

## Improving User Relationship Prediction in Twitter Metadata Using Aggregate Classification

**Mohammad Chakerolhoseini Firouzabad:** Doctoral student of Computer Engineering, Faculty of Technology and Engineering, Islamic Azad University, Neyshabur branch, Neyshabur, Iran.

**email:** m.chaker71@iau-neishabour.ac.ir

**Reza Ghaemi:** Assistant Professor, Department of Computer Engineering, Faculty of Technology and Engineering, Islamic Azad University (corresponding author), Quchan Branch, Quchan, Iran.

**email:** r.ghaemi@iauu.ac.ir

In today's world, social networks that have become a part of people's daily life, including Twitter, Telegram, Instagram, etc., are increasing and expanding day by day. Therefore, the number of their users is also increasing and as a result, a large amount of data is being exchanged and stored in these network; and this huge amount of data has turned social networks, especially Twitter, into big data. It is very important to manage, organize and prune these big data, as well as to predict the behavior of social network users.

One of the most important and effective methods for predicting user relationships in social networks is classification techniques, which in most of the applications and researches in the background of the research, are still based on criteria such as 'accuracy; and accuracy of prediction. have weakness In this article, in order to predict the user relationship in Twitter social networks, the cumulative classification method based on voting, which has two basic steps, has been used. In the first step, by using basic classification algorithms including nearest neighbor, decision tree, random forest and simple Bayesian, the outputs of each classification are obtained. In the second step, the final output of cumulative classification is calculated using the voting method. The results of the experiments on the dataset of the Twitter social network and based on the criteria of accuracy, correctness and coverage, argue that the proposed cumulative classification method based on voting has more favorable results than It has other algorithms

**Keywords:** Big Data, Social Network, Prediction of User Relationship, Cumulative Classification, Twitter.

## بهبود پیش‌بینی علاقه کاربران در کلان‌داده توییت‌ها با استفاده از طبقه‌بند تجمعی

تاریخ دریافت: ۱۴۰۱/۰۹/۲۰

تاریخ پذیرش: ۱۴۰۲/۰۷/۰۳

نوشته

محمد چاکر الحسینی فیروزآباد\*

رضا قائمی\*\*

### چکیده

در دنیای امروزی، شبکه‌های اجتماعی که بخشی از زندگی روزمره انسان‌ها شده‌اند، از جمله توییت، تلگرام، اینستاگرام و غیره، روزبه‌روز در حال افزایش و گسترش هستند. لذا تعداد کاربران آن‌ها نیز در حال افزایش است و در نتیجه، حجم داده‌های زیادی در این شبکه‌ها در حال تبادل و ذخیره‌سازی است که این حجم عظیم داده، شبکه‌های اجتماعی به‌خصوص توییت را تبدیل به کلان‌داده کرده است. مدیریت، سامان‌دهی و هرس کردن این کلان‌داده‌ها و همچنین، پیش‌بینی رفتار کاربران شبکه‌های اجتماعی امری بسیار مهم است. یکی از روش‌های مهم و تأثیرگذار برای پیش‌بینی علاقه کاربران در شبکه‌های اجتماعی، تکنیک‌های طبقه‌بندی است که در اغلب کاربردها و پژوهش‌های موجود در پیشینه تحقیق، هنوز در معیارهایی مانند دقت و صحت پیش‌بینی ضعف دارند. در این مقاله، به منظور پیش‌بینی علاقه کاربران در شبکه‌های اجتماعی توییت، از روش طبقه‌بندی تجمعی مبتنی بر رأی‌گیری که دارای دو گام اساسی است، استفاده شده است. در گام نخست، با بهره‌گیری از الگوریتم‌های طبقه‌بندی پایه‌ای شامل نزدیک‌ترین همسایه، درخت تصمیم، جنگل تصادفی و بی‌زین ساده، خروجی‌های هر طبقه‌بندی حاصل می‌شوند. در گام دوم، خروجی نهایی طبقه‌بندی تجمعی با استفاده از روش رأی‌گیری محاسبه می‌شود. نتایج آزمایش‌ها بر روی مجموعه کلان‌داده‌های شبکه اجتماعی توییت و بر اساس معیارهای دقت، صحت و پوشش، استدلال بر این دارد که روش پیشنهادی طبقه‌بندی تجمعی مبتنی بر رأی‌گیری، نتایج مطلوب‌تری را نسبت به الگوریتم‌های دیگر داشته است.

کلیدواژه: کلان‌داده، شبکه اجتماعی، پیش‌بینی علاقه کاربران، طبقه‌بندی تجمعی، توییت.

\* دانشجوی دکتری مهندسی کامپیوتر، دانشکده فنی و مهندسی، دانشگاه آزاد اسلامی، واحد نیشابور، نیشابور، ایران  
m.chaker71@iau-neishabour.ac.ir

\*\* استادیار گروه مهندسی کامپیوتر، دانشکده فنی و مهندسی، دانشگاه آزاد اسلامی (نویسنده مسئول)، واحد قوچان، قوچان، ایران  
r.ghaemi@iauc.ac.ir

نحوه استناد به این مقاله: چاکر الحسینی، محمد و قائمی، رضا (۱۴۰۳). بهبود پیش‌بینی علاقه کاربران در کلان‌داده توییت‌ها با استفاده از طبقه‌بند تجمعی. رسانه، ۳۵(۲)، ۱۰۷-۱۳۱.

## مقدمه

امروزه، بیشتر امور و کارهای ضروری انسان‌ها در بستر اینترنت و شبکه‌های اجتماعی<sup>۱</sup> قرار گرفته است. شبکه‌های اجتماعی با توجه به استقبال زیاد کاربران از این محیط‌ها در جهت تبادل پیام و اشتراک‌گذاری تصویر و ویدئو، با رشد گسترده‌ای در سال‌های اخیر مواجه بوده‌اند که توییتر<sup>۲</sup> یکی از پرکاربرترین این شبکه‌های اجتماعی است (بیهقی<sup>۳</sup> و همکاران، ۲۰۱۸: ۸). شبکه اجتماعی توییتر، فناوری مدرن برقراری ارتباط برای انسان را شکل داده و این غیرقابل انکار است که به‌طور فزاینده‌ای صرف‌نظر از موقعیت جغرافیایی، خودش را به زندگی انسان‌ها القا کرده است. طبق آخرین آمار، حدود یک میلیارد نفر از شبکه‌های اجتماعی استفاده می‌کنند که سهم توییتر، حدود ۳۳۱ میلیون کاربر فعال تا سال ۲۰۱۹ بوده است. به‌علاوه، داده‌های تولیدشده از فعالیت روزمره کاربران در شبکه‌های اجتماعی باعث ایجاد حجم زیادی داده می‌شود و لذا، کلان‌داده‌ها<sup>۴</sup> در شبکه‌های اجتماعی نقش به‌سزایی دارد و بسیار مورد استفاده قرار می‌گیرد (بیهقی و همکاران، ۲۰۱۸: ۸).

داده‌های بزرگ در کلان‌داده‌ها، از حجم عظیمی از داده‌ها به‌صورت ساخت‌یافته، نیمه‌ساخت‌یافته و یا بدون ساختار در حوزه‌های متفاوتی، مانند توییت‌ها، نظرها، پست‌ها و غیره جمع‌آوری می‌شوند. از این رو، از رایج‌ترین کاربردهای کلان‌داده‌های ایجادشده توسط شبکه‌های اجتماعی مانند توییتر، تجزیه و تحلیل نظرکاوی<sup>۵</sup> و برای نمونه، به‌دست‌آوردن بازخورد مشتریان از محصول‌ها است که برای سازمان‌ها بسیار باارزش است. بنابراین، در تجزیه و تحلیل داده‌های حجیم شبکه‌های اجتماعی از تکنیک‌های یادگیری ماشین<sup>۶</sup>، داده‌کاوی<sup>۷</sup> و متن‌کاوی<sup>۸</sup> استفاده می‌شوند. از پرکاربردترین این تکنیک‌ها به‌منظور پیش‌بینی علاقه کاربران در شبکه‌های اجتماعی، می‌توان به طبقه‌بندی<sup>۹</sup> و خوشه‌بندی<sup>۱۰</sup> اشاره کرد (مورثی<sup>۱۱</sup>، ۲۰۱۷: ۹).

تاکنون با استفاده از تکنیک‌های طبقه‌بندی، تحقیق‌های زیادی روی کلان‌داده توییتر صورت گرفته است؛ نظیر تشخیص واکسینه‌شدن افراد از کرونا، تجزیه و تحلیل تجارت الکترونیک، پیش‌بینی نتایج انتخابات ریاست‌جمهوری و میزان استقبال یا دافعه انسان‌ها در دنیا نسبت به

1. Social Networks
2. Twitter
3. Bayhaqy
4. Big Data
5. Opinion Mining
6. Machine Learning
7. Data Mining
8. Text Mining
9. Classification
10. Clustering
11. Murthy

یک تصمیم یا یک رویداد مهم که کلیه این پیش‌بینی‌ها از طریق توییت‌های شبکه اجتماعی توییت‌ر صورت گرفته است (مورثی، ۲۰۱۷: ۹). در این مقاله، به پیش‌بینی و طبقه‌بندی علاقه کاربران در شبکه اجتماعی توییت‌ر پرداخته می‌شود. در این راستا، از الگوریتم‌های طبقه‌بندی پایه و قدرتمندی شامل K- نزدیک‌ترین همسایه<sup>۱</sup>، بیزین ساده<sup>۲</sup>، جنگل تصادفی<sup>۳</sup> و درخت تصمیم<sup>۴</sup> بهره‌گیری شده است. به علاوه، به منظور بهبود معیارهای دقت<sup>۵</sup>، پوشش<sup>۶</sup> و صحت<sup>۷</sup> پیش‌بینی علاقه کاربران در شبکه اجتماعی توییت‌ر، از طبقه‌بندی تجمعی<sup>۸</sup> مبتنی بر رأی‌گیری<sup>۹</sup> استفاده شده است.

