

استخراج رابطه مبتنی بر تعبیه لغات



با فرآیند جمع‌سپاری

محمد جعفرآباد* و روح‌الله دیانت

گروه مهندسی فناوری اطلاعات، دانشکده فنی و مهندسی، دانشگاه قم، قم، ایران

چکیده

برای انجام مطالعات داده‌کاوی، تاحدودی به دلیل پیچیده‌بودن فرآیند انتخاب ویژگی در کار مورد نظر، نیاز داریم تا بخشی از برچسب زنی را به کارگران در فعالیت جمع‌سپاری واگذار کنیم. فرآیند واگذاری کارهای داده‌کاوی به کاربران، اغلب به وسیله سامانه‌های نرم افزاری و بدون اطلاع دقیق از موقعیت سنی یا جغرافیای محل سکونت کاربران صورت می‌گیرد. عدم اطمینان از عملکرد کاربران مجازی در جمع‌سپاری، میزان صحت اطلاعات دریافتی را کاهش می‌دهد. در این مقاله پیشنهاد داده‌ایم تا با استفاده از روش‌های ایجاد انگیزش، تعدادی از مردم را در محلی جمع و از آنها در جهت وظایف جمع‌سپاری استفاده کنیم. افزایش دقت در اعلام نتایج به دلیل حضور فیزیکی، سرعت بالا در گرفتن نتایج با دقت بالا در زمان تعیین‌شده، تحصیلات مناسب شرکت‌کنندگان در فعالیت و بومی بودن طرح اجرایی از ویژگی‌های این پژوهش هستند. در این پژوهش یک کار یادگیری ماشین انجام شد تا بتوانیم در ضمن آن فعالیت‌های جمع‌سپاری را با الگوریتم‌های شبکه عصبی عمیق ترکیب کنیم. وظیفه رده‌بندی برای تعبیه لغات به صورت الگوریتمی و تلفیقی با کمک جمع‌سپاری انجام می‌شود. روش پیشنهادی با افزودن داده‌های جمع‌سپار به داده‌های قبلی و تغییرات در مدل تعبیه لغات ترکیبی گلاو و وردتووک توانست نتایج مناسبی را در استخراج ویژگی به دست بیاورد.

واژگان کلیدی: جمع‌سپاری، تعبیه لغات، گلاو، وردتووک، رده‌بندی

Relation extraction based on word embedding with Crowdsourcing Process

Mohammad Jafarabad* & Rouhollah Dianat

Department of Information Technology, Faculty of Engineering,
Qom university, Qom, Iran

Abstract

For data mining studies, due to the complexity of doing feature selection process in tasks by hand, we need to send some of labeling to the workers with crowdsourcing activities. The process of outsourcing data mining tasks to users is often handled by software systems without enough knowledge of the age or geography of the users' residence. We use convolutional neural network, for doing classification in six classes: USAGE, TOPIC, COMPARE, MODEL-FEATURE, RESULT and PART-WHOLE. This article extracts the data from the abstract of 450 scientific articles and it is a total of 835 relations. One hundred of these abstracts have been selected by the crowdsourcing. Classification results in this article have been done with a slight improvement in accuracy. In this study, we computed the classification results on a combination of vocabulary vectors with using of 450 abstract relation data (100 crowd source datasets with 350 standards). The results of the implementation of the classification algorithm give us performance improvement. This paper uses the population power to perform preparing data mining works. The proposed method by adding crowdsource data to the previous data was able to obtain better results rather than the top 5 methods.

Keywords: Glove, Word2vec, Crowdsourcing, word embedding, classification

* Corresponding author

* نویسنده عهده‌دار مکاتبات

سال ۱۴۰۱ شماره ۱ پیاپی ۵۱

• تاریخ ارسال مقاله: ۱۳۹۸/۷/۵ • تاریخ پذیرش: ۱۳۹۹/۹/۱ • تاریخ انتشار: ۱۴۰۱/۰۳/۳۱ • نوع مطالعه: پژوهشی



تلاش برای ادغام استعداد، دانش و ایده‌های توزیع شده در فرایندهای نوآورانه به سوی خود جلب کرده‌اند.

هر سیستم جدیدی که تعریف می‌شود، یک سری خدمات را ارائه می‌دهد که این خدمات یک راه حل برای انجام کارهایی هستند که در قبل به‌سادگی قابل حل نبوده‌اند. سامانه جمع‌سپاری راه حلی اجتماعی برای وظایف کوچکی که قبلاً با برنامه‌نویسی قابل حل نبوده‌اند، ارائه می‌دهد. اما خود سیستم پیشنهادی نیز می‌تواند شامل ایراداتی باشد که در جمع‌سپاری بزرگترین مشکل تضمین کیفیت نتایج است. چگونه می‌توان به سمت مسیری حرکت کرد تا یک جمع‌سپار را با هزینه‌ای کمتر به مردم معرفی کرد؟ و یا وظایف جمع‌سپاری را در محیطی فیزیکی توسط جمعیت با قابلیت اطمینان بالا انجام داد؟

در ادامه مقاله به بررسی و تحلیل موضوع خواهیم پرداخت. در بخش دوم پیشینه تحقیق بیان می‌شود. در بخش سوم، مدل پیشنهادی شرح داده خواهد شد و در بخش چهارم مقایسه روش پیشنهادی با سایر روش‌ها طرح می‌شود. مقایسه رویکردی جمع‌سپاری، بحث و بررسی و نتیجه‌گیری بخش‌های پایانی این مقاله هستند.

۲- پیشینه جمع‌سپاری

این مقاله ترکیب وظیفه‌های طبقه‌بندی مبتنی بر یادگیری عمیق با جمع‌سپاری است. لذا در این بخش ابتدا پیشینه‌ای از جمع‌سپاری بیان می‌شود و در ادامه پیشینه مقالات پایه یادگیری عمیق مورد استفاده در این پژوهش بیان می‌شود.

