

Strategic data analysis model in social networks

*Rasoul Lotfi*¹
*Kamran Feizi*²
*Sayyed Mohammad Khatami Firoozabadi*³
*Mohammadreza Taghva*⁴

Type of article: Research article extracted from doctoral dissertation

Received:2021/10/18 Accepted: 2021/12/6

NAJA Strategic Studies Quarterly/Vol.7/NO.1/(serial23)Spring 2022*83-104



DOR: 20.1001.1.25381946./ssj.2022.98600

Abstract

Today, the Internet, with the creation and development of global communication platforms, has been able to play a fundamental and strategic role in governing cyberspace, and one of these spaces that has attracted a large audience is social networks. These networks have created an environment where users can overcome the limitations and barriers of their real life space and take full advantage of the opportunity provided for their interaction, communication and message transmission; Thus, the ability to analyze social networks on a large scale is a specific technical capability that each organization must implement in accordance with its mission. In the field of police actions, in order for the police management structure to be able to assess the impact of its actions on public opinion, it must be made public; therefore, in this study, in order to show how to analyze the strategic data of social networks, the pattern of data collection and processing has been drawn and in order to turn the mental abstract pattern into a tangible practical pattern, using content analysis and text analysis to examine the views of Iranian users on Twitter. Has been. In total, after deleting additional tweets, 39,000 samples were used for text analysis. In order to analyze the tweets, the samples were done based on the users' point of view. The duration of information analysis was performed in a period of 3 years (1397-1400) and during the description, data mining structure and text analysis were mentioned to understand the differences and similarities of analytical models.

Keywords: social network, twitter, text mining, feta police, strategic analysis

1. PhD Student in Information Technology Management, Faculty of Management, Allameh Tabatabaie University, Tehran, Iran (Corresponding Author), iamlotfi@yahoo.com

2. Professor, Department of Industrial Management, Faculty of Management, Allameh Tabatabaie University, Tehran, Iran

3. Professor, Department of Industrial Management, Faculty of Management, Allameh Tabatabaie University, Tehran, Iran

4. Professor, Department of Industrial Management, Faculty of Management, Allameh Tabatabaie University, Tehran, Iran

الگوی تحلیل داده‌های راهبردی در شبکه‌های اجتماعی

رسول لطفی^۱

کامران فیضی^۲

سیده محمد خاتمی فیروز آبادی^۳

محمد رضا تقوا^۴

نوع مقاله: مقاله پژوهشی مستخرج از رساله دکتری

تاریخ دریافت: ۱۴۰۰/۷/۲۶ تاریخ پذیرش نهایی: ۱۴۰۰/۹/۱۵

فصلنامه مطالعات راهبردی ناجا/سال هفتم/شماره ۱(پیاپی ۲۳)-بهار ۱۴۰۱* ۸۳-۱۰۴



DOR: 20.1001.1.25381946./ssj.2022.98600

چکیده

امروزه، اینترنت با ایجاد و توسعه پلتفرم‌های ارتباطی جهانی توانسته است نقش اساسی و راهبردی را در حکمرانی فضای مجازی بازی کند و یکی از این فضاها که مخاطبان زیادی را جذب نموده است، شبکه‌های اجتماعی است. این شبکه‌ها، فضایی ایجاد کرده‌اند که کاربران می‌توانند محدودیت‌ها و موانع فضای واقعی زندگی خود را جبران کرده و از فرصتی که برای تعاملات، ارتباطات و انتقال پیام آنها فراهم شده است، به خوبی بهره ببرند؛ از این رو، توانایی تحلیل شبکه‌های اجتماعی در مقیاس وسیع، یک توانمندی خاص فنی است که هر سازمان به تناسب مأموریت خود باید آن را پیاده‌سازی کند. در حوزه اقدامات پلیس نیز برای آنکه ساختار مدیریتی پلیس بتواند تأثیر اقدامات خود را بر افکار عمومی بسنجد، باید از نظر عموم جامعه آگاه شود؛ بنابراین، در این پژوهش، به منظور نشان دادن چگونگی تحلیل داده‌های راهبردی شبکه‌های اجتماعی، الگوی جمع‌آوری و پردازش داده‌ها ترسیم شده است و به منظور تبدیل الگوی انتزاعی ذهنی به الگوی کاربردی ملموس، با استفاده از تحلیل محتوا و متن کاوی به بررسی دیدگاه کاربران ایرانی در توییت‌ها نسبت به پلیس فتا پرداخته شده است. در مجموع، پس از حذف توییت‌های اضافی، تعداد ۳۹۰۰۰ نمونه برای متن کاوی استفاده شدند. به منظور تحلیل توییت‌ها، نمونه‌ها بر اساس دیدگاه کاربران انجام پذیرفت. مدت زمان تحلیل اطلاعات در یک بازه زمانی ۳ ساله (۱۳۹۷ تا ۱۴۰۰) انجام شد و در طی تشریح مطالب، به ساختار داده کاوی و تحلیل متن اشاره شد تا تمایز و تشابه الگوهای تحلیلی درک شود.

واژگان کلیدی: شبکه اجتماعی، توییت، متن کاوی، پلیس فتا، تحلیل راهبردی

۱. دانشجوی دکتری مدیریت فناوری اطلاعات، دانشکده مدیریت، دانشگاه علامه طباطبایی (ره)، تهران، ایران (نویسنده مسئول)، iamlotfi@yahoo.com

۲. استاد گروه مدیریت صنعتی، دانشکده مدیریت، دانشگاه علامه طباطبایی (ره)، تهران، ایران

۳. استاد گروه مدیریت صنعتی، دانشکده مدیریت، دانشگاه علامه طباطبایی (ره)، تهران، ایران

۴. استاد گروه مدیریت صنعتی، دانشکده مدیریت، دانشگاه علامه طباطبایی (ره)، تهران، ایران

