

بهبود دقت واژگان کلیدی استخراج شده از

متن فارسی با استفاده از الگوریتم Word2Vec

محمد رضا حسنی آهانگر* و علی امیری جزه

آزمایشگاه داده‌های حجیم، مرکز داده‌های حجیم و جنگ‌شناختی سایبری، دانشکده رایانه و قدرت سایبری، دانشگاه جامع امام حسین(ع)، تهران، ایران



چکیده

واژگان کلیدی لغات مهمی از سند هستند که بیان‌گر توصیفی از متن هستند و نقش بسیار مهمی در فهم دقیق و سریع از محتوا دارند. شناسایی واژگان کلیدی از متن با روش‌های معمول کاری زمان‌بر و پرهزینه است. در این مقاله ابتدا با استفاده از شبکه عصبی پیشرو و از طریق الگوریتم Word2Vec ماتریس همبستگی واژگان را به‌ازای یک سند محاسبه و سپس با استفاده از ماتریس همبستگی و یک فهرست اولیه محدود از واژگان کلیدی، نزدیک‌ترین واژگان را از نظر شباهت در قالب فهرست نزدیک‌ترین همسایگی‌ها استخراج می‌کنیم. فهرست به‌دست آمده را به‌صورت نزولی مرتب و از ابتدای فهرست، درصدهای مختلفی از واژگان را انتخاب و به‌ازای هر درصد، ده مرتبه فرایند آموزش شبکه عصبی و ساخت ماتریس همبستگی و استخراج فهرست نزدیک‌ترین همسایگی‌ها را تکرار و در نهایت میانگین دقت، فراخوانی و معیار F را محاسبه می‌کنیم. این کار را تا جایی ادامه می‌دهیم که به بهترین نتایج در ارزیابی دست یابیم؛ نتایج نشان می‌دهند که به‌ازای انتخاب حداکثر چهار درصد واژگان از ابتدای فهرست نزدیک‌ترین همسایگی‌ها، نتایج مورد قبولی به‌دست می‌آید. الگوریتم بر روی پیکره‌ای با هشتصد خبر که به‌صورت دستی واژگان کلیدی آن‌ها را استخراج کرده‌ایم، آزمایش شده است و نتایج آزمایش‌ها نشان می‌دهد که دقت روش پیشنهادی ۷۸ درصد خواهد بود.

واژگان کلیدی: واژگان کلیدی، الگوریتم word2Vec، شبکه عصبی، وزن دهی ویژگی

Improving Precision of Keywords Extracted From Persian Text Using Word2Vec Algorithm

Mohammad Reza Hasani Ahangar* & Ali Amiri Jezeh

Big Data Laboratory, Computer Science and Technology Center, Faculty of Information and Communication Technology, Imam Hossein University, Tehran, Iran

Abstract

Keywords can present the main concepts of the text without human intervention according to the model. Keywords are important vocabulary words that describe the text and play a very important role in accurate and fast understanding of the content. The purpose of extracting keywords is to identify the subject of the text and the main content of the text in the shortest time. Keyword extraction plays an important role in the fields of text summarization, document labeling, information retrieval, and subject extraction from text. For example, summarizing the contents of large texts into smaller texts is difficult, but having keywords in the text can make you aware of the topics in the text. Identifying keywords from the text with common methods is time-consuming and costly. Keyword extraction methods can be classified into two types with observer and without observer. In general, the process of extracting keywords can be explained in such a way that first the text is converted into smaller units called the word, then the redundant words are removed and the remaining words are weighted, then the keywords are selected from these words. Our proposed method in this paper for identifying keywords is a method with observer. In this paper, we first calculate the word correlation matrix per document using a feed forward neural network and Word2Vec algorithm. Then, using the correlation matrix and a limited

* Corresponding

* نویسنده عهده‌دار مکاتبات

سال ۱۴۰۰ شماره ۱ پیاپی ۴۷

تاریخ ارسال مقاله: ۱۳۹۷/۰۲/۰۲ • تاریخ پذیرش: ۱۳۹۹/۱۲/۰۹ • تاریخ انتشار: ۱۴۰۰/۰۳/۰۱ • نوع مطالعه: کاربردی

فصلنامه علمی



۵۱

initial list of keywords, we extract the closest words in terms of similarity in the form of the list of nearest neighbors. Next we sort the last list in descending format, and select different percentages of words from the beginning of the list, and repeat the process of learning the neural network 10 times for each percentage and creating a correlation matrix and extracting the list of closest neighbors. Finally, we calculate the average accuracy, recall, and F-measure. We continue to do this until we get the best results in the evaluation, the results show that for the largest selection of 40% of the words from the beginning of the list of closest neighbors, the acceptable results are obtained. The algorithm has been tested on corpus with 800 news items that have been manually extracted by keywords, and laboratory results show that the accuracy of the suggested method will be 78%.