۲-۱- پیشینه جمع‌سپاری

به‌طور سنتی، اصطلاح "جمعیت" تا حدودی به‌طور انحصاری در زمینه افرادی که در اطراف یک هدف، احساس یا تجربه مشترک هستند، استفاده می‌شود. با این حال امروزه اغلب شرکت‌ها به بحث در مورد چگونگی جمع‌آوری افراد برای اهداف سازمان می‌پردازند. جمع‌سپاری در [5] به‌عنوان استفاده از فناوری‌های اطلاعاتی برای اعطای مسئولیت‌های تجاری به جمعیت مورد استفاده قرار می‌گیرد. جمع‌سپاری یک راه مناسب برای حل مسئله است.

هم‌گرایی فناوری‌ها این قابلیت را دارد تا چشم‌انداز جدیدی از همکاری را برای انجام امورات بشر باز کند. همانند جهش‌های قبلی در قدرت فناوری، هم‌گرایی محاسبات و ارتباطات اجتماعی، بهبود در زندگی را فراهم می‌کند. زیربنای اقتصادی و علمی محاسبات فراگیر، چند دهه است که ساخته شده است. جهان فیزیکی، اجتماعی و مجازی، در حال تصادم، ادغام و توازن است. در جهان آینده، اعمال و کنش‌های ارتباطات و مجامع، متفاوت خواهند بود.

با پیشرفت فناوری، برخی از فعالیت‌های اجتماعی و اقتصادی با زیرساخت‌های فناوری اطلاعات ترکیب شده‌اند. دو اصطلاح جمعیت هوشمند و جمع‌سپاری نیز در این راستا ایجاد و تعریف می‌شوند. جمعیت هوشمند و فلش موب^۱ که نخستین بار توسط هوارد رینگلد مطرح شد، به‌عنوان سازمان‌های اجتماعی خودساختار فناوری واسطه تعریف می‌شوند [1]. اصطلاح دیگر نیز جمع‌سپاری است که به عمل مشارکت گروهی مردم در به‌دست‌آوردن اطلاعات و انجام یک کار یا پروژه با استفاده از همکاری تعداد زیادی از مردم (با پرداخت یا بدون پرداخت) گفته می‌شود [2].

آیا می‌توان روشی را برای ایجاد یک جمع‌سپار پیشنهاد داد که به‌دلیل داشتن نوآوری و جذابیت، منتهی به جمعیت فعال مناسبی برای انجام فعالیت‌های جمع‌سپاری شود؟ مقاله [3] با مطالعات انجام‌شده به این نتیجه رسیده است که تنها دو مورد از نظریه‌های انگیزشی در مطالعاتشان در ایجاد جمع‌سپار بهره‌برده‌اند. همچنین مقاله [4] با بررسی مدل‌های نوآوری باز در حوزه جمع‌سپاری می‌نویسد: دانش و خلاقیت در نژاد بشر پراکنده شده‌اند و همیشه از توانایی‌های ما در استفاده از آن فراتر رفته‌اند. به‌طور سنتی، فرایندهای نوآورانه عمدتاً (البته نه منحصر) بر تلاش درون یک سازمان و یا همکاری بین سازمان‌های نسبتاً "محدود" متکی هستند. به‌تازگی، مدل‌های نوآورانه "باز" توجه بیشتری را از سوی پژوهش‌گران و متخصصان نوآوری به‌عنوان روش‌های

^۱ فلش موب‌ها نوع خاصی از جمعیت‌های هوشمند می‌باشند که بر خلاف سایر گروه‌ها، اهدافی تفریحی و غیر معمول از تشکیل جمعیت دارند. این درحالی است که سایر جمعیت‌های هوشمند افرادی هستند دارای اهداف محکم که جهت انجام کارهای گروهی مفید با استفاده از مزایای ارتباطات شبکه‌ای، گرد هم آمده‌اند.

می‌شود. برای این ایده‌ها هزینه هم پرداخت می‌کنند. تولیدکنندگان محصولات "فکر مشتری" را درک می‌کنند که این عمل بر مبنای ایده‌هایی است که نیازهای آنها را به‌طور مؤثرتر بررسی می‌کنند. در چالش بتا کاپ^۱، تنها ۷۰۰۰ یورو پاداش تعیین شد، این به غیر از هزینه‌هایی است که صرف تبلیغات و اطلاع‌رسانی این چالش شده است.

مقاله [3] آینده این پژوهش را مربوط به فعالیت‌هایی می‌داند که به بررسی عوامل پذیرش وظایف از طرف جمعیت دارند. استفاده مؤثر از جمعیت یکی از اهداف آینده این پژوهش‌ها است. این ایده قابلیت خواهد داشت تا در آینده با انجام یادگیری نیمه‌نظارت شده، به‌صورت مداوم به بهبود نتایج دسته‌بندی کمک کند. از جمله مسائل دیگری که می‌بایست در طرح‌های بعدی مورد لحاظ قرار بگیرد این است که با چه راه‌کارهایی می‌توان یک پایداری در روند مشارکت مردمی ایجاد کرد؟ پژوهش، طراحی و توسعه در دانشگاه بردلی (پوریا، نشان می‌دهد که یادگیری فعال یا یادگیری تجربی، یادگیری مبتنی بر بازی و گیمیفیکیشن، یادگیری مبتنی بر تجزیه و سایر یادگیری‌های مبتنی بر منابع انسانی و جمع‌سپاری می‌توانند کارایی بالایی داشته باشند [19].

مقاله [20] از چندین توانایی برچسب‌زن‌ها استفاده کرده است. این مقاله به بررسی روشی برای پیدا کردن متخصصان در برچسب‌زنی می‌پردازد و با کمترین نرخ خطا بر روی پنجاه هزار مجموعه داده تصویری برچسب‌زنی را انجام داده است.

در مقاله [21] جمع‌سپاری در سندرم حاد تنفسی حاد کرونا مورد استفاده قرار گرفته است. نویسندگان با استفاده از جمع‌سپاری در شبکه‌های اجتماعی تصویری ارزشمند از شیوع کرونا در زمان واقعی ارائه می‌دهند؛ اما پوشش جغرافیایی ناهمگن است. در این روش استخراج اطلاعات برای پیش‌بینی و اطلاع‌رسانی استراتژی‌های شیوع و پس از آن تکرار این روند در چرخه‌های تکرارشونده برای نظارت و ارزیابی پیشرفت، وجود دارد. الگوریتم‌هایی نیز برای اعتبارسنجی نتایج ایجاد شده‌اند. مقاله [22] نیز بحث قیمت‌گذاری پویا را مورد بررسی قرار داده است. این مقاله همچنین بحث تجمیع و ارسال مجدد داده‌ها را برای جمع‌سپاری بررسی کرده است.