مقدمه

امروزه، اینترنت با ایجاد و توسعه پلتفرم‌های ارتباطی جهانی توانسته‌است نقش اساسی و راهبردی را در حکمرانی فضای مجازی بازی کند؛ بر این اساس، روشن است که نمی‌توان نقش فناوری اطلاعات را در تمامی عرصه‌های مرتبط بشری نادیده گرفت. یکی از این فضاها که مخاطبان زیادی را جذب نموده‌است، شبکه‌های اجتماعی است. شبکه‌های اجتماعی مجازی اغلب سرویس‌های مبتنی بر وب هستند و سرویس‌های آنلاین شامل پلتفرم‌ها یا سایت‌هایی هستند که افراد از طریق آنها نظرها و علاقه‌مندی‌های خود را بیان کرده و با دیگران به اشتراک می‌گذارند. آندریس کاپلان و مایکل هایلین^۱ رسانه اجتماعی را به‌مثابه یک دسته ابزار مبتنی بر اینترنت تعریف می‌کنند که بر بنیاد ایدئولوژی یک و فناوری وب^۲ استوار هستند و به کاربر امکان تولید محتوا و مبادله آن را می‌دهند (کاپلان و هایلین، ۲۰۱۰). این شبکه‌ها، فضایی ایجاد کرده‌اند که کاربران می‌توانند محدودیت‌ها و موانع فضای واقعی زندگی خود را جبران کرده و از فرصتی که برای تعاملات، ارتباطات و انتقال پیام آنها فراهم شده‌است، به‌خوبی بهره ببرند (اومالی، ۲۰۱۴). از جمله کارکردهای شبکه‌های اجتماعی، نقش بازخوردی آنهاست؛ به‌گونه‌ای که افراد، سازمان‌ها و دولت‌ها می‌توانند بازخورد اقدامات خود را در این شبکه‌ها جستجو کنند و از آن برای جمع‌آوری پردازش و تحلیل راهبردی داده‌ها استفاده کنند؛ چراکه سرویس‌های مبتنی بر شبکه که بتوانند داده‌های کلان^۳ از نظرات آرا و کنش و واکنش‌های کاربران را دریافت نمایند، در دسترس است و سرویس‌های تحلیل داده که بتواند با استفاده از نرم‌افزارهای تحلیلی و هوش مصنوعی^۴ آنها را تحلیل کرده و موجب خروجی تصمیم‌ساز شوند، قابل اجرا و پیاده‌سازی است. در حال حاضر، شبکه‌های اجتماعی بخش جدایی‌ناپذیر جامعه انسانی شده‌است و مردم در این شبکه‌ها عقاید خود را در خصوص موضوعات مختلف بیان می‌کنند؛ به‌عنوان مثال، روز پنجشنبه ۴ اوت ۲۰۱۱، مارک دوگان توسط یک افسر پلیس در تاتنهام کشته شد، صبح روز ششم، محتوای رسانه‌های اجتماعی خصومت فزاینده‌ای را نشان می‌دادند که تهدیدهای صریح بر ضد پلیس بود. از روز هفتم، اطلاعات رسانه‌های اجتماعی نشان‌دهنده گسترش احتمالی اعتراضات بود، اختلال در سایر نقاط لندن و سپس انگلستان را درگیر کرد. در این رخداد، داده‌های زیادی منتشر شد که می‌توانست در خصوص اطلاعات غلط در حال انتشار و افراد پشت آن با هویت مجعول به پلیس کمک کند. بعدها پلیس اذعان کرد که

1. Kaplan & Haenlein

2. Umali

3. Big Data

4. AI(Artificial Intelligence)

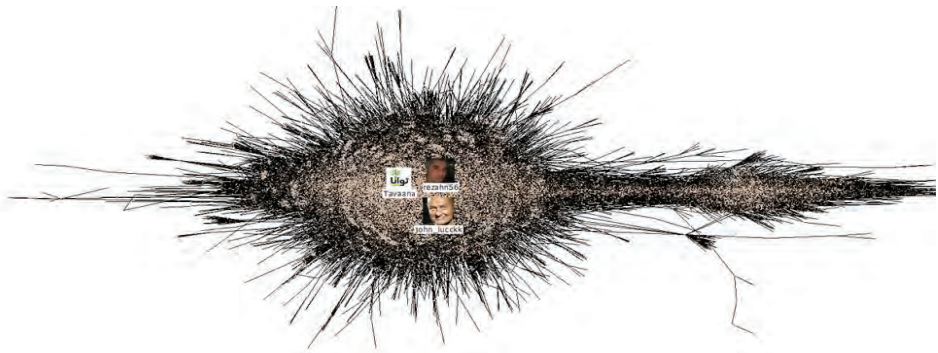
به ابزار و دانش کافی برای جمع‌آوری اطلاعات از شبکه‌های اجتماعی مسلط نبوده‌است (اومانند و همکاران، ۲۰۱۲). بر این اساس، در حوزه اقدامات پلیس، برای آنکه ساختار مدیریتی پلیس بتواند تأثیر اقداماتش را بر افکار عمومی بسنجد، باید از نظر عموم جامعه آگاه شود. توانایی تحلیل داده‌های شبکه‌های اجتماعی در مقیاس وسیع یک توانمندی است که سازمان‌های بزرگ به تناسب مأموریت خود باید از آن بهره‌مند باشند؛ بنابراین، این پژوهش در تلاش است تا الگویی برای تحلیل داده‌های راهبردی در شبکه‌های اجتماعی ارائه نماید.

ادبیات نظری پژوهش

شبکه‌های عصبی مصنوعی^۱

شبکه‌های عصبی، نسل جدید روش‌های داده‌کاوی به‌شمار می‌آیند که در دو دهه اخیر، توسعه زیادی یافته‌اند. از این روش‌ها هم می‌توان برای کشف و استخراج دانش از پایگاه داده‌ها و هم برای ایجاد الگوهای پیش‌بینی استفاده کرد. شبکه‌های عصبی به دلیل توسعه‌های اخیرشان، به‌عنوان یک ابزار داده‌کاوی بسیار متداول شده‌اند. این شبکه‌ها ابزارهایی هستند که در موارد پیش‌بینی، رده‌بندی، خوشه‌بندی و تخمین کاربرد دارند. همچنین، استفاده‌های صنعتی این ابزار نیز بسیار فراوان است که از آنها می‌توان به تشخیص سلسله‌های زمانی، تشخیص شرایط اجتماعی و تحولات آن، تشخیص خوشه‌های افراد باارزش (هر ارزشی مثبت یا منفی) و شناسایی تخلف‌هایی که در بیمه صورت می‌گیرد، اشاره نمود. تولید گراف‌هایی با حجم، جهت، گستردگی و یا پراکندگی و همچنین، اراده اجتماعی در یک موضوع و یا تشریح شبکه‌های جرایم و تروریسم یکی دیگر از کاربردهای شبکه‌های عصبی است. گراف‌هایی که در نتیجه عملکرد شبکه‌های عصبی تولید می‌شوند را می‌توان یکی از بهترین نمونه‌های آن دانست. از نرم‌افزارهایی که با استفاده از شبکه‌های عصبی از داده‌های ارتباطی، گراف‌های واضح و قابل بررسی را تولید می‌کنند، می‌توان نرم‌افزار متن باز گفی^۲ و نرم‌افزار دولتی I۲^۳ با دسترسی محدود را نام برد.

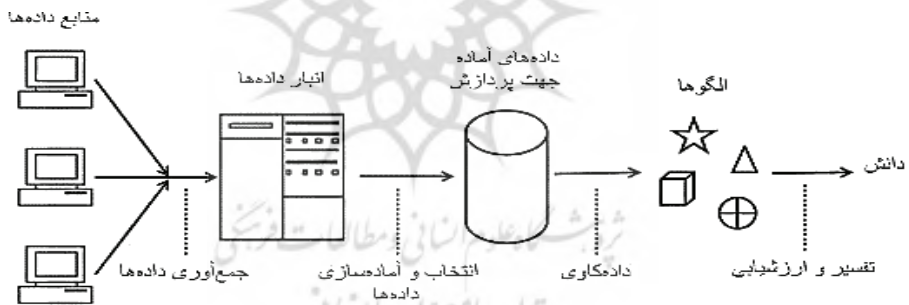
1. Artificial Neural Networks
 2. Gephi Network Analysis
 3. IBM Intelligence Analysis



شکل ۱. خروجی تحلیل روابط شبکه‌های تروریستی در توییتر با استفاده از نرم‌افزار گفی

تمایز روش‌های آماری و داده‌کاوی

به‌عنوان یک قانون کلی، فرض روش‌های آماری بر این اساس است که توزیع داده‌ها مشخص است و در بیشتر موارد، فرض بر این است که توزیع نرمال است و درنهایت، درستی یا نادرستی نتایج نهایی، به درست بودن فرض اولیه وابسته است (شکل ۲).



