

بیگلو، اعظم؛ داورپناه، محمدرضا (۱۳۹۵). محاسبه بار اطلاعاتی متون علمی فارسی براساس شاخص آنتروپی نظریه اطلاعات. پژوهشنامه کتابداری و اطلاع‌رسانی، ۶(۱)، ۸۸-۱۰۸.



محاسبه بار اطلاعاتی متون علمی فارسی براساس شاخص آنتروپی نظریه اطلاعات^۱

اعظم بیگلو^۲، دکتر محمدرضا داورپناه^۳

تاریخ پذیرش: ۹۲/۶/۳۰ تاریخ دریافت: ۹۲/۶/۲

چکیده

هدف: هدف عملده این پژوهش، محاسبه میزان بار اطلاعاتی واژه‌های متون علمی فارسی و بررسی رابطه برخی ویژگی‌های واژه و بار اطلاعاتی آن بر مبنای مقیاس آنتروپی شانون است.

روش: پژوهش حاضر با روش تحلیل محتوا و در جامعه آماری شامل ۷۵۲ مقاله برگرفته از فهرست مجلات علمی پژوهشی در سال ۱۳۸۸ صورت پذیرفت. نمونه پژوهش شامل ۳۲۰ مقاله بود که با توجه به گستردگی آن در هر حوزه تنها ۱۰ درصد از مقالات به صورت تصادفی انتخاب و مورد بررسی قرار گرفت.

یافته‌ها: پژوهش حاضر نشان داد، بار اطلاعاتی واژه با احتمال رخداد آن رابطه‌ای معکوس دارد. با افزایش تعداد حالات ممکن میزان پیش‌بینی‌پذیری و آنتروپی واژه افزایش یافته و اطلاعات کمتری منتقل می‌نماید. طول واژه نیز رابطه‌ای مستقیم با بار اطلاعاتی آن دارد. حوزه‌های مختلف علمی در میزان اطلاعاتی که انتقال می‌دهند یکسان نیستند و حوزه علوم انسانی بیشترین میزان آنتروپی و کمترین میزان اطلاعات را نسبت به سایر حوزه‌ها داراست.

کلیدواژه‌ها: نظریه اطلاعات، آنتروپی، بار اطلاعاتی واژه، متون علمی فارسی.

۱. برگرفته از پایان نامه کارشناسی ارشد

۲. کارشناس ارشد علم اطلاعات و دانش‌شناسی دانشگاه فردوسی مشهد، beiglooazam@gmail.com

۳. استاد گروه علم اطلاعات و دانش‌شناسی دانشگاه فردوسی مشهد، davarpanah@um.ac.ir

مقدمه

پیشرفت روزافرون نظام‌های ذخیره و بازیابی اطلاعات و نیاز به بهینه نمودن بازیابی اطلاعات از مدارک متین باعث شده است تا در سال‌های اخیر توجه گسترده‌تری معطوف به فنون و رویکردها در زمینه نظام‌های نمایه‌سازی خود کار شود. پیامد استفاده از نظام‌های نمایه‌سازی مبتنی بر زبان طبیعی انعکاس انواع کلمات در نمایه است. از آنجا که تمام واژگان در متن ارزش و بار اطلاعاتی یکسانی ندارند استفاده از روش‌هایی که کلمات مهم را از کلمات بی‌اهمیت تشخیص دهد همیشه در این حوزه مورد توجه بوده است. گروهی از واژگان زبان طبیعی (مانند حروف تعریف، حروف ربط، حروف اضافه، و برخی از افعال) سهم معنایی یا دستور زبانی بسیار پایینی دارند. به عبارت دیگر، حشو ویژگی بارز متون زبان طبیعی است که به منظور جلوگیری از اختلال در درک پیام متن به کار می‌رود. از سوی دیگر، بسامد واژه به تنها‌ی، به منظور اختصاص کلیدواژه‌های موضوعی مدارک چندان قابل اعتماد به نظر نمی‌رسد (داورپناه و بلندیان، ۱۳۸۶). علاوه بر این، اگر فرایند استخراج کلیدواژه بدون توجه به بار اطلاعاتی و وزن معنایی کلمه انجام پذیرد، علاوه بر حجم شدن پایگاه واژگان نمایه، ریزش کاذب و بازیابی منابع نامرتب نیز دور از انتظار نخواهد بود، چرا که واژگان اطلاعاتی بار اطلاعاتی یکسان ندارند. درواقع هر کلمه یا ترکیب به یک میزان اطلاع‌دهنده^۱ نیستند، یعنی بسیاری از ترکیبات و کلمات نباید در نظام‌های نمایه‌سازی به عنوان واژه‌نما انتخاب شوند.

بدیهی است که همه واژگان در زبان نوشتاری مقادیر یکسانی از اطلاعات را انتقال نمی‌دهند. لوهن^۲ (۱۹۵۷) نخستین کسی بود که بیان کرد واژگانی که با فراوانی بسیاری رخ می‌دهند سهم معنایی عمده‌ای در یک متن ندارند. علاوه بر این، این واژگان بخش بزرگی از متن، سهمی حدود ۲۰ تا ۳۰ درصد از نشانه‌ها در متن را، به خود اختصاص می‌دهند. از سوی دیگر، تقریباً نیمی از واژگان تنها یکبار در یک پیکره متنی^۳ رخ می‌دهند؛ درحالی که بیشتر واژگان حدود ده بار یا کمتر در متن ظاهر می‌شوند (Manning & Schutze, 1999). حتی پس از استخراج سیاهه بازدارنده، واژگان بسیار زیادی باقی خواهد ماند که همگی به عنوان واژه‌نما قابل توجه نیستند. بنابراین لازم است با استفاده از سایر روش‌ها، سودمندی واژگان باقی‌مانده تعیین شود. مطالعات فراوانی در زبان انگلیسی به رتبه‌بندی اطلاع‌دهی^۴ واژه پرداختند. منظور از اطلاع‌دهی واژه این است که آن واژه به چه میزان بیانگر ایده‌های کلیدی در مجموعه‌ای از مدارک است

1. Informative

2. Luhn

3. Corpus

4. Informativeness

(Tomokiyo & Hurst, 2003). به عبارت دیگر، اطلاع‌دهی واژه نشان‌دهنده درجه و اندازه‌ای است که عبارت کلیدی معرف و نمایانگر مدرک در دست بررسی است، و با میزان اطلاعاتی که به کاربر منتقل می‌نماید همبسته است.

تاکنون روش‌های مختلفی برای اندازه‌گیری بار اطلاعاتی واژه به کار گرفته شده است. از جمله این روش‌ها، استفاده از کمیت آنتروپی نظریه ریاضی اطلاعات^۱ شانون^۲ می‌باشد. این نظریه شاخه‌ای از نظریه آماری علوم ارتباطی است و شیوه‌ای کمی برای اندازه‌گیری محتوای اطلاعاتی پیام‌ها به دست می‌دهد. نظریه اطلاعات به طور عمده ناظر بر مسئله تعیین حداکثر ظرفیت یک کانال یا یک مجرأ برای انتقال پیام‌هاست و هدف اصلی شانون دست یافتن به شیوه‌ای بود که کارایی کانال ارتباطی را به حداکثر برساند (حری، ۱۳۸۱، ص ۲۳۲). نظریه اطلاعات با ارائه شاخصی به نام آنتروپی به اندازه‌گیری اطلاعات یک متغیر تصادفی می‌پردازد. این متغیر تصادفی می‌تواند واحدی از یک متن (حرف، کلمه، جمله، و ...) باشد. با استفاده از آنتروپی می‌توان میزان محتوای اطلاعاتی^۳ یک واژه را به عنوان متغیری تصادفی اندازه‌گرفت و کلمات با آنتروپی بالا را نادیده گرفت؛ به همین ترتیب می‌توان سیاهه‌ای از واژگان غیرمجاز ساخت. مطالعات نشان داده است که استفاده از شاخص آنتروپی نسبت به سایر معیارهای نحوی برای تشخیص و شناسایی واژگان کارکردنی و ساخت سیاهه واژگان غیرمجاز مفیدتر است (Melamed, 1997).