° Betacup: این مسابقه برای طرح ایده‌ای جهت کاهش تعداد فنجان غیر قابل بازیافت بود که هر ساله با ایجاد یک جایگزین مناسب تر برای فنجان قهوه قابل استفاده مجدد ایجاد می‌شد.

در سال‌های گذشته پردازش زبان طبیعی به سمت پژوهش‌هایی مثل سم ایول^۱ و حاشیه‌نویسی رفته است. در همین اواخر با ظهور پلتفرم‌هایی مثل ای ام تی^۲ و جمعیت گل^۳ و به‌طور هم‌زمان رشد درصد استفاده از اینترنت، جمع‌سپاری یکی از روش‌های سریع و ارزان برای پژوهش‌های پردازش متن و حاشیه‌نویسی شده است. جمع‌سپاری برای انجام ماشینی کارها به کار می‌رود، که در آن کارگران بابت کارهای انجام داده، پول دریافت می‌کنند. جمع‌سپاری کاربرد فراوانی در تحلیل احساسات دارد. یک نمونه از این فعالیت‌ها، برچسب‌زنی داده‌هاست، که کاربران برحسب حس خود، داده‌ها را برچسب‌زنی می‌کنند. برخی از پروژه‌ها در جمع‌سپاری با چندین زبان اجرا می‌شوند.

جمع‌سپاری توجه زیادی را در دنیای کسب و کار به ارمغان آورده است و بسیاری از شرکت‌ها ارزش کسب و کار بالقوه خود را با جمع‌سپاری به‌دست آورده‌اند و کارزارهایی را برای این اهداف راه‌اندازی کرده‌اند. در همین حال، جمع‌سپاری منحصر به اهداف تجاری نیست [3]. نیاز به نوآوری مداوم یکی از مهم‌ترین اولویت‌های کسب و کار در میان مدیران اجرایی و یک مسأله کلیدی در پژوهش‌های علمی است. با توجه به نیاز به جریان مستمر نوآوری در محصولات و ارائه خدمات جدید، شرکت‌ها به مخترعان حرفه‌ای برای تولید ایده‌ها متکی هستند. در همان مقاله، ۳۳ پژوهش جمع‌سپاری مطالعه شده است، تنها برخی از مطالعات بر روی یک پلتفرم خاص تمرکز کرده‌اند، مابقی به نوآوری، ارزیابی، حل مسأله، مدل، طراحی و اندازه‌گیری پرداخته‌اند.

تارنمای کیوا^۴ [8] نمونه‌ای از جمعیت‌های هوشمند است که در پی آرمان‌های بشردوستانه پدید آمده است. در اینجا تمایل به برقراری ارتباط اساس تفکر پیوند مردم برای ایجاد روابط است. در این تارنما مردم می‌توانند به یکدیگر پول قرض بدهند و یا از یکدیگر قرض بگیرند. متوسط زمان برای پیدا کردن وام در تارنمای کیوا ۹۶ ساعت است و بازپرداخت آن به‌صورت یک‌ماهه، سه ماهه و یا در یک دوره زمانی محدود است.

به‌تازگی ایده‌هایی داده می‌شود که در بازاریابی از جمع‌سپاری استفاده شود [15]. ایده‌های بسیاری مطرح می‌شود، اما نبود بازاریابی مناسب منجر به شکست آنها

¹ SEMEVAL

² AMT

³ Cf(Crowd flower)

⁴ www.kiva.org - Loans that change lives

مقاله [17] در سال ۲۰۱۶ شرکت‌کنندگان در جمع‌سپاری را اغلب افراد با اعتبار کم و نتایج آن را با کیفیت پایین گزارش می‌کند؛ لذا در بسیاری از موارد کشف حقیقت در نتایج جمع‌سپاری حاصل نمی‌شود. برخی از ویژگی‌ها که فعالیت جمع‌سپاری پیشنهادی ما را نسبت به کارهای انجام‌شده در پلتفرم‌هایی مثل آمازون و ای.تی.ام متمایز می‌کند، عبارتند از: تضمین واقعی نتایج، گزارش در لحظه، استفاده از متخصصان واقعی، موقعیت جغرافیایی مشخص، کمک به انجام سریع، کمک به معرفی و ایجاد یک جمع‌سپار، ایجاد جریان تبلیغاتی برای فعالیت جمع‌سپاری، ایجاد ایده‌های جذاب و بومی‌سازی جمع‌سپاری.

مسئله‌ای که در مقاله [16] به آن اشاره شده است، این است که ایده جمع‌سپاری باید توجه رسانه‌ای را به خودش جلب کند. این مهم برای سمینار رقابت‌های استارت‌آپی اتفاق افتاد و اخبار آن در بیش‌تر رسانه‌ها، صدا و سیما و تارنماهای دانشگاه‌های معتبری مثل صنعتی شریف و امیرکبیر منتشر شد.

۲-۲- پیشینه تعبیه لغات

زیروظیفه یک، از وظیفه هفتم سم ایول در سال ۲۰۱۸ به بررسی استخراج رابطه پرداخته است. در پژوهش حاضر همان زیروظیفه با استفاده از جمع‌سپاری انجام شده است؛ لذا دو مقاله [11] و [12] در این بخش به‌عنوان مقالات پایه بحث طبقه‌بندی مورد مطالعه قرار می‌گیرند. سایر مقالات این زیروظیفه در مقایسه‌های انجام‌شده در بخش نتایج آمده‌اند.

