



به کارگیری روش متن کاوی به منظور دسته بندی سلسله مراتبی متون با رویکرد یادگیری عمیق در داده های حجیم

بهناز اسلامی^۱، مهدی حبیب زاده مطلق^{۲*}، مجید فولادیان^۳ و زهرا رضایی^۴

^۱ گروه مهندسی کامپیوتر، دانشکده برق و کامپیوتر، واحد علوم تحقیقات، دانشگاه آزاد اسلامی، تهران، ایران،
behnazeslami30@gmail.com

^۲ استادیار گروه کامپیوتر، دانشگاه ایوانکی، سمنان، ایران، me_habi@encs.concordia.ca (نویسنده مسئول)

^۳ استادیار گروه برق و الکترونیک، دانشگاه آزاد اسلامی واحد ساوه، ساوه، ایران، majidfouladian@gmail.com

^۴ گروه مهندسی کامپیوتر، دانشکده برق و کامپیوتر، دانشگاه کاشان، اصفهان، ایران، z.rezaei2010@gmail.com

چکیده - همواره تجزیه و تحلیل متن یک روش برای استخراج دانش از متن است که جهت بازیابی دانش نهان در متون، محدودیت حافظه و زمان در پردازش داده های بزرگ به دلیل منابع داده توزیع شده در وب، موتورهای جستجو و شبکه های اجتماعی حائز اهمیت می باشد. با بهبود داده ها در اینترنت و شبکه های اجتماعی، روش های قوی تر مورد نیاز است که می تواند اطلاعات را از نزدیک و بلافاصله دسته بندی کند. از طریق انتزاع در داده های متنی، یادگیری عمیق می تواند با این چالش ها مقابله کند. در این مقاله روش یادگیری عمیق جهت دسته بندی به صورت سلسله مراتبی مورد بررسی قرار می گیرد و نتایج نشان می دهند که این روش می تواند متون شبکه های اجتماعی و وبسایت ها را با دقت ۹۸٫۸۱٪ دسته بندی کند. این مقاله همچنین نشان می دهد که این شبکه پیچیده می تواند در سطح کلمه بهتر کار کند و نیازی به دانستن ساختار نحوی و معنایی زبان ندارد. نوآوری روش پیشنهادی، اضافه کردن یک سطح به ساختار سلسله مراتبی برای تشخیص عمومی و تشخیص دقیق تر کلاس است.

کلیدواژه - پردازش زبان طبیعی، داده های حجیم، دسته بندی، متن کاوی، یادگیری عمیق

که باید اشاره کرد این است که اکثر ابزار تجزیه و تحلیل متن از پردازش خطی و خطی پشتیبانی نمی کند.

۱- مقدمه

علاقه به پردازش داده ها در دنیای امروزی توجه زیادی به مدل های داده بزرگ مانند Apache Spark [۳]، Apache Storm، [۴] و Apache Flink [۵] را به خود جلب کرده است. ابزارهای دیگر مانند Azure ML [۶] و SAMOA [۷] نیز از جمله ابزارهای بسیار کارآمد هستند که Azure ML تنها برای خدمات ابری ارائه می شود؛ که کاربران را مسدود می کند و از سوی دیگر، SAMOA فقط از داده های بلادرنگ پشتیبانی می کند و نمی تواند برای همه برنامه ها اعمال شود.

به طور کلی متن کاوی شامل بخش های مختلفی می شود که برخی از آن ها عبارتند از: مدل سازی موضوع [۸]، دسته بندی احساسات [۹] و تشخیص هرزنامه [۱۰]. در گذشته با استفاده از روش های سنتی دسته بندی متن، ابتدا متن ها به صورت نمودارهای n-gram و ویژگی های لغوی نمایش داده می شدند،

دسته بندی متن به عنوان یکی از روش های اصلی در پردازش طبیعی فرآیندهای زبان در نظر گرفته شده است و هدف این است که یک متن خاص را برچسب گذاری کند. پردازش دستی و دسته بندی داده ها برای کاربران انسان غیر عملی است، زیرا اطلاعات زیادی وجود دارد [۱].