امروزه استفاده از شاخص آنتروپی جایگزین روش‌های نمایه‌سازی واژگان براساس آمارهای فراوانی واژه شده است. روش‌های مبتنی بر فراوانی سطحی هستند و مفهوم اصلی متن را معکس نمی‌کنند، در حالی که مدل‌هایی مبتنی بر آنتروپی دقت بالایی را گزارش کردند (Lin; Fragos, Maistros & Skourlas, 2005; Yu-Shu, 2005). براساس آنچه گفته شد، مسئله اساسی این پژوهش آن است که بر مبنای آنتروپی نظریه اطلاعات بار اطلاعاتی کلمات در متون علمی فارسی چگونه است؟

هدف‌ها و ضرورت پژوهش

هدف اصلی این پژوهش محاسبه میزان آنتروپی متون علمی و تخصصی زبان فارسی است. سایر اهداف پژوهش حاضر شامل محاسبه بار اطلاعاتی واژگان متن و شناسایی واژگان کم بار اطلاعاتی (واژگان غیرمجاز)، بررسی رابطه بار اطلاعاتی یک واژه با احتمال رخداد، تعداد حالات ممکن، طول واژه

1. Mathematical Information Theory

2. CludeShanon

3. Information Content

و زمینه (حوزه) آن است. با توجه به لزوم بهبود نظام‌های بازیابی اطلاعات و نمایه‌سازی مدارک و نیز برخی ویژگی‌های خاص خط و زبان فارسی و نظر به آنکه در زبان فارسی کمتر به موضوع وزن‌دهی واژگان پرداخته شده است، لزوم پژوهشی در راستای تعیین میزان اهمیت واژه، می‌تواند در حوزه بازیابی اطلاعات و نظام‌های نمایه‌سازی راهگشا باشد.

فرضیه‌های پژوهش

۱. میزان بار اطلاعاتی واژه با احتمال رخداد آن رابطه‌ای معکوس دارد.
۲. هرچه میزان آنتروپی متن بیشتر باشد میزان حضور اطلاعات در متن کمتر است.
۳. هرچه تعداد حالات ممکن یک واژه کمتر باشد بار اطلاعاتی آن واژه بیشتر است.
۴. بین طول کلمه و بار اطلاعاتی آن رابطه وجود دارد.
۵. میزان بار اطلاعاتی متون در حوزه‌های مختلف علمی متفاوت می‌باشد.

پیشینه پژوهش

محاسبه بار اطلاعاتی واژه براساس آنتروپی در پژوهش‌های بسیار و با اهداف مختلفی مورد توجه قرار گرفته است. کاربرد آنتروپی در تحلیل و مطالعه متون سابقه‌ای طولانی دارد. در ادامه به پژوهش‌هایی که میزان اطلاعات واژه را براساس آنتروپی سنجیده‌اند اشاره می‌کنیم.

یکی از هدف‌های تعیین بار اطلاعاتی واژه، شناسایی واژه‌های عمومی و اختصاصی یک متن است؛ میزان اطلاعات و بنابراین، آنتروپی واژه‌های عمومی و تخصصی در متون یکسان نیست و این موضوع در پژوهش‌های بسیاری مورد بررسی قرار گرفته است. کاربالو و چارنیاک^۱ (۱۹۹۹) از طریق محاسبه میزان آنتروپی کلمات سطح اختصاصی بودن^۲ آن‌ها را مورد اندازه‌گیری قرار دادند. نتایج نشان داد اسم‌هایی با آنتروپی و فراوانی بیشتر عمومی تر هستند و می‌توان گفت رابطه معکوسی بین فراوانی و محتوای معنایی یک کلمه وجود دارد. پژوهش دیگری با هدف شناسایی واژگان اختصاصی، پژوهش ریو و چوی^۳ (۲۰۰۴) است که با بررسی اصطلاحات اصطلاحات اصطلاحات مش^۴ دریافتند آنتروپی با دقیقی بسیار بالا (۸۲) درصد) به شناسایی اصطلاحات تخصصی می‌پردازد. نمیروفسکی و دوبراینین^۵ (۲۰۰۸) به بررسی اهمیت

1. Caraballo & Charniak

2. Specificity

3. Ryu & Choi

4. MeSH

5. Nemirovsky & Dobrynin

واژه در بافت و زمینه‌ای که در آن رخ می‌دهد، پرداختند. نتایج نشان داد واژه‌ای مهم است که در یک حوزه معنای تخصصی داشته باشد و در تعداد مدارک اندکی رخ دهد. به عبارت دیگر واژه تخصصی در متن آنتروپی پایینی دارد. درنهایت، می‌توان به پژوهش کیریو^۱ (۲۰۰۹) اشاره کرد که او نیز مشابه پژوهش‌های قبل، اختصاصی بودن، بار اطلاعاتی و محتوای اطلاعاتی واژه بر اساس آنتروپی را مورد بررسی قرار داد. او نشان داد به کارگیری مدل آنتروپی نتایج بسیار بهتری نسبت به روش‌های سنتی وزن دهی کلمات در پی دارد.

گروه دیگری از مطالعات به بررسی بار اطلاعاتی واژه پرداختند. از جمله این پژوهش‌ها پژوهش مونتمورو و زانت^۲ (۲۰۰۱) است که بیش از سایر پژوهش‌ها به مطالعه حاضر نزدیک و با آن همسو می‌باشد. پژوهشگران در مطالعه خود به تحلیل آماری کلمات در مجموعه متون ادبی انگلیسی پرداختند. نتایج نشان داد آنتروپی با افزایش تکرار کلمه افزایش می‌یابد، به عبارتی، کلمات کمیاب آنتروپی بسیار پایینی دارند. گنzel و چارنیاک^۳ (۲۰۰۲) با بررسی متون وال استریت ژورنال^۴ از سال ۱۹۸۷ تا ۱۹۸۹ به این نتیجه که «جملات در میزان اطلاعاتی که انتقال می‌دهند متفاوت هستند»، دست یافت.

تعیین واژگان بازدارنده نیز یکی از هدف‌های محاسبه بار اطلاعاتی واژه به حساب می‌آید. زو و همکارانش^۵ (۲۰۰۶) مدلی مبتنی بر نظریه اطلاعات به منظور ساخت سیاهه واژگان بازدارنده برای متون چینی ارائه کردند. از دیدگاه نظریه اطلاعات واژگان بازدارنده واژگانی هستند که اطلاعات کمی را انتقال می‌دهند. ارزش اطلاعاتی واژه w_j توسط آنتروپی اندازه‌گیری می‌شود. احتمال $P_{i,j}$ فراوانی آن در مدرک Di تقسیم بر تعداد کل واژگان در مدرک Di است. بنابراین، ارزش آنتروپی (H) برای واژه w_j به صورت زیر محاسبه گردید:

$$H(w_j) = \sum_N^{i=1} P_{i,j} \times \log\left(\frac{1}{P_{i,j}}\right)$$

پژوهش‌هایی که به طور مستقیم به اندازه‌گیری «آنتروپی واژه» پرداخته‌اند در زبان فارسی چندان پرشمار نیستند. آنتروپی حروف فارسی در مطالعه میرزایی (۱۳۸۵) مورد توجه قرار گرفت. هدف این

-
1. Kireyev
 2. Montemurro & Zanette
 3. Genzel & charniak
 4. The Wall Street Journal
 5. Zou et al.

مطالعه بررسی حشو در زبان فارسی با رویکرد نظریه اطلاعات بود. در نهایت این گونه نتیجه‌گیری شد که آنتروپی زبان فارسی روی حرف اول کلمه ۵ بیت است و زوائد آن حدود ۸۰ درصد می‌باشد. بنابراین اختیار ما برای سخن گفتن ۲۰ درصد و احاطه زبان فارسی بر گفتار ما حدود ۸۰ درصد است.

مطالعه بر نرخ آنتروپی زبان فارسی در پژوهش هاشمی و ساوجی (۱۳۸۶) نیز صورت پذیرفت. با استفاده از نتایج به دست آمده از مدل‌سازی متن فارسی و انگلیسی مشخص شد نرخ آنتروپی متن فارسی بالاتر از متن انگلیسی و در نتیجه قابلیت فشرده پذیری متن فارسی کمتر از متن انگلیسی است. در پژوهش‌های ذکر شده اشاره‌ای به بار اطلاعاتی نشده است، در حالی که، حجازی، علی و عابد^۱ (۱۹۸۷) در مطالعه خود به اندازه‌گیری بار اطلاعاتی هر حرف و ترکیباتی از چند حرف پرداخته و سپس این میزان اطلاعات را با زبان انگلیسی مقایسه نمودند. در نهایت چنین استنتاج شد که وقتی ارزش n (طول کلمه) افزایش می‌یابد، آنتروپی کاهش می‌یابد.

