


## Evaluating the Structure of the Inverted Pyramid in the Big Persian News Corpus: News Discourse Analysis based on the Correlation Coefficient between Title and News Content

Ghayoomi, Masood<sup>1</sup> 

Assistant Professor of Linguistics, Institute for Humanities and Cultural Studies, Tehran, Iran

### Abstract


News discourse is a type of discourse analysis that deals with the analysis of news discourse. Due to the fact that in the formatting of news there are two hidden features of selection and prominence in the communication representation of news, the inverted pyramid of news is used to grade the importance of the discourse parts of the news. Although it is desirable to meet the structure of the inverted pyramid of news, sometimes this structure may change. In this article, we put an effort to analyse the discourse analysis of Persian news websites with the help of statistical analysis. To research the goal, data science can be used. This inter-discipline deals with data analysis from a scientific aspect, finding implicit concepts to be obtained from data analysis and extracting knowledge from the data. In the framework of data science, we examined the Persian news corpus and studied the existence of semantic correlation between the news title and the news content based on the structure of the news inverted pyramid. To achieve the goal, by using the crawling method, a relatively large news corpus with a volume of 14 billion words has been obtained from 24 news websites. After pre-processing and normalizing the corpus, in the framework of distributional semantics, the vector of title news and content have been created by using the Word2Vec tool for creating the vector model to have the vector representation of each news. After segmenting news content into three parts (lead, body and further explanation about the lead) according to the inverted pyramid, the Pearson correlation coefficient has been used to calculate the correlation between the title and each part of the news. Although Pearson's correlation coefficient was positive for a large number of news, zero value and no correlation was found for the news. On average, the correlation between the headline and the news lead and body was higher than the correlation between the headline and the lead development. This research can be used as a method to carefully select the title and content and filter the news according to the inverted pyramid structure.

**Keywords:** news corpus, Pearson correlation coefficient, distributional semantics, Word2Wok, news inverted pyramid

1. m.ghayoomi@ihcs.ac.ir

**How to cite:** Ghayoomi, M. (2023). Evaluating the Structure of the Inverted Pyramid in the Big Persian News Corpus: News Discourse Analysis based on the Correlation Coefficient between Title and News Content. *Language and Linguistics*, 18(35), 21-45. doi: 10.30465/lsi.2023.8497

## ارزیابی ساختار هرم وارونه در پیکره بزرگ خبری فارسی: تحلیل گفتمان خبری براساس همبستگی میان عنوان و محتوای خبر

قیومی، مسعود  | استادیار پژوهشگاه علوم انسانی و مطالعات فرهنگی، تهران، ایران

**چکیده:** گفتمان خبری گونه‌ای تحلیل گفتمان است که به تحلیل ساختار گفتمان خبری می‌پردازد. با توجه به این که در قالب‌بندی اخبار دو ویژگی انتخاب و برجستگی در نمود ارتباطی خبر نهفته است، از ساختار هرم وارونه خبر برای درجه‌بندی اهمیت بخش‌های گفتمانی خبر استفاده می‌شود. اگرچه رعایت ساختار هرم وارونه خبر مطلوب است، گاهی ممکن است در گفتمان خبری این ساختار دچار تغییر شود که در این مقاله تلاش می‌شود با کمک تحلیل آماری، به تحلیل ساختار گفتمان و بگانه‌های خبری فارسی پرداخته شود. برای رسیدن به هدف می‌توان از علم داده استفاده کرد. این بین‌رشته‌ای از جنبه علمی به تحلیل داده، یافتن مفاهیم ضمنی به‌دست‌آمده از تحلیل داده‌ها و استخراج دانش از داده‌ها می‌پردازد. در چارچوب علم داده به بررسی پیکره‌ای متون خبری فارسی پرداخته شده و وجود رابطه همبستگی معنایی میان عنوان خبر و محتوای خبر در ساختار هرم وارونه خبر مورد مطالعه قرار می‌گیرد. برای دستیابی به این هدف، با استفاده از روش خزش، یک پیکره خبری نسبتاً بزرگ با حجمی بالغ بر ۱۴ میلیارد واژه از ۲۴ وبگاه خبری به‌دست آمده است. پس از پیش‌پردازش و اعمال یکدستی نسبی در این پیکره، در چارچوب معناشناسی توزیعی، بردار عنوان خبر و متن خبر با استفاده از مدل بردارسازی واژه ورد و ک به‌دست آمده و براساس آن بردار هر خبر ساخته شده است. پس از بخش‌بندی محتوای هر خبر براساس هرم وارونه خبر به سه قسمت سرخ (لید)، بدنه و توسعه سرخ، با استفاده از

ضریب همبستگی پیرسون، میزان همبستگی میان عنوان و هریک از سه بخش خبر محاسبه شد. اگرچه ضریب همبستگی پیرسون برای حجم زیادی از خبرها مثبت بود، ارزش صفر و عدم وجود همبستگی برای خبرها یافت شد. به‌طور متوسط، همبستگی میان عنوان و بدنه خبر بیش از همبستگی میان عنوان و توسعه سرنخ بود. این پژوهش می‌تواند به‌عنوان روشی برای دقت در انتخاب عنوان و محتوا و پالایش خبر منطبق بر هرم وارونه استفاده شود.

**کلیدواژه‌ها:** پیکره خبری، ضریب همبستگی پیرسون، معناشناسی توزیعی، وردآوک، هرم وارونه خبر.

## ۱ مقدمه

تحلیل گفتمان یکی از حوزه‌های زبان‌شناسی نظری است که علاوه بر نقش تعاملی خود در جامعه برای انتقال اطلاعات و توصیف روابط اجتماعی و رویکردهای شخصی، دارای نقش ارجاعی در توصیف محتوا و تحلیل کاربردی زبان نوشتاری و گفتاری یا نشانه‌ها در زبان است (براون<sup>۱</sup> و یول<sup>۲</sup>، ۱۹۸۳: ۱-۲). اگرچه هدف اولیه زبان برای اهداف تعاملی است، می‌توان از آن برای اهداف ارجاعی و روابط اجتماعی نیز بهره برد. براون و یول (۱۹۸۳: ۶ و ۱۹۰) مفاهیم خاصی را برای متن تعریف کرده‌اند. بنا بر نظر آنان، هرگونه ثبت شفاهی عملکرد ارتباطی زبان، متن است؛ و همچنین متن چیزی جز پیوستگی معنایی میان جملات تشکیل‌دهنده آن نیست. پر واضح است الزاماً همیشه عملکرد ارتباطی زبان به‌صورت شفاهی نیست و گونه مکتوب، مانند متن خبری، کتاب و مقاله نیز وجود دارد تا ارتباط زبانی میان تولیدکننده محتوا و دریافت‌کننده محتوا حاصل شود. در نوع مکتوب متن که از طریق خط میسر می‌شود، از زبان در نقش تعاملی خود برای ارتباط فرد با رویدادهای محیط پیرامون استفاده می‌شود. دویوگراند<sup>۳</sup> و درسلر<sup>۴</sup> (۱۹۸۱: ۳) متن را نوعی نمود ارتباطی دانسته و هفت ویژگی را برای آن مطرح کرده‌اند. متن خبری کم‌وبیش حاوی ویژگی‌های ذکرشده است که مطالعه در این حوزه در چارچوب تحلیل گفتمان، با اصطلاح «گفتمان خبری»<sup>۵</sup> (بدنارک<sup>۶</sup> و کپل<sup>۷</sup>، ۲۰۱۲) شناخته می‌شود.

1. G. Brown

2. G. Yule

3. R. de Beaugrand

4. W. U. Dressler

5. News discourse

6. M. Bednarek

7. H. Caple

باتوجه به نقش تعاملی زبان در خبر، باید ارتباط متنی با مخاطب دارای ویژگی تأثیرگذاری باشد (دوبوگراند و درسلر، ۱۹۸۱: ۳۱). برای تحقق این ویژگی نیاز است عناوین خبری به گونه‌ای برای مخاطب جذاب باشد؛ به عبارتی دیگر، زبان خبر باید به گونه‌ای طراحی شود که برای مخاطب جاذبه داشته باشد که بل<sup>۱</sup> (۱۹۹۱: ۱۲۱) از آن با اصطلاح «طراحی مخاطب»<sup>۲</sup> یاد می‌کند.

برای بررسی ویژگی تعاملی زبان و رعایت تأثیرگذاری بر مخاطب از طریق ایجاد جذابیت در عنوان خبر و ارتباط عنوان با محتوای خبر، در این پژوهش که اساساً ماهیت بین‌رشته‌ای دارد می‌کوشیم با کمک تحلیل آماری در چارچوب تحلیل گفتمان خبری و «علم داده»<sup>۳</sup>، به بررسی همبستگی عنوان و بخش‌های تشکیل‌دهنده خبر پردازیم تا بتوان به این پرسش‌ها پاسخ داد که آیا به‌هنگام تهیه و انتشار خبر نکاتی که به‌صورت نظری مطرح است رعایت می‌شود؟ و برای جذب مخاطب، عناوین خبر متبلورکننده کدام بخش از خبر است؟

ساختار مقاله حاضر به این صورت است که پس از مقدمه، در بخش ۲ مطالعات انجام‌شده در حوزه تحلیل ساختار خبر مرور می‌شود. در بخش ۳، چارچوب نظری چندبعدی در حوزه گفتمان خبری، معرفی هرم وارونه خبر و همچنین علم داده ارائه می‌شود. شیوه گردآوری داده‌ها، تهیه پیکره زبانی و شیوه انجام پژوهش در بخش ۴ توضیح داده می‌شود. در بخش ۵، به تحلیل داده‌ها و پاسخ به پرسش‌ها پرداخته شده و در نهایت، در بخش ۶، مقاله با نتیجه‌گیری خاتمه می‌یابد.