مقاله پایه [11] با موضوع استخراج رابطه در مقالات علمی، وظیفه استخراج رابطه معنایی و طبقه‌بندی در چکیده مقاله‌های علمی را تشریح می‌کند. این چالش بر روابط معنایی با دامنه خاص متمرکز است و شامل سه زیرشاخه مختلف است. زیرشاخه‌ها به‌گونه‌ای طراحی شده بودند که تأثیر مراحل مختلف پیش‌پردازش بر نتایج طبقه‌بندی رابطه را با یکدیگر مقایسه کنند. انتظار می‌رود که این وظیفه برای طیف گسترده‌ای از پژوهش‌گرانی که در استخراج دانش تخصصی فعالیت دارند، مرتبط باشد.

هدف از [11] این است که به‌طور خودکار روابط معنایی را در یک مقاله از نشریات علمی شناسایی کند. در این کار، جفت‌موجودیت‌ها به‌عنوان نمونه‌های رده ایستا لحاظ می‌شوند. این وظیفه شامل شناسایی و طبقه‌بندی نمونه‌های روابط معنایی بین مفاهیم در مجموعه‌ای از

شش دسته گسسته است. این روابط خاص با حوزه علمی است و نمونه‌های آنها اغلب در چکیده یا مقدمه مقالات علمی یافت می‌شوند. این کار برای تهیه چارچوبی برای ارزیابی منظم مراحل لازم برای استخراج کامل اطلاعات از متن علمی انجام می‌شود.

مقاله پایه [12] نیز با موضوع طبقه‌بندی رابطه‌ای^۱، دو سامانه را برای طبقه‌بندی رابطه معنایی^۲ توصیف کرده است. این مقاله وظیفه طبقه‌بندی رابطه معنایی را اجرا کرده است. یک مدل SVM^۳ و یک مدل CNN^۴ ویژگی‌های متراکم^۵ وردتووک و ویژگی‌های پراکنده دستی را ترکیب می‌کنند. برای آموزش مدل‌ها، دو مجموعه داده ارائه شده برای زیرشاخه‌ها به‌منظور تعادل رده‌ها، ترکیب شدند. نتیجه این‌طور بود که عملکرد مدل SVM بهتر از CNN شد.

مقاله [6] به بررسی تجزیه و تحلیل احساسات اخبار با استفاده از جمع‌سپاری می‌پردازد. این مقاله رویکردی برای استخراج نمرات ۳۷۰۰۰ واژه داشته است. این رویکرد به‌وسیله تارنمای رپلر^۶ پیاده‌سازی شده است. این مقاله از یک روش بدون نظارت یادگیری استفاده می‌کند. بسیاری از روش‌های تجزیه و تحلیل احساسات از منابع واژگانی استفاده می‌کنند. برخی مقالات لغت را خارج از متن بررسی کرده و می‌گویند هر لغت دارای یک قطبیت مثبت و منفی است. در مقاله [7] داده‌ها با چندین رده (از خیلی منفی تا خیلی مثبت) رده‌بندی شده‌اند. که رتبه‌بندی این داده‌ها به‌وسیله جمع‌سپاری انجام شده است.

۳- روش پژوهش

استخراج رابطه، کاری برای استخراج اطلاعات است که هدف آن شناسایی و طبقه‌بندی روابط معنایی بین موجودیت‌ها در متن است. این کار به‌طوراساسی ساختار را از اطلاعات متنی بدون ساختار استخراج می‌کند و به ما این امکان را می‌دهد تا اطلاعات ارزشمندی را در مورد نحوه تعامل موجودیت‌ها به‌دست آوریم. در نتیجه ظرفیت انسانی را برای تجزیه و تحلیل (اغلب زیاد) داده‌های متنی بهبود می‌بخشد. در استخراج، جفت‌موجودیت‌ها^۷ از یک سند بررسی و نوع رابطه پیش‌بینی می‌شود.

¹ relation classification

² semantic relation classification

^۳ ماشین بردار پشتیبان

⁴ convolutional neural network

⁵ word2vec

⁶ rappler.com

⁷ pairs of entities

ممکن است در جملات، متن چندکلمه‌ای وجود داشته باشد که در یک تعبیه یا در هر دو مدل از قبل آموزش دیده نشده باشند. یک بردار صفر برای این کلمات استفاده می‌شود. برای هر مدل تعبیه‌ای که آنها را پوشش نمی‌دهد، این بردار لحاظ می‌شود. آخرین عملیاتی که در مورد بردارهای تعبیه‌شده انجام می‌دهیم، نرمال کردن آنها به مقادیر مشابه است.

۳-۲- نوآوری مدل پیشنهادی

در مرحله بعد از آماده‌سازی داده‌ها، زیرمجموعه‌ای کارآمد از ویژگی‌های ورودی را با استفاده از BPSO³ انتخاب می‌کنیم. برای اعمال BPSO، باید یک تابع برازش⁴ تعریف و به‌طور مشابه مقادیر مناسبی را برای هاپر پارامتر الگوریتم تعریف کنیم. نمره F به‌دست می‌آید که با استفاده از الگوی آموزش دیده با هر زیرمجموعه انتخابی از ورودی‌ها در مجموعه داده‌های توسعه به‌عنوان مقدار برازش برای BPSO در نظر گرفته می‌شود.

بهینه‌سازی ازدحام ذرات بدووی یا همان BPSO به یک تابع برازش⁵ نیاز دارد تا میزان برتری راه حل انتخابی را به‌دست بیاورد؛ لذا می‌بایست یک طبقه‌بند با ویژگی‌های انتخاب‌شده برای هر بردار ایجاد شود. این طبقه‌بند بر روی داده توسعه‌ای به‌وسیله SVM اعمال می‌شود. این در حالی است که در سایر مدل‌ها با CNN انجام می‌شد، که اگرچه دقت بیشتری داشت، ولی زمان بیشتری برای به نتیجه رسیدن نیاز داشت. طول دوره آموزش در CNN زیاد بود؛ لذا طبقه‌بند به SVM در BPSO تغییر داده شد؛ لذا با توجه به تعداد بسیار زیاد پارامترها لازم است که مدل برای CNN آموزش داده شود و تصمیم گرفتیم از یک فرآیند انتخاب ورودی نیز به جای استفاده از CNN استفاده کنیم. دلایل اصلی انتخاب SVM عبارتند از: تعداد بسیار کم هاپر پارامترهای قابل تنظیم و توانایی آن برای کار با تعداد زیادی از ویژگی‌های ورودی نیز بالا است. گفتنی است که ما به هاپر پارامترهای SVM⁶ اهمیت نمی‌دهیم و مقادیر پیش فرض را بدون تغییر نگه می‌داریم.

