

## تشخیص ارائه‌دهندگان خدمات در پیام‌های فارسی تلگرام مبتنی بر روش‌های انتخاب ویژگی

فرزانه ابراهیمی<sup>۱</sup>، محمد علی زارع چاهوکی<sup>۲\*</sup>، علی هاشمی<sup>۳</sup>

<sup>۱</sup> دانشجوی ارشد مهندسی کامپیوتر، دانشگاه یزد، یزد، ایران

<sup>۲</sup> استادیار دانشکده مهندسی کامپیوتر، دانشگاه یزد، یزد، ایران

<sup>۳</sup> دکتری مهندسی کامپیوتر، دانشگاه یزد، یزد، ایران

### چکیده

پیام‌رسان تلگرام بستری مناسب برای کاربرانی است که به دنبال خرید محصول یا دریافت خدمات به صورت آنلاین هستند. در این پیام‌رسان، برای درخواست محصول، امکان دسترسی مستقیم به ارائه‌دهندگان کالا و خدمات وجود ندارد. لذا باید ابتدا در گروه‌های تلگرامی مرتبط عضو شد و درخواست خود را در تک تک گروه‌ها ثبت کرد و منتظر پاسخ ماند. این امر ضمن زمان‌بر بودن، با مشکلاتی همراه است. هدف از این پژوهش تشخیص کاربران ارائه‌دهنده خدمات به منظور ارتباط مستقیم و مؤثر با مشتریان است. بدین منظور از ایده دسته‌بندی پیام‌های فارسی منتشر شده در تلگرام استفاده شد. یکی از مشکلات دسته‌بندی این پیام‌ها، ابعاد بزرگ فضای ویژگی است که سبب کاهش دقت و افزایش زمان دسته‌بندی می‌شود. برای حل این مسئله از روش‌های انتخاب ویژگی استفاده شد. روش پیشنهادی این پژوهش، بر مبنای ترکیب روش‌های انتخاب ویژگی مبتنی بر فیلتر محلی و سراسری است. نوآوری این پژوهش در استفاده از روش‌های ترکیبی انتخاب ویژگی جهت دسته‌بندی خودکار پیام‌های فارسی تلگرام، به منظور شناسایی کاربران ارائه‌دهنده خدمات است. روش پیشنهادی، ضمن کاهش تعداد ویژگی‌ها و انتخاب ویژگی‌های مرتبط، سبب بهبود عملکرد دسته‌بندی و تشخیص کاربران ارائه‌دهنده خدمات می‌شود.

### مقاله پژوهشی

تاریخ دریافت:

۱۴۰۱/۰۵/۰۹

تاریخ پذیرش:

۱۴۰۱/۱۱/۱۰

کلیدواژه‌ها:

دسته‌بندی، کاربران ارائه‌دهنده خدمات، انتخاب ویژگی، کاهش ویژگی، یادگیری ماشین

نویسنده مسئول:

chahooki@yazd.ac.ir



10.22034/ABMIR.2023.18780.1012

E-ISSN: [2821-2037](https://doi.org/10.22034/ABMIR.2023.18780.1012) /© 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/>).



## ۱- مقدمه

با شناسایی کاربران ارائه‌دهنده کالا و خدمات، این امکان وجود دارد تا بتوانند از فعالان هر حوزه، مشاوره بگیرند. جهت تشخیص کاربران ارائه‌دهنده خدمات از روش دسته‌بندی<sup>۳</sup> پیام‌های ارسالی آن‌ها در گروه‌های مختلف تلگرامی استفاده شد. عموماً دسته‌بندی خودکار متن بر مبنای ویژگی‌های استخراجی انجام می‌گیرد. یکی از مشکلات اساسی در دسته‌بندی، تعداد زیاد ویژگی‌های استخراجی است، چرا که اغلب این ویژگی‌ها نامربوط و زائد هستند و بر کارایی دسته‌بندی تأثیر منفی می‌گذارند [۱]. انتخاب ویژگی<sup>۴</sup>، راه‌حلی مناسب برای کاهش ابعاد بزرگ فضای ویژگی و افزایش کارایی عملکرد دسته‌بندی است [۲]. هدف از انتخاب ویژگی، انتخاب بهینه‌ترین زیرمجموعه ویژگی از کل فضای ویژگی‌های اصلی متن مورد نظر است [۳]، به طوری که بتوان ضمن کاهش تعداد ویژگی‌ها، به دقت دسته‌بندی قابل قبولی نیز دست یافت. با انتخاب ویژگی‌های مناسب، زمان آموزش و محاسبات در یادگیری ماشین کاهش پیدا کرده و دقت پیش‌گویی افزایش می‌یابد [۴].

این پژوهش روشی را جهت شناسایی کاربران ارائه‌دهنده خدمات در پیام‌رسان تلگرام ارائه می‌دهد. مبنای کار این روش، دسته‌بندی پیام‌های ارسالی کاربران در گروه‌های مختلف تلگرامی است. در گام اول پیام‌های ارسالی شده توسط هر کاربر در گروه‌های مختلف تلگرامی استخراج می‌شود. در گام دوم با استفاده از پرکاربردترین روش‌های انتخاب ویژگی مبتنی بر فیلتر که به دودسته محلی و سراسری تقسیم می‌شوند، ویژگی‌های غیرضروری حذف می‌شوند. در این پژوهش ویژگی‌های استخراجی، همان کلمات موجود در پیام‌ها هستند. سپس هر یک از روش‌های محلی و سراسری باهم ترکیب و ویژگی‌های بهتر انتخاب و با استفاده از این ویژگی‌ها دسته‌بندی انجام می‌شود. نتایج این آزمایش‌ها نشان داد که ترکیب روش‌های ترکیبی انتخاب ویژگی، در اکثر موارد می‌تواند سبب افزایش دقت نسبت به حالت بدون انتخاب ویژگی، در دسته‌بندی و شناسایی کاربران ارائه‌دهنده خدمات شود.

با گسترش شبکه‌های اجتماعی<sup>۱</sup>، تحولی بنیادین در نحوه ارتباط صاحبان کسب‌وکار با مشتریانشان به وجود آمد. در حال حاضر، هر فردی بدون توجه به اندازه کسب‌وکار و میزان بودجه‌ای که دارد، می‌تواند از فضای مجازی و شبکه‌های اجتماعی برای ارتباط با مشتریان، و معرفی محصولات و خدماتش به طیف وسیعی از مردم استفاده کند. تلگرام<sup>۲</sup>، یکی از پیام‌رسان‌هایی است که بین ایرانیان محبوبیت زیادی کسب کرده است. صاحبان کسب‌وکارهای مختلف با ایجاد گروه و کانال‌های تلگرامی، روزانه میلیون‌ها پیام فارسی از جنس تبلیغ و معرفی محصول و خدمات منتشر می‌کنند که بستر مناسبی جهت درآمدزایی برای آن‌ها ایجاد کرده است.

شناسایی ارائه‌دهندگان کالا و خدمات در تلگرام سبب برقراری ارتباط مستقیم بین صاحبان کسب‌وکار و مشتریانشان می‌شود. به‌طور مثال کاربری که درخواست خرید یک قطعه خودرو را دارد می‌تواند فقط با نوشتن درخواست خود و یک‌بار ارسال آن، این تقاضا را به‌طور مستقیم برای مجموعه‌ای از ارائه‌دهندگان کالا که در حوزه<sup>۵</sup> فروش قطعات خودرو فعالیت می‌کنند، ارسال کند. لذا دیگر نیاز به عضویت در تک‌تک کانال‌ها یا گروه‌های مرتبط نیست. این امر برای ارائه‌دهندگان کالا و خدمات در هر حوزه کاری دسترسی سریع فراهم می‌کند. همچنین زمان پاسخگویی به درخواست مشتریان را کاهش و احتمال پاسخ مثبت را افزایش می‌دهد. از طرفی با شناسایی کاربران ارائه‌دهنده خدمات، می‌توان کمپین‌هایی مخصوص هر حوزه کاری راه‌اندازی کرد که صاحبان کسب‌وکار در آن حوزه خاص، نظرات و تجربیات خود را بیان کنند. این کمپین‌ها کمک شایانی به افراد تازه‌وارد به هر حوزه می‌کند، زیرا می‌توانند در مدت‌زمان کمی از تجربیات پیش‌کسوتان باخبر شوند و مسیر هموارتری را برای ادامه کسب‌وکار خود ترسیم کنند. بسیاری از افراد برای راه‌اندازی یک فعالیت، خواهان اطلاعات درباره چالش‌های ورود به آن حوزه کاری هستند تا با آمادگی و بینش قوی‌تری ورود کرده و احتمال شکست خود را پایین بیاورند.