## ۲ پیشینه مطالعاتی

از آنجاکه تحلیل گفتمان خبری در چارچوب هرم وارونه خبر و همبستگی معنایی در خبر محوریت موضوع این پژوهش را شکل می‌دهد، مطالعات انجام‌شده با موضوعات متمرکز بر تحلیل گفتمان با استفاده از هرم خبر و همبستگی محتوایی بررسی می‌شود.

ژانگ<sup>۴</sup> و لیو<sup>۵</sup> (۲۰۱۶) در پژوهش خود تلاش کرده‌اند که با استفاده از قیاس بین گفتمان و درخت تجزیه وابستگی جملات به طراحی چکیده متن پردازند. برای این هدف، ۳۵۹

1. A. Bell

2. audience design

3. data science

4. H. Zhang

5. H. Liu

مقاله مجله وال استریت<sup>۱</sup> که حاوی دادگان درختی گفتمان مبتنی بر نظریه ساختار بلاغی شکل گرفته است را انتخاب کرده و هر درخت گفتمان را به سه درخت وابستگی در سطوح گفتمان، پاراگراف و جمله تبدیل کرده‌اند. آنها به‌طور تجربی ساختار «خلاصه + جزئیات» یا «هرم وارونه»<sup>۲</sup> گفتمان خبری را براساس داده‌های موجود مطالعه کرده‌اند.

پتیربرگ<sup>۳</sup> (۲۰۱۰) به مقایسه گزارش‌های خبری تلویزیونی اروپا و اسکاندیناوی در دهه‌های ۱۹۸۰ و ۱۹۹۰ با ساختار معیار هرم وارونه خبر پرداخته است. براساس یافته‌ها، در مقالات خبری تغییری صورت گرفته است به این صورت که از عنینت‌گرایی «پارادایم خبری» به سمت آرمان‌گرایی و «داستان‌سرایی» تغییر جهت صورت گرفته است و از ساختار هرم وارونه خبر خارج شده است.

امده<sup>۴</sup>، کلیمت<sup>۵</sup> و اشلوتس<sup>۶</sup> (۲۰۱۶) درک خبر توسط نوجوانان ۱۲ تا ۱۷ ساله را مورد مطالعه قرار داده‌اند. در این پژوهش شیوه سستی پخش اخبار در رسانه که مبتنی بر ساختار هرم وارونه خبر است را با شیوه روایی خبر مقایسه کرده‌اند. براساس نتایج این پژوهش، به‌طور کلی خبر روایی نتوانسته است سبب افزایش درک خبر در نوجوانان شود.

نورامبونا<sup>۷</sup>، هورنینگ<sup>۸</sup> و میترا<sup>۹</sup> (۲۰۲۰) به ارزیابی هرم وارونه خبر و امتیازدهی به بخش‌های هرم به‌هنگام استخراج اطلاعات برای دو موضوع الف) شیوه «۵ سؤال و ۱ چگونگی» در ساختار هرم وارونه خبر که با اصطلاح «5WIH» شناخته می‌شود و در برجسب‌گذاری بخش‌های مختلف خبر به‌کار رفته است و ب) خلاصه‌سازی متن پرداخته‌اند. در این پژوهش، از مقالات اسوشیتدپرس<sup>۱۰</sup> مربوط به خبرهای فوری و غیرفوری در ژانر اقتصادی و سیاسی استفاده شده است. براساس نتایج به‌دست آمده، ساختار خبر در اخبار فوری و غیرفوری متفاوت است و اخبار فوری از ساختار هرم وارونه اخبار بیشتر استفاده می‌کند.

مطالعات انجام شده در فارسی بیشتر در چارچوب تحلیل گفتمان تصویری بوده و به بررسی همبستگی معنایی عکس به‌کاررفته در متن و محتوای متن پرداخته‌اند. اقبالی (۱۳۹۰) به بررسی تصویرسازی کتاب‌های داستان کودک بین ۱۳۴۰ تا ۱۳۸۰ پرداخته و ۵۰ نمونه را بررسی کرده است. براساس یافته‌های وی، الزاماً همیشه همبستگی معنایی در کتب داستان کودکان رعایت نشده و گاهی تصویرسازی مغایر با متن داستان بوده است.

1. Wall Street Journal  
4. K. Emde  
7. B. K. Norambuena  
10. Associated Press

2. inverted pyramid  
5. C. Klimmt  
8. M. A. Horning

3. E. Ytreberg  
6. D. M. Schlütz  
9. T. Mitra

عظیمی فرد و همکاران (۱۳۹۶) در چارچوب تحلیل گفتمان تصویری، به بررسی همبستگی معنایی عکس به کاررفته در خبر و محتوای خبری پرداخته‌اند. در این پژوهش، سه وبگاه خبری خبر، العالم و پرس تی وی با یکدیگر مقایسه شده‌است. در این مطالعه، به صورت تصادفی، ۳۰ خبر از مجموع ۶۶۴ خبر به عنوان داده‌های خبر انتخاب شده و به ارزیابی کیفی آنها پرداخته شده‌است. براساس تحلیل انجام شده، در ۱۶ خبر همبستگی عکس و متن خبر رعایت شده‌است، در ۶ خبر تاحدودی این همبستگی رعایت شده و در ۸ خبر عدم وجود همبستگی دیده شده‌است. شایان ذکر است در مقایسه انجام شده، شبکه پرس تی وی عملکرد بهتری در انتخاب تصویر و ارتباط آن با محتوای خبر در مقایسه با شبکه العالم و خبر داشته‌است.

اردکانی فرد و همکاران (۱۴۰۰) به تحلیل محتوایی خبرگزاری‌های رسمی در سال ۱۳۹۹ با محوریت دین پرداخته‌اند و ویژگی‌های این دسته از اخبار را استخراج کرده‌اند. براساس نتایج به دست آمده از این پژوهش، ساختار معیار هرم وارونه خبر در این گونه خبری به این صورت رعایت شده‌است که ابتدا تازگی خبر و سپس رویکرد تبلیغی در خبر حفظ شده‌است. همچنین، ضمن رعایت گفتمان سستی نسبت به دین، به ابعاد سیاسی، اجتماعی و مناسکی دین توجه شده‌است.

### ۳ چارچوب نظری

در انجام پژوهش حاضر سه موضوع دخیل است که اساس نظری این مطالعه را شکل می‌دهد: الف) یک موضوع تحلیل گفتمان خبر و ویژگی‌های آن است؛ ب) موضوع دیگر علم داده برای پاسخ به سؤالات پژوهش در حجم زیاد داده با کمک رایانه، آمار و روش‌های پردازشی است؛ و ج) بازنمایی معنایی واژه‌ها در چارچوب «معناشناسی توزیعی»<sup>۱</sup> است. در ادامه، نظریاتی که چارچوب این پژوهش را شکل می‌دهد توضیح داده می‌شود.

#### ۳-۱ تحلیل گفتمان خبر

در زبان‌شناسی ساختگرایی، دوسوسور<sup>۲</sup> (۱۹۱۶) دو سطح صورت و معنا را برای زبان قائل شده‌است که از طریق آوا در شنیدار یا خط در نوشتار تجلی یافته و به مخاطب منتقل

1. distributional semantics

2. F. de Saussure

می‌شود. خبر نوعی متن است که از طریق نوشتار، پیامی که دارای معنا است را به مخاطب منتقل می‌کند. گفتمان خبری زیرمجموعه تحلیل گفتمان است که به بررسی زبان در خبر می‌پردازد. خبر که گونه‌ای از متن است در نقش ارتباطی خود حاوی اطلاعات خارق‌العاده و خاصی است که برای جلب توجه مردم منتشر می‌شود. برای مثال، اگر گفته شود که «سگی فردی را گاز گرفته‌است» خبر نیست؛ چون اغلب اوقات این اتفاق می‌افتد و حاوی اطلاعات جدید نیست. درحالی‌که اگر «مردی سگی را گاز بگیرد» این رویداد یک خبر است (ایتول<sup>۱</sup> و اندرسون<sup>۲</sup>، ۲۰۰۶) و حاوی اطلاعات غیرقابل پیش‌بینی است. ژنی<sup>۳</sup> (۲۰۱۸) خبر را اطلاعاتی در مورد یک رویداد می‌داند و مطرح می‌کند که خبر به‌عنوان یک پنجره در جهت به‌دست‌آوردن دانش، ایفای نقش می‌کند. بنابراین، خبر نوعی نمود زبانی در ایجاد ارتباط میان فرد و اتفاقات محیط اطراف است و از نقش تعاملی زبان بهره‌برده می‌شود. همان‌طور که در مقدمه به آن اشاره شد، دوبوگراند و درس‌لر (۱۹۸۱: ۳) هفت ویژگی را برای نمود ارتباطی متن اشاره کرده‌اند که عبارت است از الف) وجود انسجام، ب) وجود پیوستگی معنایی، پ) وجود هدفمندی در نویسنده برای خلق متن، ت) وجود تبادل اطلاعات و انتظارات میان خواننده و شنونده در اطلاعات تبادل‌شده، ث) وجود عوامل مرتبط با تولید متن و ج) وجود ارتباط در تولید درک متن حاضر با تولید و درک متن پیشین. دو ویژگی انسجام و پیوستگی معنایی متن-محور و مابقی ویژگی‌ها کاربرد-محور است (آقاگلزاده، ۱۳۹۴: ۱۰۹). در متن خبری می‌توان این هفت ویژگی را یافت؛ از این‌رو، این نوع داده از درجه اهمیت بالایی برخوردار است.

