



کاربرد تحلیل گفتمان در خلاصه سازی خودکار متون علمی^۱

پگاه تاجر، عبدالرسول جوکار، سید مصطفی فخر احمد، علیرضا خرمایی، هاجر ستوده

چکیده

هدف این پژوهش آشکارسازی کاربرد ویژه تحلیل گفتمان در خلاصه‌سازی خودکار متون علمی است و به این منظور از رویکرد مرور و تحلیل نوشتارها استفاده شده است. متن یک واحد زبانشناسی پیچیده است و ساختار گفتمانی و تئوری‌های سازماندهی متن را می‌توان شاخصی برای تفسیر متن و انتخاب جمله مناسب به منظور تولید خودکار خلاصه‌های خواناتر و منسجم‌تر در نظر گرفت. نشانه گذاری بلاغی اطلاعات، برای فرایندهای جدید دسترسی به اطلاعات علمی مفید می‌باشد. استراتژی‌های ویژه خلاصه‌سازی خودکار مقالات علمی عموماً تلفیقی از روش‌های سطحی و معنایی هستند. یکی از استراتژی‌هایی که از تجزیه و تحلیل کلام بهره می‌برد عبارت است از ناحیه بندی استدلالی. بر اساس این استراتژی طرح‌های حاشیه نویسی گفتمانی ویژه متون علمی به وجود آمده اند که بر اساس آنها می‌توان ابعاد چهارگانه حل مسئله، تخصیص فکری، استدلال علمی و نگرش به آثار دیگران را در متون علمی معلوم کرد و بر غنای خلاصه کلیدواژه‌ها: تحلیل گفتمان، ژانر علمی، خلاصه سازی خودکار متن، خلاصه سازی گفتمان مدار، خلاصه سازی استنادی، ناحیه بندی استدلالی

۱. برگرفته از رساله دکتری

های خودکار افزود. مشهورترین طرح حاشیه نویسی نواحی بلاغی مقالات علمی، طرح توپفل می باشد. از کاربردهای مهم و ویژه تحلیل گفتمان در خلاصه سازی خودکار ژانر علمی، تولید خلاصه های استنادی گفتمان مدار می باشد که بر اساس آن می توان انواع خلاصه های خودکار مبتنی بر پرسش کاربر و خلاصه های چندسندی ای که تضادها و شباهت های مدارک در خلاصه نهایی مشهود است را تولید نمود.

۱- مقدمه

رشد فزاینده اطلاعات متنی در محیط وب و نیاز به دسترسی، ارزیابی و انتخاب سریع آنها، منجر به پدیدآمدن ابزارهایی برای خلاصه سازی خودکار متن شده است. این ابزارها، در ژانرهای مختلف از جمله خبری، حقوقی و علمی و نیز در خلاصه سازی صفحات وب به منظور نمایش در تلفن همراه و سیستم های بازیابی اطلاعات، کاربردهای فراوانی دارند.

خلاصه سازی خودکار متون از تکنیک های پردازش زبان طبیعی است که همواره مورد توجه زبان شناسان رایانشی و متخصصان بازیابی اطلاعات بوده است. خلاصه سازی خودکار متن عبارت است از نمایش فشرده، دقیق و منسجم متن ورودی به طوری که متن خروجی مفاهیم مهم متن ورودی را در بر داشته باشد (میبوری^۱، ۱۹۹۵ نقل در رحمان و بورا^۲، ۲۰۱۵). پژوهش در حوزه خلاصه سازی خودکار، در اواخر دهه پنجاه میلادی، زمانی که تمایل زیادی به خودکار سازی چکیده نویسی اسناد فنی وجود داشت، مورد توجه جامعه علمی قرار گرفت (لوهن^۳، ۱۹۵۸) و با به کار گیری هوش مصنوعی به طور چشمگیری افزایش یافت (دیجانگ^۴، ۱۹۸۲). از دهه نود، با ظهور کنفرانسهای علمی متعدد مرتبط با خلاصه سازی خودکار، تمایل به پژوهش در میان پژوهشگران این حوزه رنگ و بوی تازه ای به خود گرفت که نقطه اوج آن را می توان توسعه چارچوب های ارزیابی از سال ۲۰۰۰ در آمریکا دانست. دو برنامه ارزیابی مهم عبارتند از کنفرانس های درک سند (دی. یو. سی.)^۵ و کنفرانس های تحلیل متن (تی. ای. سی.)^۶.

چگونگی انتخاب محتوای ضروری یک سند و نحوه بیان محتوای انتخاب شده به صورتی فشرده، دو عنصر اساسی فرآیند خلاصه سازی متن است (جونز^۷، ۱۹۹۳؛ جونز و اندرس- نیگمیر^۸، ۱۹۹۵). به طور کلی، خلاصه های خودکار به دو روش استخراجی^۹ و چکیده سازی^{۱۰} تولید می شود. خلاصه های استخراجی

^۱ Maybury

^۲ Borah

^۳ Luhn

^۴ DeJong

^۵ Document Understanding Conference (DUC)

^۶ Text Analysis Conference (TAC)

^۷ Jones

^۸ Endres-Niggemeyer

^۹ Extractive

^{۱۰} Abstractive

شامل مجموعه‌ای از جملات سند ورودی هستند. این نوع خلاصه‌سازی یک روش ساده اما قوی برای خلاصه‌سازی متن است. لذا بسیاری از پژوهش‌های این حوزه، از نوع استخراجی هستند. استخراج و انتخاب جملات مناسب برای خلاصه نهایی، بر اساس ویژگی‌های متون صورت می‌پذیرند. این ویژگی‌ها، طیف وسیعی از ویژگی‌های سطحی مانند فراوانی کلمات و محل قرار گرفتن جمله در متن تا ویژگی‌های معنایی و گفتمانی را شامل می‌شوند (سیگون^۱ و پویبئو^۲، ۲۰۱۳). خلاصه‌های چکیده‌ای، تفسیری از متن اصلی را ارائه می‌کنند و در تولید آنها مفاهیم جملات اصلی به شکل کوتاه‌تر بازنویسی می‌شوند (مانی^۳، ۲۰۰۱). خلاصه‌سازی از منظر تعداد اسناد، می‌تواند تک‌سندی^۴ یا چندسندی^۵ باشد. در خلاصه‌سازهای تک‌سندی، ورودی تنها یک سند است. این در حالی است که خلاصه‌سازی چندسندی بازنمون مختصر و منسجمی از محتوای کلیدی مجموعه‌ای از مدارک مربوط به یکدیگر است که به صورت خودکار تولید شده باشد. (سیگون و پویبئو، ۲۰۱۳).

