



یک چارچوب داده کاوی رسانه های اجتماعی مبتنی بر LDA برای اقتصاد چرخشی پلاستیک

یانگیمین ژو^{۱*}، چاندراکامهامپاتی، یونگ کیانگ چنگ^۱، نیشیکانت میشر^۲، نور وولندهاری^۲، پائولین دویتز^۳

دریافت: ۲۴ فوریه ۲۰۲۳ / پذیرش: ۲۱ نوامبر ۲۰۲۳ © نویسنده)
۲۰۲۴ (گان)

چکیده

تولیدانبوه زباله های پلاستیکی باعث ایجاد یک بحران فوری در حوزه سلامت عمومی در سراسر جهان شده است. اگرچه سیاست های دولتی و نوآوری های صنعتی نیروهای محرک برای مقابله با این چالش هستند، اما تلاش برای درک نگرش های عمومی می تواند کارایی این فرآیند را بهبود بخشد. رسانه های اجتماعی به راه های اصلی برای کسب اطلاعات و ابراز نظرات و احساسات توسط عموم مردم تبدیل شده اند. این امر ما را بر آن داشت تا با استفاده از داده های رسانه های اجتماعی، برداشت ها و واکنش های رفتاری نسبت به استفاده از پلاستیک را کاوش کنیم. در این مقاله، ما چارچوبی برای جمع آوری و تحلیل داده ها بر اساس رسانه های جریان اصلی در بریتانیا پیشنهاد دادیم تا نظرات عمومی در مورد پلاستیک ها را به دست آوریم. یک مدل یادگیری ماشین بدون نظارت مبتنی بر تخصیص نهفته دیریکله (LDA) برای تجزیه و تحلیل و خوشه بندی موضوعات به منظور مقابله با عدم حاشیه نویسی محتوای داده ها به کارگرفته شده است. سپس یک روش فرهنگ لغت اضافی برای ارزیابی احساسات نظرات ارائه شد. این چارچوب همچنین ابزارهایی برای تجسم مدل و نتایج برای تحریک درک های عمیق ارائه می دهد. ما اثربخشی چارچوب را با استفاده از آن برای تجزیه و تحلیل سه رسانه اجتماعی جریان اصلی، که در آن 6 دسته موضوعی سطح اول و 13 دسته موضوعی سطح دوم از متون نظرات مربوط به پلاستیک ها شناسایی شده اند، تأیید کردیم. نتایج نشان می دهد که احساسات عمومی نسبت به محصولات پلاستیکی به طور کلی پایدار است. توزیع مکانی-زمانی احساسات هر موضوع با تعداد تکرار آن همبستگی بالایی دارد.

کلمات کلیدی: مصورسازی مدل · تحلیل احساسات · طبقه بندی نظرات · LDA

اختصارات

ال دی ای	واژگان پردازش زبان طبیعی
ان ال پی	تخصیص نهفته دیریکله
ترکیبات آلی فرار (VOC)	
آران ان	حافظه کوتاه مدت بلند مدت
LSTM	شبکه عصبی بازگشتی
جی آر یو	واحد بازگشتی دروازه ای
بی بی سی	نمایه سازی معنایی پنهان شرکت
ال اس آی	پخش بریتانیا

امقدمه

محصولات پلاستیکی در دهه 1950 به طور گسترده مورد استفاده قرار گرفتند. به دلیل هزینه تولید پایین، شدت متوسط، سهولت پردازش و راحتی استفاده، پلاستیک تقریباً در همه جا مورد استفاده قرار گرفته و با اهمیت روزافزون در مدیریت پسماند شناخته شده است. گزارش بنیاد الن مک آرتور نشان می دهد که در حال حاضر سالانه حدود 500 میلیون تن پلاستیک در جهان تولید می شود، اما تنها 10٪ از آنها بازیافت می شوند. علاوه بر این، بسیاری از زباله های پلاستیکی از سیستم دفع نشت می کنند و به یک تهدید زیست محیطی تبدیل می شوند. در نتیجه، میزان پلاستیک در زباله های شهری از 1٪ به بیش از 10٪ افزایش یافته است. [۱] این مشکل نشان می دهد که طراحی، تولید و استفاده از پلاستیک ها خطی است، نه یک الگوی دایره ای. به عنوان مثال، با در نظر گرفتن بسته بندی پلاستیکی، حدود ۳۲٪ از تولید جهانی در محیط طبیعی مانند اقیانوس و هوا رها شده و حدود ۴۰٪ از آن دفن شده است. فقط

* یانگیمین ژو

y.xue-2019@hull.ac.uk

چاندراکامهامپاتی

c.kambhampati@hull.ac.uk

یونگ کیانگ چنگ

y.cheng@hull.ac.uk

^۱ گروه علوم و فناوری کامپیوتر، دانشگاه هال، هال 7RX، انگلستان

^۲ دانشکده بازرگانی، دانشگاه هال، هال 7RX، انگلستان

^۳ گروه جغرافیا، زمین شناسی و محیط زیست، دانشگاه هال، هال 7RX، انگلستان

همزمان، تعداد موضوعات ناشناخته طبقه بندی می شود و کیفیت نتایج طبقه بندی ارزیابی می شود. در این مقاله، ما چارچوبی را برای ساده سازی جنبه های مختلف نظرکاوی در مورد پلاستیک ها بر اساس تخصیص دیریکله پنهان (LDA) پیشنهاد می کنیم [۱۳، ۱۴]. یک روش مدل سازی موضوعی بدون نظارت که می تواند موضوعات بالقوه را در اسناد شناسایی کند، برای بررسی موضوعات پنهان یا کلماتی که اغلب در مجموعه ظاهر می شوند، اتخاذ شده است. این کلمات بر اساس احتمالات سند موضوعی و نزدیکی کلمه-موضوع محاسبه می شوند. برخلاف روش سنتی استفاده از تجربه برای قضاوت در مورد تعداد و نتیجه طبقه بندی متن، این مقاله دو روش تخمین را معرفی می کند و نتایج را برای انتخاب نتایج طبقه بندی متن مناسب، مصورسازی می کند. پس از تعیین جنبه های مختلف هر نظر، نتایج طبقه بندی را مصورسازی می کنیم تا درک شهودی تر و واضح تری از فرآیند و دامنه موضوعات ایجاد کنیم. این به خلاصه کردن محتوای موضوع کمک می کند. برای انجام طبقه بندی احساسات، ابتدا امتیاز احساسات هر کلمه را از لیست امتیاز کلمات تعیین شده به دست می آوریم و سپس نتیجه امتیاز احساسات نظرات کامل را برای نرمال سازی محاسبه می کنیم [۱۶] در نهایت، از طریق طبقه بندی موضوعی و امتیازدهی به احساسات، می توان یافته های آموزنده ای از برداشت های عمومی از پلاستیک های بازیافتی و تغییرات رفتاری را برای کمک به سیاست گذاری ارائه داد.

ادامه مقاله به شرح زیر سازماندهی شده است: بخش ۲. بخش مربوطه، اثر مرتبط است. ۳. جزئیات چارچوب را در بخش ... ارائه می دهد. ۴. نتایج مختلف را ارائه می دهدیم و نتایج خود را با سایر روش های پیشرفته مقایسه می کنیم. مقاله را در بخش نتیجه گیری می کنیم. ۵.

۲ کار مرتبط

در اینجا، ما کارهای مرتبط مورد استفاده در این تحقیق، به طور خاص، پیشینه LDA، NLP، ارزیابی مدل موضوعی و تحلیل احساسات را ارائه می دهیم.

۲.۱ پردازش زبان طبیعی

پردازش زبان طبیعی (NLP) پلی بین زبان ماشین و زبان انسان است. این رشته به مطالعه چگونگی فهم زبان انسان توسط کامپیوترهای پردازش [۱۷]. این فناوری اغلب بالاترین وظیفه و قلمرو هوش مصنوعی را در بر می گیرد. این شاخه ای از هوش مصنوعی است و نقطه تلاقی علوم کامپیوتر و زبان شناسی است. NLP همچنین به عنوان زبان شناسی محاسباتی در رشته زبان شناسی شناخته می شود. تحقیقات اولیه پردازش زبان طبیعی شامل تحلیل دستوری، تحلیل معنایی و درک متن است. هر تحلیل به یک مرحله پیش پردازش نیاز دارد که شامل برخی از فناوری های مهم مانند توکن سازی است. در این فرآیند، متن ورودی از سند به کوچکترین واحدها (کلمات، عبارات، شکلک ها و غیره) تقسیم می شود. در بیشتر موارد، این واحد در

۱۰٪ از بسته بندی های پلاستیکی بازیافت شده اند و تنها ۲٪ از آنها به همان سطح از بازسازی دست یافته اند. ۲. در مقابل، در یک سیستم چرخشی، محصولات (از جمله بسته بندی) به گونه ای طراحی می شوند که تا حد امکان دوام بیاورند، قابل استفاده مجدد باشند و در پایان عمر قابل بازیابی باشند. ۱. سیاست گذاران در بریتانیا، اتحادیه اروپا و جاهای دیگر در حال تدوین سیاست هایی برای ممنوعیت پلاستیک های یکبار مصرف و بهبود بازیابی پلاستیک ها هستند. ادبیات [۵، ۳] دیدگاه های عمومی در مورد محصولات پلاستیکی را به عنوان مرجعی برای سیاست گذاران جهت تدوین سیاست های مربوطه و همچنین برای شرکت ها در تدوین برنامه های تجاری، خلاصه کرده اند. [۴].

دروانشناسی، علوم شناختی و جامعه شناسی، مطالعات نشان داده اند که افراد از نظرات برای هدایت تصمیمات و رفتارهای فردی خود استفاده می کنند. [۹] در مطالعه مایا بی، نشان داده شده است که چگونه رسانه های اجتماعی برای عموم اعتیادآور هستند [۲۸] اما برای دستیابی به یک اقتصاد چرخشی پلاستیکی مؤثر، باید بتوان طیف وسیعی از نظرات عمومی را درک کرد. برای مثال، تحمل عمومی برای افزایش هزینه های محصولات زیست محیطی چقدر است؟ مردم تا چه حد از علل و پیامدهای سوءمدیریت پلاستیک آگاه هستند و با آنها همدردی می کنند؟ از آنجایی که رسانه های اجتماعی به تدریج به رسانه ای مهم برای مردم جهت کسب و انتشار اطلاعات تبدیل می شوند، روزانه حجم عظیمی از داده های رسانه های اجتماعی و نظرات در اینترنت تولید می شود. یکی از رویکردها برای دستیابی به این بینش ها، تحلیل نظرات در رسانه های اجتماعی است که بر رفتار عمومی از کسب اطلاعات گرفته تا رفتار پس از خرید نیز تأثیر گذاشته است. [۸-۶] علاوه بر این، صدها میلیون کاربر اینترنت از طریق اینترنت و رسانه های اجتماعی به انواع مختلف اطلاعات مرتبط با پلاستیک دسترسی دارند. در مقایسه با سایر کانال های ارتباطی سنتی، رسانه های اجتماعی بسیار کارآمدتر هستند. به طور کلی، با افزایش تعداد نظرات، موضوعات و احساسات در مورد آنها به سرعت گسترش یافته و تکامل یافته اند و باعث نوسانات روانی و تغییرات عاطفی در سراسر جامعه شده اند. ۱۰. چگونگی حذف مؤثر نظرات اسپم، یافتن نظرات ارزشمند و ارائه آنها به خوانندگان، یا ارائه آنها برای تحلیل بیشتر افکار عمومی و وظایف متن کاوی، ارزش های کاربردی مهمی دارند.