در بخش دوم به مرور ادبیات و پیشینه تحقیق شامل کلان‌داده و طبقه‌بندی‌ها پرداخته شده است. بخش سوم، الگوریتم‌های طبقه‌بندی پایه‌ای مورد استفاده و همچنین، الگوریتم طبقه‌بندی تجمعی پیشنهادی تشریح شده‌اند. در بخش چهارم، نتایج حاصل‌شده از الگوریتم‌های طبقه‌بندی پایه و همچنین، طبقه‌بندی تجمعی مبتنی بر رأی‌گیری پیشنهادی مورد ارزیابی قرار گرفته است. در نهایت بخش پنجم، به نتیجه‌گیری و کارهای آینده پرداخته است.

## ادبیات و پیشینه تحقیق

در این بخش، در مورد مفاهیم اصلی کلان‌داده‌ها، شبکه‌های اجتماعی و طبقه‌بندی توضیح داده شده است و سپس، عملکرد الگوریتم‌های طبقه‌بندی پایه‌ای مورد استفاده در این مقاله شامل K- نزدیک‌ترین همسایه، بیزین ساده، جنگل تصادفی و درخت تصمیم تشریح شده‌اند. در ادامه، پیشینه تحقیق در زمینه پیش‌بینی و طبقه‌بندی علاقه کاربران در شبکه‌های اجتماعی به‌خصوص رسانه توییت‌ر ارائه شده است.

### ۱. کلان‌داده‌ها

کلان‌داده اصطلاحی است برای جمع‌آوری مجموعه‌های اطلاعاتی گران و پیچیده‌ای که پردازش آن با استفاده از برنامه‌های کاربردی مرسوم پردازش اطلاعات، دشوار است. چالش‌هایی از جمله پیچیدگی حافظه و زمانی بالای ذخیره‌سازی، مدیریت، پرس و جو، اشتراک‌گذاری، تبادل، ادراک و نقض حفاظت در کلان‌داده‌ها وجود دارند. برای کاهش الگوهای تجاری، پیش‌بینی

1. K-Nearest Neighbor
2. Naive Bayesian
3. Random Forest
4. Decision Tree
5. Accuracy
6. Recall
7. Precision
8. Classification Ensemble
9. Voting

بیماری‌ها، شناسایی تضادها و غیره، در مقایسه با مجموعه داده‌های کوچک‌تر، به مجموعه‌های داده بزرگ‌تری نیاز است. کار با حجم عظیمی از اطلاعات با استفاده از چارچوب‌های مدیریت پایگاه داده اجتماعی و اندازه‌گیری‌های صفحه نمایش و بسته‌های ادراک مشکل است و نیاز به برنامه‌نویسی موازی بسیار زیادی دارد که روی ده‌ها، صدها یا حتی تعداد زیادی سرورس دهنده اجرا می‌شوند (سری و آنوشا، ۲۰۱۶: ۱).

برخلاف داده‌های سنتی، کلان‌داده به داده‌های بزرگ در حال رشد اشاره دارد و مجموعه داده‌هایی که دارای ساختارهای ناهمگن هستند، شامل داده‌های ساخت یافته، بدون ساختار و نیمه‌ساخت یافته. کلان‌داده ماهیت پیچیده‌ای دارد که به فناوری‌های قدرتمند و الگوریتم‌های پیشرفته‌ای نیاز دارد (اوسوس و بتجلن<sup>۲</sup> و همکاران، ۲۰۱۸: ۲). شکل ۱ پنج اصل در کلان‌داده‌ها شامل ظرفیت، تنوع، صحت، ارزش و سرعت را نشان داده است. از آنجایی که داده‌ها از منابع بسیار متفاوتی جمع‌آوری می‌شوند، صحت به مفهوم کیفیت، دقت و قابلیت اطمینان داده‌های جمع‌آوری شده است. همچنین، ارزش به توانایی سازمان برای تبدیل داده‌های عظیم به کسب‌وکار واقعی اشاره دارد، چراکه داده‌های دقیق باعث می‌شوند کسب‌وکارها به رفع نیازها و انتظارات مشتریان خود نزدیک‌تر شوند (پاپاکیراکو و باربوناکس<sup>۳</sup>، ۲۰۲۲: ۵).

ظرفیت	تنوع	صحت	ارزش	سرعت
ترازایی	ساخت یافته	اعتبار	آماری	روی خط و برون خط داده
سوابق	نیمه‌ساخت یافته	مسئولیت	رویدادها	زمان نزدیک و واقعی
جداول - فایل‌ها	بدون ساختار	در دسترس	همبستگی	جریان‌ها
تراکش‌ها	همه موارد بالا	قابل اعتماد	فرضی	

شکل ۱. پنج اصل در کلان‌داده

## ۲. شبکه‌های اجتماعی

شبکه و رسانه‌های اجتماعی، مجموعه‌ای از برنامه‌های کاربردی مبتنی بر اینترنت است که بر اساس ایده وب ۲ ایجاد شد و ابتدا در سال ۲۰۰۴ برای توصیف محتویات و برنامه‌های کاربردی مورد استفاده قرار گرفت و می‌توان آن را به‌جای این‌که به‌طور سنتی توسط افراد ایجاد، تهیه و منتشر شود توسط کاربران به‌طور مداوم و به‌طور مختلف از طریق مشارکت و همکاری تغییر داد. استفاده گسترده از نرم‌افزار و سخت‌افزار موجود برای دسترسی به

1. Sri & Anusha
2. Oussous & Benjelloun
3. Papakyriakou & Barbounakis

چارچوب‌های شبکه و رسانه‌های اجتماعی از طریق اینترنت، منجر به ایجاد و تبادل محتوای تولیدشده توسط کاربران شده است (قانی<sup>۱</sup> و همکاران، ۲۰۱۹: ۳).

شبکه‌های اجتماعی به‌عنوان خدمات مبتنی بر وب تعریف می‌شوند که ابتدا در آن، افراد مجاز به ایجاد نمایه عمومی یا نیمه‌عمومی خود هستند. سپس، این کاربران مجاز به ارتباط با دیگران‌اند و در نهایت برای تشکیل یک شبکه، این افراد مجاز به مشاهده و ارتباط با سایر کاربران و فعالیت‌هایی هستند که در شبکه آن‌ها منتشر می‌شوند. نمونه‌هایی از شبکه‌های اجتماعی که حجم زیادی از داده‌های بدون‌ساختار را دارا هستند عبارت است از فیس‌بوک، توئیتر، اینستاگرام، لینکدین، وب‌نوشت‌ها، ویکی‌ها و یوتیوب (قانی و حمید و همکاران، ۲۰۱۹: ۳).

### پیش‌بینی و طبقه‌بندی

داده‌کاوی فرایند استخراج دانش از داده‌های به‌نسبت عظیم است. داده‌کاوی دارای سه تکنیک اصلی طبقه‌بندی، خوشه‌بندی، الگوهای مکرر و قوانین وابستگی آن‌ها است. در خوشه‌بندی، مجموعه‌ای از داده‌ها بر اساس شباهت و تفاوت بین ویژگی‌هایشان گروه‌بندی می‌شوند. مدل خوشه‌بندی حاصل می‌تواند برای طبقه‌بندی داده‌های آینده مورد استفاده قرار گیرد. داده‌کاوی فرایند استخراج اطلاعات از یک مجموعه داده و تبدیل آن به یک ساختار قابل درک است. وظیفه داده‌کاوی، تحلیل خودکار یا نیمه‌خودکار مقادیر زیادی از داده‌ها برای استخراج الگوهای جالب ناشناخته قبلی است. داده‌کاوی شامل شش وظیفه اصلی است که عبارت‌اند از تشخیص ناهنجاری، یادگیری قوانین وابستگی، خوشه‌بندی، طبقه‌بندی، رگرسیون و خلاصه‌سازی. طبقه‌بندی یک تکنیک پرکاربرد در داده‌کاوی است و به‌طور گسترده‌ای در زمینه‌های مختلف مورد استفاده قرار گرفته است که برای پیش‌بینی عضویت گروهی از نمونه‌های آزمون در یک کلاس مشخص، استفاده می‌شود. در این مقاله، چندین الگوریتم طبقه‌بندی پایه‌ای استفاده شده است، شامل K- نزدیک‌ترین همسایه، بیزین ساده، جنگل تصادفی و درخت تصمیم (کساورج و سوکوماران، ۲۰۱۳: ۴).

طبقه‌بندی، مسئله پیش‌بینی گسسته متغیر تصادفی Y از متغیر تصادفی X است که طبقه‌بندی الگو یا تشخیص الگو نیز نامیده می‌شود. الگوریتم طبقه‌بندی، داده‌ها را به تعداد معینی از کلاس‌ها دسته‌بندی می‌کند و برچسبی را به هر کلاس اختصاص می‌دهد. ایده اصلی در الگوریتم‌های طبقه‌بندی به‌منظور پیش‌بینی کلاس هدف توسط تجزیه و تحلیل مجموعه داده‌های آموزشی، دسته‌بندی داده‌ها به تعداد معینی از کلاس‌ها است (پاپاکیراکو و باریناکس، ۲۰۲۲: ۵).