از آنجایی که متن یک واحد زبان‌شناسی پیچیده است، ساختار گفتمانی و یا تئوری‌های سازماندهی متن را می‌توان شاخصی برای تفسیر متن و انتخاب جمله مناسب به منظور تولید خودکار خلاصه‌های خواناتر و منسجم‌تر در نظر گرفت.

۲- بیان مسأله

بهترین منابع باید در مناسب‌ترین قالب و در کمترین زمان در دسترس محققان قرار گیرد. با توجه به رشد روزافزون انتشارت علمی، خلاصه‌سازی خودکار متون، راهکار مناسبی برای تسهیل و تسریع فرآیند درک محتوای اسناد علمی است.

امروزه، خلاصه‌سازی خودکار چندسندی متون علمی اجتناب‌ناپذیر است. پایگاه‌های اطلاعات علمی علاوه بر متن اصلی، چکیده‌های نویسنده را در اختیار کاربر قرار می‌دهند اما در مواجهه با انبوهی از مقالات، یک سیستم خلاصه‌ساز چندسندی می‌تواند فشرده‌ای از اهداف، روش‌ها و نتایج مجموعه‌ای از مدارک پژوهشی را در اختیار محقق قرار دهد و او را از صرف وقت زیاد برای بررسی انبوهی از چکیده مقالات و یا متن کامل آنها به ویژه در مراحل آغازین پژوهش مانند انتخاب موضوع و مرور پیشینه‌ها تا حد زیادی بی‌نیاز کند.

همانطور که در مقدمه گفته شد، فرآیند خلاصه‌سازی چندسندی دشوارتر از نوع تک‌سندی است و با وجود ارائه راهکارهای مختلف در سال‌های گذشته، همچنان چالش‌هایی وجود دارد. بیشترین تلاش‌ها برای رفع معضلات خلاصه‌سازی چندسندی در ژانر خبری انجام شده‌اند که بعضاً به نتایج

^۱ Saggion

^۲ Poibeau

^۳ Mani

^۴ Single document

^۵ Multi-document

قابل قبولی رسیده‌اند. با این حال، با توجه به اینکه متون خبری با متون علمی ماهیتاً تفاوت دارند، نتایج به دست آمده در ژانر خبری را نمی‌توان به ژانر علمی تعمیم داد. با توجه به این مهم، مسأله این پژوهش آشکارسازی کاربرد ویژه تحلیل گفتمان در خلاصه‌سازی خودکار متون علمی است و به این منظور از رویکرد مرور و تحلیل نوشتارها استفاده شده است.

۳- اهمیت پژوهش

تولید خلاصه‌های خودکار باکیفیت مطلوب همواره از چالش‌های خلاصه‌سازی خودکار متن بوده‌است. به همین دلیل، خلاصه‌سازی خودکار متن را از جمله مباحث پیچیده پردازش زبان طبیعی می‌دانند که علی‌رغم پیشرفت‌ها، همواره مسائل زیادی را پیش روی جامعه علمی قرار می‌دهد (سیگون و پویبیو، ۲۰۱۳). این در حالی است که اکثر پژوهش‌های خلاصه‌سازی خودکار بر مقالات خبری تمرکز دارند و به دیگر ژانرها از جمله متون علمی توجه کمتری شده‌است (تویفل^۱ و موئنز^۲، ۲۰۰۲).

از طرف دیگر، اطمینان از پیوستگی متن خلاصه‌های تولید شده به صورت خودکار کار دشواری است زیرا نه تنها به درک محتوای هر قطعه متنی بلکه به دانش ساختار گفتمانی متن هم نیازمند است (رادو^۳، هاوی^۴ و مک کئون^۵، ۲۰۰۲). بنابراین آشکارسازی کاربرد های ویژه تحلیل گفتمان در خلاصه‌سازی ژانر علمی می‌تواند منجر به اکتشاف و ارائه الگوریتم های نوین خلاصه‌سازی خودکار متون علمی گردد.

۴- پیشینه پژوهش

خلاصه‌سازی خودکار متن از حوزه هایی است که همواره مورد توجه پژوهشگران بوده و در ۲۵ سال گذشته نیز انواع رویکردهای خلاصه‌سازی ارائه شده‌است. تعدادی از این رویکرد ها بر تخمین توزیع موضوع/محتوای خلاصه نهایی از طریق تکنیک هایی چون تحلیل معنایی پنهان^۶، مدل سازی موضوعی^۷ و مدل های بیزی^۸ تمرکز کرده‌اند. برخی دیگر، خلاصه‌سازی را یک مسئله بهینه سازی دانسته و عده ای دیگر با اتخاذ رویکرد طبقه‌بندی نظارتی سعی کرده‌اند راه حل هایی

^۱ Teufel

^۲ Moens

^۳ Radev

^۴ Hovy

^۵ McKeown

^۶ Latent Semantic Analysis (LSA)

^۷ Topic Modelling

^۸ Bayesian models

برای خلاصه‌سازی خودکار ارائه دهند. مدل های گراف‌مدار، رویکرد دیگری است که در ادبیات خلاصه‌سازی متن، بسیار مورد توجه قرار گرفته است. در این گونه مدل‌ها هدف اصلی، یافتن مرکزی ترین جمله در گراف شباهت جملات یک سند می‌باشد. خلاصه‌سازهایی چون «لکس رنک»^۱ و «تکس رنک»^۲ از این رویکرد بهره می‌برند. علاوه بر این‌ها، حداکثرسازی تازگی و پرهیز از افزونگی در یک خلاصه سیستمی با بهره گیری از الگوریتم انتخاب حریصانه^۳ راهکار دیگری است که بررسی شده است. از جمله ایده‌هایی که اخیراً مورد توجه بوده است، بررسی ساختار بلاغی مدارک بر اساس تئوری «آر.اس. تی.» برای تجزیه گفتمانی مدارک می‌باشد (کهن و گوهریان، ۲۰۱۵).