به علاوه تهیه مجموعه از متن هایی که برای آموزش چنین مدلی استفاده می شود، به عنوان یک چالش در نظر گرفته می شود. در این میان، کارشناسان موضوعی قادر هستند تا اسناد را تشخیص داده و عمل برچسب گذاری را انجام داده و مجموعه آموزشی را ایجاد کنند. انجام این فرآیند، زمان و انرژی زیادی را صرف می کند. علاوه بر این، تعداد اسناد متنی مورد نیاز برای ساخت یک مدل ناشناخته است، در حالی که ایده کلی وجود اسناد متنی بیشتر جهت ارائه مدل بهتر می باشد [۲]. نکته دیگر



در این مقاله، شبکه‌های عصبی کانولوشن برای استفاده از استخراج ویژگی از متن استفاده شده است. در [۱۳] برخی از ویژگی‌های داده‌های متنی از طریق ترسیم شبکه عصبی کانولوشن ساخته شده‌اند و این ویژگی‌ها با استفاده از یک بردار ۳۶۰ درجه‌ای که شامل بردار کلمه و نوع ساختار کلمه می‌باشد نشان داده شده‌اند. در این مقاله، خروجی انتهایی یک لایه کاملاً متصل در شبکه عصبی کانولوشن همراه با سایر طبقه‌بندی‌ها استفاده شده است. همچنین در این مقاله مشخص شد که ویژگی‌های ایجاد شده توسط شبکه، دسته‌بندی بهتر داده‌های آموزشی را از طریق سیستم‌های مختلف یادگیری ارائه می‌دهد. علاوه بر این مشخص گردید، ویژگی‌های ایجاد شده بهتر است جهت دسته‌بندی به لایه آخر شبکه عصبی متکی باشند. در [۱۴] نشان داده شده که دقت دسته‌بندی متون با استفاده از تکنیک‌های پیشرفته پردازش زبان طبیعی به‌طور قابل توجهی بهبود بخشیده شده است. در این مطالعه، گروهی از نظرات پزشکی در مورد یک دسته از داروها برای بررسی اثرات جانبی داروها از طریق ویژگی‌های مفید آن‌ها مورد بررسی قرار گرفت. نتایج نشان داد که انتخاب دقیق مشخصات و ویژگی‌ها می‌تواند به‌طور خودکار نتایج نهایی برخی از روش‌های دسته‌بندی کلاسیک از جمله بیزین و ماشین بردار پشتیبان را بهبود بخشد.

۳- روش پیشنهادی

روش پیشنهادی شامل دو مرحله است: جمع آموری داده‌ها جهت ساخت مدل و پیش‌بینی بلادرنگ متون. در این بخش، این مراحل شرح داده می‌شوند.

۳-۱- جمع آموری داده‌ها

با استفاده از خزش گر Scrapy متون مختلف از وبسایت‌های خبری CNN و رویترز در ۱۱ دسته در بازه زمانی سال ۲۰۱۶ تا سال ۲۰۱۷ بر اساس موضوعات مختلف جمع‌آوری شدند.

سپس یک مدل خطی یا روش هسته بر روی این نمودارها اعمال می‌گردید [۸]. اما اخیراً ارائه متنوع متون با استفاده از روش‌های یادگیری عمیق همانند شبکه‌های عصبی کانولوشن [۱۱] یا بر اساس شبکه‌های عصبی مکرر که خودشان بر اساس حافظه کوتاه‌مدت [۱۲] ارائه می‌شوند ارائه شده است. دسته‌بندی و پیش‌بینی دو عملیات برای تجزیه و تحلیل داده‌ها و مشتق مدل‌ها باهدف توصیف دسته‌های اصلی داده‌ها، درک آن‌ها و پیش‌بینی رفتار آینده آن‌هاست. به‌طور کلی، دسته‌بندی داده‌ها یک فرآیند دو مرحله‌ای است. مرحله اول، ساخت مدل است و مرحله دوم، استفاده از مدل و پیش‌بینی از طریق داده‌های قبلی است. دسته‌بندی اسناد به این معنی است که دسته‌بندی این اسناد و متون را به گروه‌هایی با ویژگی‌های مشترک، به‌طوری که بعداً می‌توان از این اسناد بر اساس این نقطه مشترک استفاده کرد.

