

گزارش نهایی طرح پژوهشی

عنوان

**تحلیل و تطبیق میزان همخوانی دامنه و اهداف موضوعی نشریات فارسی وزارت عتف با محتوای**

**مقالات منتشر شده در آنها طی بازه زمانی سه ساله ۱۳۹۷-۱۳۹۹**

**Analyzing and matching the scope of MSRT Persian Journals with the content of  
the articles published in them during the period of 2019-2021**

مجری:

دکتر نرجس ورع

تیر ماه ۱۴۰۲

بِسْمِ اللَّهِ الرَّحْمَنِ الرَّحِيمِ

## چکیده

**هدف:** بر پایه توصیه کمیسیون نشریات علمی در مورد پابندی به محورهای موضوعی اعلام شده از سوی نشریات علمی، همه نشریات باید به سوی این خط مشی حرکت کرده و انتظار می‌رود محتوای مقالات، انعکاس دهنده پژوهش‌های مرتبط و مندرج با دامنه و اهداف نشریه باشد. هدف این رویکرد، فراهم ساختن بستر، به گونه‌ای است که جامعه دانشگاهی اعم از پژوهشگران، نویسندگان مقالات و خوانندگان بتوانند هر نشریه را با موضوع یا موضوعات تخصصی شناسایی کنند. از اینرو در پژوهش حاضر تلاش بر این است که با استفاده از روش تحلیل محتوای مقالات نشریات معتبر وزارت علوم در شش حوزه موضوعی سطح کلان، به تعیین میزان همخوانی مقالات با اهداف و دامنه موضوعی اعلام شده از سوی نشریات پرداخته شود.

**روش:** پژوهش از نوع کاربردی و به روش متن‌کاوی، با رویکرد تحلیلی انجام شده است. جامعه آماری پژوهش حاضر را ۱۱۵۰ عنوان نشریه فارسی مصوب وزارت علوم، پژوهش‌ها و فناوری (۱۴۰۰) بر اساس پرتال نشریات<sup>۱</sup> تشکیل می‌دهد. با توجه به محدودیت‌های زمانی، منابع انسانی و مالی، امکان بررسی اهداف و دامنه موضوعی مقالات تمام نشریات در پژوهش حاضر وجود نداشت؛ بنابراین بایستی حجم نمونه مشخص میشد. در همین راستا، حداقل حجم نمونه مورد نیاز توسط نرم افزار، ۱۱۶ عنوان نشریه تعیین شد. در ادامه با استفاده از روش نمونه‌گیری تصادفی طبقه‌ای، نمونه پژوهش به گونه‌ای انتخاب شد که تناسب بین نمونه و جامعه در هر یک از شش حوزه موضوعی سطح کلان و میانی تعیین شده توسط وزارت عتف رعایت گردد. پس از استخراج داده‌های کتابشناختی مورد نیاز، در ابتدا پیش پردازش شامل نرمال سازی، حذف کلمات اضافی صورت گرفت. سپس دومین مرحله با فنون مختلف متن‌کاوی مانند تکرار واژگان و خوشه بندی انجام شد. در مرحله سوم شباهت‌سنجی و تجزیه و تحلیل دانش انجام شد.

**یافته‌ها:** یافته‌ها نشان داد در ۸ درصد نشریات مورد بررسی، میزان همخوانی یا تطابق محتوای مقالات با دامنه موضوعی مندرج در وب سایت به میزان ۹۵ درصد رعایت شده بود. کمترین میزان همخوانی ۴۰ درصد است که در ۲۰ درصد

---

<sup>1</sup> journals.msrt.ir

نشریات مورد بررسی مشاهده گردید. در این میان ۴۷ درصد نشریات به میزان ۶۰ درصد، ۲۳ درصد نشریات به میزان ۷۰ درصد و ۲ درصد نشریات به میزان ۸۰ درصد به موضوعات اعلامی پایبند بودند.

**نتیجه گیری:** نتایج حاصل از این ارزیابی ها می تواند یاریگر مسئولان بالادستی و سیاست گذاران عرصه علم و دانش کشور در برنامه ریزی های پژوهشی، ارزیابی، رصد پیشرفتهای علمی و ارتقای کیفیت فعالیتهای پژوهشی باشد. همچنین دست اندرکاران نشریات به منظور ارتقای محتوای نشریه و تدوین شیوه نامه های استاندارد می توانند از نتایج این ارزیابی ها بهره مند گردند. در پیش گرفتن رویکرد شفاف موضوعی از سوی نشریات می تواند دستاوردهای سودمندی برای جامعه دانشگاهی داشته باشد. با مشخص شدن گرایش های موضوعی نشریه های علمی، هویت بارزتری برای هر نشریه ایجاد می شود؛ به طوری که هر متخصصی می تواند انتظار داشته باشد که نشریه یا نشریه های خاصی، زمینه های مطالعاتی وی را پوشش می دهد.

**کلیدواژه ها:** نشریه علمی، متن کاوی، دامنه موضوعی نشریه، تحلیل محتوا، مقالات نشریه، مشابهت یابی

## فهرست مطالب

### فصل اول: کلیات

- ۸ - مقدمه
- ۱۲ - بیان مسئله
- ۱۳ - اهمیت پژوهش و ضرورت
- ۱۶ - تعریف های مفهومی و عملیاتی

### فصل دوم: مبانی نظری و مرور پیشینه های پژوهش

- ۱۹ - مقدمه
- ۱۹ - مبانی نظری
- ۴۸ - پیشینه های پژوهش

### فصل سوم: روش پژوهش

- ۵۵ - مقدمه
- ۵۵ - روش پژوهش
- ۵۶ - جامعه پژوهش
- ۵۷ - روش گردآوری داده ها
- ۵۸ - روش تجزیه و تحلیل داده ها

### فصل چهارم: تجزیه و تحلیل داده ها

- ۶۲ - مقدمه
- ۶۲ - استخراج و آماده سازی داده ها

۶۵ - تجزیه و تحلیل

۶۷ - یافته‌ها

### فصل پنجم: نتیجه‌گیری و ارائه پیشنهادها

۷۶ - مقدمه

۷۶ - بحث و نتیجه‌گیری

۷۹ - پیشنهادهای پژوهش

۸۰ - فهرست منابع

# فصل اول

## کلیات پژوهش

نشریات علمی مهم‌ترین محمل برای ارائه آخرین یافته‌های پژوهشی و فناورانه است و منبع راهبردی برای مطالعه روند رشد و تغییر پژوهش و فناوری محسوب می‌گردد (بهکامی و دیم<sup>۱</sup>، ۲۰۱۲). این ابزار کارآمد اطلاع‌رسانی، از مشخصه‌های اصلی ورود یک نظام اجتماعی به دوره توسعه یافتگی تلقی می‌شود که در شکل دهی به شبکه پیچیده تبادل اطلاعات در سطح جهان از نقش ویژه‌ای برخوردار است. نشریات علمی به طور معمول نموداری از حیات علمی جامعه به شمار می‌روند که مورد توجه نویسندگان و پژوهشگران رشته‌های مختلف در سطح وسیعی است. مقاله‌های نشریات به دلیل تازگی و سرعت در ارائه آخرین دستاوردهای علمی، تنوع و کوتاه بودن مطالب، ارائه تحلیل‌های دقیق از موضوع‌های پیچیده علمی، دسترسی آسان، ارائه نقد و بررسی‌های علمی، به‌عنوان بستری برای انتشار و ارائه جدیدترین نتایج و دستاوردهای پژوهش‌های حوزه‌های مختلف علمی به شمار می‌وند. به سبب این ویژگی، نشریات در میان سایر مدارک و منابع علمی از جذابیت بیشتری برای نویسندگان و پژوهشگران بیشتری برخوردار هستند (رضایت، ۱۳۹۴)؛ چرا که مقالات علمی به پژوهشگران کمک می‌کند دانش رشته مربوطه خود را به روز نگه‌دارند و روند تحقیقات و مطالعات خود را مدیریت کنند (آبل و نیولین<sup>۲</sup>، ۲۰۰۲). بنابراین نشریات علمی و پژوهش‌های دانشگاهی همواره دارای ارتباط و مزایای متقابل برای یکدیگر هستند. از یک‌سو، یافته‌های پژوهش‌های نوین عمدتاً در مجلات علمی منتشر می‌شوند (براون<sup>۳</sup>، ۲۰۰۵) و از سوی دیگر تعداد مخاطبان و شناخت آنان از مجلات نیز مبتنی بر پژوهش‌های با کیفیت منتشر شده در آنها است. به بیان دیگر، پژوهشگران برای افزایش آگاهی و انجام پژوهش‌های جدید به مجلات علمی نیاز دارند و اشتهار مجلات علمی نیز از طریق انتشار پژوهش‌های مطرح و اثربخش، کسب می‌شود.

بر اساس مطالب پیش‌گفته و با توجه به نقش و اهمیتی که نشریات علمی در عرصه علم و دانش دارند، لازم است به صورت دوره‌ای از نظر کمی و کیفی مورد ارزیابی قرار گیرند. نتایج حاصل از این ارزیابی‌ها می‌تواند یاریگر مسئولان

---

<sup>1</sup> Behkami & Daim

<sup>2</sup> Abel & Newlin

<sup>3</sup> Brown



بالادستی و سیاستگذاران عرصه علم و دانش کشور در برنامه‌ریزی‌های پژوهشی، ارزیابی، رصد پیشرفت‌های علمی و ارتقای کیفیت فعالیتهای پژوهشی باشد. همچنین دست اندرکاران نشریات بمنظور ارتقای محتوای نشریه و تدوین شیوه نامه های استاندارد می توانند از نتایج این ارزیابی ها بهره مند گردند.

این قبیل ارزیابی ها بر اساس ملاکها، معیارها و استانداردهای علمی، یکی از فرایندهای رایجی است که وزارت علوم، تحقیقات و فناوری، پایگاه استنادی علوم جهان اسلام و مسئولان نشریات به منظور بررسی نقاط قوت و ضعف نشریات مدنظر قرار دارند. به طور معمول ارزیابی ها در دو بخش ساختاری و محتوایی انجام می‌گردد. بخش ساختاری به اطلاعات شناسنامه ای و چارچوب وبگاه نشریه برمی‌گردد؛ اما با توجه به رشد روزافزون تولیدات علمی، تحلیل محتوای مقاله ها و نشریه ها، دسته بندی موضوعی آنها و شناخت انواع مطالب منتشرشده در نشریه های تخصصی برای دسترسی سریع و آسان به اطلاعات، ضروری بنظر می‌رسد.

در فرایند ارزیابی محتوایی نشریات، یکی از روشهای علمی و معتبر تحلیل محتوایی است. تحلیل محتوا، متشکل از مجموعه ای از روندهای گردآوری و سازمان دهی اطلاعات به شکلی استاندارد است که تحلیل گر را قادر می سازد دربارهٔ ویژگی ها و مفاهیم موجود در نشریه، استنباط هایی داشته باشد. بدون شک برای پژوهشگران و دست اندرکاران نشریات این موضوع دارای اهمیت است که بدانند در زمینه علمی موضوعات چه روندی را طی کرده است. تحلیل محتوا، ابزاری قدرتمند برای بررسی روندها، الگوی اسناد، الگوی نوشتاری و غیره است. همچنین اهداف پژوهشگران را می توان بر اساس داده های موجود در نشریه ها و موضوع هایی که آنها برای مقاله خود انتخاب می‌کنند، استنباط نمود (ماجی، جال و مهارانا<sup>1</sup>، ۲۰۱۶). در این راستا انتظار می‌رود محتوای مقالات، انعکاس دهنده پژوهش های مرتبط با دامنه و اهداف نشریه و دربردارنده اولویتها، گرایش ها و موضوع های بنیادین موجود در یک مقطع زمانی باشد. از این رو، بررسی موضوع هایی که در مقالات نشریه بدان پرداخته شده، علاوه بر اینکه تصویری عینی از شرایط موجود و انعکاس پژوهش ها در آن نشریه را نشان می دهد، می تواند نشان دهنده گرایش های موضوعی و کمبودهای احتمالی موجود در حوزه های مورد

---

<sup>1</sup> Maji, jal & meharata

بررسی باشد. عبارتی در گستره وسیع منابع منتشر شده در یک حیطه علمی و پژوهشی، انجام یک مطالعه تطبیقی، علاوه بر آشکارسازی سیر تکوینی مطالعات آن حیطه و شناسایی تفاوتها و شباهتها، راهنمایی برای ساماندهی و مدیریت هرچه بهتر نشر آثار پژوهشی علمی به شمار میآید.

با هدف دستیابی به این مهم، گام نخست، تحلیل محتوای مقالات نشریه مبتنی بر تشریح اهداف و برنامه های آن است. به کارگیری این روش موجب استنتاج معتبر از اطلاعات موجود در متن، فراهم شدن دانش، تشریح وسیع اطلاعات، بینش جدید و راهنمای کاربردی برای عمل میشود (قهنویه و همکاران، ۱۳۹۰؛ هو و هئو، ۲۰۱۳). در این روش، محتوا به صورت نظامدار و کمی توصیف میشود؛ به طوریکه داده های کیفی به داده های کمی تبدیل میشوند. همچنین بررسی گزارش نتایج مطالعات تطبیقی به افراد متخصص در یک رشته علمی این امکان را میدهد که از دانش موجود در ارتباط با حیطه پژوهشی و شغلی خود بهره لازم را ببرند و نقاط ضعف و قوت را دریابند.

به منظور انجام تحلیل محتوا، استخراج خودکار داده های متنی و استفاده از فنون متن کاوی ضروری به نظر می رسد. سالوم و همکاران<sup>۲</sup> (۲۰۱۸) استفاده از فنون متن کاوی را روشی برای شناسایی موضوعات متون علمی و سیر تکاملی این موضوعات بیان می کنند. یکی از الگوریتمهای مهم و مفید متن کاوی، خوشه بندی است. با انجام عملیات خوشه بندی، حیطه گسترده ای از داده های پراکنده در گروه های مدون و سازمان یافته قرار میگیرند. گروههای متعدد ایجاد شده با برخورداری از ویژگیهای مشترک درون هر گروه دارای ارتباط ساختاری با یکدیگر هستند. با این روش داده های مربوط به مقالات درون خوشه های واحد قرار میگیرند به گونه ای که مقالات درون هر خوشه دارای حداکثر شباهت با یکدیگر و حداقل شباهت با دیگر خوشه ها هستند. روشهای متن کاوی و به طور خاص مدل سازی موضوعی قادر است علاوه بر کشف موضوع های پنهان و ارتباط آن با موضوع های آشکار، موضوع های کمتر شناخته شده یا کمتر پرداخته شده را شناسایی کند. تکنیک متن کاوی به عنوان یکی از روشهای تحلیل محتوا از طریق پردازش، استخراج و مرتب سازی اطلاعات به

---

<sup>1</sup> Woo H, Heo

<sup>2</sup> Salom et al.

ترسیم مدل مفهومی و مصورسازی اطلاعات میپردازد. همچنین امکان تحلیل، مسیریابی، نمایش و آشکارسازی ساختار و دانش مفید و ضمنی را از میان انبوهی از داده‌های ساختارنیافته آشکار میکند (شفرین و برون<sup>۱</sup>، ۲۰۰۴). این نوع شناسایی دانش ضمنی، اهمیت زیادی در تعیین اولویت‌های پژوهشی، تدوین برنامه‌های راهبردی سیاست‌گذاران، اطلاع از شکاف موضوع‌های پژوهشی پژوهشگران ایرانی، تأثیرات بین‌رشته‌ای، ارتقای جایگاه حرفه‌ای و شغلی پژوهشگران و پژوهشگران هر رشته داشته و مسیری هموار را پیش روی برنامه‌ریزان و سیاست‌گذاران علمی-پژوهشی در تدوین برنامه‌های راهبردی قرار می‌دهد.

نتایج به دست آمده از مطالعات نشان داده که مدل‌سازی موضوعی نه تنها می‌تواند یک ابزار سودمند برای استخراج اطلاعات از داده‌های متنی باشد؛ بلکه نسبت به بسیاری از رویکردهای سنتی و روش‌های مبتنی بر خوشه‌نیز عملکرد بهتری در امر بازیابی اطلاعات دارند. با توجه به کاربرد مدل‌سازی موضوعی در درک ساختار موضوعی، ارتباطات بین اسناد، گرایش‌های پژوهشی، دارای پتانسیل بالایی در تولید ایده‌های پژوهشی، تشویق به همکاری پژوهشگران و به طور کلی در حوزه سیاست‌گذاری علمی و پژوهشی دارد.

بر اساس مطالب پیش گفته، شناسایی اهداف و حوزه‌های موضوعی نشریات علمی، قلمروهای موضوعی گوناگون، دیداری‌سازی و مطالعه آنها از مسائلی است که نیازمند پژوهش است. در مجموع بر پایه توصیه‌های کمیسیون بررسی نشریات علمی کشور و نیز دفتر سیاست‌گذاری و برنامه‌ریزی امور پژوهشی وزارت علوم، تحقیقات و فناوری در مورد پابندی محورهای موضوعات اعلام شده از سوی نشریات علمی، همه نشریات باید به سوی این خط مشی حرکت کنند. هدف کلی این رویکرد، فراهم ساختن بستر به گونه‌ای که جامعه دانشگاهی اعم از پژوهشگران، نویسندگان مقالات و خوانندگان بتواند هر نشریه را با موضوع یا موضوعات تخصصی شناسایی کنند (آیین‌نامه کمیسیون نشریات وزارت عتف، ۱۳۹۸)

---

<sup>1</sup> Shefrin & brown

از اینرو در پژوهش حاضر تلاش بر این است که با استفاده از روش تحلیل محتوای مقالات نشریات معتبر وزارت علوم در شش حوزه موضوعی سطح کلان، به تعیین میزان همخوانی موضوع یا موضوعات اعلام شده از سوی نشریات از یک سو و مقالات منتشر شده در آنها از سوی دیگر پرداخته شود.

## ۱-۲. بیان مسئله

آمار نشان می دهد که در حال حاضر ۳۵۰۰۰ نشریه در سراسر جهان منتشر می شود (اولریخ<sup>۱</sup>، ۲۰۲۳). در این میان جمهوری اسلامی ایران دارای بیش از ۲۰۰۰ نشریه علمی معتبر وزارتین است که سالانه بالغ بر ۶۰۰۰۰ مقاله را در برمی گیرد (اسفند ۱۴۰۱). از اینرو، انتشار یافته های پژوهشی، برای به اشتراک گذاری دانش و همچنین اعتبار حرفه ای پژوهشگران در یک نشریه مناسب و مرتبط حائز اهمیت و از دشوارترین جنبه های انتشار نتایج پژوهش ها پژوهش ها است (وانگ و هو<sup>۲</sup>، ۲۰۱۵). انتخاب نشریه متأثر از عوامل گوناگونی است که در این میان عامل ارتباط دامنه و حوزه موضوعی نشریه با مقاله از اهمیت ویژه ای برخوردار است؛ زیرا ارسال مقالات حاصل از یک پژوهش، به نشریاتی که ربط کمتری با محتوا و نیاز خوانندگان دارند باعث ائتلاف وقت پژوهشگر می شود. همچنین حوزه سردبیری نشریه نیز مدت زمانی را صرف بررسی و احتمالاً ارسال مقاله به داوری نموده و اگر طی هر یک از مراحل مقاله رد شود این فرآیند باید مجدداً برای نشریه ای دیگر تکرار گردد. نتایج مطالعات نشان می دهد متوسط مدت زمان تصمیم گیری و اعلام نظر جهت پذیرش اولیه مقاله به نویسنده، طولانی و به طور متوسط ۴۱ روز است (ویمن و اسمیتز<sup>۳</sup>، ۲۰۱۷؛ نووین و همکاران<sup>۴</sup>، ۲۰۱۸) که این مدت در ۳۱ درصد موارد به بالای ۶ ماه نیز می رسد (مالیگان، هال و رافائل<sup>۵</sup>، ۲۰۱۳). دلیل اصلی رد مقاله در بسیاری از موارد، در مرحله بررسی اولیه توسط سردبیران و داوران، عدم تطابق و تناسب موضوعی مقاله با دامنه موضوعی نشریه است. در برخی موارد، نرخ رد اولیه

---

<sup>1</sup> Ulrichs

<sup>2</sup> Wang & Hou

<sup>3</sup> Huisman and Smits

<sup>4</sup> Nguyen et al.

<sup>5</sup> Mulligan, Hall, Raphael

حتی قبل از ارسال به داوری تا ۸۸ درصد اعلام شده است (اندرسون<sup>۱</sup>، ۲۰۱۲). بنابراین ارسال مقاله به نشریه نامرتبط، علاوه بر اینکه مزایای کمتری برای ارتقای فردی و سازمانی پژوهشگر به همراه دارد منجر به از دست رفتن زمان قابل ملاحظه، تازگی، اهمیت و تاثیرگذاری یافته‌های پژوهش‌ها پژوهش‌های خواهد شد (هو، وو، لین<sup>۲</sup>، ۲۰۱۶). این چالش با شفاف سازی حوزه های موضوعی مورد تایید نشریه جهت دریافت و پذیرش مقاله تا حد قابل توجهی بهبود خواهد یافت. همچنین تخصصی شدن نشریه های علمی موجب تخصصی شدن فرایند داوری، ویراستاری و نشر مقالات خواهد شد. از این نظر، انسجام و یکدستی بهتری در نشریه ها به وجود خواهد آمد. این وضعیت می تواند از دریافت و داوری مقاله های تکراری در دفاتر مجلات و نیز بروز سرقت علمی جلوگیری نماید.

در همین راستا و به منظور پاسخ به پرسش های پژوهش طرح حاضر بر پایه فنون متن کاوی به بررسی میزان همخوانی دامنه و اهداف موضوعی نشریات فارسی وزارت عتف با محتوای مقالات منتشر شده در آنها پرداخته است. روش های متن کاوی به منظور تحلیل خودکار انتشارات علمی به کاررفته است که به بررسی اسناد به منظور شناسایی مضامین یا موضوعات آنها می پردازد. از نتایج حاصل می توان به منظور تحلیل چگونگی ارتباط مباحث با یکدیگر و چگونگی تکامل آنها باگذشت زمان استفاده کرد.

### ۱-۳. اهمیت پژوهش و ضرورت

- در ایران ۱۴۵۰ عنوان نشریه علمی معتبر توسط وزارت علوم، تحقیقات و فناوری در راستای افزایش مشارکت علمی و انتشار نتایج مطالعات و پژوهش‌ها پژوهش‌ها در سطح ملی و بین المللی منتشر می شود (اسفند ۱۴۰۱). این نشریات به لحاظ کمی و کیفی در سطوح الف تا دال رتبه بندی می شوند. در این میان ۸۵ درصد این نشریات از لحاظ درجه علمی دارای رتبه های الف و ب هستند که این امر نشان از اعتبار علمی بالای آنها دارد. افزایش تعداد نشریه های علمی ایجاب می کند که برای حفظ و توسعه کیفی آنها تلاش شود. از اینرو ارزیابی دقیق موضوعی مقالات این نشریات می تواند شاخصی از روند و گرایش

---

<sup>1</sup> Anderson

<sup>2</sup> Hou, Wu and Lin

فعالیت‌های پژوهشی را نشان دهد. این پژوهش در جهت پایش وضعیت موجود و بهبود جایگاه نشریات کشور در حوزه مربوطه و تعیین الگوی موضوعی است و در مجموع این نکته که نشریات به چه میزان در اعلام و رعایت دامنه و حوزه های موضوعی که در آن مقاله می پذیرند شفاف عمل کرده اند قابل بررسی و حائز اهمیت است.

- در پیش گرفتن این رویکرد از سوی نشریات می تواند دستاوردهای سودمندی برای جامعه دانشگاهی داشته باشد. با تخصصی شدن گرایش های موضوعی نشریه های علمی، هویت بارزتری برای هر نشریه ایجاد می شود؛ به طوری که هر متخصصی می تواند انتظار داشته باشد که نشریه یا نشریه های خاصی، زمینه های مطالعاتی وی را پوشش خواهد داد. تخصص گرایی فرآیندی علمی و عملی است که می توان به یاری آن، پژوهش های اعضای هیات علمی و پژوهشگران کشور را معطوف به حوزه های تخصصی کرد. در تعریف حوزه تخصص باید گفت که حوزه تخصص به آن حوزه یا گرایش خاصی اطلاق می گردد که مشتمل بر مجموعه مسائل نظام مندی است که نقشه راه فعالیت نشریه را در بازه زمانی طولانی تعیین می نماید و نتیجه آن به فراوری دانش انباشته و رسیدن به نظریه های جدید در آن حوزه علمی منتهی می گردد. بدین ترتیب، نشریه و نویسندگان آن به مرجعی در جامعه علمی تبدیل می گردند (فتاحی، ۱۳۹۴).