1. Ghani
2. Kesavaraj & Sukumaran

## ۱.K- نزدیک ترین همسایه

الگوریتم طبقه بندی K- نزدیک ترین همسایه، بر اساس یادگیری توسط مقایسه عمل می کند که در آن، ویژگی های n بعدی هر نمونه داده، نشان دهنده یک نقطه در فضای n بعدی است. نمونه های آموزشی در یک فضای الگوی n بعدی ذخیره می شوند. هنگامی که یک نمونه داده آزمون توسط K- نزدیک ترین همسایه طبقه بندی می شود، فضای الگو را برای K نمونه آموزشی که نزدیک ترین به نمونه آزمون هستند، جست و جو می کند و فاصله نمونه آزمون را بر حسب فاصله اقلیدسی مطابق رابطه (۱)، با سایر نمونه های آموزشی محاسبه می کند. نمونه آزمون، دارای رایج ترین کلاس در میان K- نزدیک ترین همسایه خود است. زمانی که  $K=1$  باشد، نمونه آزمون به کلاس نزدیک ترین نمونه آموزشی به آن اختصاص داده می شود (فیو، ۲۰۰۹: ۶).

$$d(X, Y) = \sqrt{\sum_{i=1}^n (x_i - y_i)^2} \quad (1)$$

در رابطه (۱)، متغیرهای X و Y دو نقطه مورد نظر هستند که با محاسبه مجموع اختلاف آن ها، فاصله بین دو نقطه محاسبه می شود. الگوریتم طبقه بندی K- نزدیک ترین همسایه مبتنی بر نمونه و از نوع یادگیرنده های تنبل است. طبقه بندی نمونه آزمون برای طبقه بندی تنبل ممکن است با وجود تعداد زیادی همسایه بالقوه، متحمل هزینه های زمانی زیادی شود. برای رفع این چالش، از الگوریتم طبقه بندی K- نزدیک ترین همسایه وزن دار استفاده می شود که در آن، بر اساس اهمیت هر نمونه آموزشی یا هر ویژگی داده، وزنی به آن انتساب داده شود و زمانی که نمونه آموزشی یا ویژگی نامربوط زیادی وجود داشته باشد، ممکن است مشکل به وجود آورد (فیو، ۲۰۰۹: ۶).

## ۲. بیزین ساده

از تکنیک دسته بند بیزین اغلب به عنوان یک راه کار ساده برای دسته بندی و تعیین روشی برای تشخیص برچسب اشیا یا نقاط استفاده می شود. به منظور به کارگیری دسته بند بیزین ساده، الگوریتم یکتایی وجود ندارد، بلکه خانواده ای از الگوریتم ها موجود است که با فرض استقلال ویژگی ها یا متغیرها نسبت به یکدیگر عمل می کنند. برای نمونه، اندازه یک میوه و رنگ آن که متغیرهای مستقل در نظر گرفته می شوند، در تعیین نوع آن میوه مؤثر هستند. بنابراین، اگر میوه ای با رنگ قرمز و دارای اندازه حدود ۱۰ سانتی متر باشد، با احتمال زیاد سیب است. یکی از مزایای قابل توجه در دسته بند بیزین ساده، امکان برآورد پارامترهای مدل با اندازه نمونه کوچک به عنوان مجموعه داده آموزشی است (فیو، ۲۰۰۹: ۶). بیزین ساده را می توان یک مدل بر مبنای احتمال شرطی در نظر گرفت. فرض کنید  $X=(x_1, \dots, x_n)$  برداری از n ویژگی

را بیان کند که به صورت متغیرهای مستقل هستند. به این ترتیب می توان احتمال رخداد CK، یعنی  $P(CK|x_1, \dots, x_n)$  را به عنوان یکی از حالت های کلاس رخداد های مختلف به ازای Kهای متفاوت، در رابطه (۲) نمایش داد (فیو، ۲۰۰۹: ۶؛ متسیس و آندرو<sup>۱</sup> و همکاران، ۲۰۰۶: ۷).

$$P(Ck | X) = (P(Ck)P(X | Ck))/P(X) \quad (2)$$

رابطه (۲)، همان قضیه بیز است که X بردار و CK رخداد است و بر اساس احتمالات پیشامدهای پیشین<sup>۲</sup>، پسین<sup>۳</sup>، درست نمایی<sup>۴</sup> و شواهد<sup>۵</sup>، در رابطه (۳) باز نویسی شده است (فیو، ۲۰۰۹: ۶):

$$Posterior = \frac{Prior \times Likelihood}{Evidence} \quad (3)$$

به این ترتیب، برای محاسبه احتمال  $P(CK|x_1, \dots, x_n)$  کافی است از احتمال توأم کمک گرفته شود و با احتمال شرطی با توجه به استقلال متغیرها، مطابق رابطه (۴) ساده تر نمایش داده شود (فیو، ۲۰۰۹: ۶):

$$P(Ck, x_1, x_2, \dots, x_n) = p(x_1, x_2, \dots, x_n, Ck) p(x_2 | x_3, \dots, x_n, Ck) \dots p(x_n - 1 | x_n, Ck) p(x_n | Ck) p(Ck) \quad (4)$$

به این ترتیب احتمال توأم را به صورت حاصل ضرب احتمال شرطی می توان نوشت:

$$P(Ck, x_1, x_2, \dots, x_n) = p(Ck) \prod_{i=1}^n p(x_i | Ck) \quad (5)$$

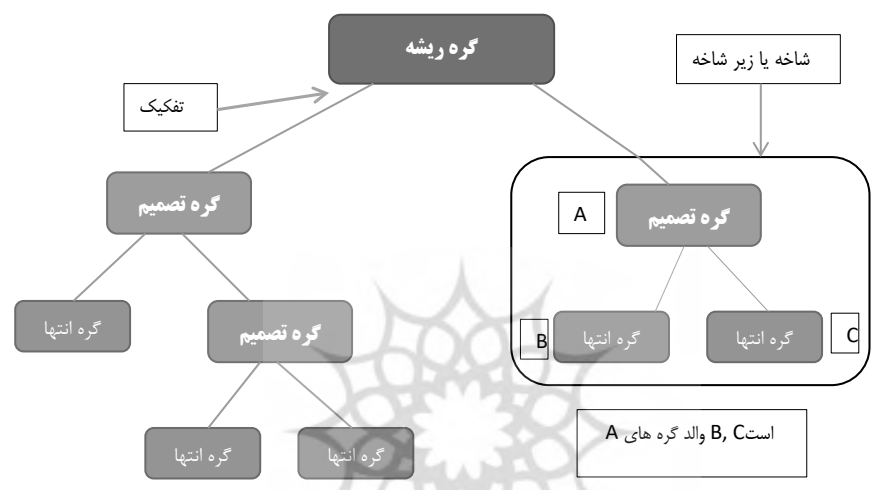
### ۳. درخت تصمیم

درخت تصمیم یکی از روش های قدرتمندی است که اغلب در زمینه های مختلفی نظیر یادگیری ماشین، پردازش تصویر و شناسایی الگوها مورد استفاده قرار می گیرد. درخت تصمیم یک مدل متوالی است که یک سری آزمون های پایه را به طور کارآمد و منسجم متحد می کند که در آن، یک ویژگی عددی با عددی دیگر که مقدار آستانه در هر آزمون است، مقایسه می شوند. درخت تصمیم یک مدل طبقه بندی است که از گره و شاخه ها تشکیل شده است و هر گره، ویژگی هایی را در یک دسته برای طبقه بندی نشان می دهد و هر زیر مجموعه مقداری را تعریف می کند که

1. Metsis & Androutsopoulos
2. Prior
3. Posterior
4. likelihood
5. Evidence



می‌تواند توسط گره گرفته شود. درخت‌های تصمیم به دلیل تجزیه و تحلیل ساده و دقت آن‌ها بر روی فرم‌های داده‌های متعدد، زمینه‌های پیاده‌سازی بسیاری را پیدا کرده‌اند. انواع متفاوتی از الگوریتم‌های درخت تصمیم وجود دارند، از جمله دوگانگی تکراری<sup>۱</sup>، جانشین درخت طبقه‌بندی و رگرسیون<sup>۲</sup>، رگرسیون تطبیقی چندمتغیره<sup>۳</sup>، تعمیم یافته<sup>۴</sup>، بی‌طرفانه<sup>۵</sup>، تشخیص و تخمین تعامل<sup>۶</sup>، درختان استنتاج شرطی<sup>۷</sup> و غیره. شکل ۲، نمونه‌ای از یک درخت تصمیم را نشان می‌دهد (چاربوتی و عبدالله زی<sup>۸</sup>، ۲۰۲۱: ۱۰).