در چارچوب گفتمان خبری، یک ساختار کلی برای خبر معرفی شده‌است که به «هرم وارونه»<sup>۴</sup> معروف است (بدنارک و کپل، ۲۰۱۲: ۱۰۰-۹۶). این چارچوب در قرن ۱۸ به‌عنوان ابزار مناسب انتشار اطلاعات معرفی شده‌است (پوتکر<sup>۵</sup>، ۲۰۰۳). کاربرد عملی این چارچوب در اوایل قرن ۱۹ هم‌زمان با اکتشاف تلگراف قابل مشاهده بود چراکه براساس اهمیت، موضوعات در ارتباطات مطرح می‌شد (کاناویلهاس<sup>۶</sup>، ۲۰۰۷). در شکل ۱ ساختار هرم وارونه خبر نمایش داده شده‌است. هر خبر متشکل از چهار بخش الف) عنوان یا سرخط خبر، ب) مقدمه یا سرخط خبر، پ) بدنه خبر و ت) اطلاعات تکمیلی جهت توسعه سرخط خبر است. در سرخط، خلاصه رویداد به‌گونه‌ای که برای خواننده ایجاد جذابیت کند

1. B. Itule

4. inverted pyramid

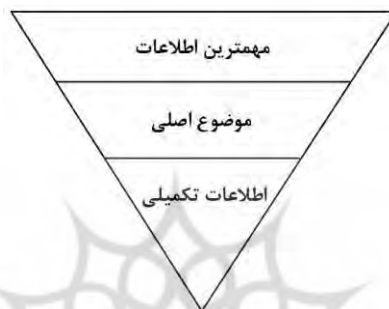
2. D. Anderson

5. H. Pöttker

3. Q. Xie

6. J. Canavilhas

بسیار کوتاه و موجز بیان می‌شود. ویژگی سرخط این است که معمولاً به صورت عبارت و به ندرت به صورت جمله بیان می‌شود. در بخش مقدمه یا سرخط، جنبه‌های جذاب رویداد در اندازه یک جمله کوتاه یک سطر بیان می‌شود. به عبارتی دیگر، در قسمت سرخط که تکمیل عنوان خبر است، مهمترین اطلاعات به صورت جمله بیان می‌شود. در بخش‌های بعدی که بدنه و توسعه سرخط است به بسط موضوع اصلی خبر و ارائه اطلاعات تکمیلی می‌پردازد.



شکل ۱- هرم وارونه خبر

اگرچه در هرم وارونه اطلاعات از نظر درجه اهمیت مرتب شده است و رعایت این ساختار برای روزنامه‌نگاران مهم است؛ نگاه زبان‌شناسان در تحلیل گفتمان خبر به این صورت است که روابط منطقی یا بلاغی میان عنوان یا سرخط و بقیه اجزای خبر (بل، ۱۹۹۱: ۱۷۰) یا پیوند میان بخش‌های عنوان و سرخط (فیزا، ادما<sup>۲</sup> و وایت<sup>۳</sup>، ۲۰۰۸) مورد توجه است. در پژوهش حاضر تلاش می‌شود پیوند و همبستگی معنایی میان عنوان و بخش‌های مختلف خبر در چارچوب علم داده بررسی شود. *گامه علوم انسانی و مطالعات فرهنگی*

## ۲-۳ علم داده

علم داده در دهه‌های ۱۹۶۰ و ۱۹۷۰ با مطرح شدن کشف اطلاعات و تحلیل محتوای معنایی شکل گرفت (کائو<sup>۴</sup>، ۲۰۱۷). روش‌شناسی علم داده در چارچوب علم آمار و علم رایانه تکامل یافت؛ تا این که در سال ۱۹۹۸ مفهوم جدید علم داده توسط «چیکو هایشی»<sup>۵</sup>

1. S. Feez  
4. L. Cao

2. R. Edema  
5. Chikio Hayashi

3. P. R. R. White



معرفی شد (مورتاگ<sup>۱</sup> و دولین<sup>۲</sup>، ۲۰۱۸). کلوند<sup>۳</sup> (۲۰۰۱) مفهوم امروزی علم داده را این‌گونه معرفی می‌کند: علم داده یک بین‌رشته‌ای است که در آن داده‌های فاقد ساختار، ساختارمند شده و با بهره‌گیری از روش‌های الگوریتمی پردازش اطلاعات می‌توان به استخراج دانش و افزایش بینش پرداخت. بنابراین، داده، اطلاعات و دانش ارکان هرم دانش را در علم داده می‌سازد که در شکل ۲ نمایش داده شده‌است. در تعریف دالکیر<sup>۴</sup> (۲۰۰۵)، «داده» محتوای قابل دیدن یا قابل تغییر، «اطلاعات» بازنمایی داده تحلیل شده و «دانش» اطلاعات نظری و مفید است. در پژوهش حاضر، داده همان متن خام خبری است؛ اطلاعات بخش‌های محتوای خبری و دانش یافتن همبستگی معنایی میان عنوان و بخش‌های خبری است.



شکل ۲- هرم دانش

در حوزه علم داده، از آمار و روش‌های «یادگیری ماشین»<sup>۵</sup> در هوش مصنوعی برای ساخت الگوریتم‌های پردازشی جدید بهره برده می‌شود؛ و براساس نیازهای تعریف‌شده در این حوزه، کار ذخیره‌سازی و استخراج اطلاعات از داده به‌صورت الگوریتمی انجام می‌پذیرد. همچنین، تحلیل داده برای حل مسائل پیچیده و نمایش بصری اطلاعات نیز به‌کار می‌رود تا به استخراج دانش و افزایش بینش از داده منجر شود. برای تحقق اهداف علم داده، به مدل‌سازی داده نیاز است. در مدل‌سازی داده، تمام فرایندهای پردازش داده به‌صورت روشمند و الگوریتمی تعریف و سازماندهی می‌شود. از آنجاکه انواع داده‌های متنی، صوتی، تصویری و

1. F. Murtagh  
4. K. Dalkir

2. K. Devlin  
5. machine learning

3. W. S. Cleveland

عددی وجود دارد، نحوه جمع‌آوری، ذخیره‌سازی، ساماندهی و استفاده از این اطلاعات متفاوت است.

### ۳-۳ معناشناسی توزیعی

مدل‌سازی داده بخش مهم علم داده است که روش‌شناسی انجام پژوهش را در دل خود دارد و بر فرایند تحلیل اثربخش است. برای بررسی وجود همبستگی معنایی میان عنوان خبر و سایر بخش‌های سازنده خبر در پژوهش حاضر نیاز است این تحلیل در چارچوب معناشناسی توزیعی انجام پذیرد. هریس<sup>۱</sup> (۱۹۵۴) معتقد است که واژه‌هایی که در یک بافت زبانی یکسان به کار می‌رود، این تمایل را دارند که از نظر معنایی به یکدیگر شبیه<sup>۲</sup> باشد. بنا بر نظر وی، معنای هر واژه منعکس‌کننده بافتی است که آن واژه در آن بافت به کار می‌رود. نظر هریس در چارچوب معناشناسی توزیعی قرار می‌گیرد که منجر به معرفی شدن «فرضیه توزیعی»<sup>۳</sup> شده است. بر اساس این فرضیه، واژه‌هایی که تمایل به حضور در یک بافت یکسان دارد، از نظر معنایی مشابه یکدیگر است. بنا بر نظر هریس، معنای واژه تأثیرپذیر از «بافت جایگاهی»<sup>۴</sup> است که آن واژه در آن بافت ظاهر شده است. در همین راستا، فرث<sup>۵</sup> (۱۹۵۷) می‌افزاید که با توجه به واژه‌های اطراف یک واژه می‌توان معنای یک واژه را مشخص کرد. نتیجه این نظرات آن است که بافت زبانی واژه، نقش بسیار مهمی در تعیین معنای یک واژه دارد.

در چارچوب معناشناسی توزیعی، میکولوو<sup>۶</sup> و همکاران (۲۰۱۳) مدلی را معرفی کرده‌اند که می‌توان معنای واژه را به صورت بردار<sup>۷</sup> بازنمایی کرد. در این شیوه بازنمایی معنایی، به جای صورت واژه از بردار واژه که براساس بافت جایگاهی واژه در یک پیکره زبانی به دست آمده است استفاده می‌شود.