Keywords: keywords, word2vec algorithm, neural network, giving weight features

واژگان در سند دیگری به عنوان واژگان کلیدی آموزنده در نظر گرفته شوند. در سند یادشده همان طور که بیان شد، واژگانی به عنوان واژگان کلیدی آموزنده در نظر گرفته می شوند که ارتباط خیلی قوی با محتویات و معنی متن داشته باشند، واژگانی از قبیل " گاز، کارخانه، گرم، زمین، ازون، آب، یخ، قطب " البته این واژگان نیز ممکن است در سند دیگری به عنوان واژگان کلیدی غیرمفید لحاظ شوند. در این مقاله استخراج واژگان کلیدی آموزنده مدنظر است. تاکنون روش های بسیاری برای استخراج واژگان کلیدی از یک سند ارائه شده است، که دارای معایب و مزایایی هستند، روشی که در این مقاله ارائه می شود یک روش مبتنی بر شبکه عصبی با در نظر گرفتن همبستگی واژگان از لحاظ معنایی و شباهت با یکدیگر است. در این مقاله ابتدا مروری بر کارهای پیشین و مشابه صورت پذیرفته است؛ سپس روش رایج استخراج واژگان کلیدی بیان و در ادامه راه حل پیشنهادی با بیان جزئیات مطرح و در پایان نتایج کار ارائه شده است.

۲- مروری بر کارهای پیشین

واژگان و عبارات کلیدی می توانند مفاهیم اصلی متن را بدون دخالت انسان بسته به مدل ارائه دهند [6]. خلاصه سازی محتویات متن های بزرگ به متن های کوچک تر برای آنکه بهتر فهمیده شوند کاری سخت است اما با داشتن واژگان کلیدی موجود در متن می توان از موضوعات داخل متن آگاه شد. از این روی کارهای بسیاری در این حوزه صورت پذیرفته است.

در [7] ابتدا اسناد خوشه بندی می شوند، سپس با استفاده از ویژگی های به دست آمده از نزدیک ترین همسایگان یک سند، K سندی که در واقع بیشترین شباهت را به سند مشخص شده دارند، به سند اضافه و سند بسط داده می شود، در نهایت الگوریتم رتبه بندی گراف بر روی سندی که گسترش یافته است، اعمال می شود و بر اساس ویژگی های سند و همسایگان سند، واژگان کلیدی استخراج شده از سند، بهبود می یابند. در [8] ابتدا عبارات

۱- مقدمه

هدف از استخراج واژگان کلیدی تشخیص موضوع متن و محتوای اصلی متن در کمترین زمان است. استخراج واژگان کلیدی نقش مهمی در زمینه های خلاصه سازی متون، برچسب گذاری اسناد، بازیابی اطلاعات و استخراج موضوع از متن دارد [1]. یکی از مشکلات مهم در تشخیص واژگان کلیدی، شناسایی واژگان غیرکلیدی است [2] به گونه ای که بیانگر موضوع اصلی متن نیستند. واژگان کلیدی، عناصر بسیار مهمی در جست و جو و دسترسی اطلاعات هستند. آن ها می توانند به عنوان مجموعه واژگان تشریح کننده سند در طی عملیات جست و جو مدنظر قرار بگیرند. به عبارت دیگر، به هر عبارت مهمی که محتویات داخل سند را تشریح کند، واژه کلیدی گفته می شود [3]. واژگان کلیدی در دو گروه طبقه بندی می شوند [4]، واژگان کلیدی تابعی^۱ و واژگان کلیدی آموزنده^۲.

واژگان کلیدی تابعی یا غیرمفید برای واژگان دستوری یا مرتبط با زبان استفاده می شوند [4] که منظور از واژگان دستوری، واژگان معناداری هستند که به اشیا و مفاهیم اشاره می کنند، مانند " کتاب، قلم و نوشتن "؛ این واژگان ارتباط کمی با محتویات سند دارند و باید حذف شوند و در نظر گرفته نشوند.

واژگان کلیدی آموزنده ارتباط خیلی قوی با محتویات و معنای متن سند دارند. تشخیص مرز بین واژگان کلیدی تابعی و آموزنده خیلی سخت نیست و می توان یک مرز فازی برای آن ها در نظر گرفت [4, 5]. برای مثال سندی را که اثرات افزایش گازهای گلخانه ای را در سوراخ شدن لایه ازون بررسی کرده است در نظر بگیرید؛ آن دسته از واژگان کلیدی که ارتباط کمتری نسبت به محتویات سند دارند، به عنوان واژگان کلیدی غیرمفید در نظر گرفته می شوند، مانند واژگانی از قبیل " سال، جهان، جبران، صدمات " البته ممکن است این

¹ Functional

² Informative

قرار می‌گیرند، واژگان کلیدی مشابه‌ای دارند و شناسایی واژگان کلیدی جدید در این اسناد ممکن نیست. در [10] هم‌رخدادی صریح واژگان مدنظر قرار گرفته‌است و واژگانی که به‌منظور جلوگیری از تکرار با ضمائر و واژگان اختصار بیان شده‌اند جزو واژگان پرتکرار قرار نمی‌گیرند. در [11] واژگانی که فقط با یکدیگر تکرار می‌شوند، تشکیل یک عبارت کلیدی را نمی‌دهند. در [12, 13] عبارات کلیدی تابعی به مراتب بیشتر از دیگر روش‌ها وجود دارد، چون از ویژگی‌های روابط معنایی بین واژگان برای شناسایی عبارات کلیدی استفاده نشده است.

