	0.0570	0.2280	0.0421	0.1815	0.1768	0.1845	-0.2752	0.4209	0.0157	0.2889	0.0089	0.1077	0.1119
موتور	-0.0189	0.3312	0.0590	0.0749	-0.2150	0.0213	0.4239	-0.1896	-0.0519	-0.2824	0.1591	0.1059	0.1959	0.1119	0.1119
کلمه‌ها	0.0120	0.2771	0.0289	0.2071	-0.1093	0.0203	0.1033	-0.2859	0.0403	0.1978	0.2209	0.0089	0.0070	0.2009	0.1119
رغم سوسوم	0.1147	0.1737	0.1479	-0.2010	0.0540	-0.3432	0.0840	0.0840	0.0112	-0.0473	0.0071	-0.2178	-0.4128	0.1119	0.1119
فرمانت تک	-0.1000	0.2739	0.1463	0.1943	0.0494	-0.0645	0.0212	0.0151	-0.0641	0.1937	-0.1193	-0.1000	-0.1486	0.0767	0.1119

شکل (۶): مدل برداری وُردتووک بر روی مجموعه داده

۳-۵ احتمال فراوانی کلمات در متن

کاربران معمولاً در نوشته‌های خود از کلمات تکراری استفاده می‌کنند. دلیل استفاده از فراوانی کلمات به این دلیل است که کاربران معمولاً برای راحتی و تسریع در نوشتار از کلمات خاصی برای نوشتن یا صحبت‌های خود استفاده می‌کنند. در این پژوهش فراوانی کلمات دوتایی را شمارش شده است. محاسبه فراوانی کلمات یکی از ویژگی‌های مؤثر در انتخاب کلمات کلیدی به شمار می‌رود.

احتمال فراوانی کلمه، به‌تنهایی نمی‌تواند در استخراج کلمات کلیدی تأثیرگذار باشد و باید با سایر روش‌های ذکر شده ترکیب شود تا بهترین نتیجه ایجاد شود.

جهت محاسبه احتمال فراوانی کلمه از رابطه (۲) استفاده می‌کنیم که عبارت $count(w)$ بیانگر تعداد تکرار عبارت و C بیانگر تعداد کل کلمات متن است.

$$f(w) = count(w)/c \quad (2)$$

۳-۶ امتیاز جایگاه کلمه در متن

اکثر افراد برای تأثیرگذاری بیشتر آگهی‌شان در ابتدای متن خود، کلمات مهم‌تر و پایه‌ای را می‌نویسند، به همین جهت ترکیب کردن امتیاز جایگاه کلمات با روش‌های ذکر شده باعث بهبود عملکرد می‌شود. به طور مثال به متن آگهی زیر توجه کنید:

«پراید دوگانه‌سوز فروشی. بدون رنگ. با شماره زیر تماس بگیرید.»

می‌بینیم که برچسب پراید دوگانه‌سوز در ابتدای جمله حاضر شده است و رفته‌رفته بر تأثیر کلمات دوتایی کاسته می‌شود.



شکل (۵): روش پیشنهادی برای استخراج کلمات کلیدی

در روش پیشنهادی بعد از پیش‌پردازش مجموعه داده‌ها، مدل وُردتووک روی کلمات دوتایی هر آگهی ساخته و میانگین شباهت کسینوسی هر کلمه با سایر کلمات محاسبه شده است. یکی دیگر از معیارهای مهم در شناسایی کلمات کلیدی، تعداد تکرار کلمات است اما این مورد به‌تنهایی نقشی در انتخاب کلمه کلیدی ندارد، زیرا ممکن است یک کلمه با بسامد کم، معنی اصلی متن را مشخص کند. به همین دلیل از مدل وُردتووک برای ساخت بردار تعبیه کلمات استفاده شده است. این مدل، معنا و مفهوم کلمات را به همراه وابستگی بین کلمات نیز درک می‌کند. در استخراج کلمات کلیدی، کلمه‌ای که فاصله کمتری با سایر کلمات دارد و در اطراف آن تراکم کلمات بیشتر است، انتخاب می‌شود. به همین علت روش پیشنهادی روشی پویا است و قادر به استخراج هر نوع کلمه کلیدی است. در این پژوهش برای بهبود استخراج کلمات کلیدی، با روش‌های آماری نیز ترکیب شده است که در ادامه به توضیح هر کدام می‌پردازیم.