### ۴ جمع‌آوری داده

انجام هر پژوهشی به داده نیاز دارد و هر قدر حجم داده‌های پژوهش بیشتر باشد، نتایج به دست آمده موثق‌تر است. هدف اصلی این پژوهش یافتن همبستگی میان عنوان و محتوای

1. Z. S. Harris

4. local context

7. vector

2. similar

5. J. R. Firth

3. distributional hypothesis

6. T. Mikolov

خبر است. از این رو، به یک پیکره خبری نیاز است که می‌توان به واسطه خزش از وبگاه‌های خبری به تهیه این پیکره اقدام کرد. دسترسی به بایگانی برخط وبگاه‌های خبری در رسیدن به هدف این پژوهش کمک شایانی می‌کند. وجود داده با حجم زیاد سبب می‌شود بتوان از روش‌های پیکره-محور<sup>۱</sup> و نه پیکره-بنیان<sup>۲</sup> (توگنینی-بونلی<sup>۳</sup>، ۲۰۰۱: ۸۴-۸۵) در تحلیل داده‌ها بهره برد. در روش پیکره-محور از علم داده استفاده می‌شود و از بررسی داده به فرضیه می‌رسیم؛ در حالی که در روش پیکره-بنیان فرضیه اولیه‌ای وجود دارد و با بررسی داده به اثبات آن فرضیه می‌پردازیم. اکثر پژوهش‌های انجام‌شده فرضیه‌محور و از نوع پیکره-بنیان است. طرح حاضر می‌تواند زمینه را برای انجام طرح‌های پیکره-محور مبتنی بر علم داده فراهم کند.

#### ۴-۱ اهمیت داده‌های خبری

امروز حجم زیادی از انواع داده‌ها در بستر وب وجود دارد؛ و از میان آنها داده‌ی متنی از حجم و تنوع بیشتری برخوردار است. بنابراین محتوای زیادی از این نوع داده تولید می‌شود. یکی از منابع این نوع داده که به‌طور روزانه به تولید حجم زیادی از محتوا می‌پردازد خبرهای منتشرشده در وبگاه‌های خبری است. این نوع داده که توجه پژوهشگران را به خود جلب کرده‌است از جنبه‌ی زبان‌شناسی و کاربردشناسی اهمیت دارد که در ادامه توضیح داده می‌شود. این نوع داده از تنوع ژانر، مانند سیاسی، اجتماعی، اقتصادی و غیره، برخوردار است که می‌توان به‌صورت الگوریتمی ژانر خبر را تعیین کرد. تولید محتوا در وبگاه‌های خبری روزانه است. بنابراین با تحلیل هم‌زمانی می‌توان به بررسی سیر تحول موضوعات پرداخت. این داده‌ها معمولاً در وبگاه‌های خبری بایگانی می‌شود. بنابراین تا زمان حیات روزنامه، امکان دسترسی به سوابق اخبار وجود دارد. گزارش‌های روزانه از تحولات سیاسی، اجتماعی، اقتصادی و غیره سبب می‌شود بتوان زبان را رصد کرد و از نواژه‌ها، تغییرات معنایی واژه‌ها و همچنین تغییرات زبانی با توجه به مؤلفه‌ی زمان آگاه شد. از آنجاکه خبرها توسط افراد مختلف نگارش می‌شود، سبک نگارشی افراد در داده‌ها مستتر است که می‌تواند در مطالعات سبک‌شناسی نیز به‌کار رود. معمولاً دسترسی به اخبار این وبگاه‌ها رایگان است. در نتیجه، برای گردآوری داده‌ها نباید هزینه‌ای پرداخت شود. از آنجاکه امروزه

1. corpus-driven

2. corpus-based

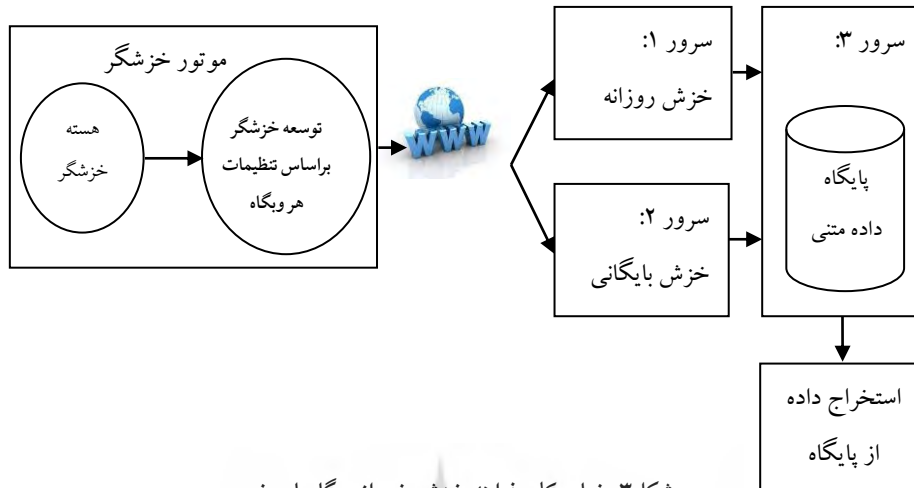
3. E. Tognini-Bonelli

به منظور پردازش زبان طبیعی به حجم مناسب داده برای مدل‌سازی زبان نیاز است، این نوع داده می‌تواند خیلی سریع با کمترین هزینه نیاز الگوریتم‌های پردازشی را برآورد. هر خبر متشکل از چند موضوع است که کنار یکدیگر قرار گرفته‌است. می‌توان موضوع‌های موجود در هر خبر را با کمک الگوریتم «مدل‌سازی موضوع»<sup>۱</sup> (بلای<sup>۲</sup> و همکاران، ۲۰۰۳) استخراج کرد و خبرهای دارای تشابه موضوعی را خوشه‌بندی کرد. علاوه بر آن می‌توان کلیدواژه‌های مرتبط با هر موضوع را نیز استخراج کرد که می‌تواند در تولید هشتگ<sup>۳</sup> برای گروه‌بندی خبرها کاربرد داشته باشد. تحلیل احساسات نهفته در هر خبر و یافتن گرایش و رویکرد خبرگزاری در انتشار خبر موضوع دیگری است که کاربرد دیگر پیکره‌خبری را مشخص می‌کند. تمیزدادن خبر جعلی که با قصد فریب افراد منتشر می‌شود از خبر موثق رویکرد دیگری در تحلیل این نوع داده است. خلاصه‌سازی الگوریتمی خبر نیز موضوعی است که به داده‌های خبری نیاز دارد. علاوه بر موارد مطرح‌شده بالا که کارایی‌های مختلف پیکره‌خبری را مشخص می‌کند می‌توان از این داده در تحلیل‌های دیگری استفاده کرد. حال که اهمیت داده‌های خبری، چه از بعد پژوهشی و چه از بعد کاربردی، مشخص شد، در ادامه، شیوه جمع‌آوری داده‌های مورد نیاز پژوهش حاضر توضیح داده می‌شود.

## ۲-۴ خزش و ذخیره‌سازی اطلاعات

مرحله اول، تهیه «پایگاه داده متنی» است که به عنوان بایگانی کلی از اخبار بتواند برای سایر پژوهش‌های علوم انسانی مبتنی بر داده مورد استفاده قرار گیرد. در این مدل خزش، اخبار به صورت روزانه از وبگاه‌های خبری خزش شده و به داده‌های قبلی در این دو پایگاه داده افزوده می‌شود؛ بنابراین یک فرایند پیوستار و پویا در طراحی جمع‌آوری داده وجود دارد و مقطعی نیست. نتیجه این عملکرد می‌تواند ضمن به‌روزرسانی روزانه داده‌ها در پایگاه، کمک کند تا از این اطلاعات برای رصد تغییرات زبانی و یافتن نواژه‌های وارد شده به زبان کمک گرفت (قیومی، ۱۳۹۹). نمای کلی فرایند خزش خبر از وبگاه‌های خبری در شکل ۲ نمایش داده شده‌است.

1. Topic modeling
2. D. M. Blei
3. Hashtag



شکل ۳- نمای کلی فرایند خزش خبر از وبگاه‌های خبری

همان‌گونه که در شکل ۳ مشخص است، برای جمع‌آوری داده به سه بخش نیاز است: الف) موتور خزشگر از وب، ب) تعریف فرایند استخراج اطلاعات از پایگاه‌های داده و ج) سرور برای خزش و ذخیره‌سازی داده در پایگاه داده نیاز است که این سه بخش در ادامه توضیح داده می‌شود.

به هنگام تهیه موتور خزشگر از وب و تعریف فرایند استخراج اطلاعات، با بررسی‌های انجام‌شده از مجموعه وبگاه‌های خبری هدف، به این نتیجه دست یافته شد که یک موتور خزشگر دو بخشی طراحی شود. این دو بخش شامل هسته خزشگر و توسعه این هسته براساس تنظیمات مورد نیاز هر وبگاه خبری است:

الف) هسته خزشگر: هسته خزشگر حاوی مفاهیم پایه‌ای است که برای خزش هر وبگاه خبری نیاز است. بنابراین با ایجاد این هسته، اطلاعات کلیدی قابل استخراج از وبگاه‌ها مشخص می‌شود که این کار ضمن افزایش سرعت برای راه‌اندازی خزش یک وبگاه، در زمان صرفه‌جویی می‌کند.