در ادبیات پردازش زبان طبیعی، مطالعات تویفل و مؤنزر در سال ۲۰۰۲ را می‌توان اولین پژوهش خلاصه‌سازی خودکار ژانر علمی دانست. اما ریشه‌های خلاصه‌سازی خودکار متون علمی را می‌توان در ادبیات علم اطلاعات و دانش‌شناسی و مباحث چکیده‌نویسی خودکار سال‌های دهه پنجاه قرن بیستم جست.

تمایل به خلاصه‌سازی خودکار متون علمی از زمانی آغاز شد که کتاب‌ها و مقالات علمی به صورت دیجیتالی ذخیره شدند. نمونه‌های اولیه چنین تمایلاتی را در دهه پنجاه میلادی و در پژوهش‌هایی که به چکیده‌نویسی و نمایه‌سازی خودکار مقالات و کتاب‌های علمی می‌پرداختند، می‌توان دید (لوهن، ۱۹۵۸؛ بکسندال^۴، ۱۹۵۸). با این وجود، جامعه پژوهشی خلاصه‌سازی خودکار، از خلاصه‌سازی مقالات علمی به خلاصه‌سازی مقالات خبری روی آورد و این روند تا دوباره فعال شدن خلاصه‌سازی خودکار ژانر علمی در اواخر دهه نود میلادی، همچنان ادامه داشت (یلوگلو^۵ و همکاران، ۲۰۱۱). رویکرد خلاصه‌سازی پیکره‌بنیاد، در سال ۱۹۹۵ توسط کوپیک^۶ و همکارانش در هشتمین کنفرانس بین‌المللی سالانه «ای.سی.ام. سیگیر»^۷ پیشنهاد شد (کوپیک و همکاران، ۱۹۹۵). لازم به ذکر است که این کنفرانس به موضوع پژوهش و توسعه در بازیابی اطلاعات می‌پرداخت.

مرور ادبیات خلاصه‌سازی نشان می‌دهد که با وجود فعال شدن این حوزه از دهه نود، تعداد پژوهش‌های آن در مقایسه با پژوهش‌های خلاصه‌سازی خودکار متون خبری اندک است. اغلب پژوهش‌های خلاصه‌سازی خودکار مقالات علمی، به خلاصه‌سازی تک‌سندی مقالات پرداخته‌اند

^۱ LexRank

^۲ TextRank

^۳ Greedy selection

این الگوریتم، بهترین انتخاب را با توجه به شرایط مسئله انجام می‌دهد به امید آنکه با ادامه همین روش بهینه‌سازی انجام شود.

^۴ Baxendale

^۵ Yeloglu

^۶ Kupiec

^۷ ACM SIGIR

و از روش‌های یادگیری ماشینی بهره برده‌اند. پژوهشگرانی چون کوپیک و همکاران (۱۹۹۵)، توپفل و موئنز (۲۰۰۲) یادگیری ماشینی نظارتی و برخی دیگر مانند قزوینیان و رادو (۲۰۰۸) و (۲۰۱۰)، الکیس^۱ و همکاران (۲۰۰۸) و ابوجبارا^۲ و رادو (۲۰۱۱) رویکرد غیرنظارتی را به کار گرفته‌اند. در این میان پژوهشگرانی هم هستند که خلاصه‌سازهای چندمنظوره برای ژانرهای مختلف و با اهداف مختلف، طراحی کرده‌اند که بعضاً ژانر علمی را هم شامل می‌شود. برای مثال، لورت^۳ و پالومار^۴ (۲۰۱۲) در پژوهشی با عنوان «کامپنیدیوم^۵: یک ابزار خلاصه‌سازی متن برای تولید خلاصه‌های چندمنظوره از ژانرها و حوزه‌های گوناگون»، یک سامانه خلاصه‌ساز طراحی کردند که قادر است متداول‌ترین انواع خلاصه‌ها از جمله تک‌چندی، چندسندی، استخراجی، چکیده‌ای، پرسش‌مدار و مبتنی بر احساس^۶ را به صورت خودکار تولید کند.

۵- چارچوب نظری پژوهش

در این بخش ضمن تعریف اجمالی گفتمان و تحلیل آن، مباحث نظری خلاصه‌سازی گفتمان مدار که اساس نظری این مطالعه را تشکیل می‌دهد، تشریح می‌گردد.

۵-۱- تعریف گفتمان

چادرون^۷ و ریچاردز^۸ (۱۹۸۵)، دو تعریف برای گفتمان ارائه کرده‌اند که عبارتند از: الف. گفتمان اصطلاحی است عام برای نمونه‌های کاربرد زبان؛ یعنی زبان که برای برقراری ارتباط تولید شده است، ب. برخلاف دستور زبان که با عبارت‌ها و جمله‌ها سر و کار دارد، گفتمان به واحدهای زبانی بزرگ‌تر چون بند، مصاحبه، مکالمه و متن نظر دارد (بشیری، ۱۳۹۴). براون^۹ و یول^{۱۰} (۱۹۸۳) معتقدند، گفتمان عبارت است از تجزیه و تحلیل زبان در کاربرد. بنابراین گفتمان نمی‌تواند به توصیف صورت‌های زبانی، جدا از اهداف و نقش‌هایی که این صورت‌ها برای پرداختن به آنها در امور انسانی به وجود آمده‌اند، بپردازد.

^۱ Elkiss

^۲ Abu-Jbara

^۳ Lloret

^۴ Palomar

^۵ COMPENDIUM

^۶ Sentiment-based

^۷ Chaudron

^۸ Richards

^۹ Brown

^{۱۰} Yule

تحلیل گفتمان برچسبی است برای آن دسته از تحلیل های متون که منبعث از نظریه ها و روال های زبان شناسی هستند. تحلیل گفتمان در وهله نخست به شرایط خوش ساختگی (پیوستگی و انسجام) و قواعد قیاسی می پردازد (البرزی ورکی، ۱۳۸۶).