- با گسترش علوم و تخصصی شدن هر حوزه علمی، جایگاه تخصص ها نیز اهمیت بیشتری می یابد. گرایش های تخصصی در رشته های دانشگاهی و رویکردهای کاملاً تخصصی در پژوهش، حکایت از روند روبه توسعه علم دارند. از سوی دیگر، توسعه گرایش های تخصصی در هر کشور نشان دهنده تلاش اندیشمندان و مدیران آن کشور در جهت توسعه علم و درنوردیدن مرزهای موجود است. به بیان دیگر، تخصص گرایی، شاخص پیشرفت علم در هر جامعه تلقی می شود. در این میان، دانشگاه ها و مراکز آموزشی و پژوهشی نقش اصلی و تعیین کننده دارند. توسعه علم و ایجاد رشته های جدید دانشگاهی، بویژه در مقاطع کارشناسی ارشد و دکترا عامل مهمی در توسعه تخصص ها به شمار می روند. در همین راستا، انتظار می رود اعضای هیات علمی و دانشجویان مقاطع تحصیلات تکمیلی عمده تلاش های خود را در بعد آموزش، پژوهش و فعالیت های علمی بر یک زمینه تخصصی محدود متمرکز و در همان زمینه رشد و متخصص شوند. این روند سال ها در دانشگاه ها و مراکز آموزشی و پژوهشی جوامع پیشرفته اعمال شده و در حال حاضر وضعیت به گونه ای است که هر مؤسسه یا هر فرد در یک یا حداکثر دو

زمینه خاص و محدود، به فعالیت اشتغال دارد. به همین دلیل، فعالیت‌های پژوهشی و تولید علم در آن جوامع از پیشرفت قابل توجهی برخوردار است. هر پژوهشگر در طول حیات علمی خود تنها یک یا دو زمینه خاص را دنبال می‌کند و در همان زمینه(ها) نیز متخصص و صاحب نظر می‌شود. تولیدات علمی (کتاب‌ها و مقاله‌ها) نیز کاملاً تخصصی است. به عبارت دیگر، کمتر مشاهده می‌کنیم که پژوهشگران آن جوامع دچار پراکنده کاری و گسیختگی علمی شوند. در مقابل، نگاهی به مقاله‌ها و کتاب‌های تألیفی یا ترجمه‌ای اعضای هیئت علمی و پژوهشگران ایرانی حاکی از آن است که بسیاری از پژوهشگران دچار نوعی پراکنده کاری و فعالیت در هر زمینه‌ای هستند. عمومیت این وضعیت در رشته‌های علوم انسانی و اجتماعی بیش از دیگر علوم است. بنظر میرسد چنانچه هر نشریه در حوزه‌ای تخصصی مربوطه به جنبه‌های نظری و مفهومی، در عرصه کار علمی طی سال‌های فعالیتش متمرکز شود، خروجی‌های قابل استفاده و عمیق تری خواهد داشت. چنین راهبردی به عمق و غنای کارهای علمی پژوهشگران بسیار خواهد افزود و باعث می‌شود که پراکنده کاری و سطحی نگری به حداقل ممکن کاهش یابد (گاهنامه عتف، ۱۳۹۵). از این رو، مراجعه منظم به نشریه‌های معین به شکل چاپی یا الکترونیکی، وی را در جریان روند پژوهش‌های تخصصی حوزه مورد علاقه او قرار خواهد داد.

- با توجه به اهمیت انتشار یافته‌های پژوهشی و هزینه و زمانی که صرف انجام آن می‌شود؛ توجه به موضوع پژوهش‌های انجام شده در دوره‌های زمانی مختلف حائز اهمیت است. چنین بررسی‌هایی نشان می‌دهد که در دوره‌های زمانی مختلف به چه موضوعاتی بیشتر یا کمتر پرداخته شده است. هنگامی که کاستی‌های پژوهش در موضوعاتی خاص تعیین شد می‌توان انتظار داشت که خط مشی پژوهش‌های آینده طوری برنامه‌ریزی شود که کاستی‌ها بمنظور شناسایی اولویت‌های پژوهشی انجام و تصویری کلی از سیر موضوعی که نشریات در پیش گرفته‌اند مشخص گردد.

- با مشخص شدن دامنه موضوعی نشریات، پژوهشگران می‌توانند در انتخاب و استفاده از منابع مرتبط در یک حوزه تخصصی بهتر عمل کنند؛ لذا نتایج این پژوهش می‌تواند در این امر نیز به کار گرفته شود. همچنین دست‌اندرکاران نشریات نیز می‌توانند در تامین نیاز جامعه هدف مخاطبان هدف دقیق‌تر عمل کنند.

- همچنین عمق بخشیدن به پژوهش‌ها و توانمندسازی نشریات در جهت شناسایی مرزهای علمی تخصصی متقن و معتبر،

زمینه سازی برای تعمق، ارتقا و شکوفایی خلاقیت و ابتکار علمی محققین، افزایش کارایی نشریات در جهت گشایی و حل مسائل واقعی کشور، امکان ایجاد شبکه های علمی تخصصی بین نویسندگان، استفاده از اطلاعات تخصصی در فعالیت تصمیم گیری مدیران ارشد مراکز پژوهشی وزارت علوم، تحقیقات و فناوری از دیگر نتایج طرح حاضر در رفع برخی چالشها خواهد بود.

#### ۴-۱. هدف پژوهش

هدف اصلی پژوهش حاضر، تحلیل میزان همخوانی مقالات منتشر شده با دامنه موضوعی نشریات وزارت علوم، تحقیقات و فناوری بر پایه فنون متن کاوی است.

در این راستا اهداف جزئی به شرح زیر تعریف می گردد:

- بررسی میزان همخوانی محتوای مقالات نشریات وزارت علوم در شش حوزه موضوعی علوم انسانی، فنی و مهندسی، علوم پایه، هنر و معماری، کشاورزی و منابع طبیعی و دامپزشکی
- فراوانی نشریات بر اساس میزان همخوانی محتوای مقالات نشریه با حوزه موضوعی مندرج در وبگاه به تفکیک حوزه موضوعی

#### ۴-۲. تعریف های مفهومی و عملیاتی

• متن کاوی

○ تعریف مفهومی

متن کاوی تکنیکی میان رشته ای است و این امکان را فراهم می کند تا بتوان از طریق شناسایی و اکتشاف الگوهای در داده های متنی از کلان داده ها به شکلی مفید استفاده نمود (تروینس و همکاران<sup>۱</sup>، ۲۰۱۴). همچنین متن کاوی می تواند در جهت درک

---

<sup>1</sup> Truyens et al.



بهتر اطلاعات موجود در اسناد به کار گرفته شود. در واقع متن کاوی امکان بازیابی و درک اطلاعات پنهان در متن‌ها را فراهم می‌کند و برای کشف ساختار، الگوها و دانش در مجموعه‌های متنی بزرگ به کار می‌رود (تسنگ<sup>۱</sup> و همکاران، ۲۰۰۷).

#### ○ تعریف عملیاتی

در این پژوهش از متن کاوی جهت شناسایی و اکتشاف محتوای موضوعی مقالات نشریات استفاده گردید. بدین ترتیب که ابتدا با مراجعه به وبگاه هر نشریه اهداف و دامنه موضوعی و همچنین اطلاعات کتابشناختی مقالات بازه زمانی مورد بررسی استخراج و با استفاده از فنون متن کاوی به داده‌هایی که قابل خواندن و پردازش برای ماشین باشد تبدیل شد.

### • خوشه‌بندی موضوعی

#### ○ تعریف مفهومی

یکی از الگوریتمهای مهم و مفید متن کاوی، خوشه بندی است. با انجام عملیات خوشه بندی، حیطه گسترده ای از داده های پراکنده در گروه های مدون و سازمان یافته قرار میگیرند. گروههای متعدد ایجاد شده با برخورداری از ویژگیهای مشترک درون هر گروه دارای ارتباط ساختاری با یکدیگر هستند. با این روش داده های مربوط به مقالات درون خوشه های واحد قرار میگیرند به گونه ای که مقالات درون هر خوشه دارای حداکثر شباهت با یکدیگر و حداقل شباهت با دیگر خوشه ها هستند (لامبا و مدهوسدان<sup>۲</sup>، ۲۰۱۸).

#### ○ تعریف عملیاتی

در این پژوهش محتوای نشریات با استفاده از الگوریتم فضای برداری مورد بررسی و تحلیل قرار گرفته است. بدین ترتیب که برای هر مدرک یک بردار مبتنی بر فرکانس تکرار واژه‌ها تشکیل و سپس با برداری که بر اساس واژه‌های کلیدی مستخرج از دامنه و حوزه‌های موضوعی ایجاد گردیده شباهت سنجی شد.

---

<sup>1</sup> Tseng et al.

<sup>2</sup> Lamba & Madhusudhan

**فصل دوم**

**مبانی نظری**

**و**

**مرور پیشینه‌های پژوهش**

## ۲-۱. مقدمه

در این فصل مبانی نظری تبیین گردیده و پیشینه‌های پژوهش مرور شده است. به بیان دیگر، ابتدا به مبانی نظری پژوهش در مباحث مرتبط با داده‌کاوی و متن‌کاوی پرداخته شده و سپس پیشینه‌های پژوهش‌های در سطح ملی و بین‌المللی مرور شده و در پایان به استنتاج و نتیجه‌گیری از مرور پیشینه‌های پژوهش پرداخته شده است.

## ۲-۲. مبانی نظری

### ۲-۲-۱. داده‌کاوی

پیشرفت‌های به وجود آمده در جمع‌آوری داده‌ها و قابلیت‌های ذخیره‌سازی در طی دهه‌های اخیر باعث شده در بسیاری از علوم با حجم بزرگی از داده‌ها روبرو شویم. داده‌کاوی کوششی برای به‌دست آوردن اطلاعات مفید از میان این داده‌هاست و رشد بی‌رویه داده‌ها در سطح بین‌المللی اهمیت داده‌کاوی را دوچندان کرده است. فناوری مدیریت پایگاه داده‌های پیشرفته انواع مختلفی از داده‌ها را می‌تواند در خود جای دهد، در نتیجه تکنیک‌های آماری و ابزار مدیریت سستی برای تحلیل این داده‌ها کافی نیست و استخراج دانش از این حجم بسیار زیاد داده‌ها چالشی اساسی تلقی می‌شود.

برای داده کاوی تعاریف متعددی وجود دارد. برخی از این تعاریف عبارت‌اند از: فرآیند به خدمت گرفتن یک روش‌شناسی رایانه‌ای که با استفاده از فنون مختلف، دانش را به طور مستقیم از داده‌ها استخراج می‌کند (بارس و کمپور<sup>۱</sup>، ۲۰۰۸). داده کاوی جستجویی است برای اطلاعات جدید و نوین از میان مقادیر کلان داده‌ها و فرآیندی است مشارکتی میان انسان و کامپیوتر (یوکسل ترک و همکاران<sup>۲</sup>، ۲۰۱۴). داده کاوی فرآیند اکتشاف و تحلیل داده‌ها به وسیله ابزار خودکار و نیمه خودکار به منظور اکتشاف الگوهای معنی‌دار و قواعد است (پراناتا و اسکینر<sup>۳</sup>، ۲۰۱۵).

اصولاً داده کاوی، پایگاه‌های داده‌ی بزرگ را به عنوان منبع بالقوه‌ای از دانش ارزشمند برای تصمیم‌گیری در نظر می‌گیرد. در فرآیند داده کاوی بهترین نتیجه زمانی حاصل می‌شود که دانش یک فرد خبره در خصوص یک مسئله با توانایی‌های کامپیوتر ترکیب شود (خان<sup>۴</sup>، ۲۰۱۶). فقط در این صورت است که می‌توان بدون نفی توانایی‌های فرد خبره، سیستم‌های کامپیوتری را در خدمت نیرومندتر ساختن او قرارداد.

## ۱-۱-۲-۲. اهداف داده کاوی

دو هدف اساسی داده کاوی، پیش‌بینی<sup>۵</sup> و توصیف<sup>۶</sup> است. در عملیات پیش‌بینی بعضی از متغیرها یا حوزه‌هایی از مجموعه‌های داده به منظور پیش‌بینی ارزش ناشناخته یا ارزش آینده داده‌های دیگر مورد استفاده قرار می‌گیرد. همچنین داده کاوی بر یافتن الگوهای تشریحی داده‌ها که به وسیله انسان می‌تواند تعبیر شود، تمرکز می‌کند. در نتیجه داده کاوی را می‌توان در یکی از دو گروه زیر جای داد:

- **پیش‌بینی:** این روش با استفاده از مجموعه داده‌ها، مدل‌هایی را برای توضیح سیستم تولید می‌کند که با استفاده از آن‌ها می‌توان عملکرد متغیرهای مختلف را پیش‌بینی کرد.

---

<sup>1</sup> Baars and Kemper

<sup>2</sup> Yukselturk et al.

<sup>3</sup> Pranata and Skinner

<sup>4</sup> Khan

<sup>5</sup> Prediction

<sup>6</sup> Description

- **توصیف:** اطلاعات جدید را بر اساس مجموعه داده‌های در دسترس تولید می‌کند این داده‌ها الگوهای رفتاری متغیرها را تشریح می‌کند.

هدف از داده کاوی پیش‌بینی‌کننده تولید مدلی است که با استفاده از یک کد اجرایی وظایفی چون پیش‌بینی، دسته‌بندی، تخمین مقدار، تخمین عملکرد و برخی دیگر از موارد را انجام دهد. هدف از داده کاوی تشریحی دستیابی به درکی کامل از سیستم تحت بررسی با استفاده از الگوهای پنهان در آن و روابط درون مجموعه‌های داده است (کانتاردزیک<sup>۱</sup>، ۲۰۱۱).

## ۲-۱-۲-۲. ریشه‌های داده کاوی

پایه‌های اصلی داده کاوی بر دو اصل آمار و یادگیری ماشین<sup>۲</sup> استوار است. آمار نیز ریشه در ریاضیات و منطق داشته و بنابراین داده کاوی نیز علاوه بر آمار ریشه در این دو علم دارد. در مقابل یادگیری ماشینی نیز علمی مبتنی بر رایانه است که اصول آن را در هوش مصنوعی می‌توان یافت. تضادی که در اینجا آشکار می‌شود این است که آمار به دلیل طبیعت ریاضی خود متمایل به فرموله کردن مسائل و مدل‌سازی است، درحالی که یادگیری ماشینی مسائل را با استفاده از الگوریتم‌ها حل می‌کند. اینجاست که باید نسبت به ترکیب این دو علم برای استفاده از آن‌ها در داده کاوی اقدام کرد.

داده کاوی رویه‌های تحلیلی را در زمینه‌های آمار، ریاضیات، تجارت و نظریه اقتصاد پیوند می‌زند. داده کاوی علاوه بر علوم فوق به‌خاطر استفاده از اصول اساسی مدل‌سازی از نظریه کنترل نیز استفاده می‌کند. این نظریه عموماً در سیستم‌های مهندسی و فرآیندهای صنعتی مورد استفاده قرار می‌گیرد. بنابراین داده کاوی یک فناوری میان‌رشته‌ای است و برای استفاده مؤثر از آن باید از علوم تشکیل‌دهنده آن شناخت کافی داشت. البته زمانی که بخواهیم از داده کاوی برای مقاصد نوآورانه و خلاقانه‌تر استفاده کنیم، نیاز به این شناخت به مراتب عمیق‌تر می‌شود.

---

<sup>1</sup> Kantardzic

<sup>2</sup> Machine Learning

علی‌رغم ارتباط میان داده‌کاوی و آمار، تفاوت‌های اساسی میان این دو علم وجود دارد. آمار، علمی تأییدی<sup>۱</sup> است؛ یعنی کوشش دارد مفروضاتی را با استفاده از فنون مختلف تصدیق یا رد کند، درحالی‌که داده‌کاوی یک علم اکتشافی<sup>۲</sup> است، بدین معنی که سعی دارد الگوهای دانشی داده‌های موجود را کشف نماید. همچنین آمار استنتاجی از نمونه‌های کوچک و تعمیم آن‌ها به جامعه استفاده می‌کند و ماهیتاً توان پردازش نمونه‌های بزرگ را ندارد درحالی‌که در داده‌کاوی از نمونه‌های بسیار بزرگ و حتی خود جامعه استفاده می‌شود. دلیل این امر استفاده این فناوری از روش‌های پیشرفته کامپیوتری است که توان پردازش بالایی دارد و نهایتاً آمار فقط می‌تواند نمونه را به جامعه‌ای که از آن انتخاب‌شده تعمیم دهد، درحالی‌که در داده‌کاوی تحلیل‌ها بر روی کل جامعه انجام می‌شود (جهرمی و همکاران، ۲۰۱۶).

### ۳-۱-۲-۲. کاربردهای داده‌کاوی

تاریخچه کشف دانش در پایگاه‌های اطلاعاتی که امروزه به داده‌کاوی مشهور است، قدمت چندانی ندارد. در اوایل دهه ۱۹۹۱، هنگامی که اصطلاح کشف دانش در پایگاه‌های اطلاعاتی برای نخستین بار مطرح شد، هجده‌ای همگانی به سمت طراحی الگوریتم‌های داده‌کاوی صورت پذیرفت (ماربان و همکاران<sup>۳</sup>، ۲۰۰۹) و این‌زمانی بود که شرکت‌ها اقدام به ذخیره‌سازی مقادیر عظیم داده‌ها کردند و به دنبال روش‌هایی برای استفاده از این انبار داده‌ها بودند (نیکلسن<sup>۴</sup>، ۲۰۰۳).  
باتوجه‌به توان تحلیل بالای فناوری داده‌کاوی و باوجود قدرت پردازش بی‌نظیر آن، از این فناوری می‌توان برای حل مسائل بی‌شماری در دنیای واقعی استفاده کرد. برخی از کاربردهای داده‌کاوی عبارت‌اند از:

- تحلیل رفتار افراد و گروه‌ها (ماربان و همکاران، ۲۰۰۹)؛

- پردازش اطلاعات پزشکی (مسروپیان و اوسیانیکوف<sup>۵</sup>، ۲۰۱۴)؛

---

<sup>1</sup> Confirmatory

<sup>2</sup> Exploratory

<sup>3</sup> Marbán et al.

<sup>4</sup> Nicholson

<sup>5</sup> Mesropyan & Ovsyannikov

- تشخیص الگوهای رفتاری مصرف کنندگان (کانتاردزیک<sup>۱</sup>، ۲۰۱۱)؛
- یافتن پروتئین‌های مختلف از نقشه ژنی<sup>۲</sup> موجودات زنده (میترا و آچاریا<sup>۳</sup>، ۲۰۰۳)؛
- هوشمندی کسب‌وکار و کاهش ابهامات ناشی از محیط (ماریناکوس و داسکالکی<sup>۴</sup>، ۲۰۱۶)؛
- مبارزه با جرم و جنایت و تشخیص الگوهای رفتاری گروه‌های تروریستی (عبدالله و همکاران<sup>۵</sup>، ۲۰۱۶) و
- بهینه‌سازی تصمیمات و تخمین‌ها در بازارهای مالی (جین و همکاران<sup>۶</sup>، ۲۰۱۳).

#### ۴-۱-۲-۲. فرآیند داده‌کاوی

باتوجه به این امر که داده‌کاوی فرآیند اکتشاف مدل‌های گوناگون، خلاصه‌ها و ارزش‌های نشأت گرفته از مجموعه خاصی از داده‌ها است، برای پیاده‌سازی چنین فرآیندی باید از روش‌شناسی خاصی استفاده کرد. در این راستا روش‌شناسی فرآیند استاندارد میان صنعتی داده‌کاوی<sup>۷</sup> به‌وسیله تحلیل نمایندگی‌های دایملر کرایسلر ایجاد شد (چاپمن<sup>۸</sup> و همکاران، ۲۰۰۰) این روش‌شناسی توانمند و منعطف جهت ارتقای شایستگی داده‌کاوی در حل مسائل سازمانی است. بر اساس این روش، یک پروژه داده‌کاوی مبتنی بر چرخه عمر، متشکل از شش گام است و این گام‌ها به‌صورت مستمر و تکراری در تمام فرآیند داده‌کاوی به کار گرفته می‌شود. گام‌های روش‌شناسی داده‌کاوی CRISP به شرح زیر است (چاپمن و همکاران، ۲۰۰۰؛ لیو و یو<sup>۹</sup>، ۲۰۱۵):

<sup>1</sup> kantardzic

<sup>2</sup> Genomic Map

<sup>3</sup> Mitra and Acharya

<sup>4</sup> Marinakos & Daskalaki

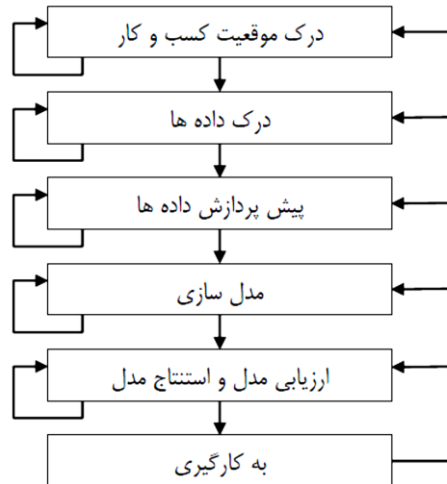
<sup>5</sup> Abdallah et al.

<sup>6</sup> Jain et al.

<sup>7</sup> Cross Industry Standard Process(CRISP)

<sup>8</sup> Chapman

<sup>9</sup> Liu and Yu



نمودار ۱-۲. مراحل داده کاوی CRISP

## ۲-۲-۲. متن کاوی

متن کاوی اکتشاف و استخراج دانش جذاب<sup>۱</sup> و غیر بدیهی<sup>۲</sup> از متن آزاد یا غیر ساختارمند است و تمامی فعالیت‌هایی که به نوعی به دنبال استخراج دانش از متن هستند را شامل می‌شود (هانگ و ژانگ<sup>۳</sup>، ۲۰۱۲؛ لی<sup>۴</sup> و همکاران، ۲۰۰۸). هدف از متن کاوی کشف اطلاعات جدید و ناشناخته، به وسیله استخراج خودکار اطلاعات از داده‌های متنی بدون ساختار یا نیمه ساختاریافته است و کاربردهای گسترده‌ای در تجزیه و تحلیل و پردازش مدارک متنی دارد (بلی<sup>۵</sup>، ۲۰۱۲؛ آبرامسون<sup>۶</sup> و همکاران، ۲۰۱۴؛ وانگ، بلی و هکرمن<sup>۷</sup>، ۲۰۱۲؛ استیورز و گریفیت<sup>۸</sup>، ۲۰۰۷؛ جلودار<sup>۹</sup> و همکاران، ۲۰۱۹).

1 Interesting  
 2 Non-trivial  
 3 Hung & Zhang  
 4 Lee  
 5 Blei  
 6 Abramson  
 7 Wang, Blei & Heckerman  
 8 Steyvers & Griffiths  
 9 Jelodar



پیشرفت پیوسته فناوری باعث افزایش شدید حجم اطلاعات، خصوصاً اطلاعات علمی و تکنیکی شده است. در نتیجه، پیگیری اطلاعات به وظیفه‌ای چالش‌برانگیز برای دانشمندان تبدیل شده است. بنابراین، تکنیک‌های متن‌کاوی برای کمک به دانشمندان در زمینه کسب اطلاعات سودمند از حجم بسیار زیادی از اطلاعات، مورد نیاز است. متن‌کاوی و یا کشف دانش از متن اشاره به فرآیندی است که موجب به دست آوردن الگوهای غیر بدیهی، جالب و باکیفیت بالا و همچنین اطلاعات و دانش از اسناد متنی ساختار نیافته می‌شود. متن‌کاوی (که به عنوان داده‌کاوی متنی<sup>۱</sup> نیز شناخته می‌شود) به جست‌وجو در میان داده‌های متنی برای استخراج اطلاعات مفید می‌پردازد که معمولاً طبیعتی ساختار نیافته دارد (الویدیان<sup>۲</sup> و همکاران، ۲۰۱۵). هدف اولیه متن‌کاوی، بازیابی اطلاعات از متون ساختار نیافته و همچنین ارائه دانش به صورت خالص و چکیده برای کاربران است (آنانیادو و مکنوت<sup>۳</sup>، ۲۰۰۶). هدف از متن‌کاوی قادر ساختن کاربران برای جمع‌آوری، ذخیره، کشف، و تفسیر دانش مورد نیاز برای پژوهش و آموزش مؤثر و نظام‌مند است. متن‌کاوی شامل سه فعالیت اصلی است: بازیابی اطلاعات که بازیابی متون مربوط به پرسش استفاده‌کنندگان است؛ خلاصه اطلاعات که شناختن و استخراج نکات ریز متون که مربوط به پرسش است؛ داده‌کاوی رابطه مستقیم یا غیرمستقیم بین قسمت‌های اطلاعات استخراجی از متون را پیدا می‌کند. به بیان دیگر متن‌کاوی شاخه‌ای از داده‌کاوی یا همان کشف دانش است (توماس، مک‌نات و آنانیادو<sup>۴</sup>، ۲۰۱۱).

هدف متن‌کاوی، کشف اطلاعات ناشناخته‌ای است که هنوز کسی نمی‌داند و بنابراین مستند نشده است (گولاب<sup>۵</sup>، ۲۰۰۶). از آنجا که بسیاری از اطلاعات به شکل متن ذخیره شده‌اند، متن‌کاوی، ارزش اقتصادی بسیار بالایی در پی خواهد داشت. دانش ممکن است از منابع گوناگون اطلاعاتی به دست آمده باشد، اما متون ساختار نیافته، بیشترین منابع دانش در دسترس را تشکیل می‌دهد. مسئله کشف دانش از متون، استخراج مفاهیم صریح و غیرصریح و روابط معنایی میان مفاهیم با استفاده از فنون پردازش زبان طبیعی است. هدف استخراج دانش، به دست آوردن آگاهی درباره کلان داده‌های متنی است. کشف دانش از

---

<sup>1</sup> Text Data Mining

<sup>2</sup> Alwidian

<sup>3</sup> Ananiadou and McNaught

<sup>4</sup> Thomas, McNaught & Ananiadou

<sup>5</sup> Golub

متن ریشه در پردازش زبان طبیعی دارد؛ اما از روش‌های آماری، یادگیری ماشینی، استدلال استخراج اطلاعات، مدیریت دانش و دیگر رشته‌های مرتبط برای فرآیند کشف، وام گرفته شده است. کشف دانش از متن، نقش مهمی در ظهور برنامه‌هایی مانند فهم متن ایفاء می‌کند. متن کاوی با استفاده از پردازش زبان طبیعی، کاربردپذیری کشف دانش از داده‌ها را به‌طور قابل توجهی افزایش داده است. این بدان معناست که نیازی نیست فرآیند کشف دانش از داده‌ها را تنها به آن دسته از اطلاعات موجود در پایگاه‌های ساختاریافته محدود کنیم.