ب) توسعه هسته خزشگر و انطباق با نیاز هر وبگاه: بخش دوم موتور خزشگر حاوی تنظیمات هر وبگاه جهت استخراج اطلاعات مورد نظر از هر صفحه خبری با در نظر گرفتن مفاهیم پایه مشخص‌شده در هسته خزشگر است. در این بخش، ابتدا باید براساس ساختاری که در هر صفحه وب جهت نمایش خبر تنظیم شده است تجزیه شود و از هر صفحه فراداده‌هایی مانند نام وبگاه خبری، تاریخ و ساعت انتشار خبر، کد خبر، عنوان خبر،

خلاصه خبر که معمولاً توسط خبرگزاری‌ها درج می‌شود، متن خبر، کلیدواژه‌ها و پیوند دسترسی به خبر استخراج شود.

ساختار سامانه‌ای که برای خزش و ذخیره‌سازی داده‌های خزش‌شده طراحی شده است به این صورت است که به سه سرور نیاز است و هرکدام از سرورها نقش مشخصی را در جمع‌آوری و ذخیره داده دارد. سرور شماره (۱) فقط برای خزش روزانه اخبار از وبگاه‌های خبری هدف تخصیص داده شده است. سرور شماره (۲) فقط برای خزش بایگانی اخبار از وبگاه‌های خبری هدف تخصیص داده شده است. وجود تعداد زیاد صفحات وب بایگانی سبب شد برای کاهش زمان در فرایند خزش از بایگانی وبگاه‌ها، از شیوه خزش هم‌زمان و موازی استفاده شود. در سرور شماره (۳) پایگاه داده به نام «پایگاه داده متنی» وجود دارد که ساختار جداول آن براساس اطلاعات خزش‌شده طراحی شده است. از آنجا که حجم اخبار هر وبگاه نسبتاً زیاد بوده و خزش چندین وبگاه خبری موجب حجیم‌شدن و زمان‌برشدن استخراج داده و فراداده‌های مرتبط شد، داده‌های خزش‌شده از هر وبگاه در یک پایگاه مستقل سازماندهی شده است. بنابراین، در این پژوهش ۲۴ پایگاه داده تهیه شده است. استخراج اطلاعات پایگاه داده مربوط به هر وبگاه خبری و تجمیع آنها می‌تواند به یک پیکره بزرگ خبری منجر شود. برای کاربردی‌شدن این پیکره بزرگ که از منابع مختلف گردآوری شده است در پژوهش‌های مبتنی بر آمار و پردازش رایانشی نیاز است یکپارچگی نسبی در پیکره ایجاد گردد. از این رو، با در نظر داشتن اشکالات متداول در تهیه پیکره برای زبان فارسی (قیومی<sup>۱</sup> و همکاران، ۲۰۱۰) می‌بایست کار پیش‌پردازش انجام پذیرد.

در فرایند خزش، ۲۴ وبگاه خبری را انتخاب کردیم. ویژگی این وبگاه‌ها دسترسی به «هم‌نشری بیش‌ساده»<sup>۲</sup> در وبگاه است. ویژگی هم‌نشری بیش‌ساده این است که اجازه دسترسی به روزرسانی‌های وبگاه‌ها را در ساختار معیاری که قابلیت خوانش توسط رایانه دارد می‌دهد. اطلاعات آماری موجود در «پایگاه داده متنی» از بایگانی ۲۴ وبگاه خبری خزش‌شده از تاریخ ۱۳۶۸/۴/۱ تا ۱۴۰۱/۲/۲۱ در جدول ۱ گزارش شده است. همانگونه که مشخص است، تعداد کل اسناد خزش‌شده، اعم از اخبار بایگانی و روزانه، بالغ بر ۴۴ میلیون خبر است که این حجم داده بیش از ۱۴ میلیارد واژه را شامل شده است. شایان ذکر است بعضی از وبگاه‌های خبری، علاوه بر ارائه محتوای خبری به زبان فارسی، اخبار گزینش‌شده را به زبان‌های غیرفارسی نیز ارائه می‌کنند. حجم داده گزارش‌شده در جدول ۱ متعلق به بخش فارسی وبگاه‌های خبری منتخب است.

جدول ۱- اطلاعات آماری استخراج‌شده از بایگانی وبگاه‌های منتخب خزش‌شده

تاریخ شروع	تعداد زبان‌ها	تعداد واژه‌ها	تعداد اسناد خبری	وبگاه خبری
۱۳۶۸/۴/۱	فارسی، انگلیسی، عربی، ترکی، اسپانیایی، اردو، روسی، آلمانی، فرانسه، چینی	۲۳۵۴۵۴۴۲۵۰	۷۴۸۹۸۱۳۷	ایرنا
۱۳۸۲/۴/۱۱	فارسی، انگلیسی، عربی، اردو، ترکی، کردی	۱۹۵۲۷۶۲۰۹۳	۵۸۰۷۷۰۴	مهرنیوز
۱۳۷۸/۱/۲۲	فارسی، انگلیسی، عربی، فرانسه	۲۳۸۳۲۸۱۴۳۸	۶۲۳۶۱۶۳	ایسنا
۱۳۹۰/۲/۴	فارسی، انگلیسی، عربی	۹۴۳۵۸۷۲۶۱	۶۲۹۶۱۱۵	خبرگزاری جوان
۱۳۹۱/۸/۲۰	فارسی، انگلیسی، عربی	۱۲۱۲۱۵۸۷۵۷	۲۸۹۱۰۳۶	تسنیم نیوز
۱۳۹۰/۱۲/۹	فارسی، انگلیسی، عربی	۵۷۳۵۱۷۷۳۳	۲۲۴۷۵۷۳	خبرگزاری صداوسیما
۱۳۸۹/۱/۲۲	فارسی	۴۹۳۸۱۳۰۱۶	۱۲۰۰۱۴۷	مشرق نیوز
۱۳۸۳/۱/۱۱	فارسی، عربی	۴۶۴۷۶۲۵۶۸	۱۰۷۸۷۸۶	تابناک
۱۳۸۷/۵/۲۵	فارسی، انگلیسی، عربی	۳۹۰۱۶۷۹۷۹	۱۰۱۷۰۲۲	خبرآنلاین
۱۳۹۴/۱۰/۲۴	فارسی، انگلیسی، عربی، ترکی، اسپانیایی، اردو، فرانسه	۱۰۳۰۰۱۲۲۱	۹۴۴۹۳۵	شفقنا
۱۳۸۵/۱/۱	فارسی	۳۶۰۶۹۶۰۹۱	۸۵۲۰۸۳	همشهری آنلاین
۱۳۸۵/۴/۱۳	فارسی، عربی	۴۶۶۱۳۴۷۴۱	۱۱۹۰۴۴۶	عصرایران
۱۳۸۱/۴/۱۱	فارسی، انگلیسی، عربی، ترکی	۱۴۰۰۳۳۵۴۴۵	۳۹۹۷۱۸۳	فارس نیوز
۱۳۹۱/۶/۵	فارسی	۵۶۹۱۵۹۲۳	۶۰۰۶۱۵	نامه نیوز
۱۳۸۷/۱۱/۲۶	فارسی	۱۶۵۹۱۰۳۳۱	۵۲۹۶۱۷	شفاف
۱۳۹۱/۶/۱۴	فارسی	۱۹۶۱۴۱۳۵۳	۴۶۱۴۷۴	روزنو
۱۳۹۰/۲/۲۷	فارسی	۱۴۵۳۳۳۲۰۹	۴۵۶۴۰۸	ایران اکونومیست
۱۳۸۷/۲/۷	فارسی	۱۳۸۹۴۵۵۳۰	۲۵۲۷۰۱	رجانیوز
۱۳۹۲/۲/۱۷	فارسی	۱۱۰۴۲۸۵۶۵	۲۶۳۹۸۴	صبحانه آنلاین
۱۳۹۰/۱۰/۲۰	فارسی	۷۵۵۸۱۴۲۴	۴۴۲۹۷۱	شبکه خبر
۱۳۹۰/۱۰/۱۴	فارسی	۶۴۵۲۲۳۲۲	۱۷۶۹۸۰	اخبار بانک
۱۳۹۲/۲/۲۳	فارسی	۲۸۴۸۶۹۳۳	۱۰۷۱۶۱	تیترنیوز
۱۳۹۳/۸/۲۶	فارسی، انگلیسی	۱۲۳۸۴۵۷۷	۱۸۵۴۷	وزارت بهداشت
۱۳۹۲/۳/۴	فارسی	۳۱۵۹۵۷۱	۹۹۶۹	سازمان مدیریت بحران کشور
از ۱۳۶۸/۴/۱ تا ۱۴۰۱/۲/۲۱	۱۱ زبان	۱۴۰۹۶۶۳۹۳۳۱	۴۴۵۶۹۴۵۷	جمع کل

در جدول ۲ اطلاعات آماری استخراج شده از تاریخ ۱۳۶۸/۴/۱ تا ۱۴۰۱/۲/۲۱ براساس تفکیک زبانی از پایگاه داده متنی گزارش شده است. شایان ذکر است پژوهش حاضر محدود به بررسی خبرهای فارسی شده است.