امروزه، افراد کمی در مورد موضوعات خاص مشابه دیدگاه ها در مورد محصولات پلاستیکی، به تحلیل افکار عمومی می پردازند. مطالعات اغلب به این نتیجه می رسند که بخش خصوصی باید اقدامات بیشتری برای کاهش ضایعات پلاستیکی انجام دهد. به دلیل تفاوت در تخصص و انگیزه مفسران آنلاین، نظرات آنلاین اغلب از نظر محتوا و احساسات غیرقابل پیش بینی هستند. [۱۱]. این نظرات حاوی اطلاعات متنی می توانند موضوع و احساس نظر را بررسی کنند تا دیدگاه نظر را تحلیل کنند. هدف از این مطالعه، به دست آوردن دیدگاه های عمومی در مورد موضوع از بررسی های مرتبط با پلاستیک است تا واکنش عموم به اقتصاد پلاستیک آشکار شود. ما نظرات مبتنی بر داده های رسانه های اجتماعی را تجزیه و تحلیل می کنیم تا ویژگی های زمانی و مکانی و ویژگی های توزیع مکانی و زمانی موضوعات مختلف را با استفاده از یک مدل استخراج و طبقه بندی موضوع نامشخص دریافت کنیم [۱۲]. در

روش محاسبه نیز بسیار شهودی است. VOC پی دلیومن = دلیومن /
ن، کجادیلیومن تعداد است من کلمه ام در VOC مربوط به موضوع،
و تعداد کل کلمات مربوط به موضوع.

شکل کلمات یا جملات یا پاراگراف ها. فرآیند نرمال سازی می تواند تمام حروف را به حروف کوچک تبدیل کند. علائم نگارشی و ریشه یابی را حذف کنید، که یک متن کامل را به کلمات اساسی تبدیل می کند [۸] مرحله بعدی تولید یک مدل، یک فرهنگ لغت و یک پیکره زبانی برای فیلتر کردن کلمات با فرکانس پایین و ساخت بردارها است.

۲.۲ فاصله نزدیک (LDA)

[LDA^{۱۳}] یک مدل احتمالی مولد است که می تواند برای طبقه بندی پیکره ها استفاده شود. LDA یک مدل احتمال بیزی سلسله مراتبی سه سطحی است که شامل کلمه، موضوع و سند می شود. اصل این روش این است که فرض کنیم هر کلمه از یک موضوع بالقوه پنهان در پشت آن استخراج می شود. در طول فرآیند مولد، انتخاب موضوع و انتخاب کلمه دو فرآیند تصادفی هستند. در فرآیند تولید، انتخاب موضوع و انتخاب کلمه دو فرآیند تصادفی هستند. ابتدا، برای هر سند، یک موضوع را از توزیع موضوع استخراج کنید. سپس، یک کلمه از توزیع کلمه مربوط به موضوع انتخاب شده استخراج می شود. فرآیند فوق را تا زمانی که هر کلمه در سند پیمایش شود، تکرار کنید. فرض می شود که هم سندی که متعلق به یک موضوع است و هم موضوعی که می تواند توسط یک کلمه نمایش داده شود، از یک توزیع چندجمله ای پیروی می کنند. به طور خلاصه، هدف LDA شناسایی موضوعات از اسناد است، یعنی تبدیل ماتریس سند-کلمه به یک ماتریس سند-موضوع (توزیع).

مزیت LDA این است که یک یادگیری ماشین بدون نظارت است و به مجموعه داده های آموزشی که به صورت دستی حاشیه نویسی شده اند، متکی نیست، یعنی تنها ورودی ها مجموعه اسناد و تعداد موضوعات هستند. علاوه بر این، LDA همیشه می تواند کلمات نماینده ای برای توصیف هر موضوع پیدا کند.

برای مدل متن پنهان LDA، عیب اصلی این است که موقعیت مجموعه ای از کلمات را در متن در نظر نگرفته است. از این رو، نمی تواند معانی مختلفی را که کلمات یکسان با ترتیب های مختلف می توانند بیان کنند، تشخیص دهد. علاوه بر این، از آنجا که متون طولانی حاوی کلمات بیشتری هستند، تطبیق موضوعات آنها دشوار است. مشکل دیگری که باید حل شود این است که کلماتی که از موضوعات مختلف در LDA تشکیل شده اند، دوباره استفاده می شوند. این امر منجر به همپوشانی موضوعات به جای استقلال آنها می شود.

اول، تعریف می کنیم دی همانطور که سند تنظیم شده است. تی مجموعه ای از موضوعات است. در مجموعه سندی، سند می تواند به عنوان یک دنباله کلمه در نظر گرفته شود (\dots, α_n). د شامل n کلمه است. تمام کلمات مختلف در این بی در یک مجموعه بزرگ واژگان (VOC) ترکیب می شوند. در D، احتمال هر سند مربوط به موضوعات مختلف پی تی ا پی تی ... α_2 پی تی ... β_k کجا پی تی من احتمال این است که d متناظر با من موضوع هفتم در تی نگاه شهودی به محاسبه پی تی من تی من /ن، گجانی من تعداد کلمات مربوط به من موضوع هفتم در د، و ن در د تعداد کل کلمات است. برای هر تی در تی احتمال تولید کلمات متفاوتی > پی دبلیو، پی دبلیو ... α_2 پی دبلیو < گجایی دبلیو من احتمال آن است که تی تولید می کند من کلمه هفتم در

$$(1) \quad \text{ص} = \int d w (x_p \times w | t_p w | \text{زمان} | \text{زمان})$$

این فرمول، فرآیند اصلی است که لایه موضوع را به عنوان لایه میانی که کلمات و سند را به هم متصل می کند، در نظر می گیرد و سپس احتمال ... را می دهد. دلیلی در سند از طریق دوتی با استفاده از جریان دوتی، می توانیم محاسبه کنیم ص (دلیلی) یک کلمه در یک سند را وقتی که با هر موضوعی مطابقت دارد، بررسی کنید و سپس موضوع مربوط به کلمه را بر اساس این نتایج به روزرسانی کنید. سپس، اگر به روزرسانی، موضوع مربوط به کلمه را تغییر دهد، به نوبه خود تأثیر خواهد گذاشت. دو

در زیر فرآیند یادگیری الگوریتم LDA آمده است: در ابتدا، دوتایی به صورت تصادفی اختصاص داده می شود (برای همه دوتایی سپس، فرآیند اصلی فوق را به طور مداوم تکرار کنید، تا نتیجه همگرایی نهایی، که خروجی LDA است.

برای من کلمه هفتم در یکی از موضوعات سند ۵، هنگامی که مربوط به کلمه، فرمول می تواند به صورت زیر جی، سپس موارد فوق اصلاح شود:

(2) $\left(\begin{smallmatrix} 1 & 0 \\ 0 & 1 \end{smallmatrix} \right) \left(\begin{smallmatrix} 1 & 0 \\ 0 & 1 \end{smallmatrix} \right) \left(\begin{smallmatrix} 1 & 0 \\ 0 & 1 \end{smallmatrix} \right)$ پڄي ديليومن |دها $\left(\begin{smallmatrix} 1 & 0 \\ 0 & 1 \end{smallmatrix} \right)$ پڄي ديليومن |دها $\left(\begin{smallmatrix} 1 & 0 \\ 0 & 1 \end{smallmatrix} \right)$ پڄي ديليومن |دها.

هنگام برشمردن موضوعات درتی برای بدست آوردن همه‌چی (که \mathcal{H})، آن مقدارچی است؟ یک سپس، طبق این مقادیر احتمال، برابر است بامن کلمه هفتم در \mathbf{h} یک موضوع انتخاب کنید. ساده ترین ایده این است که من که به حد اکثر می رساند \mathbf{h} (من \mathbf{h}) (توجه داشته باشید که فقط چی یک متغیر در این فرمول است)، که $\arg\max$ است.

در مرحله بعد، اگر متن کلمه اصلی، *مرد* را موضوعی را از بین موضوعات انتخاب می کند آن را تحت تأثیر قرار خواهد داد. *مرد* می تواند از آن استنباط کرد فرمول محاسبه دو بردار فوق). این تأثیر همچنین به محاسبه ... منتقل خواهد شد. *پ (دلیو اد)* بالا. محاسبه کنید *پ (دلیو اد)* برای همه در همه حال در درونی و موضوع را به عنوان یک موضوع دوباره انتخاب کنید. این روش را تا بعد از آن اجرا کنید تکرارهای حلقه، نتیجه نهایی مورد نیاز LDA را پوشش داده و دریافت می کند. در اینجا، ما دو احتمال را به دست آوردیم *پ (کلمه | موضوع)* و *پ (موضوع | سند)*.

۲.۳ ارزیابی مدل موضوعی

اکثروش های سنتی از بازرسی بصری یا دانش قبلی برای ارزیابی عملکردمدل های انتخاب موضوع استفاده می کنند. شهودی ترین روش، قضاوت دستی در مورد موضوع استخراج شده است، اما بدیهی است که این کار زمان بر است. قضاوت بصری دستی عمدتاً شامل ارزیابی نتایج خوشه بندی با استفاده از فناوری تجسم ومعرفی معنای موضوعات توسط کلمات کلیدی تولید شده از مدل، موضوع است [۱۵].

ازسوی دیگر، روش های ارزیابی خودکار شامل ارزیابی اثر
خوشه بندی توسط ضریب سیلوئیت و انسجام برای اندازه گیری
اینکه آیا کلمات موجود در موضوع ... هستند یا خیر، می شوند.

ضریب سیلوئت (Silhouette Coefficient) یک روش ارزیابی عملکرد خوشه بندی است که اولین بار توسط پیتر جی. روسو در سال ۱۹۸۶ پیشنهاد شد. این روش ترکیبی از انسجام (cohesion) و تفکیک پذیری (resolution) است. این روش یکی از روش های ارزیابی تأثیر الگوریتم های مختلف یا حالت های عملیاتی مختلف بر نتایج خوشه بندی بر اساس داده های اصلی یکسان است. برای هر نمونه من در داده ها، میانگین فاصله از من به نمونه های دیگر در این خوشه است (الف/ا) حداقل میانگین فاصله از I تا تمام نمونه های خوشه های دیگر، b است. (من) فرمول می تواند به صورت زیر باشد:

بازنویسی شده به صورت:

$$(من) \text{ (من) } - \text{الف (من) حد اکثر } \{ \text{الف (من) } , \text{ب (من) } \}$$

$$(من) = \frac{1 - \text{الف (من) (ب (من))}}{\text{الف (من) (ب (من))}}$$

(4)

- (1) میانگین فاصله را محاسبه کنید/الف (i) بین نمونه i و سایر نمونه های همان خوشه. چه زمانی/الف (i) کوچکتر باشد، اثر گروه بندی خوشه من بهتر است. (من) به عنوان عدم تشابه درون خوشه ای نمونه در نظر گرفته می شود.
- (2) میانگین فاصله بین نمونه I و سایر نمونه ها در خوشه دیگر b برابر است با (من). این به عنوان عدم تشابه بین خوشه ای نمونه تعریف می شود. من:

من = دقیقه ب من ۱، ب من ۲، ...، ب من i کی

- (3) چه زمانی من نزدیک به ۱ باشد، خوشه بندی نمونه من منطقی است. اگر من نزدیک به ۱- باشد، عکس آن صادق است.

انسجام موضوعی به این صورت است که با مقایسه شباهت معنایی بین کلمات با امتیاز بالاتر، به موضوع امتیاز داده می شود. در موضوع.

$$\text{هسته } i \text{ من، وی جی، } = \text{ورود به سیستم وی، وی جی، } + \text{لاچی وی من وی جی، } - \text{(()) (())}$$

پنجم گروهی از کلمات است که برای توصیف موضوع استفاده می شود. E به معنای اطمینان از این است که نمره یک عدد حقیقی را برمی گرداند. در اینجا، هر چه اپسیلون کوچکتر باشد، نتیجه کوچکتر خواهد بود. احتمال کلمه ص با شمارش فراوانی کلمات روی مجموعه داده محاسبه می شود. این الگوریتم بر اساس مجموعه داده اصلی مدل موضوع آموزشی است و به مجموعه داده خارجی متکی نیست.