<sup>3</sup> classification

<sup>4</sup> Feature selection

<sup>1</sup> Social network

<sup>2</sup> Telegram

طبقه‌بندی از پیش تعریف شده به جای اندازه‌گیری مستقل برای ارزیابی زیرمجموعه ویژگی‌ها استفاده می‌کنند. روش‌های یادگیر مینا در مقایسه با روش‌های فیلتر، عملکرد بهتری دارند، اما زمانی که تعداد ویژگی‌ها بسیار زیاد شود، از نظر محاسباتی گران‌تر می‌شوند. دو روش انتخاب متوالی روبه‌جلو (SFS) و حذف متوالی به عقب (SBS)، دو روش متداول حریصانه هستند که اغلب برای انتخاب ویژگی استفاده می‌شوند. [۵].

در روش‌های یادگیر حاصل، انتخاب ویژگی و یادگیری مدل به صورت هم‌زمان انجام می‌شود. روش‌های یادگیر حاصل زمان محاسبه برای دسته‌بندی مجدد زیرمجموعه‌های مختلف در روش یادگیر مینا را کاهش می‌دهند [۶]. در مجموع روش‌های فیلتر محاسبات ساده‌تر و زمان پردازش کوتاه‌تری دارند. لذا جهت کار با مجموعه داده‌های با ابعاد بالا مناسب هستند. اما روش‌های یادگیر مینا با وجود کارایی قابل قبول، از نظر محاسباتی پیچیده‌تر و پرهزینه‌تر هستند [۵]. از این رو روش‌های فیلتر در مقایسه با روش‌های یادگیر مینا سریع‌تر عمل می‌کنند. ضعف روش‌های یادگیر حاصل در نادیده گرفتن ارتباط بین ویژگی‌ها است که سبب مشکلاتی در انتخاب ویژگی‌های نهایی می‌شود [۷]. آقای یوسال ترکیب یک روش انتخاب ویژگی سراسری مبتنی بر فیلتر و یک روش انتخاب ویژگی محلی را برای بهبود انتخاب ویژگی، استفاده کرده است که در واقع به دنبال ایجاد مجموعه‌ای از ویژگی‌هایی هستند که تقریباً به طور یکسان تمام کلاس‌ها را در برگیرند. نتایج تجزیه و تحلیل تجربی نشان می‌دهد که این روش ترکیبی، عملکرد دسته‌بندی بهتری نسبت به عملکرد تکی روش‌های انتخاب ویژگی سراسری دارد؛ اما این روش در برخورد با مجموعه داده‌های نامتعادل دارای تعداد زیادی کلاس، مشکل دارد [۷].

پژوهش آقای اگنی‌هتری طرح انتخاب ویژگی متغیر سراسری (VGFSS) را که ترکیبی از یک روش سراسری و روش نسبت شانس است، جهت طبقه‌بندی خودکار اسناد متنی ارائه می‌دهد. این

در گام سوم خروجی ویژگی‌های روش ترکیبی گام دوم، به عنوان ورودی برای روش انتخاب ویژگی یادگیر مینا استفاده می‌شود. با این روش بار دیگر ویژگی‌های کم‌اهمیت حذف می‌شود. نتایج آزمایش‌ها نشان می‌دهد که استفاده از روش ترکیبی فیلتر و یادگیر مینا عملکرد بهتری نسبت به روش‌های فیلتر دارد. در ادامه در بخش دوم به بررسی مطالعات پیشین پرداخته می‌شود. روش پیشنهادی در بخش سوم به تفصیل بیان می‌شود. نتایج استفاده از روش پیشنهادی در بخش چهارم ارائه و در بخش پنجم نیز به بیان خلاصه کارهای انجام شده و نتیجه‌گیری پرداخته می‌شود.

## ۲- مطالعات پیشین

تکنیک‌های انتخاب ویژگی به طور کلی به سه دسته فیلتر<sup>۱</sup>، یادگیر مینا<sup>۲</sup> و یادگیر حاصل<sup>۳</sup> طبقه‌بندی می‌شوند. روش‌های فیلتر، مستقل از الگوریتم یادگیری، ویژگی‌ها را بر اساس یک مرحله پیش‌پردازش انتخاب می‌کنند. در این روش ابتدا به تک‌تک ویژگی‌ها امتیازی داده می‌شود و بر اساس این امتیاز مرتب می‌شوند و ویژگی‌هایی با بالاترین امتیاز انتخاب می‌شوند. در نهایت این مجموعه به عنوان ورودی به الگوریتم دسته‌بندی داده می‌شود. روش‌های فیلتر شامل دودسته محلی<sup>۴</sup> و سراسری<sup>۵</sup> می‌شوند. روش‌های انتخاب ویژگی سراسری معروف عبارت‌اند از: فراوانی سند<sup>۶</sup> (DF)، بهره<sup>۷</sup> اطلاعات (IG)، شاخص‌های جینی<sup>۸</sup> (GI) و انتخاب‌گر ویژگی متمایز<sup>۹</sup> (DFS). از روش‌های انتخاب ویژگی محلی پرکاربرد می‌توان نسبت شانس<sup>۱۰</sup> (OR) و ضریب همبستگی<sup>۱۱</sup> (CC) را نام برد. روش فیلتر از الگوریتم یادگیری خارجی برای ارزیابی عملکرد ویژگی‌های انتخاب شده استفاده می‌کند [۱].

در روش‌های یادگیر مینا، زیرمجموعه‌ای از ویژگی‌ها بر اساس یک الگوریتم یادگیری به گونه‌ای انتخاب می‌شوند که سبب بیشینه شدن کارایی آن روش یادگیری شود. روش‌های یادگیر مینا تقریباً شبیه به روش‌های فیلتر عمل می‌کنند با این تفاوت که از یک الگوریتم

<sup>7</sup> Information Gain (IG)

<sup>8</sup> Gini Index (GI)

<sup>9</sup> Distinguishing Feature Selector (DFS)

<sup>10</sup> Odd Ratio (OR)

<sup>11</sup> Correlation Coefficient (CC)

<sup>1</sup> Filter

<sup>2</sup> Wrapper

<sup>3</sup> Embed

<sup>4</sup> Local

<sup>5</sup> Global

<sup>6</sup> Document Frequency (DF)

ویژگی‌های زیاد استفاده کردند. این روش، ترکیبی از رویکرد مبتنی بر فرکانس و رویکرد مبتنی بر خوشه است [۱۲].

آقای اسیف و همکارانش یک روش انتخاب ویژگی جدید مبتنی بر بهینه‌سازی ازدحام ذرات تطبیقی با وزن خود اینرسی (SIW-APSO) برای افزایش عملکرد سیستم‌های طبقه‌بندی متن پیشنهاد می‌کنند. این روش به دلیل صلاحیت جستجوی بالا و توانایی یافتن زیرمجموعه ویژگی‌ها به‌طور کارآمد، دارای پدیده همگرایی سریع است. تجزیه و تحلیل تجربی این مقاله نشان می‌دهد که روش پیشنهادی نسبت به الگوریتم‌های پیشرفته موجود در مجموعه داده‌های روترز-۲۱۵۷۸ عملکرد بهتری دارد [۱۳].