۳-۴ آموزش وُردتووک روی مجموعه داده

در این قسمت به نحوه آموزش مدل وُردتووک روی کل آگهی‌های املاک می‌پردازیم. مدل‌های مرسوم ساخته شده همگی روی عبارت‌های تک‌کلمه‌ای آموزش دیده‌اند. مدلی که در این پژوهش آماده کرده‌ایم روی عبارت‌های دو کلمه‌ای آموزش دیده است. جهت ایجاد عبارت‌های دو کلمه‌ای از کتابخانه شمارش بردار^۱ استفاده کرده‌ایم. بعد از تقسیم جملات به عبارت‌های دو کلمه‌ای، مدل وُردتووک را با کمک کتابخانه جنسیم^۲ آموزش داده‌ایم. مدل ساخته شده بر اساس حالت اسکپ-گرام آموزش دیده است.

² Gensim

¹ CountVectorizer

Pseudo-code of Keyword Extraction

Input: U (sentence)

Output: R (list of keywords)

Variables:

var u_i : each bigram word in U

var M : model of word2vec

var F : frequency of bigram

var P : position of bigram

var D : dictionary of bigram with sum of score

var S : similarity of vector

var $score$: final score for bigram

Pseudo-code:

for each bigram u_i in U :

$P=1/(position+1)$ of u_i in U ;

$F=count(u_i)/count$ of all word in U ;

set $S=0$;

for each bigram u_j in U :

if u_i in M AND u_j in M :

set $S += cosine$ similarity (u_i, u_j);

set $score = S + F + P$;

$D.append(\{u_i: Score\})$;

set $R = sort(D)$;

return R ;

شبه‌کد (۱): روش استخراج کلمات کلیدی

۴- آزمایش‌های تجربی

برای حفظ اطلاعات و مدیریت بهتر نتایج، وب‌سرویس بر پایه جنگو طراحی شده‌است. جهت ارزیابی، از روش‌های آماری همچون TF-IDF، YAKE، عبارت اول، تکست‌رنک، پوزیشن‌رنک و تاپیک‌رنک و از روش‌های درون‌سازی همچون پارس‌برت، ویکی‌پارسی و مدل‌آی‌کی استفاده شده‌است. جهت ارزیابی مدل‌آی‌کی، جملات به عبارت‌های دوتایی و جهت ارزیابی سایر روش‌ها، جملات به تعداد کلمات تشکیل‌دهنده تقسیم شده‌اند. برای ساخت مدل‌آی‌کی، به آموزش کلمات کلیدی دو کلمه‌ای و تک کلمه‌ای پرداخته شده است که نتایج مدل دو کلمه‌ای، بهبود قابل توجهی نسبت به مدل تک کلمه‌ای داشته است. جهت سهولت و دقت در ارزیابی، داده‌ها در پایگاه‌داده ذخیره شده و بعد از پیش‌پردازش، ارزیابی روی آن‌ها انجام شده‌است.

در این پژوهش از معکوس جایگاه هر عبارت دوکلمه‌ای، استفاده شده است. این عملکرد از رابطه (۳) محاسبه می‌شود.

$$p(w) = \frac{1}{position(w) + 1} \quad (3)$$

در رابطه (۳) مقدار $position(w)$ به جایگاه عبارت دو کلمه‌ای در متن اشاره می‌کند. این مقدار از ۱ شروع می‌شود و به تعداد کلمات موجود در متن ادامه پیدا می‌کند.

۳-۷ محاسبه شباهت کلمات

پس از محاسبه جایگاه کلمات و احتمال فراوانی عبارت دو کلمه‌ای در متن، شباهت عبارت با تمامی عبارت‌های موجود در متن، با استفاده از بردار وردتووک و شباهت کسینوسی، محاسبه می‌شود. جهت محاسبه میانگین شباهت کسینوسی عبارت با تمام عبارت‌های موجود در متن، از معادله (۴) استفاده می‌شود.

$$sim(w) = \frac{\sum_i^n (1 - (w \cdot w_i) / |w| \cdot |w_i|)}{n} \quad (4)$$

شباهت کسینوسی، برای تعیین شباهت اسناد استفاده می‌شود. از نظر ریاضی شباهت کسینوسی، کسینوس زاویه بین دو بردار پیش‌بینی‌شده در یک فضای چندبعدی است. در رابطه (۴) مقدار i به‌ازای تمام عبارت‌های متن ورودی تغییر می‌کند. مقدار w بیانگر کلمات دوتایی هست که قصد داریم امتیاز آن را محاسبه کنیم و w_i بیانگر سایر عبارت‌های منتخب است. مقدار n در این الگوریتم بیانگر تعداد کل عبارت‌های موجود در جمله است.