شکل ۲. فرآیند داده‌کاوی

در مقابل روش‌های یادگیری، یادگیری ماشین از هیچ فرضی در مورد داده‌ها استفاده نمی‌کند و همین باعث تفاوت‌هایی بین این دو روش می‌شود (مانیلا^۱ و همکاران، ۲۰۱۱). در جدول شماره ۱ تفاوت‌های روش آنالیز آماری و داده‌کاوی مشاهده می‌شود.

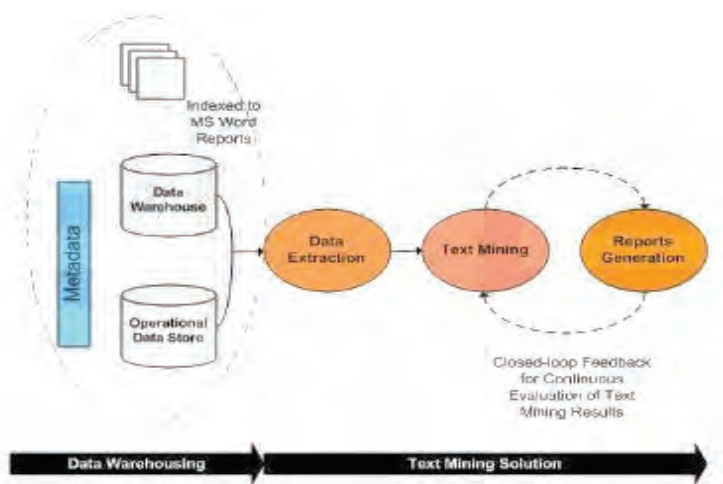
روش	آنالیز آماری	داده کاوی
فرضیه	آمار شناسان همیشه یا یک فرضیه شروع به کار می کنند.	به فرضیه احتیاجی ندارد.
نوع داده ها	آنها از داده های عددی استفاده می کنند.	ابزارهای داده کاوی از انواع مختلف داده، نه تنها عددی می توانند استفاده کنند.
ایجاد روابط	آمارشناسان باید رابطه هایی را ایجاد کنند که به فرضیه آنها مربوط است.	الگوریتمهای داده کاوی به طور اتوماتیک روابط را ایجاد می کنند.
صحت داده ها	آنها می توانند داده های ناپجا و نادرست را در طول آنالیز مشخص کنند.	داده کاوی به داده های صحیح و درست نیاز دارد.
قابلیت تفسیر	آنها می توانند نتایج کار خود را تفسیر و برای مدیران بیان کنند.	نتایج داده کاوی نسبتا پیچیده می باشد و نیاز به متخصصانی جهت بیان آنها به مدیران دارد.

جدول ۱. تفاوت های روش آنالیز آماری و داده کاوی (حسین زاده، ۱۳۹۲)

تحلیل محتوا

در خصوص تحلیل محتوا تعریف های مختلفی وجود دارد؛ اما یکی از تعریف های جامع، تعریف برنارد برلسون^۱ است که پس از گذشت سال ها، هنوز بسیاری از پژوهشگران به آن استناد می کنند. برلسون، تحلیل محتوا را یک شیوه تحقیقی تعریف کرده است که برای تشریح عینی، منظم و کمی محتوای آشکار پیام های ارتباطی به کار می رود (ضغیمی و همکاران، ۱۳۹۷). به عبارت دیگر، برنارد برلسون تحلیل محتوا را روش مطالعه و تجزیه و تحلیل ارتباط به شیوه های نظام مند، عینی، کمی و آشکار برای اندازه گیری متغیرها توصیف می کند؛ اما در کل، تحلیل محتوا به فرآیند تحلیل متن برای استخراج یا کشف اطلاعات و واقعیت های معتبر، جدید و از پیش ناشناخته، پنهان، مفید و قابل درک از داده های ساخت نیافته و نیمه ساخت یافته به صورت خود کار گفته می شود.

در شکل ۳، معماری مفهومی فرآیند تحلیل محتوا شامل مرحله انبار داده، استخراج ویژگی و الگوسازی توصیف شده است. در بخش انتهایی (تولید گزارش که توسط کاربر سیستم اعمال می شود)، می توان ایرادها را بهبود بخشید یا نظرات متخصص حوزه را اعمال کرد.



شکل ۳. معماری مفهومی تحلیل محتوا

شباهت‌ها و تفاوت‌های داده‌کاوی و تحلیل محتوا

متن‌کاوی برخی از رهیافت‌های خود را از پژوهش‌های داده‌کاوی به دست آورده‌است؛ بنابراین، شباهت‌های زیادی بین معماری داده‌کاوی و تحلیل محتوا وجود دارد. در هر دو روش، رویه‌هایی مبنی بر پیش‌پردازش، الگوریتم‌های کشف الگو و مؤلفه‌های لایه‌بازنمایی وجود دارد (مانند ابزارهای بازنمایی بصری برای بهبود نتیجه‌ها)؛ همچنین، تحلیل محتوا بسیاری از انواع الگوها را در عملیات اکتشاف دانش هسته‌ای که ابتدا در تحقیقات داده‌کاوی معرفی شده‌اند، اتخاذ می‌کند. در تفاوت داده‌کاوی و تحلیل محتوا نیز می‌توان مطرح کرد که داده‌کاوی شامل روش‌هایی برای اکتشافی از انواع داده‌ها مانند داده‌های عددی، چندرسانه‌ای و... است (شیخ و شامبیاتی، ۱۳۹۴).

این امکان وجود دارد که تمام انواع داده با استفاده از الگوریتم‌های عمومی داده‌کاوی پردازش شوند و تحلیل محتوا به‌عنوان زیرمجموعه‌ای از داده‌کاوی عمومی تلقی شود. تفاوت اصلی داده‌کاوی و تحلیل محتوا، در نوع داده و الگوریتم‌های مربوط به آنهاست. همان‌طور که در جدول شماره ۲ مشاهده می‌شود، داده‌کاوی بر داده‌های ساختاریافته تمرکز دارد؛ این درحالی است که در تحلیل محتوا بر روی اطلاعات غیرساخت یافته یا نیمه‌ساخت یافته پردازش صورت می‌گیرد.

متن کاوی	داده کاوی	
متون (غیر ساخت یافته و نیمه ساخت یافته)	داده‌های عددی (ساختار یافته)	موضوع مورد بررسی
متون فاقد شکل	بانک‌های اطلاعاتی رابطه‌ای	ساختار موضوع
باز یابی اطلاعات، دسته‌بندی محتوا، مقایسه متون و...	توصیف، پیش‌بینی و...	هدف
الگوریتم‌های خوشه‌بندی، طبقه‌بندی، زبان‌شناسی، هستان‌شناسی و...	یادگیری ماشین، درخت تصمیم‌گیری، شبکه‌های عصبی، رگرسیون و...	روش‌ها
پیاده‌سازی وسیع از سال ۲۰۰۰	پیاده‌سازی وسیع از سال ۱۹۹۴	دوره توسعه