آقای افروسیدینیز و همکارانش از ۱۲ روش تکی و ۶ روش گروهی برای انتخاب ویژگی جهت طبقه‌بندی داده‌های محیطی، به‌ویژه در حوزه مدل‌سازی توزیع گونه‌ها استفاده می‌کنند. این روش‌ها هر ۳ دسته روش‌های انتخاب ویژگی فیلتر، یادگیرمبنا، و یادگیر حاصل را شامل می‌شوند. آزمایش آن‌ها بر روی ۸ مجموعه داده‌های محیطی نشان می‌دهد که توضیحات افزودنی (SHAP) و اهمیت جایگشت مؤثرترین روش‌های فردی هستند که هر دو در دسته روش‌های یادگیرمبنا قرار می‌گیرد. به‌طور کلی، روش‌های فیلتر عملکرد ضعیفی دارند و روش‌های یادگیر حاصل در این بین قرار می‌گیرند. از ۶ روش گروهی در نظر گرفته‌شده، Reciprocal Ranking، بهترین عملکرد را دارد [۱۴].

## ۲-۱ روش‌های انتخاب ویژگی مورد استفاده در این

### پژوهش

روش‌های انتخاب ویژگی محلی مورد استفاده در این پژوهش شامل نسبت شانس و ضریب همبستگی است که در ادامه به توضیح عملکرد این روش‌ها پرداخته می‌شود. معیار نسبت شانس، عضویت در کلاس خاصی را با صورت کسر و عدم عضویت را با مخرج کسر اندازه‌گیری می‌کند. امتیاز عضویت و عدم عضویت، با تقسیم به یکدیگر نرمال می‌شوند [۱۰]. طبق رابطه (۱) به‌صورت مجزا امتیاز نسبت شانس هر اصطلاح برای تمام کلاس‌ها محاسبه و بالاترین امتیاز را برای اختصاص یک برچسب کلاس به اصطلاح انتخاب می‌شود. برای تعیین برچسب، مقدار امتیاز مطلق در نظر گرفته می‌شود.

روش تعداد متغیری از ویژگی‌ها را از هر کلاس بر اساس توزیع عبارات در کلاس‌ها انتخاب می‌کند که این امر تضمین می‌کند حداقل تعداد اصطلاحات از هر کلاس انتخاب می‌شود [۸].

آقای ملو و همکارانش با مقایسه و ارزیابی تجربی دو روش انتخاب ویژگی محلی و سراسری بر روی مجموعه داده‌های طبقه‌بندی چند برچسبی مسطح و سلسله مراتبی رایج (که معمولاً در محک زدن عملکرد طبقه‌بندی کننده‌های چند برچسبی استفاده می‌شوند) نشان می‌دهند که انتخاب ویژگی محلی عملکرد بهتری نسبت به انتخاب ویژگی سراسری در این نوع داده‌ها دارند. [۹]

در برخی پژوهش‌ها، روش‌های محلی برای تعداد ویژگی‌های کم و روش‌های سراسری برای تعداد ویژگی‌های زیاد نتایج بهتری داشته‌اند [۵]. آقای کرس و همکارانش از روش‌های فیلتر به‌صورت مستقل و ترکیبی استفاده کردند [۱۰]. اندازه‌گیری دقیق روابط ویژگی‌های نامزد، ویژگی‌ها و دسته‌بندی‌های انتخاب‌شده در فرآیند انتخاب، به‌ویژه برای داده‌های با ابعاد بالا و اندازه نمونه کوچک، یک کار چالش‌برانگیز است. برای این منظور در پژوهش آقای وی و همکارانش معیار جدیدی به نام اهمیت ویژگی پویا (DFI) و همچنین الگوریتم انتخاب ویژگی متناظر آن بانام انتخاب ویژگی مبتنی بر اهمیت ویژگی پویا (DGIFS) پیشنهاد شده است. این الگوریتم ترکیبی از روش انتخاب ویژگی مبتنی بر اهمیت ویژگی پویا با روش‌های فیلتر است. هدف آن‌ها رسیدن به دقت دسته‌بندی بالاتر با تعداد کمتری از ویژگی‌ها است. [۱۱].

تکنیک کیسه کلمات اغلب برای ارائه سند در طبقه‌بندی متن استفاده می‌شود. با این حال، برای مجموعه بزرگی از اسناد که ابعاد بردار کیسه کلمات بسیار بزرگ است، دسته‌بندی متن به دلیل داده‌های پراکنده، برآزش بیش‌ازحد، و ویژگی‌های نامربوط، به چالشی جدی تبدیل می‌شود. روش انتخاب ویژگی فیلتر با حذف ویژگی‌های نامربوط از بردار کیسه کلمات، تعداد ویژگی‌ها را کاهش می‌دهد. آقای نام و همکارانش نقاط ضعف و نقاط قوت دو رویکرد انتخاب ویژگی فیلتر که رویکرد مبتنی بر فرکانس و رویکرد مبتنی بر خوشه است را تحلیل کردند، سپس از روش‌های انتخاب ویژگی فیلتر ترکیبی برای بهبود طبقه‌بندی متن با تعداد



$$IG(t) - \sum_{i=1}^M P(C_i) \log P(C_i) + P(t) \sum_{i=1}^M P(C_i|t) \log P(C_i|t) + P(\bar{t}) \sum_{i=1}^M P(C_i|\bar{t}) \log P(C_i|\bar{t}) \quad (3)$$

$M$  تعداد کلاس‌ها،  $P(C_i)$  احتمال کلاس  $C_i$ ،  $P(t)$  احتمال وجود اصطلاح  $t$  و  $P(\bar{t})$  احتمال عدم وجود اصطلاح  $t$  را نشان می‌دهد.  $P(C_i|t)$  احتمال کلاس  $C_i$  به شرط وجود اصطلاح  $t$  است.  $P(C_i|\bar{t})$  احتمال کلاس  $C_i$  به شرط عدم وجود اصطلاح  $t$  است. طریقه محاسبه روش شاخص جینی در رابطه (۴) بیان شده است. این روش به شکل بهبودیافته در دسته‌بند درخت تصمیم استفاده می‌شود [۱۰].

$$GI(t) = \sum_{i=1}^M P(t|C_i)^2 P(C_i|t)^2 \quad (4)$$

$M$  تعداد کلاس‌ها،  $P(t|C_i)$  احتمال وجود اصطلاح  $t$  به شرط کلاس  $C_i$  است.  $P(C_i|t)$  احتمال کلاس  $C_i$  به شرط وجود اصطلاح  $t$  است.

روش انتخاب ویژگی مبتنی بر فیلتر زمانی مطلوب است که امتیاز بالایی را به اصطلاحات متمایز و امتیاز پایینی را به موارد نامربوط اختصاص دهد. یکی از روش‌های انتخاب ویژگی مناسب برای دسته‌بندی متن، انتخاب‌گر ویژگی متمایز است [۱۰]. این روش اصطلاحات متمایز را انتخاب می‌کند [۸]. طریقه محاسبه این روش در رابطه (۵) قابل مشاهده است.

$$DFS(t) = \sum_{i=1}^M \frac{P(C_i|t)}{P(\bar{t}|C_i) + P(t|\bar{C}_i) + 1} \quad (5)$$

$P(t|C_i)$  احتمال وجود کلاس  $C_i$  به شرط وجود اصطلاح  $t$  است.  $P(\bar{t}|C_i)$  احتمال عدم وجود اصطلاح  $t$  به شرط کلاس  $C_i$  است.  $P(t|\bar{C}_i)$  احتمال شرطی اصطلاح  $t$  برای تمام کلاس‌ها به جز کلاس  $C_i$  است.

فراوانی سند مربوط به یک اصطلاح، تعداد اسناد آموزشی است که اصطلاح موردنظر در آن اسناد وجود دارد. از این روش به دلیل

$$OR(t|C_i) = \log \frac{P(t|C_i)[1 - P(t|\bar{C}_i)]}{[1 - P(t|C_i)]P(t|\bar{C}_i)} \quad (1)$$

$P(t|C_i)$  احتمال وجود اصطلاح  $t$  به شرط کلاس  $C_i$  است.  $P(t|\bar{C}_i)$  احتمال شرطی اصطلاح  $t$  برای تمام کلاس‌ها به جز کلاس  $C_i$  است.