تحلیل گفتمان، رویکردی میان رشته ای است که ریشه در زبان شناسی دارد. هدف عمده تحلیل گفتمان این است که تکنیک و روش جدیدی را در مطالعه متون، رسانه ها، فرهنگ ها، علوم، سیاست، اجتماع و مواردی مانند آن به دست دهد. تحلیل گفتمان جزء روش های تحقیق کیفی بوده که جهت کشف معنای به کار رفته در متن یا سخن به کار می رود (کلانتری و دیگران، ۱۳۸۸).

۵-۲- نشانگرهای گفتمانی

اولین تلاش ها برای معرفی نشانگرهای گفتمانی به عنوان یک موجودیت زبانی توسط لابو^۱ و فانشل^۲ در سال ۱۹۷۷ م. صورت گرفت. این دو، نشانگر گفتمانی را موجودیتی می دانستند که به موضوع^۳ دانش توزیع شده بین افراد شرکت کننده در بحث ارجاع می دهد (فراسر^۴، ۱۹۹۹). نشانگرهای گفتمانی وابسته به بافت هستند و بر اساس نگرش ها، شرکت کنندگان و متن طبقه بندی می شوند. بنابراین دارای کارکردهای بلاغی هم در سطح متن و هم در سطح بین فردی می باشند. این نشانگرها، نقش مهمی در درک گفتمان و پیشرفت اطلاعات دارند (شیفرین^۵، ۱۹۸۷ نقل در یانگ^۶، ۲۰۱۱). کلمات و عباراتی چون 'all in', 'well', 'so', 'yeah', 'right', 'like', 'But' okay و 'after all', 'you know' نمونه هایی از نشانگرهای گفتمانی در زبان انگلیسی هستند (یانگ، ۲۰۱۱). شیفرین (۱۹۸۷) اهمیت نشانگرهای گفتمانی را آشکار ساخت و مدل انسجامی خود را ارائه داد. مدل او شامل سه سطح معنا، نحو^۷ و سازماندهی گفتمان بود. وی با ارائه این مدل درصد بود تا بررسی کند که نشانگرهای گفتمانی چگونه می توانند به انسجام شفاهی کمک کنند (آرچاکیس^۸، ۲۰۰۱ نقل در یانگ، ۲۰۱۱). او ۱۱ نوع از نشانگرهای گفتمانی را در مدل خود توصیف نمود (شیفرین، ۱۹۸۷). در حالی که بسیاری از زبان شناسان به تحلیل توصیفی نشانگرهای گفتمانی در زبانی خاص و کلام بومی پرداختند، مولر^۹ (۲۰۰۴) بر نشانگرهای گفتمانی زبان مقصد که غیربومیان در کلام به کار می گیرند، تمرکز نمود (یانگ، ۲۰۱۱).

^۱ Labov

^۲ Fanshel

^۳ Topic

^۴ Fraser

^۵ Schiffrin

^۶ Yang

^۷ Syntactic

^۸ Archakis

^۹ Muller

در ادبیات تجزیه و تحلیل کلام، عبارت های متنوعی وجود دارند که در برابر عبارت «نشانگرهای گفتمانی» به کار رفته اند. برخی از آنها عبارتند از: عبارت های نشانه، ذرات گفتمانی، عبارات پراگماتیک، نشانه های پراگماتیک، عملگرهای پراگماتیک، اتصال دهنده های جمله و عملگرهای گفتمانی (فراسر، ۱۹۹۹).

۵-۳- خلاصه سازی خودکار گفتمان مدار

انسجام و خوانایی خلاصه های سیستمی از چالشهای خلاصه سازی خودکار به ویژه نوع چندسندی آن است که برای مقابله با آن راهکار تعیین ناحیه های بلاغی متن^۱ ارائه شده است. این روش، با تعیین ساختار گفتمانی متن تلاش می کند جملات را به نحوی طبقه بندی کند که خلاصه هایی منسجم و خوانا تولید شود. بر این اساس، نوعی از خلاصه سازی به نام خلاصه سازی گفتمان مدار پیشنهاد شده است. ساختار گفتمانی یک متن شامل تمام روابط بین قطعه های یک متن است که ساختار آن متن را می سازد. ایده به کار گیری ساختار گفتمانی متن برای خلاصه سازی متون در سال ۱۹۹۷ توسط مارکو^۲ برای تولید خلاصه های کلی پیشنهاد شد (باسما^۴، ۲۰۰۸) و حوزه ای از خلاصه سازی خودکار به نام خلاصه سازی گفتمان مدار را پدید آورد. این نوع خلاصه سازی از نظریه های سازماندهی گفتمان بهره می برد. یکی از مشهورترین این نظریه ها نظریه ساختار بلاغی^۵ است. بر اساس این نظریه می توان متن را به صورت ساختاری درختی یا درخت بلاغی بازنمون نمود که در آن محدوده های متنی^۶ با به کارگیری مجموعه ای از روابط گفتمانی از پیش تعیین شده به هم پیوند داده می شوند (مان^۷ و تامپسون^۸، ۱۹۸۸).

بر اساس نظریه ساختار بلاغی، استراتژی ناحیه بندی استدلالی متون علمی به منظور خلاصه سازی خودکار متون علمی شکل گرفت. در این استراتژی، فرض بر این است که گفتمان علمی در بردارنده توصیف هایی از نظرات مثبت و منفی راجع به سهم پژوهشگران از یک حوزه علمی است. این نوع گفتمان، حاصل یک بازی بلاغی است که سعی دارد جایگاه و سهم پژوهشگران از یک حوزه علمی را ارتقاء دهد. بر این اساس، متون علمی باید ادعاهای دانش را در تقابل با آثار پژوهشی قبلی به

^۱ Rhetorical Zones

^۲ Passage

^۳ Marcu

^۴ Bosma

^۵ Rhetorical Structure Theory (RST)

^۶ Text spans

^۷ Mann

^۸ Thompson

وضوح نشان دهند (فیساس^۱، رونزانو^۲ و سیگون، ۲۰۱۵). از این منظر، مدل گفتمانی ادعاهای دانش^۳ معروف به «مدل سیمون تویفل^۴» یا همان ناحیه بندی استدلالی به وجود آمده است (تویفل، ۱۹۹۹؛ تویفل و موئنز، ۲۰۰۲؛ تویفل، سیدهارتان^۵ و بچلور^۶، ۲۰۰۹) در ناحیه بندی استدلالی، وضعیت بلاغی هر جمله به منظور بازنمون بافت گفتمانی جملات استخراج شده، تعیین می شود. به عبارت دیگر در این استراتژی، فرایند استخراج جمله با تحلیل گفتمان ترکیب شده است (تویفل و موئنز، ۲۰۰۲).