پس از محاسبه $sim(w)$ ، $f(w)$ و $p(w)$ امتیاز کلمه w مطابق رابطه (۵) محاسبه می‌شود.

$$score(w) = f(w) + p(w) + (2 * sim(w)) \quad (5)$$

جهت تأثیرگذاری بیشتر مقدار وردتووک، مقدار $sim(w)$ را در ۲ ضرب کرده‌ایم. امتیاز تمامی عبارت‌ها محاسبه شده و بر اساس بیشترین امتیاز مرتب می‌شوند. در این پژوهش از حد آستانه‌های ۵، ۱۰ و ۱۵ برای انتخاب کلمات کلیدی استفاده شده است. برای بیان واضح‌تر مراحل ذکر شده، تمامی الگوریتم‌ها در شبه‌کد ۱ بیان شده‌اند.

۱-۴ معیارهای ارزیابی

در این قسمت به معرفی معیارهای ارزیابی می‌پردازیم. معیارهای ارزیابی در اغلب تحقیقات شامل دقت، صحت، فراخوانی و امتیاز F1 هستند. این معیارها توسط چهار مؤلفه ماتریس اغتشاش محاسبه می‌شوند. هر یک از چهار مؤلفه عبارت‌اند از: مثبت‌های درست، منفی‌های درست، مثبت‌های نادرست و منفی‌های نادرست. با توجه به رابطه (۶)، جهت محاسبه مقدار دقت، نسبت عبارت‌هایی که به‌درستی به‌عنوان کلمه کلیدی تشخیص داده شده‌اند به همه عبارت‌هایی که به‌عنوان کلمه کلیدی تشخیص داده شده‌اند، محاسبه می‌شود.

$$(6) \quad \frac{TP}{TP + FP}$$

نسبت عبارت‌هایی که به‌درستی به‌عنوان کلمه کلیدی تشخیص داده شده‌اند به کل کلمات کلیدی با عنوان مقدار فراخوانی، از رابطه (۷) محاسبه می‌شود.

$$(7) \quad \frac{TP}{TP + FN}$$

در نهایت جهت محاسبه امتیاز F_1 از رابطه (۸) استفاده می‌کنیم. این رابطه مشخص می‌کند دقت و فراخوانی در مجموع به چه صورت هستند.

$$(8) \quad \frac{2 \times Precision \times Recall}{Precision + Recall}$$

علاوه بر معیار امتیاز F1، معیار افت همینگ نیز محاسبه شده است. معیار افت همینگ، معیار ارزیابی برای داده‌های چند برچسبی است. این معیار با این فرض که هر برچسب یک اهمیت مساوی دارد، نسبت برچسب‌هایی را که به اشتباه پیش‌بینی شده‌اند به کل تعداد برچسب‌ها محاسبه می‌کند. به عبارتی دیگر، نمونه‌ای که به برچسب نادرست نسبت داده شده باشد یا برچسب صحیحی که به یک نمونه متعلق باشد اما پیش‌بینی نشده باشد نسبت به کل برچسب‌های موجود محاسبه می‌شود.

این معیار برای هر نمونه ورودی تعداد برچسب‌های نادرست را شمارش کرده و سپس با تقسیم بر تعداد کل برچسب‌ها، خطای همینگ را به‌عنوان افت همینگ ارائه می‌دهد. فرمول این معیار از رابطه (۹) محاسبه می‌شود.

$$(9) \quad Hamming Loss(y, \hat{y}) = \frac{1}{N} \sum_{i=1}^N \frac{|y_i \Delta \hat{y}_i|}{L}$$

در رابطه (۹) عملگر N تعداد کل نمونه‌ها و L تعداد کل برچسب‌ها را نمایش می‌دهد. $|y_i \Delta \hat{y}_i|$ بیانگر برچسب واقعی و \hat{y}_i بیانگر برچسب پیش‌بینی شده نمونه i ام است. $|y_i \Delta \hat{y}_i|$ تابع نشانگر است که اگر برچسب‌های واقعی و پیش‌بینی شده برای نمونه i ام یکسان باشند ۱ را برمی‌گرداند و در غیر این صورت ۰ را برمی‌گرداند.