باتوجه به اینکه بیشتر اطلاعات ارزشمند برای استخراج هم‌اکنون در متون زبان طبیعی وجود دارد، پردازش زبان طبیعی می‌تواند فنون موردنیاز برای متن کاوی را فراهم کرده و دانش را به‌طور خودکار از این متون استخراج کند. متن کاوی درباره جستجوی انگاره‌ها در متن زبان طبیعی است و عبارت است از فرآیند تحلیل متن به‌منظور استخراج اطلاعات از حجم عظیمی از متون غیر ساختاریافته. متن کاوی نوعی فناوری است که امکان کشف انگاره‌ها و گرایش‌ها را به‌طور نسبتاً خودکار از درون متن غیر ساختاریافته آزاد فراهم می‌کند. همان‌طور که گفته شد، متن کاوی به‌عنوان امتداد طبیعی داده کاوی به شمار می‌رود و استفاده از همان فناوری‌های داده کاوی در قلمرو اطلاعات متنی است. متن کاوی در مقایسه با داده کاوی فرآیند پیچیده‌تری است، زیرا با آن دسته از داده‌های متنی سروکار دارد که غیر ساختاریافته هستند. متن کاوی قلمرویی چند رشته‌ای است و شامل بازیابی اطلاعات، تحلیل متن، دسته‌بندی اطلاعات، طبقه‌بندی اطلاعات، دیداری‌سازی اطلاعات، فناوری پایگاه‌های اطلاعاتی، یادگیری ماشینی و داده کاوی است. علاوه بر این از ابزارهای متن کاوی برای کمک به پژوهشگران در زمینه‌های مختلف استفاده می‌شود. به‌عنوان مثال، یک اخترشناس که علائیم اشعه ایکس را در ناحیه‌ای از فضا کشف می‌کند ممکن است تمایل داشته باشد که در پیشینه‌های و منابع الکترونیکی جستجو نماید تا ببیند که آیا تاکنون هیچ نوع اشعه مادون قرمز در همان ناحیه کشف شده است یا خیر؟ یک زیست‌شناس که دارای فهرستی از تعداد ۱۹۹ ژن است که در مقالات مختلف شناسایی شده‌اند ممکن است تمایل داشته باشد که به سراغ پژوهش‌های منتشرشده موجود برود و در جستجوی مقالاتی باشد که درباره عملکرد این ژن‌ها توضیح داده‌اند (شارما<sup>۱</sup>، ۲۰۰۵).

---

<sup>1</sup> Sharma

متن کاوی که به تحلیل هوشمند متن، داده کاوی متنی یا کشف دانش در متن، نیز مشهور است عموماً به فرآیند استخراج دانش و اطلاعات موردعلاقه و مهم از مجموعه متنی غیر ساختاریافته اشاره دارد. به عبارت دیگر، متن کاوی فرآیند تحلیل طبیعی متن به منظور کشف و ثبت اطلاعات معنایی برای درونداد و ذخیره سازی در یک ساختار سازمان یافته دانش است. یکی از تکنیک های مفید متن کاوی دسته بندی است که دسته بندی برای کشف توزیع ها و انگاره های داده های موردعلاقه در کلان داده ها مورد استفاده قرار می گیرد. با استفاده از این تکنیک می توان بدون اتکاء بر هیچ دانش پیش زمینه ای، ساختارها یا دسته های موردعلاقه را به طور مستقیم از داده ها شناسایی نمود (پونر - پوراتا<sup>۱</sup> و همکاران، ۲۰۰۷).

از مهم ترین مزایای متن کاوی می توان به موارد زیر اشاره نمود:

- ۱- استخراج اطلاعات مفید و سودمند از کلان داده ها در زمان کوتاه؛
- ۲- پیش بینی جنبه های آینده بر اساس مهیا کردن آمار و ارقام و مشاهدات؛
- ۳- طراحی و ارائه الگوهایی از داده های در دست بررسی؛
- ۴- مشاهده و تحلیل اطلاعات متنی سازمان های امنیتی در منابع اینترنتی و شبکه های اجتماعی؛
- ۵- تحلیل، ذخیره سازی و دسترس پذیری اطلاعات وبسایت ها و موتورهای جست و جو به منظور پردازش و جست و جوی مؤثرتر و با دقت تر
- ۶- تحلیل لغوی و تشخیص الگوها به منظور مطالعه توزیع فرکانس واژگان (هاشمی<sup>۲</sup> و همکاران، ۲۰۱۵؛ ربهولز-شومان<sup>۳</sup> و همکاران، ۲۰۱۲؛ راجمان و وسلی<sup>۴</sup>، ۲۰۰۴؛ سالوم<sup>۵</sup> و همکاران، ۲۰۱۷؛ اوکالهان<sup>۶</sup> و همکاران، ۲۰۱۵).

## ۱-۲-۲-۲. رابطه متن کاوی و داده کاوی

---

<sup>1</sup> Pons-Porrata

<sup>2</sup> Hashimi

<sup>3</sup> Rebholz-Schuhmann

<sup>4</sup> Rajman and Vesely

<sup>5</sup> Salloum

<sup>6</sup> O'callaghan

متن کاوی شامل طیف گسترده‌ای از وظایف است که می‌تواند اطلاعاتی را در مورد جنبه‌های مختلف متون به ارمغان

آورد. وظایف معمول استخراج متن شامل موارد زیر است: (ماینر<sup>۱</sup> و همکاران، ۲۰۱۲):

- طبقه‌بندی اسناد<sup>۲</sup> - اختصاص سند به یک یا چند طبقه از پیش تعریف شده (به‌عنوان مثال، اختصاص مقاله روزنامه به یک یا

چند طبقه، برچسب‌گذاری نامه‌های الکترونیکی به‌عنوان هرزنامه)؛

- خوشه‌بندی<sup>۳</sup> - گروه‌بندی اسناد با توجه به شباهت آن‌ها، به‌عنوان مثال، به‌منظور شناسایی اسنادی که دارای یک موضوع

مشترک هستند؛

- جمع‌بندی<sup>۴</sup> (خلاصه‌سازی) - یافتن مهم‌ترین قسمت‌ها در یک یا چند سند و ایجاد متنی که به‌طور قابل توجهی کوتاه‌تر از

سند اصلی باشد؛

- بازیابی اطلاعات<sup>۵</sup> - بازیابی اسنادی که با کوئری مطابقت دارند و اطلاعات موردنیاز را از مجموعه بزرگی از اسناد نشان

می‌دهد؛

- استخراج معنی<sup>۶</sup> - اسناد یا قسمت‌هایی از آن را با شناسایی موضوعات پنهان، تجزیه و تحلیل احساسات، عقاید نشان می‌دهد؛

- استخراج اطلاعات<sup>۷</sup> - استخراج اطلاعات ساختار یافته مانند موجودیت‌ها، رویدادها یا روابط از متون بدون ساختار؛

- استخراج روابط<sup>۸</sup> - یافتن ارتباط بین مفاهیم یا اصطلاحات در متون؛

- تحلیل روند<sup>۹</sup> - بررسی چگونگی تغییر مفاهیم موجود در اسناد در گذر زمان؛

- ترجمه ماشینی<sup>۱۰</sup> - تبدیل متن نوشته‌شده از یک زبان به زبان دیگر.

---

<sup>1</sup> Miner

<sup>2</sup> categorization of documents

<sup>3</sup> clustering

<sup>4</sup> summarization

<sup>5</sup> information retrieval

<sup>6</sup> extracting the meaning

<sup>7</sup> information extraction

<sup>8</sup> association mining

<sup>9</sup> trend analysis

<sup>10</sup> machine translation

داده کاوی فرآیند خودکار یا نیمه خودکار یافتن دانش ضمنی، قبلاً ناشناخته و بالقوه مفید در مجموعه داده‌های ذخیره شده الکترونیکی است. دانش دارای شکلی از الگوهای ساختاری در داده‌ها است که می‌تواند برای پیش‌بینی یا ارائه پاسخ در آینده نیز مورد استفاده قرار گیرد (ویتن و فرانک<sup>۱</sup>، ۲۰۰۲). داده کاوی شامل روش‌ها، ابزارها، الگوریتم‌ها یا مدل‌های مختلف است. جهت انجام عملیات داده کاوی نیاز هست که داده‌ها به شکل ساختاریافته باشند. این بدان معنی است که داده‌ها را می‌توان به صورت جدول یا یک پایگاه داده رابطه‌ای نشان داد. داده‌ها به شکل مجموعه‌ای از مثال‌ها (نمونه‌ها) با مقادیر خاص ویژگی‌های آن‌ها (ویژگی‌ها، متغیرها، زمینه‌ها) توصیف می‌شود. ویژگی‌ها می‌تواند انواع مختلفی داشته باشد که عبارتند از (دانگ و لیو<sup>۲</sup>، ۲۰۰۴؛ لیو و متودا<sup>۳</sup>، ۲۰۱۲):

- طبقه‌ای (اسمی)<sup>۴</sup> - دامنه مجموعه‌ای گسسته از مقادیر است که در آن ترتیب معنی ندارد؛

- دو دویی<sup>۵</sup> - نوع خاصی از ویژگی‌های طبقه‌ای که فقط با دو مقدار ممکن است؛

- ترتیبی<sup>۶</sup> - دامنه یک مجموعه گسسته از مقادیر است که می‌تواند دارای یک ترتیب خاص باشد؛

- عددی<sup>۷</sup> - مقدار یک عدد که این عدد صحیح یا پیوسته است.

یک واحد از متن می‌تواند یک جمله، چند جمله باهم ترکیب شده در یک پاراگراف یا متن‌های بسیار طولانی‌تری مانند صفحات وب، ایمیل، مقاله یا کتاب باشد. گاهی اوقات، یک متن می‌تواند فقط چند کلمه باشد که جمله معتبری نیست که کاملاً معمول است، به عنوان مثال، برای پست‌های کوتاه در شبکه‌های اجتماعی، برخی قوانین (نحو) ترکیب می‌شوند؛ بنابراین برای اینکه بتوان از روش‌های داده کاوی در متن استفاده کرد، باید آن‌ها را به نمایشی ساختاری تبدیل کرد.

---

<sup>1</sup> Witten and Frank

<sup>2</sup> Dong and Liu

<sup>3</sup> Liu H, Motoda

<sup>4</sup> categorical (nominal)

<sup>5</sup> binary

<sup>6</sup> ordinal

<sup>7</sup> numerical

به‌طور کلی می‌توان گفت که داده‌کاوی با داده‌های ساختاریافته و نرمال‌شده سروکار دارد و با پایگاه داده‌های رابطه‌ای کار می‌کند؛ اما در عوض متن‌کاوی با داده‌های ساختاریافته یا نیمه‌ساختاریافته یعنی متون موجود در مقالات، اسناد و غیره سروکار دارد. علاوه بر این دسترسی کم به ساختار در متون دلیلی دیگر بر این است که متن‌کاوی کاری بسیار مشکل است. مفاهیم موجود در متون معمولاً بسیار انتزاعی بوده و به‌سختی می‌توان آن‌ها را مدل‌سازی موضوعی کرد. همچنین وقوع کلمات مترادف (کلمات متفاوت از نظر نوشتاری اما هم‌معنی) یا کلمات با تلفظ یکسان (کلمات با تلفظ یکسان اما معنی متفاوت) پیدا کردن رابطه منطقی بین بخش‌های مختلف متن را مشکل می‌کند (اسپیناکیس و پرسترا، ۲۰۰۴).

فرق داده‌کاوی با متن‌کاوی این است که در متن‌کاوی الگوها از متن زبان طبیعی استخراج می‌شوند درحالی‌که در داده‌کاوی الگوها از پایگاه‌های داده ساخت‌یافته به دست می‌آید. متن‌کاوی، متصل کردن اطلاعات استخراج‌شده به یکدیگر برای تشکیل حقایق یا فرضیه‌های جدید است تا پس از آن به کمک روش‌های متعارف آزمایش، بررسی بیشتری شوند (لیو و متودا، ۲۰۱۲). به‌عبارت‌دیگر تفاوت میان داده‌کاوی و متن‌کاوی در این است که داده‌کاوی، اطلاعات را گردآوری و فهرست‌نویسی کرده، سپس اقدام به تولید دانش از بین حجم عظیمی از داده‌ها می‌کند اما متن‌کاوی، حوزه‌ای نو و میان‌رشته‌ای است که از رشته‌های بازبایی اطلاعات، داده‌کاوی، یادگیری ماشینی، آمار و زبان‌شناسی محاسباتی مشتق شده است و عمدتاً بر مستندات متنی تکیه دارد (رضضانی و همکاران، ۱۳۹۳).

## ۲-۲-۲-۲. تاریخچه متن‌کاوی

آنچه امروزه به‌طور عام متن‌کاوی نامیده می‌شود، مجموعه‌ای متشکل از علوم مختلف از قبیل زبان‌شناسی، آمار، کامپیوتر، مدیریت، هوش مصنوعی و دیگر حوزه‌ها است. در اصل باید ریشه اصلی روش‌های متن‌کاوی را در تلاش‌ها و اقدامات صورت گرفته برای استفاده از فنون کمی و آماری در علم زبان‌شناسی جست که به این مجموعه اقدامات،

---

<sup>1</sup> Spinakis and Peristera

زبان‌شناسی کمی<sup>۱</sup> می‌گویند. تاریخچه استفاده از زبان‌شناسی کمی به حداقل قرن نوزدهم میلادی بازمی‌گردد. اما فعالیت کلاسیک تئوری گونه‌زیم ۱۹۴۹ به‌عنوان یکی از مهم‌ترین پیش‌گامان عرصه تحلیل کمی زبان‌شناسی، شناخته می‌شود (گلیسون<sup>۲</sup> و همکاران، ۲۰۰۵). از دهه هفتاد میلادی تاکنون، افزایش چشمگیری در میزان علاقه به این حوزه از علوم اطلاعات مشاهده شده است. پژوهش ویلی یکی از نخستین کاربردهای این فنون در مطالعه منابع و ادبیات علمی، محسوب می‌شود (جانسنز<sup>۳</sup> و همکاران، ۲۰۰۶).

نخستین بار لان<sup>۴</sup> در سال ۱۹۵۸ میلادی، مفهوم متن‌کاوی را در مقاله خود مطرح نمود. سپس در سال ۱۹۶۱ نیز دویله در مقاله‌ای به متن‌کاوی و روش‌های مرتبط با آن اشاره کرد و بیان داشت که «طبقه‌بندی و سازمان‌دهی اطلاعات» می‌تواند از تجزیه و تحلیل تکرار و توزیع کلمات به‌کاررفته در آن صورت پذیرد. اگر این دو مورد را به‌عنوان نخستین موارد مطرح‌شدن متن‌کاوی در نظر بگیریم، می‌توان گفت که متن‌کاوی ممکن است مفهومی جدید باشد، اما رویای استخراج خودکار اطلاعات از متون، هم‌زمان با پیدایش کامپیوتر به‌وجود آمده است. سوانسون<sup>۵</sup> در سال ۱۹۹۱ در مقاله خود بیان کرد که متون علمی باید به‌عنوان پدیده‌های ارزشمند طبیعی، کشف، اصلاح و تحلیل شوند و به‌این ترتیب نظر دانشمندان به استفاده از اطلاعات با تحلیل هوشمندانه را جلب کرد.

سوانسون با طراحی نرم‌افزاری نخستین گام در جهت عملی‌سازی متن‌کاوی را طی کرد. او از این نرم‌افزار در استخراج اطلاعات مفید از متون پزشکی استفاده نمود. این نرم‌افزار ارو اسمیس نام دارد. این نرم‌افزار که به‌صورت تخصصی در ارتباط با پردازش متون پزشکی تهیه شده قادر است پس از دریافت متون، کلمات کلیدی و متون مرتبط با یکدیگر را مشخص کند. سوانسون هیچ‌گاه از اصطلاح متن‌کاوی برای این نرم‌افزار استفاده نکرد، اما به نظر می‌رسد که این نرم‌افزار نخستین نرم‌افزار متن‌کاوی بوده است. از این رو می‌توان او را پدر متن‌کاوی مدرن نامید.

---

<sup>1</sup> Quantitative linguistics

<sup>2</sup> Glenisson

<sup>3</sup> Janssens

<sup>4</sup> Luhn

<sup>5</sup> Swanson

لیندزی<sup>۱</sup> و گوردون<sup>۲</sup> ۱۹۹۹ کارهای سوانسون را بدون آنکه نام متن کاوی بر آن نهند، ادامه دادند. نرم افزار آن‌ها دو واژه را به صورت هم‌زمان در میان متون جستجو می‌کرد و نتیجه را در فهرستی قرار می‌داد تا کاربران آن‌ها را به عنوان مقاله‌های تکمیلی مورد مطالعه قرار دهند. لیندزی و گوردون در همان سال روش TF-IDF<sup>۳</sup> را برای رتبه‌بندی متون و واژه‌ها به این نرم‌افزار افزودند.

هیرست<sup>۴</sup> (۱۹۹۷) نیز متن کاوی را با دسترسی به اطلاعات (بازیابی سنتی اطلاعات) متمایز ساخته است. بازیابی سنتی اطلاعات، بیشتر روی بازیابی متون مرتبط باهم تأکید می‌کند. متونی که با نیازهای اطلاعاتی کاربر مرتبط است. بر اساس تعریف هیرست داده کاوی- که متن کاوی نیز نوعی از داده کاوی به شمار می‌آید- تنها با اطلاعات و بازیابی آن سروکار ندارد، بلکه تلاش می‌کند تا اطلاعات جدیدی را از میان داده‌ها کشف نماید که پیش‌ازین حتی برای ایجادکننده داده‌ها هم مشخص نبود. او معتقد بود که داده کاوی و متن کاوی مبتنی بر شانس هست، درحالی‌که بازیابی اطلاعات مبتنی بر هدف است؛ بنابراین کارهایی نظیر جستجوی واژه‌ها برای پاسخ به پرسش‌ها، متن کاوی محسوب نمی‌شود. او بر این باور بود که بازیابی و دستیابی به اطلاعات می‌تواند به عنوان کاری تکمیلی و پشتیبان برای متن کاوی به شمار رود.

نخستین محصول نرم‌افزاری<sup>۵</sup> در زمینه متن کاوی در سال ۱۹۹۹ توسط شرکت آی بی ام<sup>۶</sup> به بازار عرضه گردید. نرم‌افزار مذکور شامل مجموعه ابزارهایی است که به استخراج اطلاعات از متن پرداخته و موجب غنای متن می‌شود. اطلاعاتی که از محتوای متن استخراج می‌شود می‌تواند ویژگی‌های متن نظیر زبان متن، اسامی افراد، تاریخ‌ها، مقادیر پول و ... باشد. استخراج این ویژگی‌ها به صورت خودکار بوده و بر اساس واژه‌نامه از پیش تعریف‌شده‌ای انجام نمی‌شود. این نرم‌افزار قادر است تا روی یک متن یا مجموعه‌ای از متون پردازش‌های لازم را انجام دهد. در این نرم‌افزار که شمارش واژه‌های مورد استفاده در متن،

---

<sup>1</sup> Lindsay

<sup>2</sup> Gordon

<sup>3</sup> Term Frequency–Inverse Document Frequency

<sup>4</sup> Hearst

<sup>5</sup> Intelligent Miner for Text

<sup>6</sup> IBM



اساس پردازش‌های بعدی آن است قسمتی برای تشخیص اصطلاح‌ها دارد و متون را خوشه‌بندی و دسته‌بندی می‌کند (در<sup>۱</sup> و همکاران، ۱۹۹۹).

با توسعه مفهوم متن کاوی، مفاهیم دیگری نیز پایه‌پای آن رشد کرد: بازیابی اطلاعات از مجموعه متون، بازیابی اطلاعات از یک متن، کشف دانش از بانک‌های اطلاعاتی، مدیریت دانش در سازمان‌ها و مصورسازی داده‌ها و اطلاعات، این مفاهیم توسط کوستوف<sup>۲</sup>، اوراد<sup>۳</sup> و لوزیویچ<sup>۴</sup> در سال ۲۰۰۰ در چند مقاله منتشر شد تا میان این مفاهیم تمایز قائل شوند. در سال ۲۰۰۱ میلادی کوستوف و دمارکو<sup>۵</sup> متن کاوی را با «استخراج اطلاعات از متون فنی» تعریف کردند. بر اساس این تعریف، متن کاوی شامل سه بخش است: بازیابی اطلاعات، پردازش اطلاعات و یکپارچگی اطلاعات. پردازش اطلاعات به استخراج الگوهای موجود در متون بازیابی شده اطلاق می‌شود و یکپارچگی اطلاعات ترکیبی از اطلاعات پردازش شده توسط کامپیوتر است که انسان اطلاعات را بازیابی نموده و مطالعه می‌کند. این مرحله مفهوم سیستم انسان- ماشین را تداعی می‌کند.

### ۲-۲-۲-۳. روش‌شناسی متن کاوی

متن کاوی فرآیند نیمه خودکاری است که برای استخراج الگوها و کشف دانش از روی انبوهی از منابع داده‌ای غیر ساختاریافته مورد استفاده قرار می‌گیرد (دلن<sup>۶</sup>، ۲۰۱۴). متن کاوی با داده کاوی رابطه تنگاتنگی دارند و هدف کلی و برخی فرآیندهای عملیاتی آن‌ها یکسان است. همچنین در متن کاوی ورودی فرآیندها مجموعه‌ای از فایل‌های متنی غیر ساختاریافته بوده که از آن جمله می‌توان به فایل‌های Word، PDF، گزیده‌ای از متون و فایل‌های XML اشاره نمود. مزایای متن کاوی بیشتر در قلمروهایی است که انبوهی از داده‌های متنی قرار دارد. از جمله مهم‌ترین این قلمروها می‌توان به

---

<sup>1</sup> Dörre

<sup>2</sup> Kostoff

<sup>3</sup> Oard

<sup>4</sup> Losiewicz

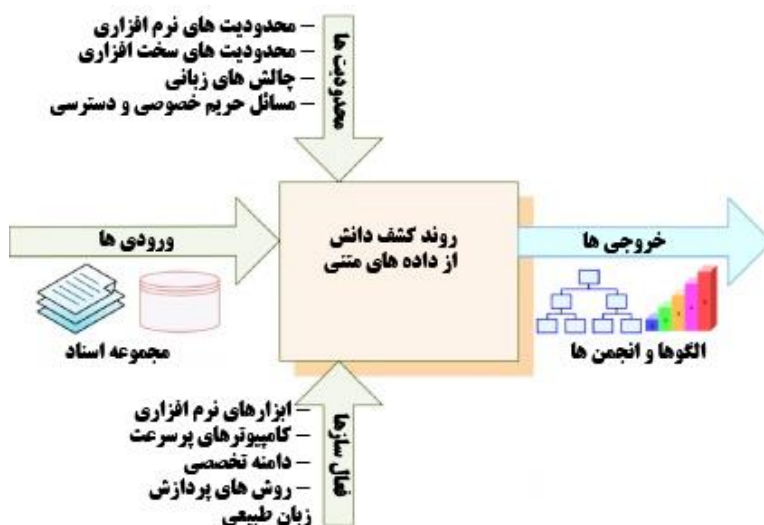
<sup>5</sup> DeMarco

<sup>6</sup> Delen

بخش مرور پیشینه مقالات مختلف، بخش مالی (گزارش‌های فصلی، تفسیرهای رسانه‌ای)، خدمات بهداشتی (گزارش مرخصی بیمار، یادداشت‌های پزشک)، حقوق و مسایل قانونی (دستورات دادگاه)، زیست‌شناسی (فعل و انفعالات مولکولی)، فناوری (فایل‌های پروانه ثبت اختراع) و بازاریابی (نظرات مشتریان) اشاره نمود (کیم و دلن، ۲۰۱۸).

#### ۲-۲-۲-۴. فرآیند متن کاوی

مطالعات متن کاوی با فرآیند اصولی و کارآمد و با بهره‌گیری از بهترین رویه‌های کاربردی بنا نهاده می‌شود. هر پروژه‌ی متن کاوی به فرآیند استاندارد و موردپذیرش مجریان این قلمرو، همچون فرآیند استاندارد صنعتی Cross برای داده کاوی (CRISP-DM) نیاز دارد. در سطوح بالاتر فرآیند متن کاوی از طریق نمودار نمایش داده می‌شود نمونه چنین نموداری به همراه اجزاء تشکیل دهنده شامل ورودی‌ها، خروجی‌ها، کنترل‌گرها و سازوکارها (توانمندسازی‌ها) با پیکان‌های جهت‌دار در نمودار ۲-۲-۲ نمایش داده شده‌اند (کیم و دلن، ۲۰۱۸).



نمودار ۲-۲-۲- فرآیند متن کاوی در قالب یک مدل مفهومی کلی (کیم و دلن، ۲۰۱۸)

نمودار در شرایط مختلف به سه مرحله تقسیم می‌شود. هر یک از این مراحل دارای ورودی‌های خاصی بوده که برای تولید خروجی‌های اختصاصی بکار گرفته می‌شود. در صورتی که به هر دلیلی خروجی بخشی از این فرآیندها با انتظارات سازگاری نداشته باشد، اعمال بازخورد (بازگشت رو به عقب) نسبت به وظیفه‌ی در دست اجرا ضروری خواهد بود. در نمودار ۲-۲ نمای گرافیکی از این وضعیت نمایش داده شده است.