۲.۴ تحلیل احساسات

تحلیل احساسات فرآیند تحلیل، پردازش، القا و استدلال ذهنی متن با رنگ های احساسی است [۱۹، ۲۰]. تحلیل احساسات می تواند از روش های سنتی مبتنی بر فرهنگ لغت احساسات یا روش های مبتنی بر یادگیری عمیق استفاده کند. [۲۵] روش مبتنی بر فرهنگ لغت عمدتاً شامل تدوین مجموعه ای از فرهنگ های احساسات و قوانین، تجزیه جملات، تجزیه و تحلیل و تطبیق است.

فرهنگ لغت های متن (به طور کلی تحلیل اجزای کلام، تحلیل وابستگی نحوی)، محاسبه مقادیر احساسی، و در نهایت استفاده از مقادیر احساسی به عنوان گرایش احساسی متن. [مبنای قضاوت] [۲۲، ۲۱، ۱۶]. بر اساس طبقه بندی احساسات مبتنی بر یادگیری عمیق، جمله ابتدا بر اساس طبقه بندی احساسات مبتنی بر یادگیری عمیق است. ابتدا جمله پیش پردازش می شود، مانند تقسیم بندی کلمات، کلمات توقف و تبدیل ساده شده و سنتی. سپس، کدگذاری بردار کلمات و همچنین استخراج ویژگی با استفاده از RNN (شبکه عصبی بازگشتی) مانند LSTM (حافظه کوتاه مدت بلند) یا GRU (واحد بازگشتی دروازه دار) انجام می شود.

(3) [۲۳] مراحل عملیاتی عبارتند از: سند باید از جملات و برجسب ها

تشکیل شده باشد. جمله را توکنیزه کنید،

بنابراین لیستی از کلمات آن را نشان می دهد. سپس، از طریق ویژگی های ساده کلمات یونیگرام، به ترتیب از نمونه های ذهنی و عینی برای حفظ توزیع متعادل و یکنواخت کلاس در مجموعه آموزش و مجموعه تست استفاده می شود. از این ویژگی ها برای به دست آوردن یک نمایش مقدار ویژگی از مجموعه داده خود استفاده می شود. سپس باید طبقه بندی کننده را روی مجموعه آموزش آموزش دهیم و در نهایت نتایج را خروجی دهیم.

۳ چارچوب پیشنهادی

این چارچوب شامل چهار بخش است: جمع آوری و پیش پردازش داده ها، طبقه بندی موضوعات، انتخاب مدل و مصورسازی، و تحلیل احساسات، همانطور که در شکل نشان داده شده است. ۱.

در مرحله اول، داده های نظرات را می توان از رسانه های اجتماعی آنلاین بررسی کرد. در مرحله بعد، نظرات پیش پردازش می شوند تا نویز محتوای تکراری یا نظرات با ساختار غیرمتنی برای پردازش در نظرات فیلتر شوند. پس از پیش پردازش، یادگیری ماشین بدون نظارت LDA بر روی مجموعه داده های آموزشی آموزش داده شد. پس از آن، یک ارزیابی انجام شد.

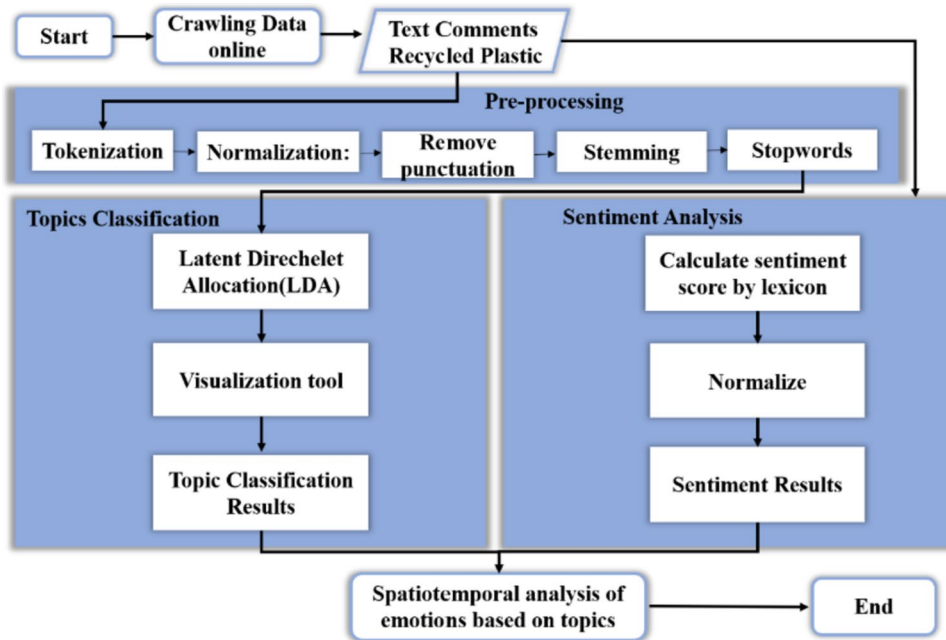
(5) فرآیند و یک ابزار تجسم برای کمک به ما در انتخاب تعداد موضوعات

و خلاصه کردن موضوعات. در عین حال، مدل ارزیابی احساسات به هر نظر یک امتیاز احساسی می دهد. این امتیاز می تواند برای طبقه بندی احساسات یا مطالعه تغییرات در احساسات استفاده شود. در نهایت، با ترکیب طبقه بندی موضوعی و نمرات احساسی، می توانیم یک سری زمانی از تغییرات احساسی در جنبه های مختلف محصولات پلاستیکی به دست آوریم. از آنجایی که داده های مورد بررسی از نظرات خوانندگان رسانه های اصلی بریتانیا از جمله بی بی سی، گاردین و میل آنلاین در زیر اخبار مربوط به پلاستیک قابل بازیافت گرفته شده است، داده ها دارای ویژگی های غیرقابل پیش بینی بودن و بدون برجسب های طبقه بندی موجود هستند. بنابراین، ما باید از تکنیک هایی استفاده کنیم که نیازی به مجموعه آموزشی برای طبقه بندی نظرات ندارند. در اینجا، محقق یک روش یادگیری بدون نظارت مبتنی بر کلمه به نام LDA را معرفی می کند.

۳.۱ جمع آوری و پیش پردازش داده ها

چند رسانه مانند بنگاه سخن پراکنی بریتانیا (بی بی سی)، گاردین و میل آنلاین تا حد زیادی

شکل ۱ روش تحقیق



رتبه بندی 30 کلمه ای که بیشترین تکرار را در داده های جمع آوری شده ما داشته اند را نشان می دهد.

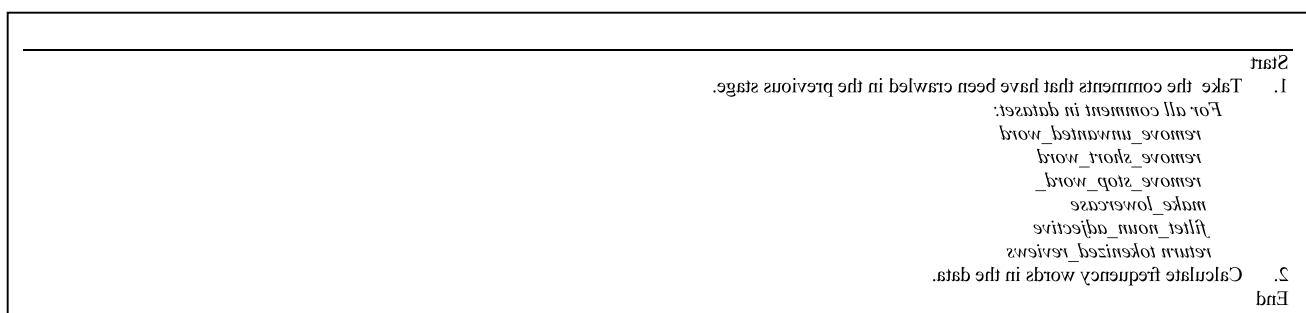
۳.۲ طبقه بندی موضوعات

چارچوب استخراج و طبقه بندی موضوع بر اساس مدل موضوعی LDA ساخته شده است و احساسات عمومی در مورد موضوع به صورت لایه ای از متون مرتبط در رسانه های اجتماعی به دست می آید. از مدل موضوعی LDA برای استخراج موضوع استفاده می شود تا توزیع احتمال موضوع هر متن و توزیع احتمال کلمه هر موضوع تولید شود. پس از آن، داده های نمونه از موضوع حاشیه نویسی شده در کل مجموعه داده ها طبقه بندی می شوند.

مصرف اخبار آنلاین در بریتانیا را تحت سلطه خود قرار داده است [۲۴] اعتماد خوانندگان به این رسانه ها به این معنی است که صرف نظر از اینکه خبر چیست، آنها مستقیماً و مکرراً جستجو و مطالعه می کنند. نظرات بیان شده پس از خواندن خبر اغلب ارتباط زیادی با موضوع خبر دارد.

این مرحله، فرآیندی است که توسط الگوریتم ۱ انجام می شود، از جمله: متن نظرات اصلی حاوی اطلاعات تداخلی مانند فاصله ها، لینک های http و علائم نگارشی است. برای حذف نویز و بهبود کارایی تقسیم بندی کلمات، داده های اصلی باید فیلتر متن شوند. ما از عبارات منظم پایتون برای فیلتر کردن متن اصلی رسانه های اجتماعی و حذف اطلاعات تداخلی (مانند لینک های http، علائم نگارشی)، کلمات بی کیفیت، متن بی کیفیت و متن تکراری استفاده می کنیم. شکل ۲

الگوریتم ۱

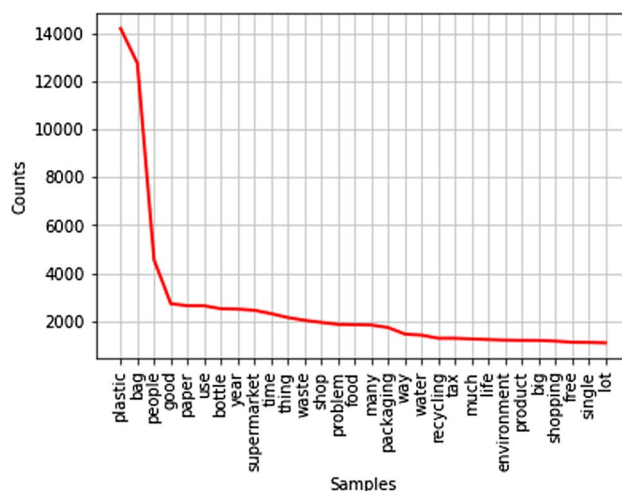


ساختارهای موضوعی را از تعداد زیادی سند استخراج می کند و هر سند را به یک یا چند موضوع اختصاص می دهد (شکل ۱). [۲۵] هنگام ساخت یک مدل موضوعی، برخی پارامترهای اصلی باید از قبل تنظیم شوند. این پارامترها شامل تعداد موضوعات K ، آلفای پیشینی توزیع موضوع، بتای پیشینی توزیع کلمات و تعداد اسنادی است که باید در هر اندازه قطعه بلوک آموزشی استفاده شوند. تعداد کل قبولی های ارزیابی آموزشی. در این آزمایش حفظ شده است. $(\alpha=0.1, \beta=5, \gamma=1)$.