شکل ۲. درخت تصمیم

#### ۴. جنگل تصادفی

همان‌طور که از نام این الگوریتم پیداست، جنگل تصادفی یک الگوریتم طبقه‌بندی نظارت شده است که داده‌ها را با ساخت تعدادی طبقه‌بندی‌کننده با هدف دستیابی به دقت بالاتر پیش‌بینی، طبقه‌بندی می‌کند (شایک و سرین واسان<sup>۹</sup>، ۲۰۱۹: ۱۵). جنگل تصادفی برای رتبه‌بندی اهمیت متغیرها در یک مسئله رگرسیونی یا طبقه‌بندی می‌تواند استفاده شود. نخستین گام، اندازه‌گیری

1. Repetitive Duality
2. Regression
3. Multivariate Adaptive Regression
4. Generalized
5. Impartially
6. Estimation of Interaction
7. Conditional Inference Trees
8. Charbuty & Abdulazeez
9. Shaik & Srinivasan

اهمیت یک متغیر در مجموعه داده  $D_n(X_i, Y_i)=1$  است که برازش جنگل تصادفی است. در طول فرایند، برازش خطای خارج از محدوده هر نقطه از داده ثبت می‌شود و میانگین آن در جنگل محاسبه می‌شود (زکریا، ۲۰۱۴: ۱۲). از مزایای مهم الگوریتم جنگل تصادفی می‌توان به دقت بالا، کارآمدی در مدیریت پایگاه‌های داده‌های بزرگ، و کنترل داده‌های از دست‌رفته بدون به‌خطر انداختن دقت، اشاره کرد. نمونه‌های اولیه برای دادن اطلاعات یا داده‌های متا در رابطه بین متغیرهای مختلف استفاده می‌شود (شایک و سرین واسان، ۲۰۱۹: ۱۵).

## پیشینه تحقیق

در این بخش، به مرور تحقیق‌های موجود از چند سال گذشته در خصوص طبقه‌بندی در شبکه اجتماعی توئیتر، پرداخته می‌شود. مقاله هارجلو و گورجار<sup>۲</sup> و همکاران (۲۰۲۰: ۱۷) با استفاده از روش‌های طبقه‌بندی رگرسیون، بیزین ساده و ماشین بردار پشتیبان<sup>۳</sup> بر روی مجموعه داده Sentiment140 مربوط به دانشگاه استفورد توانستند احساساتی که در توئیتهای تبادلی می‌شوند را بر مبنای دو دسته مثبت و منفی طبقه‌بندی کنند. به دلیل محبوبیت زیاد توئیتر، افراد زیادی هرزنامه‌هایی را به منظور ایجاد حساب‌های جعلی و انتشار پیام‌های مخرب ایجاد می‌کنند. در تحقیق آدوله و هان<sup>۴</sup> و همکاران (۲۰۲۰: ۱۸) با تلفیق الگوریتم خوشه‌بندی K-Means و روش‌های طبقه‌بندی ماشین بردار پشتیبان و جنگل تصادفی، برای تشخیص حساب‌های جعلی به میزان ۹۶/۳ درصد دقت، حساسیت و پوشش بهبود داشته است.

یکی از موارد مهم رفع ابهام برای کاربران، حمایت از تصمیم آن‌ها و سطح رضایتشان است. در پژوهش فان و نگوین<sup>۵</sup> و همکاران (۲۰۲۱: ۱۹) سطح رضایت کاربران در سه مرحله محاسبه شده است. مرحله اول، موضوعی را مشخص می‌کند که کاربر به آن علاقه‌مند است. مرحله دوم، جنبه‌هایی از موضوع و احساسات آن‌ها در توئیتهای استخراج می‌شود. مرحله سوم، سطح رضایت کاربر با توجه به هر نوع احساسات شناسایی و محاسبه می‌شود. در نهایت، با استفاده از درخت تصمیم، بهبود لازم حاصل شده است. یکی از عوامل مهم تجزیه و تحلیل احساسات کاربران در مواجهه با تجارت الکترونیک است. در تحقیق بایاق و سفنیرانتو<sup>۶</sup> و همکاران (۲۰۱۸: ۸) با استفاده از روش‌های طبقه‌بندی نزدیک‌ترین همسایه، درخت تصمیم و بیزین ساده، احساسات کاربران در تجارت الکترونیک تجزیه و تحلیل شد. بیزین ساده عملکرد بهتری در طبقه‌بندی احساسات بر اساس معیارهای دقت، حساسیت و پوشش داشته است.

1. Zakariah
2. Harjule & Gurjar
3. Support Vector Machine
4. Adewole & Han
5. Phan & Nguyen
6. Bayhaqy & Sfenrianto

طبق پژوهش نئوجی و گرج<sup>۱</sup> و همکاران (۲۰۲۱: ۲۱) آموزش در اندونزی امری واجب و پرداخت هزینه آن الزامی است. به علت این که کلیه هزینه‌های مدارس از طریق شهریه انجام می‌پذیرد، اما بسیاری از دانش‌آموزان تأخیر در پرداخت دارند. به همین منظور، با استفاده از روش طبقه‌بندی نزدیک‌ترین همسایه و بیزین ساده نشان داده شد که بیزین ساده در پیش‌بینی دقت بالاتری داشته است.

در تحقیق (عباس و سلیه<sup>۲</sup> و همکاران، ۲۰۲۰: ۲۲) در سال ۲۰۲۱، روی موضوع اعتراض کشاورزان بیش از ۲۰ هزار توئیتر برای شناسایی احساس اعتراضی از طریق روش‌های طبقه‌بندی بیزین ساده، درخت تصمیم، ماشین بردار پشتیبان و جنگل تصادفی انجام شد که روش جنگل تصادفی دقت بالاتری نسبت به سایر روش‌ها داشت. در داده‌های رسانه‌های اجتماعی، توئیتر اغلب از متن مبهم استفاده می‌کند که می‌تواند در شناسایی احساسات مثبت یا منفی مشکل ایجاد کند. بیش از یک میلیارد پیام رسانه‌های اجتماعی وجود دارد که باید در یک پایگاه داده مناسب ذخیره‌سازی شوند و برای تجزیه و تحلیل، به درستی پردازش شوند. در پژوهش شمتر و چاکرابورتی<sup>۳</sup> (۲۰۲۱: ۲۳) با ترکیب روش طبقه‌بندی تجمعی با بیزین ساده، درخت تصمیم، رگرسیون و شبکه عصبی پرسپترون<sup>۴</sup> نشان داده شد که عملکرد بهتر شده است. در دو سال اخیر، به دلیل همه‌گیری ویروس کرونا در جهان، به منظور رفع این ویروس تلاش شد و سه واکسن مدرنا، فایزر و آسترازنکا از مشهورترین واکسن‌ها تولید شدند. در تحقیق شمتر و چاکرابورتی در سال ۲۰۲۱، با استفاده از روش طبقه‌بندی K- نزدیک‌ترین همسایه نظرات توئیتری راجع به واکسن‌های مذکور، به سه دسته مثبت، منفی و خنثی طبقه‌بندی شدند. مهم بودن ویروس کرونا به عنوان یک مسئله حیاتی و همچنین، تجزیه و تحلیل احساسات و نظرسنجی کاربران از طریق توئیترها اهمیت زیادی پیدا کرده است، اما تجزیه و تحلیل‌ها بر روی توئیترهای انگلیسی بوده است.

در پژوهش الحاشدی و الفوحیدی<sup>۵</sup> و همکاران (۲۰۲۰: ۲۴) با استفاده از طبقه‌بندی بیزین ساده برای تحلیل بر روی زبان عربی صورت گرفت. یکی از مسائل مهمی که جنبه تحقیقات جالبی به خصوص برای پزشکی قانونی ارائه می‌دهد، تعیین جنسیت بر مبنای توئیترها است، که در تحقیق و شیس و میحان<sup>۶</sup> (۲۰۲۰: ۲۵) با استفاده از روش پردازش زبان طبیعی<sup>۷</sup> و همچنین، روش‌های طبقه‌بندی بیزین ساده، ماشین بردار پشتیبان و رگرسیون خطی در طبقه‌بندی جنسیت

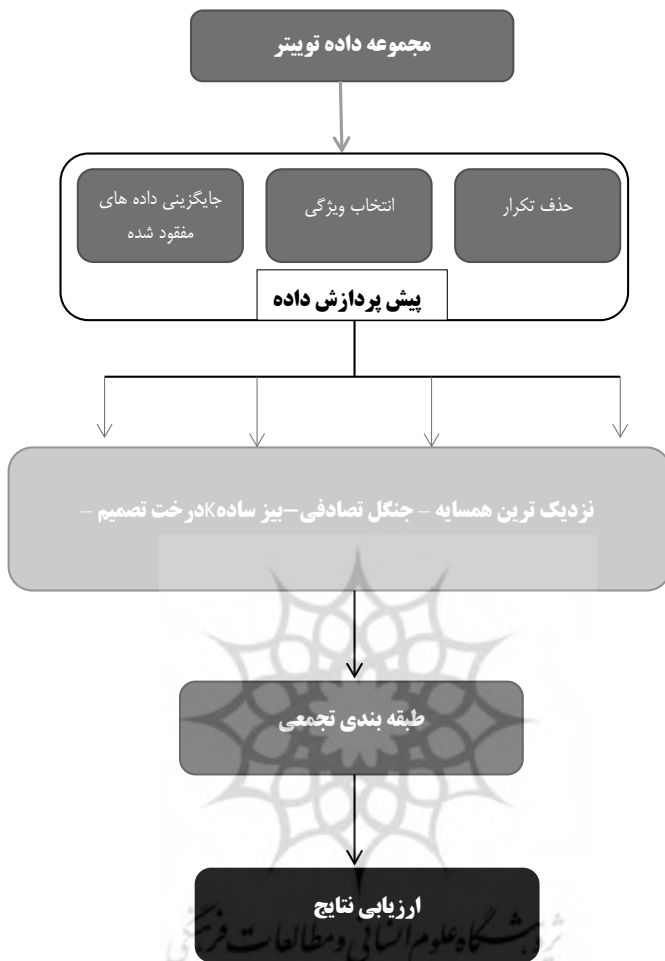
1. Neogi, Garg
2. Abbas, Salih
3. Shamrat, Chakraborty
4. Perceptron Neural Network
5. Al-Hashedi, Al-Fuhaidi
6. Vashisth & Meehan
7. Natural Language Processing

کمک‌کننده بوده است. در رویکرد مشتق‌شده پژوهش پاتل و پسی<sup>۱</sup> (۲۰۲۰: ۲۶)، تحلیلی بر روی داده‌های توئیتر برای جام جهانی فوتبال ۲۰۱۴ برزیل انجام شد تا احساسات مردم در سراسر جهان تجزیه و تحلیل شود. با استفاده روش‌های طبقه‌بندی بیزین ساده، ماشین بردار پشتیبان، جنگل تصادفی و K- نزدیک‌ترین همسایه نشان داده شد که بیزین ساده دقت بالاتری داشته است.