۳- روش استخراج واژگان کلیدی

روش‌های استخراج واژگان و عبارات کلیدی را می‌توان در دو نوع با ناظر^۳ و بدون ناظر^۴ دسته‌بندی کرد [7]. اگر در مجموعه‌ای از اسناد واژگان کلیدی مشخص شده باشند، فرایند استخراج واژگان کلیدی با ناظر و در غیر این صورت بدون ناظر خواهد بود؛ اما به‌طور کلی می‌توان فرایند استخراج واژگان کلیدی از متن را این‌گونه توضیح داد که در ابتدا متن به واحدهای کوچک‌تری به نام واژه تبدیل می‌شود، بعدازآن واژگان زائد حذف و واژگان باقیمانده وزن‌دهی می‌شوند، سپس واژگان کلیدی از بین این واژگان انتخاب می‌شوند.

۴- بردار ویژگی‌ها

واژه‌ها و واژگان برای چندین دهه به‌عنوان واحد اصلی در زبان طبیعی مورد مطالعه قرار گرفته‌اند [14]. در پردازش زبان طبیعی یکی از مهم‌ترین پارامترها در تمام وظایف، چگونگی ارائه واژگان به‌عنوان ورودی به هر یک از مدل‌هاست، به بیانی ساده‌تر ابتدا باید مفهوم شباهت و تفاوت بین واژگان را داشته باشیم و این در قالب بردار ویژگی‌های متن قابل بیان است [15]. آنچه در انتخاب واژگان کلیدی حائز اهمیت است، وزن‌دهی واژگان نامزد است. به‌منظور وزن‌دهی واژگان برای استخراج واژگان کلیدی [16]، می‌توان از روش‌هایی اعم از تعداد تکرار واژه (TF) در سند و اسناد مختلف (IDF) [17]، روش‌های یادگیری ماشین [18, 19] و ترکیب روش‌های آماری و زبان‌شناختی [20] استفاده کرد. از آن‌جایی که تعداد قابل‌توجهی از روش‌های قدیمی از معیار TF-IDF استفاده می‌کنند [21] و به‌منظور مقایسه با روش پیشنهادی، در ادامه این روش را توضیح خواهیم داد.

اسمی از سند استخراج می‌شوند، سپس عباراتی که دست‌کم یک واژه مشابه دارند در داخل یک خوشه قرار می‌گیرند؛ با ویژگی تعداد تکرار عبارات اسمی و واژگان داخل این عبارات، خوشه‌ها رتبه‌بندی می‌شوند؛ درنهایت عباراتی به‌عنوان عبارات کلیدی انتخاب می‌شوند که خوشه‌ی مربوطه بالاترین رتبه را داشته باشد. در [9] برای استخراج واژگان کلیدی از الگوریتم‌های یادگیری با ناظر استفاده شده است، به‌گونه‌ای که ابتدا عبارات کلیدی در یک سند را به‌صورت فهرستی از عبارات کلیدی مشخص می‌کنند، سپس الگوریتم باید یاد بگیرد که اسناد را در دو دسته منفی و مثبت طبقه‌بندی کند. در [10] ابتدا واژگان پرتکرار از متن استخراج می‌شوند، سپس هم‌رخدادی این واژگان با تمام واژگان موجود در متن در یک ماتریس محاسبه می‌شود، در این روش ادعا شده است که هر واژه‌ای که با واژگان پرتکرار هم‌رخدادی بیشتری داشته باشد، بهترین نامزد برای واژه کلیدی است. در [11] ابتدا نامزدهای واژگان کلیدی از متن استخراج می‌شوند، سپس هم‌رخدادی بین این نامزدها محاسبه و به‌ازای فراوانی به هر یک امتیازی نسبت داده می‌شود، نامزدهایی به‌عنوان عبارات کلیدی انتخاب می‌شوند که بیشترین امتیاز را داشته باشند. در [12] از شبکه عصبی برای استخراج عبارات کلیدی از سند استفاده می‌شود، ابتدا شبکه عصبی با استفاده از ویژگی‌هایی همچون فراوانی واژه و موقعیت واژه در سند آموزش داده می‌شود، سپس عبارات کلیدی نامزدشده با درنظرگرفتن یک حد آستانه تعریف شده در دو دسته عبارات کلیدی و عبارات غیرکلیدی دسته‌بندی می‌شوند. در [13] از شبکه رقابتی LVQ^۱ و شبکه عصبی MLP^۲ برای استخراج واژگان کلیدی متون استفاده شده است، این الگوریتم‌ها درواقع واژگان موجود در یک متن را به دو خوشه واژگان کلیدی و غیرکلیدی دسته‌بندی می‌کنند، در این روش ادعا شده که نتیجه‌ها بر روی هشتاد سند متنی نشان داده است که روش MLP نسبت به روش دیگر نتایج بهتری را به دست می‌دهد. در [7] واژگان کلیدی که استخراج می‌شوند، نمی‌توانند نمایندگان مناسبی برای اسنادی باشند که بیشترین فاصله را از مرکز خوشه دارا هستند. روش [8] وابسته به زبان است و شناسایی وابسته‌های پسین و پیشین به‌ازای زبان‌های مختلف تفاوت دارند، همچنین استخراج عبارات اسمی، بازنمایی مناسبی برای شناسایی تمام واژگان کلیدی یک سند نیستند. روش [9] نیاز به داده قابل توجهی برای یادگیری و آموزش دارد و فقط اسنادی که در یک دسته