چگونه تعداد مناسب موضوعات را انتخاب کنیم؟ هرچه موضوعات بیشتری انتخاب کنیم، موضوعات خاص تر خواهند بود. با این حال، انتخاب موضوعات زیاد ممکن است تشخیص نظرات را به روشی معنادار غیرممکن کند. در عین حال، تعداد کم موضوعات منجر به ترکیب نظرات به سمت جنبه های یکسانی می شود که باید به دسته های مختلفی تعلق داشته باشند.

بامقایسه کیفیت موضوعات تولید شده، تعداد بهینه موضوعات انتخاب می شود. مرحله بعدی محاسبه و مقایسه احتمالات موضوعاتی است که یک بررسی واحد به آنها تعلق دارد و تعیین موضوعات نظر. در نهایت، از این موضوع برای طبقه بندی همه اسناد استفاده می شود. مباحث مختلفی از مدل LDA ساخته شده اند [۲۶]، که در آن هر موضوع ترکیبی از کلمات کلیدی است و هر کلمه کلیدی وزن خاصی به موضوع می دهد. از لیست کلمات کلیدی، کلماتی که به درک و خلاصه سازی موضوعات کمک می کنند. این کلمات کلیدی وزن ها برای خلاصه سازی محتوای موضوع استفاده خواهند شد.

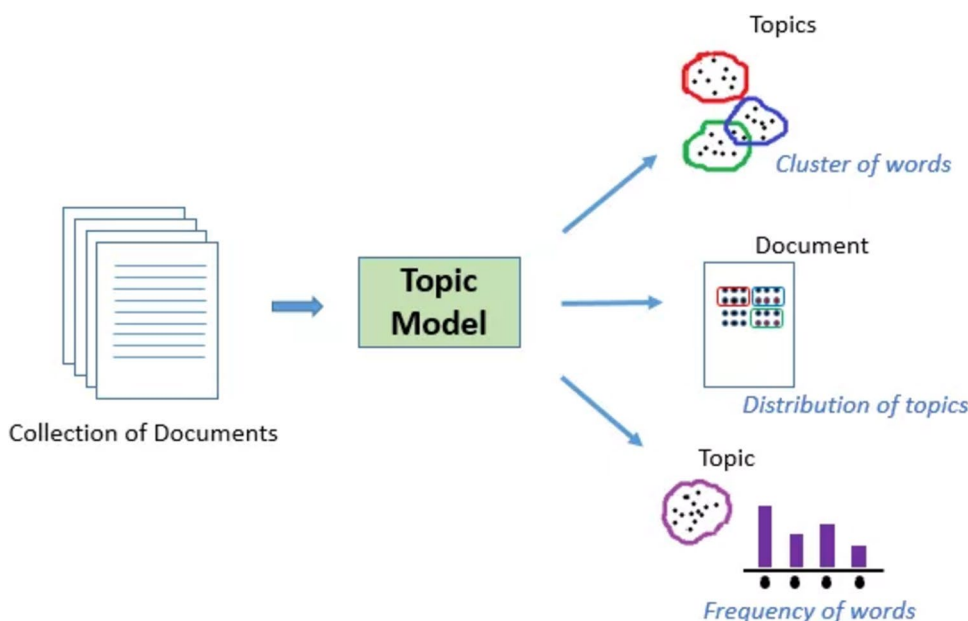
برای هر فرآیند LDA، مراحل زیر به عنوان الگوریتم ۲ دنبال می شوند:



شکل ۲ 30 کلمه پرتکرار در داده ها

مدل های LDA قادر به تشخیص موضوعات در اسناد و کاوش در پیکره های اطلاعات پنهان هستند و طیف گسترده ای از کاربردها را در سناریوهایی مانند تجمیع موضوعات، استخراج اطلاعات از متن بدون ساختار و انتخاب ویژگی دارند. این مدل می تواند (1) الگوهای موضوعی پنهان در پیکره را کشف کند؛ (2) اسناد را بر اساس موضوعات حاشیه نویسی کند؛ و (3) از حاشیه نویسی برای سازماندهی، جمع آوری، خلاصه سازی و بازیابی اسناد استفاده کند. LDA یک مدل احتمالاتی برای حل مسئله مدل سازی موضوع متن است. این یک روش یادگیری بدون نظارت است که به طور خودکار الگوهای پنهان را کشف می کند.

شکل ۳ یک گردش کار از مدل قالب معمولی [۲۵]



الگوریتم ۲

2	End
4	Summarize topics.
3	Display the results of topic classification (key words and word frequency under Hidden topic content)
	$w_{m,n}$: a particular word for word classification $[m, n]$
	$z_{m,n}$: topic index of the word in comments m
	N_m : the length of comments m
	\hat{v}_m : topic distribution for comments m
	N : vocabulary size
	M : the total number of comments
	$\hat{\phi}$: Dirichlet parameters
	$\hat{\phi}_k$: word distribution for topic k
	K : the number of topics
	Parameters and variables:
	end for
	end for
	sample term for word $w_{m,n} \sim Mult(\hat{\phi}_{z_{m,n}})$
	sample topic index $z_{m,n} \sim Mult(\hat{v}_m)$
	for all words $n \in [1, N_m]$:
	sample document length $N_m \sim Potts(\hat{\alpha})$
	sample mixture proportion $\hat{v}_m \sim Dir(\hat{\alpha})$
	end for
	sample mixture components $\hat{\phi}_k \sim Dir(\hat{\beta})$
	for all topics $k \in [1, K]$:
2	and word-topic.
2	Selected LDA model, calculates the probability of each word appearing in the document and compares also the probability of document-topic
1	From the previous pre-processed step, get the term list.

۳.۳ انتخاب مدل و تجسم نه

در مرحله بعد، ما از یک ابزار تعاملی مبتنی بر وب برای نمایش نتایج مدل LDA استفاده می کنیم [۱۴]، از جمله معنی هر موضوع، شیوه هر موضوع و ارتباط هر موضوع.

یک سیستم تجسم تعاملی مبتنی بر وب است. این سیستم نه تنها ترتیب جهانی کلمات را در یک موضوع ارائه می دهد، بلکه توزیع موضوعی کلمات خاص و کلمات انحصاری را نیز نمایش می دهد. به نوبه خود، یک معیار جدید برای ارتباط پیشنهاد می کند. این مقاله همچنین از نتایج یک مطالعه کاربری استفاده می کند تا نشان دهد که انتخاب کلمات به ترتیب نزولی احتمال برای تفسیر موضوع بهینه نیست. برای خلاصه سازی موضوع، بهینه ترین توضیح صرفاً کلمات با بیشترین فراوانی یا کلمات کاملاً خاص نیست. بلکه باید از کلماتی که مرتبط تر هستند برای توضیح طبقه بندی موضوع استفاده شود LDAvis

رابط کاربری عملیات در شکل نشان داده شده است. ۴ می توان بانگه داشتن ماوس روی دایره سمت چپ، موضوع خاص را مشاهده کرد. پس از انتخاب، سمت راست 30 کلمه برتر مرتبط با این موضوع را نشان می دهد. می توان از این کلمات برای خلاصه کردن میانگین این موضوع استفاده کرد. اندازه دایره، فراوانی این موضوع را نشان می دهد. در اینجا، ما از تحلیل چندبعدی استفاده می کنیم و مؤلفه های اصلی را استخراج می کنیم.

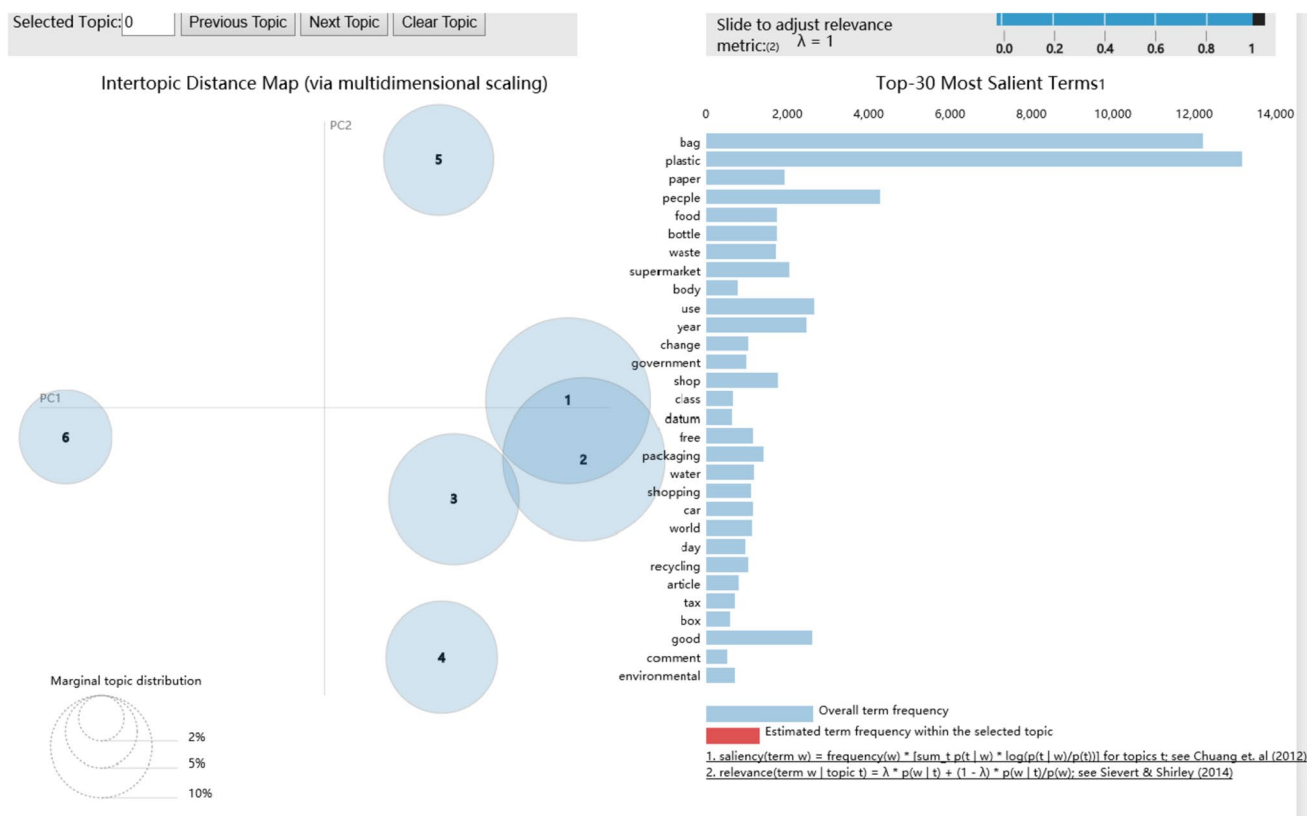
به عنوان ابعاد، و توزیع موضوع را روی این دو بعد قرار دهید. فاصله بین موضوع نشان دهنده ارتباط آنها است. ترتیب کلمات در سمت راست، وزن کلماتی را که در موضوع نقش دارند نشان می دهد. ارتباط را می توان با پارامتر تنظیم کرد $\lambda \in [0, 1]$. تغییر مقدار λ ترتیب وزن کلمات در موضوع را تغییر می دهد، که به تجسم محتوای اصلی موضوع کمک می کند. هرچه بزرگتر باشد λ ، هرچه تعداد کلمات بیشتر باشد، فرکانس آنها کمتر است. λ ، کلمات خاص تر.

$$\text{ارتفاع} = p(w|t) + \lambda * p(w) + (1 - \lambda) * \frac{p(w|t)}{p(w)} \quad (6)$$

هدف از طبقه بندی موضوعات، داشتن حداقل تعداد ممکن از دسته بندی های موضوعی و نسبتاً مستقل است. زیرا تعداد زیاد دسته بندی های موضوعی بی معنی است. همپوشانی موضوعات، نتایج تحقیق را گیج کننده خواهد کرد. این چارچوب دارای توابع انعطاف پذیری برای تنظیم تعداد دسته بندی ها و موضوعات است.