جدول ۲. تفاوت تحلیل محتوا و داده کاوی

روش‌شناسی پژوهش

قلمرو مکانی این پژوهش جهت ارائه الگوی تحلیل داده‌های راهبردی شبکه‌های اجتماعی، توییت‌ها بوده و جامعه آن، کاربران ایرانی توییت‌هاست. روش گردآوری، پیمایش شبکه اجتماعی بوده و سعی شده‌است تا با استفاده از ابزارهای موجود، نسبت به دریافت توییت‌های کاربران اقدام شود. در این تحقیق، از طریق طراحی یک برنامه کاربردی سعی شده‌است تا داده‌های هدف از شبکه اجتماعی توییت‌ها دریافت شود؛ بدین صورت که داده‌های دریافت شده پس از خلاصه‌سازی، مضمون‌بندی (کد) شده و پژوهشگر با مرور متون، معانی را کدبندی و استخراج کرده و پس از آن نسبت به کمی‌سازی اقدام شد. در مرحله بعد، شمارگان کلمات به دست آمد و با رویکردی جدید گرافی از شبکه کلمات حاصل شد و در کنار آن، به تحلیل محتوای کیفی پرداخته شد. تحلیل محتوای کیفی را می‌توان روش تحقیقی برای تفسیر ذهنی محتوای داده‌های متنی از طریق فرآیندهای طبقه‌بندی نظام‌مند، کدبندی و تم‌سازی یا طراحی الگوهای شناخته‌شده دانست (هسیه و شانون، ۲۰۰۵). تحلیل محتوا نه تنها به خلاصه‌سازی متن اصلی کمک می‌کند؛ بلکه نگرش‌ها و ادراک صاحب پیام را نیز منتقل می‌نماید (کریستوفر^۲ و همکاران، ۲۰۱۳).

1. Hsieh & Shannon
2. Christopher

چارچوب عملی و کارکرد پژوهش

در این پژوهش، نشان داده می‌شود که فرآیند پیاده‌سازی تحلیل محتواهای راهبردی چگونه است. در ابتدا، با استفاده از سرویس‌های نرم‌افزاری منابع باز^۱ و با استفاده از API توییتر^۲ (ای.پی.ای که به برنامه‌های دیگر اجازه استفاده از داده‌های خود را می‌دهد) شبکه محتوایی هدف - که توییتر است - بررسی شده و با تعریف کلیدواژه‌های مرتبط با جستجو^۳ تعیین محل ذخیره‌سازی^۴ جمع‌آوری داده‌های مرتبط آغاز می‌شود.

۱. به‌کارگیری الگو

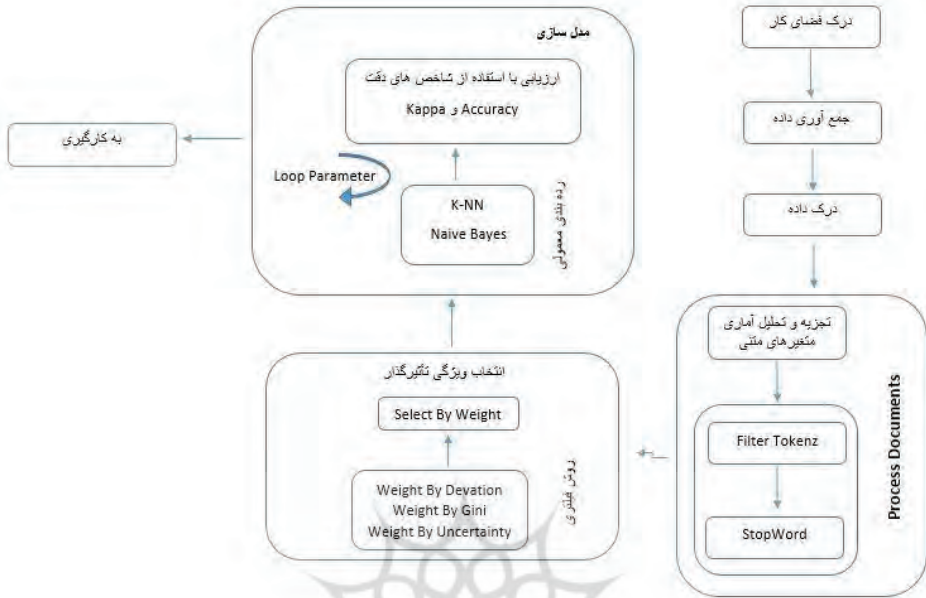
دانش حاصل شده به‌سازماندهی نیاز دارد و باید به‌شکلی ارائه شود که بهره‌وران بتوانند از آن استفاده کنند. بسته به ملزومات کار، فاز به‌کارگیری می‌تواند به‌سادگی ایجاد یک گزارش یا به پیچیدگی اجرای یک فرآیند قابل تکرار داده‌کاوی باشد. با توجه به اینکه این پژوهش، یک پژوهش کاربردی است، نتایج می‌تواند به‌صورت یک سیستم سایه‌ای، به‌موازات کار کارشناسان کاربرد داشته باشد؛ به این معنا که نتایج طبقه‌بندی و ارزیابی سیستم داده‌کاوی، بر ارزیابی کارشناسان تصمیم‌ساز در سازمان تأثیرگذار باشد.

۲. درک مسئله کاری

این مرحله، ابتدا بر درک اهداف و ملزومات طرح از دیدگاه سازمان متمرکز می‌شود. هدف داده‌کاوی در این پژوهش، "پیش‌بینی‌کننده" و نوع آن، "طبقه‌بندی" است؛ بدین‌گونه که با توجه به بررسی محتوای کاربران شبکه‌های اجتماعی، نظرات آنها را در طیف‌های مختلف جمع‌آوری، تحلیل و دسته‌بندی می‌نماید.

پژوهشگاه علوم انسانی و مطالعات فرهنگی
پرتال جامع علوم انسانی

1. Open source
2. Application Programming Interface
3. Keyword
4. Storage



شکل ۴. الگوی کلی فرآیند تحلیل داده‌های شبکه‌های اجتماعی

۳. درک داده‌ها

این مرحله با جمع آوری اولیه داده‌ها شروع می‌شود و به توصیف داده‌ها و تعیین کیفیت آنها می‌پردازد. متن (کامنت)‌ها به صورت انسانی کدبندی شده و هر کاربر با توجه به دیدگاهش بر چسب (لیبل) زده می‌شود. در این پژوهش، دسته‌بندی‌ها به صورت زیر تعریف شدند:

* فکاهی؛

* مطالبه‌ای؛

* خبری؛

* انتقادی؛

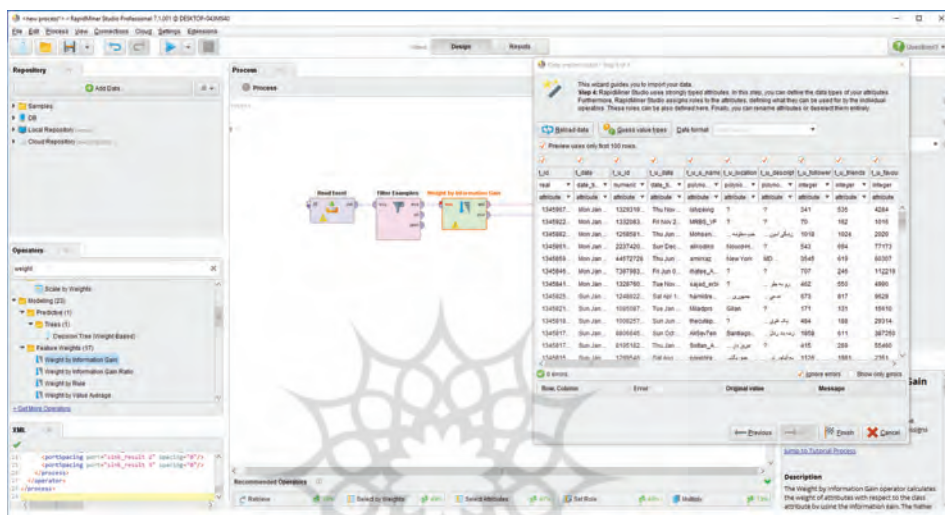
* اعتراضی؛

* تمجید؛

* نامرتبط.