³ Superwise

⁴ Unsuperwise

¹ Learning Vector Quantization

² Multilayer Perceptron

$$p(w_o | w_l) = \frac{\exp(v'_{wo} v_{wl})}{\sum_{w=1}^W \exp(v'_w v_{wl})} \quad (1)$$

یکی از جدیدترین و پیشرفته‌ترین روش‌ها برای بازنمایی^۱ واژگان و استخراج بردار ویژگی‌ها از متن، الگوریتم Word2Vec است. در واقع خروجی Word2Vec بازنمایی بردار چگال واژگان است که می‌تواند مفاهیم معنایی و همبستگی موضوعی را نشان دهد [22]. سطرهای ماتریس بردار ویژگی‌ها در این روش، بیان‌گر روابط معنایی و روابط موضوعی واژگان با یکدیگر هستند [14, 15].

۱-۴- وزن دهی به روش TF-IDF^۲

در این روش تعداد تکرار یک کلمه در یک سند را در مقابل تعداد تکرار آن کلمه در مجموعه تمام اسناد در نظر می‌گیریم. در این روش اگر واژه‌ای در اسناد زیادی ظاهر شود، وزن واژه کمتر خواهد شد و واژه اهمیت کمتری پیدا خواهد کرد.

$$a_{td} = tf_{td} * \log \frac{N}{n_d} \quad (2)$$

در فرمول (۱) tf_{td} نمایانگر تعداد دفعات ظاهر شدن واژه t در سند d و N بیان‌گر تعداد تمامی اسناد موجود در مجموعه داده است و n_d نشان‌دهنده تعداد اسنادی است که واژه t در آن وجود دارد [23]. پس از محاسبه TF-IDF آن دسته از واژگانی که بیشتری نسبت به سایر واژگان دارند، به عنوان واژگان کلیدی انتخاب می‌شوند [21].

۲-۴- وزن دهی به روش skip-gram

در الگوریتم Word2Vec هدف از آموزش مدل skip-gram این است که واژه‌هایی را پیدا کنیم که برای پیش‌بینی واژگان همسایه در یک جمله یا سند مفید باشند [24]. همان‌طور که در فرمول (۳) مشاهده می‌شود هدف از روش skip-gram پیشینه‌کردن میانگین لگاریتم احتمال آمدن واژگان در اطراف یک واژه است.

$$\frac{1}{T} \sum_{t=1}^T \sum_{-c \leq j \leq c, j \neq 0} \log p(w_{t+j} | w_t) \quad (3)$$

در فرمول (۳) c تعداد مجموعه داده ورودی را برای آموزش نشان می‌دهد که با افزایش این مقدار، دقت و هزینه بالاتر می‌رود و $p(w_{t+j} | w_t)$ بر اساس فرمول (۱) به دست می‌آید که v_w و v'_w به ترتیب بردارهای ورودی و خروجی واژگان هستند و W تعداد تمام واژگان است.

شکل (۲) معماری مدل skip-gram را نشان می‌دهد، در این مدل با داشتن یک واژه باید واژگان بعدی

¹ Representation

² Term Frequency inverse Document Frequency

را پیش‌بینی کنیم [24]. از آنجایی که واژگان دورتر نسبت به واژگان نزدیک‌تر رابطه موضوعی یا مفهومی کمتری دارند [25]، پس لازم است که در ابتدا بیشینه فاصله بین واژه فعلی و پیش‌بینی شده در یک جمله را مشخص کنیم، ما این فاصله را با نام اندازه پنجره^۳ می‌شناسیم. در شکل (۱) این فاصله برابر دو در نظر گرفته شده است، یعنی با داشتن یک واژه باید دو لغت قبل‌تر و دو لغت بعدتر را پیش‌بینی کنیم.