جدول ۲- اطلاعات آماری از پیکره بزرگ متنی براساس تفکیک زبانی

زبان	تعداد وبگاه	تعداد اسناد خبری	تعداد واژه‌ها
فارسی	۲۴	۴۴۵۶۹۴۵۷	۱۴۰۹۶۶۳۹۳۳۱
عربی	۱۱	۷۰۱۴۰۷	۱۴۱۳۵۳۷۵۹
انگلیسی	۱۰	۸۱۳۰۳۱	۲۳۲۲۱۵۶۷۴
ترکی	۴	۸۱۶۹۱	۱۲۴۸۹۶۶۸
اردو	۳	۸۹۸۱۴	۲۱۷۰۵۶۹۶
فرانسه	۳	۳۸۱۱۳	۸۲۷۸۶۱۹
اسپانیایی	۲	۱۲۲۶	۳۳۱۵۸۶
کردی	۱	۶۰۲۹۰	۹۳۱۶۷۱۱
روسی	۱	۲۹۰۵۵	۳۸۱۷۹۸۰۶
آلمانی	۱	۱۱۸۱۴	۱۴۴۶۰۳۸
چینی	۱	۲۴۷۲	-

## ۵ تحلیل داده

در بخش ۳، گفتمان خبری و ساختار خبر معرفی شد. توضیح داده شد که خبر متشکل از عنوان و محتوا است و براساس هرم وارونه محتوا به سه قسمت سرنخ، بدنه و توسعه سرنخ تقسیم می‌شود. در این پژوهش می‌خواهیم میزان همبستگی معنایی میان عنوان و سه بخش محتوای خبری را مشخص کنیم تا از نظر آماری بینیم ارتباط عنوان با کدام بخش از محتوای خبر ارتباط بیشتر دارد. از این رو، تلاش می‌کنیم با کمک ابزار آماری ضریب همبستگی پیرسون به پرسش‌های این پژوهش پاسخ دهیم.

همبستگی یک شیوه آماری است که برای مقایسه دو متغیر کمی پیوسته استفاده می‌شود. به عبارتی دیگر، در همبستگی، شدت پیوستگی دو متغیر سنجیده می‌شود. ضریب همبستگی عددی بین ۱ تا -۱ است. اگر ضریب همبستگی پیرسون عدد ۱ باشد، بیانگر رابطه مستقیم بین دو متغیر است به این مفهوم که هرگونه تغییر در یک متغیر، متغیر دیگر نیز تغییر می‌کند. چنانچه این ضریب عدد -۱ باشد، رابطه معکوس بین دو متغیر وجود دارد



که با افزایش یک متغیر، متغیر دیگر کاهش می‌یابد. اگر این ضریب صفر باشد بیانگر عدم وجود رابطه خطی بین این دو متغیر است (بادی<sup>۱</sup> و اسمیت<sup>۲</sup>، ۲۰۰۹: ۹۴-۹۲). از تساوی (۱) برای محاسبه ضریب همبستگی پیرسون استفاده می‌شود:

(۱)

$$r = \frac{N \sum xy - (\sum x)(\sum y)}{\sqrt{[N \sum x^2 - (\sum x)^2][N \sum y^2 - (\sum y)^2]}}$$

که در این تساوی،  $N$  تعداد جفت‌های مورد نظر برای مقایسه است،  $\sum xy$  مجموع تعداد توالی جفت‌های مورد نظر است،  $\sum x$  مجموع تعداد متغیر  $x$  و  $\sum y$  مجموع تعداد متغیر  $y$  است،  $\sum x^2$  مجموع مجذور تعداد متغیر  $x$  و  $\sum y^2$  مجموع مجذور تعداد متغیر  $y$  است. از آنجاکه در این پژوهش به دنبال وجود رابطه همبستگی بین عنوان و محتوای خبری هستیم، برای مقایسه‌پذیری این دو به عنوان دو متغیر در یک متن خبری، نیاز است ساختار متنی داده از شکل اولیه خود خارج و در چارچوب معناشناسی توزیعی که براساس آرای هریس (۱۹۵۴) و فرث (۱۹۵۷) شکل گرفته است در قالب بردار بازنمایی شود. در این پژوهش تلاش کرده‌ایم با استفاده از ابزار ورد ۲ و ک<sup>۳</sup> که توسط میکولو و همکاران (۲۰۱۳) تهیه شده است برای بردارسازی داده‌های خزش شده از وب استفاده کنیم.

یکی از متغیرهای این پژوهش عدم توازن بین عنوان خبر و متن خبر است؛ چراکه عنوان خبر به طور متوسط حدود ۱۰ واژه بوده و محتوای خبر از این تعداد بیشتر است. در محاسبه همبستگی میان دو متغیر، تعداد دویه‌دوی متغیرها یکسان است. بنابراین، بردارسازی داده در قالب معناشناسی توزیعی سبب یکسان شدن این متغیر شده و مقایسه‌پذیری داده را ساده می‌سازد؛ به این صورت که می‌توان براساس تعداد مشخص بُعد در بردارها، به مقایسه کنترل شده داده پرداخت.

از جمله دیگر متغیرهایی که در فرایند بازنمایی معنایی واژه‌ها باید تثبیت شود تعداد بُعد بردارها و همچنین تعیین تعداد واژه‌های همسایه در بافت جایگاهی برای محدود کردن بافت و تعیین معنای آن واژه در آن بافت است. برای این هدف، تعداد ۱۰۰ بُعد و ۸ واژه در بافت جایگاهی (۴ واژه قبل و ۴ واژه بعد از واژه هدف) را در فرایند بردارسازی هر واژه مورد استفاده قرار می‌دهیم. حجم پیکره حاصل از تجمیع ۲۴ وبگاه خبری که بیش از

۱۴ میلیارد واژه است توسط یک سرور که دارای دو پردازنده ۱۲ هسته‌ای و ۲۵۶ گیگابایت حافظه رم است پردازش شده و با استفاده از وردوک بردار واژگان پیکره به دست آمده است. پس از بازنمایی معنایی واژه‌ها در قالب بردار، عنوان و محتوای هر خبر از صورت واژه به بردارهای ۱۰۰ بُعدی تبدیل شده است؛ به این صورت که بردار واژه‌های یک متن با یکدیگر جمع شده و سپس براساس تعداد واژگان متن، بردار میانگین واژه‌ها به دست آمده است. بردار میانگین سبب حذف اثرگذاری متغیر تعداد واژه‌های متن بر بردار می‌شود.

از آنجاکه در این پژوهش قصد بررسی رابطه همبستگی معنایی میان عنوان و بخش‌های سرنخ، بدنه و توسعه سرنخ را داریم، هر متن خبری خزش شده را به سه بخش تقطیع کرده‌ایم؛ به این صورت که جمله اول متن خبر به عنوان سرنخ خبر، جمله آخر به عنوان توسعه سرنخ و جملات بینابین سرنخ و توسعه سرنخ به عنوان بدنه خبر تلقی شده است. بر اساس این تقسیم‌بندی، متن‌هایی که حداقل حاوی سه جمله است را پالایش کرده و در این پژوهش استفاده کرده‌ایم. فرایند بازنمایی برداری را علاوه بر عنوان، به‌طور مجزا برای این بخش‌ها تکرار می‌کنیم.

در جدول ۳ رابطه همبستگی پیرسون میانگین بین عنوان و سرنخ، عنوان و بدنه و همچنین عنوان و توسعه سرنخ برای هر روزنامه گزارش شده است. براساس امتیاز ضریب همبستگی پیرسون میانگین در وبگاه‌های خبری، همبستگی عنوان-سرنخ و همچنین عنوان-بدنه تقریباً یکسان است و براساس آزمون تی<sup>۱</sup> دو دنباله‌ای، تفاوت معناداری بین این دو بخش خبری وجود ندارد. ولی همبستگی عنوان-توسعه سرنخ نسبت به دو بخش دیگر خبری با تفاوت معناداری کمتر است ( $p < 0/01$ ). نتایج به دست آمده حاکی از آن است که اگرچه عنوان خبر باید بیشترین ارتباط معنایی را با محتوای خبر داشته باشد و براساس هرم وارونه این نقش به بخش سرنخ خبر واگذار شده است، در بررسی وبگاه‌های خبری می‌بینیم که همبستگی عنوان-سرنخ و همچنین عنوان-بدنه حدود 0/58 است. این عدد بیانگر این است که سرنخ و بدنه خبر هر دو به یک اندازه با عنوان خبر همبستگی دارد که با ساختار هرم وارونه کمی متفاوت است؛ چراکه بدنه باید ارتباط معنایی کمتری را با عنوان داشته باشد. شایان ذکر است امتیاز به دست آمده ضریب همبستگی خیلی بالا نیست؛ چراکه با ۱ فاصله معناداری دارد. به نظر می‌رسد برای رفع این مشکل و انتقال اطلاعات بیشتر همسو بین عنوان و سرنخ باید دقت بیشتری توسط خبرنگاران صورت پذیرد.