پس از دریافت داده‌ها با کلیدواژه مورد نظر (که در این نمونه، پلیس فتا مد نظر قرار گرفته و تمامی متن‌هایی که پلیس فتا در آن درج شده است، جمع‌آوری شده است) کاربران و متن‌های

آنها در دسته‌بندی‌های فوق تقسیم می‌شوند. این مرحله، در واقع، "آماده‌سازی داده" است. مرحله آماده‌سازی داده‌ها، شامل کلیه فعالیت‌هایی است که برای ساختن مجموعه داده‌های نهایی (داده‌هایی که برای الگوسازی آماده شده‌اند) از داده‌های خام اولیه به کار می‌رود. هر چه کیفیت این آماده‌سازی بهتر باشد، الگوسازی نیز بهتر خواهد بود (مجیبی‌منش، ۱۳۹۰).



شکل ۵. ورود داده‌ها به نرم‌افزار رپید ماینر

وظیفه آماده‌سازی داده‌ها در چند دوره انجام می‌گیرد و هیچ ترتیب از پیش تعریف شده‌ای ندارد. این وظیفه شامل انتخاب جداول^۱، رکوردها^۲، خصیصه‌ها^۳ و نیز انتقال^۴ و پاک‌سازی داده^۵ برای الگوسازی است. در این مرحله، داده‌هایی که در مرحله قبل دریافت شدند، استخراج و ثبت شدند. سپس داده‌ها در یک پایگاه داده جامع و یکپارچه (البته پالایش نشده) به‌منزله یک پایگاه داده رابطه‌ای قرار گرفتند. سپس داده‌ها پالایش شده و ساختار موردنظر برای الگوسازی روی آنها اعمال شد. برای پاک‌سازی و پیش‌پردازش داده‌ها، دو عملیات کاهش داده^۶ و اعمال تغییرات^۷ در شکل داده‌ها، روی پایگاه داده رابطه‌ای^۸ صورت گرفت.

1. Tables
2. Records
3. Attributes
4. Transferring
5. Data Cleaning
6. Reduction
7. Change assigning
8. Relational Database

۴. الگوسازی

در این مرحله، انواع روش‌های الگوسازی انتخاب شده و به کار گرفته می‌شوند. در مجموع، برای یک نوع مسئله داده کاوی، چندین روش وجود دارد. برخی از روش‌ها نیازمند شکل ویژه‌ای از داده‌ها هستند؛ بنابراین، به منظور آماده‌سازی داده برای اعمال روش‌ها در پژوهش حاضر، با استفاده از نرم‌افزار رپیدماینر^۱ الگوسازی داده‌ها، به‌طور جداگانه روی هر الگوریتم طبقه‌بندی انجام گرفت و الگوریتمی که بالاترین صحت را داشت، مبنای الگوسازی و استخراج دانش از آن قرار گرفت. هدف از این کار، استخراج دانش با توجه به داده‌های آموزشی از الگوریتم انتخابی است.

۵. ارزیابی نتایج

در این مرحله از طرح، الگویی که از دیدگاه تحلیل داده، کیفیت بالایی دارد، ساخته شده است. پیش از اقدام برای به‌کارگیری الگو، باید الگو به‌طور کلی ارزیابی شده^۲ و گام‌های اجرایی برای تطابق با اهداف سازمانی مرور شوند. در اینجا، با توجه به مقایسه‌های صورت گرفته، روشی که از بقیه پیش‌بینی دقیق‌تری را انجام می‌دهد، برای استفاده بهره‌وران ارائه می‌شود.

۶. ورود داده و پیش‌پردازش

داده‌های آماده‌شده در نرم‌افزار رپیدماینر وارد شده و در مراحل ابتدایی، تعیین می‌شود که کدام ویژگی‌های داده از اهداف داده‌پردازی موردنظر محسوب می‌شوند. انتخاب ویژگی را می‌توان به‌عنوان فرآیند شناسایی ویژگی‌های مرتبط و حذف ویژگی‌های غیرمرتبط و تکراری با هدف مشاهده زیرمجموعه‌ای از ویژگی‌ها که مسئله را به‌خوبی و با حداقل کاهش درجه کارایی تشریح می‌کنند، تعریف کرد. این کار مزایای گوناگونی دارد؛ از جمله:

الف. بهبود کارایی الگوریتم‌های یادگیری ماشین؛

ب. درک داده، کسب دانش درباره فرآیند و کمک به بصری‌سازی آن؛

ج. کاهش داده کلی، محدود کردن نیازمندی‌ها ذخیره‌سازی و احتمالاً کمک به کاهش هزینه‌ها؛

د. کاهش مجموعه ویژگی‌ها، ذخیره‌سازی منابع در دور بعدی گردآوری داده یا در طول بهره‌برداری؛

ه. سادگی و قابلیت استفاده از الگوهای ساده‌تر و کسب سرعت.

در کنار همه دلایل پیشین، در سناریوهای "تحلیل کلان‌داده"، انتخاب ویژگی نقشی اساسی ایفا می‌کند. با توجه به حجم داده‌های سه‌ساله، برای بالابردن دقت الگوها، از انواع الگوریتم‌های کاهش

1. RapidMiner
2. Evaluation

داده^۱ که یک ویژگی تکنیکی است، برای کاهش داده‌ها در فرآیند داده‌کاوی استفاده می‌شود. در کاهش داده‌ها حجم داده‌ها را کاهش می‌دهند تا بتوان از آنها برای تجزیه و تحلیل کارآمدتر استفاده کرد؛ چراکه مجموعه داده ممکن است دارای تعداد زیادی ویژگی باشد (felici & others, 2006).

۷. یادگیری نظارت شده

در یادگیری نظارت شده^۲ کار با ایمپورت کردن مجموعه داده‌های شامل ویژگی‌های خصیصه‌های آموزش^۳ خصیصه‌های هدف^۴ آغاز می‌شود. الگوریتم یادگیری نظارت شده رابطه بین مثال‌های آموزش و متغیرهای هدف مختص آنها را به دست می‌آورد و آن رابطه یادگرفته شده را برای دسته‌بندی ورودی‌های کاملاً جدید (بدون هدف‌ها) مورد استفاده قرار می‌دهد. برای نمایش اینکه یادگیری نظارت شده چگونه کار می‌کند، یک مثال از پیش‌بینی نمرات دانش‌آموزان بر پایه ساعات مطالعه آنها ارائه می‌شود:

$$Y = f(X) + C$$

که در آن:

- F رابطه بین نمرات و تعداد ساعاتی است که مدیران راهبردی به منظور آماده‌شدن برای امتحانات به مطالعه می‌پردازند؛
- X ورودی است (تعداد ساعاتی که دانش‌آموز خود را آماده می‌کند)؛
- Y خروجی است (نمراتی که دانش‌آموزان در آزمون کسب کرده‌اند)؛
- C یک خطای تصادفی است.

هدف نهایی یادگیری نظارت شده، پیش‌بینی Y با حداکثر دقت برای ورودی جدید داده شده X است. چندین راه برای پیاده‌سازی یادگیری نظارت شده وجود دارد که برخی از متداول‌ترین رویکردها در ادامه مورد بررسی قرار می‌گیرند.