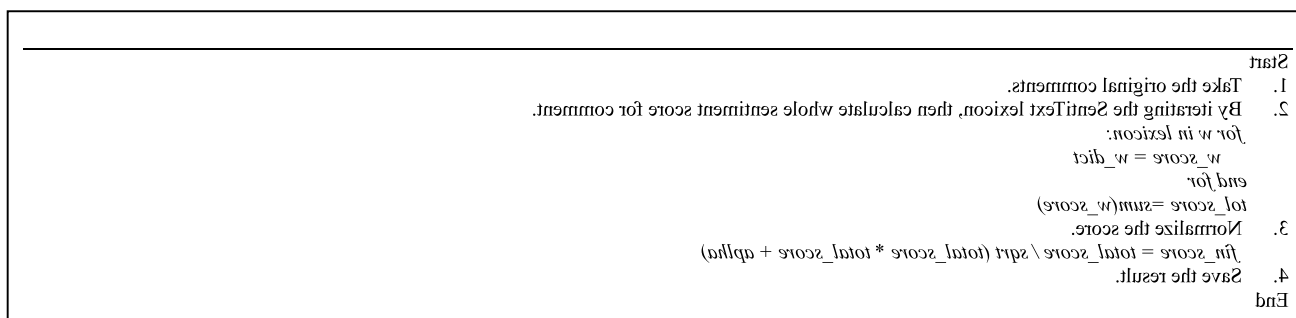
۳.۴ تحلیل احساسات

ما از ابزار تحلیل Vader متعلق به NLTK استفاده می کنیم [۱۶] در این چارچوب به عنوان الگوریتم ۳. Vader (فرهنگ لغت و استدلال احساسات آگاه از ظرفیت)، یک ابزار تحلیل احساسات مبتنی بر فرهنگ لغت و قانون است که نیازی به آموزش یا ... ندارد.



شکل ۴ تجسم مدل موضوعی

الگوریتم ۳



(7)

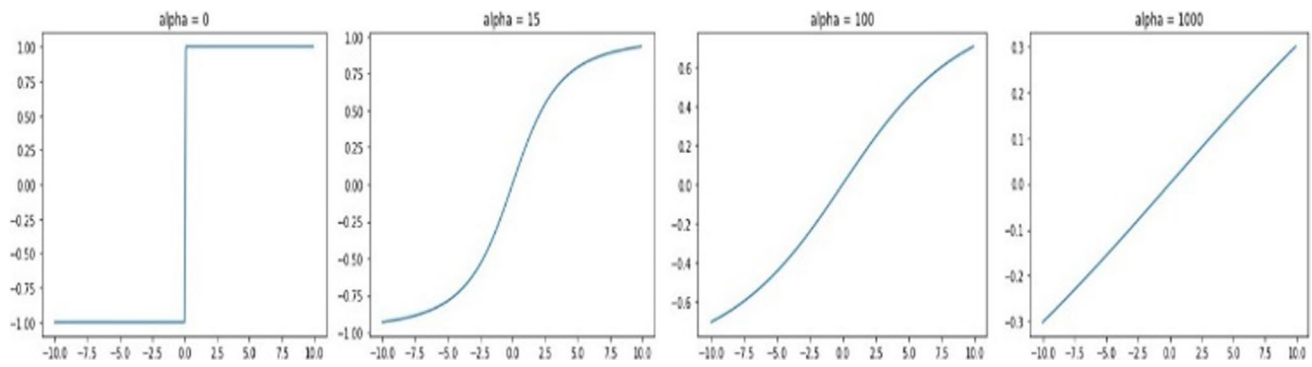
$$\text{کوری هنجار} = \sqrt{\frac{\text{امتیاز} +}{\text{امتیاز}}}$$

درفرمول، پارامتری است که برای تنظیم توزیع فاصله ای نتایج نرمال شده استفاده می شود. تأثیر معادله (۷) در شکل نشان داده شده است. برای طبقه بندی راحت تر احساسات، در مقاله ما آلفا = ۱۵ انتخاب شده است. این امر باعث می شود نتیجه نرمال شده در محدوده مورد نظر یکنواخت تر باشد. در نهایت، یک آستانه استاندارد برای طبقه بندی جملات تعیین کنید:

احساسات مثبت: امتیاز ترکیبی ≤ 0.05 .

به منظور کسب امتیاز ارزیابی احساسات، سفارشی سازی شده است و به طور ویژه با احساسات بیان شده در رسانه های اجتماعی تنظیم شده است. این ابزار برای محتوای رسانه های اجتماعی طراحی شده است، بنابراین در محتوایی که می توانید در رسانه های اجتماعی پیدا کنید، بهترین عملکرد را دارد.

این ابزار برای هر نظر یک امتیاز ترکیبی ارائه می دهد. امتیاز ترکیبی با جمع کردن امتیازهای متناظر هر کلمه در فرهنگ لغت، سپس تنظیم آنها طبق قوانین و در نهایت نرمال سازی آنها بین 1- و 1+ محاسبه می شود. این فرهنگ لغت شامل بیش از 9000 کلمه است که به صورت دستی با امتیازهای احساسی مشخص شده اند. فرمول نرمال سازی به صورت زیر است:



شکل ۵: چگونه آلفا بر نتیجه نرمال شده تأثیر می گذارد

و دیگران مردم به محصولات پلاستیکی، خرید و سیاست بیشتر از محتوای خانواده، غذا و غیره توجه می کنند، به ترتیب ۱۰، ۱۶۸، ۹۲۸۱ و ۵۷۰۰. ترتیب کلمات کلیدی به ترتیب بزرگی ضرایب است. ضریب فراوانی وقوع کلمه است. کلماتی که با رنگ قرمز مشخص شده اند، کلمات کلیدی هستند که برای تسهیل درک خواننده از موضوع خلاصه شده اند.

برای درک جزئیات دقیق، ما یک طبقه بندی ثانویه برای سه موضوع برتر با بیشترین تعداد پیام، یعنی محصولات پلاستیکی، خرید و سیاست گذاری، به صورت جدول انجام دادیم.

۴.۱ مجموعه داده ها

داده ها از فید خوانندگان مقالات خبری از بی بی سی، گاردین و دیلی میل آنلاین گرفته شده اند. به طور خاص، ما ۷۹۲۴ نظر از ۲۶ خبر از بی بی سی، با بازه زمانی از آگوست ۲۰۱۷ تا اکتبر ۲۰۱۹؛ ۲۲۸۵۷ نظر از ۹۷ خبر از گاردین، با بازه زمانی از جولای ۲۰۱۶ تا اکتبر ۲۰۱۹؛ و ۳۳۲۱ نظر از ۲۲ خبر از دیلی میل آنلاین، با بازه زمانی از ژانویه ۲۰۱۸ تا اکتبر ۲۰۱۹ داریم.

موضوع محصولات پلاستیکی شامل هفت زیرموضوع است: استفاده مجدد، منابع، بازیافت، بسته بندی آب، آلودگی و موارد دیگر. سه زیرموضوع در موضوع خرید وجود دارد: بسته بندی قابل استفاده مجدد، بسته بندی پولی و بسته بندی تحویل. علاوه بر این، کسب و کار، رسانه و محیط زیست زیرموضوعات موضوع سیاست گذاری هستند.

۴.۳ انتخاب مدل

مانمرات انسجام شماره های مختلف موضوع را در شکل ۱ مقایسه کرده ایم. می توانیم ببینیم که وقتی تعداد موضوعات ۶ است، امتیاز انسجام بالاترین است. جدول ۳

۴.۲ دسته بندی موضوعات

میزان نتایج دسته بندی نظرات بر اساس موضوع را نشان می دهد. مانظرات را به شش دسته تقسیم کردیم: محصول پلاستیکی، خرید، سیاست، خانواده، غذا

جدول ۱: شرح مباحث سطح اول

اعداد	کلمات کلیدی
۱۰,۱۶۸	محصول: ۰.۰۹۰*پلاستیک+ ۰.۰۲۰*بطری+ ۰.۰۲۰*ضایعات+ ۰.۰۱۳*آب+ ۰.۰۱۲*روغن+ ۰.۰۱۲*محصول+ ۰.۰۱۲*بازیافت
۹۲۸۱	موضوع ۲: خرید ۰.۱۴۵*کیسه+ ۰.۰۶۴*پلاستیک+ ۰.۰۲۳*کاغذ+ ۰.۰۲۱*سوپرمارکت+ ۰.۰۲۱*استفاده+ ۰.۰۱۸*فروشگاه+ ۰.۰۱۴*افراد-لطفاً+ ۰.۰۱۴*رایگان+ ۰.۰۱۳*خرید+ ۰.۰۱۱*سال
۵۷۰۰	موضوع ۳: سیاست ۰.۰۳۴*مردم+ ۰.۰۲۰*خوب+ ۰.۰۱۹*تغییر+ ۰.۰۱۸*دولت+ ۰.۰۱۷*چیز+ ۰.۰۱۵*جهان+ ۰.۰۱۳*مشکل+ ۰.۰۱۳*مالیات+ ۰.۰۱۳*محیط زیست+ ۰.۰۱۱*کشور+ ۰.۰۱۳*lem
۳۵۵۰	موضوع ۴: خانواده نفر+ ۰.۰۲۲*سال+ ۰.۰۱۸*روز+ ۰.۰۱۶*زمان+ ۰.۰۱۵*ماشین+ ۰.۰۱۴*مرد+ ۰.۰۱۳*بچه+ ۰.۰۱۲*فرزند+ ۰.۰۱۱*کفش+ ۰.۰۲۲*هزینه+ ۰.۰۱۰*بسیاری
۲۷۳۳	موضوع ۵: غذا ۰.۰۲۷*غذا+ ۰.۰۱۶*جعبه+ ۰.۰۱۴*بسته-میوه+ ۰.۰۱۲*گوشت+ ۰.۰۱۱*چیز+ ۰.۰۱۱*خوب+ ۰.۰۱۰*محل+ ۰.۰۰۹*بی کیفیت+ ۰.۰۰۹*کوچک+ ۰.۰۱۴*ing
۲۱۱۹	موضوع: سایر ۰.۰۰۹*under+ ۰.۰۰۹*demand+ ۰.۰۱۰*article+ ۰.۰۱۶*comment+ ۰.۰۱۹*datum+ ۰.۰۲۳*class+ ۰.۰۲۴*body+ ۰.۰۲۸*سوال+ ۰.۰۰۹*متشکرم+ ۰.۰۰۹*درسته+ ۰.۰۰۸*سوال