مرحله نخست. تعیین بخش‌های منتخب از متن: هدف اصلی فعالیت نخست گردآوری تمامی اسناد مرتبط با قلمروی در دست مطالعه است. در ادامه‌ی کار اسناد گردآوری شده به قالب خاصی تبدیل و در همان شرایط بایگانی شده تا همگی در ساختاری واحد تحت پردازش‌های رایانه‌ای متعاقب قرار بگیرند (برای مثال فایل‌های متنی در فرمت ASCII)؛

مرحله دوم. پیش‌پردازش داده‌ها (تشکیل ماتریس واژه / تفکیک برحسب سند): پس از نهایی سازی وضعیت متون، ماتریس واژه‌های تفکیک شده برحسب سند (TDM) بر اساس اطلاعات در دسترس تشکیل می‌شود. در این ماتریس سطرها نشان‌دهنده اسناد و ستون‌ها بیانگر واژه‌ها (کلیدواژه‌های مدنظر، هدف، جستجو شده) است. روابط برقرار شده بین واژه‌ها و اسناد به کمک شاخص‌ها نمایش داده می‌شوند (برای مثال شاخص‌های نسبی که ساده‌ترین آن‌ها تعداد تکرارپذیری هر واژه در اسناد مربوطه است)؛

مرحله سوم. استخراج دانش: هنگامی که به ماتریس TDM با ساختاردهی بهینه دسترسی باشد، به استخراج الگوها مبادرت ورزیده و آن‌ها به‌عنوان خوشه‌ها تعریف می‌شود. خوشه‌بندی یک فرآیند غیر نظارتی بوده که طی آن اشیاء یا رخدادها در گروه‌بندی‌های نرمال (معنادار) دسته‌بندی می‌شوند. فرآیند غیر نظارتی فرآیندی است که از هیچ‌گونه الگو یا دانش قبلی برای هدایت فرآیند خوشه‌بندی سود نخواهد بُرد. در فرآیند خوشه‌بندی غیر نظارتی، مجموعه اشیاء / مؤلفه‌های فاقد برچسب‌گذاری (اسناد، دیدگاه‌های مشتری، صفحات وب) بدون هرگونه دانش قبلی به خوشه‌های معنادار منتقل می‌شوند. فرضیه بنیادی بیانگر این نکته بوده که اسناد مرتبط کاملاً شبیه یکدیگر هستند، این در حالی است که اسناد غیر مرتبط از کمترین شباهت نسبت به یکدیگر سود می‌برند. در صورتی که این فرضیه معتبر باشد، آنگاه خوشه‌بندی اسناد بر مبنای میزان شباهت محتوی به ارتقای کیفیت جستجو می‌انجامد.

یافتن تعداد بهینه‌ای از خوشه‌ها به هیچ‌عنوان وظیفه‌ای ساده نخواهد بود. در واقع هیچ‌گونه فرمول ریاضیاتی برای انجام این عمل پیش‌بینی نشده است (یک الگوریتم با ساختار بسته). تعیین تعداد بهینه‌ای از خوشه‌ها همچنان در قلمروی فرآیندهای شهودی-آزمایشی قرار گرفته و طی آن تعداد خوشه‌ها به تدریج از تعداد کم تا تعداد بیشتر افزایش یافته (یا برعکس) تا زمانی که تعداد خوشه‌های حاصله به شیوه‌ای بهینه بیانگر مجموعه داده‌ی چندبعدی در دست ارزیابی باشند (کیم و دلن، ۲۰۱۸).

### ۳-۲-۲. خوشه‌بندی

در سال‌های اخیر، به دلیل پیشرفت‌هایی که در داده‌کاوی، قدرت رایانه‌ها و بسته‌های نرم‌افزاری آماری حاوی الگوریتم‌های تحلیل خوشه، صورت گرفته است، تحلیل خوشه به موضوعی مهم مبدل شده است. اغلب نیاز است تا مجموعه‌ای از اسناد به گروه‌ها یا خوشه‌هایی همگن تقسیم‌بندی شوند. اگر این مجموعه حاوی تعداد اندکی مستند باشد، این کار می‌تواند به صورت دستی انجام گیرد. ولی اگر با حجم زیادی از اسناد سروکار داشته باشیم، انجام دستی این فرآیند، زمان‌بر و غیر مؤثر، خواهد بود. داده و الگو یکی از شاخص‌های بسیار مهم در دنیای اطلاعات است. خوشه‌بندی یکی از بهترین روش‌هایی است که برای کار با داده‌ها ارائه شده است. قابلیت آن در ورود به فضای داده و تشخیص ساختار آن‌ها، خوشه‌بندی را یکی از ایده‌آل‌ترین سازوکارها برای کار با دنیای عظیم داده‌ها کرده است. ایده آن نخستین بار در دهه ۱۹۳۰ ارائه شد و امروزه با پیشرفت‌ها و جهش‌های عظیمی که پدید آمده، خوشه‌بندی در کاربردها و جنبه‌های مختلفی حضور یافته است. هدف نهایی خوشه‌بندی این است که داده‌های موجود را به چند گروه تقسیم کنند و در این تقسیم‌بندی داده‌های گروه‌های مختلف باید حداکثر تفاوت ممکن را به هم داشته باشند و داده‌های موجود در یک گروه باید بسیار به هم شبیه باشند. البته کیفیت نتایج خوشه‌بندی به روش اندازه‌گیری شباهت و توانایی و قدرت الگوریتم در کشف الگوهای مخفی میان داده‌ها بستگی دارد (بورگلاند<sup>۱</sup>، ۲۰۱۳).

---

<sup>1</sup> Borglund

خوشه‌بندی یکی از کاربردهای متن‌کاوی است. خوشه‌بندی به فرآیند تقسیم مجموعه‌ای از داده‌ها (یا اشیاء) به زیر کلاس‌هایی با مفهوم خوشه اطلاق می‌شود. به این ترتیب یک خوشه شامل یک سری داده‌های مشابه است که همانند یک گروه واحد، رفتار می‌کنند (بورگلاند، ۲۰۱۳؛ هانگ<sup>۱</sup>، ۲۰۰۸). تکنیک خوشه‌بندی تکنیکی قابل‌اعتماد است که به‌طور کلی برای تحلیل کلان داده‌ها مورد استفاده قرار می‌گیرد. ثابت شده است که خوشه‌بندی متن یکی از مؤثرترین ابزارهای مورد استفاده برای تحلیل عنوان‌ها است (کلیفتون<sup>۲</sup> و همکاران، ۲۰۰۴). هر مجموعه که موجودیت نامیده می‌شود، با یک خوشه مشخص می‌شود که مربوط به یکی از عنوان‌ها در مجموعه نوشته‌هاست (سالوم و همکاران، ۲۰۱۸).

خوشه‌بندی به عنوان یکی از روش‌های داده‌کاوی توصیفی، تکنیکی برای گروه‌بندی مشاهدات به  $K$  گروه (خوشه) مختلف است؛ به طوری که داده‌های هر خوشه بالاترین درجه شباهت را داشته باشند و داده‌های متعلق به خوشه‌های مختلف از حداکثر درجه عدم شباهت برخوردار باشند (کافمن و راسو<sup>۳</sup>، ۲۰۰۹). خوشه‌بندی یکی از تکنیک‌هایی است که به‌طور گسترده به منظور متن‌کاوی، شناسایی الگو، تحلیل صفحات وب و تحلیل بازار مورد استفاده قرار می‌گیرد. در تعریفی دیگر، خوشه‌بندی برای جدا کردن یک جمعیت غیر همگن، به تعدادی از زیرگروه‌های همگن، بدون طبقه‌های از پیش تعریف شده، استفاده می‌شود. گروه‌ها با توجه به شباهت مفاهیم و موجودیت‌های درون خوشه‌ها و بر اساس متغیرهایی مشخص و محتوای تحلیل خوشه‌بندی می‌شوند (تراپی<sup>۴</sup> و همکاران، ۲۰۱۱). مویی<sup>۵</sup> و سارستد<sup>۶</sup> ۲۰۱۱ مراحل شش‌گانه‌ای برای انجام فرآیند خوشه‌بندی معرفی کرده‌اند که در نمودار ۲-۳ نشان داده شده است.

---

<sup>1</sup> Huang

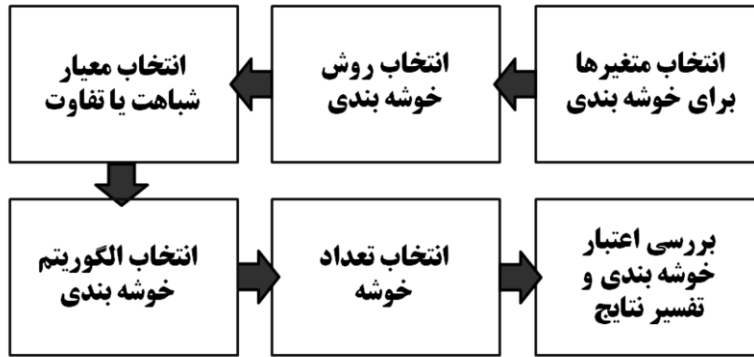
<sup>2</sup> Clifton

<sup>3</sup> Kaufman and Rousseeuw

<sup>4</sup> Trappey

<sup>5</sup> Mooi

<sup>6</sup> Sarstedt



نمودار ۳-۲. مراحل شش گانه خوشه بندی (مویی و سارستد، ۲۰۱۱)

### ۱-۳-۲-۲. مراحل فرآیند خوشه بندی

- **انتخاب متغیرها برای خوشه بندی.** الگوریتم های خوشه بندی نتایج را فارغ از اینکه داده ها پیش فرض های لازم را برای خوشه بندی دارا باشند در اختیار کاربر قرار می دهد؛ بنابراین انتخاب متغیرهای صحیح مرتبط برای کسب نتایج معتبر بسیار حائز اهمیت است.
- **انتخاب روش خوشه بندی.** در این مرحله به طور معمول بهینه کردن مشخصه ای مانند به حداقل رساندن واریانس داخل خوشه ها یا حداکثر کردن فاصله بین خوشه ها انجام می شود.
- **انتخاب معیار شباهت یا تفاوت.** گام بعد برای انجام فرآیند خوشه بندی، انتخاب معیار تفاوت یا تشابه است تا بتوان فاصله دو مشاهده را به صورت یک مقدار عددی مشخص نمود (کافمن و راسو<sup>۱</sup>، ۲۰۰۹). معیارهای متفاوتی برای خوشه بندی مجموعه داده ها وجود دارد که به ماهیت داده ها (کمی، اسمی و غیره) وابسته است (اوریت<sup>۲</sup> و همکاران، ۲۰۱۱).

<sup>1</sup> kaufman & rasoi

<sup>2</sup> Everitt

- انتخاب الگوریتم خوشه‌بندی. بسته به نوع خوشه‌بندی (سلسله مراتبی، تقسیمی و غیره) از الگوریتم‌های مختلفی همچون الگوریتم پیوند، پیوند کامل، پیوند متوسط، الگوریتم مرکز، الگوریتم K-Means و غیره استفاده می‌شود.
- **انتخاب تعداد خوشه.** معیارهای اعتبار مختلفی برای ارزیابی خوشه‌بندی و انتخاب تعداد خوشه بر اساس محاسبات آماری وجود دارد. این معیارها بسته به روش خوشه‌بندی متفاوت است.
- **بررسی اعتبار خوشه‌بندی و تفسیر نتایج.** بررسی اعتبار خوشه‌بندی برای جلوگیری از ایجاد نتایج نامعتبر و خوشه‌های غیرحقیقی لازم است. توصیه می‌شود که پژوهشگران الگوریتم‌های مختلف خوشه‌بندی را برای داده‌های خود انجام داده و نتایج حاصل را بر اساس معیارهای اعتبار مقایسه نمایند تا با اطمینان بیشتری خوشه‌ها را انتخاب کنند. همچنین توصیه می‌شود تا خوشه‌بندی با تعداد خوشه‌های مختلف صورت گرفته و نتایج مقایسه شوند. استفاده از نظرات خبرگان حوزه موردنظر نیز در بررسی اعتبار نتایج خوشه‌بندی مفید خواهد بود (اوریت<sup>۱</sup> و همکاران، ۲۰۱۱).

#### ۲-۲-۴. مدل‌سازی موضوعی

به منظور برخورداری از روشی بهتر در جهت مدیریت بهینه اسناد الکترونیکی، استفاده از روش‌ها یا ابزارهایی که به‌طور خودکار عمل سازمان‌دهی، جستجو، نمایه‌سازی و مرور مجموعه‌های عظیم را انجام می‌دهد، ضروری است. براساس یافته‌های پژوهش‌های کنونی در مورد یادگیری ماشین و آمارها، تکنیک‌های جدیدی برای یافتن الگوهای واژه‌ها در مجموعه اسناد با استفاده از مدل‌های احتمالی سلسله مراتبی، توسعه یافته‌اند. این مدل‌ها، «مدل‌های موضوعی» نامیده می‌شود. کشف الگوها به‌طور معمول موضوعات اساسی به‌منظور شکل‌دادن اسنادی را منعکس کرده و به‌آسانی به سایر انواع داده‌ها تعمیم داده می‌شود. مدل‌های موضوعی به‌جای تحلیل واژه‌ها، تصاویر، داده‌های بیولوژیکی و اطلاعات و داده‌های نظرسنجی را تحلیل می‌کند (دیوید و همکاران، ۲۰۰۶).

اهمیت اصلی مدل‌سازی موضوعی، کشف الگوهای کاربرد واژه‌ها و چگونگی ارتباط با اسنادی است که الگوهای مشابه را ارائه می‌دهند (استیورز<sup>۱</sup> و همکاران، ۲۰۰۷). در مدل‌سازی موضوعی هر سند موجود در یک مجموعه توسط هیستوگرامی که دربرگیرنده میزان وقوع واژه‌ها است، نمایش داده می‌شود. هیستوگرام با توزیع روی تعداد موضوع‌های معینی مدل‌سازی شده است که هر یک از آن‌ها توزیعی روی کلمات موجود در واژه‌نامه‌اند. با درک توزیع‌ها، می‌توان از هر سند، هیستوگرامی را ترسیم نمود. مدل‌های موضوعی مختلف از جمله تحلیل معنایی پنهان<sup>۲</sup>، تحلیل معنایی پنهانی احتمالاتی<sup>۳</sup>، تخصیص پنهان دیریکله، مدل موضوعی هم‌بسته<sup>۴</sup> طبقه‌بندی دقیقی در قلمرو کشف مدل‌سازی موضوعی ارائه داده‌اند (هافمن<sup>۵</sup>، ۲۰۰۱). باگذشت زمان، موضوعات موجود در مجموعه اسناد رشد داشته و باید گفت مدل‌سازی موضوعی بدون در نظر گرفتن متغیر زمان باعث ابهام در کشف موضوع می‌شود. مدل‌سازی موضوعی با لحاظ نمودن متغیر زمان، مدل‌سازی تکاملی موضوعی نامیده می‌شود. این مدل‌سازی اطلاعات پنهان و مهم در مجموعه اسناد را افشاء نموده و امکان شناسایی موضوعات را با توجه به زمان و بررسی تکاملی‌شان در طول زمان ارائه می‌دهد. قلمروهای مختلفی می‌توانند از مدل‌های تکاملی موضوعی استفاده کنند.

#### ۱-۴-۲. روش‌های مدل‌سازی موضوعی

در این بخش، در خصوص برخی از روش‌های مدل‌سازی موضوعی که با کلمات، اسناد و موضوعات سروکار دارند، بحث می‌شوند. به‌علاوه، ایده کلی این روش‌ها و مثال‌هایی در صورت لزوم ارائه شده است. این روش‌ها، کاربردهای زیادی دارند در اینجا چگونگی کاربرد هر روش به‌طور مختصر بیان می‌گردد.

---

<sup>1</sup> Steyvers

<sup>2</sup> Latent Semantic Analysis (LSA)

<sup>3</sup> Probabilistic Latent Semantic Analysis (PLSA)

<sup>4</sup> Correlated Topic Model (CTM)

<sup>5</sup> Hofmann



## • تحلیل معنایی پنهان

تحلیل معنایی پنهان (LSA) روش یا تکنیک مربوط به پردازش زبان طبیعی است. هدف اصلی تحلیل معنایی پنهان، ایجاد بردار بازنمایی برای متون و نشانه‌گذاری محتوای معنایی است. با استفاده از بازنمایی برداری، شباهت میان متون محاسبه می‌شود. در گذشته، تحلیل معنایی پنهان، نمایه‌سازی معنایی پنهان<sup>۱</sup> (LSI) نام داشت و برای بازیابی اطلاعات، ارائه شده بود. تحلیل معنایی پنهان بایستی ابعادی برای روشی مانند تطبیق کلمات کلیدی، تطبیق وزن کلمات کلیدی و بازنمایی بردار مربوط به رخداد کلمات در اسناد، داشته باشد (سیلوا و دی کاروالهو<sup>۲</sup>، ۲۰۲۱). همچنین، تحلیل معنایی پنهان از تجزیه مقادارهای تکی<sup>۳</sup> (SVD) برای مرتب‌سازی مجدد داده‌ها استفاده می‌نماید. تجزیه مقادارهای تکی روشی است که از یک ماتریس برای پیکربندی مجدد و محاسبه نقصان‌های فضای بردار، استفاده می‌شود. نقصان‌های فضای بردار از بالاترین تا کمترین اهمیت سازمان‌دهی می‌شوند. برای توصیف اساسی‌ترین مراحل تحلیل معنایی پنهان ابتدا، مجموعه عظیمی از متون مرتبط گردآوری شده و سپس توسط اسناد، تقسیم می‌شود. سپس، ماتریس هم‌رخدادی برای عبارات و اسناد می‌گردد. در ادامه، هر سلول انتخاب و محاسبه می‌شود. در نهایت باید افزود تجزیه مقادارهای تکی نقش مهمی برای محاسبه ابعاد و ایجاد ماتریس سه‌بعدی ایفاء می‌کند (القمدی و الفلقی<sup>۴</sup>، ۲۰۱۵).

## • تحلیل معنایی پنهان احتمالی

تحلیل معنایی پنهان احتمالی (PLSA) روشی است که پس از روش تحلیل معنایی پنهان و برای برطرف نمودن معایب تحلیل معنایی پنهان ارائه شد. تحلیل معنایی پنهان احتمالی روشی است که امکان نمایه‌سازی خودکار سند را بر اساس مدل کلاس پنهانی آماری افزایش می‌دهد و درصدد بهبود تحلیل معنایی پنهان به صورت احتمالاتی و با استفاده از مدل مولد

---

<sup>1</sup> Latent Semantic Indexing

<sup>2</sup> Silva & de Carvalho

<sup>3</sup> Singular Value Decomposition

<sup>4</sup> Alghamdi and Alfalqi

است. هدف اصلی تحلیل معنایی پنهان احتمالی، شناسایی و تمیز میان مفاهیم مختلف واژه‌های به‌کاررفته بدون استفاده از واژه‌نامه است. این روش دو نتیجه مهم دارد: نخست اینکه ابهام‌زدایی از واژه‌های دارای چند معنا را ممکن می‌سازد و دوم با گروه‌بندی واژه‌هایی که چارچوب مشترکی دارند، شباهت‌های معمول را نشان می‌دهد (هافمن، ۲۰۰۱).

تحلیل معنایی پنهان احتمالی مبتنی بر مدلی آماری است که در قالب مدلی ابعادی شناخته‌شده است. مدل ابعادی مدل متغیر پنهان برای داده‌های بارخداد مشترک است که دسته متغیرهای مشاهده نشده را با هر مشاهده، مرتبط می‌کند. روش تحلیل معنایی پنهان احتمالی برای بهبود تحلیل معنایی پنهان ارائه‌شده و بنابراین مشکلاتی که تحلیل معنایی پنهان نمی‌تواند حل کند را آدرس‌دهی می‌نماید. تحلیل معنایی پنهان احتمالی کاربردهای موفقی از جمله بینایی رایانه‌ای<sup>۱</sup> و سیستم‌های توصیه‌گر<sup>۲</sup> در دنیای واقعی داشته است.

کاربردهای تحلیل معنایی پنهان احتمالی حوزه‌های مختلفی از جمله بازیابی و فیلترسازی اطلاعات، پردازش زبان طبیعی و یادگیری ماشین از متن را شامل می‌شود. به‌ویژه، برخی از این کاربردها، طبقه‌بندی پردازش خودکار، دسته‌بندی، ردگیری موضوعی، بازیابی تصویر و توصیه خودکار پرسش هستند. در ادامه دو مورد از این کاربردها توضیح داده شده است.

○ بازیابی تصویر. مدل تحلیل معنایی پنهان احتمالی ویژگی‌های بصری دارد که باعث می‌شود برای نمایش هر تصویر به‌عنوان مجموعه واژه‌های بصری از یک واژه‌نامه بصری به کار رود. وقوع واژه‌های بصری در یک تصویر در بردار رخداد مشترک محاسبه شده است. هر تصویر، بردارهای رخداد مشترکی دارد که به ساخت جدول رخداد مشترک برای به‌دست آوردن مدل تحلیل معنایی پنهان احتمالی کمک می‌کند. پس از شناخت مدل تحلیل معنایی پنهان احتمالی می‌توان آن را به همه تصاویر پایگاه داده اعمال نمود. سپس آرایش برداری برای نمایش هر تصویر به کار می‌رود و هر عنصر بردار درجه‌ای که هر تصویر، موضوع معینی را ترسیم می‌کند را نشان می‌دهد (رامبرگ<sup>۳</sup> و همکاران، ۲۰۰۸)

---

<sup>1</sup> Computer vision

<sup>2</sup> recommender systems

<sup>3</sup> Romberg

○ توصیه خودکار پرسش. یکی از کاربردهای مهمی که تحلیل معنایی پنهان احتمالی با آن سروکار دارد، وظیفه توصیه پرسش است. در این کاربرد، واژه مستقل از کاربرد است. اگر کاربر معنای خاصی مدنظر داشته باشد و نیز زمانی که کاربر، پاسخها و معناهای نهفته مربوط به پرسش را دریافت می کند، می تواند توصیه هایی براساس شباهت های معناهای نهفته ارائه دهد. وو<sup>۱</sup> و همکاران (۲۰۰۸) اظهار داشتند که تحلیل معنایی پنهان احتمالی به منظور مدل سازی پروفایل کاربران استفاده می شود. همچنین پرسشها با حذف احتمالات موضوعات پنهان پشت واژهها مدل سازی می شوند. دلیل این امر آن است که پروفایل کاربر توسط پرسشها و پاسخها نمایش داده شده است. از این رو، فقط چگونگی مدل سازی درست پرسش مدنظر قرار می گیرد.

### • تخصیص پنهانی دیریکله

بهبود روش مدل سازی ترکیبی علت اصلی بروز مدل تخصیص پنهانی دیریکله (LDA) بود که قابلیت تبادل واژهها و اسناد را از روش قدیمی و توسط تحلیل معنایی پنهان احتمالی و تحلیل معنایی پنهان ثبت می کند. این امر همزمان با نظریه بازنمون کلاسیک<sup>۲</sup> در سال ۱۹۹۰ روی داد (بلی و همکاران، ۲۰۰۳).

مجموعه اسناد الکترونیکی متعدد از جمله صفحات وبی، وبلاگ های علمی و مقالات اخباری چالش های جدیدی برای پژوهشگران در جامعه داده کاوی ایجاد کرده اند. ضرورت روش های خود کار بصری سازی، تحلیل و خلاصه سازی این مجموعه اسناد افزایش یافته است. اخیراً، مدل سازی موضوعی پنهان به عنوان تکنیکی بدون نظارت، محبوبیت زیادی برای کشف موضوع در مجموعه کلان داده های متنی کسب کرده است. این مدل، الگوریتمی برای متن کاوی است که براساس مدل های موضوعی آماری (بیزین) بوده و کاربرد بسیار گسترده ای داشته است. مدل تخصیص پنهانی دیریکله یک مدل مولد است که در صدد تقلید پردازش نوشتن است؛ بنابراین سعی می کند یک سند را براساس موضوع مورد نظر، تولید کند. این مدل همچنین می تواند برای سایر انواع داده ها به کار رود. ده ها مدل مبتنی بر مدل تخصیص پنهانی دیریکله وجود دارد از جمله متن کاوی موقتی،

---

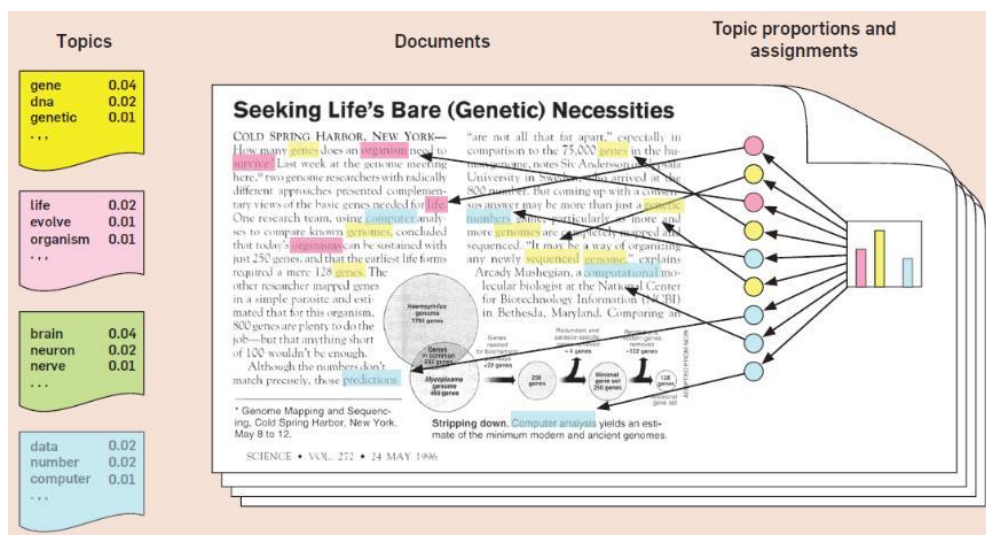
<sup>1</sup> Wu

<sup>2</sup> Classic representation theorem

تحلیل موضوع - نویسنده، مدل های موضوعی نظارت شده، خوشه بندی مشترک دیریکله و مدل تخصیص پنهانی دیریکله مبتنی بر بیوانفورماتیک (شن<sup>۱</sup> و همکاران، ۲۰۰۸؛ وانگ و همکاران، ۲۰۱۶).