جدول ۳- رابطه همبستگی میانگین پیرسون بین عنوان-سرnx، عنوان-بدنه و عنوان-توسعه سرnx

همبستگی پیرسون			وبگاه خبری
عنوان-توسعه سرnx	عنوان-بدنه	عنوان-سرnx	
0/131376	0/566799	0/618047	ایرنا
0/498211	0/608299	0/56198	مهرنیوز
0/180871	0/62791	0/681766	ایسنا
0/0597137	0/234171	0/227679	خبرگزاری جوان
0/17796	0/627867	0/537342	تسنیم نیوز
0/350428	0/48171	0/54448	خبرگزاری صداوسیما
0/453529	0/573402	0/557071	مشرق نیوز
0/460634	0/580151	0/602778	تابناک
0/158178	0/603539	0/585131	خبرآنلاین
0/408475	<b>0/64996</b>	<b>0/704899</b>	شفقنا
0/467625	0/568211	0/536267	همشهری آنلاین
0/44618	0/583441	0/630682	عصرایران
0/144334	0/624237	0/572369	فارس نیوز
0/498552	0/612639	0/585108	نامه نیوز
0/420367	0/569097	0/569855	شفاف
0/476263	0/613549	0/623696	روز نو
0/363095	0/6193358	0/594269	ایران اکونومیست
0/441498	0/561825	0/593515	رجانیوز
0/317141	0/578347	0/568287	صبحانه آنلاین
0/394299	0/563152	0/551843	شبکه خبر
0/490011	0/600346	0/596059	اخبار بانک
<b>0/501147</b>	0/604544	0/580048	تیترنیوز
0/24137	0/589484	0/637157	وزارت بهداشت
0/452316	0/550338	0/561328	سازمان مدیریت بحران کشور
<b>0/355566</b>	<b>0/574682</b>	<b>0/575902</b>	میانگین وبگاه‌های خبری

از میان ۲۴ وبگاه خبری و براساس محاسبات انجام‌شده در جدول ۳، ۲۳ بگه خبری (۹۶٪ وبگاه‌ها) ضریب همبستگی نسبتاً بالایی را در مقایسه بین عنوان و سرnx به‌دست آورده‌است. از میان این وبگاه‌های خبری، وبگاه «شفقنا» بالاترین و وبگاه «خبرگزاری جوان» پایین‌ترین ضریب همبستگی پیرسون را در عنوان-سرnx و همچنین عنوان-بدنه به‌دست

آورده‌است. وبگاه‌های خبری «ایسنا»، «وزارت بهداشت» و «ایرنا»، به ترتیب، جایگاه‌های دوم تا چهارم در مقایسه عنوان-سرنخ به‌دست آورده‌است.

۲۲ بگه خبری (۹۲٪ وبگاه‌ها) ضریب همبستگی نسبتاً بالایی را در مقایسه بین عنوان و بدنه به‌دست آورده‌است و مابقی کمتر بوده‌است. این نکته بیان می‌کند در این وبگاه‌ها مفهوم عنوان در بدنه خبر متبلور است. چنین می‌توان نتیجه گرفت که رابطه عنوان-سرنخ و عنوان-بدنه که بخش‌های اصلی خبر در هرم وارونه را شامل می‌شود در اکثر خبرگزاری‌ها رعایت می‌شود.

از میان ۲۴ وبگاه خبری، در ۱۳ وبگاه (۵۴٪ وبگاه‌ها) همبستگی عنوان-بدنه از عنوان-سرنخ بیشتر بود. این نتیجه بیان می‌کند که در این وبگاه‌های خبری، بدنه خبر از اهمیت بیشتری نسبت به سرنخ خبر برخوردار است که این یافته با هرم وارونه در تضاد است. شایان ذکر است در مورد همبستگی عنوان-توسعه سرنخ به این نتیجه رسیدیم که در همه وبگاه‌های خبری این همبستگی در مقایسه با عنوان-سرنخ و عنوان-بدنه کاهش یافته‌است که منطبق با هرم وارونه خبر است.

در جدول ۴، فراوانی نسبی تعداد خبرهایی که براساس ضریب همبستگی پیرسون، امتیاز بالاتر یا پایین‌تر از امتیاز ۰/۵ را به‌دست آورده، یا دارای امتیاز ضریب همبستگی صفر یا ۱- است گزارش شده‌است. همان‌طور که مشخص است، با توجه به خبرگزاری، نتایج متفاوتی از نظر تعداد اسناد خبری برای تقسیم‌بندی چهارگانه مبتنی بر ضریب همبستگی پیرسون به‌دست آمده‌است. از آنجاکه در پژوهش حاضر از داده‌های خبری استفاده می‌شود، ضریب همبستگی صفر یا ۱- به یک صورت تعبیر می‌شود و به مفهوم عدم وجود رابطه همبستگی است. به‌طور کلی، ۷۰/۶۱٪ از تمامی مستندهای خبری رابطه همبستگی نسبتاً بالا ( $r \geq 0/5$ ) بین عنوان و سرنخ وجود دارد؛ ۷۲/۵۹٪ از خبرها رابطه همبستگی نسبتاً بالا بین عنوان و بدنه وجود دارد؛ و ۴۲/۵۳٪ از خبرها رابطه همبستگی نسبتاً بالا بین عنوان و توسعه سرنخ وجود دارد. همچنین، ۲۶/۳۰٪ از متون خبری رابطه همبستگی نسبتاً پایینی ( $r < 0/5$ ) بین عنوان و سرنخ وجود دارد؛ ۲۳/۹۱٪ از خبرها همبستگی نسبتاً پایینی بین عنوان-بدنه وجود دارد؛ و همبستگی نسبتاً پایینی بین عنوان و توسعه سرنخ در ۴۴/۳۰٪ از خبرها اتفاق افتاده‌است.

براساس نتایج جدول ۴، بعضی از خبرهای وبگاه‌های خبری همبستگی بین عنوان با سرنخ، بدنه و توسعه سرنخ وجود ندارد. به‌طور متوسط، در ۱/۶۱٪ و ۱/۶۰٪ از خبرهای وبگاه‌ها، به ترتیب، هیچ رابطه همبستگی در عنوان-سرنخ و عنوان-بدنه ندارد و در ۱۲/۴۷٪ از خبرها، هیچ رابطه همبستگی در عنوان-توسعه سرنخ خبرها ندارد.

بنا بر نتایج گزارش شده در جدول ۴، در اکثر وبگاه‌های خبری، کمتر از یک درصد از اخبار، همبستگی در عنوان-سرنخ و همچنین عنوان-بدنه وجود ندارد؛ با این وجود، در بیش از ۳۰٪ از خبرهای وبگاه «خبرگزاری جوان»، همبستگی در عنوان-سرنخ و عنوان-بدنه وجود ندارد که بیانگر عدم دقت در انتخاب عنوان یا عدم رعایت مدل وارونه خبر به‌هنگام تنظیم خبر است. عدم همبستگی در عنوان-توسعه سرنخ، در بیش از ۳۲٪ از خبرهای این خبرگزاری وجود دارد که بالاترین تعداد خبر در مقایسه با سایر خبرگزاری‌ها است. در وبگاه‌های «شفقنا» و «ایسنا»، همبستگی عنوان-سرنخ و عنوان-بدنه در بیش از ۸۵٪ از خبرهای منتشرشده وجود دارد. اگرچه همبستگی عنوان-سرنخ در خبرهای وبگاه «ایران اکنومیست» در ۷۴/۵۶٪ از خبرها مشاهده شده است، در بیش از ۸۰٪ از خبرهای این خبرگزاری، همبستگی در عنوان-بدنه وجود دارد. نکته قابل توجه در بررسی همبستگی عنوان-توسعه سرنخ این است که در بعضی از وبگاه‌ها، این رابطه در حجم بسیار کمی از خبرها، کمتر از ۱۰٪، اتفاق افتاده است؛ مانند «ایرنا»، «خبرگزاری جوان»، «تسنیم‌نیوز» و «فارس‌نیوز». ولی در خبرگزاری‌هایی مانند «مهرنیوز»، «نامه‌نیوز»، «روز نو»، «اخبار بانک» و «تیترنیوز» همبستگی عنوان-توسعه سرنخ در بیش از ۵۰٪ از خبرها وجود دارد.

جدول ۴- فراوانی نسبی خبرهای دارای ضریب همبستگی پیرسون بالاتر و پایین‌تر از ۰/۵، صفر و منفی یک

ضریب همبستگی	عنوان-سرنخ (%)			عنوان-بدنه (%)			عنوان-توسعه سرنخ (%)		
	$r \leq 0$	$r < 0/5$	$r \geq 0/5$	$r \leq 0$	$r < 0/5$	$r \geq 0/5$	$r \leq 0$	$r < 0/5$	$r \geq 0/5$
ایرنا	2/30	18/60	72/41	0/90	26/33	70/08	11/74	53/71	1/40
مهرنیوز	0/23	30/10	68/54	0/10	20/93	78/43	0/80	41/78	54/46
ایسنا	0/27	13/60	85/21	0/10	17/38	81/88	8/93	39/71	13/87
خبرگزاری جوان	30/65	10/17	28/26	30/65	8/56	29/95	32/09	23/68	1/51
تسنیم نیوز	0/19	34/62	64/01	0/06	17/98	81/58	3/92	63/37	6/25

6/08	1/93	8/46	1/78	2/23	1/20	3/31	1/75	1/19	1/44	16/13	1/71	2/23	5/09
45/00	49/33	33/37	41/57	48/46	39/08	64/84	47/31	46/65	51/52	38/00	45/63	44/94	40/08
29/49	43/23	40/18	51/51	41/29	55/69	1/37	44/93	48/12	38/66	10/88	47/67	47/07	35/32
0/12	0/30	0/12	0/09	0/27	0/12	0/10	0/66	0/36	0/06	0/24	0/65	0/53	2/26
26/43	30/34	18/66	20/92	28/74	21/23	17/96	23/12	28/27	14/02	21/94	24/80	27/13	33/38
72/91	68/15	80/71	78/54	69/98	78/00	81/50	74/09	70/16	85/64	77/07	73/09	71/09	56/23
0/13	0/32	0/15	0/12	0/32	0/17	0/19	0/50	0/30	0/07	0/18	0/90	0/25	0/86
29/92	25/42	24/54	20/57	29/31	26/40	28/67	17/98	35/49	10/07	26/50	20/99	31/91	28/65
69/14	73/03	74/56	78/75	69/17	72/56	70/15	79/54	62/71	89/59	72/36	75/50	66/64	66/62
صبحانه آنلاین	رجانیوز	ایران اکونومیست	روز نو	شفاف	نامه نیوز	فارس نیوز	عصرایران	همشهری آنلاین	شفقنا	خبرآنلاین	تابناک	مشرق نیوز	خبرگزاری صداوسیما