هاشمی	تصریح	کرد	وظیفه	مطبوعات	نوشتن	است
هاشمی	تصریح	کرد	وظیفه	مطبوعات	نوشتن	است
هاشمی	تصریح	کرد	وظیفه	مطبوعات	نوشتن	است
هاشمی	تصریح	کرد	وظیفه	مطبوعات	نوشتن	است

شکل (۱)- اندازه پنجره در مدل Skip-gram
(Figure-1): Window Size in Skip-gram Model

در شکل (۲) و در لایه ورودی، x نشان‌دهنده واژه ورودی در نمونه آموزشی است که به صورت یک بردار $1 \times V$ به شبکه عصبی داده می‌شود و $[y_1, \dots, y_c]$ یک بردار به طول $1 \times V$ است که مربوط به واژگان خروجی در نمونه آموزشی با هدف کمینه‌سازی مجموع خطای پیش‌بینی همه واژگان همسایه در لایه خروجی است، ماتریس $W_{V \times N}$ ، ماتریس وزن بین لایه ورودی و لایه پنهان است که سطر i^{th} نشان‌دهنده وزن‌های مربوط به کلمه i^{th} در بین مجموعه لغات است. ماتریس وزن W همان چیزی است که ما علاقه‌مند به یادگیری آن هستیم؛ زیرا پس از پایان آموزش شبکه عصبی شامل بردار تمام واژگان موجود در مجموعه واژگان است. هر بردار کلمه خروجی دارای یک ارتباط $N \times V$ با ماتریس خروجی W دارد.

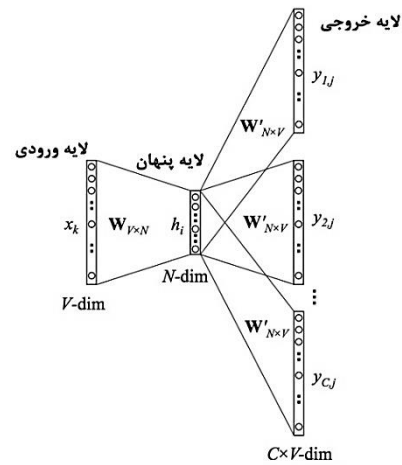
برخی از پارامترهایی که در این روش مورد استفاده قرار گرفته‌اند بدین شرح است:

۱. به صورت پیش‌فرض بیشینه فاصله بین واژه فعلی و پیش‌بینی شده در یک جمله در الگوریتم مقدار سه در نظر گرفته شده است؛ اما با توجه به این که در زبان فارسی عباراتی با طول چهار واژه داریم و فروانی عباراتی با بیشتر از چهار واژه به کمینه می‌رسد، بنابراین بیشینه فاصله بین واژه فعلی و پیش‌بینی شده در یک جمله چهار در نظر گرفته شده است.
۲. کمینه فرکانس واژگان به منظور نادیده گرفته شدن در فرایند آموزش عدد یک در نظر گرفته شده است.
۳. کمینه فرکانس واژگان به منظور نادیده گرفته شدن در فرایند آموزش عدد یک در نظر گرفته شده است.
۴. طول ابعاد بردار ویژگی به اندازه تعداد واژگان موجود در متن در نظر گرفته شده است.

³ Window Size

۶- آموزش شبکه عصبی

در آموزش شبکه عصبی نمی‌توانیم یک واژه را به‌عنوان یک رشته متنی به شبکه عصبی بدهیم، بدین منظور لازم است ابتدا لغت‌نامه واژگان را به‌دست آوریم، برای مثال V واژه منحصره‌فرد داریم که می‌خواهیم از بین آن‌ها واژه "ایران" را به‌عنوان ورودی به شبکه عصبی بدهیم، تنها کافی است که یک بردار را به طول $1 \times V$ که تمام خانه‌های آن با عدد صفر پر شده‌اند، به‌جز خانه‌ای را که مربوط به واژه "ایران" است، به‌عنوان ورودی به شبکه عصبی بدهیم.



شکل (۲): معماری مدل skip-gram [26]

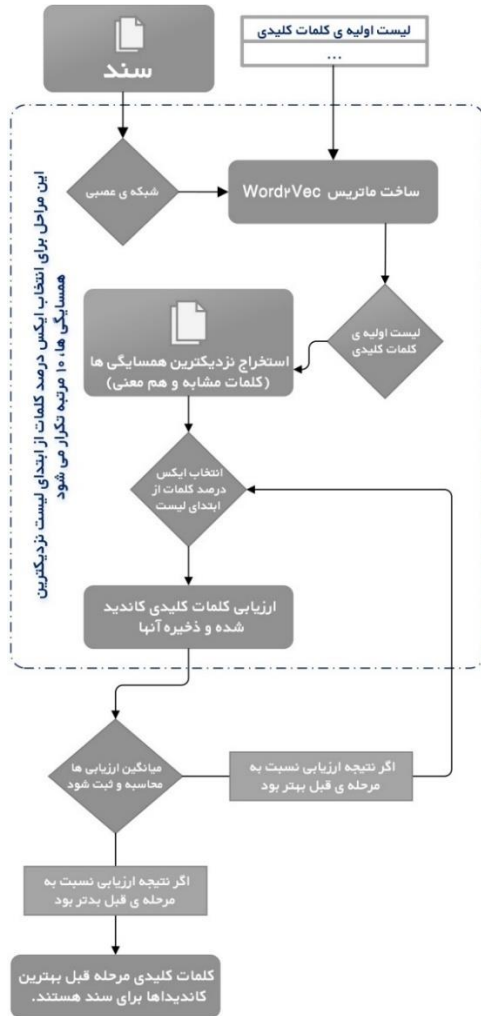
(Figure-2): Skip-gram Model Architecture [26]

۵. کمینه فرکانس واژگان به‌منظور نادیده‌گرفته‌شدن در فرایند آموزش عدد یک در نظر گرفته شده است.
۶. طول ابعاد بردار ویژگی به اندازه تعداد واژگان موجود در متن در نظر گرفته شده است.