با نگاهی به مطالعات مربوط به آنتروپی واژه و بار اطلاعاتی آن می‌توان دریافت که هدف عمدۀ این پژوهش‌ها شناسایی واژه‌های عمومی و اختصاصی، واژگان بازدارنده (غیرمجاز) و تعیین رابطه بار اطلاعاتی واژه با ویژگی‌های آن از جمله طول واژه، فراوانی و ... است؛ البته قابل ذکر است که این اهداف جملگی، بر بازیابی اطلاعات بهینه تأکید دارند. البته بررسی بار اطلاعاتی واژه، صرف نظر از هدف آن، موضوعی بود که در پژوهش‌های انجام شده در ایران به چشم نخورد و به نظر می‌رسد تنها، آنتروپی به عنوان کلیتی در ارتباط با زبان مورد توجه قرار گرفته است، در حالی که بازیابی اطلاعات و نمایه‌سازی کارآمد، به انجام مطالعات گستره‌تر و دقیق‌تری بر بار اطلاعاتی واژه متکی است.

روش‌شناسی پژوهش

این پژوهش با روش تحلیل محتوا انجام پذیرفت. ویر^۲ (۱۹۹۰) تحلیل محتوا را روش تحقیقی قاعده‌مند به منظور تحلیل اطلاعات متنی طبق یک روش استاندارد و تعریف شده معرفی می‌نماید که به محقق اجازه می‌دهد تا براساس این اطلاعات، نتیجه‌گیری کند. رویکردهای کمی تحلیل متن با استفاده از روش‌های آماری و ریاضی صورت می‌پذیرد و از آنجا که هدف این پژوهش استفاده از شاخص آنتروپی به منظور تحلیل متون فارسی است، روش تحلیل محتوا روشی مناسب می‌باشد.

1. Hegazi, Ali & Abed
2. Weber

جامعه آماری این پژوهش مقالات مندرج در آخرین شماره منتشر شده در مجلات علمی-پژوهشی حوزه‌های ادبیات و علوم انسانی، علوم پایه، فنی و مهندسی، و کشاورزی است. این مجلات که از فهرست نشریات مورد تائید وزارت علوم، تحقیقات و فناوری در سال ۱۳۸۸ استخراج شد، شامل ۲۶۱ عنوان مجله در چهار حوزه است. تعداد مقالات منتشر شده در آخرین شماره این مجلات نیز حدود ۲۶۵۰ مقاله بود که از این تعداد مقاله تنها عناوینی مورد توجه قرار گرفت که به صورت الکترونیکی در دسترس بودند. به این صورت تعداد ۸۳ مجله و ۷۵۲ مقاله در حوزه علوم انسانی، ۲۲ مجله و ۱۷۰ مقاله در حوزه فنی و مهندسی، ۳۸ مجله و ۴۹۵ مقاله در حوزه کشاورزی و منابع طبیعی و ۹ مجله و ۹۸ مقاله در حوزه علوم پایه به عنوان جامعه آماری پژوهش مورد توجه قرار گرفت، سپس با استفاده از فرمول زیر تعداد مقالاتی که باید در هر حوزه مورد بررسی قرار گیرند محاسبه شد:

$$\frac{\text{حجم جامعه آن حوزه}}{\text{حجم کل جامعه}} = \frac{\text{حجم نمونه} \times \text{تعداد مقالات هر حوزه}}{\text{حجم کل مقالات هر حوزه}}$$

نتایج حاصل از این معادله تعیین کرد که تعداد ۱۵۴ مقاله در حوزه علوم انسانی، ۳۵ مقاله در حوزه فنی و مهندسی، ۱۰۱ مقاله در حوزه کشاورزی و منابع طبیعی، و ۲۰ مقاله در حوزه علوم پایه باید مورد بررسی قرار گیرند. با توجه به لزوم مقایسه متون در حوزه‌های مختلف و به کارگیری روش‌های آماری، حجم نمونه در حوزه علوم پایه نیز به ۳۰ عنوان افزایش یافت. شناسایی این مقالات از میان کل تعداد مقالات در هر حوزه به صورت تصادفی ساده انجام پذیرفت. قابل ذکر است در فرضیاتی که در سطح مقاله باید آزمون می‌شد تنها از ده درصد مقالات یعنی ۳۲ مقاله استفاده شد و این تعداد نیز به صورت تصادفی ساده انتخاب شدند.

گردآوری داده‌ها

به منظور تحلیل متون مورد مطالعه، ابتدا متن الکترونیکی متون به قالب Word تبدیل و در این محیط نرم‌افزاری تحلیل و پردازش واژگان امکان‌پذیر شد. پس از تهیه متون الکترونیکی، واژگان تفکیک گردید. تفکیک واژگان در دو مرحله صورت پذیرفت:

۱. اول واژگان هر متن از نظر شیوه نگارش مورد ویرایش قرار گرفت. در این مرحله براساس آئین نگارش زبان فارسی، واژگان ساده و مرکب تشخیص داده شده و متمایز شدند. معیارهای مورد

استفاده در این مرحله معیارهای مورد استفاده در پژوهش‌های پیشین (داورپناه و بلندیان، ۱۳۸۶؛ سنجی و داورپناه، ۱۳۸۸) و قواعد نگارش متون فارسی (وحیدیان کامیار، ۱۳۷۹) است.

۲. در مرحله دوم، متن مورد نظر را در محیط نرم‌افزاری Word و با استفاده از گزینه Table، به جدول تبدیل شد.

پس از تفکیک واژگان مقالات، برای هر مقاله یک جدول با چند ستون طراحی گردید. ستون اول مربوط به فراوانی واژه، ستون دوم طول واژه، ستون سوم حالات ممکن، ستون چهارم احتمال رخداد واژه و ستون پنجم آنتروپی بود. شمارش فراوانی واژگان به منظور محاسبه احتمال رخداد آن لازم است. تعداد کل واژه‌های هر مقاله و طول واژه در نوار وضعیت با استفاده از فعال نمودن قابلیت Word Count قابل مشاهده است. حالات ممکن یک واژه، شمارش یک واژه با کلیه کلمه‌های خانواده و هم‌ریشه بود. آنتروپی واژه نیز به منظور برآورد بار اطلاعاتی واژگان مورد توجه قرار گرفت.

یافته‌های پژوهش

با استفاده از داده‌های گردآوری شده فرضیات پژوهش مورد بررسی قرار گرفت که توضیحات آن ذیل فرضیات ارائه می‌شود:

فرضیه شماره ۱. میزان بار اطلاعاتی واژه با احتمال رخداد آن رابطه‌ای معکوس دارد.

هدف از طرح این فرضیه بررسی رابطه بین احتمال رخداد یک واژه در یک متن و میزان اطلاعات آن است، به این صورت که واژه‌ای با احتمال وقوع مشخص نسبت به سایر واژه‌ها، بار اطلاعاتی کمتر یا بیشتری خواهد داشت. متغیرهای مورد بررسی در این فرضیه شامل احتمال رخداد (P_i) و بار اطلاعاتی است. به منظور بررسی بار اطلاعاتی واژگان در این پژوهش از مفهوم کمیت آنتروپی استفاده گردید. از آنجا که رابطه آنتروپی و اطلاعات یک رابطه معکوس می‌باشد آنتروپی بالای هر واژه نشانگر بار اطلاعاتی اندک آن است. پس از ویرایش متون، شمارش واژه‌های هر مقاله و محاسبه فراوانی (f_i)، احتمال رخداد تک تک واژه‌ها محاسبه گردید به این صورت که فراوانی و تعداد تکرار یک واژه در یک مقاله به تعداد کل واژگان آن مقاله (N) تقسیم شد:

$$P_i = \frac{f_i}{N}$$

داده‌های مربوط به هر مقاله شامل فراوانی واژه، وارد برنامه اکسل شد و سپس با استفاده از نرم‌افزار آماری SPSS احتمال رخداد، لگاریتم احتمال رخداد و حاصل ضرب این دو کمیت محاسبه و طبق معادله

زیر آنتروپی واژه به دست آمد:

$$H = -P(i) \text{Log} P(i)$$

در مرحله آزمون فرض به منظور بررسی این رابطه از آزمون همبستگی پیرسون استفاده شد.

ضریب همبستگی نشانه وجود ارتباط بین دو متغیر است. قابل ذکر است که هرچه ارتباط دو متغیر شدید باشد مقدار ضریب به $+1$ و -1 - نزدیک تر خواهد بود و با کاهش ارتباط بین دو متغیر مقدار ضریب به صفر نزدیک می شود. همچنین اگر دو متغیر مستقل باشند مقدار ضریب همبستگی برابر صفر خواهد بود.