از آنجاکه اسناد ممکن است بر اساس محتوا یا بر اساس ویژگی دیگری مانند تاریخ اسناد دسته‌بندی شوند. در دسته‌بندی محتوا اسناد، ابتدا برخی اسناد با موضوع پیش تعیین شده به‌عنوان نمونه برای یادگیری به دستگاه داده می‌شوند. سپس، بر اساس کلمات در هر سند، شبکه یاد می‌گیرد و با توجه به احتمال کلمات در هر سند، دستگاه موضوع سایر اسناد را پیش‌بینی می‌کند. سپس، با استفاده از یادگیری عمیق دسته‌بندی به‌صورت سلسله مراتبی انجام شده و نتایج ارزیابی خواهد شد. نکته مهم در این مقاله معرفی این است که دسته‌بندی فارغ از در نظر گرفتن چالش‌های داده‌های بزرگ و گفتار محاوره‌ای انجام می‌پذیرد.

مراحل کلی این مقاله به شرح زیر است:

- جمع‌آوری داده‌ها
- پیش‌پردازش
- دسته‌بندی داده‌های حجیم به روش سلسله مراتبی

۲- مطالعات قبلی

در [۱۲] نشان داده شده است که آموزش شبکه عصبی کانولوشن بر بردارهای کلمات که قبلاً ایجاد شده است، برای دسته‌بندی جمله بسیار مؤثر است. علاوه بر این، یک شبکه عصبی کانولوشن تک لایه نیز عملکرد بهتری در مقایسه با الگوریتم‌های بیزین و ماشین بردار پشتیبان دارد. علاوه بر این،



۲-۳- پیش پردازش متون

۳-۳- معماری روش پیشنهادی

در این مقاله، داده‌های ورودی به صورت جداگانه به معماری پیشنهادی افزوده می‌شوند. معماری پیشنهادی بر اساس ساختار شبکه‌های عصبی است و ورودی آن لزوماً فرمت متن است و بر اساس مکانیزم داخلی آن متون را به بردارهای عددی تغییر می‌دهد و فرایند خود را دنبال می‌کند. به منظور کاهش اشتباهات از روش سلسله مراتبی مبتنی بر یادگیری عمیق در شبکه عصبی استفاده می‌شود. در حقیقت هدف معرفی روشی است که در آن خطای دسته‌بندی به حداقل برسد. اگرچه روش‌های دسته‌بندی متون بر مبنای شبکه‌های عصبی دارای راندمان بالایی هستند [۱۷-۱۵]. لذا فرضیه تحقیق حاضر این است که اگر از دانش ساخت متون در ساختار مدل تولید شده استفاده کنیم، متن بهتری به دست می‌آید. ایده ما در ارائه این مدل این است که تمام اجزای متن در پاسخ به جستجو یک وابستگی یکسان ندارند. علاوه بر این، شناسایی بخش‌های مختلف روابط بین کلمه هم وابسته است.