جدول ۲ شرح موضوع در

موضوع سطح مشت	مبحث سطح دوم	وزن هاو کلمات کلیدی
محصولات پلاستیکی یوس نی	دیگر	کلاس+ "0.010* پترو+ "0.024* com+ "0.033* لیمک+ "0.032* بدنه+ "0.028* داده+ "0.025*+ "0.035* زیر خط+ "0.009* بزرگ+ "0.007* یاد+ "0.009* leum
استفاده مجدد		0.105* بطری+ "0.052+ "لیوان+ "0.024* پلاستیک+ "0.020* شیر+ "0.019* با "فروشگاه+ "0.018* سپرده+ "0.015* فروشگاه+ "0.012* افراد+ "0.011* طرح+ "0.009* سوپرمارکت "0.083* پلاستیک+ "0.024* روغن+ "0.021* مصرف+ "0.018* انرژی+ "0.018* محصول+ "0.016* کیسه+ "0.011* جایگزین "0.010* کاغذ+ "0.008* تکی+ "0.008* ماده
باز یافت		0.۰۷۴* پلاستیک+ "۰.۲۸* ضایعات+ "۰.۱۸* کیسه+ "۰.۱۷* باز یافت+ "۰.۱۶* غذا+ "۰.۱۶* بسته بندی+ "۰.۱۲* احتمالا "مردم+ "0.011+ باز یافت+ "0.009* زیاد+ "0.012* lem
بسته آب		0.076* آب+ "0.021+ پلاستیک "نفت+ "0.012+ تقاضا+ "0.012* شیر+ "0.011* سال+ "0.009* رشد+ "0.009* زغال سنگ+ "0.009* مواد+ "0.009* بطری شده+ "0.014* tic
آلودگی		0.113* پلاستیک+ "0.032* کیسه+ "0.018* اقیانوس+ "0.016* زیاله+ "0.012* غذا+ "0.012* کشور+ "0.011* دریا+ "0.010* مالش- "مشکل+ "0.009* چیز+ "0.009* bish
سیاست	کسب و کار	0.032* پلاستیک+ "0.022* دولت+ "0.022* مالیات+ "0.018* کیسه+ "0.016* مشکل+ "0.015* پول+ "0.014* مردم لطفاً+ "0.010* محیطی+ "0.010* خوب+ "0.009* بزرگ
	رسانه	0.031* خوب+ "0.022* تغییر+ "0.018* چیز+ "0.017* آب و هوا+ "0.017* نکته+ "0.013* مقاله+ "0.011* افراد لطفاً+ "0.011* روش+ "0.010* محیط+ "0.010* عالی
	محیط زیست	"مردم+ "0.023+ جهان+ "0.014+ انسان+ "0.014* کشور+ "0.040 "امتحان کنید+ "0.013+ بسیاری+ "0.009+ جهانی+ "0.008* فقیر+ "0.008* جمعیت+ "0.008* سیاره+ "0.008* چیز "0.015* نفر+ "0.009+ زمان+ "0.008* وسایل+ "0.008* شخص+ "0.007* بطری+ "0.007* ماشین+ "0.006* سال+ "0.006* کفش هزینه+ "0.006* خوب+ "0.005* چیز
خرید	بسته استفاده مجدد	0.028* کیسه+ "0.025+ فروشگاه+ "0.016* پلاستیک+ "0.016* کاغذ+ "0.014* شارژ+ "0.014* سوپرمارکت+ "0.013* مشتری- "0.013+ افراد+ "0.011+ پول+ "0.011* رایگان
	بسته پولی	"0.150* کیسه+ "0.082* پلاستیک+ "0.025* استفاده+ "0.020* کاغذ+ "0.017* سوپرمارکت+ "0.015* خرید+ "0.014* افراد لطفاً+ "0.013* سال+ "0.013* مجرد+ "0.011* حامل
	بسته تحویل	

چهارزشی دارد که کلمات رایج عبارتند از تغییر، دولت، مالیات و محیط زیست. اینها کلمات مهمی هستند که به ما در خلاصه کردن موضوع کمک می کنند. علاوه بر این، آب و هوا، فقیر، جهانی، شغل و عمل در کلمات خاص ظاهر می شوند و مردم و جهان در کلمات با فرکانس بالا ظاهر می شوند. بنابراین، ما فکر می کنیم این موضوع در مورد سیاست است. باتوجه به توزیع کلمات تحت مدل موضوعی تولید شده توسط R های مختلف، این طرح نامگذاری طبیعتاً مناسب است: فقط تعداد کلمات محتمل (مثلاً ۵ تا ۱۰) و خاص ترین کلمات در توزیع باید به عنوان توصیفگر موضوع استفاده شوند. این روش معمولاً خوب عمل می کند.

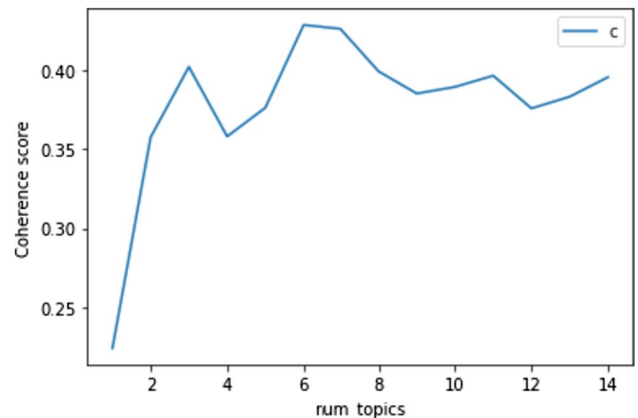
۴.۴ تحلیل احساسات موضوعات

میانگین نمرات احساسات همانطور که در جدول نشان داده شده است ۴ میانگین کلی احساسات خنثی است. نظرات مردم در مورد خرید، غذا و خانواده نسبتاً مثبت است. سایر موارد خنثی هستند.

ارقام ۷ و ۸ رایج ترین کلماتی که در نظرات مثبت و منفی ظاهر می شوند را نشان می دهد. در میان آنها، مردم عموماً نسبت به اسراف دیدگاه منفی تری دارند.

۴.۵ مقایسه عملکرد

مانتایج LSI، LDA (شاخص گذاری معنایی پنهان) و HDP (فرآیند دیریکله سلسله مراتبی) را مقایسه کردیم [۲۷]. برای بدست آوردن موضوع ... عمل می کند (SVD) بر اساس تجزیه مقدار تکین LSI یک مدل موضوعی ساده و کاربردی است LSI.



شکل ۶ نمره انسجام موضوعات مختلف عددی

جزئیات ترتیب کلمات کلیدی در موضوعات مختلف را خلاصه می کند.

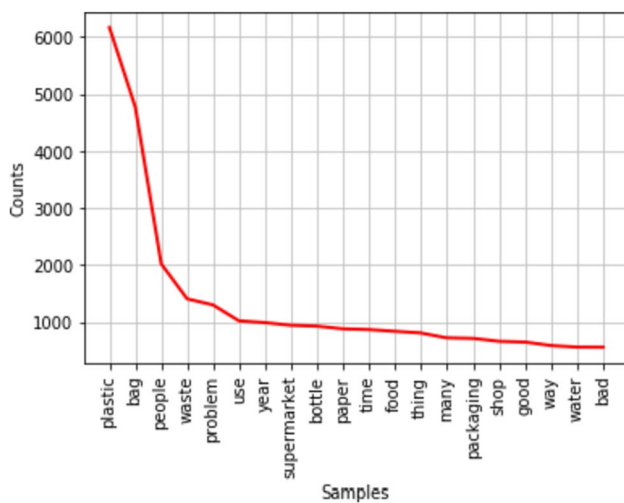
روش اصلی ما برای خلاصه سازی خوشه ها، استفاده از تابع تجسم برای یافتن کلمات رایج با فراوانی بالا ($\lambda=1$) و ویژگی های خاص ($\lambda=0$) و همچنین یک نشانگر موجودی ($\lambda=0.5$) بین فراوانی بالا و ویژگی های خاص بر اساس شش موضوع شامل محصول، غذا، خانواده، سیاست، خرید و موارد دیگر. همانطور که در جدول مشاهده می شود، چه زمانی اگر ۱ باشد، کلماتی که بیشتر ظاهر می شوند، ارتباط بیشتری با موضوع دارند. وقتی ۰ است، تنظیم شده باشد، کلماتی که خاص تر و انحصاری تر هستند، ارتباط بیشتری با موضوع دارند. برای مثال، وقتی به درون سیاست موضوع نگاه می کنیم، مهم نیست

جدول ۳ ترتیب کلمات کلیدی در موضوعات مختلف

$\lambda=0$	$\lambda=0.5$	$\lambda=1$
محصول بطری، زباله، آب، بازیافت، شیشه، نام، پیوند، بازیافت، انرژی، com	پلاستیک، بطری، زباله، آب، بازیافت، نفت، محصول، شیشه، نام، لینک	پلاستیک، بطری، زباله، آب، روغن، محصول، بازیافت، استفاده از بسته بندی، شیشه
غذا جعبه، میوه، گوشت، فله، مقوا، ماهی، سبزیجات، زباله، قانون، تازه	غذا، جعبه، میوه، گوشت، فله، مقوا، ماهی، سبزیجات، زباله، بسته بندی	غذا، جعبه، بسته بندی، میوه، گوشت، چیزهای دیگر، خوب، محلی، گشاد، کوچک
خانواده مرد، بچه، کودک، قهوه، فنجان، زن، مدرسه، خود، سگ، خود	مرد، روز، بچه، کودک، قهوه، فنجان، ماشین، زن، سال، مدرسه	مردم، سال، روز، زمان، ماشین، مرد، بچه، کودک، قهوه، زیاد، قدیمی
سیاست تغییر، دولت، مالیات، محیط زیست، آب و هوا، فقیر، جهانی، شغل، مدت، اقدام	تغییر، دولت، مردم، مالیات، جهان، خوب، محیط زیست، مسئله، سیاره، انسان	مردم، خوبی، تغییر، حکومت، چیز، جهان، مشکل، مالیات، محیط زیست، مسئله
خرید کیف، کاغذ، رایگان، خرید، هزینه، حامل، لاینر، خیریه، خواربارفروشی، کول	کیسه، پلاستیک، کاغذ، سوپرمارکت، مغازه، رایگان، خرید، استفاده، هزینه، اپراتور	کیسه، پلاستیک، کاغذ، سوپرمارکت، استفاده، مغازه، مردم، رایگان، خرید، سال
دیگر بدنه، کلاس، داده، توضیح، زیرخط دار، تشکر، سوال، حیوان، زمین، نرم	بدنه، کلاس، داده، توضیح، زیرخط دار، تشکر، مقاله، سوال، زمین، حیوان	بدنه، کلاس، داده، نظر، مقاله، تقاضا، زیر خط کشیدن، تشکر کردن، درست گفتن، سوال کردن

جدول ۴ میانگین امتیاز احساسات

میانگین احساسات امتیاز	
۰.۰۴۵۲۴	کل نظرات
۰.۰۴۶۶۴	سیاست
۰.۰۸۴۹۷	خرید
۰.۰۶۰۷۴	غذا
۰.۰۲۷۴۶	دیگر
۰.۰۰۵۰۲	خانواده
۰.۰۱۷۱۲	محصول پلاستیکی



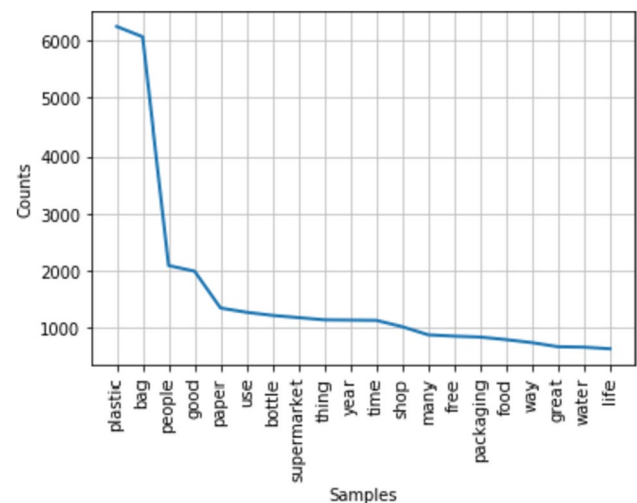
شکل ۸ بیست کلمه ای که بیشترین تکرار را در نظرات منفی داشته اند

جدول ۵ عملکرد مدل تاپیکز

عملکرد مدل موضوعی	الیدی ای	الاس ای	اچدی بی
ضریب سیلوئیت	۰.۰۵۵۳۰۷	۰.۰۶۵۸۹۵	۰.۰۲۱۳۱۴
نمره انسجام	۰.۳۳۰۱۵	۰.۲۷۴۷۰۱	۰.۱۷۴۸۷۲

نمره ضریب سیلوئیت LDA برابر با ۰.۰۵۵ است که به ... نزدیک تر است. ۱، نشان می دهد که خوشه بندی LDA معقول است. دو مدل تم دیگر منفی هستند، که نشان می دهد نتایج طبقه بندی آنها بیش از حد انتظار دقیق نیست. در امتیاز انسجام، LDA بالاترین امتیاز یعنی ۰.۳۲ را کسب کرد که نشان می دهد کلمات با معانی مشابه در مدل LDA تمایل دارند در زمینه های مشابه ظاهر شوند. اکثر کلمات تحت طبقه بندی موضوعی آن ارتباط نزدیکی با هم دارند، بنابراین موضوع منسجم تر در نظر گرفته می شود. در جدول ۵ نتایج نشان می دهد که LDA در اثر خوشه بندی و انسجام موضوعی اسناد متنی کوتاه، مزایایی دارد.