۲-۴ معرفی روش‌های استفاده شده در ارزیابی

برای ارزیابی کیفیت برچسب‌زنی اسناد و ارزیابی هر یک از الگوریتم‌های درون‌سازی، ۴۰۰ سند به‌صورت تصادفی انتخاب و توسط فرد خبره برچسب‌گذاری شده است. برای هر آگهی، برچسب‌ها به‌صورت چند برچسبی انتخاب شده‌اند. در مدل چند برچسبی، برای هر آگهی می‌توان به تعداد برچسب‌های موردنیاز، برچسب اضافه کرد. به طور مثال برای آگهی «فروش املاک با نازل‌ترین قیمت. جهت هماهنگی، با املاکی سهند تماس حاصل فرمایید»، دو برچسب فروش املاک و املاکی سهند انتخاب می‌شود. جهت ارزیابی، روش آی‌کی با روش‌های بدون ناظر استخراج کلمات کلیدی همچون TF-IDF، YAKE، عبارت اول، تکست‌رنک، پوزیشن‌رنک، تاپیک‌رنک، برت و وردتووک مقایسه شده است که نتایج روش پیشنهادی عملکرد بهتری را ثبت کرده است.

۳-۴ ارزیابی روش‌های استخراج کلمات کلیدی

جهت بررسی عملکرد هر یک از روش‌های ذکر شده، آزمایش‌هایی را روی آگهی‌های تلگرام انجام داده‌ایم. ارزیابی انجام شده با استفاده از دو معیار امتیاز F1 و افت همینگ انجام شده است. جهت انجام آزمایش‌ها از سه آستانه ۵، ۱۰ و ۱۵ استفاده کرده‌ایم. در جدول ۱ درصد امتیاز ماکرو-F1 روش‌های مختلف ثبت شده است. با توجه به نتایج، مدل آی‌کی دو کلمه‌ای (وردتووک) عملکرد بهتری در مقایسه با روش‌های درون‌سازی دیگر و روش‌های آماری داشته است.

جدول (۲): مقایسه‌ی افت همینگ روش‌های پیشنهادی

افت همینگ			نام روش
p@15	p@10	p@5	
۰/۰۲۴۲	۰/۰۲۴۲	۰/۰۲۴۲	TF-IDF
۰/۰۲۷۱	۰/۰۲۷۱	۰/۰۲۷۱	عبارت اول
۰/۰۲۶۹	۰/۰۲۶۹	۰/۰۲۷۳	تکسترنک
۰/۰۲۶۷	۰/۰۲۶۷	۰/۰۲۶۷	پوزیشنرنک
۰/۰۲۶۸	۰/۰۲۶۸	۰/۰۲۶۸	تاپیکرنک
۰/۰۲۴۸	۰/۰۲۴۳	۰/۰۲۴۹	YAKE
۰/۰۱۸۷	۰/۰۱۸۴	۰/۰۱۹۲	مدل آی‌کی دو-گرم (وُردتووک)
۰/۰۲۰۱	۰/۰۱۹۵	۰/۰۲۲۰	مدل آی‌کی تک-گرم (وُردتووک)
۰/۰۲۰۵	۰/۰۲۰۲	۰/۰۲۱۵	مدل وی‌کی پدیا (وُردتووک)
۰/۰۲۰۴	۰/۰۱۹۶	۰/۰۱۹۳	مدل آی‌کی دو-گرم + مدل وی‌کی پدیا (وُردتووک)
۰/۰۲۴۰	۰/۰۲۳۲	۰/۰۲۳۱	برت

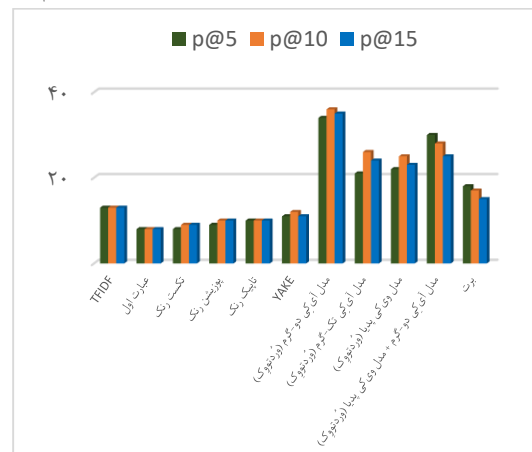
با توجه به آمارهای به‌دست‌آمده، الگوریتم وُردتووک نسبت به سایر روش‌های استخراج کلمات کلیدی کیفیت بهتری داشته است. در میان روش‌های درون‌سازی، روش مدل آی‌کی که با آموزش عبارت‌های دو کلمه‌ای آموزش‌دیده است کیفیت بالاتری داشته است. بهترین حد آستانه تعلق دارد به حد آستانه ۱۰، که از ماکرو-F1 و افت همینگ بهتری برخوردار است.