همان طور که در جدول شماره ۱ مشاهده می شود، مقدار P در همه مقالات کمتر از 0.05 است

بنابراین فرض آماری رد و فرض پژوهش پذیرفته می شود، به این معنی که میزان بار اطلاعاتی واژه متن با احتمال رخداد آن رابطه دارد. از طرفی مقدار ضریب همبستگی پیرسون نزدیک به 1 و علامت آن مثبت است که نشان دهنده همبستگی بالا میان دو متغیر احتمال رخداد و آنتروپی می باشد.

جدول شماره ۱. بررسی ضریب همبستگی آنتروپی واژه با احتمال رخداد

تعداد واژه	Sig. (2-tailed)	ضریب همبستگی	کد مقاله	ردیف
۹۳۹	.	.۰/۹۸۴	۳	۱
۱۱۹۲	.	.۰/۹۷۱	۴	۲
۸۷۳	.	.۰/۹۷۲	۶	۳
۸۵۵	.	.۰/۹۷۴	۷	۴
۷۰۰	.	.۰/۹۷۷	۱۰	۵
۵۰۰	.	.۰/۹۸۲	۱۷	۶
۸۵۸	.	.۰/۹۷۹	۲۰	۷
۱۵۶۶	.	.۰/۹۸۲	۲۱	۸
۵۴۹	.	.۰/۹۷۴	۲۳	۹
۶۰۹	.	.۰/۹۷۳	۲۷	۱۰
۶۸۶	.	.۰/۹۳۸	۲۸	۱۱
۱۲۲۴	.	.۰/۹۷۸	۳۰	۱۲
۴۸۲	.	.۰/۹۸	۳۳	۱۳
۵۶۴	.	.۰/۹۸	۳۵	۱۴
۱۴۲۴	.	.۰/۹۷	۴۴	۱۵
۱۶۱۵	.	.۰/۹۸	۴۸	۱۶
۷۰۲	.	.۰/۹۷۹	۴۹	۱۷
۸۵۲	.	.۰/۹۷۳	۵۲	۱۸

۶۵۴	.	۰/۹۷	۵۳	۱۹
۵۹۳	.	۰/۹۸	۵۴	۲۰
۶۸۸	.	۰/۹۶۷	۷۳	۲۱
۷۵۹	.	۰/۹۸۱	۷۸	۲۲
۱۵۴۲	.	۰/۹۶۶	۸۹	۲۳
۵۷۷	.	۰/۹۷۹	۹۲	۲۴
۵۷۴	.	۰/۹۷۹	۹۶	۲۵
۸۸۸	.	۰/۹۷۶	۱۰۰	۲۶
۱۸۳۴	.	۰/۹۶۹	۱۰۸	۲۷
۱۱۹۴	.	۰/۹۶۷	۱۱۰	۲۸
۱۲۵۶	.	۰/۹۸۴	۱۲۸	۲۹
۸۶۱	.	۰/۹۸	۱۳۳	۳۰
۱۵۷۲	.	۰/۹۶۵	۱۵۰	۳۱
۱۲۲۹	.	۰/۹۸۷	۱۵۴	۳۲

فرضیه شماره ۲. هرچه میزان آنتروپی متن بیشتر باشد میزان حضور اطلاعات در متن کمتر است.

همان‌طور که پیش از این ذکر شد، آنتروپی کمیتی قابل محاسبه برای هر واحد زبانی است.

آنتروپی متن برابر با مجموع آنتروپی واژگان آن متن است. به این ترتیب برای هر مقاله عددی مشبت نشانگر آنتروپی آن مقاله به دست آورده‌یم. به منظور بررسی میزان حضور اطلاعات در متن کلمات پربار و کم‌بار اطلاعاتی مورد توجه قرار گرفتند. شناسایی این کلمات به این صورت تحقق یافت که ابتدا واژگان هر مقاله به ترتیب آنتروپی (از بزرگ‌ترین به کوچک‌ترین) مرتب شدند، سپس با توجه به میانگین آنتروپی در هر مقاله دو طیف واژه شناسایی گردید. واژگان کم‌بار همان واژگانی بودند که به واسطه عدد آنتروپی بزرگ در ردیف‌های بالای جدول قرار گرفتند. این واژگان بزرگ‌ترین فراوانی را نیز دارا بودند. پس از آن واژگانی که با آنتروپی و فراوانی کم، در زیر میانگین قرار گرفتند، به عنوان واژگان پربار در نظر گرفته شدند. جدول شماره ۲ مقایسه آنتروپی مقاله و تعداد واژگان پربار و کم‌بار اطلاعاتی را نشان می‌دهد:

جدول شماره ۲. میزان حضور اطلاعات در متن

ردیف	کد	میانگین آنتروپی	متن	تعداد واژگان	بالای میانگین	پایین میانگین	فراوانی واژگان کمبار	نسبت فراوانی واژگان بربار	نسبت کم- بار به بربار
۱	۳	۰/۰۰۸۶	۰/۸۳۲۸	۴۹۵۱	۳۹۲۵	۱۰۲۶	۰/۷۹۲۷۶۹	۰/۲۰۷۲۳۱	۳/۸۲۵۵۳۳
۲	۴	۰/۰۰۷	۸/۳۴۴۷۹	۴۱۸۲	۳۰۴۹	۱۱۳۶	۰/۷۲۹۰۷۷	۰/۲۷۱۶۴	۲/۶۸۳۹۸۲
۳	۶	۰/۰۰۹	۷/۸۲۸۲۱۶	۴۵۷۴	۳۴۵۷	۱۱۱۷	۰/۷۵۵۷۹۴	۰/۲۴۴۲۰۶	۳/۰۹۴۹۰۳
۴	۷	۰/۰۱۲	۸/۱۱۷	۳۳۲۹	۲۳۸۴	۹۴۵	۰/۷۱۶۱۳۱	۰/۲۸۳۸۶۹	۲/۵۲۲۷۵۲
۵	۱۰	۰/۰۱۱۷	۸/۰۸۴۷	۲۰۳۹	۱۳۴۴	۶۹۵	۰/۶۵۹۱۴۷	۰/۳۴۰۸۵۳	۱/۹۳۳۸۱۶
۶	۱۷	۰/۰۱۵۳	۷/۷۶۵۴۷	۱۸۳۰	۱۲۶۵	۵۶۵	۰/۶۹۱۲۵۷	۰/۳۰۸۷۴۳	۲/۲۳۸۹۴
۷	۲۰	۰/۰۱۳۴	۷/۸۵۴۸۵	۲۴۵۳	۱۷۹۰	۶۶۳	۰/۷۲۹۷۱۹	۰/۲۷۰۲۸۱	۲/۶۹۹۸۵۳
۸	۲۱	۰/۰۰۵۵	۸/۰۵۵۶۲۱	۶۷۴۱	۵۰۳۹	۱۷۰۲	۰/۷۴۷۵۱۵	۰/۲۵۲۴۸۵	۲/۹۶۰۶۳۱
۹	۲۳	۰/۰۱۳۱	۷/۱۸۵۲۲	۲۸۵۲	۲۲۰۶	۶۴۶	۰/۷۷۳۴۹۲	۰/۲۲۶۵۰۸	۳/۴۱۴۸۵۵
۱۰	۲۷	۰/۰۱۲۹	۷/۸۲۷۰۷	۲۲۲۴	۱۵۰۶	۷۱۸	۰/۶۷۷۱۵۸	۰/۳۲۲۸۴۲	۲/۰۹۷۴۹
۱۱	۲۸	۰/۰۲۷۶	۷/۹۸۷۱	۲۲۹۰	۱۶۷۳	۶۱۷	۰/۷۳۰۵۶۸	۰/۲۶۹۴۳۲	۲/۷۱۱۵۱۲
۱۲	۳۰	۰/۰۰۷۱	۸/۶۹۴۲۳	۳۷۸۵	۲۶۱۰	۱۱۷۵	۰/۶۸۹۵۶۴	۰/۳۱۰۴۳۶	۲/۲۲۱۲۷۶
۱۳	۳۳	۰/۰۱۵۶	۷/۰۵۳۳۴۲	۱۸۵۲	۱۳۲۵	۵۲۷	۰/۷۱۵۴۴۳	۰/۲۸۴۵۵۷	۲/۵۱۴۲۳۴
۱۴	۳۵	۰/۰۱۴	۷/۹۱۲۹۸	۱۷۶۲	۱۲۶۱	۵۰۱	۰/۷۱۰۵۶۴	۰/۲۸۴۳۴۶	۲/۵۱۶۹۶۶
۱۵	۴۴	۰/۰۰۶	۸/۸۶۰۳۵	۴۷۲۶	۳۳۰۹	۱۴۱۷	۰/۷۰۰۱۶۹	۰/۲۹۹۸۳۱	۲/۳۳۵۲۱۲
۱۶	۴۸	۰/۰۰۵۳	۸/۸۶۰۰۲	۵۵۸۳	۴۰۹۳	۱۴۹۰	۰/۷۳۳۱۱۸	۰/۲۶۹۸۸۲	۲/۷۴۶۹۷۴
۱۷	۴۹	۰/۰۱۱۶	۸/۱۵۹۹۹	۲۲۵۳	۱۶۲۹	۶۲۴	۰/۷۲۳۰۳۶	۰/۲۷۶۹۶۴	۲/۶۱۰۵۷۸
۱۸	۵۲	۰/۰۰۹۵	۸/۰۹۳۴۴	۳۷۱۹	۲۸۳۲	۸۸۷	۰/۷۶۱۴۹۵	۰/۲۳۸۵۰۵	۳/۱۹۲۷۸۴
۱۹	۵۳	۰/۰۱۱۸	۸/۰۴۷۰۸	۲۷۸۲	۲۰۶۶	۷۱۶	۰/۷۴۲۶۳۱	۰/۲۵۷۳۶۹	۲/۸۰۵۴۷۲
۲۰	۵۴	۰/۰۱۳۶	۷/۵۹۸۸۲	۱۷۴۴	۱۲۱۱	۵۲۳	۰/۶۹۴۳۸۱	۰/۳۰۵۶۱۹	۲/۲۷۲۰۴۸
۲۱	۷۳	۰/۰۱۱۴	۷/۸۶۹۶۸	۲۶۵۶	۱۹۰۹	۷۴۷	۰/۷۱۸۷۵	۰/۲۸۱۲۵	۲/۵۵۵۵۵۶
۲۲	۷۸	۰/۰۱۰۲	۸/۳۱۷۱۱	۳۵۴۴	۲۷۴۱	۸۰۳	۰/۷۷۳۴۲	۰/۲۲۶۵۸	۳/۴۱۳۴۵۲
۲۳	۸۹	۰/۰۰۵۶	۸/۶۴۷۹۱	۵۰۶۱	۳۶۲۶	۱۴۳۵	۰/۷۱۶۴۵۹	۰/۲۸۳۵۴۱	۲/۵۲۶۸۲۷
۲۴	۹۲	۰/۰۱۳۸	۷/۹۵۱۷۶	۲۰۹۲	۱۴۵۷	۶۳۵	۰/۶۹۶۴۶۳	۰/۳۰۴۵۳۷	۲/۲۹۴۴۹۱
۲۵	۹۶	۰/۰۱۳۹	۷/۹۷۳۷۷	۱۶۶۰	۱۱۳۵	۵۲۵	۰/۶۸۳۷۳۵	۰/۳۱۶۲۶۵	۲/۱۶۱۹۰۵
۲۶	۱۰۰	۰/۰۰۹۵	۸/۴۲۸۴۵	۲۷۷۶	۱۹۴۳	۸۲۳	۰/۶۹۹۹۲۸	۰/۳۰۰۰۷۲	۲/۳۳۲۵۴۴