متن. مدل HDP از ویژگی تقسیم نامتناهی دسته فرآیند دیریکله در فضای محدود استفاده می کند و به طور تطبیقی تعداد موضوعات را برای بدست آوردن مجموعه موضوعات با تخصیص بهینه ساختار مجموعه اسناد تطبیق می دهد.



شکل ۷ بیست کلمه ای که بیشترین تکرار را در نظرات مثبت داشته اند

جدول ۶ نتایج احساسات در مورد تحلیل موضوعات

نتایج احساسات در مورد موضوعات		
موضوع	احساسات	نتایج (به درصد)
سیاست	منفی	۷.۸۴۶۷
	خنثی	۶.۴۵۹۰
	مثبت	۲.۶۶۸۲
خرید	منفی	۱۲.۵۳۶۹
	خنثی	۸.۸۱۱۵
	مثبت	۶.۲۸۹۳
غذا	منفی	۳.۴۰۶۷
	خنثی	۲.۷۲۷۷
	مثبت	۲.۰۹۳۴
دیگر	منفی	۲.۴۴۴۸
	خنثی	۲.۰۲۲۰
	مثبت	۱.۸۴۳۳
خانواده	منفی	۳.۹۴۵۷
	خنثی	۳.۷۴۰۲
	مثبت	۲.۸۸۵۶
محصول	منفی	۱۱.۷۷۷۵
	خنثی	۱۰.۷۳۲۳
	مثبت	۷.۷۶۹۳
درصد کل		۱۰۰

۲۶۳۵ نظر مثبت و ۸۹۶ نظر منفی بود. از نظر خرید، ۴۲۱۰ نظر مثبت و ۲۱۱۲ نظر منفی بود. از نظر محصولات، ۳۹۵۵ نظر مثبت و ۲۶۰۹ نظر منفی بود.

از میانگین امتیاز احساسات، عموم مردم نسبت به وضعیت فعلی پلاستیک محافظه کار و خوش بین هستند. در میان آنها، خرید، غذاخوری و خانواده مثبت هستند. با این حال، با قضاوت از تعداد نظرات طبقه بندی شده بر اساس احساسات، تعداد نظرات منفی جزء اصلی را تشکیل می دهد. نظرات منفی بیشتری در موردشش موضوع نسبت به نظرات مثبت وجود داشت.

۴.۷ تحلیل احساسات مبتنی بر زمان

تحلیل سری های زمانی شش دسته موضوعی سطح اول در شکل نشان داده شده است. ۹ به جز موضوع خرید که پس از سال ۲۰۱۸ کمی کاهش یافت، بقیه موضوعات همگی در حال افزایش هستند. به خصوص پس از سال ۲۰۱۷، حجم نظرات به شدت افزایش یافته است.

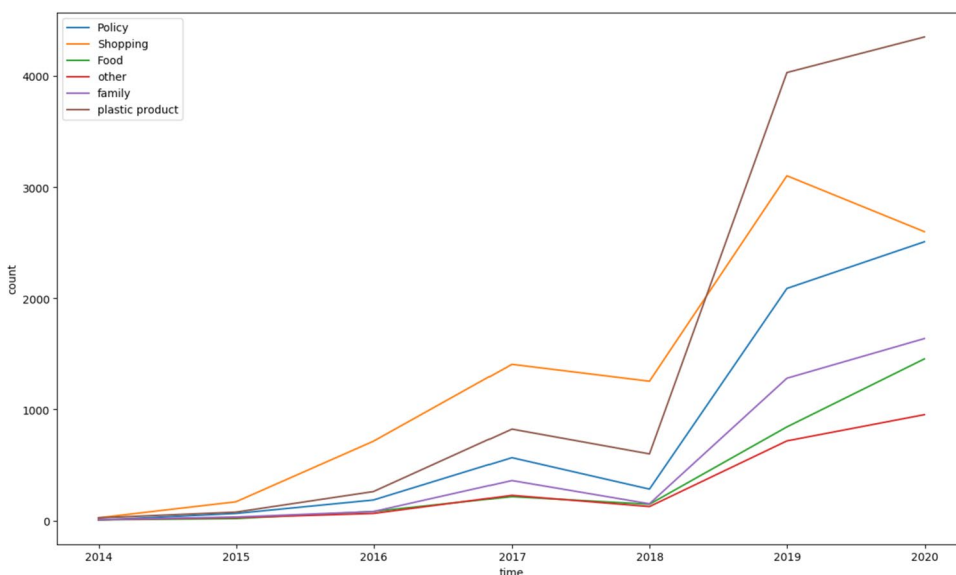
در شکل ۱۰، این نمودار سری زمانی تغییرات احساسی مردم را نشان می دهد. نگرش مردم نسبت به سیاست به تدریج از منفی به مثبت تغییر کرد. احساسات مردم در مورد خرید مثبت شد. رضایت مردم از مواد غذایی و محصولات پلاستیکی کاهش یافته است. احساسات در مورد موضوع خانواده نوسان دارد، اما خنثی است. احساسات عمومی نسبت به محصولات پلاستیکی از منفی در سال ۲۰۱۳ به مثبت در سال ۲۰۱۴ تغییر کرد و به تدریج خنثی شد.

۴.۶ تحلیل احساسات مبتنی بر موضوع

از جدول می توان نتیجه گرفت ۶ اینکه خوانندگان اغلب در مورد اخبار پلاستیک نظر می دهند. نتیجه تجزیه و تحلیل، ۱۴۰۹۰ نظر مثبت، ۱۱۵۸۳ نظر منفی و ۷۹۰۸ نظر خنثی است. اگرچه نظرات مثبت زیادی وجود دارد، اما هنوز نظرات منفی زیادی وجود دارد، به این معنی که عموم مردم از وضعیت فعلی محصولات پلاستیکی راضی نیستند یا برای حذف نظرات منفی به بهبود نیاز دارند. از نظر سیاست ها،

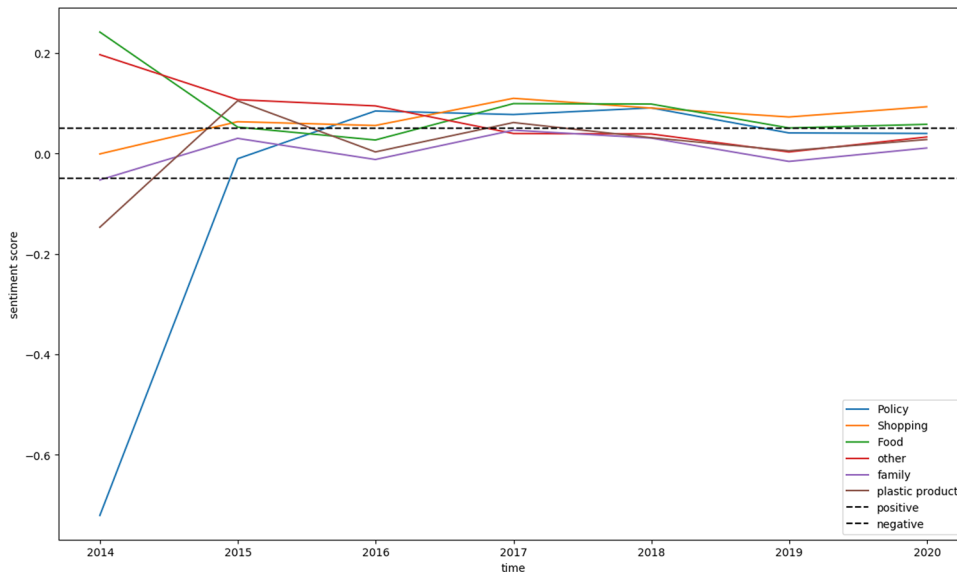
۵ نتیجه گیری

مایک چارچوب استخراج و طبقه بندی موضوع برای به دست آوردن و تحلیل برداشت های عمومی مبتنی بر رسانه های اجتماعی پیشنهاد دادیم. بر اساس این چارچوب، ما



شکل ۹ سری زمانی موضوعی سطح اول مرتبط

شکل ۱۰ احساسات مرتبط با سری زمانی موضوعی سطح اول



وضعیت فعلی بدتر شده است. مشکلات بیشتری کشف شده اند که منجر به ابراز احساسات منفی بیشتری شده است. مشکل پلاستیک توجه فزاینده ای را به خود جلب کرده و منجر به افزایش چشمگیر توجه و همدردی عمومی برای تجزیه و تحلیل محیط زیست شده است.

تقدیرنامه ماایلیم از استاد راهنمایم، چاندراسخار، به خاطر راهنمایی هایش در تمام مراحل این فرایند تشکر کنم. ماایلیم از پروفیسور نیشیکانت به خاطر ارائه سرنخ های اولیه برای داده های تحقیق تشکر کنم.

مشارکت های نویسندگان. نوشته شد و همه نویسندگان در مورد نسخه های قبلی نسخه خطی اظهار نظر کردند. همه نویسندگان نسخه خطی نهایی را خواندند و تأیید کردند YX انجام شد. اولین پیش نویس نسخه خطی توسط YX نگارش - بررسی و ویرایش، مدیریت پروژه، همه نویسندگان در مفهوم سازی و طراحی مطالعه مشارکت داشتند. تهیه مطالب، جمع آوری داده ها و تحلیل توسط PD: تحلیل نتایج، نگارش - بررسی NW: گردآوری داده ها، جمع آوری منابع داده NM: گردآوری داده ها، نگارش - بررسی و ویرایش، نظارت، مدیریت پروژه YC: گردآوری داده ها، نگارش - بررسی و ویرایش، نظارت، مدیریت پروژه CK: مفهوم سازی، روش شناسی، تحلیل قالب، نرم افزار، اعتبارسنجی، نگارش - پیش نویس اصلی، مصورسازی، نگارش - بررسی و ویرایش YX:

بودجه نویسندگان اعلام می کنند که ما هیچ گونه منافع مالی رقابتی یا روابط شخصی شناخته شده ای نداریم که بتواند بر کار گزارش شده در این مقاله تأثیر بگذارد.

در دسترس بودن داده ها و مواد تمام داده های زیربنایی این مقاله بنا به درخواست معقول نویسنده مسئول به اشتراک گذاشته خواهد شد.

اعلامیه ها

تضاد منافع هیچ تضاد منافی در ارسال این مقاله وجود ندارد و انتشار آن توسط همه نویسندگان تأیید شده است. من از طرف نویسندگان همکارم اعلام می کنم که کار شرح داده شده، تحقیقی اصیل بوده که قبلاً منتشر نشده و برای انتشار در جای دیگری، چه به طور کامل و چه جزئی، در دست بررسی نیست. همه نویسندگان ذکر شده، مقاله پیوست را تأیید کرده اند.