یکی دیگر از مواردی که در استخراج کلمات کلیدی مفید است، مدت زمانی هست که یک الگوریتم نیاز دارد تا یک کلمه کلیدی را استخراج کند. در جدول ۳ مقایسه زمان برای استخراج کلمات کلیدی انجام شده است. این میانگین زمان بر حسب میلی‌ثانیه بوده و مشخص‌کننده زمان طول‌کشیدن برای یافتن یک کلمه کلیدی است. در ارزیابی‌های انجام شده بهترین زمان ثبت شده برای مدل تکسترنک و در روش‌های درون‌سازی، بهترین زمان ثبت شده متعلق به مدل آی‌کی دو کلمه‌ای است.

جدول (۱): مقایسه‌ی F1 روش‌های پیشنهادی

ماکرو-F1			نام روش
p@15	p@10	p@5	
۱۳	۱۳	۱۳	TF-IDF
۸	۸	۸	عبارت اول
۹	۹	۸	تکسترنک
۱۰	۱۰	۹	پوزیشنرنک
۱۰	۱۰	۱۰	تاپیکرنک
۱۱	۱۲	۱۱	YAKE
۳۵	۳۶	۳۴	آی‌کی دو-گرم (وُردتووک)
۲۴	۲۶	۲۱	آی‌کی تک-گرم (وُردتووک)
۲۳	۲۵	۲۲	مدل وی‌کی پدیا (وُردتووک)
۲۵	۲۸	۳۰	مدل آی‌کی دو-گرم + مدل وی‌کی پدیا (وُردتووک)
۱۵	۱۷	۱۸	برت

برای مقایسه و نمایش بهتر ارزیابی‌های انجام شده، اطلاعات را در قالب نمودار در سه آستانه ۵، ۱۰ و ۱۵ در شکل ۷ ترسیم کرده‌ایم.



شکل (۷): نمودار ارزیابی روش‌های پیشنهادی

در جدول ۲ مقدار افت همینگ را برای تمامی معیارهای ارزیابی محاسبه کرده‌ایم که نشان می‌دهد بهترین مقدار افت همینگ به معیار مدل آی‌کی دو کلمه‌ای با حد آستانه ۱۰ متعلق است.

توجه شده است استخراج کلمات دوکلمه‌ای بوده است زیرا اکثر کلمات کلیدی از دو کلمه و بیشتر تشکیل می‌شوند. از نکات دیگری که در این پژوهش استفاده شده است، امتیاز دهی به واژگان ابتدای جمله است زیرا مفهوم اصلی در آگهی‌ها معمولاً در ابتدا بیان می‌شود. و مورد آخر با توجه به اینکه اکثر آگهی‌ها به زبان عامیانه نوشته می‌شوند، ترکیب روش‌های درون‌سازی با روش‌های آماری کیفیت خوبی را در تشخیص کلمات کلیدی به وجود آورده است.

پژوهش‌های انجام شده در این پژوهش مربوط به آگهی‌های تلگرام در حوزه املاک بوده‌اند. در آینده با گسترش استخراج کلمات کلیدی در حوزه‌های دیگر مانند خودرو، کاریابی، لوازم‌خانگی، لوازم الکترونیکی و غیره می‌توانیم به مدل عمومی برای استخراج کلمات کلیدی دست پیدا کنیم تا مدل بتواند در تمام دسته‌بندی‌ها به‌خوبی عمل کند.

از دیگر اقداماتی که در آینده انجام می‌شود، ساخت مدل درون‌سازی از عبارت‌های سه و چهار کلمه‌ای می‌باشد تا سامانه بتواند عبارت‌های پویا و متناسب را استخراج کند.