روش پیشنهادی یک ساختار عصبی جدید یا شبکه سلسله مراتبی است که برای انجام دو ویژگی اصلی در مورد ساختار متن طراحی شده است. نخست آن که چون متون ساختار سلسله مراتبی دارند (از جمله کلمات عبارت و عبارات یک متن و در مقیاس بزرگ متن از مجموعه‌ای از اسناد)، به همین دلیل نمایش متون ابتدا با عبارات آغاز می‌شود و از کنار هم قرار گرفتن عبارات متن نهایی حاصل می‌شود. همچنین، مشاهده شده است که کلمات و عبارات متن متفاوت متون و بارهای مختلف هستند. علاوه بر این، اهمیت کلمات و عبارات عمیقاً وابسته به زمینه‌ای است که از آن‌ها استفاده می‌شود. به طوری که همان کلمه ممکن است به طور کاملاً متفاوت در یک زمینه دیگر بامعنای دیگری به کار برود. به همین دلیل به منظور در نظر گرفتن این حساسیت، مدل پیشنهادی در دو سطح مکانیسم توجه قرار گرفته است: یکی در سطح کلمه و دیگری در سطح عبارت.

ساختار در روش پیشنهادی سلسله مراتبی است و مکانیزم آن در ابتدا در سطح کلمه و سپس در سطح جمله و در نهایت در سطح سند است. در واقع، این به این دلیل است که اهمیت کلمات و جملات وابسته به زمینه‌ای است که در آن‌ها استفاده

- نشانه‌گذاری: تبدیل یک سند به یک گروه از کلمات (W_0, W_1, \dots, W_n) که در آن n فرکانس کلمات منحصر به فرد در یک داکيومنت متنی است.
- تبدیل حروف بزرگ به کوچک: تمام کلمات را از حروف بزرگ به حروف کوچک تبدیل می‌کند. بر این اساس، کلمات مانند "Runner" یا "runner" به معنای دونده، هر دو به عنوان کلمه "runner" شناسایی می‌شوند
- حذف علائم زائد: به معنای حذف علائمی همانند: ، ؛ " ' ! @ # \$ % ^ * < > ؟ و موارد دیگر.

جدول ۱: اطلاعات دسته‌ها به تفکیک موضوع

شماره برچسب	دسته	زبان
۰	هنر	انگلیسی
۱	ورزشی	انگلیسی
۲	پزشکی	انگلیسی
۳	مذهبی	انگلیسی
۴	خبری	انگلیسی
۵	سیاسی	انگلیسی
۶	فرهنگی	انگلیسی
۷	اقتصادی	انگلیسی
۸	اجتماعی	انگلیسی
۹	تاریخی	انگلیسی
۱۰	سایر	انگلیسی

جدول ۲: نمونه‌ای از پیش پردازش

اعمال پیش پردازش	
متن	"Christopher Columbus", the earliest European visitor, made his first land Farsill in the New World on the Bahamian island of San Salvador in 1492!!!
تبدیل حروف بزرگ به کوچک	"christopher columbus", the earliest european visitor, made his first land farsill in the new world on the bahamian island of san salvador in 1492!!!
حذف علائم زائد	christopher Columbus the earliest european visitor made his first landFarsill in the new world on the bahamian island of san salvador in 1492



بصری در انسان است، به این معنی که توجه بصری انسان می تواند در یک نقطه خاص از یک تصویر با "وضوح بالا" متمرکز شود درحالی که در اطراف آن با "وضوح پایین" درک می شود و سپس / او می تواند نقطه کانونی را در طول زمان تنظیم کند. مکانیزم توجه در NLP، مدل را با قابلیت یادگیری آنچه مطابق با متن ورودی و آنچه تاکنون تولید کرده است، یاد می گیرد.

بر این اساس، از طریق مکانیزم توجه، ما کلمات را که در درک متن بسیار مهم هستند، شناسایی می کنیم، سپس واژه هایی را که از لحاظ معنایی غنی هستند را به شکل یک بردار تبدیل می کنیم.

به منظور تشخیص جملات که سرخ برای دسته بندی یک متن هستند، اهمیت هر جمله از طریق یک بردار و با استفاده از یک بردار مشخص می شود. سپس، با توجه به جملات در اسناد متنی، احتمال قوی تر از دسته بندی متن چند متن می تواند به دست آید، که چقدر از متن یک موضوع خاص را ارائه می دهد و چقدر از آن موضوع دیگری را نشان می دهد. هدف از انجام این کار یافتن متون مبتنی بر وب است که از کلاس های مختلف تشکیل شده و گاهی درون یک متن، مثلاً متن اقتصادی، مسائل سیاسی مطرح می شود. که در این حالت سطح وابستگی متن به هر کلاس مهم است.