³ binary particle swarm optimization

⁴ برای اینکه خوب یا بد بودن جواب را تعیین Fitness Function کنیم، از مفهومی به نام تابع برازش استفاده می‌کنیم. اگر تابع نتواند میزان خوب بودن جواب را به‌درستی نشان دهد، الگوریتم معیار ارزیابی دقت عملکرد نخواهد داشت.

⁵ fitness function

⁶ support vector machine

مطالب علمی با ساختارهای پیچیده و تکیه بر اصطلاحات تخصصی پیش‌بینی نوع رابطه صحیح¹ را دشوارتر می‌کنند. در پژوهش پیش روی مشکل تشخیص رابطه به‌عنوان یک کار طبقه‌بندی با ناظر کار شده است. در اینجا ناظر جمع‌سپاری است.

۳-۱- مدل طرح پیشنهادی

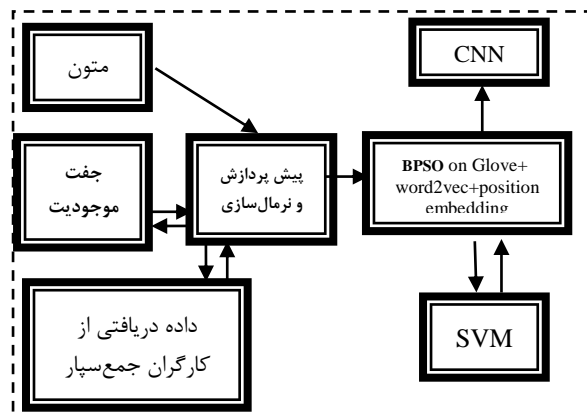
دو مدل تحلیل شده‌اند. یک ماشین بردار پشتیبانی (SVM)، که از مجموعه‌ای غنی از ویژگی‌ها در ترکیب ویژگی‌های متراکم وردتووک و ویژگی‌های پراکنده دستی بهره برده است و در روشی دیگر از یک مدل عصبی پیچشی (CNN) استفاده می‌شود. بهترین نتیجه با مدل SVM به‌دست آمد.

روش پیشنهادی ما برای استخراج رابطه معنایی شامل سه مرحله اساسی است: تهیه داده‌ها، تصمیم‌گیری در مورد ورودی‌ها و آموزش مدل CNN. شکل (۱) نمای کلی مدل پیشنهادی را نشان می‌دهد. ما ابتدا بر روی اسناد ورودی کار می‌کنیم تا برخی از اقدامات پیش پردازش از جمله نشانه‌گذاری را انجام دهیم. ما به دنبال بردارهای تعبیه لغات هستیم؛ لذا از مدل‌های پیش‌آموزش داده‌شده وردتووک و گلاو²[28] بر روی مجموعه‌متن‌های بزرگ استفاده کرده‌ایم. همچنین تعبیه موقعیت برای هر نشان در یک جمله معین می‌شود.

در پایان مرحله پیش‌پردازش و در مرحله آموزش، هر سه بردار جاسازی‌شده با هم جمع می‌شوند تا یک ماتریس ورودی اصلی با اندازه $(e1 + e2 + d) * n$ ایجاد شود که در آن n بیشینه طول جملات نشانه‌گذاری شده است و e1 و e2 به ترتیب طول بردارهای تعبیه‌شده به‌ازای هر نشانه از وردتووک و گلاو هستند. طول تعبیه موقعیت برای هر کلمه به‌وسیله d نشان داده شده است.

(شکل-۱): مدل پیشنهادی برای طبقه‌بندی موجودیت‌ها

(Figure-1): Proposed model for classifying entities



استفاده از پتانسیل
اتحادیه‌های شغلی

رویداد جمعیت هوشمند در هر جلسه یک حوزه شغلی را مورد مطالعه قرار داد. اتحادیه‌ها و اصناف کمک به‌سزایی در ارسال پیامک برای مخاطبین حوزه‌های شغلی خود داشتند و بدین شکل جمعیت هوشمند هر جلسه تا حدی نیز تخصصی و با حضور اعضاء آن صنف تشکیل شد. ارسال پیامک تبلیغاتی به گروه‌های هدف نیز تأثیر به‌سزایی در افزایش جمعیت همگون داشت. برای مثال در رقابت‌های استارت آپی املاک به ۱۳ هزار نفر از مشاوران املاک استان تهران پیامک ارسال شد.

ایجاد حس نوآوری

مردم علاقه‌مند به کسب اطلاعات نوین در حوزه کسب و کار هستند. از طرفی میل ذاتی انسان به فراگیری، کشف و تحصیل رضایت می‌تواند با تحول و رشد جمعیت‌های هوشمند در دهه اخیر مرتبط باشد. میل به کسب دانش و اطلاعات نوین نقش مؤثری در تشکیل این جمعیت هوشمند داشت. با توجه به اینکه ذات استارت‌آپ‌ها، شامل مفاهیم خلاقیت و ایده‌پردازی است، استفاده از آنها کمک به‌سزایی در تشکیل جمعیت دارد.

انجام تبلیغات انگیزشی
صوتی/ تصویری

ایجاد پادکست‌های تبلیغاتی از استارت آپ‌ها و ساخت کلیپ تبلیغاتی از استارت‌آپ‌ها

اجرا در محیط دانشگاهی

بنابر مقاله [16] از مطالعه میان ۵۵ نشریه، ۳۵ مقاله (۶۴٪) تنها توسط دانشگاهیان، ۱۰ مقاله (۱۸٪) تنها از صنعت و ۱۰ مقاله (۱۸٪) توسط همکاری دانشگاهیان و صنعت نوشته شده است. این نشان می‌دهد که دانشگاهی‌ها نقش مهمی در تحقق اهداف پژوهش دارند.