۵- راه‌حل پیشنهادی

در روش پیشنهادی، واژگان کلیدی از بردار ویژگی‌های یک سند و همبستگی معنایی و موضوعی واژگان به‌دست می‌آیند. تمام ویژگی‌های یک سند تحت بردار ویژگی‌ها با استفاده از شبکه عصبی پیشرو ۱ و از طریق الگوریتم Word2Vec در قالب یک ماتریس نمایش داده می‌شوند. با توجه به نتایجی که در [25] آمده است، در این مقاله برای ساخت بردار واژگان از روش Skip-gram استفاده می‌کنیم. از مزایای این روش برای ساخت بردار ویژگی‌ها نسبت به روش TF-IDF این است که از ماتریس نهایی ایجادشده با این روش می‌توان میزان ارتباط بین واژگان را محاسبه کرد. محاسبه میزان ارتباط بین واژگان با ترکیب بردارهای آنها با یکدیگر امکان‌پذیر است. برای مثال سه واژه "فرانسه"، "ایران" و "تهران" را در نظر بگیرید با توجه به رابطه "بردار(فرانسه) - بردار(ایران) + بردار(تهران) = ؟" به جواب "بردار(پاریس)" خواهیم رسید.

معماری روش پیشنهادی در شکل (۱) بیان شده است. در روش پیشنهادی فهرست اولیه واژگان کلیدی به شرط داشتن ارتباط قوی با محتوای سند که همان واژگان کلیدی آموزنده هستند، توسط خبره مشخص می‌شود. مجموعه اسنادی که توسط خبره واژگان کلیدی آنها مشخص شده حدود هشتصد سند است. به‌منظور ارزیابی واژگان کلیدی نامزدشده بعد از استخراج واژگان کلیدی با استفاده از معیار F ارزیابی صورت می‌پذیرد.



شکل (۳): معماری روش پیشنهادی

(Figure 3): Suggested Architecture for Improving Keywords

$$[0 \ 0 \ 0 \ 1 \ 0] \times \begin{bmatrix} 17 & 24 & 1 \\ 23 & 5 & 7 \\ 4 & 6 & 13 \\ 10 & 12 & 19 \\ 11 & 18 & 25 \end{bmatrix} = [10 \ 12 \ 19]$$

شکل (۴): بردار ویژگی‌ها در آموزش شبکه عصبی

(Figure-4): Features Vector in Neural Network Training

برای مثال در جمله "بزرگ‌ترین دشمن اسرائیل ایران است"، واژه "ایران" چهارمین واژه است، در شکل

^۱ Feed forward Neural Network

(۴) نیز در چهارمین خانه از ماتریس 1×5 عدد یک بیان گر این موضوع است. تعداد ویژگی‌ها، پارامتری هستند که می‌توان برحسب نیاز کم و زیاد کرد؛ در این مثال ما سه ویژگی را برای ماتریس دوم در نظر گرفته‌ایم، درواقع این ماتریس همان وزن‌های لایه نخست است که تعداد سطرهای آن بیان گر تعداد واژگان لغت‌نامه است که سطر چهارم از این ماتریس درواقع ویژگی‌های استخراج‌شده برای واژه "ایران" است. ماتریس نخست به‌صورت یک شاخص عمل می‌کند و درواقع بیان گر شماره سطر واژه مورد نظر در ماتریس وزن است که ویژگی‌های متناظر با آن در خروجی نمایش داده شده است.

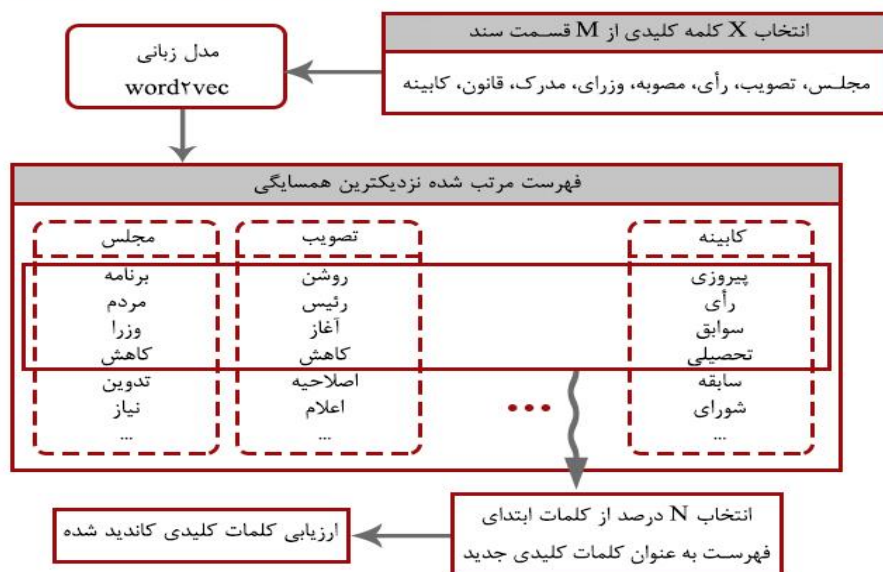


(شکل-۵): نحوه تولید بردار خروجی در فرایند یادگیری شبکه عصبی

(Figure-5): How to generate an output vector in the process of learning the neural network