۶- کاربرد تحلیل گفتمان در خلاصه‌سازی ژانر علمی

در راستای مباحثی که در بخش چارچوب نظری مطرح شد، استراتژی‌های ویژه خلاصه‌سازی خودکار مقالات علمی توسعه یافته‌اند که عموماً تلفیقی از روش‌های سطحی و معنایی هستند. یکی از استراتژی‌هایی که از تجزیه و تحلیل کلام بهره می‌برد عبارت است از ناحیه بندی استدلالی^۷ بر اساس این استراتژی طرح‌های حاشیه نویسی گفتمانی ویژه متون علمی به وجود آمده‌اند که در ادامه مورد بحث قرار می‌گیرند.

۶-۱- طرح‌های مبتنی بر ناحیه بندی استدلالی

نشانه‌گذاری بلاغی اطلاعات، برای فرایندهای جدید دسترسی به اطلاعات علمی مفید می‌باشد. مثلاً در حوزه بازبانی اطلاعات علمی، می‌توان از اطلاعات بلاغی به عنوان گزاره‌های بیان‌گر تغییر پارادایم استفاده کرد (لیساکیک^۸ و همکاران، ۲۰۰۵ نقل در تویفل، سیدهارتان و بچلور، ۲۰۰۹). زیرا مقالاتی که شامل این نوع گزاره‌ها هستند، مقالات تأثیرگذار یک حوزه محسوب می‌شوند. برای مثال، ۷۵ درصد از مقالات «اف»^۹ که یک سرویس توصیه‌گر مقالات حوزه‌های پزشکی و زیست‌شناسی است، حاوی جملات تغییر پارادایم است. این مقالات از منظر خبرگان،

^۱ Fisas

^۲ Ronzano

^۳ Knowledge Claim Discourse Model (KCDM)

^۴ Simone Teufel's model

^۵ Siddharthan

^۶ Batchelor

^۷ Argumentative Zoning (AZ)

^۸ Lisacek

^۹ Biology ۱۰۰۰ Faculty of

از اهمیت خاصی در حوزه برخوردارند (تویفل، سیدهارتان و بچلور، ۲۰۰۹). بنابراین ساختار بلاغی مقالات علمی در خلاصه‌سازی مقالات علمی نیز حائز اهمیت می‌باشد. زیرا بر اساس آن، می‌توان خلاصه‌های تک‌سندی و چندسندی تولید کرد که به تفاوت‌ها و شباهت‌های آثار مورد استناد قرار گرفته هم توجه می‌کنند. برای تعیین وضعیت بلاغی جملات منتخب، ابتدا لازم است ساختار یک مقاله علمی تشریح شود. ساختار یک مقاله علمی ابعاد مختلفی دارد که عبارتند از:

حل مسئله: پژوهش عموماً به صورت فعالیتی برای حل مسئله توصیف می‌شود (جردن^۱، ۱۹۸۴؛ زاپن^۲، ۱۹۸۳ نقل در تویفل و موئنز، ۲۰۰۲). انتظار می‌رود در هر مقاله پژوهشی سه نوع اطلاعات وجود داشته باشد که عبارتند از: مسئله و اهداف پژوهش، راه حل‌ها (روش‌ها) و نتایج. این ساختار که ساختار مسئله راه حل نام دارد در بسیاری از رشته‌ها به ویژه در علوم تجربی ساختاری غالب است که به صورت ساختار ثابت مقدمه/روش/نتایج/بحث مشهود است (ون دیجک^۳، ۱۹۸۰ نقل در تویفل و موئنز، ۲۰۰۲). این در حالی است که متونی علمی هم وجود دارند که عیناً از این ساختار پیروی نمی‌کنند.

تخصیص فکری^۴: متون علمی باید سهم خود از دنیای علم و اصالت اثر^۵ را به صورت واضح نشان دهند و جنبه‌های جدید کار را در مقابله با رویکردهای پژوهشگران دیگر و بیانیه‌های علمی عموماً پذیرفته شده، آشکار سازند. نویسندگان عموماً دلیل محکمی برای نوشتن مقاله خود دارند و آن را به گونه‌ای می‌نگارند که تفکر خاص شان آشکارا بیان شود.

جدول ۱: طرح حاشیه نویسی نواحی بلاغی تویفل و موئنز (۲۰۰۲)

توصیف	ناحیه بلاغی
ناحیه هدف خاص پژوهشی مقاله	AIM
ناحیه گزاره‌هایی درباره ساختار بخش	TEXTUAL
ناحیه توصیف اثر خود: روش‌شناسی، نتایج و بحث (ناحیه خنثی)	OWN
ناحیه پیشینه علمی عموماً پذیرفته شده	BACKGROUND
ناحیه گزاره‌هایی در مقایسه یا تضاد با اثر دیگر؛ ضعف اثر دیگر	CONTRAST
ناحیه گزاره‌هایی در موافقت با اثر دیگر یا ادامه اثر دیگر	BASIS
ناحیه توصیف اثر دیگران (ناحیه خنثی)	OTHER

^۱ Jordan

^۲ Zappen

^۳ van Dijk

^۴ Intellectual attribution

^۵ Contribution

خوانندگان نیز معمولاً آن را به آسانی درک می‌کنند (تویفل و موئنز، ۲۰۰۲). استدلال علمی^۱: علم تنها کارخانه‌ای از واقعیت‌ها نیست بلکه جنبه‌های اجتماعی قوی‌ای دارد. زیرا موفقیت یک پژوهشگر، به توانایی او برای متقاعد کردن دیگران در مورد اعتبار استدلال‌هایش وابسته است (اسویلز^۲، ۱۹۹۰ نقل در تویفل و موئنز، ۲۰۰۲). به عبارت دیگر، نویسنده تلاش می‌کند اعتبار «دعای دانش»^۳ خود را نشان دهد. دانشی که در واقع قطعه دانش جدیدی است که از طریق داوری و انتشار با مخازن دانش آن حوزه یکپارچه خواهد شد (تویفل، سیدهارتان و بچلور، ۲۰۰۹). او این کار را از طریق یک قطعه متنی که از نظر بلاغی انسجام دارد، انجام می‌دهد (میرز^۴، ۱۹۹۲ نقل در تویفل و موئنز، ۲۰۰۲).