## روش پیشنهادی

از چالش‌های موجود در شبکه‌های اجتماعی به‌خصوص توئیتر، عدم استفاده مطلوب در استفاده از طبقه‌بندها است. انتخاب غیر صحیح و عدم بررسی صحیح در این مورد، در معیارهایی چون دقت، حساسیت، پوشش و غیره، باعث خروجی ناکارآمدی می‌شود، به‌خصوص در کلان‌داده‌ای چون توئیتر که انتخاب بهترین طبقه‌بند در فرایند بهینه‌سازی می‌تواند بسیار مؤثر باشد. شکل ۳ عملکرد روش پیشنهادی پیش‌بینی و طبقه‌بندی علاقه‌کاربران را در کلان‌داده رسانه اجتماعی توئیتر نشان داده است.

گام اول، فرایند جمع‌آوری داده‌ها از توئیتر و برچسب‌گذاری توئیتهای است. پس از برچسب‌گذاری داده‌ها، گام دوم پیش‌پردازش توئیتهای و آماده‌سازی آن‌ها برای تجزیه و تحلیل و تبدیل شدن است که دارای چندین مرحله است. پیش‌پردازش توئیتهای شامل پاک‌سازی، تبدیل نفی، تبدیل شکلک‌ها (ایموجی)، فیلتر کردن و غیره است. پس از پیش‌پردازش توئیتهای، در گام سوم، توئیتهای با استفاده از الگوریتم‌های پایه‌ای، اما قدرتمندی همچون بیزین ساده، K- نزدیک‌ترین همسایه، جنگل تصادفی و درخت تصمیم، طبقه‌بندی می‌شوند. در این مقاله و در گام چهارم، برای پیش‌بینی علاقه‌کاربران، روش طبقه‌بندی تجمعی استفاده شده است که خروجی الگوریتم‌های طبقه‌بندی پایه‌ای را دریافت کرده است و مبتنی بر روش رأی‌گیری، خروجی نهایی را تعیین می‌کند. در نهایت در گام پنجم، روش‌های طبقه‌بندی پایه‌ای و طبقه‌بند تجمعی پیشنهاد شده، بر اساس معیارهای دقت، حساسیت و پوشش ارزیابی شده‌اند.



شکل ۳. روش پیشنهادی پیش‌بینی و طبقه‌بندی علاقه‌کاربران در کلان‌داده تویتر

### ۱. پیش‌پردازش داده‌ها

در پیش‌پردازش، ابتدا داده خروجی از مجموعه داده تویتر، توسط روش جایگزینی داده‌های مفقود شده، مقادیر می‌شوند، به این صورت که داده‌های مفقود شده با مقادیر مشخص شده، جایگزین و مقادیر می‌شوند. در اصل مقادیرهای مفقود شده در تویتر یا بر مبنای سیستم یا بر مبنای کاربر ثبت صورت نمی‌گیرد، و یا در داده‌های آماری باید لحاظ شوند، در حالی که داده‌ای ثبت نشده است که این ثبت نشدن یا حاصل نویز و یا خطا در روال سیستم بوده است و

یا حاصل عدم درج کاربرد و از این رو، در مجموعه داده باعث ایجاد داده مفقود شده می شود. بعد از مرحله جایگزینی، انتخاب ویژگی<sup>۱</sup> صورت می پذیرد و خصوصیت های مهم و مؤثر توییت که به منظور تجزیه و تحلیل مورد نظر است، انتخاب می شود. در نهایت در مرحله پیش پردازش، حذف تکرار<sup>۲</sup> انجام می گیرد، به طوری که هر رکورد که دارای خصوصیت تکراری باشد را حذف می کند و از مقادیر زائد داده و حجم پردازش بیهوده آن ها جلوگیری می شود.

### طبقه بندی های منفرد پایه

در این بخش، به تشریح و توسعه هر یک از الگوریتم های طبقه بندی پایه ای و قدرتمند شامل بیزین ساده، K- نزدیک ترین همسایه، جنگل تصادفی و درخت تصمیم، روی توییت های کلان داده رسانه اجتماعی توییت پرداخته شده است. به عنوان اولین الگوریتم طبقه بندی پایه به منظور پیش بینی علاقه کاربران، از طبقه بندی بیزین ساده استفاده شده است که در شکل ۴ نشان داده شده است (فیو، ۲۰۰۹: ۶).

**ALGORITHM NB TRAINING**

1. Let  $V$  be the vocabulary of ALL words in the documents in  $D$
2. For each category  $c_i \in C$ 
  - Let  $D_i$  be the subset of documents in  $D$  in category  $c_i$
  - $P(c_i) = |D_i| / |D|$
  - Let  $T_i$  be the concatenation of all the documents in  $D_i$
  - Let  $n_i$  be the total number of word occurrences in  $T_i$
  - For each word  $w_j \in V$ 
    - Let  $n_{ij}$  be the number of occurrences of  $w_j$  in  $T_i$
    - Let  $P(w_j | c_i) = (n_{ij} + 1) / (n_i + |V|)$

---

**ALGORITHM NB TESTING**

1. Given a test document  $X$
2. Let  $n$  be the number of word occurrences in  $X$
3. Return the category:
  - $$\arg \max_{c_i \in C} P(c_i) = \prod_{i=1}^n P(a_i | c_i)$$

where  $a_i$  is the word occurring at the  $i^{\text{th}}$  position in  $X$

شکل ۴. الگوریتم طبقه بندی بیزین ساده (پاتیل و پاوار<sup>۳</sup>، ۲۰۱۲: ۲۸)

1. Feature Selection
2. Remove Duplicate
3. Patil & Pawar

دومین الگوریتم طبقه‌بندی پایه برای پیش‌بینی علاقه کاربران، طبقه‌بند K- نزدیک‌ترین همسایه است که مبتنی بر نمونه است و از نوع یادگیرنده‌های تنبیل است که در شکل ۵ نشان داده شده است (فیو، ۲۰۰۹: ۶).

```

K ← number of nearest neighbors
For each object X in the test set do
    calculate the distance D(X,Y) between X and every object Y in the training set
    neighborhood ← the k neighbors in the training set closest to X
    Xclass ← SelectClass (neighborhood)
End for

```

شکل ۵. الگوریتم طبقه‌بندی K- نزدیک‌ترین همسایه (آرچانا و الانگوان، ۲۰۱۴: ۱۳)

به‌عنوان سومین الگوریتم طبقه‌بندی پایه به‌منظور پیش‌بینی علاقه کاربران، از طبقه‌بند درخت تصمیم که در مقابل داده‌های نویزدار حساس نیست، استفاده شده است و در شکل ۶ نشان داده شده است (هان و کامبر، ۲۰۰۶: ۱۶).

```

(1) create a node N;
(2) if tuples in D are all of the same class, C then
(3)     return N as a leaf node labeled with the class C;
(4) if attribute_list is empty then
(5)     return N as a leaf node labeled with the majority class in D; // majority voting
(6) apply Attribute_selection_method(D, attribute_list) to find the "best" splitting_criterion;
(7) label node N with splitting_criterion;
(8) if splitting_attribute is discrete-valued and
    multiway splits allowed then // not restricted to binary trees
(9)     attribute_list ← attribute_list - splitting_attribute; // remove splitting_attribute
(10) for each outcome j of splitting_criterion
    // partition the tuples and grow subtrees for each partition
(11)     let D_j be the set of data tuples in D satisfying outcome j; // a partition
(12)     if D_j is empty then
(13)         attach a leaf labeled with the majority class in D to node N;
(14)     else attach the node returned by Generate_decision_tree(D_j, attribute_list) to node N;
    endfor;
(15) return N;

```

شکل ۶. الگوریتم طبقه‌بندی درخت تصمیم

چهارمین الگوریتم طبقه‌بندی پایه برای پیش‌بینی علاقه کاربران، طبقه‌بند جنگل تصادفی است و در رتبه‌بندی اهمیت متغیرها در یک مسئله رگرسیون یا طبقه‌بندی کاربرد فراوانی داشته است. این الگوریتم در شکل ۷ نشان داده شده است (زکریا، ۲۰۱۴: ۱۲).