بر پایه مجموعه داده موجود، مسئله یادگیری ماشین در دو نوع دسته‌بندی^۵ و رگرسیون^۶ قرار می‌گیرد. اگر داده‌های موجود دارای مقادیر ورودی (آموزش) و خروجی (هدف) باشند، مسئله از نوع دسته‌بندی است. اگر مجموعه داده دارای مقادیر عددی پیوسته^۷ بدون هرگونه برچسب هدفی

1. Feature subset selection
2. Supervised Training
3. Training attributes
4. Target attributes
5. Classification
6. Regression
7. Continuous numerical values

باشد، مسئله از نوع رگرسیون محسوب می شود.

۸. دسته بندی

در این پژوهش، داده های توییت به منظور پیش بینی اینکه کاربران در خصوص اقدامات پلیس فتا چه گفته اند، مد نظر قرار گرفته است؛ با این هدف که نتیجه این تحقیق در دسته بندی می تواند برای مدیران راهبردی کارگشا باشد. دسته بندی^۱ یک مسئله پیش بینی^۲ است که برچسب های^۳ کلاس دسته ای را که گسسته یا بدون ترتیب هستند، پیش بینی می کند. این یک فرآیند دو مرحله ای است که شامل مرحله یادگیری و دسته بندی می شود.

روش های دسته بندی و انتخاب بهترین آنها

الگوریتم های متداول داده کاوی برای تحلیل داده استفاده می شود و بسته به نوع داده ها و اهداف مد نظر و بر اساس منابع علمی و تجربی دانشمندان علم داده^۴، می توان الگوریتم مورد نظر را در نرم افزار به کار گرفت تا به نتایج دقیق تری حاصل شود. در اینجا، به برخی از این الگوریتم ها و ویژگی های آنها پرداخته می شود.

الف. نزدیک ترین همسایگی (KNN^۵):

روشی از دسته بندی است و با وجود سادگی، عملکرد خوبی را برای مجموعه های آموزشی بزرگ ارائه می دهد که بر اساسی ترین فرض نهفته در پیش بینی ها تکیه می کند؛ اینکه مشاهدات با ویژگی های مشابه نتایج مشابهی خواهند داشت. روش KNN یک مقدار پیش بینی شده را برای یک مشاهده جدید بر اساس تعدد یا میانگین در مجموعه آموزشی^۶ تعیین می کند؛ به این معنا که با مقداری از داده با الگوریتم KNN به سیستم یاد می دهیم که داده های دسته بندی نشده جدید را با توجه به خصیصه های آنها (همسایگی) در کنار هم دسته بندی نماید. با توجه به حجم نامحدودی از داده ها، هر مشاهده ای دارای "همسایه" های زیادی خواهد بود که با توجه به همه ویژگی های اندازه گیری شده به طور خودکار نزدیک هستند (Joshua S. Richman, 2013).

ب. درخت تصمیم (Decision Tree)

منظور از "درخت تصمیم" الگوهای پشتیبانی تصمیم هستند که الگوها را با استفاده از یک سلسله

1. Classification
2. Prediction
3. Label
4. Data Science
5. K- Nearest Neighborhood
6. Training dataset

قوانین کاملاً مشخص^۱ طبقه‌بندی می‌کنند. آنها نمودارهای درختی^۲ هستند که در آنها هر گره شاخه^۳ یک دسته بین تعدادی از دسته‌ها را نشان می‌دهد و هر گره برگ^۴ نشان‌دهنده نتیجه‌ای از انتخاب‌های این دسته‌ها است (Ranganathan, 2013).

ج. نایو بیز (Naive Bayes)

در مرحله یادگیری، الگوی دسته‌بندی، دسته‌بند را با تحلیل مجموعه داده آموزش می‌سازد. در مرحله دسته‌بندی، برچسب‌های کلاس برای داده‌های موجود پیش‌بینی می‌شوند. تاپل‌های مجموعه داده و برچسب‌های کلاس مربوط در حال تحلیل، به دو دسته "مجموعه آموزش" و "مجموعه آزمون" تقسیم می‌شوند. مجموعه داده آزمون به منظور تخمین پیش‌بینی دسته‌بندی استفاده می‌شود. دقت دسته‌بندی در صدی از آزمون است که توسط دسته‌بند به درستی دسته‌بندی شده‌اند. به منظور کسب دقت بیشتر، بهترین راه آزمون الگوریتم‌های گوناگون و آزمودن شاخص‌های مختلف در هر الگوریتم است. بهترین حالت با اعتبارسنجی متقابل^۵ قابل انتخاب است. برای انتخاب یک الگوریتم خوب برای مسئله، شاخص‌هایی مانند دقت، زمان آموزش، خطی بودن^۷ تعداد شاخص‌ها و شرایط خاص باید برای الگوریتم‌های متفاوت در نظر گرفته شود.

طرح عملی تحلیل داده‌های راهبردی

در این تحقیق، جملات انتشار یافته توسط کاربران و نظرات آنها نسبت به پلیس فتا دسته‌بندی شد. داده‌ها پس از متن‌کاوی به الگویی برای کل داده‌های سه‌ساله تبدیل شد که دقت‌های به دست آمده با سه الگوریتم پیش‌گفته، در جدول‌های زیر قابل مشاهده است.

پژوهشگاه علوم انسانی و مطالعات فرهنگی
پرتال جامع علوم انسانی

1. Well-Defined
2. Tree-like graphs
3. Branch node
4. leaf node
5. Training set
6. Cross-validation
7. Linearity

دقت با الگوریتم k نزدیک ترین همسایگی

accuracy: 92.41% +/- 1.60% (mikro: 92.41%)

	true فکاهی	true مضایقه	true خبری	true نامرغبت	true انتقادی	true اعتراضی	true تصحیح	class precision
pred. فکاهی	159	62	19	54	2	22	14	47.89%
pred. مضایقه	0	700	6	3	0	0	0	98.73%
pred. خبری	0	4	646	1	0	0	3	98.78%
pred. نامرغبت	4	9	5	402	0	0	0	95.71%
pred. انتقادی	0	0	0	1	14	4	0	73.68%
pred. اعتراضی	1	2	0	2	6	599	0	98.20%
pred. تصحیح	3	0	4	0	0	0	293	97.67%
class recall	95.21%	90.09%	95.00%	86.83%	63.64%	95.84%	94.52%	

دقت با الگوریتم درخت تصمیم

accuracy: 38.04% +/- 10.15% (mikro: 38.04%)

	true فکاهی	true مضایقه	true خبری	true نامرغبت	true انتقادی	true اعتراضی	true تصحیح	class precision
pred. فکاهی	75	0	0	2	0	0	24	74.26%
pred. مضایقه	92	777	676	461	12	348	266	29.52%
pred. خبری	0	0	3	0	0	0	0	100.00%
pred. نامرغبت	0	0	0	0	0	0	0	0.00%
pred. انتقادی	0	0	0	0	6	0	0	100.00%
pred. اعتراضی	0	0	0	0	4	277	0	98.58%
pred. تصحیح	0	0	1	0	0	0	20	95.24%
class recall	44.91%	100.00%	0.44%	0.00%	27.27%	44.32%	6.45%	