۲/۹۵۷۰۸۹	۰/۲۵۲۷۱۱	۰/۷۴۷۲۸۹	۲۰۵۱	۶۰۶۵	۸۱۱۶	۸/۷۴۰۱۵	۰/۰۰۴۸	۱۰۸	۲۷
۲/۳۲۹۸۲۶	۰/۳۰۰۳۱۶	۰/۶۹۹۶۸۴	۱۳۳۱	۳۱۰۱	۴۴۳۲	۸/۴۱۱۰۴	۰/۰۰۷۱	۱۱۰	۲۸
۱/۸۹۷۲۲۴	۰/۳۴۵۱۵۸	۰/۶۵۴۸۴۲	۱۲۲۶	۲۳۲۶	۳۵۵۲	۸/۶۵۲۱۹	۰/۰۰۶۹	۱۲۸	۲۹
۲/۴۹۳۰۱	۰/۲۸۶۲۸۶	۰/۷۱۳۷۱۴	۷۸۷	۱۹۶۲	۲۷۴۹	۸/۴۰۳۹۸	۰/۰۰۹۸	۱۳۳	۳۰
۲/۳۷۱۶۵۸	۰/۲۹۶۵۹	۰/۷۰۳۴۱	۱۴۹۶	۳۵۴۸	۵۰۴۴	۸/۸۸۷۸	۰/۰۰۵۷	۱۵۰	۳۱
۲/۲۳۵۵۹۷	۰/۳۰۹۰۶۲	۰/۶۹۰۹۳۸	۱۱۶۳	۲۶۰۰	۳۷۶۳	۸/۵۴۸۸۵	۰/۰۰۷	۱۵۴	۳۲

به منظور آزمون فرض فوق ابتدا نسبت واژگانی که آنتروپی آنها بالای میانگین بود به تعداد کل واژگان برای هر مقاله محاسبه شد، این محاسبه برای تعداد واژگانی که آنتروپی آنها پایین میانگین بود نیز حساب شد. سپس بین این دو تعداد آزمون همبستگی انجام شد.

جدول شماره ۳. اختلاف معناداری تعداد واژگان بالا و پایین میانگین

N	Sig. (2-tailed)	ضریب همبستگی
۳۲	.	-۱

نتایج آزمون نشان داد که بین تعداد واژگان بالای میانگین (نمایانگر واژگان کم بار اطلاعاتی) و واژگان پایین میانگین (نمایانگر واژگان پربار اطلاعاتی) در همه مقالات رابطه معنادار معکوسی وجود دارد. مرحله دوم آزمون فرض، بررسی رابطه این دو سطح از متغیر میزان اطلاعات متن با متغیر آنتروپی متن بود. به این منظور از آزمون رگرسیون چند متغیری^۱ استفاده شد. جدول شماره ۴ حاصل انجام این آزمون در نرم افزار اس.بی.اس.^۲ است.

جدول شماره ۴. رابطه اطلاعات متن با آنتروپی

Sig.	T	ضرایب استاندارد Beta	ضرایب غیراستاندارد		Model
			Std. Error	B	
.	۸/۷۳۷		۰/۹۴۳	۸/۲۴۲	(Constant)
.	۷/۷۵۷	۱/۲۵	۰/۴۰۶	۷/۴۵۶	بالای میانگین
۰/۶۶۱	-۰/۴۴۳	-۰/۰۹۲	۰/۰۰۱	۰	پایین میانگین

ارزش P به دست آمده در دو سطح واژگان بالا و پایین میانگین آنتروپی، نشان‌دهنده ارتباط معنی‌دار بین اطلاعات متن و آنتروپی آن در سطح واژگان بالای میانگین است. این ارتباط در سطح واژگان پایین میانگین معنادار نیست. درواقع آنتروپی متن با تعداد واژگان کم بار اطلاعاتی رابطه دارد و می‌توان گفت هرچه تعداد این واژگان بیشتر باشد آنتروپی متن نیز بالاتر است.

1. Multiple Regression

2. SPSS

فرضیه شماره ۳. هرچه تعداد حالات ممکن یک واژه کمتر باشد بار اطلاعاتی آن واژه بیشتر است. تعداد حالات ممکن یک واژه، شکل‌های مختلف ظهور یک واژه در متن اعم از حالات اسمی، فعلی، قیدی و ... است. به طور مثال حالات ممکن واژه تدارک، در یک مقاله شامل واژه‌های تدارک، تدارک دیدن، تدارکات، تدارک دید، خواهد بود. پس از شمارش حالات ممکن هر واژه یک ستون به آن اختصاص یافت و سپس فراوانی کلیه حالات ممکن یک واژه تجمعی گردید. به این ترتیب واژه‌ای با ۶ حالت ممکن تنها به یک شکل در جدول ظاهر می‌شود، در حالی که فراوانی سایر حالات به فراوانی آن افزوده گشته است. در مرحله بعد واژگان غیرمجاز از جداول حذف شد. چشم‌بوشی از واژگان غیرمجاز به این دلیل انجام شد که این واژگان بیشترین فراوانی را دارا بوده در عین حال حالات ممکن متعددی ندارند و این امر می‌توانست صحت آزمون همبستگی را تحت الشاعر قرار دهد. پس از این مرحله به بررسی رابطه این ستون با ستون احتمال رخداد پرداختیم.