هدف نهایی، رسیدن به یک مدل عمومی برای استخراج کلمات کلیدی از هر نوع و دسته می‌باشد. از اقداماتی که می‌توان انجام داد، تشخیص موضوع و دسته‌بندی آگهی و استخراج کلمه کلیدی متناسب با دسته است. به طور مثال آگهی حوزه خودرو توسط مدل خودرو و آگهی بدون دسته با مدل عمومی باید استخراج شود. این اقدام در بالابردن کیفیت کلمات استخراج شده مؤثر است.

سپاس‌گزاری

سپاس خداوند یکتای عزت‌مندی که رحمت و دانش او در سراسر گیتی گسترده شده، آسمان‌ها و زمین همه از آن اوست و علم و دانش حقیقی را بر هر که بخواهد موهبت می‌فرماید. رحمت و لطف او را بی‌نهایت سپاس می‌گویم چرا که فهم و درک مطالب این پژوهش را بر من ارزانی داشت و مرا به این اصل رساند که علم و ایمان دو بال یک پروازند. توفیق تلاش به من داد و هر بار که خطا کردم فرصتی دوباره، تا با امید، تلاشی تازه را آغاز می‌کنم و به خواست او به نتیجه مطلوب نائل آیم. به‌راستی که همه چیز از آن اوست و همه چیز به خواست اوست. همچنین به طور ویژه

جدول (۳): مدت زمان طول کشیدن استخراج یک کلمه

نام روش	مدت زمان (میلی ثانیه)
TF-IDF	۳۳۸
عبارت اول	۱۲۵
تکست‌رنک	۱۲۴
پوزیشن‌رنک	۱۳۳
تاپیک‌رنک	۱۴۳
YAKE	۲۵۷
مدل آی‌کی دو-گرم (وُردتووک)	۱۵۴
مدل آی‌کی تک-گرم (وُردتووک)	۱۵۶
مدل وی‌کی پدیا (وُردتووک)	۲۶۰
مدل آی‌کی دو-گرم + مدل وی‌کی پدیا (وُردتووک)	۲۵۰
برت	۶۰۰

۵- نتیجه‌گیری و پژوهش‌های آتی

به‌طور کلی، استخراج کلمات کلیدی از داده‌های متنی، یکی از مهم‌ترین و پرکاربردترین مسائل در متن‌کاوی و بازیابی اطلاعات بوده است که در این مورد تشخیص کلمات کلیدی نقش تعیین‌کننده‌ای در پردازش اطلاعات دارد.

اغلب روش‌های استخراج کلمات کلیدی، روی عبارت‌های فارسی نتیجه مطلوبی ندارند. یکی دیگر از مشکلات، استفاده از افراد از کلمات عامیانه در آگهی‌ها هست که تشخیص کلمات کلیدی را با چالش همراه می‌کند.

روش پیشنهادی در این پژوهش، روشی ترکیبی از مدل‌های آماری و درون‌سازی متن جهت استخراج کلمات کلیدی است. این روش نیاز به دانش خارجی نداشته و با داده‌های بدون برچسب به‌خوبی کار می‌کند. با توجه به اینکه روش‌های درون‌سازی معنا و مفهوم هر کلمه را تشخیص می‌دهند، امتیاز هر کلمه را با توجه به مکانش محاسبه می‌کنند. روش‌های درون‌سازی توانایی خوبی در تشخیص کلمات عامیانه دارند. از نکاتی که در روش‌های پیشنهادی به آن

- [11] Y. Goldberg, Neural network methods in natural language processing. San Rafael, CA: Morgan & Claypool, 2017.
- [12] Kedia and M. Rasu, Hands-On Python Natural Language Processing: Explore tools and techniques to analyze and process text with a view to building real-world NLP applications. Birmingham, England: Packt Publishing, 2020.
- [13] T. Haarman, B. Zijlema, and M. Wiering, "Unsupervised keyphrase extraction for web pages," Multimodal technol. interact., vol. 3, no. 3, p. 58, 2019.
- [14] K. Bennani-Smires, C. Musat, A. Hossmann, M. Baeriswyl, and M. Jaggi, "Simple Unsupervised Keyphrase Extraction using Sentence Embeddings," in Proceedings of the 22nd Conference on Computational Natural Language Learning, 2018.
- [15] M. Saraswathi and V. Balu, "Preprocessing Techniques for Effective Data Extraction and Computation," IUP Journal of Computer Sciences, Volume 7(3), p. 27, 2013.
- [16] O. Hajipoor, et al., "Determine the Sentiment for Persian Words and Phrases Using Deep Learning," Computer Society of Iran Conference, Volume 24, 2019.
- [17] Harris Hawks' Optimization algorithm," Knowledge-Based Systems,

از جناب آقای دکتر محمدعلی زارع چاهوکی بابت ارشادات ارزنده، ثمربخش و دلسوزانه کمال امتنان را دارم.