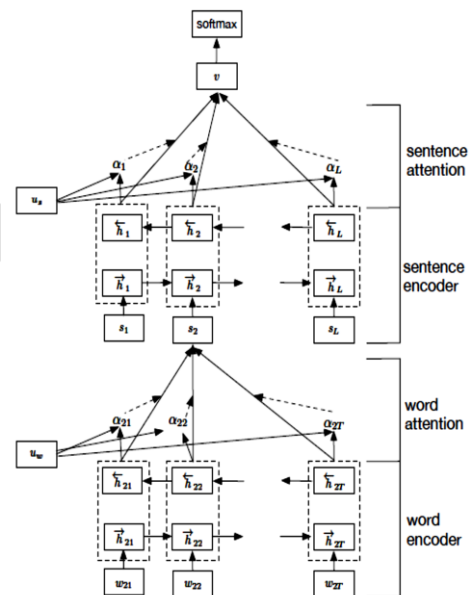
۴- نتایج

در این مقاله مجموع کل داده ها در بخش آموزش شامل ۱۹۶۰۱۰ سند و ۲۲۰۰۰ سند برای بخش آزمون بود. تمامی داده های آموزش و آزمون به صورت دوره ای به صورت جداگانه به شبکه وارد می شدند که سیستم فرآیند دسته بندی را در کمتر از ۰٫۱۴ انجام داد. مرحله آزمایش روی GEFORCE GTX 1080 انجام گردید، که یک میلیون سند متنی در هر ثانیه دسته بندی شدند. این پژوهش بر اساس زبان پایتون با استفاده از ابزارهای ویژه مانند کتابخانه های جستجوی متن مانند NLTK، Genism و TensorFlow صورت گرفته است.

اثر پیش پردازش داده ها در تصاویر زیر نشان داده شده است، همان طور که مشخص است، با توجه به این واقعیت که داده ها از تمام وب سایت های گرفته شده اند، بنابراین آن ها شامل برخی بخش های غیر مرتبط هستند. از آنجا که تنها بخش های

می شود، به طوری که یک کلمه یا جمله ممکن است یک مقدار کاملاً متفاوت در یک زمینه متفاوت داشته باشد. بنابراین، سه سطح مکانیزم در نظر گرفته شده است، و این به ما امکان می دهد بخش های مهم تر متن را از قسمت های مهم کمتر درحالی که ساخت یک متن خاص ارائه می دهد، تمایز بخشد.

ساختار کلی شبکه سلسله مراتبی مورد استفاده (HAN) در نمودار زیر ارائه شده است که شامل اجزای مختلف است: رمزگذار و توجه در سطح کلمه و رمزگذار و توجه در سطح جمله و در نهایت اضافه کردن رمزگذار و توجه به کل سند.



شکل ۱: معماری روش HAN [۱۸]

از طریق روش سلسله مراتبی استفاده از بردارهای کلمات (Word2vec) بر اساس عملکرد مکانیک رمزگذار و توجه در سطوح کلمه و جمله به صورت متقابل (یکبار از ابتدا به ابتدای یکبار از ابتدا تا انتها)، شباهت ها شناسایی می شوند و در زمان آموزش، دقت آن در تشخیص افزایش می یابد. به طور کلی می توان گفت که آموزش شبکه HAN، به شدت به بردارهای کلمات (Word2Vec) بستگی دارد و این بردارها هر قدر غنی تر باشند آموزش بهتر انجام می گیرد. در واقع، آنچه در این روش پیشنهادی مهم است، جایگزینی عملکرد پیش فرض آن از tanh به sigmoid است که بر دقت یا خطای خروجی تأثیر می گذارد [۱۸].