6/67	40/61	40/42	0/30	29/05	69/37	0/33	32/67	65/54	شبکه خبر
1/28	40/58	53/78	0/19	21/84	77/06	0/15	24/44	74/66	اخبار بانک
0/87	40/71	55/20	0/19	22/14	76/84	0/24	27/30	71/18	تیترنیوز
10/26	36/08	21/94	0/27	25/99	72/50	0/23	19/03	77/71	وزارت بهداشت
1/65	47/15	45/26	0/34	31/40	66/57	0/19	31/97	66/83	سازمان مدیریت بحران کشور
12/47	44/30	42/53	1/60	23/91	72/56	1/61	26/30	70/61	متوسط کل

## ۶ جمع‌بندی و نتیجه‌گیری

آنچه در این مقاله توضیح داده شد، تحلیل گفتمان خبری و ارزیابی همبستگی معنایی میان عنوان خبر و قسمت‌های تشکیل‌دهنده خبر براساس ساختار هرم وارونه با استفاده از روش آماری بود. برای رسیدن به هدف، در چارچوب علم داده ابتدا پیکره نسبتاً بزرگی در بازه زمانی ۱۳۶۸ تا ۱۴۰۱ با حجم بیش از ۱۴ میلیارد واژه به‌واسطه خزش ۲۴ وبگاه خبری به‌دست آمد و برای این پژوهش مورد استفاده قرار گرفت. از معیار آماری همبستگی پیرسون برای یافتن همبستگی معنایی بین عنوان و بخش‌های خبر استفاده شد. برای یافتن همبستگی معنایی عنوان و بخش‌های خبر، در چارچوب معناشناسی توزیعی، بازنمایی معنایی عنوان و بخش‌های خبر با استفاده از ابزار ورد۲وک به بردار تبدیل شد و کار تحلیل انجام شد.

براساس نتایج به‌دست‌آمده، در ۹۶٪ از خبرهای وبگاه‌ها، همبستگی معنایی عنوان و سرنخ بیش از ۰/۵ بود. همبستگی نسبتاً بالای معنایی بین عنوان و بدنه خبر در ۹۲٪ از وبگاه‌ها اتفاق افتاده بود؛ و در ۸٪ از وبگاه‌ها، همبستگی نسبتاً بالای معنایی ( $r \geq 0/5$ ) بین عنوان و توسعه سرنخ خبر وجود داشت. این نتیجه بیانگر این نکته است که اکثر خبرگزاری‌ها ساختار وارونه خبر را رعایت می‌کند؛ ولی می‌توان نمونه‌هایی از وبگاه‌های خبری یافت که این ساختار را رعایت نمی‌کند. همچنین، علاوه بر وجود همبستگی معنایی

بین عنوان و سرنخ، تقریباً به یک میزان این همبستگی میان عنوان و بدنه نیز وجود دارد. دستاورد کاربردی این پژوهش این است که می‌توان با کمک روش‌شناسی معرفی‌شده، یک ارزیابی اولیه بر ساختار محتوایی انتشار خبر نمود و آن را با ساختار معیار هرم وارونه خبر مقایسه کرد.

## منابع

- آقاگلزاده، فردوس (۱۳۹۴). *تحلیل گفتمان انتقادی*. تهران: انتشارات علمی و فرهنگی.
- اردکانی‌فرد، زهرا، محمدمهدی فرقانی و مناسلگی (۱۴۰۰). «روزنامه‌نگاری دین در ایران: تحلیل محتوای خبرگزاری‌های رسمی در سال ۱۳۹۹». *فصلنامه مطالعات فرهنگ - ارتباطات*. س ۲۳، ش ۵۹، ۱۵۰-۱۲۳.
- اقبال، پرویز (۱۳۹۰). *تبیین رابطه تصویر و متن در تصویرسازی کتاب‌های داستانی کودک در ایران از سال ۱۳۴۰ تا ۱۳۸۰*. رساله دکتری. تهران: دانشکده هنر دانشگاه شاهد.
- عظیمی‌فرد، فاطمه، سیاوش صلواتیان و علی‌رضا عمادالدین (۱۳۹۶). «مقایسه کاربرد عکس خبری در وبگاه شبکه‌های خبر، العالم و پرس تی‌وی». *فصلنامه مطالعات رسانه‌های نوین*. س ۳، ش ۱۱، ۱۳۶-۹۵.
- قیومی، مسعود (۱۳۹۹). «چالش سامانه‌های پردازش طبیعی در مواجهه با نوآژه‌های «کرونایی»». *همایش ملی ابعاد انسانی-اجتماعی کرونا در ایران*. تهران: پژوهشگاه علوم انسانی و مطالعات فرهنگی.

- Bednarek, M., & H. Caple (2012). *News Discourse*. London: Continuum International Publishing Group.
- Bell, A. (1991). *The Language of News Media*. Oxford: Blackwell.
- Blei, D. M., et al. (2003). "Latent Dirichlet allocation". *Journal of Machine Learning Research*. 3: 993-1022.
- Boddy, R., & G. Smith (2009). *Statistical Methods in Practice: For scientists and Technologists*. Chichester, U.K.: Wiley.
- Brown, G., & G. Yule (1983). *Discourse Analysis*. Cambridge: Cambridge University Press.
- Cao, L. (2017). "Data science: A comprehensive overview". *ACM Computing Surveys*. 50.
- De Beaugrande, R., & W. U. Dressler (1981). *Introduction to Text Linguistics*. London: Longman.
- De Saussure, F. (1916). *Cours de linguistique générale*. Lausanne, Paris: Payot.
- Canavilhas, J. (2007). "Web journalism: From the inverted pyramid to the tumbled pyramid". *Biblioteca on-line de ciências da comunicação*.
- Cleveland, W. S. (2001). "Data science: an action plan for expanding the technical areas of the field of statistics". *International Statistical Review*. 69, 21-26.
- Dalkir, K. (2005). *Knowledge Management in Theory and Practice*. Amsterdam: Elsevier Science Ltd.
- Emde, K., & C. Klimmt, & D. M. Schlütz (2016). "Does storytelling help adolescents to process the news?". *Journalism Studies*. 17(5): 608-627.



- Feez, S., & R. Iedema, & P. R. R. White (2008) *Media Literacy*. Surry Hills, NSW: NSW Adult Migrant Education Service.
- Firth, J.R. (1957). "A synopsis of linguistic theory 1930-1955," *Studies in Linguistic Analysis (special volume of the Philological Society)*, chapter 1, pp: 1-32, Oxford: Blackwell.
- Ghayoomi, M., S. Momtazi, and M. Bijankhan (2010). "A study of corpus development for Persian". *International Journal on Asian Language Processing*. 20(1): 17-33.
- Harris, Z. S. (1954). "Distributional structure". *Word*. 23 (10): 146-162.
- Itule, B., & D. Anderson (2006). *News Writing and Reporting for Today's Media*. 7th Edition, McGraw-Hill.
- Mikolov, T., et al. (2013). "Distributed representations of words and phrases and their compositionality". *Advances in Neural Information Processing Systems*. C. J. C. Burges, et al. (eds.), 26: 3111-3119.
- Murtagh, F., & K. Devlin (2018). "The development of data science: Implications for education, employment, research, and the data revolution for sustainable development". *Big Data and Cognitive Computing*. 2.
- Norambuena, B. K., & M. A. Horning, & T. Mitra (2020). "Evaluating the inverted pyramid structure through automatic 5W1H extraction and summarization". *Computational Journalism Symposium*. Boston, MA, USA, March 20 – 21.
- Pöttker, H. (2003). "News and its communicative quality: the inverted pyramid—when and why did it appear?". *Journalism Studies*. 4(4): 501-511.
- Tognini-Bonelli, E. (2001). *Corpus Linguistics at Work, Studies in Corpus Linguistics*. Amsterdam, The Netherlands: John Benjamins Publishing Company.
- Xie, Q. (2018). "Critical discourse analysis of news discourse". *Theory and Practice in Language Studies*. 8 (4): 399-403.
- Zhang, H., & H. Liu (2016). "Visualizing structural "inverted pyramids in English news discourse across levels". *Text & Talk*. 36: 110-89.
- Ytreberg, E. (2010). "Moving out of the inverted pyramid: Narratives and descriptions in television news". *Journalism Studies*. 2(3): 357-371.

استناد به این مقاله: قیومی، مسعود (۱۴۰۱). ارزیابی ساختار هرم وارونه در پیکره بزرگ خبری فارسی: تحلیل گفتمان خبری براساس همبستگی میان عنوان و محتوای خبر. *زبان و زبان‌شناسی* ۱۸ (۳۵)، ۲۱-۴۵. doi: 10.30465/lsi.2023.849.۴۵-۲۱