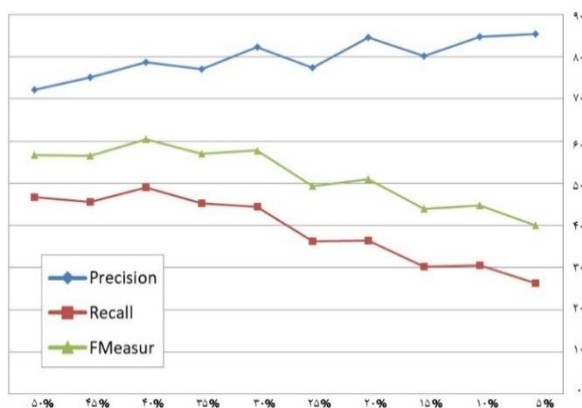
در شکل (۵) از چپ به راست به ترتیب بردار ورودی ($1 \times V$)، ماتریس وزن لایه ورودی به لایه پنهان و همچنین تولید بردارهای نهایی متن ($V \times D$)، لایه پنهان ($1 \times D$)، ماتریس وزن لایه پنهان به لایه خروجی ($D \times V$) و بردار خروجی ($1 \times V$) در فرایند یادگیری شبکه عصبی نشان داده شده است. این شکل به خوبی نشان می‌دهد که چگونه ابعاد ماتریس ورودی در لایه پنهان کمتر شده است، D تعداد ابعاد بردارهای واژه را نشان می‌دهد، در واقع در این مثال در لایه پنهان یک بردار به طول V به یک بردار به طول تعداد ویژگی‌ها تبدیل شده است. در مرحله بعد به منظور تولید بردار خروجی لازم است که ماتریس $1 \times D$ را در ترانهاده ماتریس ویژگی‌ها ضرب کنیم. نکته‌ای که لازم است، گفته شود این است که مجموع احتمالاتی که در بردار خروجی می‌آیند در هر سطر برابر یک است. در پایان کار بردارهای نهایی واژگان، همان بردار مورد نظر ما است.

نایب رئیس اول کمیسیون تدوین آیین نامه داخلی مجلس از شروط مجلس برای وزرای پیشنهادی و تغییر نحوه رأی اعتماد به کابینه خبر داد. طبق این مصوبه و در صورتی که در صحن علنی هم به تصویب برسد، وزرای پیشنهادی باید شرایطی که در قانون ذکر می‌شود را داشته باشند و همچنین نحوه رأی اعتماد به کابینه نیز تغییر می‌کند. طبق مصوبه کمیسیون، وزرای پیشنهادی باید حداقل ۱۰ سال سابقه تجربی مرتبط با وزارتخانه مربوطه را داشته باشند. نایب رئیس اول کمیسیون تدوین آیین نامه داخلی مجلس در ادامه «دارا بودن حداقل مدرک فوق لیسانس یا معادل آن» برای وزرای پیشنهادی را از دیگر مصوبات کمیسیون متبوعش اعلام کرد و گفت: در سنوات و دوره‌های گذشته دارا بودن مدرک تحصیلی شرط نبود. وی تأکید کرد: مجلس به دنبال روشن شدن تکلیفی است که دولت باید در معرفی و پیشنهاد وزرا آن حداقل ها را رعایت کند. کاتب تأکید کرد: به دنبال آن هستیم بعد از پایان تعطیلات تابستانی و آغاز به کار مجلس تا قبل از دوازده مردادماه که مراسم تنفیذ رئیس جمهور منتخب برگزار می‌شود، این اصلاحیه در صحن علنی مطرح شده و به تصویب برسد و شورای نگهبان آن را تأیید کند تا این قانون به معرفی کابینه یازدهم برسد و دولت مکلف به رعایت حداقل شاخص‌ها در معرفی کابینه شود.