جدول شماره ۵. بررسی ضریب همبستگی آنتروپی واژه با حالات ممکن آن

ردیف	کد مقاله	ضریب همبستگی	Sig. (2-tailed)	تعداد واژه
۱	۳	۰/۵۲۴	.	۹۳۹
۲	۴	۰/۴۲۴	.	۱۱۹۲
۳	۶	۰/۴۷۸	.	۸۷۳
۴	۷	۰/۴۸۶	.	۸۵۵
۵	۱۰	۰/۵۰۹	.	۷۰۰
۶	۱۷	۰/۵۰۲	.	۵۰۰
۷	۲۰	۰/۵۵۸	.	۵۸۵
۸	۲۱	۰/۵۶۲	.	۱۵۶۶
۹	۲۳	۰/۳۳۸	.	۵۴۹
۱۰	۲۷	۰/۴۷	.	۶۰۹
۱۱	۲۸	۰/۴۹۲	.	۶۸۶
۱۲	۳۰	۰/۵۷۷	.	۱۲۲۴
۱۳	۳۳	۰/۴۶۱	.	۴۸۲
۱۴	۳۵	۰/۶۰۳	.	۵۶۴
۱۵	۴۴	۱	.	۱۴۲۴
۱۶	۴۸	۰/۳۴۲	.	۱۶۱۵
۱۷	۴۹	۰/۶۶۸	.	۷۰۲

ردیف	کد مقاله	ضریب همبستگی	Sig. (2-tailed)	تعداد واژه
۱۸	۵۲	۰/۴۴	.	۸۵۲
۱۹	۵۳	۰/۳۸۹	.	۶۵۴
۲۰	۵۴	۰/۵۴۹	.	۵۹۳
۲۱	۷۳	۰/۴۶۱	.	۶۸۸
۲۲	۷۸	۰/۳۸۷	.	۷۵۹
۲۳	۸۹	۰/۵۵۵	.	۱۵۴۲
۲۴	۹۲	۰/۵۹	.	۵۷۷
۲۵	۹۶	۰/۵۸۶	.	۵۷۴
۲۶	۱۰۰	۰/۵۲۲	.	۸۸۸
۲۷	۱۰۸	۰/۴۵۵	.	۱۸۳۴
۲۸	۱۱۰	۰/۶۰۲	.	۱۱۹۴
۲۹	۱۲۸	۰/۵۱۸	.	۱۲۵۶
۳۰	۱۳۳	۰/۶۰۳	.	۸۶۱
۳۱	۱۵۰	۰/۴۹۵	.	۱۵۷۲
۳۲	۱۵۴	۰/۴۴۷	.	۱۲۲۹

با توجه به مقدار P به دست آمده از تک تک مقالات فرض پژوهش پذیرفته می شود:

$$\text{Sig.(2-tailed)} = \text{P-Value} < 0.05$$

بررسی جدول شماره ۵ نشان دهنده پذیرش فرض آزمون است. طبق یافته های این جدول ارتباط معنی داری میان تعداد حالات ممکن و آنتروپی وجود دارد. ضریب همبستگی پرسون 0.3 تا 1 می باشد و این مطلب نشان دهنده همبستگی متوسط تا قوی است.

فرضیه شماره ۴. بین طول کلمه و بار اطلاعاتی آن رابطه وجود دارد.

هدف از طرح این فرضیه بررسی این موضوع است که آیا واژگان با طول های متفاوت، بار اطلاعاتی یکسانی دارند یا نه؟ به عبارت دیگر آیا تک تک واژگان صرف نظر از متنی که در آن ظاهر می شوند به میزان یکسانی اطلاعات منتقل می نمایند و آنتروپی مشابهی دارند. این فرضیه نیز باید در سطح تک تک واژه های مقالات نمونه بررسی شد. متغیرهای مورد بررسی در این فرضیه طول کلمه و آنتروپی آن می باشد. طول کلمه همان تعداد کاراکترهای یک واژه است که پس از انجام عملیات ویرایش و آماده سازی متون به صورت جداول با استفاده از قابلیت Word Count محاسبه و در ستونی جداگانه در صفحه اکسل به نمایش درآمد. آنتروپی کلیه واژگان نیز همان گونه که قبله توضیح داده شد محاسبه

گردید.

جدول شماره ۶. رابطه آنتروبی واژه با طول آن

ردیف	کد مقاله	ضریب همبستگی	Sig. (2-tailed)	تعداد واژه
۱	۳	-۰/۲۳۹	.	۹۳۹
۲	۴	-۰/۱۹۷	.	۱۱۹۲
۳	۶	-۰/۲۱۶	.	۸۳۷
۴	۷	-۰/۲۱۱	.	۸۵۷
۵	۱۰	-۰/۲۳۳	.	۶۹۹
۶	۱۷	-۰/۲۱۱	.	۵۰۰
۷	۲۰	-۰/۱۹۷	.	۵۸۶
۸	۲۱	-۰/۲۱۳	.	۱۵۶۶
۹	۲۳	-۰/۰۹۴	.۰/۰۲۸	۵۴۹
۱۰	۲۷	-۰/۲۱۹	.	۶۰۹
۱۱	۲۸	-۰/۲۰۲	.	۶۸۶
۱۲	۳۰	-۰/۱۲۶	.	۱۲۲۳
۱۳	۳۳	-۰/۲۵۲	.	۴۸۲
۱۴	۳۵	-۰/۲۱۲	.	۵۶۳
۱۵	۴۴	-۰/۱۸۶	.	۱۴۸۱
۱۶	۴۸	-۰/۱۸۳	.	۱۶۱۵
۱۷	۴۹	-۰/۲۲۶	.	۷۰۲
۱۸	۵۲	-۰/۲۲۶	.	۸۰۲
۱۹	۵۳	-۰/۲۳۴	.	۶۴۵
۲۰	۵۴	-۰/۲۲۳	.	۵۹۲
۲۱	۷۳	-۰/۱۸۶	.	۶۸۸
۲۲	۷۸	-۰/۱۰۶	.۰/۰۰۳	۷۵۹
۲۳	۸۹	-۰/۱۸	.	۱۵۴۲
۲۴	۹۲	-۰/۲۲۲	.	۵۷۷
۲۵	۹۶	-۰/۲۱۲	.	۵۷۴
۲۶	۱۰۰	-۰/۱۸۳	.	۸۸۸
۲۷	۱۰۸	-۰/۱۸۸	.	۱۸۳۴
۲۸	۱۱۰	-۰/۱۷۹	.	۱۱۹۴

۱۲۵۶	.	-۰/۲	۱۲۸	۲۹
۸۶۱	.	۰/۱۹۷	۱۳۳	۳۰
۱۵۷۲	.	۰/۱۹۳	۱۵۰	۳۱
۱۲۲۹	.	۰/۲۱۹	۱۵۴	۳۲

با توجه به جدول شماره ۶ در تمامی مقالات مورد بررسی ارزش P (Sig.) کوچک‌تر از ۰/۰۵ است که این مطلب بیانگر وجود رابطه میان طول کلمه و بار اطلاعاتی آن می‌باشد. ضریب همبستگی به دست آمده نیز در هر ۳۲ نمونه مورد بررسی عددی منفی است که رابطه‌ای معکوس بین آنتروپی واژه و طول آن را نشان می‌دهد. به این ترتیب، هرچه واژه‌ای کوتاه‌تر باشد آنتروپی واژه بیشتر و بنابراین اطلاعات واژه کمتر خواهد بود.

فرضیه شماره ۵. مقدار بار اطلاعاتی متون در حوزه‌های مختلف علمی متفاوت می‌باشد.

هدف از مطرح نمودن فرضیه فوق، بررسی این موضوع است که آیا بافت‌های واژگانی با درون‌ماهی متفاوت، میزان اطلاعات یکسانی منتقل می‌سازند یا برخی متون با موضوعی خاص نسبت به متون دیگر با موضوعی متفاوت اطلاعات بیشتری دارا هستند؟ همان‌طور که پیش از این گفته شد کمیت آنتروپی قابل محاسبه برای کلیه واحدهای زبانی است؛ پس از محاسبه آنتروپی تک تک واژگان یک مقاله با استفاده از فرمول می‌توان آنتروپی متن را قابل محاسبه ساخت:

$$H = - \sum P(i) \log P(i)$$

علامت سیگما (Σ) در فرمول نشان‌دهنده این است که آنتروپی متن حاصل جمع آنتروپی کلیه واژگان آن متن است. به این ترتیب برای هر مقاله یک عدد مثبت نشانگر آنتروپی متن به دست می‌آید. بدیهی است که متغیر مورد بررسی در این فرضیه آنتروپی H هر مقاله است. سطح مورد آزمون نیز چهار حوزه است که میانگین آنتروپی هر حوزه در جدول شماره ۷ ارائه شده است. قابل ذکر است که سطح مورد توجه در این فرضیه حوزه است، به این ترتیب ۳۲۰ مقاله اولیه مورد بررسی قرار گرفتند.