داده های رسانه های اجتماعی از نظرات مردم در مورد محصولات پلاستیکی را به طور جامع تجزیه و تحلیل کردیم. نتایج ارزیابی نشان می دهد که روش استخراج موضوع و طبقه بندی پیشنهادی در این مقاله امکان پذیر است. در عین حال، همراه با امتیاز احساسی رسانه های اجتماعی، در بعد زمان تجزیه و تحلیل می شود و نتایج می تواند تغییرات احساسات را تحت طبقه بندی پیدا کند.

پس از مقایسه نتایج الگوریتم مدل های موضوعی مختلف، با مقایسه دو شاخص مدل طبقه بندی موضوعی، یعنی کارایی سیلوئیت و امتیاز انسجام، مدل LDA مناسب تری را برای طبقه بندی نظرات انتخاب کردیم. در ترکیب با امتیاز ارزیابی احساسات، احساسات را تحت طبقه بندی موضوعی طبقه بندی می کنیم. نتایج طبقه بندی نشان می دهد که تعداد نظرات منفی از سوی عموم مردم بیشتر از نظرات مثبت است. در نهایت، در ترکیب با تحلیل سری های زمانی، مشخص شد که بحث در مورد پلاستیک ها سال به سال در حال افزایش است. با این حال، نگرش عموم مردم نسبت به طبقه بندی موضوعات پلاستیکی به تدریج صلح آمیزتر و خنثی تر شده است.

از منظر نمرات احساسات، نظرات عمومی عموماً معتقدند که پایین ترین امتیاز محصولات پلاستیکی به این معنی است که پلاستیک ها جای بیشتری برای بهبود در محصولات دارند. با توجه به نسبت نظرات طبقه بندی شده بر اساس احساسات، نظرات منفی بیشتری وجود دارد و عموم مردم نسبت به این سیاست منفی تر هستند. در عین حال، نتایج تحقیق نشان می دهد که اگرچه مردم از سال ۲۰۱۸ توجه بیشتری به اطلاعات مربوط به پلاستیک نشان داده اند (تعداد نظرات به شدت افزایش یافته است)، اما ابراز احساسات عمومی درجات مختلفی از کاهش را نشان داده است. این ممکن است به این دلیل باشد که در سال ۲۰۱۸، دولت بریتانیا یک استراتژی جدید منابع و زباله را که به کاهش ضایعات پلاستیکی اختصاص داده شده بود، منتشر کرد. در نتیجه، توجه مردم به پلاستیک ها...

دسترسی آزاد این مقاله تحت مجوز بین المللی Attribution 4.0 Creative Commons منتشر شده است که استفاده، اشتراک گذاری، اقتباس، توزیع و تکثیر در هر رسانه یا قالبی را مجاز می داند، مادامی که به نویسنده(گان) اصلی و منبع، اعتبار کافی داده شود، پیوندی به مجوز Creative Commons ارائه شود و در صورت ایجاد تغییرات، مشخص شود که آیا تغییراتی ایجاد شده است یا خیر. تصاویر یا سایر مطالب شخص ثالث در این مقاله در مجوز Creative Commons مقاله گنجانده شده اند، مگر اینکه در خط اعتباری مطلب، خلاف آن ذکر شده باشد. اگر مطلبی در مجوز Creative Commons مقاله گنجانده نشده باشد و استفاده مورد نظر شما طبق مقررات قانونی مجاز نباشد یا از حد مجاز تجاوز کند، باید مستقیماً از دارنده حق چاپ اجازه بگیرید. برای مشاهده نسخه ای از این مجوز، به <http://creativecommons.org/licenses/by/4.0/>.

منابع

۱. استاهل، دبلیو آر: اقتصاد چرخشی. طبیعت ۵۳۱(7595)، 438-435 (2016)
۲. بنیاد الن مک آرتور، مرکز مک کینزی برای تجارت و محیط زیست. رشد درون: چشم انداز اقتصاد چرخشی برای یک اروپای رقابتی [M]. بنیاد الن مک آرتور (2015)
۳. اشنور، آر ای جی، آلبیوی، وی، چوداری، ام. و همکاران: کاهش آلودگی دریایی ناشی از پلاستیک های یکبار مصرف (SUPS): یک بررسی. مار. آلوده کننده. بول. آلوده کننده. بول. ۱۳۷، ۱۵۷-۱۷۱ (۲۰۱۸)
۴. دیلکس-هافمن، ال. اس.، پرات، اس.، لیکاک، بی. و همکاران: نگرش عمومی نسبت به پلاستیک ها، منابع، حفاظت، بازیافت، بازیافت. ۲۳۷، ۱۴۷-۲۳۵ (۲۰۱۹)
۵. برزنو، ن. بازده مثبت و تعادل: بازخورد همزمان بین افکار عمومی و سیاست اجتماعی. مطالعه سیاست. مجله مطالعات سیاست. مجله ۴۵ (4)، 583-612 (2017)
۶. منگولد، دبلیو جی، فالدرز، دی جی: رسانه های اجتماعی: عنصر ترکیبی جدید از ترکیب تبلیغات، اتوبوس، افق. افق. ۵۲(4)، 365-357 (2009)
۷. فیشر، ای.، روبر، ای. آر. تعامل اجتماعی از طریق رسانه های اجتماعی جدید: (چگونه) تعاملات در توئیتر می تواند بر تفکر و رفتار مؤثر تأثیر بگذارد؟ جی. یاس. و نتور. ۲۴(1)، 1-18 (2011)
۸. ویلیامز، آر ال، کاترل، جی. چهار روش هوشمندانه برای اداره جوامع آنلاین. مدیریت اسلون ام آی تی. مدیریت تجدیدنظر. تجدیدنظر. ۴۱(4)، 81 (2000)
۹. مارگتس، اچ. رفتار سیاسی و اکوستیک رسانه های اجتماعی. نات. هوم. بها. و بها. ۱(4)، 1-3 (2017)
۱۰. لاروش، م.، حبیبی، م.، ریچارد، م. او. و همکاران: تأثیرات جوامع برند مبتنی بر رسانه های اجتماعی بر نشانگرهای جامعه برند، شیوه های خلق ارزش، اعتماد به برند و وفاداری به برند. محاسبات. رفتار. رفتار. رفتار. رفتار. ۲۸(5)، 1755-1767 (2012)
۱۱. اسپارکس، بی. ای.، براونینگ، وی. شکایت در فضای مجازی: انگیزه ها و اشکال شکایات آنلاین مهمانان هتل. مجله بیمارستان. مارکت. مدیریت. مدیریت. ۱۹(7)، 797-818 (2010)
۱۲. ساندرسان، ن.، ژانگ، ی.، باودین، س. و همکاران: سیستم و روشی برای استخراج موضوع و نظرکاوی: ثبت اختراع ایالات متحده ۸,۵۳۳,۲۰۸، ۱۰ سپتامبر ۲۰۱۳
۱۳. هدایت الله، آ. ف.، معارف، م. ر. وظایف پیش پردازش در پیام های توئیتر اندونزیایی. مجله کنفرانس فیزیک. فصلنامه ۸۰۱(1)، 012072 (2017)
- Res. تخصیص دیریکله نهفته. جی. ماخ. یاد بگیرید. ۳۱4. Blei, DM, Ng, AY, Jordan, MI: ۱۰۲۲-۹۹۳ (۲۰۰۳) CAD/CAM و Komputerowego Wspomagania Medycyny 2018(Zakład Projektowania Systemów. تشخیص و تجسم خودکار. موضوعات در مجموعه داده های متنی بزرگ: 15. Romaszko KP: 16. باکشی، آر کی، کاور، ان.، کاور، آر. و همکاران: نظرکاوی و تحلیل احساسات. در: سومین کنفرانس بین المللی محاسبات برای توسعه پایدار جهانی (IEEE INDIACOM). 2016. صفحات 455-452 (2016)
۱۷. هوتو، سی جی، گیلبرت، ای. ویدر: یک مدل مبتنی بر قاعده و صرفه جو برای تحلیل احساسات متن رسانه های اجتماعی. در: هشتمین کنفرانس بین المللی AAAI در مورد وبلاگ ها و رسانه های اجتماعی (۲۰۱۴)
۱۸. چوداری، کی. آر. پردازش زبان طبیعی. در: مبانی هوش مصنوعی، صفحات ۶۰۳-۶۴۹. اشپرینگر، دهلی نو (۲۰۲۰)
۱۹. جیان کیانگ، ز.، شیائولین، گ. تحقیقات مقایسه ای در مورد روش های پیش پردازش متن در تحلیل احساسات توئیتر. IEEE Access. ۵، ۲۸۷۰-۲۸۷۹ (۲۰۱۷)
۲۰. پریانتینا، ر.، نوپمبر، ITS، سارنو، ر. و همکاران: تحلیل احساسات نظرات هتل ها با استفاده از تخصیص پنهان دیریکله، شباهت معنایی و LSTM. مجله بین المللی مهندسی سیستم ها. ۱۲(4)، 142-155 (2019)
۲۱. لو، ل. روش تحلیل احساسات متن شبکه ای با ترکیب نمایش متن LDA و GRU-CNN. Pers. Ubiquit. Comput. Ubiquit. Comput. ۲۳(۳-۴)، ۴۱۲-۴۰۵ (۲۰۱۹)
۲۲. همتیان، ف.، سهرابی، م. ک. مروری بر تکنیک های طبقه بندی برای نظرکاوی و تحلیل احساسات. Artif. Intell. Rev.. Intell. Rev. ۵۲(3)، 1545-1495 (2019)
۲۳. جونگلینگ، ر.، داتا، س.، سربرنیک، الف. انتخاب سلاح های شما: در مورد ابزارهای تحلیل احساسات برای تحقیقات مهندسی نرم افزار. در: کنفرانس بین المللی IEEE در مورد نگهداری و تکامل نرم افزار (ICSME) 2015. صفحات 535-531 (2015)
- موجود در [آنلاین] Ofcom.org.uk [۲۴. Ofcom.org.uk report.pdf/news-consumption-2019- _data/assets/pdf_file/0027/157914/uk-https://www.ofcom.org.uk (۲۰۲۱). دسترسی در ۱۶ نوامبر ۲۰۲۱]
۲۵. دو، اچ اچ، پراساد، پی دبلیو سی، ماگ، ای. و همکاران: یادگیری عمیق برای تحلیل احساسات مبتنی بر جنبه: یک بررسی مقایسه ای. سامانه تخصصی. کاربرد. ۱۱۸، ۲۷۲-۲۹۹ (۲۰۱۹)
۲۶. هدایت الله، ای اف، آدیتیا، اس کی، کریمه، اس تی جی و همکاران: مدل سازی موضوعی شرایط آب و هوایی و اقلیمی در توئیتر با استفاده از تخصیص دیریکله پنهان (LDA). کنفرانس IOP، سری مواد، علوم، مهندسی. ۴۸۲(1)، 012033 (2019)
۲۷. genism.org معنانشناسی آماری در پایتون. برگرفته از - ۲۷. Řehůřek, R., Sojka, P.: Gensim
۲۸. مایا، ب.، آرکوب، او. ای: راه حل های تقریبی هیلبرت و رفتارهای هندسی کسری یک مدل کسری دینامیکی از اعتماد به رسانه های اجتماعی که توسط عملگر دیفرانسیلی کاپوتو کسری تأیید شده است. آشوب، سالیتون ها، فراکتال ها: ۱۰X، ۱۰۰۹۲ (۲۰۲۳)
- یادداشت ناشر** اشپرینگر نیچر در مورد ادعاهای مربوط به صلاحیت قضایی در نقشه های منتشر شده و وابستگی های سازمانی بی طرف باقی می ماند.