روش ضریب همبستگی، شایستگی یک زیرمجموعه از اصطلاحات را با در نظر گرفتن قدرت پیش‌بینی تکی هر اصطلاح و درجه افزونگی بین آن‌ها، ارزیابی می‌کند. زیرمجموعه مناسب از ویژگی‌ها، همبستگی بالایی در کلاس و همبستگی پایینی بین کلاس‌های مختلف دارد. در رابطه (۲) ضریب همبستگی  $CC$  برای اصطلاح  $t$  در یک کلاس  $C_i$  برای مجموعه‌ای با  $N$  نمونه تعریف شده است [۳].

$$CC(t, C_i) = \frac{1}{\sqrt{N}} \frac{P(t, C_i)P(\bar{t}, \bar{C}_i) - P(t, \bar{C}_i)P(\bar{t}, C_i)}{\sqrt{P(t)P(\bar{t})P(C_i)P(\bar{C}_i)}} \quad (2)$$

$N$  تعداد کل اسناد را نشان می‌دهد.  $P(t|C_i)$  احتمال وجود اصطلاح  $t$  به شرط کلاس  $C_i$  است.  $P(\bar{t}|\bar{C}_i)$  احتمال عدم وجود اصطلاح  $t$  برای تمام کلاس‌ها به جز کلاس  $C_i$  است.  $P(\bar{t}|C_i)$  احتمال عدم وجود اصطلاح  $t$  به شرط کلاس  $C_i$  است.  $P(t)$  احتمال وجود اصطلاح  $t$  و  $P(\bar{t})$  احتمال عدم وجود اصطلاح  $t$  را نشان می‌دهد.  $P(C_i)$  احتمال وجود کلاس  $C_i$  و  $P(\bar{C}_i)$  احتمال عدم وجود کلاس  $C_i$  را نشان می‌دهد.

روش‌های انتخاب ویژگی سراسری این پژوهش شامل فراوانی سند، بهره اطلاعات، شاخص‌های جینی و انتخاب‌گر ویژگی متمایز است که در ادامه به معرفی این روش‌ها پرداخته می‌شود. روش بهره اطلاعات، طبق رابطه (۳) سهم وجود یا عدم وجود اصطلاح برای دسته‌بندی صحیح اسناد متنی را محاسبه می‌کند. روش بهره اطلاعات یک استراتژی نظارت‌شده است. یعنی از برچسب اسناد در انتخاب اصطلاحات برای حفظ آن‌ها استفاده می‌کند [۵].

<sup>1</sup> Term

ارائه‌دهنده خدمات شناخته می‌شود. در غیر این صورت، سند برچسب کلاس منفی می‌گیرد. در طی این پژوهش پیام‌های مربوط به ۳۰ هزار کاربر تلگرام برچسب‌زنی شد. از این تعداد ۹,۵۰۶ کاربر دارای برچسب مثبت، ۱۷,۸۲۸ کاربر دارای برچسب منفی و ۲,۶۶۶ کاربر دارای برچسب نامشخص بودند. ابتدا داده‌های با برچسب‌های نامشخص کنار گذاشته شدند. لذا ۲۷,۳۳۴ کاربر با برچسب مثبت و منفی داریم. داده‌های به‌دست‌آمده را به دو قسمت تقسیم می‌کنیم و مراحل آماده‌سازی داده‌ها را بر روی این تقسیم‌بندی انجام می‌دهیم. ۸۰ درصد اول داده‌ها را به‌عنوان مجموعه آموزشی و ۲۰ درصد بعدی را به‌عنوان مجموعه آزمون در نظر می‌گیریم. در داده‌های آموزش و آزمون، نمونه‌های با برچسب منفی بیشتر از نمونه‌های با برچسب مثبت است. در این پژوهش جهت دسته‌بندی از چهار روش ماشین بردار پشتیبان<sup>۳</sup> (SVM)، بیزین ساده<sup>۴</sup> (NB)، شبکه عصبی<sup>۵</sup> (MLP) و درخت تصمیم<sup>۶</sup> (DT) استفاده کردیم. همچنین از معیار میانگین میکرو<sup>۷</sup> (Micro-F1)، برای ارزیابی نتایج استفاده شد.

### ۲-۳ انتخاب ویژگی

دسته‌بندی به‌صورت دستی، بسیار زمان‌بر است، لذا برای تشخیص کاربران ارائه‌دهنده خدمات، از روش‌های دسته‌بندی خودکار اسناد استفاده شد. روش‌های دسته‌بندی بر مبنای ویژگی‌های استخراجی از متن، عمل می‌کنند. در این پژوهش ویژگی‌های استخراجی، همان کلمات موجود در پیام‌ها هستند. لذا در گام بعد بر روی همه اسناد جداسازی کلمات<sup>۸</sup> انجام شد و در نهایت برای هر کلمه، فرکانس تکرار آن در کل اسناد محاسبه شد. بعد از حذف کلمات توقف<sup>۹</sup>، برای رسیدن به ویژگی‌های بهتر، کلماتی که فرکانس بسیار بالا یا بسیار پایین داشتند حذف شدند زیرا کلماتی که نادر و یا بسیار پرتکرار هستند معمولاً ارزش بالایی در دسته‌بندی ندارند. سپس کلمات باقیمانده، به‌عنوان ویژگی و واژگان موجود در کیسه کلمات<sup>۱۰</sup> در نظر گرفته شدند.

هزینه کم محاسبات، در دسته‌بندی متن بیشتر استفاده می‌شود [۸]. رویکرد روش فراوانی سند این است که در دسته‌بندی متن، اصطلاحات نادر غیر آموزنده هستند و باید در هنگام انتخاب ویژگی حذف شوند. نحوه محاسبه این روش در رابطه (۶) نشان داده شده است.  $N$  تعداد کل اسناد،  $P(t)$  احتمال رخداد اصطلاح  $t$  را نشان می‌دهد.

$$DF(t) = N \cdot P(t) \quad (6)$$

### ۳- روش تحقیق

برای شناسایی کاربران ارائه‌دهنده خدمات، از دسته‌بندی پیام‌های ارسالی آن‌ها در گروه‌های تلگرامی استفاده شد. روش پیشنهادی در چندین مرحله انجام شد که به‌طور خلاصه شامل جمع‌آوری اطلاعات، پیش‌پردازش، انتخاب ویژگی، اجرا روش‌های یادگیری ماشین و ارزیابی نتایج می‌شود. در ادامه جزئیات هر کدام از مراحل به‌تفصیل بیان می‌شود.

### ۱-۳ جمع‌آوری داده‌ها

در گام نخست یک میلیون پیام منتشرشده در گروه‌های مختلف تلگرامی به همراه شناسه ارسال‌کننده آن‌ها، از سامانه ایده‌کاوا<sup>۱</sup> که یک موتور جست‌وجوی تلگرام است، جمع‌آوری شد. سپس این داده‌ها برای نمایه‌سازی به الاستیک سرچ<sup>۲</sup> منتقل شدند. الاستیک سرچ یک موتور جست‌وجوگر متن است. در گام بعد با جست‌وجو در نمایه ساخته‌شده مبتنی بر شناسه تعلق‌گرفته به هر کاربر، مجموع پیام‌های ارسالی هر کاربر، جمع‌آوری و مجدداً، نمایه جدیدی ایجاد شد. در نمایه جدید هر کاربر با مجموع پیام‌های ارسالی خود، به‌عنوان یک سند لحاظ می‌شود. تعداد اسناد ایجادشده در این مرحله به ۳۰ هزار می‌رسد.