مکانیزم توجه در شبکه های عصبی یک کپی از مکانیزم توجه



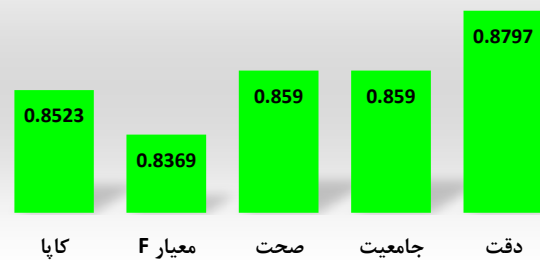
سلسله مراتبی گسترده استفاده شد. در ابتدا کلمات موجود در سطح جملات متون موجود به صورت (BOW) استخراج شدند، سپس با استفاده از قدرت بالای ترکیب معنایی برای دسته بندی استفاده شد. در مقایسه با روش های سنتی و مدرن دسته بندی متن، رویکرد کنونی می تواند نتایج دسته بندی اسناد به زبان انگلیسی را بهبود بخشد. در این مقاله شبکه های توجه سلسله مراتبی معرفی گردید. این مدل در نهایت بردارهای سند را از طریق ایجاد یک بردار متشکل از کلمات در بردارهای جمله می سازد که نتیجه همگرایی بردارهای کلمات مهم بوده و در نهایت، بردارهای سند تشکیل می شوند. مدل پیشنهادی روش های قبلی را به طور قابل توجهی بهبود می بخشد. ارائه لایه های به دست آمده توجه نشان می دهد که مدل عملکرد مناسبی و قابل توجهی در انتخاب کلمات و جملات داشته است. علاوه بر این، در متن هایی که احتمال آن ها کمتر از آستانه تشخیص است، استفاده از حافظه کوتاه مدت در RNN می تواند به مدل کمک کند تا مدل ها را به طور مؤثرتر تشخیص دهد. نقطه نهایی این است که در این روش، تعداد پارامترهایی که در مرحله آموزش به روز می شوند کمتر گردیده و این خود باعث می شود مدت زمان فاز آموزش کاهش یافته و با سرعت بالاتری صورت پذیرد. از این رو این روش می تواند جهت دسته بندی متون در مباحث داده های حجیم راه حل مناسبی باشد.

مراجع

- [1] Stahl F, Gaber MM, Adedoyin-Olowe M. "A Survey of Data Mining Techniques for Social Media Analysis". *Journal of Data Mining & Digital Humanities*. 2014.
- [2] Kwon O, Sim JM. "Effects of data set features on the performances of classification algorithms". *Expert Systems with Applications*. Vol. 40, No. 5, pp. 1847-1857, 2013.
- [3] Meng X, Bradley J, Yavuz B, Sparks E, Venkataraman S, Liu D, Freeman J, Tsai DB, Amde M, Owen S, Xin D. "Mllib: Machine learning in apache spark". *The Journal of Machine Learning Research*. Vol. 17, No. 1, pp. 1235-1241, 2016.
- [4] Allen ST, Jankowski M, Pathirana P. "Storm Applied: Strategies for real-time event processing". *Manning Publications Co.*; 2015.
- [5] Carbone P, Katsifodimos A, Ewen S, Markl V, Haridi S, Tzoumas K. "Apache flink: Stream and batch processing in a single engine". *Bulletin of the IEEE Computer Society Technical Committee on Data Engineering*. Vol. 36, No. 4, 2015.
- [6] Barga R, Fontama V, Tok WH, Cabrera-Cordon L. "Predictive analytics with Microsoft Azure machine learning". *Apress*, 2015.
- [7] De Francisci Morales G. "SAMOA: A platform for mining big data streams". In: *Proceedings of the 22nd International Conference on World Wide Web*. pp. 777-778, 2013.
- [8] Wang, S., & Manning, C. D. "Baselines and bigrams: Simple, good sentiment and topic classification", In: *Proceedings of the 50th Annual Meeting of the Association for Computational Linguistic*. Vol. 2, pp. 90-94, 2012.