دقت با الگوریتم نیویز

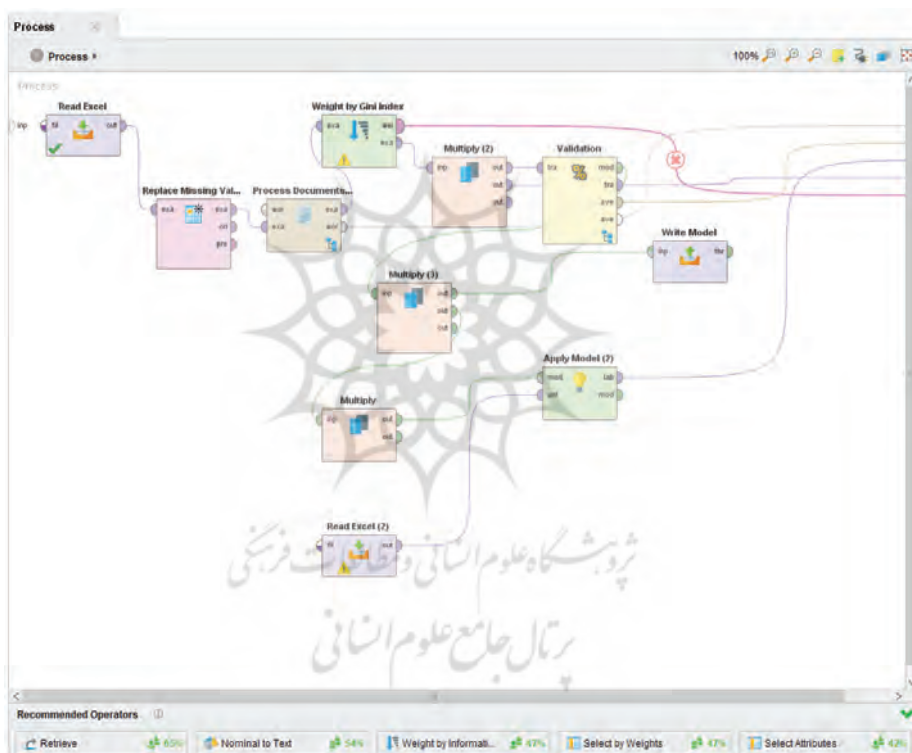
accuracy: 93.23% +/- 1.13% (mikro: 93.23%)

	true فکاهی	true مضایقه	true خبری	true نامرغبت	true انتقادی	true اعتراضی	true تصحیح	class precision
pred. فکاهی	155	3	0	5	0	0	6	91.72%
pred. مضایقه	7	733	10	20	0	9	4	93.61%
pred. خبری	1	12	656	9	2	4	2	95.63%
pred. نامرغبت	0	16	2	411	0	1	4	94.70%
pred. انتقادی	2	0	0	7	20	42	0	28.17%
pred. اعتراضی	1	13	1	7	0	569	0	96.28%
pred. تصحیح	1	0	11	4	0	0	294	94.84%
class recall	92.81%	94.34%	96.47%	88.77%	90.91%	91.04%	94.84%	

جدول ۳. دقت های به دست آمده در تحلیل داده های توییت با کلیدواژه فتا

نوع الگوریتم	دقت
Naive Bayes	۹۳.۲
K-nn	۹۲.۴
Decision Tree	۳۸.۰۴

دسته‌بندی توپیت‌ها به معنای برچسب‌زدن یک مفهوم به هر توپیت و جای‌گذاری توپیت‌های مشابه در یک دسته مشابه است. به عبارتی، هر توپیتی با توجه به دارا بودن کلمات و واحدهای معنایی مشابه مستتر در آن، به یک دسته خاص اختصاص داده می‌شود. تمام توپیت‌های مربوط به پلیس فتا، پس از خوانش اولیه در پنج دسته و مقوله کلی شامل خبری، انتقادی، اعتراضی، مطالبات و فکاهی، بر روی تعداد ۳۹۹۴۵۶ توپیت الگوهای یادگیری ماشین با سه الگوریتم پیش‌گفته^۱ برابر شکل زیر در نرم‌افزار ریپدیمانر محاسبه گردید. در گام بعدی، به کمک دو اپراتور log و loop و parameter دقت الگو با تغییر شاخص‌های مختلف بررسی گردید.



سلسله‌زمانی در داده‌کاوی

ارزش و اهمیت مجموعه داده‌ها و انواع مختلف داده، به دلیل کاربردهای فراوان بر کسی پوشیده نیست. یکی از مهم‌ترین و پرکاربردترین انواع داده، داده‌های سلسله‌زمانی^۲ است. مجموعه‌های

1. Knn, Naive bayse, Decision Tree
2. Time Series

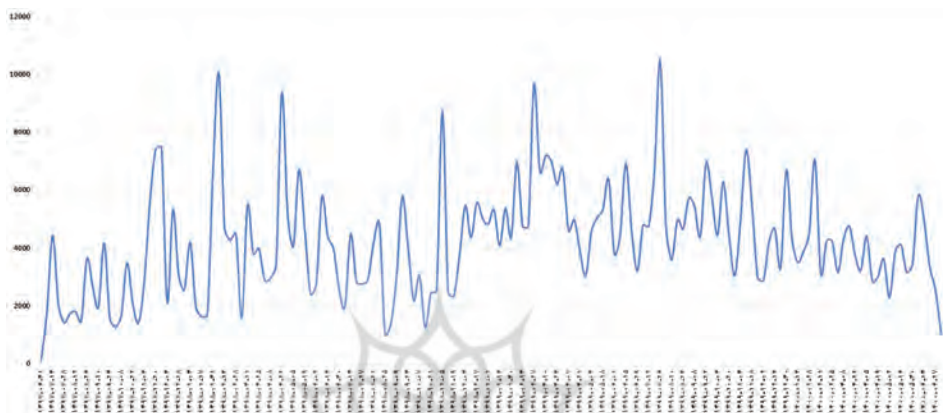
داده‌ای وجود دارد که در آنها ویژگی هدف وابسته به زمان است؛ زیرا با دنباله‌ای از دوره‌های متوالی در طول بعد زمان مرتبط است. در چنین شرایطی، گفته می‌شود که مقادیر متغیر هدف نشان‌دهنده یک سلسله‌زمانی است. از این نوع داده، برای پیش‌بینی استفاده می‌شود. در واقع، سلسله‌زمانی به کشف الگوی رفتار داده‌ها در گذشته و پیش‌بینی رفتار آینده از روی این نوع داده اطلاق می‌شود (Brockwell, 2010). این داده‌ها در داده‌کاوی برای پیش‌بینی، بسیار با اهمیت هستند. یک سلسله‌زمانی مجموعه‌ای از مشاهدات درباره یک متغیر است که در نقاط گسسته‌ای از زمان که معمولاً فاصله‌های مساوی دارند، اندازه‌گیری شده و بر حسب زمان، مرتب شده‌اند؛ بنابراین، یک سلسله‌زمانی از مشاهده یک پدیده در طول زمان به دست می‌آید. سلسله‌زمانی یک کلاس مهم از داده‌های زمانی است و مجموعه مشاهداتی است که به ترتیب زمان انجام شده‌است. داده‌های سلسله‌زمانی دارای ویژگی‌هایی مثل حجم بالای داده، ابعاد بالا و به‌روزرسانی مداوم است؛ علاوه بر این، داده‌های سلسله‌زمانی با ماهیت عددی و پیوسته مشخص می‌شوند و به‌عنوان یک کل در نظر گرفته می‌شود. افزایش استفاده از داده‌های سلسله‌زمانی، تلاش‌های تحقیق و توسعه زیادی را در زمینه داده‌کاوی آغاز کرده‌است. تولید حجم بالای داده‌ها در طول زمان، برای داده‌کاوی بسیار اهمیت دارد. استخراج الگوها و اطلاعات مفید از دل این حجم اطلاعات، اهمیت سلسله‌زمانی را نشان می‌دهد (Rojas, 2019). اگر یک بردار را در نظر بگیریم و این بردار نشان‌دهنده زمان باشد، می‌توان سلسله‌زمانی را این‌گونه تعریف نمود:

$$X(t), \quad t = 0, 1, 2, 3, \dots$$

بنابر معادله فوق، زمان می‌تواند زمان شروع یک پدیده یا زمان ثبت اولین اطلاعات مربوط به حوزه پژوهش باشد؛ بنابراین، یک متغیر تصادفی است و مقدار را در لحظه‌ای نشان می‌دهد.