نگرش به آثار دیگران: نویسندگان برای مستدل کردن افکار خود به آثار دیگران ارجاع می‌دهند. هر استدلال دلیلی دارد. به عبارت دیگر، نویسنده هر ارجاع را بر اساس هدفی ذکر می‌کند. در یک متن خوب نگارش شده، این گونه اهداف قابل درک است (تویفل و موئنز، ۲۰۰۲). حوزه تحلیل محتوای استدلال^۵، به چگونگی و چرایی استنادی نویسنده‌ها می‌پردازد. در تحلیل محتوای استدلال با بررسی بافتار استنادی تلاش می‌شود تا چگونگی ارتباط بین مقالات آشکار شود. هدف از استناد انواع گوناگونی چون علمی (نقد مثبت و منفی و...) و اجتماعی (رعایت ادب، سنت و...) دارد (زیمان^۶، ۱۹۶۹ نقل در تویفل و موئنز، ۲۰۰۲).

ابعاد مختلف ساختار یک مقاله علمی، راهنمایی جهت تعیین وضعیت بلاغی جملاتی است که قرار است برای درج در خلاصه‌نهایی انتخاب شوند. جملات دارای وضعیت بلاغی یکسان ناحیه‌های بلاغی متن را تشکیل می‌دهند. کوچک‌ترین ناحیه فقط یک جمله را شامل می‌شود و ناحیه‌های بزرگ‌تر جملات بیشتری را شامل می‌شوند (تویفل و موئنز، ۲۰۰۲). یکی از مشهورترین طرح‌های حاشیه‌نویسی نواحی بلاغی مقالات علمی، طرحی است که تویفل و موئنز (۲۰۰۲) ارائه داده‌اند. این طرح، ۷ ناحیه بلاغی^۸ را برای مقالات علمی حوزه زبان‌شناسی رایانشی شناسایی کرده است (جدول ۱). سه ناحیه OWN، OTHER و BACKGROUND بر این اساس تعریف می‌شوند که در بخش مربوطه، چه کسی صاحب آن «دعای دانش» است. دو ناحیه BASIS و CONTRAST با توجه به ارتباطشان با آثار موجود تعریف می‌شوند. بنابراین بخشی از این طرح به طرح‌های

۱ Scientific argumentation

۲ Swales

۳ Knowledge claim

۴ Myers

۵. متن پیرامونی یک استناد را گویند. به عبارت دیگر جمله‌ای است که شامل استناد به اثر است.

۶ Content citation analysis

۷ Ziman

۸ Rhetorical Zones

جدول ۲: طرح حاشیه نویسی نواحی بلاغی توپفل، سیدهارتان و باتچلور (۲۰۰۹)

توصیف	ناحیه بلاغی	توصیف	ناحیه بلاغی
نتیجه گیری (غیر قابل اندازه گیری) نویسنده در مقاله حاضر	OWN-CONC	اهداف یا فرضیه های پژوهش در مقاله حاضر	AIM
مقایسه، تضادها و تفاوت های مقاله با دیگر راه حل ها (توصیف و مقایسه خنثی)	CODI	تازگی و مزیت رویکرد خاص پژوهشگر در مقاله حاضر	NOV-ADV
کمبود راه حل در حوزه، مشکل با راه حل های دیگر (ترکیب شکاف های دانش با رویکرد انتقادی)	GAP-WEAK	زمینه های مشترک، پیشینه علمی	CO-GRO
مقایسه و برخورد با نتایج یا نظریه دیگران، برتری اثر خود	ANTISUPP	ادعای دانش قابل توجهی که توسط افراد دیگر در مقاله هاشان بیان شده است (توصیف خنثی)	OTHER
آثار دیگر، مقاله حاضر را پشتیبانی می کنند یا اثر دیگری، این مقاله را حمایت می کند	SUPPORT	ادعای دانش قابل توجهی که توسط نویسنده در مقالات گذشته اش آورده است (توصیف خنثی)	PREV-OWN
کارهای دیگران که در اثر نویسنده مقاله حاضر به کار گرفته شده است	USE	ادعای دانش جدید نویسنده مقاله: روش شناسی	OWN-MTHD
محدودیت ها و پیشنهادات پژوهش، پژوهش های آینده	FUT	راه حل، روش و آزمایش نویسنده مقاله که کارساز نبوده است (خطاهای قابل بازیافتی نویسنده)	OWN-FAIL
		بروندهای عینی و قابل اندازه گیری نویسنده در مقاله حاضر (یافته ها)	OWN-RES

رده بندی عملکرد استناد^۱ مربوط است که بر گرفته از حوزه تحلیل محتوای استناد است. AIM به اصلی ترین ادعای دانشی مقاله که معمولاً در مقدمه و نتیجه گیری تکرار می شود اشاره می نماید و بالاخره TEXTUAL، محل فیزیکی اطلاعات را از طریق مرور کلی یک بخش و یا ارائه خلاصه ای از بخش های فرعی توضیح می دهد (توپفل، سیدهارتان و بچلور، ۲۰۰۹).