شکل ۱-۲ ایده نهفته در منطق الگوریتم تخصیص پنهان دیریکله را نشان داده است. این روش فرض می کند، تعدادی موضوع که هر کدامشان توزیعی از کلمات است، در کل مجموعه مستندات موجود است (سمت چپ شکل). حال می توان تصور نمود که هر مستند این گونه شکل می گیرد:

ابتدا توزیعی از موضوعات انتخاب می شود، سپس برای هر واژه یک تخصیص موضوعی (دایره های رنگی در شکل) گزینش شده و در نهایت نیز واژه مورد نظر از موضوع مربوط انتخاب می گردد. شایان ذکر است که موضوعات و تخصیص های موضوعی به نمایش درآمده در این تصویر تنها به عنوان مثالی توضیح دهنده ارائه شده و از اعمال بر داده های واقعی به دست نیامده اند (بلی، ۲۰۱۲).



شکل ۴-۲. ایده نهفته در الگوریتم تخصیص پنهان دیریکله (گرفیت و استیورز، ۲۰۰۴)

<sup>1</sup> Shen

ایده اصلی این فرآیند این است که هر سند به صورت ترکیبی از موضوعات به صورت مدل درآمده است و هر موضوع یک توزیع احتمالی گسسته است که تعیین می کند چگونه احتمال حضور هر کلمه در موضوع مورد نظر وجود دارد. این احتمالات موضوعی، باز نمود دقیقی از سند را ارائه می کنند (برگولز<sup>۱</sup> و همکاران، ۲۰۰۸). در اینجا، یک «سند» در واقع «کیسه‌ای از کلمات»<sup>۲</sup> است که ساختاری فراتر از آماره های کلمه و موضوع ندارد. از کاربردهای مدل مبتنی بر روش مدل تخصیص پنهانی دیریکله می توان موارد زیر را نام برد:

○ **کشف نقش**<sup>۳</sup>: تحلیل شبکه اجتماعی (SNA) بررسی مدل های ریاضی برای تعامل میان افراد، سازمان ها و گروه ها است. به دلیل بروز ارتباطاتی میان مهاجمان امنیتی یازدهم سپتامبر و مجموعه داده های عظیم انسانی در سرویس های وب محبوب از جمله فیسبوک و مای اسپیس<sup>۴</sup>، علاقه به تحلیل شبکه های اجتماعی، افزایش یافته است. این امر منجر به مدل موضوع - گیرنده - نویسنده<sup>۵</sup> (ART) برای تحلیل شبکه های اجتماعی شد. مدل تخصیص پنهانی دیریکله و مدل - موضوع - نویسنده را ترکیب می کند. ایده مدل - موضوع - گیرنده - نویسنده، درک توزیع های موضوعی براساس پیام های حساس به جهت که بین دریافت کنندگان و فرستندگان ارسال شده است، دارای اهمیت می باشد (مک کلم<sup>۶</sup> و همکاران، ۲۰۰۷).

○ **موضوع احساسی**<sup>۷</sup>: مدل جفت - لینک<sup>۸</sup> - مدل تخصیص پنهانی دیریکله که بر مسئله مدل سازی اتصال متن و ارجاعات در قلمرو مدل سازی موضوعی متمرکز است. این براساس ایده مدل تخصیص پنهانی دیریکله و مدل های بلوک احتمالی عضویت ترکیبی<sup>۹</sup> (MMSB) است و امکان مدل سازی دلخواه لینک را می دهد (بائو<sup>۱۰</sup> و همکاران، ۲۰۰۹).

---

<sup>1</sup> Bergholz

<sup>2</sup> bag of words

<sup>3</sup> Role discovery

<sup>4</sup> My space

<sup>5</sup> Author-Recipient-Topic

<sup>6</sup> McCallum

<sup>7</sup> Emotion topic

<sup>8</sup> Pairwise-Link-LDA

<sup>9</sup> Mixed Membership Stochastic Block Models

<sup>10</sup> Bao

○ **نمره‌دهی خودکار متن**<sup>۱</sup>: مسئله نمره‌دهی خودکار متن همبستگی نزدیکی با طبقه‌بندی متن دارد و از دهه ۱۹۶۰ مورد بررسی قرار گرفته است. مدل تخصیص پنهان دیریکله در مقایسه با روش‌های کاهش ابعاد<sup>۲</sup> برای نمره‌دهی خودکار متن، نشان داده که رویکردهای قابل اطمینانی برای وظایف مربوط به بازیابی اطلاعات وجود دارد (کاکونن<sup>۳</sup> و همکاران، ۲۰۰۶).

#### • مدل موضوعی هم‌بسته

مدل موضوعی هم‌بسته<sup>۴</sup> نوعی مدل آماری مورد استفاده در پردازش زبان طبیعی و یادگیری ماشین است. مدل موضوعی هم‌بسته برای کشف موضوعاتی که در گروه اسناد نشان داده‌اند، به کاررفته است. کلید مدل موضوعی هم‌بسته توزیع نرمال لجستیک است. مدل‌های موضوعی هم‌بسته به مدل تخصیص پنهانی دیریکله متکی است (القمدی و والفلقی، ۲۰۱۵).

#### ۲-۲-۴-۲. ماهیت بین‌رشته‌ای متن کاوی

منظور از بین‌رشته‌ای بودن به این معنی است که با به کارگیری اصول و قواعد دانش‌های موجود (دانشی بیشتر از دو قلمرو علمی) و ترکیب آن‌ها با هم دانش جدید ایجاد می‌شود (وانگ و همکاران، ۲۰۲۰). متن کاوی زمینه‌ی بین‌رشته‌ای است که

---

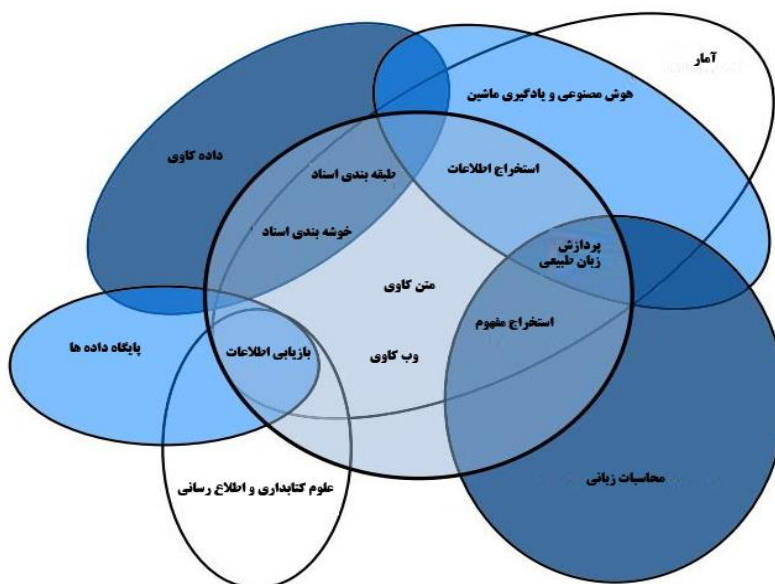
<sup>1</sup> Automated essay scoring

<sup>2</sup> Dimension Reduction Methods

<sup>3</sup> Kakkonen

<sup>4</sup> Correlated Topic Model (CTM)

به‌طور کلی از شش قلمرو دانشی و اشتراک آنها بهره می‌گیرد. آمار، هوش مصنوعی و یادگیری ماشین، داده‌کاوی، محاسبات زبانی، پایگاه داده و علوم کتابداری و اطلاع‌رسانی رشته‌هایی هستند که متن‌کاوی از آنها استفاده می‌کند. شکل ۵-۲ نمودار اشتراکات متن‌کاوی و سایر رشته‌ها را نشان می‌دهد (ماینر و همکاران، ۲۰۱۲).



شکل ۵-۲. نمودار ون اشتراکات متن‌کاوی و سایر قلمروهای علمی (ماینر و همکاران، ۲۰۱۲)

همان‌طور که در شکل ۵-۲ مشاهده می‌شود از اشتراک رشته‌های مذکور، هفت روش عملی برای رسیدن به اهداف متن‌کاوی ایجاد می‌شود که عبارتند از پردازش زبان طبیعی، طبقه‌بندی اسناد، خوشه‌بندی اسناد، بازیابی اطلاعات، وب‌کاوی و استخراج مفهوم هستند.

- طبقه‌بندی اسناد<sup>۱</sup>: در این روش بهترین برچسب از مجموعه برچسب‌های از قبل مشخص شده برای مدارک بدون برچسب<sup>۲</sup> مورد استفاده قرار می‌گیرد.

- خوشه‌بندی اسناد<sup>۳</sup>: به فرآیند گروه‌بندی مستندات مشابه درون خوشه‌های مختلف خوشه‌بندی می‌گویند.

<sup>1</sup> Document Classification

<sup>2</sup> untagged documents

<sup>3</sup> Document Clustering

- بازیابی اطلاعات<sup>۱</sup>: فرآیند نمایه‌سازی<sup>۲</sup>، جستجو و بازیابی و بازیابی مستندات از بین مجموعه داده‌های عظیم متنی با استفاده از کلیدواژه‌ها را بازیابی اطلاعات گویند. مهم‌ترین کاربرد بازیابی اطلاعات در موتورهای جستجوی دیده می‌شود.

- وب کاوی<sup>۳</sup>: فرآیند داده کاوی و متن کاوی بر روی محتوای صفحات وبی و لینک‌ها را وب کاوی می‌گویند. صفحات وب حالتی نیمه‌ساختاریافته دارند و از مجموعه‌ای از متون و لینک‌ها به صفحات دیگر تشکیل شده‌اند. لذا یک رویکرد متداول برای وب کاوی، بازنمایی صفحات وب در قالب گراف و تحلیل گراف‌های وبی است. امروزه زیرشاخه‌ای از وب کاوی برای تحلیل شبکه‌های اجتماعی مانند فیس‌بوک، توئیتر و ... بسیار مورد توجه پژوهشگران قرار گرفته است.

- استخراج اطلاعات<sup>۴</sup>: فرآیند شناسایی و استخراج موجودیت‌های مناسب و همچنین روابط بین آن‌ها از درون متن (غیر ساخت یافته) را استخراج اطلاعات می‌گویند. به عبارت دیگر، استخراج اطلاعات فرآیندی برای پردازش داده‌های غیر ساخت یافته (مثل متن، تصویر، صوت و ...) و نیمه‌ساختاریافته (مثل صفحات وب، XML، ...) و تبدیل آنها مجموعه داده ساختاریافته (از قبیل جداول پایگاه داده) است. استخراج اطلاعات به دو نوع باز (عمومی) و بسته (خاص و در حوزه‌ای مشخص) تقسیم می‌شود.

- پردازش زبان طبیعی<sup>۵</sup>: هدف آن، پردازش سطح پایین و فهم (درک) زبان و بخصوص متن توسط کامپیوترها است. معمولاً معادل با اصلاح زبان‌شناسی رایانشی<sup>۶</sup> بکار گرفته می‌شود، هرچند که اغلب زبان‌شناسان زبان‌شناسی رایانشی را کلی‌تر از پردازش زبان طبیعی می‌دانند.

---

<sup>1</sup> Information Retrieval-IR

<sup>2</sup> indexing

<sup>3</sup> Web Mining

<sup>4</sup> Information Extraction-IE

<sup>5</sup> Natural Language Processing-NLP

<sup>6</sup> Computational linguistics



- استخراج مفاهیم<sup>۱</sup>: فرآیند گروه‌بندی واژه‌ها و عبارات درون گروه‌های متنی که دارای مشابهت معنایی هستند را استخراج مفاهیم گویند. معمولاً از تکنیک‌های آماری (مانند n-grams)، مدل‌سازی موضوعی<sup>۲</sup> و خوشه‌بندی متون و واژه‌ها برای استخراج مفاهیم استفاده می‌شود.

### ۲-۳. مرور پیشینه‌های پژوهش

در دهه‌های گذشته، چالش اصلی در فرایند تصمیم‌گیری و انتخاب، کمبود اطلاعات بود؛ اما امروزه اشباع اطلاعاتی و انتخاب اطلاعات مناسب، چالش اصلی به شمار می‌آید. فضای تولیدات علمی به صورت مقاله، ثبت اختراع و اکتشاف، چاپ کتاب، انتشار مجلات علمی-تخصصی، چه از نظر اقتصادی و چه از جنبه‌های مدیریتی، رعایت اخلاق و امانتداری علمی، تسهیم منافع علمی، آموزش آداب و اخلاق علمی عدم ضابطه‌مندی را نشان می‌دهد. این عدم ضابطه‌مندی، بر آشفتگی در تولید علم دامن می‌زند (کینی، ۲۰۰۷). از طرفی از نظر عملیاتی، نویسندگان با شناسایی و طبقه‌بندی میدانی نشریات و انتساب ارزش و ارتباط موضوعی هر نشریه اقدام به ارسال و انتشار یافته‌های پژوهش خود می‌کنند. از اینرو این نکته که نشریات تا چه اندازه در اعلام حوزه‌های موضوعی که مقالات را می‌پذیرند شفاف عمل کرده و به آن پایبند بوده حائز اهمیت است. گاهی موضوعات اعلامی از سوی نشریات بسیار کلی است و کاربر را در انتخاب نشریه مرتبط با دشواری مواجه می‌کند. از طرفی با توجه به محدودیت‌های زمانی و همچنین حجم اطلاعات، جستجو و تحلیل دستی منابع علمی عملاً امکان‌پذیر نیست و آنچه حائز اهمیت است استفاده از روش‌های یادگیری ماشین مانند متن کاوی برای پردازش داده‌های بزرگ است. از جمله کاربردهای متن کاوی می‌توان به دسته‌بندی، خوشه‌بندی، خلاصه‌سازی و یافتن روابط میان مفاهیم در متون اشاره کرد. متن کاوی دارای دو روش یادگیری «با نظارت»<sup>۳</sup> و «بدون نظارت»<sup>۴</sup> است (الهیاری، ۲۰۱۷). در روش یادگیری با نظارت، هر یک از داده‌های آموزشی به دسته‌ای خاص که از ابتدا مشخص هستند نسبت داده می‌شود. بعبارتی در هنگام آموزش ناظری

<sup>1</sup> Concept Extraction

<sup>2</sup> topics modeling

<sup>3</sup> supervised

<sup>4</sup> unsupervised

وجود دارد که اطلاعاتی علاوه بر داده‌های آموزش در اختیار یادگیرنده قرار می‌دهد. در یادگیری بدون نظارت، یادگیرنده است که باید درون داده‌ها به دنبال ساختاری خاص بگردد؛ زیرا به‌جز داده‌های آموزشی، هیچ اطلاعاتی در اختیار یادگیرنده قرار ندارد (کراپر و همکاران<sup>۱</sup>، ۲۰۱۹).

در پیوند با پژوهش حاضر طبقه‌بندی یا دسته‌بندی موضوعی مقالات نشریه‌ها به گونه‌ای مد نظر است که بتواند مفهوم را به درستی منتقل نماید. از اینرو از روشهای مدل‌سازی موضوعی بهره‌برده خواهد شد. مدل‌سازی موضوعی تکنیکی است که ساختار موضوع را در مجموعه‌ای از اسناد کشف و تفسیر مینماید (بلی<sup>۲</sup>، ۲۰۱۲). به عبارت دیگر، در حالت کلی روشها و الگوریتمهایی هستند که متن را پردازش کرده و موضوعات مختلف موجود در آن را استخراج مینمایند.

الگوریتم‌ها برای تولید و ارائه پیشنهاد از عناصر متنی مقالات مانند عنوان، چکیده و واژه‌های کلیدی استفاده می‌کنند (آگاروال، ۲۰۱۶). که به شکل انبوهی از داده‌های ساخت‌نیافته در متون وجود دارند و بایستی به نحوی که برای ماشین قابل استفاده و درک باشد استخراج شود (کوچوکتانچ<sup>۳</sup> و همکاران، ۲۰۱۲). بنابراین میتوان گفت که در مدل‌سازی موضوعی هر سند با توجه به موضوعات موجود در آن تفسیر و سازماندهی میگردد. در این مدلها، هر متن یا سند به صورت توزیعی از موضوعات ارائه میگردد؛ در حالی که هر موضوع هم توزیعی روی واژگان میباشد (کروا و بانسل، ۲۰۲۰). در مدل‌سازی موضوعی موضوعات با تشخیص الگوهایی مثل کلمات و تکرارشان تشخیص داده می‌شوند. مدل فضای برداری (VSM) اولین مدل جبری ساده‌ای بود که مستقیماً بر اساس ماتریس اصطلاح-سند برای استخراج اطلاعات معنایی ارائه گردید (سالتون و همکاران، ۱۹۷۵). یک مدل فضای برداری پایه میتواند متن را به کاراکترهای یونی‌گرم تا ان-گرم تقسیم کند که بر اساس روش کیف کلمات (BOW) میباشد. به عبارت دیگر، ترتیب دقیق اصطلاح در یک سند نادیده گرفته شده؛ اما تکرار وقوع هر اصطلاح به عنوان یک فاکتور مهم نگهداری میگردد (بری و همکاران، ۱۹۹۹). کاربرد مدل فضای برداری در حوزه‌هایی است که نیاز به محاسبه میزان شباهت میان کلمات و اصطلاحات موجود در اسناد دارند که از آن

---

<sup>1</sup> Kruber et al.

<sup>2</sup> Blei

<sup>3</sup> Küçükünç

جمله میتوان به موتورهای جستجو، پردازش زبان طبیعی و ماشین اشاره کرد (گرب و همکاران، ۲۰۱۳).

برای مدل‌سازی سامانه بازایی مانند موتورهای جستجو یا مجموعه دیجیتالی، همه کلمات در یک سند به یک اندازه مهم نیستند بنابراین به هر اصطلاح در سند یک وزن بر اساس تعداد وقوع اصطلاح در آن سند داده میشود. این وزن بهتر است بر اساس وقوع یک اصطلاح در سند و عدم تکرار در اسناد دیگر باشد. بنابراین در متون معمولاً از وزندهی Tf\_IDF استفاده میگردد (تورنی و پاتل، ۲۰۱۰).

الگوریتمهای مدل‌سازی موضوعی با استفاده از تکنیکهای مختلف تلاش میکنند تا موضوعات خوشه بندی شده را تحت عنوان یک موضوع ارائه دهند. مدل توزیعی مانند فضای برداری، آنالیز پنهان مفهوم، آنالیز پنهان مفهومی احتمالی و تخصیص دیریکله پنهان میتواند معنای کلمات را از متن با استفاده از روشهای آماری استخراج کند (کرین و همکاران، ۲۰۱۲، زمانی و همکاران، ۱۳۹۳). به منظور انجام مدل‌سازی موضوعی در متن، نیاز به پردازش زبان میباشد. ابزارهای استاندارد پیش پردازش و نرمالسازی ایجاد شده برای متون زبان فارسی به صورت رایگان منتشر نشده یا دقت مناسب را ندارند (کامیابی و همکاران، ۱۳۹۷). به دلیل نزدیکی دبیره زبان فارسی با عربی، همواره در نگارش تعدادی از حرف‌ها مشکل نویسه‌های عربی معادل وجود دارد. از جمله آنها می‌توان به حروف «ک»، «ی»، «هزّه و... اشاره کرد. علاوه بر این، وجود نویسه «ـ»، تشدید، تنوین و موارد مشابه (کامیابی و همکاران، ۱۳۹۷؛ سراجی<sup>۱</sup>، ۲۰۱۰) ایست واژه‌ها<sup>۲</sup> از جمله حروف اضافه، بسیاری از قیده‌ها، حروف ربط و افعال از جمله اقدام‌های لازم قبل از شروع پردازش متن است که این پردازشها در هر سطح، نیازمند دانش، منابع و پیکره‌های مورد نیاز آن سطح و سطوح پایینتر است. در دسترس بودن منابع و دانش برای انجام تحقیق در حیطه‌ی پردازش زبان طبیعی از جمله چالشهای پردازش زبان طبیعی است.

در بحث متن کاوی و پردازش زبان طبیعی از شیوه‌ها و الگوریتم‌های متفاوتی به منظور وزندهی به واژگان استفاده می‌شود.

هر یک از این الگوریتم‌ها دارای نقاط قوت و ضعفی هستند که براساس روش‌شناسی پژوهش‌های گوناگون و نیز شرایط حاکم بر پژوهش‌ها برای وزندهی به واژگان مورد استفاده قرار می‌گیرند. مرور پژوهش‌هایی که از روش TF-IDF برای وزندهی

---

<sup>1</sup> Seraji

<sup>2</sup> Stop word

استفاده کرده‌اند، حاکی از آن است که این روش از نقاط قوت بسیار زیادی برخوردار است (راموس<sup>۱</sup>، ۲۰۰۳؛ ایلیاس<sup>۲</sup>، ۲۰۱۹؛ کیم و دلن، ۲۰۱۸). راموس (۲۰۰۳) یکی از مهم‌ترین مزایا و کاربرد TF-IDF را در حجم بالای داده‌ها می‌داند. به بیان دیگر، استفاده از TF-IDF در پژوهش‌هایی که دارای داده‌هایی با حجم بالا هستند، توصیه می‌شود.

پژوهش‌های مرتبط با بررسی میزان تخصص‌گرایی یا تنوع موضوعی مقالات نشریات با استفاده از روش‌های متن‌کاوی به صورت موردی انجام شده است. در پیوند با پژوهش حاضر به دلیل تفاوت‌های ذاتی زبان فارسی نسبت به سایر زبانتها، در این بخش به پژوهش‌های مرتبط با زبان فارسی پرداخته شده است.

تحلیل محتوای موضوعی یا ترسیم نقشه موضوعی مقالات با استفاده از روش‌های متن‌کاوی به صورت موردی برای برخی از نشریه‌های علمی انجام شده است. نتایج مطالعه‌ی فضلی پور و جمالی مهموئی در بررسی میزان تخصص‌گرایی مجله‌های حوزه علم اطلاعات و دانش‌شناسی فارسی طی سالهای ۱۳۸۷ الی ۱۳۹۱ نشان داد که ۹۳ درصد مقالات با موضوع‌های مرتبط و ۷ درصد مقالات با موضوع‌های غیرمرتبط با علم اطلاعات در این نشریات منتشر شده‌اند. همچنین نتایج نشان داد در غالب نشریات تناسبی بین موضوعی و نام نشریه وجود ندارد. در این پژوهش از طبقه‌بندی موضوعی نیز استفاده شده است (۱۳۹۳). ناغانی و عابسی (۱۳۹۵) از روش متن‌کاوی به منظور آنالیز محتویات مقالات علمی در زمینه مهندسی صنایع استفاده و تکنیک‌هایی را ارائه دادند.

شکوهیان و همکاران (۱۳۹۸) به ارائه مدل دسته‌بندی موضوعی تولیدات علمی حوزه سلامت با استفاده از روش‌های متن‌کاوی پرداختند. در این مطالعه از مدلی ترکیبی برای دسته‌بندی موضوعی تولیدات علمی استفاده شده است. داده‌ها مربوط به سال ۲۰۰۹ الی ۲۰۱۹ مستخرج از پایگاه پابمد می‌باشد. نتایج نشان داد ترکیب دسته‌بندی و خوشه‌بندی می‌تواند دقت دسته‌بندی متون سلامت را افزایش دهد.

---

<sup>1</sup> Ramos

<sup>2</sup> Ilias

دهدشتی شاهرخ (۱۳۹۸) به تحلیل متن کاوی نشریه برنامه ریزی و توسعه گردشگری با استفاده از تکنیک متن کاوی پرداخته است. در این پژوهش ساختار محتوایی مقالات بررسی و حوزه های مطالعاتی این فصلنامه مورد بررسی قرار گرفته اند.

تجزیه و تحلیل موضوعی مقالات منتشرشد کتابداری و اطلاع رسانی پزشکی در ایران با استفاده از فنون متن کاوی نیز توسط داستانی و همکاران (۱۳۹۹) مورد مطالعه قرار گرفت. این مطالعه با استفاده از رویکرد اکتشافی و توصیفی به تجزیه و تحلیل مقالات کتابداری و اطلاع رسانی منتشرشده در مجلات تخصصی این حوزه در ایران از سال ۱۳۷۶ تا ۱۳۹۸ با استفاده از فنون متن کاوی پرداخته است. جهت تعیین موضوعات منتشرشده، الگوریتم TF-IDF انتخاب گردیده است. جهت شناسایی مهمترین واژگان به کار رفته در مقالات از الگوریتم وزن دهی و از زبان برنامه نویسی پایتون نیز جهت اجرای الگوریتمهای متن کاوی استفاده شده است.