در گام بعد، دسته‌بندی در قالب برچسب‌زنی دستی بر روی اسناد انجام شد. فرایند برچسب‌زنی به این صورت است که، اگر مجموع پیام‌های ارسالی یک کاربر، از نوع تبلیغات و معرفی خدمات باشد به آن سند برچسب کلاس مثبت تعلق می‌گیرد و آن کاربر به‌عنوان

<sup>6</sup> Decision Tree (DT)

<sup>7</sup> Micro-averaging

<sup>8</sup> Tokenization

<sup>9</sup> Stop words

<sup>10</sup> Bag Of Words (BOW)

<sup>1</sup> Idekav.com

<sup>2</sup> Elastic search

<sup>3</sup> Support Vector Machine (SVM)

<sup>4</sup> Naive Bayes (NB)

<sup>5</sup> Multi-Layer Perceptron (MLP)

ویژگی‌های انتخاب‌شده سراسری تعریف می‌کنیم که شامل تعداد  $G$  ویژگی است.

### ترکیب مجموعه ویژگی‌ها

از دو مرحله قبل، دو مجموعه ویژگی به دست آمد که هر کدام مقدار تعیین‌شده‌ای دارند. در این مرحله با توجه به شکل (۱) مجموعه ویژگی خروجی مرتب‌شده از هر یک از روش‌های محلی را با مجموعه ویژگی خروجی مرتب‌شده از هر یک از روش‌های سراسری ترکیب می‌کنیم. از این دو مجموعه، ویژگی‌هایی انتخاب می‌شوند که امتیازات بالاتری داشته باشند. در نهایت یک مجموعه ویژگی نهایی  $F_{Final}$  به دست می‌آید.

$F_{Final} \subseteq \{F_{Local} \cup F_{Global}\}$  مجموعه ویژگی نهایی  $F_{Final}$  مقدار تعیین‌شده‌ای دارد و در آن بهترین ویژگی‌ها با مقدار بالاتر انتخاب‌شده‌اند. بعد از انتخاب ویژگی در هر روش، نوبت به ساخت بردار ویژگی می‌رسد. بردار ویژگی مورد استفاده برای روش‌های یادگیری به صورت ورودی است و هر ویژگی حضور یا عدم حضور یک کلمه را نشان می‌دهد. مجموع بردارهای ویژگی به عنوان ماتریس در نظر گرفته می‌شود. لذا برای هر ترکیب، ماتریس ویژگی مربوط به آن روش ساخته می‌شود. در گام بعد، عملکرد روش‌های انتخاب ویژگی در دسته‌بندی، توسط روش‌های یادگیری ماشینی ارزیابی می‌شود.

### انتخاب ویژگی با ترکیب روش‌های فیلتر و یادگیرمبنا

در گام بعد به منظور افزایش دقت، از ترکیب روش‌های فیلتر و یادگیرمبنا استفاده می‌شود. بدین منظور قبل از این که ویژگی‌های اصلی به روش یادگیرمبنا داده شود، باید از فیلتر روش پیشنهادی اول عبور کنند. در گام بعد ویژگی‌های مناسب انتخاب‌شده توسط روش‌های فیلتر، به عنوان ویژگی‌های ورودی، به روش یادگیرمبنا داده می‌شوند. با ترکیب روش‌های انتخاب ویژگی فیلتر و روش یادگیرمبنا ( $SFS$ ) و بررسی نتایج، مشاهده می‌شود که استفاده از این روش سبب کاهش چشمگیر ویژگی‌ها و افزایش دقت در دسته‌بندی می‌شود.

بعد از مراحل پیش‌پردازش تعداد کل ویژگی‌های استخراج‌شده از متن اصلی برابر  $7435$  عدد است. مجموعه ویژگی‌های به دست آمده از مرحله پیش‌پردازش را با  $F = \{F_1, F_2, \dots, F_N\}$  نشان داده و آن را به عنوان مجموعه‌ای از تمام ویژگی‌های اولیه تعریف می‌کنیم که در این مجموعه  $N = 7435$  است.

دسته‌بندی اسناد با استفاده از این تعداد ویژگی، زمان آموزش مدل را در یادگیری ماشینی، بسیار طولانی می‌کند. لذا برای کاهش این زمان و افزایش دقت دسته‌بندی، نیاز به اعمال روش‌های انتخاب ویژگی بر روی ویژگی‌های اولیه است تا ضمن حذف ویژگی‌های زائد، زمان دسته‌بندی اسناد را کاهش و دقت دسته‌بندی را افزایش دهد. برای اعمال روش‌های انتخاب ویژگی نیاز به ماتریس کلاس-اصطلاح است. ماتریس کلاس-اصطلاح ماتریسی است که سطرهای آن نماینده دو کلاس مثبت و منفی و ستون‌های آن، ویژگی‌های اولیه هستند. درایه  $a_{ij}$  از این ماتریس، تعداد اسناد کلاس  $j$  ام که ویژگی  $i$  ام در آن وجود داشته را نشان می‌دهد. با استفاده از این ماتریس می‌توان هر یک از روش‌های انتخاب ویژگی محلی و سراسری را اجرا کرد.

### انتخاب ویژگی با روش‌های محلی و سراسری

در این مرحله با استفاده از روش‌های انتخاب ویژگی محلی  $CC$  و  $OR$  و روش‌های انتخاب ویژگی سراسری  $GI$ ،  $IG$ ،  $DF$  و  $DFS$  به ویژگی‌های هر کلاس  $C_i$  موجود در مجموعه  $F$  یک امتیاز داده می‌شود. این امتیازها نشان‌دهنده تفاوت میان ویژگی‌ها در یک مجموعه داده است. در مرحله بعد همه ویژگی‌ها با توجه به امتیازی که در مرحله انتخاب ویژگی کسب کرده‌اند به ترتیب نزولی مرتب می‌شود. سپس در هر یک از روش‌های انتخاب ویژگی محلی،  $L$  ویژگی با بالاترین امتیاز به عنوان ویژگی‌های محلی و در هر یک از روش‌های انتخاب ویژگی سراسری،  $G$  ویژگی با بالاترین امتیاز به عنوان ویژگی‌های سراسری نهایی انتخاب می‌شوند. مقدار  $L$  و  $G$ ، یک تعداد تعیین‌شده است که به صورت تجربی حاصل شده است. در هر روش محلی،  $F_{Local} \subseteq F$  به عنوان مجموعه‌ای از ویژگی‌های انتخاب‌شده محلی، شامل تعداد  $L$  ویژگی است. در هر روش سراسری،  $F_{Global} \subseteq F$  را به عنوان مجموعه‌ای از

#### ۴- نتایج

در این بخش ابتدا نتایج استفاده از روش‌های انتخاب ویژگی تکی و ترکیبی فیلتر محلی و سراسری در دسته‌بندی‌های مختلف بر روی داده‌های آموزشی بررسی می‌شود. با توجه به نتایج به‌دست آمده، مجموعه ویژگی بهینه انتخاب می‌شود. در گام بعد، نتایج حاصل از هر دسته‌بندی با استفاده از مجموعه بهینه، بر روی داده‌های آزمون مورد بررسی قرار می‌گیرد.

#### ۴-۱ ارزیابی روش‌های انتخاب ویژگی تکی فیلتر

##### محلی و سراسری

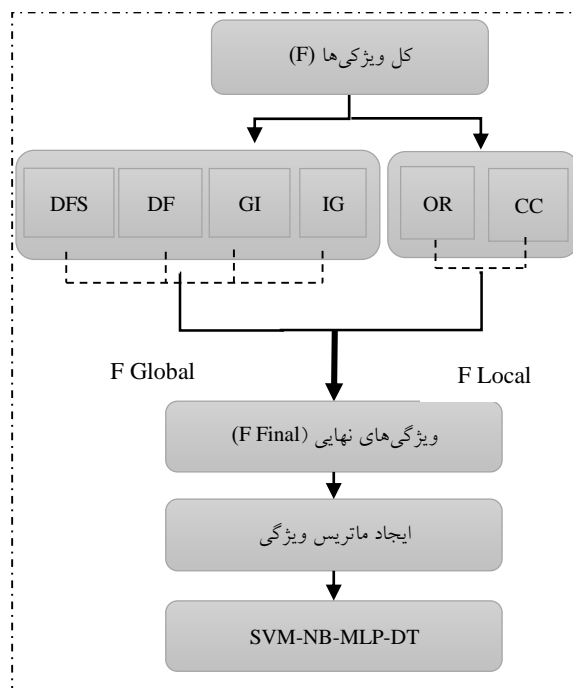
استفاده از روش‌های انتخاب ویژگی تکی بر روی داده‌های آموزشی، در بیشتر موارد سبب افزایش دقت شد. طبق نتایج ذکر شده در جدول (۱) روش انتخاب ویژگی *CC* در میان همه روش‌های انتخاب ویژگی محلی و سراسری، بهترین عملکرد را داشت به طوری که در تمامی روش‌های دسته‌بندی، بالاترین میزان دقت را حاصل کرد. دسته‌بندی‌های ماشین بردار پشتیبان و درخت تصمیم عملکرد بهتر و درصد افزایش دقت بیشتری نسبت به سایر روش‌ها داشته‌اند.