در طرح پیش رو جمعیتی حدود ۵۵۰ نفر از صاحبان کسب و کارها، مدیران کانال‌ها و گروه‌ها در شبکه‌های اجتماعی و علاقه‌مندان به توسعه کسب و کار در همایش رایگان رقابت‌های استارت آپی حاضر شدند، این جمعیت در سمینار به‌طور اشتراکی موجودیت‌های صد چکیده مطلب را برچسب‌زنی کردند.

در طرح پیشنهادی جمعیتی شامل فارغ‌التحصیل و دانشجو در بیش از پنجاه دانشگاه در فرآیند جمع‌سپاری شرکت کرده بودند، که حدود سی درصد از دانشگاه‌های دولتی ایران و نزدیک به یک درصد جمعیت نیز از دانشگاه‌های خارج از ایران بودند. بیش از هشتاد درصد شرکت‌کنندگان تحصیلات دانشگاهی داشتند و در حدود ۳۸ درصد از شرکت‌کنندگان در همایش تحصیلات کارشناسی ارشد به بالا داشتند.

۴- سناریوی پیشنهادی تعبیه برای استخراج رابطه با داده جمع‌سپار

سناریوهای مختلفی برای فعال کردن جمع‌سپاری می‌توان نوشت؛ به‌عنوان مثال در [10] بر اساس سناریو پاسخ‌های

کم و پاداش کم پژوهش انجام گرفته است. نکته‌ای که می‌بایست در فعالیت‌های جمع‌سپاری به آن توجه شود این است که مردم در انجام فعالیت‌های جمع‌سپاری احساس خستگی نکنند. بهترین سناریوها، آنهایی خواهند بود که از تأثیرگذاری مردم بر مردم در جهت انجام وظایف کمک بگیرند و در راستای فعالیت‌های مثبت و ارزنده اجتماعی تلقی شوند. سناریونویسی از مهم‌ترین بخش‌های مدیریت جمعیت‌های هوشمند است.

در این پژوهش با استفاده از اطلاع‌رسانی پیامکی، تبلیغات در شبکه‌های اجتماعی و قدرت خبری رسانه‌ای، یک رویداد کسب و کار به نام "رقابت‌های استارت آپی" به‌وسیله "مجله صدای زنده" در دانشگاه شمس‌پور تهران تشکیل شد. در هر جلسه از این رویداد، یک حوزه کسب و کاری مورد مطالعه قرار می‌گیرد. با توجه به مطالعات انجام‌شده، در کشور ایران استارت آپ‌ها در حوزه‌های مختلف شغلی فعالیت می‌کنند که از این بین در ۴۲ مورد از آن حوزه‌ها، رقابتی استارت آپی وجود دارند؛ به‌عنوان مثال ۵۶۰ تاکسی برخط، ۵ آزمایشگاه پزشکی برخط، ۱۴ معماری برخط و ۱۲ خشک‌شویی آنلاین در ایران مشغول به فعالیت هستند. این رویداد پس از یازده جلسه به جهت برندسازی به دانشگاه امیرکبیر منتقل شد.

در سناریوی پیشنهادی تصمیم گرفتیم تا برای ایجاد جمعیتی با سطح سواد بالا جهت مطالعات متن‌کاوی آینده، همایش آموزشی رایگان رقابت‌های استارت آپی را با حضور مردم برگزار کنیم. بنا شد تا در زمان برگزاری این همایش از مردم در جهت فعالیت جمع‌سپاری کمک بگیریم. هزینه‌ای به‌عنوان جایزه معادل صد هزار تومان برای دو نفر به‌طور قرعه در نظر گرفته شد. هزینه‌های اجرایی برگزاری جلسه و تبلیغات برای جمع‌آوری این جمعیت نیز از جمله هزینه‌های این فرآیند جمع‌سپاری بود.

۱-۴- وظیفه رده‌بندی در تعبیه لغات ترکیبی

به‌وسیله فرآیند جمع‌سپار

مقاله [11] با استفاده از شبکه عصبی پیچشی^۱ به‌عنوان یک وظیفه، رده‌بندی را برای شش رده کاربردی، نتیجه‌گیری^۲، جزء-کل^۱، موضوع^۲، مدل-ویژگی^۲ و

¹ Convolutional Neural Network

² USAGE

³ RESULT

مقایسه^۴ تشخیص می‌دهد. این مقاله داده مورد نیاز خود را از چکیده ۳۵۰ مقاله علمی استخراج کرده است. از این تعداد چکیده در جمع ۸۳۵ رابطه استخراج شد.

در جدول (۲) این رده‌ها بر اساس تعداد در هر رده و همچنین نوع رابطه مشخص شده‌اند. به‌عنوان مثال جمله "حافظه فیزیکی بخشی از رایانه است" یک رابطه جزء به کل عادی است و جمله "رایانه دارای مادربورد است" یک رابطه جزء به کل معکوس است. ترتیب آمدن موجودیت‌ها در جمله، جزء به کل عادی یا جزء به کل معکوس بودن را مشخص می‌کنند.

(جدول-۲): رده توزیع روی داده‌های یادگیری [11]
(Table-2): Distribution class on learning data [11]

نوع رابطه	عادی	معکوس	کل
کاربرد	296	187	483
موضوع	8	10	18
مقایسه	95	-	95
مدل-ویژگی	226	100	326
نتیجه‌گیری	52	20	72
جزء-کل	158	76	234
مجموع	835	393	1228

(جدول-۳): نتایج تفصیلی مبتنی بر اندازه‌گیری در مرحله انتخاب ویژگی ورودی. هر ردیف طبقه‌بند آموزش داده شده با SVM را با استفاده از ویژگی‌های ورودی یادشده نشان می‌دهد.

(Table-3): Detailed measurement-based results in the input feature selection step. Each row shows SVM-trained classifiers using the listed input features.