نتایج پژوهش حوزه های موضوعی مقالات را شناسایی نمود. توسعه و روند موضوعی حوزه علم اطلاعات و دانش شناسی نیز بر اساس مدل موضوعی LDA توسط باغ محمد و همکاران (۱۳۹۹) مورد بررسی قرار گرفته است. به منظور دستیابی به اهداف پژوهش، داده ها از سال ۲۰۰۸ تا ۲۰۱۹ از پایگاه اسکوپوس استخراج شدند. سپس با استفاده از الگوریتمهای متن کاوی و به طور خاص الگوریتم مدل سازی موضوعی LDA با استفاده از نرم افزار R مورد تحلیل قرار گرفتند. نتایج نشان داد که پژوهشهای موضوعی رشته علم اطلاعات و دانش شناسی در ایران همگام با رشد فناوری ها و موضوعهای جهانی توسعه یافته است. قناد و همکاران (۱۴۰۲) به کاربست فنون متن کاوی در تحلیل جریان موضوعی مقالات منتشر شده در مجلات حسابداری ایران پرداختند. در این پژوهش به منظور شناسایی مهمترین واژگان بکاررفته در مقالات از الگوریتم وزن دهی TF-IDF و جهت کاربست الگوریتمهای متن کاوی از زبان برنامه نویسی پایتون استفاده شده است. نتایج نشان داد که توجه به موضوعاتی نظیر فناوری و سیستم های اطلاعاتی حسابداری، حسابداری بخش عمومی و آموزش و پژوهش حسابداری در مقایسه با سایر حوزه های پژوهشی کمتر بوده است

در مجموع نیز یافته های مطالعه آرامو (۲۰۱۸) نشان داده است تخصص در مقابل تنوع در پژوهش‌ها پژوهش‌ها علمی نتایج بهتری را به دنبال دارد. همچنین نتایج مطالعه حری (۱۳۷۴) حاکی از آنست که بین مجلات و مقالات منتشره حوزه‌های تخصصی و ارتقا علمی متخصصان رابطه مثبت وجود دارد. میزان تخصص‌گرایی و تنوع موضوعی در مقالات نشریات با استفاده از روشهای متن کاوی در برخی نشریات مورد مطالعه قرار گرفته است.

با توجه به مرور پیشینه‌ها، برای بررسی متون علمی، شاهد حجم عظیمی از داده هستیم که برای بررسی ساختار و تحلیل روند موضوعی آنها استفاده از روش‌های سنتی و دستی، امکان‌پذیر نبوده یا با خطا مواجه خواهد بود. بنابراین همگام با توسعه فناوری‌های پردازش متن و یادگیری ماشین در تحلیل داده‌های بزرگ رویکردهای متن کاوی برای شناسایی روند موضوعی پژوهش‌ها مورد استفاده قرار گرفته است (لی و همکاران، ۲۰۱۸)؛ اما همانگونه که مشاهده شد غالب پژوهش‌ها به بررسی و مرور محتوای مقالات نشریات به صورت خاص پرداخته و نگرش جامعی در این خصوص وجود ندارد. بنابراین پژوهش حاضر بر آنست که میزان همخوانی دامنه موضوعی نشریات معتبر وزارت علوم، پژوهش‌ها پژوهش‌ها و فناوری در حوزه‌های شش‌گانه را با محتوای مقالات آنها و بر اساس فنون متن کاوی مورد تجزیه و تحلیل قرار دهد.

# فصل سوم

## روش شناسی پژوهش

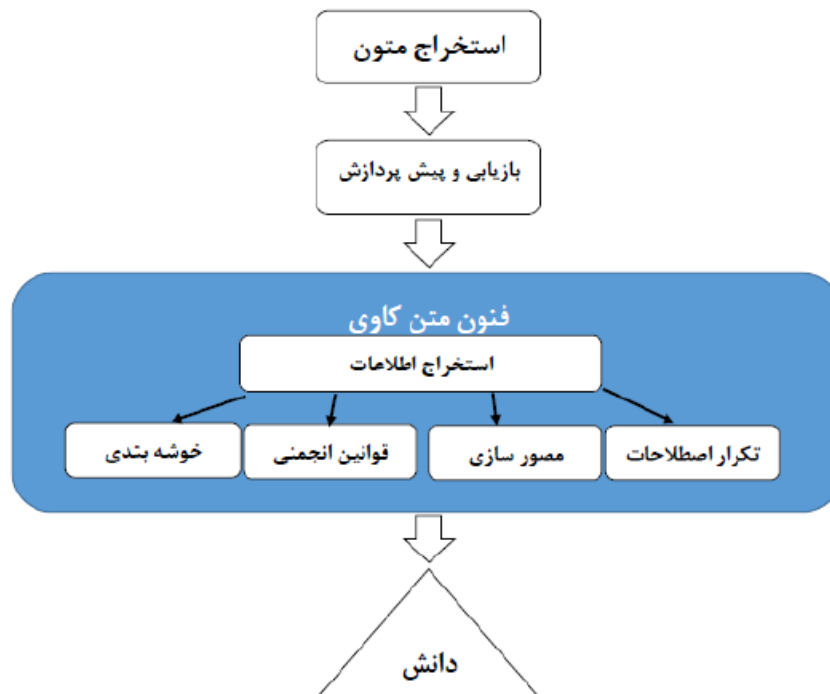
در این فصل روش‌شناسی پژوهش شامل نوع، روش و رویکرد، جامعه پژوهش، روش‌های گردآوری و تجزیه و تحلیل داده‌ها، مراحل اجرای پژوهش تبیین شده است.

### ۳-۲. روش پژوهش

پژوهش از نوع کاربردی و به روش متن کاوی با رویکرد تحلیلی انجام شده است. کاربردی بودن این پژوهش، از آن جهت است که وضع موجود قلمرو پژوهش را بررسی کرده و از نتایج به دست آمده می‌توان در تصمیم‌گیری‌ها و سیاست‌گذاری‌ها و همچنین برنامه‌ریزی‌ها استفاده کرد. همچنین در این پژوهش با استفاده از رویکرد اکتشافی-توصیفی به تجزیه و تحلیل مقالات منتشر شده در نشریات معتبر وزارت علوم پرداخته شد. این مطالعه به دلیل استفاده از فنون متن کاوی یک مطالعه اکتشافی است. متن کاوی یک روش اکتشافی داده محور است که جهت یافتن الگوها و روندها در مجموعه‌های داده‌های بزرگ استفاده می‌شود.

روش متن کاوی به کار رفته در این پژوهش برگرفته از چارچوب طراحی شده توسط ژانگ و چن و توسعه یافته توسط سالوم و همکاران است که شامل پیش پردازش متن، عملیات متن کاوری و پسا پردازش است. پیش پردازش متن شامل عملیات زیر است: انتخاب داده‌ها، دسته بندی، استخراج ویژگی، نرمال سازی، حذف کلمات زائد. دومین مرحله با فنون مختلف متن کاوی مانند خوشه بندی، بصری سازی و تکرار واژگان سروکار دارد. در خلال مرحله سوم تغییر و تبدیلهایی روی داده‌ها (مقالات علمی) از طریق توابع متن کاوی مانند ارزیابی و انتخاب دانش؛ تجزیه و تحلیل و بصری سازی دانش صورت می‌گیرد.





### ۳-۳. جامعه پژوهش

در ابتدا فهرستی از نشریات فارسی مصوب وزارت علوم، تحقیقات و فناوری مشتمل بر ۱۱۵۸ عنوان نشریه دارای رتبه طبق پرتال نشریات<sup>۱۵۶</sup> (ابان ۱۴۰۰) تهیه و سپس نمونه گیری به روش طبقه ای-تناسبی از بین شش حوزه موضوعی سطح کلان از پیش تعیین شده توسط وزارت عتف شامل علوم انسانی، فنی و مهندسی، علوم پایه، هنر و معماری، کشاورزی و منابع طبیعی و دامپزشکی انجام شد.

در جدول ۱ تعداد تمامی مجلات حوزه های مورد مطالعه و نمونه پژوهش به تفکیک حوزه موضوعی نشریه درج شده است. با توجه به محدودیت های زمانی، منابع انسانی و مالی امکان بررسی اهداف و دامنه موضوعی مقالات تمام نشریات در پژوهش حاضر وجود ندارد. بنابراین در این پژوهش بایستی حجم نمونه مشخص گردد. در همین راستا، با مشورت گرفتن از مشاور آماری طرح جهت تعیین حجم نمونه، از نرم افزار <sup>۱۵۷</sup>PASS نسخه ۱۱ استفاده شد (هینتز<sup>۱۵۸</sup>، ۲۰۱۳). با در نظر گرفتن فاصله

<sup>156</sup> journals.msrt.ir

<sup>157</sup> - Power Analysis Sample Size

158 - Hintze

اطمینان ۰/۹۵، توان آزمون ۰/۸۰ و همبستگی مورد انتظار ۰/۲۰، حداقل حجم نمونه مورد نیاز توسط نرم افزار، ۱۱۶ عنوان نشریه تعیین شد. حجم نمونه مورد نظر با استفاده از روش نمونه گیری تصادفی طبقه‌ای انتخاب شد به گونه‌ای که تناسب بین نمونه و جامعه در هر یک از شش حوزه موضوعی سطح کلان تعیین شده توسط وزارت عتف رعایت شود.

توزیع فراوانی نشریات و تعداد نمونه پژوهش در حوزه های موضوعی شش گانه وزارت عتف

ردیف	حوزه موضوعی	تعداد نشریه فارسی در حوزه (عنوان)	تعداد نمونه (عنوان)
۱	علوم انسانی	۷۲۵	۷۲
۲	علوم کشاورزی و منابع طبیعی	۱۵۸	۱۶
۳	فنی و مهندسی	۱۳۴	۱۳
۴	علوم پایه	۸۸	۹
۵	هنر و معماری	۴۲	۴
۶	دامپزشکی	۱۱	۲
	جمع	۱۱۵۸	۱۱۶

### ۳-۴. گردآوری داده ها

نکته ای که در این نوع مطالعات بایستی مورد توجه قرار گیرد و در حجم نمونه لحاظ شود؛ آن است که نمونه پژوهش در واقع شامل اطلاعات بخش حوزه موضوعی مندرج در وب سایت رسمی نشریه از یک سو و عنوان، کلیدواژه و چکیده مقالات دوره سه ساله هر یک از نشریات از سوی دیگر است. بنابراین در مرحله اول، وب سایت ۱۱۶ نشریه به منظور تعیین حوزه موضوعی درج و اعلام شده توسط نشریه بررسی و اطلاعات لازم گردآوری شد. در ادامه نشریات نمونه، در موتور

جستجوی گوگل به منظور بازیابی وب سایت اصلی مجله مورد جستجو قرار گرفت. سپس عنوان، چکیده و واژه های کلیدی مقالات هر نشریه طی یک بازه زمانی سه ساله استخراج شد که بالغ بر ۱۵۵۰۰۰ عنوان مقاله را دربرگرفت و در قالب بیب اکسل آماده سازی شد.

### ۳-۵. تجزیه و تحلیل داده‌ها

تجزیه و تحلیل داده‌ها و انجام عملیات متن کاوی در این پژوهش، شامل سه مرحله اصلی پیش پردازش،<sup>۱۵۹</sup> انجام عملیات و فنون متن کاوی و در پایان تحلیل و نتایج اصل از متن کاوی است. عملیات پاک‌سازی داده‌ها<sup>۱۶۰</sup> روی داده‌های متنی اجرا شد. پاک‌سازی داده‌ها، کیفیت داده‌ها، اعتبار الگوها و روابط استخراج شده را افزایش می دهد. پاک‌سازی داده‌ها، تنها داده‌های متنی مورد نیاز مرتبط را نگه می دارد (ایندوخیا و دامرو<sup>۱۶۱</sup>، ۲۰۱۰).

الگوریتمهای مدل‌سازی موضوعی با استفاده از تکنیکهای مختلف تلاش میکنند تا موضوعات خوشه بندی شده را تحت عنوان یک موضوع ارائه دهند. مدل توزیعی مانند فضای برداری، آنالیز پنهان مفهومی، آنالیز پنهان مفهومی احتمالی و تخصیص دیریکله پنهان میتواند معنای کلمات را از متن با استفاده از روشهای آماری استخراج کند (زمانی و همکاران، ۱۳۹۳). به منظور انجام مدل‌سازی موضوعی در متن، نیاز به پردازش زبان میباشد. ابزارهای استاندارد پیش‌پردازش و نرم‌السازی ایجاد شده برای متون زبان فارسی به صورت رایگان منتشر نشده یا دقت مناسب را ندارند (کامیابی و همکاران، ۱۳۹۷). به دلیل نزدیکی دبیره زبان فارسی با عربی، همواره در نگارش تعدادی از حروفها مشکل نویسه‌های عربی معادل وجود دارد. از جمله آنها می توان به حروف «ک»، «ی»، همزه و... اشاره کرد. علاوه بر این، حذف نویسه «ـ»، تشدید، تنوین و موارد مشابه (کامیابی و همکاران، ۱۳۹۷؛ سراجی<sup>۱۶۲</sup>، ۲۰۱۰) ایست واژه‌ها<sup>۱۶۳</sup> از جمله حروف اضافه، بسیاری از قیده‌ها، حروف ربط و افعال از جمله

---

<sup>159</sup> Preprocessing

<sup>160</sup> data cleaning

<sup>161</sup> Indurkhyia and Damerou

<sup>162</sup> Seraji

<sup>163</sup> Stop word

اقدام‌های لازم قبل از شروع پردازش متن است. سپس تنوع ریخت‌شناسی واژه‌ها در زبان فارسی طبق آخرین یافته‌ها از متون استخراج و مورد استفاده قرار گرفت (هماوند و همکاران ۱۳۹۷؛ پرئی و حمیدی ۱۳۹۶؛ ستوده و هنرجویان ۱۳۹۱). که این پردازش در هر سطح، نیازمند دانش، منابع و پیکره‌های مورد نیاز آن سطح و سطوح پایینتر است. در دسترس بودن منابع و دانش برای انجام تحقیق در حیطه‌ی پردازش زبان طبیعی از جمله چالشهای پردازش زبان طبیعی است. در زبان فارسی، این مشکلات و چالشها به مراتب بیشتر هم میشود؛ چرا که زبان فارسی به صورت ماهوی از پیچیدگیهای بیشتری برخوردار است و پژوهشهای به نسبت بسیار کمتری روی آن انجام گرفته است. پیوسته بودن نویسه‌ها، شباهت کاراکترها، کاراکترهای هم‌آوا و وجود برخی اسامی مرکب دو یا چند کلمه‌ای موجب افزایش خطا در متون فارسی میشود. ناهماهنگیهای گوناگونی که در نگارش خط فارسی دیده میشود، همچنین در نگارش رایانه‌ای متن فارسی، قالبها، ابزارها، سیستم عاملهای گوناگون و روشهای گوناگون کد کردن نوشته‌ی فارسی دیده میشود. ویژگیهای منحصر به فرد خط فارسی موجب بروز چالشهایی برای فرآیند خطایابی و تصحیح خطا میشود که برای دیگر زبانهایی که این خصوصیات و استثناءها را ندارند، مطرح نیست. عدم دسترسی به پیکره‌ها، دانش مورد نیاز و برخی چالشهای وابسته به رسم الخط زبان فارسی موجب میشود، پژوهش در این زمینه با مشکلات عدیده‌ای روبرو شده و پیشرفت در این زمینه به کندی پیش رود

در این مرحله، پیش پردازش داده‌ها با استفاده از برنامه نویسی پایتون بمنظور جداسازی واژه‌ها، یکدستی و نرمالسازی داده‌ها و به شرح زیر صورت گرفت.

- حذف نویسه‌های غیر مهم مانند فضاها، خالی اضافی، تگ‌های قالب‌بندی متن، حذف نویسه‌های غیر الفبایی شامل حذف علائم نگارشی و اعداد از متن است (بانکس<sup>۱۶۴</sup> و همکاران، ۲۰۱۸). بنابراین، این روش برای تمرکز روی کلمات و عبارات موجود در متن انجام شده است.

- شکستن اجزای متن<sup>۱۶۵</sup> به کلمات و واژه‌ها، شکستن اجزای متن به فرایند جداسازی و شکستن هر کلمه در سند براساس واحدهایی با معنی، شکستن یا تقسیم‌بندی متن گفته می‌شود (ویجایارانی<sup>۱۶۶</sup> و همکاران، ۲۰۱۵) و یکی از فرآیندهای ضروری پیش‌پردازش متن در متن‌کاوی است (اشمیدل<sup>۱۶۷</sup> و همکاران، ۲۰۱۹). که در این پژوهش داده‌های متنی براساس کلمات<sup>۱۶۸</sup> جداسازی شده است.

- نرمالسازی و یکدست‌سازی حروف مانند حروف «ی»، «ک»، همزه و... جهت یکدست‌سازی متن
- حذف ایست‌واژه‌ها: به‌منظور بازیابی یا تجزیه و تحلیل درست، حروف یا واژه‌هایی که ارزشی معنایی ندارند، با استفاده از سیاهه ایست‌واژه‌ها حذف می‌شوند (بانکس و همکاران، ۲۰۱۸). ایست‌واژه‌ها واژه‌هایی هستند که محتوای اطلاعاتی کم‌ارزشی یا بی‌ارزشی دارند و به معنادار کردن متن کمکی نخواهند کرد. افزون بر ایست‌واژه‌ها، داده‌های متنی بررسی شد و واژگان دیگری که در متن تکرار شده و معنی خاصی در متن نداشتند به سیاهه ایست‌واژه‌ها اضافه و از متن حذف شدند.

- در این پژوهش از الگوریتم وزن‌دهی واژگان TF-IDF استفاده گردید.

در مرحله بعدی تعیین فراوانی و وزن‌دهی واژگان، اجرای الگوریتم شباهت سنجی با استفاده از مدل فضای برداری انجام شد. بدین ترتیب که عبارات کلیدی مندرج در بخش دامنه موضوعی نشریه و همچنین عبارات کلیدی مستخرج از محتوای مقالات به یک سری بردار عددی تبدیل شدند. این بردارها دارای  $n$  مولفه بوده و وزندهی به آنها بر اساس تعداد تکرار کلمات لحاظ شده است. در نهایت با استخراج دانش از متون و تفسیر آن به پایان رسید.

- مقدار TF-IDF به تناسب تعداد تکرار واژه در سند افزایش می‌یابد و توسط تعداد اسنادی که در مجموعه هستند، متعادل می‌گردد (اوکوهارا<sup>۱۶۹</sup> و همکاران، ۲۰۱۸؛ وی جیرانی و جنانی، ۲۰۱۶؛ وو<sup>۱۷۰</sup> و همکاران، ۲۰۰۸).

<sup>165</sup> Tokenization

<sup>166</sup> Vijayarani

<sup>167</sup> Schmiedel

<sup>168</sup> Word tokenize

<sup>169</sup> Okuhara

<sup>170</sup> Wu

## فصل چهارم

# تجزیه و تحلیل داده‌ها

#### ۴-۱. مقدمه

در این فصل یافته های حاصل از تحلیل داده های پژوهش توصیف شده و با ترتیبی منطقی ارائه گردیده است. به این صورت که داده های تحلیل شده در قالب جدولها و نمودارها گزارش و توصیف های مربوط به هر کدام ارائه گردیده است.

#### ۴-۲. استخراج و آماده سازی دامنه موضوعی نشریات

به منظور رسیدن به هدف پژوهش ابتدا دامنه موضوعی نشریات استخراج و عبارات کلیدی پس از پیش پردازش ماشینی توسط پژوهشگر نیز کنترل و استخراج شد. به دلیل حجم زیاد داده ها جدول شمایی از فرآیند را نشان می دهد.

استخراج عبارت کلیدی از دامنه و اهداف موضوعی درج شده در وبگاه هر نشریه

ردیف	عنوان نشریه	حوزه موضوعی مندرج در وبگاه	عبارات کلیدی
۱	پردازش علائم و داده ها	- انتشار دانش علمی و تخصصی در زمینه های پردازش صوت، تصویر، متن، رمز، امنیت اطلاعات و مهندسی پزشکی - توسعه پل ارتباطی بین پژوهشگران و پژوهشگران کشور - کمک به ارتقای سطح دانش در زمینه های مورد بحث دو فصل نامه پردازش علائم و داده ها - تلاش در جهت رفع نیازهای علمی، پژوهشی و پژوهش های پژوهش های دانشجویان، پژوهشگران و	- پردازش صوت - پردازش تصویر - پردازش متن - پردازش رمز - امنیت اطلاعات - مهندسی پزشکی - پردازش علائم و داده ها

	پژوهشگران مراکز پژوهشی، پژوهش‌ها پژوهش‌های و دانشگاه		
۲	مطالعات الگوی پیشرفت اسلامی ایرانی	<p>- شناخت و معرفی و ارائه تحلیل، تبیین نظری و الگوهای عملی و همچنین نگرش‌های علمی در رابطه با موضوع الگوی پیشرفت</p> <p>- انتشار مقالات تخصصی علمی - پژوهشی در گرایش های مرتبط و ارائه نتایج پژوهش‌ها پژوهش‌ها نظری و عملی پژوهشگران.</p> <p>- فراهم کردن زمینه‌های گسترش دانش مطالعات الگوی پیشرفت اسلامی ایرانی در وجوه مختلف فرهنگی، اقتصادی و... منطبق با نگرش‌های تبیین کننده گفتمان انقلاب اسلامی، از طریق انتشار آخرین یافته‌های پژوهشی، تحلیل جدیدترین نظریه‌ها و معرفی آخرین نظر و پژوهش‌های منتشر شده مرتبط با موضوع</p>	<p>- الگوی پیشرفت اسلام</p> <p>- پژوهش‌ها پژوهش‌ها نظری</p> <p>- تحقیقات عملی</p> <p>- الگوی پیشرفت فرهنگی</p> <p>- الگوی پیشرفت اقتصادی</p> <p>- گفتمان انقلاب اسلامی</p>
۳	پژوهش‌های پولی - بانکی	<p>انتشار مطالعات و پژوهش‌ها پژوهش‌ها علمی اقتصادی در زمینه‌های اقتصاد کلان مرتبط با سیاست پولی، ارزی، بانکی و اعتباری و اقتصاد خرد مرتبط با بانکداری است. الزامی بر ارتباط مستقیم مطالعات با اقتصاد ایران وجود ندارد، اما انتظار می‌رود یافته‌های پژوهش‌ها به</p>	<p>- اقتصاد کلان</p> <p>- سیاست پولی</p> <p>- سیاست ارزی</p> <p>- سیاست بانکی و اعتباری</p>



<p>- اقتصاد خرد</p> <p>- بانکداری</p> <p>- اقتصاد ایران</p> <p>- پول و بانک</p> <p>- نظام بانکی</p> <p>- بانکداری مرکزی</p> <p>- بانکداری تجاری</p> <p>- اقتصاد پولی و بانکی</p>	<p>درک بهتر مسائل مبتلابه اقتصاد ایران به خصوص در حوزه پول و بانک کمک کند. بر این اساس، خط مشی مجله ترغیب پژوهشگران علاقه مند به مباحث اقتصادی مرتبط با حوزه های پولی و بانکی برای انجام مطالعات نظری و کاربردی در این زمینه ها و نیز آشنا کردن جامعه علمی کشور با پژوهش ها پژوهش ها ناب و عمیق اقتصادی در زمینه های مرتبط با نظام بانکی اعم از بانکداری مرکزی و تجاری است. نشر یافته ها و نظریه های جدید در زمینه های اقتصاد پولی و بانکی و فراهم آوردن زمینه های تبادل نظر در این حوزه از دیگر اهداف مجله تعریف شده است.</p>		
--	--	--	--

### ۳-۴. استخراج و آماده‌سازی عناصر مورد نیاز از محتوای مقالات نشریات

در این مرحله اطلاعات بالغ بر ۱۵۵۰۰۰ عنوان مقاله از نشریات مورد بررسی استخراج گردید. جدول ۲ شمایی از این

فرآیند را نشان می دهد.

جدول ۲- استخراج عناصر مقالات از هر نشریه

id	key_word	text	label
1	ابتکار و خلاقیت در علوم انسانی	طراحی مدل یکپارچه توسعه سطح نوآوری و تجاری سازی نوآوری	ابتکار و خلاقیت در علوم انسانی
2	ابتکار و خلاقیت در علوم انسانی	تاثیر خلاقیت و شیره های نوآورانه دانش پلیسی در پیش	ابتکار و خلاقیت در علوم انسانی
3	ابتکار و خلاقیت در علوم انسانی	رابطه هوش های چندگانه با ویژگی شخصیتی کارآفرینی و	ابتکار و خلاقیت در علوم انسانی
4	ابتکار و خلاقیت در علوم انسانی	طراحی مدل پرورش تفکر استراتژیک و خلاق در مدیران	ابتکار و خلاقیت در علوم انسانی
5	ابتکار و خلاقیت در علوم انسانی	تربیتی آموزش تفکر انتقادی در نگرش به خلاقیت و نشا	ابتکار و خلاقیت در علوم انسانی
6	ابتکار و خلاقیت در علوم انسانی	نوآوری و خلاقیت در اجرای اقدامات مدیریت منابع انسانی	ابتکار و خلاقیت در علوم انسانی
7	ابتکار و خلاقیت در علوم انسانی	الگوی کاربردی پرورش خلاقیت کودکان از طریق موسیقی	ابتکار و خلاقیت در علوم انسانی
8	ابتکار و خلاقیت در علوم انسانی	نگارش مشارکت الکترونیکی	ابتکار و خلاقیت در علوم انسانی
9	ابتکار و خلاقیت در علوم انسانی	خلاقیت	ابتکار و خلاقیت در علوم انسانی
10	ابتکار و خلاقیت در علوم انسانی	تأثیر بانگویی مشارکتی در پرورش خلاقیت دانش آموزان	ابتکار و خلاقیت در علوم انسانی
11	ابتکار و خلاقیت در علوم انسانی	مطالعه مشارکتی	ابتکار و خلاقیت در علوم انسانی
12	ابتکار و خلاقیت در علوم انسانی	تحزیم تحویل نقش واسطه ای بازخورد کار از سرپرست د	ابتکار و خلاقیت در علوم انسانی
13	ابتکار و خلاقیت در علوم انسانی	نوآوری در جذب منابع مالی: تأملی در چالش های مالی مؤثر	ابتکار و خلاقیت در علوم انسانی
14	ابتکار و خلاقیت در علوم انسانی	بررسی مؤلفه های مؤثر بر مصرف رسانه ای جوانان که	ابتکار و خلاقیت در علوم انسانی
15	ابتکار و خلاقیت در علوم انسانی	بررسی محتوای کتاب ریاضی پایه دهم دوره متوسطه	ابتکار و خلاقیت در علوم انسانی
16	ابتکار و خلاقیت در علوم انسانی	بررسی ویژگی های روانسنجی مقیاس دامنه های خلاقیت	ابتکار و خلاقیت در علوم انسانی
17	ابتکار و خلاقیت در علوم انسانی	بررسی تربیتی	ابتکار و خلاقیت در علوم انسانی
18	ابتکار و خلاقیت در علوم انسانی	تبیین علی خلاقیت: نقش سبک های دلپسندی و هیجانات مت	ابتکار و خلاقیت در علوم انسانی
19	ابتکار و خلاقیت در علوم انسانی	ساختار دانشکده معماری، زمینه پرورش خلاقیت دانشجو	ابتکار و خلاقیت در علوم انسانی
20	ابتکار و خلاقیت در علوم انسانی	اجرای استراتژی و عملکرد با توجه به نقش تحلیلگر های	ابتکار و خلاقیت در علوم انسانی

داده های بدست آمده از این مرحله مورد پیش پردازش و آماده سازی برای استفاده توسط ماشین قرار گرفت (جدول ۳).