1. Archana & Elangovan
2. Han & Kamber

```

To generate  $c$  classifiers:
for  $i = 1$  to  $c$  do
    Randomly sample the training data  $D$  with replacement to produce  $D_i$ 
    Create a root node,  $N_i$  containing  $D_i$ 
    Call BuildTree( $N_i$ )
end for

BuildTree( $N$ ):
if  $N$  contains instances of only one class then
    return
else
    Randomly select  $\kappa\%$  of the possible splitting features in  $N$ 
    Select the feature  $F$  with the highest information gain to split on
    Create  $f$  child nodes of  $N_1, N_2, \dots, N_f$ , where  $F$  has  $f$  possible values ( $F_1, \dots, F_f$ )
    for  $i = 1$  to  $f$  do
        Set the contents of  $N_i$  to  $D_i$ , where  $D_i$  is all instances in  $N$  that match
             $F_i$ 
        Call BuildTree( $N_i$ )
    end for
end if

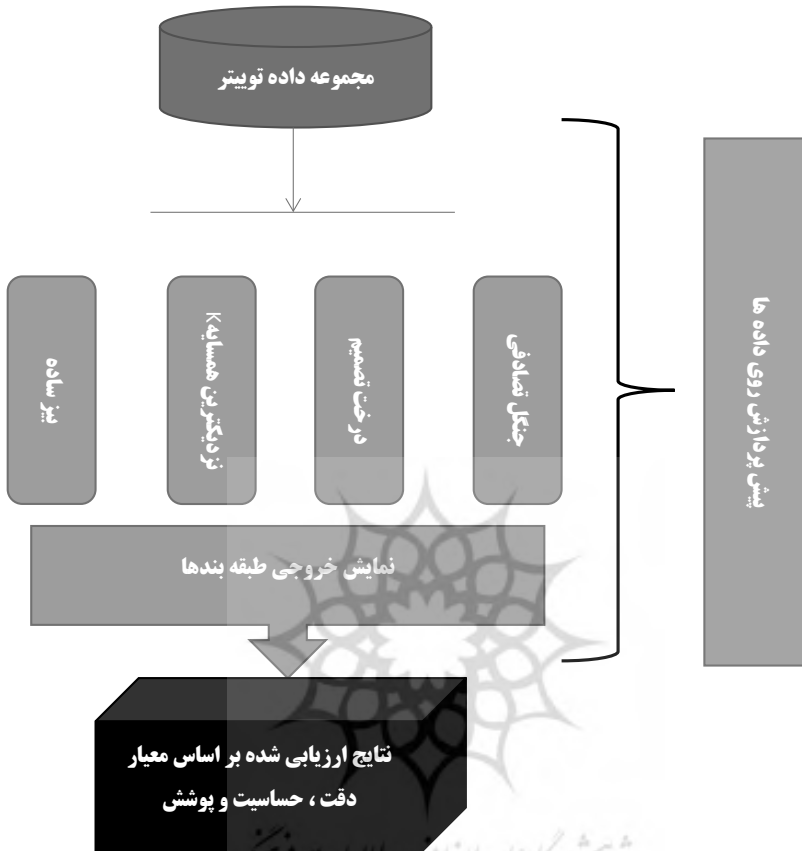
```

شکل ۷. الگوریتم طبقه‌بندی جنگل تصادفی (پینتو و کریا، ۲۰۱۸: ۲۳)

پژوهشگاه علوم انسانی و مطالعات فرهنگی  
پرتال جامع علوم انسانی



## مدل پیشنهادی طبقه‌بند تجمعی



شکل ۸. مدل پیشنهادی

طبقه‌بندهای ترکیبی از ترکیب چندین طبقه‌بند استفاده می‌کنند. در واقع این طبقه‌بندها، هرکدام مدل خود را بر روی داده‌ها می‌سازند و این مدل را ذخیره‌سازی می‌کنند. در نهایت برای طبقه‌بندی نهایی، یک رأی‌گیری در میان این طبقه‌بندها انجام می‌شود و آن طبقه‌ای که بیشترین میزان رأی را بیاورد، طبقه نهایی محسوب می‌شود. به‌طور کلی، این روش تلاش می‌کند تا از خطاهای طبقه‌بندی‌کننده‌های پایه‌ای، با هدف دستیابی به خروجی دقیق‌تر انجام پذیرد (مارتینز، واسکز: ۳۰)، (تروساس، کروسا و همکاران، ۲۰۱۶). روش تجمع خودرانداز<sup>۲</sup> نمونه‌های تصادفی از مجموعه داده‌های آموزشی

1. Martinez-Cámara, Gutiérrez-Vázquez  
2. Bootstrap Aggregation

را تولید می‌کند و الگوریتم‌های یادگیری پایه‌ای را روی هر نمونه اعمال می‌کند. سپس نتایج این طبقه‌بندی‌کننده‌های چندگانه با استفاده از رأی‌گیری میانگین یا اکثریت ترکیب می‌شوند. این روش واریانس رویه‌های ناپایدار نظیر درخت‌های تصمیم را کاهش می‌دهد و در نتیجه، معیار دقت را بهبود می‌بخشد (تروساس، کروسا و همکاران، ۲۰۱۶: ۳۱).

روش تقویت<sup>۱</sup>، یک روش مجموعه‌ای است که مدل‌های مکمل با آموزش هر مدل جدید بر روی مدل‌های قبلی که به اشتباه طبقه‌بندی شده‌اند، تولید می‌شوند. این روش تا زمانی که به محدودیتی در تعداد مدل‌ها یا معیار دقت برسد، تکرار می‌شود. اگرچه روش تقویت در برخی موارد بهتر از کیسه‌کردن عمل کرده است، اما به احتمال زیاد بیش از حد داده‌های آموزشی را تطبیق می‌دهد. به عنوان طبقه‌بندی‌کننده ضعیف، به‌طور معمول از طبقه‌بندی‌کننده‌های مبتنی بر قانون، درخت‌های یک تا دو سطحی، شبکه‌های عصبی بدون لایه‌های پنهان و غیره استفاده می‌شود (تروساس، کروسا و همکاران، ۲۰۱۶: ۳۱).

روش پشته<sup>۲</sup>، الگوریتم‌های یادگیری متفاوتی را بر روی داده‌های آموزشی اجرا می‌کند و سپس، از یک متا طبقه‌بند استفاده می‌کند که پیش‌بینی‌های هر طبقه‌بندی‌کننده را به عنوان ورودی اضافی می‌گیرد. این می‌تواند منجر به کاهش خطای سوء‌گیری یا واریانس بسته به یادگیرنده ترکیبی مورد استفاده قرار گیرد. انباشتن اغلب عملکرد بهتری نسبت به هر یک از مدل‌های آموزش دیده دارد. یک مدل رگرسیون لجستیک تک‌لایه اغلب به عنوان ترکیب‌کننده استفاده می‌شود (تروساس، کروسا و همکاران، ۲۰۱۶: ۳۱).

روش رأی‌گیری<sup>۳</sup> تلاش می‌کند تا مشکل تطابق را حل کند. سپس می‌توان از یک تولید رأی‌گیری ساده برای تخصیص داده‌ها برای اجماع نهایی استفاده کرد. رأی‌گیری می‌تواند برای تعیین عضویت هر دسته برای هر شیء اعمال شود. با این حال، در روش رأی‌گیری داده‌های مختلفی که از مدل‌های متفاوتی تولید شده است را بر اساس رأی‌گیری اکثریت بر مبنای معیارهای مورد نظر از جمله دقت انتخاب می‌کند و به عنوان خروجی نهایی ارائه می‌دهد (قائمی، سلیمان<sup>۲</sup> و همکاران، ۲۰۰۹: ۲۷)، (تروساس، کروسا و همکاران، ۲۰۱۶: ۳۱).

در این مقاله، با استفاده از روش طبقه‌بندی تجمعی و تکنیک رأی‌گیری بر مبنای اکثریت آرا از میان خروجی هر یک از طبقه‌بند‌های به‌کارگرفته شده شامل  $K$ - نزدیک‌ترین همسایه، جنگل تصادفی، درخت تصمیم، و بیز ساده، و بر مبنای معیارهای دقت، صحت و پوشش آن طبقه‌بندی که اکثریت آرا را داشته است، یعنی بهترین خروجی ممکن به منظور پیش‌بینی رفتار کاربران در مجموعه داده توپوتر را ارائه می‌دهد.

1. Boosting
2. Staking
3. Voting
4. Ghaemi & Soliman

## ارزیابی نتایج آزمایش

در این بخش، تنظیمات محیط سخت افزاری و شبیه سازی آزمایش، مجموعه داده های آزمایش شده و معیارهایی که مورد ارزیابی قرار گرفته شده است، توصیف شده اند. در نهایت، نتایج آزمایش مورد ارزیابی قرار گرفته شده است و بین نتایج حاصل از الگوریتم های طبقه بندی، مقایسه صورت پذیرفته است.

### ۱. تنظیمات محیط سخت افزاری و شبیه سازی آزمایش

کلیه شبیه سازی های آزمایش ها در این مقاله، روی سیستمی با پردازنده ۵ هسته ای و سرعت ۴/۴ گیگاهرتز، حافظه اصلی (RAM) ۸ گیگابایت و حافظه فیزیکی ۱ ترابایت، و سیستم عامل ویندوز ۷ نسخه ۶۴ بیتی، اجرا و آزمایش شده است. به علاوه، کلیه شبیه سازی های الگوریتم های طبقه بندی پایه و مدل پیشنهادی طبقه بندی تجمعی، با استفاده از شبیه ساز داده کاوی Rapid Miner نسخه ۹/۰ توسعه و پیاده سازی شده است.

### ۲. مجموعه داده آزمایش شده

کلیه الگوریتم های طبقه بندی پایه ای و مدل پیشنهادی طبقه بندی تجمعی شبیه سازی شده، بر روی مجموعه داده Tweets.csv که مربوط به کلان داده رسانه اجتماعی توئیتر است، اجرا و آزمایش شده است. مجموعه داده Tweets دارای تعداد ۵۲۵۴۳ نمونه و همچنین ۱۰ ویژگی محتوای موجود شامل محتوا، کشور، تاریخ، زمان، شناسه، زبان، تعداد مورد علاقه ها، تعداد اشتراک گذاشته شده و نوع محتوا (رشته ای یا عددی) است. به علاوه، مجموعه داده Tweets دارای ۱۰ کلاس مختلف است.