توپفل، سیدهارتان و بچلور (۲۰۰۹) این طرح را به ۱۵ ناحیه برای مقالات دو حوزه شیمی و زبان شناسی رایانشی گسترش دادند (جدول ۲). در این طرح گسترش یافته، ناحیه AIM مانند طرح قبل حفظ شده است. نام ناحیه BACKGROUND به CO-GRO COMMON GROUND تغییر یافته است. ناحیه OTHER به دو ناحیه OTHER و PREV-OWN، ناحیه BASIS به دو ناحیه USE و SUPPORT، ناحیه CONTRAST به سه ناحیه CODI، ANTISUPP و

^۱ Citation function classification schemes

OWN-MTHD، OWN-RES، OWN- چهار ناحیه OWN به بالاخره ناحیه GAP-WEAK و CONC و OWN-FAIL تقسیم شده است. با توجه به این که ناحیه TEXTUAL نسبت به دیگر نواحی اطلاعات کمتری در اختیار قرار می دهد، در این طرح حذف شده است. با توجه به اهمیت بیش از پیش تمایزات در حوزه های مختلف علم و با هدف نزدیک تر شدن به طرح های عملکرد استناد موجود، طرح تویفل، سیدهارتان و بچلور (۲۰۰۹) جزئی تر و مفصل تر از طرح تویفل و موئنز (۲۰۰۲) می باشد. لازم به ذکر است که ادعاهای دانشی که شامل استناد هستند و استنادهایی که در مثال های متعدد به وقوع می پیوندند، یکی از عوامل تصمیم گیری برای تعیین ناحیه بلاغی یک جمله است. اما، استنادها لزوماً عامل کلیدی برای چنین تصمیم گیری هایی نیستند. این در حالی است که تأثیرگذارترین عامل در این فرآیند، تمایز خاص نویسنده در آن حوزه از علم می باشد که از طریق نواحی چهارگانه OWN به خوبی قابل درک خواهد بود (تویفل، سیدهارتان و بچلور، ۲۰۰۹).

۷- بحث

با وجود اینکه بیش از ۵۰ سال است که پژوهش در حوزه خلاصه‌سازی خودکار متن آغاز شده است اما هنوز چالش اصلی، تولید خلاصه‌های با کیفیت است. بنابراین اکثر پژوهش‌ها بر بهبود نتایج تمرکز دارند. یکی از رویکردهایی که به منظور بهبود کیفیت خلاصه‌ها مورد توجه جامعه علمی قرار گرفته است، خلاصه‌سازی گفتمان مدار می باشد که همان گونه که در بخش پیشین تشریح شد، تلاش می کند اطلاعات ساختارمند مفیدی از نوشتارهای علمی استخراج کند تا در مرحله بعد بتوان آنها را در ابزارهای خودکار خلاصه‌سازی متن مورد استفاده قرار داد. با اختصاص برچسب های گفتمانی به جملات و با بهره گیری از روش های یادگیری ماشینی نظارتی می توان خلاصه هایی از این دست را تولید نمود.

در ژانر علمی، می توان رویکرد خلاصه‌سازی گفتمان مدار را با رویکرد استنادی تلفیق کرد آن را به عنوان راهکاری برای تولید الگوریتم هایی که قادر هستند خلاصه هایی با خوانایی بالاتر که محتوای اطلاعاتی مفیدتر و بیشتری دارند، ارائه داد.

در خلاصه‌سازی استنادی از مجموعه‌ای از استنادها به یک مقاله مرجع، برای تولید خلاصه بهره می برند (قزوینیان و رادو، ۲۰۰۸؛ قزوینیان و همکاران، ۲۰۱۳). مجموعه استنادها به یک مقاله، را در واقع می توان خلاصه ای از آن مقاله دانست که توسط جامعه علمی تولید شده است (الکیس و همکاران، ۲۰۰۸؛ قزوینیان و همکاران، ۲۰۱۳؛ کهن، سولداینی^۱ و گوهریان، ۲۰۱۴) که در بردارنده

^۱ Soldaini

مهمترین نکات مقاله اصلی و در واقع سهم آن اثر در جامعه علمی است. به این ترتیب، یک خلاصه استنادی با فراهم آوردن اطلاعات اضافی نسبت به چکیده نویسنده، می‌تواند بینش عمیق‌تری راجع به تأثیر آن مقاله در جامعه علمی را فراهم آورد.

اختصاص چهریزه‌های گفتمانی به جملات استنادی مقالات استنادکننده می‌تواند در طبقه‌بندی و رتبه‌بندی آنها به منظور درج در خلاصه‌نهایی مفید واقع شود. به عبارت دیگر، از طریق تحلیل گفتمان جملات استنادی مقالات استنادکننده می‌توان عملکرد استنادی را نیز معلوم کرد و بر اساس آنها جهت تولید خلاصه‌های منسجم‌تر اقدام نمود.

۸- نتیجه‌گیری

در سال‌های اخیر، رویکردهای گفتمانی خلاصه‌سازی ژانر علمی گسترش یافته‌اند. در این رویکردها فرض بر این است که با تکیه بر ویژگی‌های گفتمانی متن، نه تنها می‌توان خلاصه‌هایی با خوانایی بالاتر تولید کرد بلکه امکان بازنمون خودکار اطلاعات ساختارمند مفید برای پژوهشگران نیز از متون علمی فراهم می‌شود. پژوهش‌های انجام شده با این رویکرد، از نشانه‌های گفتمانی برای ناحیه‌بندی بلاغی متون علمی بهره می‌برند.

از کاربردهای مهم و ویژه تحلیل گفتمان در خلاصه‌سازی خودکار ژانر علمی، تولید خلاصه‌های استنادی گفتمان‌مدار می‌باشد. از آنجایی که تحلیل گفتمان جملات پیکره‌های علمی و اختصاص برچسب‌های گفتمانی و ناحیه‌بندی استدلالی آنها توسط متخصصان، به صرف وقت و هزینه زیاد نیازمند است. لازم است مؤسسات تحقیقاتی پشتیبان این مهم باشند. زیرا در صورت تحقق آن برای پیکره‌های عظیم ژانر علمی می‌توان انواع خلاصه‌های خودکار مبتنی بر پرسش‌کاربر، خلاصه‌های چندسندی که تضادها و شباهت‌هاشان در خلاصه‌نهایی مشهود است و نیز خلاصه‌های استنادمدار را تولید نمود و بدین وسیله خدمت بزرگی به بازیابی اطلاعات علمی مورد نیاز پژوهشگران ارائه داد.