جدول ۳- یکدست سازی و تعیین فراوانی عبارات

index	key_word
51	خلاقیت
8	نوآوری
5	دانشگاه آزاد اسلامی
3	دانش آموزان
3	نوآوری سازمانی
3	تفکر خلاق
3	تفکر انتقادی
3	رفتار نوآورانه
3	دانش آموزان
3	کارکنان
2	معلمان

#### ۴-۴. خوشه بندی داده‌ها

در این مرحله خوشه‌بندی داده‌ها انجام شد (جدول ۴).

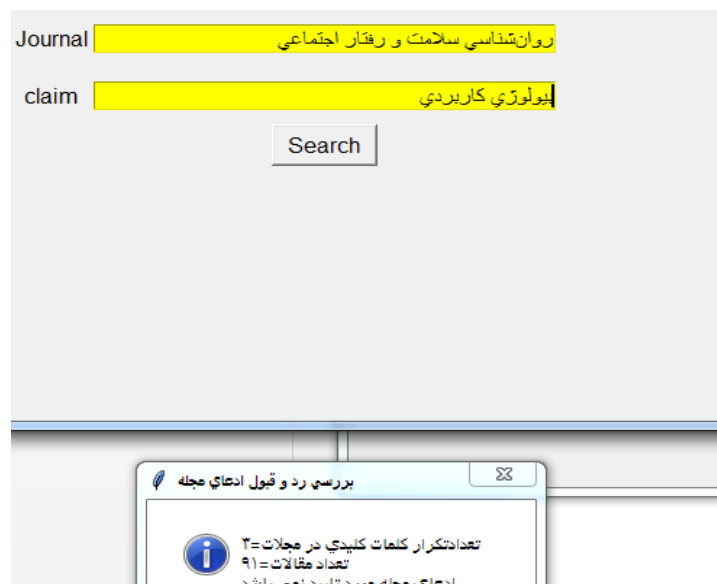
جدول ۴- خوشه بندی داده ها

انجمن‌ها	آموزش	دانشگاه	نوعی
هوش انجمن‌ها	دانش آموزش	دانشگاه آزاد اسلامی	نوعی
گزارشگری انجمن‌ها	دانش آموزش نیزهوش محاسبی	دانشگاه های بوقالی	نوعی مشاوره‌ای
افتتاح مسئولیت انجمن‌ها	دانش آموزش انجمنی	دانشگاه فرهنگیان	نوعی باز
کاهش ترس مبتنی بر شبکه های انجمن‌ها	دانش آموزش دختران	دانشگاه لرستان	نوعی منطقه ای
پیرامون شبکه انجمن‌ها	دانش آموزش نیزهوش	نوسه ی کارآموزی دانشگاهی	نوعی بر مسافتی
شبکه های انجمن‌ها	دانش آموزش نیروء متوسطه اول منطقه ۳ شهر اهواز	وزیرگی‌دان دانشگاهی	نوعی نظام رفاه
مهاجرت های انجمن‌ها		معداری دانشگاه	نوعی راهبردی
بند انجمن‌ها - فرهنگ		حمایت خانواده و دانشگاه	نوعی وب سایت ها
			نوعی کاری
			خلاقیت و نوعی
			نوعی انجمن‌ها
			صنعتی نوعی
			نوعی بر اعتماد مدیریت منابع انسانی
			نوعی بر جناب منابع مالی

#### ۴-۵. شباهت سنجی

در مرحله پایانی شباهت سنجی انجام شد (عکس ۱).

عکس ۱- شباهت سنجی حوزه موضوعی نشریه با محتوای مقالات



## ۵-۵ یافته ها

میزان همخوانی یا شباهت موضوعات مندرج در وبگاه هر نشریه و محتوای مقالات به شرح جداول ۷ تا ۱۲ محاسبه گردید.

جدول ۷- میزان همخوانی محتوای مقالات نشریات حوزه دامپزشکی با حوزه موضوعی درج شده در وبگاه

ردیف	حوزه موضوعی کلان	حوزه موضوعی میانی	عنوان نشریه	درصد همخوانی
۱	دامپزشکی	دامپزشکی	التیام	۶۰
۲	دامپزشکی	دامپزشکی	بهداشت مواد غذایی	۷۰

همانگونه که مشاهده می شود در نشریات مورد بررسی حوزه موضوعی دامپزشکی، همخوانی به میزان ۶۰ و ۷۰ درصد رعایت شده است.

جدول ۸- میزان همخوانی محتوای مقالات نشریات حوزه علوم انسانی با حوزه موضوعی درج شده در وبگاه

ردیف	حوزه موضوعی کلان	حوزه موضوعی میانی	عنوان نشریه	درصد همخوانی
۱	علوم انسانی	علوم تربیتی	ابتکار و خلاقیت در علوم انسانی	۹۵
۲	علوم انسانی	جغرافیا	اقتصاد فضا و توسعه روستایی	۹۵
۳	علوم انسانی	روانشناسی	پژوهش های روانشناسی اجتماعی	۹۵
۴	علوم انسانی	علوم اجتماعی	اعتیاد پژوهی	۹۵
۵	علوم انسانی	علوم سیاسی	پژوهش های انقلاب اسلامی	۹۵
۶	علوم انسانی	مدیریت	اندیشه آماد	۹۵
۷	علوم انسانی	علوم قرآن و حدیث	آموزه های تربیتی در قرآن و حدیث	۸۰
۸	علوم انسانی	زبان و ادبیات	ادبیات عرفانی	۷۰
۹	علوم انسانی	اخلاق	پژوهش های اخلاقی	۷۰
۱۰	علوم انسانی	ادیان، مذاهب و عرفان	اندیشه نوین دینی	۷۰
۱۱	علوم انسانی	تاریخ	پژوهش های باستان شناسی ایران	۷۰
۱۲	علوم انسانی	تربیت بدنی	پژوهش در توانبخشی ورزشی	۷۰
۱۳	علوم انسانی	جغرافیا	پژوهش های جغرافیای طبیعی	۷۰

۷۰	پژوهش های ژئومورفولوژی کمی	جغرافیا	علوم انسانی	۱۴
۷۰	بررسی های حسابداری و حسابرسی	حسابداری	علوم انسانی	۱۵
۷۰	اندازه گیری تربیتی	روانشناسی	علوم انسانی	۱۶
۷۰	پژوهش های ادبیات تطبیقی	زبان و ادبیات	علوم انسانی	۱۷
۷۰	پژوهش های اطلاعاتی و جنایی	علوم اجتماعی	علوم انسانی	۱۸
۷۰	آموزش محیط زیست و توسعه پایدار	علوم تربیتی	علوم انسانی	۱۹
۷۰	پژوهش های روابط بین الملل	علوم سیاسی	علوم انسانی	۲۰
۷۰	آموزه های حقوق کیفری	فقه و حقوق	علوم انسانی	۲۱
۷۰	پژوهش های علم و دین	فلسفه و کلام	علوم انسانی	۲۲
۷۰	آینده پژوهی ایران	مدیریت	علوم انسانی	۲۳
۶۰	اخلاق و حیانی	اخلاق	علوم انسانی	۲۴
۶۰	ادیان و عرفان	ادیان، مذاهب و عرفان	علوم انسانی	۲۵
۶۰	پژوهش های ادیانی	ادیان، مذاهب و عرفان	علوم انسانی	۲۶
۶۰	اقتصاد انرژی ایران	اقتصاد	علوم انسانی	۲۷
۶۰	اقتصاد مالی	اقتصاد	علوم انسانی	۲۸
۶۰	اقتصاد و برنامه ریزی شهری	اقتصاد	علوم انسانی	۲۹
۶۰	اقتصاد و تجارت نوین	اقتصاد	علوم انسانی	۳۰
۶۰	برنامه ریزی و بودجه	اقتصاد	علوم انسانی	۳۱
۶۰	پژوهش های تاریخی ایران و اسلام	تاریخ	علوم انسانی	۳۲
۶۰	پژوهش در طب ورزشی و فناوری	تربیت بدنی	علوم انسانی	۳۳
۶۰	پژوهش در مدیریت ورزشی و رفتار حرکتی	تربیت بدنی	علوم انسانی	۳۴
۶۰	آمایش جغرافیایی فضا	جغرافیا	علوم انسانی	۳۵
۶۰	برنامه ریزی توسعه کالبدی	جغرافیا	علوم انسانی	۳۶
۶۰	پژوهش های اقلیم شناسی	جغرافیا	علوم انسانی	۳۷
۶۰	پژوهش های حسابداری مالی و حسابرسی	حسابداری	علوم انسانی	۳۸
۶۰	پژوهش در سلامت روانشناختی	روانشناسی	علوم انسانی	۳۹

۶۰	ادب فارسی	زبان و ادبیات	علوم انسانی	۴۰
۶۰	ادبیات پارسی معاصر	زبان و ادبیات	علوم انسانی	۴۱
۶۰	ادبیات تطبیقی	زبان و ادبیات	علوم انسانی	۴۲
۶۰	ادبیات و زبانهای محلی ایران زمین	زبان و ادبیات	علوم انسانی	۴۳
۶۰	پژوهش های زبان شناسی	زبان و ادبیات	علوم انسانی	۴۴
۶۰	بازیابی دانش و نظام های معنایی	علم اطلاعات و دانش شناسی	علوم انسانی	۴۵
۶۰	بررسی مسایل اجتماعی ایران	علوم اجتماعی	علوم انسانی	۴۶
۶۰	برنامه ریزی رفاه و توسعه اجتماعی	علوم اجتماعی	علوم انسانی	۴۷
۶۰	پژوهش های انسان شناسی ایران	علوم اجتماعی	علوم انسانی	۴۸
۶۰	الگوی پیشرفت اسلامی ایرانی	علوم تربیتی	علوم انسانی	۴۹
۶۰	اندیشه های نوین تربیتی	علوم تربیتی	علوم انسانی	۵۰
۶۰	پژوهش در مسائل تعلیم و تربیت	علوم تربیتی	علوم انسانی	۵۱
۶۰	پژوهش های برنامه درسی	علوم تربیتی	علوم انسانی	۵۲
۶۰	پژوهش سیاست نظری	علوم سیاسی	علوم انسانی	۵۳
۶۰	پژوهش های سیاست اسلامی	علوم سیاسی	علوم انسانی	۵۴
۶۰	پژوهش دینی	علوم قرآن و حدیث	علوم انسانی	۵۵
۶۰	پژوهش های تفسیر تطبیقی	علوم قرآن و حدیث	علوم انسانی	۵۶
۶۰	پژوهش حقوق خصوصی	فقه و حقوق	علوم انسانی	۵۷
۶۰	پژوهش حقوق کیفری	فقه و حقوق	علوم انسانی	۵۸
۶۰	اندیشه دینی	فلسفه و کلام	علوم انسانی	۵۹
۶۰	پژوهش های اعتقادی کلامی	فلسفه و کلام	علوم انسانی	۶۰
۶۰	بورس اوراق بهادار	مالی	علوم انسانی	۶۱
۶۰	اندیشه مدیریت راهبردی	مدیریت	علوم انسانی	۶۲
۶۰	آینده پژوهی مدیریت	مدیریت	علوم انسانی	۶۳
۴۰	اقتصاد باثبات	اقتصاد	علوم انسانی	۶۴
۴۰	پژوهش های تاریخی	تاریخ	علوم انسانی	۶۵
۴۰	پژوهش های حسابداری مالی	حسابداری	علوم انسانی	۶۶

۴۰	پژوهش های روانشناسی بالینی و مشاوره	روانشناسی	علوم انسانی	۶۷
۴۰	ادبیات عرفانی و اسطوره شناختی	زبان و ادبیات	علوم انسانی	۶۸
۴۰	آینه میراث	علم اطلاعات و دانش شناسی	علوم انسانی	۶۹
۴۰	برنامه ریزی درسی نظریه و عمل	علوم تربیتی	علوم انسانی	۷۰
۴۰	پژوهش در یادگیری آموزشگاهی و مجازی	علوم تربیتی	علوم انسانی	۷۱
۴۰	آینه معرفت	فلسفه و کلام	علوم انسانی	۷۲

همانگونه که مشاهده می شود در نشریات مورد بررسی حوزه علوم انسانی، همخوانی موضوعی در ۸ درصد نشریات ۹۵ درصد و پس از آن به ترتیب در ۱ درصد نشریات ۸۰ درصد، در ۲۲ درصد نشریات ۷۰ درصد، در ۵۶ درصد نشریات ۶۰ درصد و در ۱۳ درصد نشریات ۴۰ درصد رعایت شده است.

جدول ۹- میزان همخوانی محتوای مقالات نشریات حوزه علوم پایه با حوزه موضوعی درج شده در وبگاه

ردیف	حوزه موضوعی کلان	حوزه موضوعی میانی	عنوان نشریه	درصد همخوانی
۱	علوم پایه	زمین شناسی	اکتشاف و تولید نفت و گاز	۸۰
۲	علوم پایه	زیست شناسی	اکوفیتوشیمی گیاهان دارویی	۷۰
۳	علوم پایه	محیط زیست	انسان و محیط زیست	۷۰
۴	علوم پایه	زمین شناسی	بلورشناسی و کانی شناسی ایران	۷۰
۵	علوم پایه	زیست شناسی	اقیانوس شناسی	۶۰
۶	علوم پایه	آمار	اندیشه آماری	۴۰
۷	علوم پایه	فیزیک	پژوهش سیستم های بس ذره ای	۴۰
۸	علوم پایه	ریاضی	پژوهش های ریاضی	۴۰
۹	علوم پایه	شیمی	پژوهش های شیمی (CR)	۴۰

همانگونه که مشاهده می شود در نشریات مورد بررسی حوزه علوم پایه، میزان همخوانی موضوعی در ۱۱ درصد نشریات ۸۰ درصد، در ۳۳ درصد نشریات ۷۰ درصد، در ۱۱ درصد نشریات ۶۰ درصد و در ۴۵ درصد نشریات ۴۰ درصد رعایت شده

است.

جدول ۱۰- میزان همخوانی محتوای مقالات نشریات حوزه فنی و مهندسی با حوزه موضوعی درج شده در وبگاه

ردیف	حوزه موضوعی کلان	حوزه موضوعی میانی	عنوان نشریه	درصد همخوانی
۱	فنی و مهندسی	میان رشته ای	انرژی ایران	۷۰
۲	فنی و مهندسی	میان رشته ای	پدافند الکترونیکی و سایبری	۷۰
۳	فنی و مهندسی	مهندسی شیمی	پژوهش نفت	۷۰
۴	فنی و مهندسی	پلیمر	علوم و فناوری نساجی	۶۰
۵	فنی و مهندسی	مکانیک	انرژی های تجدیدپذیر و نو	۶۰
۶	فنی و مهندسی	پلیمر	بسپارش	۶۰
۷	فنی و مهندسی	برق و کامپیوتر	الکترو مغناطیس کاربردی	۴۰
۸	فنی و مهندسی	برق و کامپیوتر	پردازش سیگنال پیشرفته	۴۰
۹	فنی و مهندسی	برق و کامپیوتر	پردازش علائم و داده ها	۴۰
۱۰	فنی و مهندسی	عمران	پژوهش های زیرساخت های عمرانی	۴۰
۱۱	فنی و مهندسی	نفت	پژوهش های سیاستگذاری و برنامه ریزی انرژی	۴۰
۱۲	فنی و مهندسی	عمران	اساس	۴۰
۱۳	فنی و مهندسی	مواد و متالورژی	علوم و مهندسی سطح	۴۰

همانگونه که مشاهده می شود در نشریات مورد بررسی حوزه فنی و مهندسی، میزان همخوانی موضوعی در ۲۳ درصد نشریات ۷۰ درصد، در ۲۳ درصد نشریات ۶۰ درصد و در ۵۴ درصد نشریات ۴۰ درصد رعایت شده است.

جدول ۱۱- میزان همخوانی محتوای مقالات نشریات حوزه کشاورزی و منابع طبیعی با حوزه موضوعی مندرج در وبگاه

ردیف	حوزه موضوعی کلان	حوزه موضوعی میانی	عنوان نشریه	درصد همخوانی
۱	کشاورزی و منابع طبیعی	زراعت	اکوفیزیولوژی گیاهی	۹۵



۹۵	به زراعی کشاورزی	زراعت	کشاورزی و منابع طبیعی	۲
۹۵	پژوهش آب ایران	آب و خاک	کشاورزی و منابع طبیعی	۳
۷۰	اقتصاد کشاورزی و توسعه	اقتصاد کشاورزی	کشاورزی و منابع طبیعی	۴
۷۰	بوم شناسی کشاورزی	اکولوژی	کشاورزی و منابع طبیعی	۵
۷۰	بیماریهای گیاهی	گیاهپزشکی	کشاورزی و منابع طبیعی	۶
۷۰	پژوهش های حبوبات ایران	میان رشته‌ای	کشاورزی و منابع طبیعی	۷
۶۰	اکوبیولوژی تالاب	محیط زیست	کشاورزی و منابع طبیعی	۸
۶۰	آبزیان زینتی	شیلات	کشاورزی و منابع طبیعی	۹
۶۰	آفات و بیماریهای گیاهی	گیاهپزشکی	کشاورزی و منابع طبیعی	۱۰
۶۰	بهره برداری و پرورش آبزیان	شیلات	کشاورزی و منابع طبیعی	۱۱
۶۰	بوم شناسی جنگل های ایران	جنگلداری	کشاورزی و منابع طبیعی	۱۲
۶۰	پژوهش های ژنتیک گیاهی	بیوتکنولوژی و ژنتیک گیاهی	کشاورزی و منابع طبیعی	۱۳
۴۰	پژوهش های آبخیزداری	آبخیزداری	کشاورزی و منابع طبیعی	۱۴
۴۰	پژوهش های صنایع غذایی	صنایع غذایی	کشاورزی و منابع طبیعی	۱۵

۱۶	کشاورزی و منابع طبیعی	اکولوژی	بوم‌شناسی کاربردی	۴۰
----	-----------------------	---------	-------------------	----

همانگونه که مشاهده می‌شود در نشریات مورد بررسی حوزه کشاورزی و منابع طبیعی، میزان همخوانی موضوعی در ۱۹ درصد نشریات ۹۵ درصد، در ۲۵ درصد نشریات ۷۰ درصد، در ۳۷ درصد نشریات ۶۰ درصد و در ۱۹ درصد نشریات ۴۰ درصد رعایت شده است.

جدول ۱۲- میزان همخوانی محتوای مقالات نشریات حوزه هنر و معماری با حوزه موضوعی مندرج در وبگاه نشریه

ردیف	حوزه موضوعی کلان	حوزه موضوعی میانی	عنوان نشریه	درصد همخوانی
۱	هنر و معماری	هنر و معماری	اثر	۶۰
۲	هنر و معماری	هنر و معماری	اندیشه معماری	۶۰
۳	هنر و معماری	هنر و معماری	باغ نظر	۶۰
۴	هنر و معماری	هنر و معماری	مبانی نظری هنرهای تجسمی	۶۰

همانگونه که مشاهده می‌شود میزان همخوانی موضوعی در همه نشریات مورد بررسی حوزه هنر و معماری به میزان ۶۰ درصد رعایت شده است.

جدول ۱۳- فراوانی نشریات بر اساس میزان همخوانی محتوای مقالات نشریه با حوزه موضوعی مندرج در وبگاه به تفکیک

حوزه موضوعی

حوزه موضوعی	میزان همخوانی				
	%۹۵	%۸۰	%۷۰	%۶۰	%۴۰
فراوانی (تعداد نشریه)					
علوم انسانی	۶	۱	۱۶	۴۰	۹
عوم پایه	۰	۱	۳	۱	۴
فنی و مهندسی	۰	۰	۳	۳	۷
کشاورزی و منابع طبیعی	۳	۰	۴	۶	۳
دامپزشکی	۰	۰	۱	۱	۰
هنر و معماری	۰	۰	۰	۴	۰

همانگونه که مشاهده می شود از مجموع نشریات مورد بررسی ۸ درصد نشریات میزان همخوانی یا تطابق محتوای مقالات با دامنه موضوعی مندرج در وب سایت را به میزان ۹۵ درصد رعایت کرده بودند. کمترین میزان همخوانی ۴۰ درصد است که در ۲۰ درصد نشریات مورد بررسی وجود داشت. در این میان میزان همخوانی در ۴۷ درصد نشریات ۶۰ درصد، ۲۳ درصد نشریات ۷۰ درصد و در ۲ درصد نشریات به میزان ۸۰ درصد رعایت شده بود.

# فصل پنجم

## نتیجه‌گیری و ارائه پیشنهادها

در بخش نخست ابتدا نتایج حاصل از تجزیه و تحلیل داده‌ها ارائه می‌گردد، سپس پیشنهادهای پژوهش تبیین

می‌شود.

## ۵-۲. بحث و نتیجه‌گیری

پژوهشگران، از مهمترین ارکان دانشی تأثیرگذار در نظام علم و فناوری محسوب می‌شوند. این افراد، علاوه بر فراهم آوردن موجبات توسعه و ترقی نظام آموزش عالی و ایجاد زمینه‌های لازم برای توسعه و پویایی علمی و فرهنگی جامعه، با ارائه دستاوردهای ارزشمند علمی و فناورانه خود، چرخ علم و فناوری را به گردش درمی‌آورند و سکان داران اصلی توسعه علم و فناوری محسوب می‌شوند.

باتوجه به اینکه نشریات علمی آخرین دستاوردهای علمی و افق‌های پژوهشی را در کوتاه‌ترین زمان ممکن منتشر می‌کنند و همچنین فرآیند داوری دقیقی را طی می‌کنند؛ یکی از مهمترین قالب‌های انتشار تولیدات علمی هستند. بنابراین انتظار می‌رود مقالات علمی منتشرشده در مجلات معتبر علمی، نماینده مناسبی از تولیدات پژوهشی یک رشته علمی باشند و تحلیل و مطالعه محتوای این مقالات بتواند در برنامه‌ریزی‌ها و سیاست‌گذاری‌های آینده آن رشته علمی نقش بسزایی داشته باشد. در حال حاضر نزدیک به ۱۵۰۰ نشریه علمی توسط وزارت علوم، تحقیقات و فناوری اعتبارگذاری شده که بیش از ۸۵ درصد آنها دارای رتبه‌های الف و ب هستند که این امر نشانگر اعتبار علمی بالای آنها دارد. اما این افزایش تعداد نشریه‌های علمی ایجاب می‌کند که برای حفظ و توسعه کیفی آنها تلاش شود. از این‌رو، ارزیابی این نشریات از جهات گوناگون، می‌تواند شاخصی در جهت پایش وضعیت موجود و بهبود جایگاه نشریات کشور باشد. اما از طرفی با توجه به محدودیت‌های زمانی و همچنین حجم اطلاعات، جستجو و تحلیل دستی منابع علمی عملاً امکان‌پذیر نیست و آنچه حائز اهمیت است استفاده از روش‌های یادگیری ماشین مانند متن‌کاوی برای پردازش داده‌های بزرگ است که در پژوهش حاضر مورد توجه قرار گرفته است. همچنین انتخاب نشریه جهت انتشار مقاله، متأثر از عوامل گوناگونی است که در این میان عامل ارتباط دامنه و حوزه

موضوعی نشریه با مقاله از اهمیت ویژه‌ای برخوردار است؛ زیرا ارسال مقالات حاصل از یک پژوهش، به نشریاتی که ربط کمتری با محتوا و نیاز خوانندگان دارند باعث اتلاف وقت پژوهشگر می‌شود. همچنین حوزه سردبیری نشریه نیز مدت زمانی را صرف بررسی و احتمالاً ارسال مقاله به داوری نموده و اگر طی هر یک از مراحل مقاله رد شود این فرآیند باید مجدداً برای نشریه‌ای دیگر تکرار گردد. از اینرو این نکته که نشریات تا چه اندازه در اعلام حوزه‌های موضوعی که مقالات را می‌پذیرند شفاف عمل کرده و به آن پایبند بوده‌اند حائز اهمیت است. گاهی موضوعات اعلامی از سوی نشریات بسیار کلی بوده و کاربرد را در انتخاب نشریه مرتبط با دشواری مواجه می‌کند.