### ۳. معیارهای ارزیابی شده

در این مقاله، صحت، دقت و پوشش معیارهایی هستند که الگوریتم های طبقه بندی پایه ای و مدل پیشنهادی طبقه بندی تجمعی، بر اساس این سه معیار مورد ارزیابی قرار گرفته اند. مطابق رابطه (۶)، معیار صحت برابر است با حاصل تقسیم تعداد نمونه هایی که به طور صحیح پیش بینی شده اند، بر تعداد کل پیش بینی های انجام شده (هان، کامبر، ۲۰۰۶: ۱۶).

$$Accuracy = \frac{TP+TN}{TP+TN+FP+FN} \quad (6)$$

که در آن، TP تعداد نمونه های مثبت صحیح پیش بینی شده، TN تعداد نمونه های منفی صحیح پیش بینی شده، FP تعداد نمونه های مثبت غلط پیش بینی شده و FN تعداد نمونه های منفی غلط پیش بینی شده هستند.

در معیار دقت، هرچه میزان پیش‌بینی غلط از میزان پیش‌بینی صحیح بیشتر باشد، مقدار دقت کمتر است. دقت برابر است با حاصل تقسیم تعداد نمونه‌های مثبت صحیح پیش‌بینی شده بر تعداد نمونه‌های مثبت صحیح پیش‌بینی شده و تعداد نمونه‌های منفی غلط پیش‌بینی شده که در رابطه (۷) نشان داده شده است (هان، کامبر، ۲۰۰۶: ۱۶).

$$Precision = \frac{TP}{TP+FP} \quad (۷)$$

مطابق با رابطه (۸) (هان، کامبر، ۲۰۰۶: ۱۶)، معیار پوشش برابر است با تعداد نمونه‌های مثبت صحیح پیش‌بینی شده، تقسیم بر مجموع تعداد نمونه‌های مثبت صحیح پیش‌بینی شده و تعداد نمونه‌های منفی غلط پیش‌بینی شده. هرچه میزان تعداد نمونه‌هایی که غلط پیش‌بینی شده است، بیشتر از تعداد نمونه‌های صحیح پیش‌بینی شده باشد، میزان پوشش کمتر است.

$$Recall = \frac{TP}{TP+FN} \quad (۸)$$

#### ۴. ارزیابی نتایج آزمایش

در این بخش، مطابق با آزمایش‌های انجام‌شده و اجرای چهار الگوریتم طبقه‌بندی پایه‌ای و قدرتمند شامل K- نزدیک‌ترین همسایه، بیزین ساده، جنگل تصادفی و درخت تصمیم و همچنین روش پیشنهادی طبقه‌بندی تجمعی مبتنی بر رأی‌گیری، روی مجموعه داده توییت‌ها به نام Tweets، به ارزیابی این الگوریتم‌ها بر اساس معیارهای دقت، صحت، و پوشش پرداخته شده است. در روش پیشنهادی طبقه‌بندی تجمعی، از روش رأی‌گیری برای به‌دست آوردن خروجی نهایی طبقه‌بندی با استفاده از خروجی هر یک از طبقه‌بندهای پایه‌ای، استفاده شده است و انتظار می‌رود که معیارهای صحت، دقت، و پوشش بهبود حاصل شود.

در جدول ۱ نتایج آزمایش اجرای چهار الگوریتم طبقه‌بندی پایه‌ای شامل K- نزدیک‌ترین همسایه، بیزین ساده، جنگل تصادفی و درخت تصمیم و همچنین، روش پیشنهادی طبقه‌بندی تجمعی مبتنی بر رأی‌گیری و بر مبنای حداکثر آرا، روی مجموعه داده توییت‌ها Tweets و بر اساس معیار دقت نشان داده شده است. طبق نتایج به‌دست آمده از خروجی هر یک از طبقه‌بندهای پایه‌ای بر اساس معیار دقت روی مجموعه داده Tweets، بدترین عملکرد متعلق به طبقه‌بند بیزین ساده با مقدار ۳۷/۹۲ درصد و بهترین عملکرد متعلق به طبقه‌بند جنگل تصادفی با مقدار ۹۲/۹۷ درصد بوده است. با استفاده از روش طبقه‌بند تجمعی پیشنهادی مبتنی بر روش رأی‌گیری، بهترین خروجی از میان خروجی‌های کلیه طبقه‌بندها انتخاب می‌شود. همان‌طور که مشاهده می‌شود، بهترین عملکرد بر اساس معیار دقت با مقدار ۹۲/۹۷، متعلق به خروجی

طبقه‌بند پایه جنگل تصادفی و همچنین، روش طبقه‌بند تجمعی پیشنهادی مبتنی بر روش رأی‌گیری بوده است.

### جدول ۱. نتایج الگوریتم‌های طبقه‌بندی پایه‌ای و تجمعی براساس معیار دقت

K- نزدیک‌ترین همسایه	بیزین ساده	جنگل تصادفی	درخت تصمیم	طبقه‌بند تجمعی
۹۲/۵۳ درصد	۳۷/۹۲ درصد	۹۲/۹۷ درصد	۹۲/۹۶ درصد	۹۲/۹۷ درصد

در جدول ۲ نتایج آزمایش اجرای چهار الگوریتم طبقه‌بندی پایه‌ای شامل K- نزدیک‌ترین همسایه، بیزین ساده، جنگل تصادفی و درخت تصمیم و همچنین، روش پیشنهادی طبقه‌بندی تجمعی مبتنی بر رأی‌گیری و بر مبنای حداکثر آرا، روی مجموعه داده توئیتر Tweets و بر اساس معیار صحت نشان داده شده است. طبق نتایج به دست آمده از خروجی هر یک از طبقه‌بندهای پایه‌ای بر اساس معیار صحت روی مجموعه داده Tweets، بدترین عملکرد متعلق به طبقه‌بند بیزین ساده با مقدار ۱۵/۱۴ درصد، و بهترین عملکرد متعلق به طبقه‌بندهای پایه‌ای جنگل تصادفی و درخت تصمیم با مقدار ۶/۰۶ درصد بوده است. با استفاده از روش طبقه‌بند تجمعی پیشنهادی مبتنی بر روش رأی‌گیری و بر مبنای حداکثر آرا، بهترین خروجی از میان خروجی‌های کلیه طبقه‌بندها انتخاب می‌شود. همان‌طور که مشاهده می‌شود، از آنجایی که خروجی طبقه‌بندهای جنگل تصادفی و درخت تصمیم با یکدیگر برابر بودند، لذا بعد از رأی‌گیری از میان طبقه‌بندها، خروجی‌های طبقه‌بندهای پایه‌ای جنگل تصادفی و درخت تصمیم را به عنوان خروجی بهینه و نهایی در نظر گرفته است. بنابراین، عملکرد روش طبقه‌بند تجمعی پیشنهادی مبتنی بر روش رأی‌گیری بر اساس معیار صحت، مقدار ۶/۰۶ درصد بوده است.

### جدول ۲. نتایج الگوریتم‌های طبقه‌بندی پایه‌ای و تجمعی براساس معیار صحت

K- نزدیک‌ترین همسایه	بیزین ساده	جنگل تصادفی	درخت تصمیم	طبقه‌بند تجمعی
۶/۷۳ درصد	۱۵/۱۴ درصد	۶/۰۶ درصد	۶/۰۶ درصد	۶/۰۶ درصد

در جدول ۳ نتایج آزمایش اجرای چهار الگوریتم طبقه‌بندی پایه‌ای شامل K- نزدیک‌ترین همسایه، بیزین ساده، جنگل تصادفی و درخت تصمیم و همچنین، روش پیشنهادی طبقه‌بندی تجمعی مبتنی بر رأی‌گیری و بر مبنای حداکثر آرا، روی مجموعه داده توئیتر Tweets و بر اساس معیار پوشش نشان داده شده است. طبق نتایج به دست آمده از خروجی هر یک از طبقه‌بندهای پایه‌ای بر اساس معیار پوشش روی مجموعه داده Tweets، بدترین عملکرد متعلق به طبقه‌بند پایه‌ای بیزین ساده با مقدار ۷/۴۲ درصد و بهترین عملکرد متعلق به طبقه‌بندهای پایه‌ای جنگل

تصادفی و درخت تصمیم با مقدار ۵/۷۴ درصد بوده است. با استفاده از روش طبقه‌بند تجمعی پیشنهادی مبتنی بر روش رأی‌گیری، بهترین خروجی از میان خروجی‌های کلیه طبقه‌بندها انتخاب می‌شود. همان‌طور که مشاهده می‌شود، از آنجایی که خروجی طبقه‌بندهای جنگل تصادفی و درخت تصمیم با یکدیگر برابر بودند، لذا بعد از رأی‌گیری از میان طبقه‌بندها، خروجی‌های طبقه‌بندهای پایه‌ای جنگل تصادفی و درخت تصمیم را به‌عنوان خروجی بهینه و نهایی در نظر گرفته است. بنابراین، عملکرد روش طبقه‌بند تجمعی پیشنهادی مبتنی بر روش رأی‌گیری بر اساس معیار صحت، مقدار ۵/۷۴ درصد بوده است.