جدول (۱): تعداد ویژگی‌های بهینه در بهترین روش‌های تکی انتخاب

ویژگی محلی و سراسری در معیار *Micro\_F1*

روش دسته‌بندی	روش انتخاب ویژگی بهتر	تعداد ویژگی	دقت با انتخاب ویژگی	دقت بدون انتخاب ویژگی	درصد کاهش ویژگی	درصد افزایش دقت
<b>SVM</b>	CC	۴۰۰	۰/۸۳۵	۰/۷۰۵	۹۴/۶۲	۱۸/۵۱
<b>DT</b>	CC	۴۰۰	۰/۸۰۳	۰/۷۰۲	۹۴/۶۲	۱۴/۳۰
<b>NB</b>	CC	۴۰۰	۰/۸۳۵	۰/۷۵۵	۹۴/۶۲	۱۰/۵۳
<b>MLP</b>	CC	۳۰۰	۰/۸۲	۰/۷۹۶	۹۵/۹۶	۲/۹۵

در دسته‌بندی ماشین بردار پشتیبان، روش انتخاب ویژگی *CC* مقدار دقت بالاتری با تعداد ۴۰۰ ویژگی داشت. دقت این روش ۰/۸۳۵ است. مقدار دقت بدون کاهش ویژگی برای این دسته‌بندی ۰/۷۰۵ بود. بنابراین استفاده از روش‌های کاهش ویژگی دقت را تا ۱۸/۵۱



شکل (۱): روند اجرایی روش پیشنهادی

#### ۳-۳ ارزیابی

در این پژوهش عملکرد روش‌های یادگیری ماشین با تعداد متفاوتی از ویژگی‌ها، بر روی نمونه‌های ارزیابی، مورد بررسی قرار گرفت و در هر روش دسته‌بندی، بهینه‌ترین تعداد ویژگی انتخاب شد. مجموعه بهینه از ویژگی‌ها، مجموعه‌ای است که با تعداد ویژگی کمتر، بیش‌ترین افزایش دقت را در دسته‌بندی حاصل کند. برای انتخاب زیرمجموعه بهینه، روند کار بدین صورت است که در هر مرحله تعداد ۵۰، ۱۰۰، ۲۰۰، تا ۴۰۰ ویژگی انتخاب می‌شود. سپس این ویژگی‌ها در یادگیرنده‌های متفاوت قرار گرفته و ارزیابی می‌شوند. به عبارتی تعداد ویژگی‌های بهینه را با استفاده از دادگان آموزش به دست می‌آوریم.

بعد از به دست آوردن مجموعه ویژگی بهینه مربوط به هر روش دسته‌بندی، در گام بعد، عملکرد روش‌های متفاوت دسته‌بندی بر روی داده‌های آزمون با این مجموعه ویژگی بهینه، مورد بررسی قرار گرفت که نتایج آن در بخش ۳ بیان شده است.



روش‌های کاهش ویژگی دقت را تا ۱۴ درصد افزایش داده است که در مقایسه با دسته‌بند ماشین بردار پشتیبان درصد کمتری است. جدول (۲): تعداد ویژگی‌های بهینه در بهترین ترکیب روش‌های انتخاب

ویژگی محلی و سراسری در معیار **Micro\_F1**

درصد افزایش دقت	درصد کاهش ویژگی	دقت بدون انتخاب ویژگی	دقت با انتخاب ویژگی	تعداد ویژگی بهینه	روش انتخاب ویژگی بهتر	روش دسته‌بندی
۱۶/۵۹	۹۴/۶۲	۰/۷۰۵	۰/۸۲۲	۴۰۰	CC_GI	SVM
۱۴/۷۴	۹۴/۶۲	۰/۷۰۲	۰/۸۰۵	۴۰۰	CC_DFS	DT
۱۰/۰۶	۹۴/۶۲	۰/۷۵۵	۰/۸۳۱	۴۰۰	CC_DFS	NB
۳/۳۸	۹۴/۶۲	۰/۷۹۶	۰/۸۲۳	۴۰۰	CC_GI	MLP

پس از بررسی نتایج مختلف، تعداد ویژگی بهینه برای داده‌های این پژوهش، ۴۰۰ ویژگی در نظر گرفته شد. لذا در مرحله بعد یادگیرنده‌های مختلف را بر روی داده‌های آزمون با زیرمجموعه ویژگی بهینه به دست آمده، موردسنجش قرار می‌دهیم.

در جدول (۳) و (۴) مقادیر نتایج ارزیابی یادگیرنده‌های مختلف بر روی داده‌های آزمون قابل مشاهده است که به طور میانگین در روش‌های انتخاب ویژگی متفاوت، ترکیب روش‌های انتخاب ویژگی مبتنی بر فیلتر محلی و سراسری، عملکرد دسته‌بندی بهتری نسبت به روش‌های انتخاب ویژگی تکی داشته است. این ترکیب‌ها، با تولید زیرمجموعه بهینه از تمامی ویژگی‌های مهم و مؤثر، ضمن کاهش ابعاد ویژگی‌ها، سبب افزایش دقت شدند.

جدول (۳): دقت دسته‌بند **NB** و **SVM** با استفاده از روش‌های

انتخاب ویژگی در معیار **Micro\_F1**

روش انتخاب ویژگی	NB	SVM
No FS	۰/۶۴۸۵	۰/۶۰۳۲
OR_DFS	۰/۷۰۵۳	۰/۷۰۲۵
OR_IG	۰/۵۷۵۶	۰/۶۰۲۵
OR_DF	۰/۷۲۳	۰/۷۰۵
OR_GI	۰/۷۰۳۵	۰/۷۰۲۲
CC_DFS	۰/۷۱۲۵	۰/۷۱۱۴
CC_IG	۰/۶۹۴۳	۰/۶۹۵۸

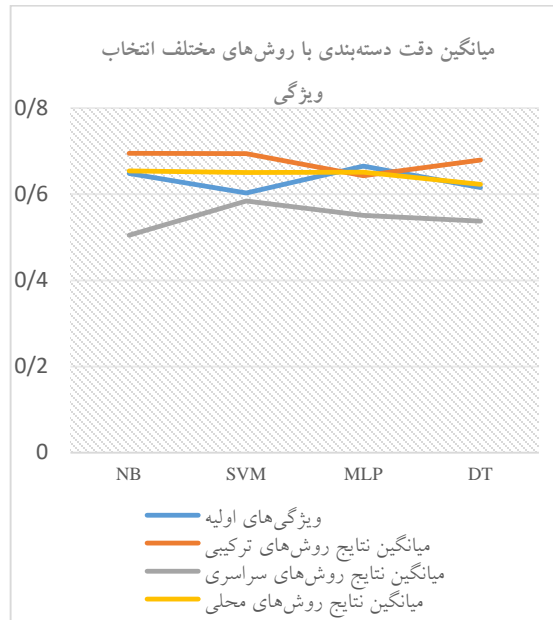
درصد افزایش داده است. دسته‌بند درخت تصمیم برای روش انتخاب ویژگی **CC** مقدار دقت بالاتری با ۴۰۰ ویژگی داشت. دقت این روش ۰/۸۰ است. مقدار دقت بدون کاهش ویژگی برای این دسته‌بند ۰/۷۰۲ بود. بنابراین استفاده از روش‌های کاهش ویژگی دقت را تا ۱۴/۳۰ درصد افزایش داده است که در مقایسه با دسته‌بند ماشین بردار پشتیبان درصد کمتری است.