جدول شماره ۷. میانگین آنتروپی در حوزه‌های علمی

ردیف	حوزه علمی	تعداد مقالات	میانگین آنتروپی
۱	ادیات و علوم انسانی	۱۵۴	۸/۴۵
۲	علوم پایه	۳۰	۷/۹۸۹
۳	فنی و مهندسی	۳۵	۷/۵۹۰۶۴
۴	کشاورزی	۱۰۱	۷/۹۲۳

به منظور بررسی این فرضیه از آزمون تحلیل واریانس^۱ یک طرفه استفاده گردید. از این آزمون در مواقعی استفاده می‌شود که آمارگر قصد دارد میانگین‌های سه جامعه یا بیشتر را با یکدیگر مقایسه کند یا به برآورد و مقایسه میانگین یک صفت در چند جامعه و همچنین وجود یا عدم وجود تفاوت معنی دار در بین نمونه‌های یک جامعه بپردازد (هویدا، ۱۳۷۸). با توجه به اینکه هدف از این فرضیه بررسی آنتروپی در چهار حوزه ادبیات و علوم انسانی، فنی و مهندسی، علوم پایه و کشاورزی می‌باشد بنابراین فرضیه با آزمون تحلیل واریانس بررسی شد.

جدول شماره ۸. آزمون تحلیل واریانس میانگین آنتروپی حوزه‌های علمی

Sig.	F	Df	مجموع مربعات خطا	ANOVA
.	۱۰/۸۸	۳	۳۰/۶۹۳	بین گروهی
		۳۱۶	۲۹۷/۱۴۶	درون گروهی

با توجه به جدول حاصل از آزمون ANOVA از $P\text{-Value}=0.000<0.05$ کوچک‌تر است و فرض پژوهش پذیرفته می‌شود (Sig.= P-Value=0.000<0.05)؛ همچنین نسبت F در صورتی که مساوی یا کمتر از یک باشد نتیجه معنی دار نیست و بالعکس هرچه بزرگ‌تر باشد اثر متغیر مستقل بر داده‌ها بیشتر است. البته میزان بزرگی نسبت F به حدی که بتوانیم آن را معنادار بدانیم بستگی به ارزش P دارد و باید کمتر از 0.05 باشد تا بتوان F را معنادار دانست (بریس، کمپ و سنلگار، ۱۳۹۱). با این شرح نسبت F بین گروه‌ها نشانگر معناداری اختلاف میان آنتروپی حوزه‌هاست. به منظور بررسی رابطه بین حوزه‌ها به صورت دو به دو از آزمون تعقیبی LSD استفاده شد. جدول شماره ۹ نشان‌دهنده میزان اختلاف در اطلاعات چهار حوزه علمی است.

جدول شماره ۹. اختلاف آنتروپی در حوزه‌های علمی

Sig.	انحراف استاندارد	اختلاف میانگین	حوزه	
۰/۰۱۸	۰/۱۹۳۵۲	۰/۴۶۰۷۴	علوم پایه	ادبیات و علوم انسانی
.	۰/۱۸۱۵۸	۰/۸۵۸۴۳	مهندسی	
.	۰/۱۲۴۱۶	۰/۵۲۵۷۳	کشاورزی	
۰/۰۱۸	۰/۱۹۳۲۵	۰/۴۶۰۷۴	ادبیات و علوم انسانی	علوم پایه
۰/۱	۰/۲۴۱۲۷	۰/۳۹۷۶۹	مهندسی	
۰/۷۴۷	۰/۲۰۱۶۳	۰/۰۶۴۹۹	کشاورزی	
.	۰/۱۸۱۵۸	۰/۸۵۸۴۳	ادبیات و علوم انسانی	مهندسي

1. Analysis of Variance (ANOVA)

۰/۱	۰/۲۴۱۲۷	۰/۳۹۷۶۹	علوم پایه	
۰/۰۸۱	۰/۱۹۰۲	۰/۳۳۲۷	کشاورزی	
۰	۰/۱۲۴۱۶	۰/۵۲۵۷۳	ادبیات و علوم انسانی	
۰/۷۴۷	۰/۲۰۱۶۳	۰/۰۶۴۹۹	علوم پایه	کشاورزی
۰/۰۸۱	۰/۱۹۰۲	۰/۳۳۲۷	مهندسی	

با بررسی ارزش P حاصل مقایسه میانگین حوزه‌ها پژوهش پذیرفته می‌شود. در نهایت می‌توان گفت که حوزه ادبیات و علوم انسانی با سه حوزه دیگر از نظر اطلاعاتی که منتقل می‌نماید متفاوت است.

نتیجه

بررسی آنتروپی و بار اطلاعاتی واژه در متون فارسی از چند جنبه در پژوهش حاضر مورد توجه قرار گرفت. اولین ویژگی مورد مطالعه یک واژه در متن احتمال رخداد واژه و رابطه آن با میزان اطلاعات واژه بود. شانون به صراحت در نظریه معروف خود تعریف کرد که اطلاعات I در یک پیام به طور معکوس با احتمال آن رابطه معکوس دارد (Bartlett, 2007). نتایج حاصل از پژوهش‌های مختلف نشان‌دهنده رابطه معکوس باراطلاعاتی یک واژه با احتمال رخداد آن است. مک‌کی^۱ (۲۰۰۳) محتوای اطلاعاتی ۲۷ رخداد ممکن را هنگامی که یک کاراکتر به طور تصادفی از متون انگلیسی انتخاب می‌شود، مورد محاسبه قرار داد. به این ترتیب محتوای اطلاعاتی رخداد حرف Z برابر ۱۰.۴ و حرف E برابر ۳.۵ بیت است. در این رابطه احتمال رخداد با میزان اطلاعات شanon رابطه معکوس دارد. بدیهی است هرچه احتمال وقوع یک پیشامد بیشتر باشد آن پیشامد اطلاعات کمتری را انتقال می‌دهد و این نتیجه در آزمون رابطه بین دو متغیر آنتروپی و احتمال رخداد در پژوهش حاضر نیز به دست آمد، به این معنا که واژه‌هایی که با آنتروپی بالایی در متن ظاهر شدن بیشترین احتمال رخداد را داشتند. فیریسون، راینر و پیکرینگ^۲ (۲۰۰۵) در پژوهش خود با بررسی هم وقوعی واژگان (احتمال رخداد دو واژه با هم)، معتقدند یک واژه با قابلیت پیش‌بینی‌پذیری بالا واژه‌ای است که با احتمال رخداد بالایی در متن ظاهر می‌شود. همان‌طور که گفته شد طبق اصول شanon واژه‌ای که پیش‌بینی‌پذیری زیادی دارد در حالی که آنتروپی بالایی دارد اطلاعات اندکی خواهد داشت. بنابراین طبق رابطه سه مفهوم مذکور خواهیم داشت:

1. MacKay

2. Frisson, Rayner & Pickering

پیش‌بینی پذیری زیاد ↔ احتمال رخداد بالا ↔ افزایش آنتروپی ↔ اطلاعات اندک

بنابر آنچه گفته شد میزان اطلاعات واژه با احتمال رخداد آن رابطه‌ای معکوس دارد.

در فرضیه دوم پژوهش حاضر، به میزان آنتروپی متن و رابطه آن با حضور اطلاعات در همان متن پرداختیم. همان‌گونه که در رابطه با آزمون فرض دوم پژوهش گفته شد، آنتروپی یک مجموعه حاصل مجموع آنتروپی تک‌تک اجزای آن مجموعه است. درواقع اگر آنتروپی معادل عدم وجود اطلاعات در یک سامانه باشد، بنابراین آنتروپی بالا نشان‌دهنده اطلاعات اندک آن خواهد بود. آزمون فرض بین متغیر آنتروپی متن و اطلاعات آن نشان داد که آنتروپی متن با تعداد واژگان بالای میانگین (واژگان کم‌بار) رابطه دارد؛ درحالی که با تعداد واژگان پایین میانگین (واژگان پر‌بار) این رابطه تأیید نشد. البته شناسایی واژه‌های پر‌بار و کم‌بار اطلاعاتی متن نیازمند تعریف شاخص دقیق‌تر دیگری است؛ ولی با توجه به یافته‌های توکانی گفت آنتروپی متن تنها با واژگانی که به طور قطعی واژگان کم‌بار مقاله هستند رابطه داشت.