طول ورودی	امتیاز ^۶ F	فراخوانی ^۵	دقت ^۵	W2vec
300	68.9	69.6	68.4	W2vec
300	67.7	68.4	67.2	گلاو
50	22.3	23.6	21.3	تعبیه موقعیت
600	70.8	71.6	70.2	W2vec + گلاو
650	70.4	71.6	69.3	W2vec + گلاو + تعبیه موقعیت
368	74.7	76.2	73.4	ورودی های نخبه

گلاو^۸ و وردتووک^۹ دو مورد از محبوب‌ترین الگوریتم‌های تعبیه لغات هستند. در کاربردهای مختلف

یادگیری ماشین از این دو الگوریتم به‌طور مجزا استفاده شده است. ما با استفاده از مقاله [12] و با ترکیب گلاو و وردتووک بردار تعبیه‌ای جدید برای لغات ایجاد کرده‌ایم. این تغییر در جهت تکمیل وظیفه سم ایول سال ۲۰۱۸ در جهت ارزیابی با هدف کشف ماهیت معانی نوشته شده است، ابتدا دو بردار گلاو و وردتووک، هر یک به طول سیصد برای کلیه کلمات استخراج شده و سپس با اضافه‌شدن پنجاه بیت برای تعبیه موقعیت لغات، کلیه لغات به‌صورت برداری به‌طول ۶۵۰ تبدیل می‌شوند؛ در نهایت جملات براساس ترکیب لغات ساخته می‌شوند. نتایج اولیه تغییرات مدل (الگوریتم) در جدول (۳) نشان داده شده است.

برای تغییر در مقاله [12] از کدهای CNN و SVM همراه با BPSO استفاده شده است؛ به‌نحوی که SVM برای برازش تابع BPSO خواهد بود. جدول (۳) نتیجه استفاده از ورودی‌های نخبه را نشان می‌دهد.

حال می‌بایست از ورودی‌های نخبه در استخراج رابطه استفاده کنیم. از نتایج به‌دست‌آمده حاصل از بردارهای لغات نخبه، زوج موجودیت‌های متن مطابق جدول (۲) با CNN پیش‌بینی می‌شوند. بهبود در جدول (۴) آمده است.

(جدول-۴): مقایسه رده‌بندی‌های مبتنی بر تلفیق و داده

استاندارد با حالت اضافه‌شدن داده جمع‌سپاری

(Table-4): Comparison of standard based classifications and data aggregation with additional data supplied by crowdsourcing

دقت ^{۱۰}	فراخوانی ^{۱۱}	امتیاز ^{۱۲} F	رده‌بندی در
79.2	84.4	81.7	ETH-DS3Lab [13]
-	-	78.9	UWNLP [14]
-	-	76.7	SIRUS_LTG_UiO [23]
-	-	74.9	ClaiRE [24]
-	-	74.2	Talla [25]
-	-	72.7	MIT-MEDG [26]
64.64	75.57	69.7	TakeLab [27]
70.20	72.80	71.20	تعبیه مبتنی بر گلاو و وردتووک با ۳۵۰ چکیده اولیه و ۱۰۰ چکیده جمع‌سپار

۵- نتیجه‌گیری

جمعیت هوشمند می‌تواند متخصص یا غیرمتخصص باشد. در جمعیت هوشمند غیر متخصص افراد از رده‌های سنی مختلفی حضور دارند. فعالیت حاضر با جمعیت

¹⁰ Precision

¹¹ Recall

¹² F-score

¹ PART-WHOLE

² TOPIC

³ MODEL-FEATURE

⁴ COMPARE

⁵ Precision

⁶ Recall

⁷ F-score

⁸ <https://nlp.stanford.edu/projects/glove/>

⁹ <http://code.google.com/archive/p/word2vec>

- ACM conference on Online social networks, 2013, pp. 27-38.
- [8] P. Belleflamme, T. Lambert, & A. Schwienbacher, "Crowdfunding: Tapping the right crowd", *Journal of business venturing*, vol. 29(5), pp. 585-609, 2014.
- [9] J. Daniels, & J. R. Feagin, "The (coming) social media revolution in the academy," *Fast Capitalism*, vol.8(2), 2019.
- [10] T. A. Gautre, & T. H. Khan, "An analysis of question answering system for education empowered by crowdsourcing" In 2018 2nd International Conference on Inventive Systems and Control (ICISC), IEEE, 2018.
- [11] K. Gábor, D. Buscaldi, A. K. Schumann, B. QasemiZadeh, H. Zargayouna, & T. Charnois, "Semeval-2018 Task 7: Semantic relation extraction and classification in scientific papers," In Proceedings of The 12th International Workshop on Semantic Evaluation, pp. 679-688, 2018.
- [12] M. Gluhak, M. P. di Buono, A. Akkasi, & J. Šnajder, "TakeLab at SemEval-2018 Task 7: Combining Sparse and Dense Features for Relation Classification in Scientific Texts", In Proceedings of The 12th International Workshop on Semantic Evaluation, pp. 842-847, 2018.
- [13] J. Rotsztein, N. Hollenstein, & C. Zhang, Eths3lab at semeval-2018 task 7: Effectively combining recurrent and convolutional neural networks for relation classification and extraction. arXiv preprint arXiv:1804.02042, 2018.
- [14] Y. Luan, M. Ostendorf, & H. Hajishirzi, "The unlp system at semeval-2018 task 7: Neural relation extraction model with selectively incorporated concept embeddings", In Proceedings of The 12th International Workshop on Semantic Evaluation, pp. 788-792, 2018, June.
- [15] S. A. Lazarus, *Cyber Mobs: A Model for Improving Protections for Internet Users* (Doctoral dissertation, Utica College), 2017.
- [16] B. L. Bayus, "Crowdsourcing new product ideas over time: An analysis of the Dell IdeaStorm community", *Management science*, vol. 59(1), pp. 226-244, 2013.
- [17] R. W. Ouyang, M. Srivastava, A. Toniolo, & T. J. Norman, "Truth discovery in crowdsourced detection of spatial events", IEEE, 2016.
- [18] B. Xiang, "The psychological effects of participation in crowdsourcing on customer's willingness to pay and recommend a brand, 2016.
- به‌طور تقریبی متخصص انجام شد. همچنین می‌تواند از اعضای داخل سازمان یا غیر وابسته به سازمان بیشتر جمعیت حاضر از افراد غیر وابسته به سازمان بودند. همچنین نوع فعالیت می‌توانست فعالیت وابسته به یکدیگر در جمعیت یا فعالیت مستقل باشد در این پژوهش راه حلی نوین برای بومی‌سازی جمع‌آوری داده‌ها و کمک به مسائل داده‌کاوی با کمترین هزینه و کمک جمعیت متخصص ارائه شد. استراتژی‌های ابتکاری آموزش و یادگیری، که از فن‌آوری کمک می‌گیرند به راه‌کارهای نوینی در تجزیه و تحلیل داده تبدیل شده‌اند.
- در این مقاله با استفاده از BPSO و SVM بردار تعبیه بهینه اشتراکی گلاو و ورد ۲۰ک انتخاب شد. در پژوهش‌های قبلی این بردار مبتنی بر PSO و CNN بود؛ سپس با استفاده از CNN استخراج ویژگی موجودیت‌ها در چکیده متون علمی انجام شد، در این پژوهش با کمک ۴۵۰ داده (صد مجموعه داده جمع‌سپار به‌همراه ۳۵۰ مجموعه استاندارد قبلی) دقت به‌نحوی افزایش پیدا کرد که روش پیشنهادی ضمن افزایش سرعت، در بین پنج روش برتر قرار گرفت.