## ۴-۲ ارزیابی ترکیب روش‌های انتخاب ویژگی محلی و سراسری

در این بخش، عملکرد ترکیبی روش‌های انتخاب ویژگی محلی و سراسری بر روی داده‌های آموزشی مقایسه می‌شود. روش‌های موردبررسی را با تعداد ویژگی‌های متفاوت اجرا کرده و نتایج را مقایسه کردیم. طبق نتایج جدول (۲) استفاده از روش‌های انتخاب ویژگی ترکیبی در اکثر موارد سبب افزایش دقت شد. میزان افزایش دقت با استفاده از روش‌های کاهش ویژگی نسبت به دقت‌های بدون کاهش ویژگی به ترتیب مربوط به دسته‌بند‌های ماشین بردار پشتیبان، درخت تصمیم و در آخر شبکه عصبی است، که از این میان میزان درصد افزایش دقت در دسته‌بند ماشین بردار پشتیبان قابل ملاحظه است.

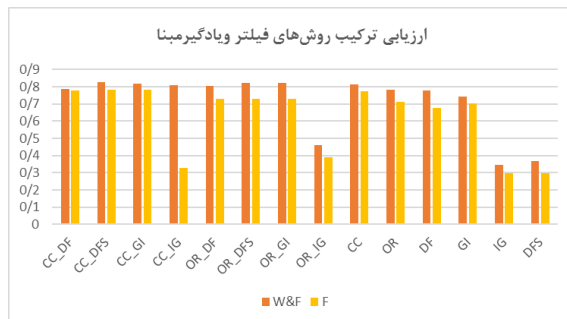
همچنین میزان افزایش دقت با روش‌های انتخاب ویژگی ترکیبی شامل **CC** بهتر است. در دسته‌بند ماشین بردار پشتیبان روش انتخاب ویژگی **CC\_GI** میزان دقت بالاتری داشت. میزان دقت این روش ۰/۸۲ است. میزان دقت، بدون کاهش ویژگی برای این دسته‌بند ۰/۷۲ است. بنابراین استفاده از روش‌های کاهش ویژگی دقت را تا ۱۶ درصد افزایش داده است. در بخش قبل برای روش‌های تکی مشاهده شد که دقت روش **CC** به‌تنهایی نیز دارای مقادیر بالایی بود. دسته‌بند ماشین بردار پشتیبان در مقایسه با دسته‌بند‌های دیگر، نسبت به مقدارهای بدون کاهش ویژگی دارای درصد افزایش دقت بیشتری برای روش‌های ترکیبی است. دسته‌بند درخت تصمیم برای روش انتخاب ویژگی ترکیبی **CC\_DFS** مقدار دقت بالاتری با ۴۰۰ ویژگی داشت. مقدار دقت این روش ۰/۸۰ است. درحالی‌که مقدار دقت بدون کاهش ویژگی برای این دسته‌بند ۰/۷۰ است. بنابراین در درخت تصمیم استفاده از

سراسری است. در میان روش‌های تکی، روش‌های محلی بهتر از روش‌های سراسری در دسته‌بندی عمل کردند.



شکل (۲): مقایسه نتایج دسته‌بندی با روش‌های انتخاب ویژگی

#### فیلتر تکی و ترکیبی



شکل (۳): مقایسه نتایج دسته‌بندی با استفاده از روش فیلتر و روش ترکیبی فیلتر و یادگیرمبنا با ۵۰ ویژگی در دسته‌بند NB و

معیار **Micro-F1**

### ۴-۳ ارزیابی ترکیب روش‌های انتخاب ویژگی فیلتر و یادگیرمبنا

همان‌طور که در شکل (۳) مشاهده می‌شود، ترکیب روش‌های فیلتر تکی و ترکیبی با روش یادگیرمبنا *SFS* دارای دقت بالاتری است. بالاترین دقت مربوط به ترکیب روش *SFS* با روش‌های فیلتری است که شامل ضریب همبستگی یا *CC* هستند. همچنین از میان روش‌های تکی، ترکیب *SFS* با *CC* یا ضریب همبستگی که از

۰/۷۱۶۵	۰/۷۲۵۴	CC_DF
۰/۷۱۶۵	۰/۷۱۹۵	CC_GI
۰/۵۹۶۴	۰/۴۲۵۶	DFS
۰/۵۳۲۸	۰/۴۰۲۵	IG
۰/۶۰۹۵	۰/۵۹۴	GI
۰/۵۹۸۸	۰/۵۹۶۴	DF
۰/۵۸۶۴	۰/۵۶۴۲	OR
۰/۷۱۵۲	۰/۷۴۴۵	CC

جدول (۴): دقت دسته‌بند MLP و DT با استفاده از روش‌های

#### انتخاب ویژگی در معیار **Micro\_F1**

DT	MLP	روش انتخاب ویژگی
۰/۶۱۵۴	۰/۶۶۵	No FS
۰/۷۰۲۹	۰/۷۰۳۴	OR_DFS
۰/۵۶۳۵	۰/۴۷۲۵	OR_IG
۰/۷۱۳۵	۰/۶۹۴	OR_DF
۰/۶۹۷۷	۰/۷۰۱۳	OR_GI
۰/۶۹۶۳	۰/۷۱۴۵	CC_DFS
۰/۶۶۸۵	۰/۳۹۶۸	CC_IG
۰/۶۹۹۴	۰/۷۳۶۸	CC_DF
۰/۶۹۴۲	۰/۷۲۶۱	CC_GI
۰/۴۸۵۶	۰/۵۴۸۲	DFS
۰/۴۹۶۴	۰/۵۰۳۵	IG
۰/۵۹۹۴	۰/۵۹۶۴	GI
۰/۵۶۸	۰/۵۵۶	DF
۰/۵۶۲۵	۰/۵۶۷۳	OR
۰/۶۸۳۴	۰/۷۳۵۶	CC

شکل (۲): میانگین نتایج روش‌های انتخاب ویژگی متفاوت در انواع دسته‌بندها را نشان می‌دهد. نتایج حاصل، بیانگر بیشترین تأثیر روش ترکیبی انتخاب ویژگی به ترتیب بر روی دسته‌بندهای ماشین بردار پشتیبان، بیزین ساده درخت تصمیم بوده است. در دسته‌بند شبکه عصبی نیز با وجود کاهش دقت نسبت به حالت بدون انتخاب ویژگی، عملکرد این دسته‌بند با کاهش تعداد ویژگی‌ها از ۷۰۰۰ ویژگی به ۴۰۰ ویژگی سنجیده شده است که این کاهش ویژگی سبب کاهش چشم‌گیر زمان پردازش می‌شود و نتیجه بااهمیتی است. همچنین در میان روش‌های انتخاب ویژگی، عملکرد روش‌های ترکیبی به‌طور میانگین بهتر از روش‌های محلی و



پرکاربردترین روش‌های محلی، روش  $CC$  یا ضریب همبستگی دارای دقت بالاتری بود. در روش‌های ترکیبی که با استفاده ترکیب پرکاربردترین روش‌های ترکیبی استفاده شده در این پژوهش به دست آمد، روش  $CC_{DF}$  دقت بالاتری را در دسته‌بندی پیام‌های تلگرام نشان داد. در نهایت بالاترین دقت مربوط به ترکیب روش  $SFS$  با روش‌های فیلتری است که شامل ضریب همبستگی یا  $CC$  هستند. به‌طور کلی استفاده از روش‌های انتخاب ویژگی جهت دسته‌بندی پیام‌های تلگرام به‌منظور شناسایی کاربران ارائه‌دهنده خدمت، ضمن کاهش قابل توجه ویژگی‌ها و انتخاب ویژگی‌های مرتبط، سبب بالا رفتن دقت عملکرد روش‌های یادگیری ماشین و کاهش زمان آموزش مدل‌های یادگیری می‌شود.