### جدول ۳. نتایج الگوریتم‌های طبقه‌بندی پایه‌ای و تجمعی بر اساس معیار پوشش

K- نزدیک‌ترین همسایه	بیزین ساده	جنگل تصادفی	درخت تصمیم	طبقه‌بند تجمعی
۶/۷۳ درصد	۷/۴۲ درصد	۵/۷۴ درصد	۵/۷۴ درصد	۵/۷۴ درصد

## نتیجه‌گیری و کارهای آینده

امروزه استقبال شدید از رسانه‌ها و شبکه‌های اجتماعی و همچنین افزایش حجم ذخیره‌سازی و مبادلات پیام‌ها، باعث ایجاد کلان‌داده‌ها شده است. از کلان‌داده‌های شبکه‌های اجتماعی برای پیش‌بینی علاقه کاربران، شناسایی رفتار آن‌ها و پیش‌بینی موارد کاربردی دیگر استفاده می‌شود. اما سامان‌دهی این کلان‌داده‌ها به‌منظور استفاده بهینه، خود یک چالش بزرگ است. یکی از تکنیک‌های سامان‌دهی این داده‌ها برای آماده‌سازی داده‌ها به‌منظور پیش‌بینی علاقه یا رفتار کاربران و غیره، تکنیک طبقه‌بندی است. در این مقاله، با استفاده از الگوریتم‌های پایه‌ای طبقه‌بندی شامل K- نزدیک‌ترین همسایه، بیزین ساده، جنگل تصادفی و درخت تصمیم و همچنین، روش پیشنهادی طبقه‌بندی تجمعی مبتنی بر روش رأی‌گیری، روی مجموعه داده Tweets و بر اساس معیارهای دقت، صحت و پوشش نشان داده شد که روش پیشنهادی طبقه‌بندی تجمعی مبتنی بر روش رأی‌گیری نسبت به چهار الگوریتم دیگر طبقه‌بندی پایه‌ای، بر اساس هر سه معیار دقت، صحت و پوشش، دارای نتایج مطلوب‌تری بوده است که نشان از کارآمدی عملکرد روش پیشنهادی طبقه‌بندی تجمعی مبتنی بر روش رأی‌گیری است.

برای کارهای آتی، می‌توان از الگوریتم‌های دیگر طبقه‌بندی نظیر ماشین بردار پشتیبان استفاده کرد. علاوه بر این، روش‌های خوشه‌بندی و تلفیق کاربرد روش‌های خوشه‌بندی با طبقه‌بندی به‌منظور تجزیه و تحلیل داده‌ها و پیش‌بینی علاقه یا رفتار کاربران در شبکه‌های اجتماعی می‌تواند از کارهای آینده در این مقاله به‌شمار رود.

## منابع

- Abbas, A. Kh., Ali Khalil, S., Harith, A. H., Qasim, M. H., & Saba Alaa, A. (2020). Twitter sentiment analysis using an ensemble majority vote classifier. *Journal of Southwest Jiaotong University*, 55(1).
- Adewole, K. S., Tao, H., Wanqing, W., Houbing, S., & Arun Kumar, S. (2020). Twitter spam account detection based on clustering and classification methods." *The Journal of Supercomputing*, 76(7), 4802-4837.
- Al-Hashedi, A., Belal A. F., Abdulqader, M., Mohsen, Y. A., Hasan Ali, G., A.K, Wedad, A., & Naseebah, M. (2022). Ensemble classifiers for Arabic sentiment analysis of social network (Twitter data) towards COVID-19-related conspiracy theories. *Applied Computational Intelligence and Soft Computing* 2022.
- Archana, S., & Elangovan, K. (2014). Survey of classification techniques in data mining." *International Journal of Computer Science and Mobile Applications*, 2(2), 65-71.
- Bayhaqy, A., Sfenrianto, S., Kaman, N., & Emil R. K. (2018). Sentiment analysis about E-commerce from tweets using decision tree, K-nearest neighbor, and naive bayes. In 2018 *international conference on orange technologies (ICOT)*, 1-6.
- Charbuty, B., & Abdulazeez, A. (2021). Classification based on decision tree algorithm for machine learning. *Journal of Applied Science and Technology Trends* 2, 1, 20-28.
- Ghaemi, R., Md Nasir, S., Hamidah, I., & Norwati M. (2009). "A survey: clustering ensembles techniques. *International Journal of Computer and Information Engineering* 3, 2, 365-374.
- Ghani, Norjihani Abdul, Suraya Hamid, Ibrahim Abaker Targio Hashem, and Ejaz Ahmed. "Social media big data analytics: A survey." *Computers in Human Behavior* 101 (2019): 417-428.
- Han, J., & Micheline K. (2006). Data mining: concepts and techniques, 2nd. *University of Illinois at Urbana Champaign: Morgan Kaufmann*.
- Harjule, P., Astha, G., Harshita, S., & Priya T. (2020). Text classification on Twitter data. In 2020 *3rd International Conference on Emerging Technologies in Computer Engineering: Machine Learning and Internet of Things (ICETCE)*, 160-164.
- Kesavaraj, G., & Sreekumar S. (2013). A study on classification techniques in data mining. In 2013 *fourth international conference on computing, communications and networking technologies (ICCCNT)*, 1-7.
- Kesavaraj, G., & Sreekumar S. (2013). A study on classification techniques in data mining. In 2013 *fourth international conference on computing, communications and networking technologies (ICCCNT)*, 1-7.
- Martinez-Cámara, E., Gutiérrez-Vázquez, Y. Javi, F, Montejo-Ráez, A., & Muñoz-Guillena, R. (2015). Ensemble classifier for twitter sentiment analysis. *NLP Applications: completing the puzzle*, 1-12.
- Metsis, V., Ion, A., & Georgios, P. (2006). Spam filtering with naive bayes-which naive bayes? In CEAS, 17, 28-69.
- Muqorobin, M., Siti, R., Isnawati, M., & Nendy, A. R. R. (2020). Classification of Community Complaints Against Public Services on Twitter. *International Journal of Computer and Information System (IJCIS)*, 1(1), 7-10.
- Murthy, D. (2017). The ontology of tweets: Mixed methods approaches to the study of Twitter. *The SAGE handbook of social media research methods*, 559-572.
- Neogi, A. S., Kirti Anilkumar, G., Ram Krishn, M., & Yogesh, K. D. (2021). Sentiment analysis and classification of Indian farmers' protest using twitter data. *International Journal of Information Management Data Insights*, 1(2).
- Oussous, A., Benjelloun, F., Ait Lahcen, R., & Belfkih, S. (2018). Big Data technologies: A survey. *King Saud University-Computer and Information Sciences* 30, 4, 431-448.
- Papakyriakou, D., & Ioannis S. (2022). Data Mining Methods: A Review. *International Journal of Computer Applications* 975: 8887.
- Patel, R., & Kalpdrum P. (2020). Sentiment analysis on twitter data of world cup soccer tournament using machine learning. *IoT* 1, 2 (14).
- Patil, A. S., & Pawar, B. V. (2012) Automated classification of web sites using Naive Bayesian algorithm.

- In *Proceedings of the international multiconference of engineers and computer scientists*, 1, 519-523.
- Phan, H. T., Ngoc Th. N., Van Cuong T., & Dosam H. (2021). An approach for a decision-making support system based on measuring the user satisfaction level on twitter. *Information Sciences 561*, 243-273.
- Phyu, Th.N. (2009). Survey of classification techniques in data mining. In *Proceedings of the international multiconference of engineers and computer scientists*, 1(5), Citeseer.
- Pinto, L. R. M., & Henrique Andrade Correia, L. (2018). Analysis of machine learning algorithms for spectrum decision in cognitive radios. In *2018 15th International symposium on wireless communication systems (ISWCS)*, 1-6.
- Sen, P. Ch., Mahimarnab, H., & Mitadru, Gh. (2020). "Supervised classification algorithms in machine learning: A survey and review." In *Emerging technology in modelling and graphics*, 99-111.
- Shaik, A. B., & Sujatha S. (2019). A brief survey on random forest ensembles in classification model. In *International Conference on Innovative Computing and Communications*, 253-260.
- Shamrat, F. M. J. M., Sovon Chakraborty, M. M., Imran, J., Naeem, M., Md Masum, B, Protiva, D., & Rahman, O. M. (2021). Sentiment analysis on twitter tweets about COVID-19 vaccines using NLP and supervised KNN classification algorithm. *Indonesian Journal of Electrical Engineering and Computer Science*, 23(1), 463-470.
- Sri, P.A., & Anusha, M. (2016). Big data-survey. *Indonesian Journal of Electrical Engineering and Informatics (IJEI)* 4 (1), 74-80.
- Troussas, Ch., Akrivi, K., & Virvou, M. (2016). Evaluation of ensemble-based sentiment classifiers for Twitter data. In *2016 7th international conference on information, intelligence, systems & applications (IISA)*, 1-6.
- Vashisth, P., & Kevin, M. (2020). Gender classification using twitter text data. In *2020 31st Irish Signals and Systems Conference (ISSC)*, 1-6.
- Zakariah, M. (2014). Classification of large datasets using Random Forest Algorithm in various applications: Survey. *International Journal of Engineering and Innovative Technology, (IJJEIT)* 4(3).

© Authors, Published by Bureau of Media Studies and Planning. This is an open-access paper distributed under the CC BY (license <http://creativecommons.org/licenses/by/4.0/>).



پژوهشگاه علوم انسانی و مطالعات فرهنگی  
رتال جامع علوم انسانی