بار اطلاعاتی واژه از جنبه تعداد حالات ممکن آن در متن نیز قابل بررسی است. هدف از بررسی رابطه تعداد حالات ممکن واژه و بار اطلاعاتی آن، بررسی این موضوع است که آیا واژه‌ای که به شکل‌های مختلف (واژه‌های هم‌خانواده و هم‌ریشه) در یک متن ظاهر می‌شود میزان اطلاعات بیشتری منتقل می‌نماید یا خیر. یو، وو، و هاوی^۱ (۲۰۰۸) در فرآیند ابهام‌زدایی معنایی از واژه به رابطه حالات مختلف یک واژه و آنتروپی اشاره کردند. به نظر آنان هرچقدر ابهام معنایی یک واژه بالا باشد و به عبارت بهتر واژه دارای معانی و اشکال متعددی باشد آنتروپی بالاتری برای آن واژه خواهیم داشت. همچنین بار اطلاعاتی آن واژه بیشتر است چرا که آن واژه غیرقابل پیش‌بینی تر است و محتوای اطلاعاتی بیشتری به همراه دارد.

بررسی رابطه طول واژه و میزان اطلاعات آن نشان داد که هرچه یک واژه کوتاه‌تر باشد میزان اطلاعات کمتری نیز خواهد داشت. به این ترتیب رابطه آنتروپی و طول واژه رابطه‌ای معکوس است. طبق نظریه شانون یک زبان کارآمد میزان اطلاعاتی نزدیک به ظرفیت کanal انتقال می‌دهد، به این ترتیب یک رابطه غیرخطی بین منفی لگاریتم احتمال واژه و طول آن وجود دارد و این رابطه همان فرمول آنتروپی شانون است. کانچو و مارتین^۲ (۲۰۱۱) نیز با بررسی متونی که به طور تصادفی و بی‌هدف تایپ می‌شوند به رابطه خطی بین طول واژه^۳ و محتوای اطلاعاتی (I) آن اشاره کردند. بنابر آنچه گفته شد نتیجه حاصل حکایت از وجود رابطه بین طول کلمه و بار اطلاعاتی آن داشت.

1. Yu, Wu & Hovy

2. Cancho & Martin

درنهایت، نتایج محاسبه بار اطلاعاتی متون حوزه‌های علمی مختلف در چهار حوزه علمی ادبیات و علوم انسانی، علوم پایه، فنی و مهندسی و کشاورزی نشان داد بیشترین آنتروپی و کمترین اطلاعات در حوزه علوم انسانی نسبت به سایر حوزه‌هایی باشد. پس از حوزه علوم انسانی، حوزه علوم پایه، کشاورزی، و فنی و مهندسی به ترتیب از بالاترین آنتروپی تا کمترین میزان آن قرار می‌گیرند. این مسئله به دلیل وجود مقالات طولانی و پر تکرار حوزه علوم انسانی است. آنچه در مقالات این حوزه به چشم می‌خورد فراوانی افراش طول مقاله احتمال رخداد واژگان غیرمجاز نیز افزایش یافته و این خود عاملی در افزایش میانگین آنتروپی حوزه است. بنابراین می‌توان گفت بار اطلاعاتی واژه وابسته به متن است و مقدار اطلاعات متون در حوزه‌های مختلف علمی متفاوت می‌باشد.

کتابنامه

بریس، نیکلا، کمپ ریچارد، و ستلگار، رزمی (۱۳۹۱). تحلیل داده‌های روانشناسی با برنامه اس پی اس اس. ترجمه خدیجه علی‌آبادی و علی صمدی. تهران: دوران

حری، عباس (۱۳۸۱)، دائره‌المعارف کتابداری و اطلاع‌رسانی. تهران: مرکز اسناد و کتابخانه ملی ایران
داورپناه، محمد رضا، و بلندیان، صدیقه (۱۳۸۶). تحلیل متن مقالات فارسی و امکان نمایه سازی ماشینی آن‌ها براساس قانون زیف. فصلنامه پژوهش در مسائل تعلیم و تربیت: ویژه نامه کتابداری و اطلاع‌رسانی، دور

دوم

سنجری، مجیده، و داورپناه، محمد رضا (۱۳۸۸). شناسایی واژه‌های غیرمفهومی (رایج) در نمایه سازی خودکار مدارک فارسی. فصلنامه کتابداری و اطلاع‌رسانی، ۴۸(۱۲)، ۲۳-۳۵

میرزایی، اردوان. (۱۳۸۵). حشو در زبان با رویکرد نظریه اطلاعات. پیک نور، ۴(۴)، ۴۰-۴۸
وحیدیان کامیار، تقی (۱۳۷۹). دستور زبان فارسی ۱. تهران: سمت
هاشمی، محسن، و ساوجی، محمد حسن (۱۳۸۶). فشرده سازی متن فارسی با استفاده از الگوریتم‌های حسابی و هافمن و مقایسه آن با فشرده سازی متن انگلیسی. مجموعه مقالات پانزدهمین کنفرانس مهندسی برق ایران.
تهران

هویدا، علیرضا (۱۳۷۸). آمار و روش‌های کمی در کتابداری و اطلاع‌رسانی. تهران: سمت.

Bartlett, M. S. (2007). Information Maximization in Face Processing. *Neurocomputing*, 70(13-15), 2204-2217.

Cancho, R., & Martín, F. M. d. P. (2011). Information content versus word length in random typing. *Journal of Statistical Mechanics: Theory and Experiment*, 2011(12), L12002.

- Caraballo, S. A., & Charniak, E. (1999). Determining The Specificity Of Nouns From Text. *Proceedings of the 1999 Joint SIGDAT Conference on Empirical Methods in Natural Language Processing and Very Large Corpora*.
- Fragos, K., Maistros, Y., & Skourlas, C. (2005).A weighted maximum entropy language model for text classification. *NLUCS*, 55-67.
- Frisson, S., Rayner, K., & Pickering, M. J. (2005). Effects of contextual predictability and transitional probability on eye movements during reading. *JExpPsychol Learn MemCogn*, 31(5), 862-877.
- Genzel, D., & Charniak, E. (2002).Entropy rate constancy in text. *Paper presented at the Proceedings of the 40th Annual Meeting on Association for Computational Linguistics*, Philadelphia, Pennsylvania.
- Hegazi, N., Ali, N., & Abed, E. (1987). Information content in textual data: Revisited for Arabic text. *Journal of the American Society for Information Science*, 38(2), 133-137.
- Kireyev, K. (2009). Semantic-Based Estimation of Term Informativeness. *The 2009 Annual Conference of the North American Chapter of the ACL*. Boulder, Colorado
- Lin, L., & Yu-Shu, L. (2005). Research and realization of naive Bayes English text classification method based on base noun phrase identification. *Paper presented at ITI 3rd International Conference*, Cairo.
- Luhn, H. P. (1957). A statistical approach to mechanized encoding and searching of literary information. *IBM Journal of Research and Development*, 1, 309-317.
- MacKay, D. J. C. (2003). *Information theory, inference and learning algorithms*. Cambridge: Cambridge University Press.
- Manning, C. D. & Schutze, H. (1999). *Foundations of Statistical Natural Language Processing*, Cambridge: MIT Press.
- Melamed, I.D. (1997). Measuring Semantic Entropy. *Proceedings of the SIGLEX Workshop on Tagging Text with Lexical Semantics*.
- Montemurro, M. A. & Zanette, D. H. (2001).Entropic analysis of the role of words in literary texts. *Advances in Complex Systems*.5(1), 7-17.
- Nemirovsky, D., & Dobrynin, V. (2008).Word importance discrimination using context information. *In Proceedings TREC*
- Ryu, P. M., & Choi, K. S. (2004). Determining the Specificity of Terms based on Information Theoretic Measures. *Paper presented at Poster Session of 3rd International Workshop on Computational Terminology*. Daejeon, Korea
- Tomokiyo, T., & Hurst, M. (2003). A language model approach to keyphrase extraction. *Paper presented at the Proceedings of the ACL 2003 workshop on Multiword expressions: analysis, acquisition and treatment - Volume 18*, Sapporo, Japan.
- Weber, P. R. (1990). *Basic content analysis*. Newbury Park, Calif.
- Yu, L.-C., Wu, C.-H., & Hovy, E. (2008). OntoNotes: corpus cleanup of mistaken agreement using word sense disambiguation. *Paper presented at the Proceedings of the 22nd International Conference on Computational Linguistics - Volume 1*, Manchester, United Kingdom.
- Zou F., Wang F.L., Deng X., & Han S. (2006), Automatic Identification of Chinese Stop Words. *A special issue on Advances in Natural Language Processing of the journal Research on Computing Science* (p 151-162).