## References

- [1] X. Deng, Y. Li, J. Weng, and J. Zhang, (2019), Feature selection for text classification: A review, *Multimedia Tools & Applications*, vol. 78, no. 3, pp. 113-115
- [2] M. Nekkaa and D. Boughaci, (2016), Hybrid harmony search combined with stochastic local search for feature selection, *Neural Processing Letters*, vol. 44, no. 1, pp. 199-220.
- [3] G. BİRİCİK, B. Dri, and A. C. SÖNMEZ, (2012), Abstract feature extraction for text classification, *Turkish Journal of Electrical Engineering & Computer Sciences*, vol. 20, no. Sup. 1, pp. 1137-1159.
- [4] L. M. Abualigah, A. T. Khader, M. A. Al-Betar, and O. A. Alomari, (2017), Text feature selection with a robust weight scheme and dynamic dimension reduction to text document clustering, *Expert Systems with Applications*, vol. 84, pp. 24-36.
- [5] D. Ö. Şahin and E. Kılıç, (2019), "Two new feature selection metrics for text classification", *Automatika: časopis za automatiku, mjerenje, elektroniku, računarstvo i komunikacije*, vol. 60, no. 2, pp. 162-171.
- [6] G. Chandrashekar and F. Sahin, (2014), "A survey on feature selection methods, *Computers & Electrical Engineering*, vol. 40, no. 1, pp. 16-28.
- [7] A. K. Uysal, (2016), An improved global feature selection scheme for text classification, *Expert systems with Applications*, vol. 43, pp. 82-92.
- [8] D. Agnihotri, K. Verma, and P. Tripathi, (2017), Variable global feature selection scheme for automatic classification of text documents, *Expert Systems with Applications*, vol. 81, pp. 268-281.

روش‌های محلی است دارای دقت بالاتری است. از میان روش‌های ترکیبی، ترکیب این روش یادگیرمنا با روش  $CC_{DFS}$  دارای دقت بالاتری است.

## ۵- بحث و نتیجه‌گیری

هدف از این پژوهش شناسایی کاربران ارائه‌دهنده خدمات در تلگرام با استفاده از دسته‌بندی پیام‌های ارسالی آن‌ها در گروه‌های مختلف تلگرامی است. برای شروع کار ابتدا پیام‌های ارسال شده توسط هر کاربر در گروه‌های مختلف تلگرامی از سامانه ایده‌کاو استخراج شد. سپس برچسب‌زنی دستی انجام شد. دسته‌بندی دستی اسناد، کاری زمان‌بر و پرهزینه است. لذا استفاده از روش‌های دسته‌بندی خودکار متن، امری ضروری است. روش‌های دسته‌بندی خودکار بر مبنای ویژگی‌های استخراجی از متن، عمل می‌کنند. در این پژوهش ویژگی‌های استخراجی، همان کلمات موجود در پیام‌ها هستند. اما ابعاد بزرگ فضای ویژگی یکی از مشکلات اصلی در دسته‌بندی خودکار است. لذا با استفاده از روش‌های انتخاب ویژگی سعی در کاهش ویژگی‌ها و انتخاب ویژگی‌های مرتبط شد.

در مرحله اول با استفاده از پرکاربردترین روش‌های انتخاب ویژگی مبتنی بر فیلتر که به دودسته محلی و سراسری تقسیم می‌شوند ویژگی‌های غیرضروری حذف شد. در مرحله دوم، هر یک از این روش‌های محلی و سراسری باهم ترکیب و ویژگی‌های بهتر انتخاب شد. نتایج این آزمایش‌ها نشان داد که ترکیب روش‌های ترکیبی انتخاب ویژگی، در اکثر موارد می‌تواند سبب افزایش دقت نسبت به حالت بدون انتخاب ویژگی، در شناسایی کاربران ارائه‌دهنده خدمات شود.

در مرحله سوم خروجی ویژگی‌های روش ترکیبی مرحله دوم به‌عنوان ورودی برای روش انتخاب و ویژگی یادگیرمنا استفاده شد. با این روش بار دیگر ویژگی‌های کم‌اهمیت حذف می‌شوند. این مراحل با تعداد ویژگی‌های مختلف انجام شد و در مجموع نتایج نشان داد که روش ترکیبی فیلتر و یادگیر مینا عملکرد بهتری نسبت به روش‌های فیلتر دارد.

به‌طور میانگین، از میان پرکاربردترین روش‌های سراسری که در این پژوهش استفاده شده است، روش  $DF$  یا فرکانس سند و  $GI$  جهت دسته‌بندی پیام‌های تلگرام دارای دقت بالاتری بودند. از میان



- [9] A. Melo and H. Paulheim, (2019), Local and global feature selection for multilabel classification with binary relevance, *Artificial intelligence review*, vol. 51, no.1, pp. 33-60.
- [10] A. Kursat. Uysal, (2018), "On two-stage feature selection methods for text classification", *IEEE Access*, vol. 6, pp.43233-43251.
- [11] G. Wei, J. Zhao, Y. Feng, A. He, and J. Yu, (2020), A novel hybrid feature selection method based on dynamic feature importance, *Applied Soft Computing*, vol.93, p. 106337.
- [12] L. N. H. Nam and H. B. Quoc, (2017), The hybrid filter feature selection methods for improving high-dimensional text categorization, *International Journal of Uncertainty, Fuzziness and Knowledge-Based Systems*, vol. 25, no. 02, pp.235.
- [13] Asif, M., Nagra, A. A., Ahmad, M. B., & Masood, K, (2022),.. Feature selection empowered by self-inertia weight adaptive particle swarm optimization for text classification. *Applied Artificial Intelligence*, vol. 36, no. 01, p.2004345.
- [14] Effrosynidis, D., & Arampatzis, A. (2021). An evaluation of feature selection methods for environmental data. *Ecological Informatics*, vol. 61, p.101224.

## service providers Identification in Telegram Persian messages based on feature selection methods

Farzaneh Ebrahimi<sup>1</sup>, Mohammad Ali Zare Chahooki<sup>2\*</sup>, Ali Hashemi<sup>3</sup>

<sup>1</sup>Senior student of Computer Engineering, Yazd University

<sup>2</sup>Assistant Professor, Faculty of Computer Engineering, Yazd University

<sup>3</sup>PhD in Computer Engineering, Yazd University

### Article Information

#### Original Research Paper

#### Received:

2022 July 31

#### Accepted:

2023 January 30

#### Keywords:

classification, Service Provider Users, Feature Selection, Feature Reduction, Machine Learning.

#### Corresponding Author\*:

chahooki@yazd.ac.ir

### Abstract

Telegram Messenger is a suitable platform for users who are looking to buy a product or receive services online. In these messengers, it is not possible to have direct access to the providers of goods and services, and in order to request the product, one must first become a member of the related groups and channels telegram. The purpose of this study is to directly identify users of service providers using the classification of Persian messages published in Telegram. One of the problems with categorizing these messages is the large size of the feature space, which reduces accuracy and increases classification time. Feature selection methods were used to solve this problem. The proposed method of this research is based on a combination of feature selection methods based on local and global filters. In this regard, in the first step, using the most widely used methods for selecting local and global filter feature, related features are selected. In the second step, a combination of local and global filtering methods is used to identify better features and increase classification accuracy. The innovation of this research is in using the combined methods of feature selection for automatic classification of Telegram Persian messages, in order to identify the users of the service provider. The proposed method, while reducing the number of features and selecting related features, improves the performance of classification and Identification of service providers.



: 10.22034/ABMIR.2023.18780.1012

E-ISSN: [2821-2037](https://doi.org/10.22034/ABMIR.2023.18780.1012) /© 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/>).