۹- پیشنهادات

با توجه به مباحث مطرح شده در این مقاله پیشنهاد می‌شود گروه‌های تحقیقاتی متشکل از متخصصین علم اطلاعات و دانش‌شناسی، زبان‌شناسان و مهندسی نرم‌افزار در دانشگاه‌ها و مؤسسات تحقیقاتی کشور برای تولید پیکره‌های حاشیه‌نویسی شده ژانر علمی به زبان فارسی ایجاد شود. همچنین، در پایان نامه‌های رشته زبان‌شناسی که به تجزیه و تحلیل کلام در متون

علمی می‌پردازند، رویکرد کیفی بیشتر از پیش تقویت گردد تا بتوان با اجتناب از گزارش‌های صرفاً آماری نشانه‌های گفتمانی موجود در متون علمی به سمت کاربردی کردن تحلیل گفتمان علمی در حوزه‌هایی چون بازیابی اطلاعات خودکار حرکت نمود و از این طریق مطالعات بین‌رشته‌ای پرثمرتری را انجام داد.

مراجع

البرزی ورکی، پرویز (۱۳۸۶). «مبانی زبان‌شناسی متن». تهران: انتشارات امیر کبیر.
بشیری، حسن (۱۳۹۴). زبان‌شناسی و تحلیل گفتمان اخبار: رویکردی عملیاتی، فصلنامه نقد کتاب: اطلاع‌رسانی و ارتباطات، شماره ۷ و ۸، ۲.
کلاتری، صمد و دیگران (۱۳۸۸). تحلیل گفتمان با تاکید بر گفتمان انتقادی به عنوان روش تحقیق کیفی. جامعه‌شناسی. شماره ۱، ۴، ۲۸-۷.

Abu-Jbara, A., and Radev, D. (2011, June). Coherent citation-based summarization of scientific papers. In Proceedings of the 49th Annual Meeting of the Association for Computational Linguistics: Human Language Technologies-Volume 1. Association for Computational Linguistics. pp. 500-509.

Baxendale, P. B. (1958). Machine-made index for technical literature—an experiment. IBM Journal of Research and Development, 2(4). pp. 354-361.

Bosma, W. E. (2008). Discourse oriented summarization. University of Twente.

Brown, G and Yule, G. (1983) Discourse Analysis, Cambridge University Press.

Cohan, A., and Goharian, N. (2015). Scientific article summarization using citation-context and article's discourse structure. Proceedings of the 2015 Conference on Empirical Methods in Natural Language Processing. Lisbon, Portugal, 17-21 September. pp. 390-400.

Cohan, A., Soldaini, L., and Goharian, N. (2014). Towards citation-based summarization of biomedical literature. In Proceedings of the Text Analysis Conference (TAC'14).

DeJong, G. (1982). An overview of the FRUMP system. Strategies for natural language processing, 113. pp. 149-176.



Elkiss, A., Shen, S., Fader, A., Erkan, G., States, D., and Radev, D. (2008). Blind men and elephants: What do citation summaries tell us about a research article?. *Journal of the American Society for Information Science and Technology*, 59(1). pp. 51-62.

Fisas, B., Ronzano, F., and Saggion, H. (2015, June). On the discursive structure of computer graphics research papers. In *The 9th Linguistic Annotation Workshop held in conjunction with NAACL 2015*. pp. 42.

Fraser, B. (1999). What are discourse markers?. *Journal of Pragmatics*, 31. pp. 931-952.

Jones, K. S. (1993). What might be in a summary?. *Information retrieval*, 93. pp. 9-26.

Jones, K. S., and Endres-Niggemeyer, B. (1995). Automatic summarizing. *Information Processing & Management*, 31(5). pp. 625-630.

Kupiec, J., Pedersen, J., and Chen, F. (1995, July). A trainable document summarizer. In *Proceedings of the 18th annual international ACM SIGIR conference on Research and development in information retrieval*. ACM. pp. 68-73.

Lloret, E., and Palomar, M. (2013). COMPENDIUM: a text summarisation tool for generating summaries of multiple purposes, domains, and genres. *Natural Language Engineering*, 19(2). pp.147-186.

Luhn, H. P. (1958). The automatic creation of literature abstracts. *IBM Journal of research and development*, 2(2). pp.159-165.

Mani, I. (2001). *Automatic summarization (Vol. 3)*. John Benjamins Publishing.

Mann, W. C., and Thompson, S. A. (1988). Rhetorical structure theory: Toward a functional theory of text organization. *Text-Interdisciplinary Journal for the Study of Discourse*, 8(3), 243-281.

Qazvinian, V., and Radev, D. R. (2008, August). Scientific paper summarization using citation summary networks. In *Proceedings of the 22nd International Conference on Computational Linguistics-Volume 1 Association for Computational Linguistics*. pp. 689-696.

Qazvinian, V., and Radev, D. R. (2010, July). Identifying non-explicit citing sentences for citation-based summarization. In *Proceedings of the 48th annual meeting of the association for computational linguistics*. Association for Computational Linguistics. pp. 555-564.



Radev, D. R., Hovy, E., and McKeown, K. (2002). Introduction to the special issue on summarization. *Computational linguistics*, 28(4). pp. 399-408.

Rahman, N., and Borah, B. (2015, September). A survey on existing extractive techniques for query-based text summarization. In *Advanced Computing and Communication (ISACC), 2015 International Symposium on IEEE*. pp. 98-102.

Saggion, H., and Poibeau, T. (2013). Automatic text summarization: Past, present and future. In *Multi-source, multilingual information extraction and summarization* (pp. 3-21). Springer Berlin Heidelberg.

Schiffrin, D. (1987). *Discourse Markers*. Cambridge: Cambridge University Press.

Teufel, S. (1999). *Argumentative Zoning: Information Extraction from Scientific Text*. School of Cognitive Science, University of Edinburg, UK..

Teufel, S., and Moens, M. (2002). Summarizing scientific articles: experiments with relevance and rhetorical status. *Computational linguistics*, 28(4). pp.409-445.

Teufel, S., Siddharthan, A., and Batchelor, C. (2009, August). Towards discipline-independent argumentative zoning: evidence from chemistry and computational linguistics. In *Proceedings of the 2009 Conference on Empirical Methods in Natural Language Processing: Volume 3-Volume 3* . Association for Computational Linguistics. pp. 1493-1502.

Yang, S. H. A. N. R. U. (2011). Investigating discourse markers in pedagogical settings: a literature review. *ARECLS*, 8. pp. 95-108.

Yeloglu, O., Milios, E., and Zincir-Heywood, N. (2011, March). Multi-document summarization of scientific corpora. In *Proceedings of the 2011 ACM Symposium on Applied Computing* .ACM. pp. 252-258.