تحلیل سلسله‌های زمانی یک ابزار راهبردی است Y با این تفسیر که تحلیل داده‌ها کمک می‌کنند تا تصمیم‌گیران اقداماتی را که از منظر کاربران قابل توجه بوده و یا دارای اهمیت فراوان بوده‌است، درک کنند. این نوع فهم از داده‌ها کمک می‌کند تا خطاهای سامانه و سازمان کاهش یابد و اقدامات خود را با اطلاع‌رسانی قبلی و پیوست رسانه‌ای، از گزند شایعات و هجوم اخبار منفی مصون دارد.



شکل ۵. سلسله‌زمانی توجه کاربران به پلیس فتا در توییتر

در شکل بالا، نقاط اوج، نشان دهنده اقدام فتا در موضوعات مختلف و نقاط نزول، زمان‌های عدم توجه به فتا بوده‌است. این نمودار زمانی کمک می‌کند تا برنامه‌ریزی رسانه‌ای به‌گونه‌ای صورت پذیرد که کاربران هدف - که مردم جامعه هستند - مدام در معرض داده‌های منتشره از فتا باشند. مدیران می‌توانند با ترکیب این زمان‌ها با وقایع دیگر، مفاهیم پنهان را کشف و اقدامات موثر را تقویت و یا اقدامات مضر هویت اجتماعی پلیس را کاهش دهند.

ابر کلمات

ابر کلمات^۱ یک روش تجسم است که مجموعه‌ای از برچسب‌های مربوط به یک مورد خاص را خلاصه می‌کند. در ابر کلمات، متخصصین تحلیل داده می‌توانند به‌صورت بصری از رشد یک کلیدواژه و هم‌جواری آن با سایر کلمات یک پدیده را تحلیل کنند، بیشترین کلمات تکرار شده درشت‌تر به نظر می‌رسند و کلمات کنار هم می‌توانند معنای خاصی داشته باشند (Rokne & Alhajj, 2014).



شکل ۵. ابر کلمات حاصل از تحلیل تکرار متن های تولیدشده در توییت کاربران در خصوص پلیس فنا

نتیجه گیری

تحلیل داده های راهبردی نقشه ی موثر را در خدمات سازمانی و شکل گیری هویت سازمانی ایفا می کند. بسته به هدف تحلیل، داده ها به خاطر دارا بودن اصالت می توانند تصمیم گیری ها را به سمت عاری از اشتباه سوق دهند؛ به همین منظور، لازم است که مدیران راهبردی در خصوص تحلیل داده ها و نقش آن در پیشبرد اهداف سازمانی از درک درستی برخوردار باشند. پیاده سازی روش های جمع آوری و تحلیل داده ها برای تمامی ابعاد سازمانی امری الزامی است که بی توجهی به آن موجب خدشه دار شدن هویت و اقدامات مفید سازمان خواهد شد.

در حقیقت، مهم ترین موضعی که می تواند یک سازمان را از لحاظ هویتی (محبوبیت و یا تنفر) نشان دهد، موضوع جهت گیری افکار عمومی نسبت به اقدامات آن است. در این پژوهش، به منظور ارائه الگوی تحلیل داده های راهبردی شبکه های اجتماعی، الگوی جمع آوری و پردازش داده ها ترسیم گردید و به منظور تبدیل الگوی انتزاعی ذهنی به الگوی عملی ملموس، با استفاده از تحلیل محتوا و متن کاوی به بررسی دیدگاه کاربران ایرانی در توییت نسبت به یکی از پلیس های فرماندهی انتظامی

جمهوری اسلامی ایران پرداخته شد. در این مسیر، طی تشریح موضوع، به ساختار داده کاوی و تحلیل متن اشاره شد تا تمایز و تشابه الگوهای تحلیلی درک شود.

فهرست منابع

منابع فارسی

- ضیغمی و همکاران (۱۳۹۷)، "تحلیل محتوا"، نشریه پرستاری ایسران، سال بیست و یکم، شماره ۵۳، صص ۷۶ - ۸۸
- کرپیندورف، کلاوس (۱۳۸۳)، تحلیل محتوا: مبانی روش‌شناسی، ترجمه هوشنگ نایی، تهران، همزه
- ادوارد والتز (۱۳۸۶)، عملیات و اصول جنگ اطلاعات، ترجمه غلامعلی جانگداز، تهران، انتشارات دانشکده امام محمد باقر (ع)
- محبی منش، امید (۱۳۹۰)، "بررسی تأثیر مدیریت کیفیت زنجیره تامین (SCQM) بر رضایت مشتریان با استفاده از مدل‌سازی معادلات ساختاری"، پایان‌نامه کارشناسی ارشد، تهران، دانشکده مدیریت دانشگاه تهران

منابع انگلیسی

- Chianese, A. & Piccialli, F. , 2016. International workshop on Data Mining of Iot Systems (DaMIS): A service oriented framework for analysing social network activities
- Christopher, A. A. , & alias Balamurugan, S. A. (2013, March). Data mining approaches for aircraft accidents prediction: An empirical study on Turkey Airline In Emerging Trends in Computing, Communication and Nanotechnology (ICE-CCN), 2013 International Conference on (pp. 739-745). IEEE
- Kaplan, A. M. , & Haenlein, M. (2010). Users of the world, unite! The challenges and opportunities of Social Media. Business horizons, 53(1), 59-68
- Joshua S. Richman, , 2011, Pages 397-408, Multivariate Neighborhood Sample Entropy: Methods in Enzymology, Volume 487 A Method for Data Reduction and Prediction of Complex Data

- Joo Chuan Tong, Shoba Ranganathan, Computational T cell vaccine design in Computer-Aided Vaccine Design, 2013
- Triantaphyllou, E. and G. Felici(Eds.), Data Mining and Knowledge Discovery Approaches based on Rule Induction Techniques, Massive Computing Series, Springer, Heidelberg, Germany, pp. 227–252, 2005
- P.J. Brockwell, Time Series Analysis, International Encyclopedia of Education (Third Edition), 2010
- Rojas, F. , Herrera, L. J. , Pomares, H. , Rojas, I. (Eds.) Theory and -Applications of Time Series Analysis, Springer, Selected Contributions from ITISE 2019Valenzuela, 2019
- Reda Alhadj, Jon Rokne , Tag Clouds In book: Encyclopedia of Social Network Analysis and Mining Edition: 1, Springer New York, 2014