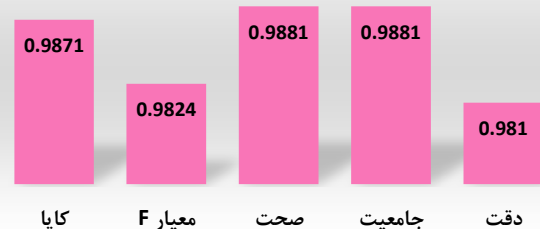
کارآمد در هر خبر تنها همان متن خبرمی باشد بخش های دیگر تأثیر منفی بر خروجی دارند به همین دلیل در مرحله پیش پردازش این موارد زائد شناسایی شده و حذف می شوند.

قبل از پیش پردازش



شکل ۲: نمودار میله ای معیارهای ارزیابی قبل از پیش پردازش

بعد از پیش پردازش



شکل ۳: نمودار میله ای معیارهای ارزیابی بعد از پیش پردازش

با توجه به زمان صرف شده برای دسته بندی متن، با توجه به این که یک صفحه استاندارد دارای ۱۵۰۰ کلمه است، سند متون مربوط به آزمون حاوی ۲۰۵۶ صفحه بود که در نهایت به منظور آزمون سرعت تجزیه و تحلیل، سرعت دسته بندی کل متون مورد آزمایش ۲۰،۱۴۸ ثانیه بود که به طور تقریبی می توان گفت که هر صفحه در ۰،۱۴۱۷۷۰۴۳ ثانیه برای تجزیه و تحلیل مورد دسته بندی قرار می گرفت.

۵- نتیجه گیری و بحث

در این مقاله از یک مدل یادگیری عمیق برای دسته بندی



- [9] Pang and Lillian Lee, "Opinion mining and sentiment analysis", *Foundations and trends in information retrieval*. Vol. 2, No. 1-2, pp. 1-135, 2008.
- [10] Sahami, M., Dumais, S., Heckerman, D., & Horvitz, E, "A Bayesian approach to filtering junk e-mail. In Learning for Text Categorization". *Papers from the 1998 workshop*, Vol. 62, pp. 98-105, 1998.
- [11] Blunsom, P., Grefenstette, E., & Kalchbrenner, N., "A convolutional neural network for modelling sentences", In: *Proceedings of the 52nd Annual Meeting of the Association for Computational Linguistics. Proceedings of the 52nd Annual Meeting of the Association for Computational Linguistics*, 2014.
- [12] Hochreiter, S., & Schmidhuber, J. "Long short-term memory", *Neural computation*. Vol. 9, pp. 1735-1780, 1997.
- [13] Y. Kim, "Convolutional Neural Networks for Sentence Classification", *Empirical Methods in Natural Language Processing (EMNLP)*. Doha, 2014.
- [14] S. Poria, E. Cambria and A. Gelbukh, "Deep Convolutional Neural Network Textual Features and Multiple Kernel Learning for Utterance-Level Multimodal Sentiment Analysis". *Empirical Methods in Natural Language Processing*. Lisbon, 2015.
- [15] A. Sarker and G. Gonzalez, "Portable automatic text classification for adverse drug reaction detection", *Journal of Biomedical Informatics*. Vol. 53, pp. 196-207, 2015.
- [16] Zhang X, Zhao J, LeCun Y. "Character-level convolutional networks for text classification". In: *Advances in neural information processing systems*. pp. 649-657, 2015.
- [17] Johnson R, Zhang T. "Effective use of word order for text categorization with convolutional neural networks". ArXiv preprint arXiv:1412.1058. 2014.
- [18] Z. Yang, D. Yang, C. Dyer, X. He, A. Smola and E. Hovy, "Hierarchical Attention Networks for Document Classification", in *NAACL-HLT*, San Diego, 2016.

Archive