8- References

۸- مراجع

- [1] H. Rheingold, *Smart mobs: The next social revolution*. Basic books, 2007.
- [2] D. Zhou, Q. Liu, J. C. Platt, C. Meek, & N. B. Shah, Regularized minimax conditional entropy for crowdsourcing. arXiv preprint arXiv:1503.07240, 2015.
- [3] Y. Zhao, & Q. Zhu, "Evaluation on crowdsourcing research: Current status and future direction", *Information Systems Frontiers*, vol. 16(3), pp. 417-434, 2014.
- [4] S. Marjanovic, C. Fry, & J. Chataway, "Crowdsourcing based business models: In search of evidence for innovation 2.0", *Science and public policy*, vol. 39(3), pp. 318-332, 2012.
- [5] J. Prpić, P. P. Shukla, J. H. Kietzmann, & I. P. McCarthy, "How to work a crowd: Developing crowd capital through crowdsourcing", *Business Horizons*, vol. 58(1), pp. 77-85, 2015.
- [6] J. Staiano and M. Guerini, "DepecheMood: a Lexicon for emotion analysis from crowd-annotated news," *arXiv preprint arXiv1405*, pp. 1605, 2014.
- [7] P. Gonçalves, M. Araújo, F. Benevenuto, and M. Cha, "Comparing and combining sentiment analysis methods," in Proceedings of the first



محمد جعفرآباد مدرک کارشناسی و کارشناسی ارشد خود را در رشته مهندسی رایانه گرایش معماری سامانه های رایانه دریافت کرده و هم‌اکنون دانشجوی مقطع دکترای دانشگاه قم در رشته مهندسی فناوری اطلاعات است. ایشان در زمینه‌های داده‌کاوی و رمزنگاری مطالعاتی داشته است.

jafarabadm@yahoo.com



روح‌الله دیانت فارغ‌التحصیل دکترای معماری سامانه‌های رایانه‌ای دانشگاه شریف است. ایشان هیأت علمی دانشگاه قم هستند و پژوهش‌هایی را در زمینه داده‌کاوی و پردازش صوت و تصویر داشته‌اند.

rouhollah.dianat@gmail.com

- [19] M. A. Rashid, K. Deo, D. Prasad, K. Singh, S. Chand, & M. Assaf, TEduChain: A platform for crowdsourcing tertiary education fund using blockchain technology. arXiv preprint arXiv:1901.06327, 2019.
- [20] P. Welinder, & P. Perona, "Online crowdsourcing: rating annotators and obtaining cost-effective labels", In 2010 IEEE Computer Society Conference on Computer Vision and Pattern Recognition-Workshops pp. 25-32, IEEE, 2010.
- [21] G. M. Leung, & K. Leung, "Crowdsourcing data to mitigate epidemics", *The Lancet Digital Health*, vol. 2(4), e156-e157, 2020.
- [22] A. Druksa, V. Fedorova, D. Ustalov, O. Megorskaya, E. Zermirnova, & D. Baidakova, "Crowdsourcing Practice for Efficient Data Labeling: Aggregation, Incremental Relabeling, and Pricing", In *Proceedings of the 2020 ACM SIGMOD International Conference on Management of Data*, pp. 2623-2627, 2020.
- [23] F. Nooralahzadeh, & L. Øvrelid, "Syntactic dependency representations in neural relation classification", arXiv preprint arXiv:1805.11461, 2018.
- [24] L. Hettinger, A. Dallmann, A. Zehe, T. Niebler, & A. Hotho, "Claire at semeval-2018 task 7: Classification of relations using embeddings", In *Proceedings of The 12th International Workshop on Semantic Evaluation*, pp. 836-841, 2018.
- [25] B. Pratap, D. Shank, O. Ositelu, & B. Galbraith, "Talla at SemEval-2018 task 7: Hybrid loss optimization for relation classification using convolutional neural networks", In *Proceedings of The 12th International Workshop on Semantic Evaluation*, pp. 863-867, 2018.
- [26] D. Jin, F. Dernoncourt, E. Sergeeva, M. McDermott, & G. Chauhan, "MIT-MEDG at SemEval-2018 task 7: Semantic relation classification via convolution neural network", In *Proceedings of The 12th International Workshop on Semantic Evaluation*, pp. 798-804, 2018.
- [27] M. Gluhak, M. P. di Buono, A. Akkasi, & J. Šnajder, "TakeLab at SemEval-2018 Task 7: Combining Sparse and Dense Features for Relation Classification in Scientific Texts, In *Proceedings of The 12th International Workshop on Semantic Evaluation*, pp. 842-847, 2018.
- [28] J. Pennington, R. Socher, & C. D. Manning, "Glove: Global vectors for word representation", In *Proceedings of the 2014 conference on empirical methods in natural language processing (EMNLP)*, pp. 1532-1543, 2014.