در این راستا پژوهش حاضر با هدف تحلیل و تطبیق میزان همخوانی دامنه و اهداف موضوعی نشریات فارسی وزارت عتف با محتوای مقالات منتشر شده در آنها طی بازه زمانی سه ساله ۱۳۹۷-۱۳۹۹ انجام شد تا مشخص گردد که نشریات به چه میزان در اعلام و رعایت دامنه و حوزه‌های موضوعی که در آن مقاله می‌پذیرند رویکرد شفافی داشته‌اند. در این پژوهش مجموعاً ۱۱۶ نشریه و بالغ بر ۱۵۰۰۰ مقاله مورد بررسی و تحلیل قرار گرفت و از روش فضای برداری به منظور تعیین میزان شباهت استفاده شد. از نظر دوره انتشار، نشریات مورد مطالعه به صورت فصلنامه و دو فصلنامه منتشر می‌شوند. همچنین نشریات مورد مطالعه از نظر پوشش موضوعی به ترتیب از حوزه‌های علوم انسانی، فنی و مهندسی، علوم پایه، هنر و معماری، کشاورزی و منابع طبیعی و دامپزشکی می‌باشند.

یافته‌های پژوهش در میزان رعایت و پایبندی نشریات به اهداف و دامنه موضوعی اعلام شده و محتوای مقالات منتشر در هر یک از آنها نشان داد تنها در ۸ درصد نشریات مورد بررسی به میزان ۹۵ درصد تطابق وجود داشت و نویسندگان را به طور شفاف از موضوعات مورد پذیرش نشریه آگاه کرده و خود نیز پایبند به حوزه‌های موضوعی اعلام شده بوده‌اند. در ۱۹ درصد نشریات اعلام موضوعات بسیار کلی بوده و تطابق نیز تنها در حد ۴۰ درصد می‌باشد.

در مطالعات پیشین نیز پژوهشگران به این نتیجه رسیده‌اند که درصدی از مقالات منتشر شده در نشریات با موضوع های اعلامی همخوانی ندارند از جمله نتایج مطالعه افصلی پور و جمالی مهموئی در بررسی میزان تخصص گرایی مجله های حوزه علم اطلاعات و دانش شناسی فارسی طی سالهای ۱۳۸۷ الی ۱۳۹۱ نشان داد ۷ درصد مقالات با موضوعای غیرمرتبط با

علم اطلاعات در این نشریات منتشر شده اند. همچنین تحلیلی محتوایی مقالات نشریه اخلاق در علوم و فناوری نشان داد موضوع هایی از جمله کدهای اخلاقی، رهبری اخلاقی، تعهد سازمانی و .. که بیشترین تعداد مقالات را طی سالهای ۸۵ تا ۹۴ به خود اختصاص داده اند در اداف و دامنه موضوعی نشریه ذکر نشده اند (غلامی، ۱۳۹۵). بیشترین سهم موضوعی مقالات در نشریه آموزش و توسعه منابع انسانی طی سالهای ۹۳ تا ۹۷ مربوط به موضوع اندازه گیری اثرات و نتایج برنامه های آموزش و توسعه و سپس مباحثی از قبیل سنجش اثربخشی دوره های آموزشی، آموزش و یادگیری الکترونیکی، مدیریت دانایی بود (روحانی راد و همکاران، ۱۳۹۹) که در اهداف و دامنه موضوعی نشریه هیچ اشاره ای بدانها نشده است.

هرچند در برخی نشریات مانند نشریه مطالعات مدیریت گردشگری، گرایشهای موضوعی مقالات طی سالهای ۱۳۹۱ تا ۱۳۹۸ در نشریه مطالعات مدیریت گردشگری نشان داد مدیریت توریسم، تبلیغات توریسم، اقتصاد توریسم و مدیریت هتلها به ترتیب دارای بیشترین تعداد مقاله هستند (رضانیا و همکاران، ۱۴۰۰) و با دامنه و حوزه موضوعی مدرج در وبگاه نشریه مطابقت دارد. همچنین در نشریه علمی راهبرد با توجه به یافته های پژوهش آخوندی (۱۳۹۸) همخوانی لازم به میزان مطلوبی رعایت شده است. بر اساس یافته های مطالعه خاشعی و همکاران (۱۳۹۸) و مطابق با مقالات منتشر شده در این نشریه طی سالهای ۹۰ تا ۹۶، نشریه مطالعات مدیریت بهبود و تحول نیز به صورت شفاف موضوعات مورد پذیرش نشریه را در وبگاه اطلاع رسانی نموده است.

در پیوند با مطالب پیش گفته می توان اینچنین مطرح کرد که در پیش گرفتن رویکرد شفاف موضوعی از سوی نشریات می تواند دستاوردهای سودمندی برای جامعه دانشگاهی داشته باشد. با مشخص بودن گرایش های موضوعی نشریه های علمی، هویت بارزتری برای هر نشریه ایجاد می شود؛ به طوری که هر متخصصی می تواند انتظار داشته باشد که نشریه یا نشریه های خاصی، زمینه های مطالعاتی وی را پوشش خواهد داد. تخصص گرایی فرآیندی علمی و عملی است که می توان به یاری آن، پژوهش های اعضای هیات علمی و پژوهشگران کشور را معطوف به حوزه های تخصصی کرد. در تعریف حوزه تخصص باید گفت که حوزه تخصص به آن حوزه یا گرایش خاصی اطلاق می گردد که مشتمل بر مجموعه مسائل نظام مندی است که نقشه راه فعالیت نشریه را در بازه زمانی طولانی تعیین می نماید و نتیجه آن به فراوری دانش انباشته و رسیدن به

نظریه های جدید در آن حوزه علمی منتهی می گردد. بدین ترتیب، نشریه و نویسندگان آن به مرجعی در جامعه علمی تبدیل می گردند.

### ۵-۳ پیشنهادهای کاربردی

پیشنهاد می شود گروه دبیران، ضمن بررسی منظم روند انتشار مقالات در مجلات خود، از سیر موضوعی مقالات همراستا با پژوهش های جهانی اطمینان حاصل نمایند و ضمن انعطاف پذیری در موضوعات مورد پذیرش، از انحراف مجله تحت مدیریت خود از خط مشی آن جلوگیری کنند. همچنین با توجه به گرایش های موضوعی در گذر زمان، می توانند از همپوشانی موضوعی مجلات ممانعت به عمل آورده و موضوعات مغفول را در دامنه موضوعی مجلات بگنجانند. در این راستا به منظور تعدیل وضعیت موجود پیشنهاد میشود ضمن استخراج بیشترین واژه های کلیدی به کار رفته شده در مقالات، فهرست حوزه های موضوعی مندرج در وبگاه نشریه را روزآمد نمایند. این امر علاوه بر اینکه از ارسال مقالات غیر مرتبط با موضوع نشریه و اتلاف زمان و هزینه پژوهشگر و دست اندرکارن نشریه جلوگیری می شود، تخصصی شدن و عدم پراکندگی موضوعی نشریه را نیز به دنبال خواهد داشت. همچنین پیشنهاد می شود پژوهشگران نیز پیش از انجام پژوهش، ضمن بررسی موضوعات روز، به مجلاتی که در زمینه های موضوعی مرتبط و تخصصی اقدام به انتشار مقاله می کنند، توجه خاص تری داشته باشند.

### ۵-۴ پیشنهادهای پژوهش

- پیشنهاد می شود پژوهشی مشابه در خصوص مجلات ایرانی بین المللی انجام و نتایج پژوهش حاضر با پژوهش پیشنهادی مقایسه گردد.
- پیشنهاد می شود پژوهش با بررسی دیگر نشریات مجدداً تکرار و در صورت کسب نتایج مشابه، شفاف سازی ذکر دامنه موضوعی نشریات از سوی کمیسیون به آنها ابلاغ گردد



## فهرست منابع

- آخوندی، عباس. (۱۳۹۸). تحلیل محتوای مقالات فصلنامه راهبرد. فصلنامه علمی راهبرد، ۲۸(۲): ۱۶۷-۱۹۶. doi: 20.1001.1.10283102.1398.28.2.7.0
- استادزاده، زیبا؛ داوودی، حسین؛ حیدری، حسن؛ میرمهدی، سیدرضا (۱۳۹۷). تحلیل محتوای مقالات «دوفصلنامه مطالعات اسلام و روان‌شناسی (۱۳۹۵-۱۳۸۶)». مطالعات اسلام و روان‌شناسی. ۱۲ (۲۳): ۱۴۱-۱۵۸.
- افضل‌پور حدیثه، جمالی مهموئی حمید رضا (۱۳۹۳). میزان تخصص‌گرایی مجله‌های علمی- پژوهشی علم اطلاعات و دانش‌شناسی فارسی بر اساس موضوع مقالات منتشر شده بین سال‌های ۱۳۸۷ تا ۱۳۹۱. تحقیقات اطلاع‌رسانی و کتابخانه‌های عمومی. ۲۰ (۳): ۳۹۷-۴۲۲.
- اللهیاری، محمدصادق؛ دقیقی ماسوله، زهرا؛ افتخاری، هاجر (۱۳۹۴). تحلیل محتوای مقالات تخصصی علوم ترویج و آموزش کشاورزی در ایران (۱۳۹۲-۱۳۸۸). علوم ترویج و آموزش کشاورزی، ۱۱(۱): ۲۲۹-۲۴۹.
- باغ محمد، مریم، منصوری، علی، چشمه سهرابی، مهرداد. (۱۴۰۱). بررسی توسعه و روند موضوعی حوزه علم اطلاعات و دانش‌شناسی بر اساس مدل موضوعی LDA. پژوهشنامه پردازش و مدیریت اطلاعات، ۳۶(۲): ۲۹۷-۲۹۸. doi: 10.35050/JIPM010.2020.001۳۲۸
- توکلی زاده راوری، محمد؛ دهقانی، نجابتیان؛ سهیلی. (۱۳۹۴). تحلیل محتوای مقالات فارسی نشریات علمی ایران در زمینه ازدواج و طلاق با روش خوشه‌بندی سلسله مراتبی. فصلنامه فرهنگی تربیتی زنان و خانواده، ۱۰(۳۲).
- جانانی، پیمان؛ رودباری؛ تهمتن؛ صدقی (۱۳۹۱). تحلیل محتوای مقالات نشریات دانشکده‌های پرستاری و مامائی دانشگاه‌های علوم پزشکی ایران. مراقبتهای پرستاری و مامایی، ۲(۱): ۵۳-۶۱.
- حاتمی ناغانی، بهمنو عباسی، مسعود. (۱۳۹۵). تحلیل محتوایی مقالات علمی با استفاده از متن کاوی. مطالعات مدیریت کسب و کار هوشمند 5(18), 137-167. doi: 10.22054/ims.2017.7014
- حری، عباس. (۱۳۷۴). بررسی رابطه میان مجلات منتشره حوزه‌های تخصصی و ارتقاء علمی متخصصان کشور.

پژوهش‌های کاربردی روانشناختی ۱ (۲): ۲۳-۴۴.

- حسینی نسب، اعظم؛ محمدی، مهدی؛ طالعی، عبدالحسین (۱۳۹۵). تحلیل موضوعی مقالات علمی پژوهشی حوزه علوم قرآن؛ نشر یافته بین سال‌های ۱۳۸۶-۱۳۹۰. دو فصلنامه کتاب‌قیم. ۶ (۱۵): ۱۴۱-۱۶۸.
- خاشعی ورنامخواستی، وحید؛ طیبی ابوالحسنی، سیدامیرحسین؛ اسدی خانقاه، شیرین (۱۳۹۸). تحلیل محتوای مقالات فصلنامه مطالعات مدیریت (بهبود و تحول) طی دوره ۷ ساله. (۱۳۹۰-۹۶). مطالعات مدیریت (بهبود و تحول): ۲۸ (۹۳): ۱۲۹-۱۶۰.
- داستانی، میثم؛ موسوی چلک، افشین؛ ضیائی، ثریا؛ دلقندی، فائزه. (۱۳۹۹). تجزیه و تحلیل موضوعی مقالات منتشرشده ی کتابداری و اطلاع‌رسانی پزشکی در ایران با استفاده از فنون متن‌کاوی. تصویر سلامت ۱۱ (۴): ۳۵۵-  
doi: 10.34172/doh.2020.43.۳۶۷
- دهدشتی شاهرخ، زهره. (۱۳۹۹). تحلیل ساختار محتوایی فصلنامه مطالعات مدیریت گردشگری با استفاده از تکنیک متن‌کاوی. مطالعات مدیریت گردشگری ۱۵ (۱۵): ۹۷-۱۲۷. doi: 10.22054/tms.2020.42515.2141
- رضایانیا، درنا، مقتنی پور، مجید رضا و ظفرمند، سید جواد. (۱۴۰۰). تحلیل محتوای مقالات فصلنامه مطالعات مدیریت گردشگری (۱۳۹۱-۱۳۹۸) بر مبنای الگوی مشارکت نویسندگان و ویژگی‌های روش‌شناسی، استنادی و عملکردی مقالات. فصلنامه بازیابی دانش و نظام‌های معنایی. doi: 169-201. 8(28),  
10.22054/jks.2021.60526.1433
- سلک، محسن، بزرگی، اشرف السادات. (۱۳۸۹). تحلیل محتوای مقالات منتشر شده در دو نشریه فصلنامه "کتابداری و اطلاع‌رسانی" و "فصلنامه کتاب" در سال‌های ۱۳۸۵ و ۱۳۸۶. دانش‌شناسی. ۳ (۱۰): ۲۵-۴۰.
- شکوهیان، محبوبه؛ عاصمی، عاصفه؛ شعبانی، احمد؛ چشمه سهرابی، مظفر. (۱۴۰۱). ارائه مدل دسته‌بندی موضوعی تولیدات علمی حوزه سلامت با استفاده از روش‌های متن‌کاوی. پژوهشنامه پردازش و مدیریت اطلاعات ۳۵ (۲):  
doi: 10.35050/JIPM010.2020.061.۵۷۴-۵۵۳

- صفوی جهرمی، گلایل، طباطبائیان، سید حبیب الله، حنفی زاده، پیام، حاجی میرزائی، حامد. (۱۳۹۹). نقشه موضوعی مقالات حوزه تولید محتوای دیجیتال برای کودکان و نوجوانان. مطالعات کتابداری و علم اطلاعات. ۱۴ (۳).
- غلامی، طاهره. (۱۳۹۵). تحلیل محتوای مقالات فصلنامه اخلاق در علوم و فناوری (۱۳۸۵ - ۱۳۹۴). اخلاق در علوم و فناوری. ۱۱ (۱): ۲۹-۳۸.
- فتاحی، رحمت الله. (۱۳۹۴) تخصصی کردن گرایش موضوعی مجله های علمی و چالشهای آن. پژوهشنامه پردازش و مدیریت اطلاعات. ۳۰ (۳): ۶۰۱-۶۰۲.
- قتاد، مصطفی؛ عرب مازار یزدی، محمد؛ صفرزاده بندری؛ محمدحسین؛ حصارزاده، رضا. (۱۴۰۲). کاربرد فنون متن کاوی در تحلیل جریان موضوعی مقالات منتشره در مجلات حسابداری ایران. پژوهش های حسابداری مالی و حسابرسی ۱۵ (۵۸): ۱-۳۸.
- کاظمی، عبدالحسن. (۱۳۸۸) مشکلات سنجش تولیدات علمی کشور. مدیریت اطلاعات سلامت. فروردین ۱۳۸۸.
- Abel, R.E. & Newlin, L. W. (2002). *Scholarly Publishing: Books Journals, Publishers, and Libraries in the Twentieth Century*. N.Y.: ISBN 0-471-21929-0.Y.: Wiley.
- Abramson, D., Lees, M., Krzhizhanovskaya, V. V., Dongarra, J. J., & Sloot, P. M. (2014). Big Data Meets Computational Science, Preface for ICCS. In *ICCS*. 1-7.
- Alghamdi, R., & Alfalqi, K. (2015). A survey of topic modeling in text mining. *Int. J. Adv. Comput. Sci. Appl.(IJACSA)*, 6(1):147-153.
- Alwidian, S. A. A., Bani-Salameh, H. A., Alslaity, A. A. N. (2015). Text data mining: a proposed framework and future perspectives. *International Journal of Business Information Systems*, 18(2): 127-140.
- Anderson, K. (2012). Editorial Rejection — Increasingly Important, Yet Often Overlooked Or Dismissed, in The Scholarly Kitchen. <https://scholarlykitchen.sspnet.org/2012/04/19/post-publication-peerreview-what-value-do-usage-based-metrics-offer>.

- Ananiadou, S., McNaught, J. (2006). *Text mining for biology and biomedicine*. London: Artech House.
- Baars, H., Kemper, H. G. (2008). Management support with structured and unstructured data—an integrated business intelligence framework. *Information Systems Management*, 25(2): 132-148.
- Blei DM. (2012). Probabilistic topic models. *Communications of the ACM*. 55(4): 77-84.
- Chapman, P., Clinton, J., Kerber, R., Khabaza, T., Reinartz, T., Shearer, C., & Wirth, R. (2000). CRISP-DM 1.0: Step-by-step data mining guide. SPSS inc, 9, 13.
- Chen, H., Huang, X., & Li, Z. (2022). A content analysis of Chinese news coverage on COVID-19 and tourism. *Current Issues in Tourism*, 25(2), 198-205.
- Choudhary AK, Oluikpe PI, Harding JA, Carrillo PM. (2009). The needs and benefits of Text Mining applications on Post-Project Reviews. *Computers in Industry*. 1;60(9):728-40.
- Clifton, C., Cooley, R., & Rennie, J. (2004). Topcat: Data mining for topic identification in a text corpus. *IEEE transactions on knowledge and data engineering*, 16(8): 949-964.
- Davies, K. (2012). Content analysis of research articles in information systems (LIS) journals. *Library and Information Research*, 36(112), 16-28.
- Delen, D. (2014). *Real-world data mining: applied business analytics and decision making*. New Jersey: Financial Times Press
- Dong, G., Liu, H. (Eds.). (2018). *Feature engineering for machine learning and data analytics*. CRC Press.
- Elo S, Kyngäs H. (2008). The qualitative content analysis process. *Journal of advanced nursing*. 62(1):107-15.
- Glenisson, P., Glänzel, W., Janssens, F., & De Moor, B. (2005). Combining full text and bibliometric information in mapping scientific disciplines. *Information processing & management*, 41(6): 1548-1572.

- Golub, K. (2006). Using controlled vocabularies in automated subject classification of textual web pages, in the context of browsing. *IEEE TCDL Bulletin*, 2(2): 1-11.
- Hofmann T. (1999). Probabilistic latent semantic indexing. *Proceedings of the 22nd annual international ACM SIGIR conference on Research and development in information retrieval*; Berkeley, California, USA: Association for Computing Machinery; 50–57.
- Huisman, J., and J. Smits. (2017). Duration and quality of the peer review process: the author’s perspective. *Scientometrics*, 113 (1): 633-650.
- Ilias, K. (2019). A distributed text clustering approach for detecting duplicate questions (Doctoral dissertation, Aristotle University of Thessaloniki Thessaloniki, Greece).
- Indurkha, N., Damerau, F. J. (Eds.). (2010). Handbook of natural language processing (Vol. 2). CRC Press.
- Jain, N., & Srivastava, V. (2013). Data mining techniques: a survey paper. *IJRET: International Journal of Research in Engineering and Technology*, 2(11): 2319-1163.
- Janssens, F., Leta, J., Glänzel, W., & De Moor, B. (2006). Towards mapping library and information science. *Information processing & management*, 42(6): 1614-1642.
- Kantardzic, M. (2011). *Data mining: concepts, models, methods, and algorithms*. John Wiley & Sons.
- Khan, R. A. (2016). Data Mining: A Tool for Customer Relationship Management. *Data Mining and Knowledge Engineering*, 8(4), 95-99.
- Kim, Y., M., & Delen, D. (2018). Medical informatics research trend analysis: A text mining approach. *Health informatics journal*, 24(4): 432-452.
- Lamba, M., & Madhusudhan, M. (2019). Mapping of topics in DESIDOC Journal of Library and Information Technology, India: a study. *Scientometrics*, 120(2): 477-505
- Lin, Z., S. Hou, and J. Wu. (2016). The correlation between editorial delay and the ratio of highly cited papers in Nature, Science and Physical Review Letters.

Scientometrics 107 (3): 1457-1464.

- Majhi, S., Jal, Ch. & Maharana, B. (2016). Content analysis of Journal articles on Wiki in Science Direct Database. *Library Philosophy and Practice (e-journal)*, 1 – 15.
- Marbán, O., Segovia, J., Menasalvas, E., & Fernández-Baizán, C. (2009). Toward data mining engineering: A software engineering approach. *Information systems*, 34(1): 87-107.
- Marinakos, G., Daskalaki, S. (2016). Viability prediction for retail business units using data mining techniques: a practical application in the Greek pharmaceutical sector. *International Journal of Computational Economics and Econometrics*, 6(1): 1-12.
- Mesropyan, V. R., Ovsyannikov, M. V. (2014). Prospects for the application of scientometric methods for forecasting. *Scientific and Technical Information Processing*, 41(1): 38-46.
- Mitra, S., Acharya, T. (2003). *Data Mining: Concepts and Algorithms From Multimedia to Bioinformatics*. John Wiley & Sons.
- Mulligan, A., Hall, L. & Raphael, E. (2013). Peer review in a changing world: An international study measuring the attitudes of researchers. *Journal of the American Society for Information Science and Technology*, 64(1), 132-161.
- Nguyen, T. T., F. Maxwell Harper, L. Terveen, et al. (2018). User Personality and User Satisfaction with Recommender Systems. *Information Systems Frontiers* 20 (6): 1173-1189.
- Nicholson, S., Hwang, S. Y., Keezer, P., & O'Neill, E. T. (2003). The bibliomining process: Data warehousing and data mining for libraries. Sponsored by SIG LT. *Proceedings of the American Society for Information Science and Technology*, 40(1): 478-479.
- Pranata, I., Skinner, G. (2015). Segmenting and targeting customers through clusters selection & analysis. In *2015 International Conference on Advanced Computer Science and Information Systems (ICACISIS)*, 303-308.

- Ramos, J. (2003). Using tf-idf to determine word relevance in document queries. In Proceedings of the first instructional conference on machine learning. 242 (1): 29-48.
- Saheb T, Saheb M. (2019). Analyzing and visualizing knowledge structures of health informatics from 1974 to 2018: A bibliometric and social network analysis. *Healthc Inform Res.* 25(2): 61-72.
- Salloum SA, Al-Emran M, Monem AA, Shaalan K. (2017). A survey of text mining in social media: facebook and twitter perspectives. *Adv. Sci. Technol. Eng. Syst. J.* 2(1):127-33.
- Sharma, N. (2005). *Discovering knowledge with text mining*. A&M University-Kingsville.
- Silva, R. O., de Carvalho, R. L. (2021). Using Latent Semantic Indexing as a metric for evaluating research potentialities through Innovation Public Policies. *Academic Journal on Computing, Engineering and Applied Mathematics*, 2(2): 10-15.
- Spinakis, A., Peristera, P. (2004). Text mining tools: Evaluation methods and criteria. In *Text Mining and its Applications* (pp. 131-149). Springer, Berlin, Heidelberg.
- Srivastava AN, Sahami M. (2009). *Text Mining: Classification, Clustering, and Applications* (Chapman & Hall/CRC Data Mining and Knowledge Discovery Series) 1st Ed. Uk: Chapman and Hall/CRC.
- Smith SD, Ng K-M, Brinson J, Mityagin E. (2008). Multiculturalism, diversity, and social advocacy: A 17-year content analysis of counselor education and supervision. *Counselor Education and Supervision*.47(4):249.
- Soleimani Nezhad A, Salajegheh M, Tayyebi Nia E. (2019). Clustering scientific articles based on the k\_means algorithm Case Study: Iranian Research Institute for information Science and Technology . *Iranian Journal of Information Processing and Management*. 34(2):871-96.
- Soler, V. (2007). Writing titles in science: An exploratory study. *English for Specific Purposes*, 26(1), 90–102.

- Stroud, D., Pennington, P., Cleaver, C., Collins, J. R., & Terry, N. (2017). A content analysis of research articles in *The Journal for Specialists in Group Work*: 1998–2015. *The Journal for Specialists in Group Work*, 42(2), 194-210.
- Wang C, Blei D, Heckerman D. (2008). Continuous time dynamic topic models. *Proceedings of the Twenty- Fourth Conference on Uncertainty in Artificial Intelligence*; Helsinki, Finland: AUAI Press; 579–86.
- Wang, W. T., and Y. P. Hou. (2015). Motivations of employees’ knowledge sharing behaviors: A selfdetermination perspective. *Information and Organization* 25 (1): 1-26.
- Wei X, Croft WB. (2006). LDA-based document models for ad-hoc retrieval. *Proceedings of the 29<sup>th</sup> annual international ACM SIGIR conference on Research and development in information retrieval*; Seattle, Washington, USA: Association for Computing Machinery; 178–85.
- Witten, I. H., Frank, E. (2002). Data mining: practical machine learning tools and techniques with Java implementations. *Acm Sigmod Record*, 31(1): 76-77.
- Woo H, Heo N. (2013). A content analysis of qualitative research in select ACA journals (2005–2010). *Counseling Outcome Research and Evaluation*. 4(1):13-25.
- Yukselturk, E., Ozekes, S., & Turel, Y. K. (2014). Predicting dropout student: An application of data mining methods in an online education program. *European Journal of Open, Distance and e-learning*, 17(1): 118-133.
- Yau C-K, Porter A, Newman N, Suominen A. (2014). Clustering scientific documents with topic modeling. *Scientometrics*. 100(3): 767-86.
- Yao, W., J. He, H. Wang, Y. Zhang, & J. Cao. (2015). Collaborative topic ranking: Leveraging item metadata for sparsity reduction. In *Proceedings of the AAAI Conference on Artificial Intelligence* 29 (1): 374-380.
- Zeigler EF. (1987). Sport Management: Past, present, future. *Journal of sport management*. 1(1):4-24.