References

- [1] A. Onan, S. Korukoğlu, and H. Bulut, "Ensemble of keyword extraction methods and classifiers in text classification," Expert Syst. Appl., vol. 57, pp. 232–247, 2016.
- [2] "Tagging posts properly for users and SEO," Yoast, 11-Apr-2019. [Online]. Available: <https://yoast.com/tagging-posts-properly-for-users-and-seo/>. [Accessed: 23-Jan-2022].
- [3] Williams, WordPress for beginners 2021: A visual step-by-step guide to mastering WordPress. Independently Published, 2020.
- [4] Moz, "What are keywords?," Moz, 28-Mar-2017. [Online]. Available: <https://moz.com/learn/seo/what-are-keywords>. [Accessed: 23-Jan-2022].
- [5] "Telegram revenue and usage statistics (2022)," Business of Apps, 08-Aug-2017. [Online]. Available: <https://www.businessofapps.com/data/telegram-statistics>. [Accessed: 23-Jan-2022].
- [6] H. Tankovska, "Most popular global mobile messenger apps as of January 2023, based on number of monthly active users," 20022. [Online]. Available: <https://www.statista.com/statistics/258749/most-popular-global-mobilemessenger-apps/>. [Accessed: 2-May-2023].
- [7] Y. Sun, H. Qiu, Y. Zheng, Z. Wang, and C. Zhang, "SIFRank: A new baseline for unsupervised keyphrase extraction based on pre-trained language model," IEEE Access, vol. 8, pp. 10896–10906, 2020.
- [8] Won, M., Martins, B. and Raimundo, F. (2019) Automatic extraction of relevant keyphrases for the study of issue competition. In Proceedings of the 20th International Conference on Computational Linguistics and Intelligent Text Processing, Berkeley, La Rochelle, France, April 7-13, 2019.
- [9] C. Florescu and C. Caragea, "PositionRank: An unsupervised approach to keyphrase extraction from scholarly documents," in Proceedings of the 55th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers), 2017.
- [10] R. Ruder, "word-embeddings-1," 1142016. [Online]. Available: <https://www.ruder.io/word-embeddings-1/>.

Extracting keywords from Telegram ads based on the integration of word embedding and statistical features

Mohammad Ali Zare Chahooki*, Mohammad Mahdi Sedghian

Computer Engineering Department, Yazd University, Yazd, Iran

Article Information

Original Research Paper

Received:

2023 November 16

Accepted:

2024 April 8

Keywords:

Tag extraction, search engine optimization (SEO), deep learning, word embedding

Corresponding Author*:

chahooki@yazd.ac.ir

Abstract

In today's society, news and advertisements have a special place in people's communication with each other. By specifying the main words of the ad, you can understand the general concept. One of these uses of keywords is to use them in SEO. Keywords indicate the main purpose of the ad. Since keywords have many uses in recognizing the meaning of text, identifying automatic and improved methods for extracting this category of words has always been of interest to researchers. In terms of SEO, keywords are queries searched by users, for this reason, keyword extraction will play an effective role in increasing site SEO. Ideakav site is a system that collects Telegram messages and ads. According to the requirement of the idekav system in extracting keywords from advertisements published in Telegram, in this research we seek to provide a method that can identify and extract keywords. The quality of extracted keywords plays a significant role in improving SEO and advertising statistics. The data used in this research is from the data collection of Ideakav system. The approach of the current research is an unsupervised integrated method for extracting keywords that uses a combination of statistical methods and word internalization to extract keywords. The word embedding algorithm is able to understand colloquial conversations and the semantic structure of the text, therefore, it is useful in identifying keywords in Telegram ads that are often published in popular form. The obtained results show that the combination of the two-word model word embedding method along with the statistical methods presented in this research has produced a better performance in extracting keywords than other proposed methods.

 : 10.22034/ABMIR.2024.20879.1041

E-ISSN: [2821-2037](https://doi.org/10.22034/ABMIR.2024.20879.1041) /© 2023. Published by Yazd University This is an open access article under the CC BY 4.0 License (<https://creativecommons.org/licenses/by/4.0/>).