(شکل-۶): نحوه شناسایی واژگان کلیدی جدید با استفاده از روش پیشنهادی

(Figure-6): How to identify new keywords using the suggested method



(نمودار ۱-): ارزیابی روش پیشنهادی (محور افقی درصد واژگان انتخاب شده از نزدیک ترین فهرست

همسایگی و محور عمودی درصد ارزیابی)

(Chart-1): Assessment of the suggested method (horizontal axis The percentage of words selected from the list of nearest neighbors and the vertical axis is the percentage of evaluation)

(جدول ۱-): ارزیابی روش پیشنهادی

(Table-1): Assessment of the suggested method

میانگین	ارزیابی به ازای تعداد دفعات تکرار فرایند انتخاب واژگان کاندید										درصد انتخاب
	10	9	8	7	6	5	4	3	2	1	
85.35	76.92	78.57	91.67	76.92	91.67	81.82	80.00	91.67	90.91	93.33	5%
84.78	80.00	92.31	87.50	80.00	78.57	85.71	88.24	92.31	84.62	78.57	10%
80.24	62.50	86.67	75.00	81.25	84.62	78.57	80.00	86.67	73.33	93.75	15%
84.53	88.24	77.78	87.50	84.21	83.33	73.68	88.24	88.24	80.00	94.12	20%
77.31	88.89	68.42	73.68	73.68	70.00	72.22	80.95	88.24	83.33	73.68	25%
82.28	81.82	81.82	71.43	75.00	85.71	85.71	83.33	91.30	80.95	85.71	30%
77.03	66.67	81.82	86.36	69.57	80.00	72.73	79.17	75.00	75.00	84.00	35%
78.75	80.77	83.33	81.48	75.00	82.61	79.17	73.08	76.00	76.92	79.16	40%
75.11	77.78	66.67	75.00	80.00	77.27	74.07	70.00	83.33	72.00	75.00	45%
72.11	60.71	76.92	70.37	76.92	70.37	73.91	69.23	70.37	74.07	78.26	50%
26.25	25.00	27.50	27.50	25.00	27.50	22.50	20.00	27.50	25.00	35.00	5%
30.50	30.00	30.00	35.00	30.00	27.50	30.00	37.50	30.00	27.50	27.50	10%
30.25	25.00	32.50	30.00	32.50	27.50	27.50	30.00	32.50	27.50	37.50	15%
36.50	37.50	35.00	35.00	40.00	37.50	35.00	37.50	37.50	30.00	40.00	20%
36.25	40.00	32.50	35.00	35.00	35.00	32.50	42.50	37.50	37.50	35.00	25%
44.50	45.00	45.00	37.50	37.50	45.00	45.00	50.00	52.50	42.50	45.00	30%
45.25	40.00	45.00	47.50	40.00	50.00	40.00	47.50	45.00	45.00	52.50	35%
49.00	52.50	50.00	55.00	45.00	47.50	47.50	47.50	47.50	50.00	47.50	40%
45.50	52.50	40.00	45.00	50.00	42.50	50.00	35.00	50.00	45.00	45.00	45%
46.75	42.50	50.00	47.50	50.00	47.50	42.50	45.00	47.50	50.00	45.00	50%
40.05	37.74	40.74	42.31	37.74	42.31	35.29	32.00	42.31	39.22	50.90	5%
44.79	43.64	45.28	50.00	43.64	40.74	44.44	52.63	45.28	41.51	40.74	10%
43.90	35.71	47.27	42.86	46.43	41.51	40.74	43.64	47.27	40.00	53.57	15%
50.94	52.63	48.28	50.00	54.24	51.72	47.46	52.63	52.63	43.64	56.14	20%
49.32	55.17	44.07	47.46	47.46	46.67	44.83	55.74	52.63	51.72	47.45	25%
57.73	58.06	58.06	49.18	50.00	59.02	59.02	62.50	66.67	55.74	59.02	30%
56.98	50.00	58.06	61.29	50.79	61.54	51.61	59.37	56.25	56.25	64.61	35%
60.38	63.64	62.50	65.67	56.25	60.32	59.37	57.58	58.46	60.61	59.37	40%
56.58	62.69	50.00	56.25	61.54	54.84	59.70	46.67	62.50	55.38	56.25	45%
56.67	50.00	60.61	56.72	60.61	56.72	53.97	54.55	56.72	59.70	57.14	50%

دقت

بازخوانی

معیار F



۷- نزدیک ترین همسایگی

با استفاده از مدل word2Vec که خروجی یادگیری شبکه‌ی عصبی است، می‌توان واژگانی را که رابطه‌ی معنایی و همبستگی موضوعی با یکدیگر دارند، را پیدا کرد. در این مقاله یک پیکره‌ هشتصدتایی را از اخبار که واژگان کلیدی آنها توسط خبره استخراج شده‌اند، آماده کردیم، سپس با استفاده از روش اعتبارسنجی متقابل^۱ که یک روش برای ارزیابی مدل است و بیان‌گر میزان مستقل بودن و تعمیم‌پذیری نتایج یک تحلیل آماری بر روی یک مجموعه‌داده نسبت به داده‌های آموزشی است، X واژه کلیدی را از M قسمت سند شناسایی و انتخاب می‌کنیم. برای این منظور ابتدا متن سند به صورت تصادفی به M قسمت تقسیم می‌شود که با توجه به طول سندهای مورد استفاده در این مقاله، بیشینه مقدار M در روش پیشنهادی مقدار پنج در نظر گرفته شده است. انتخاب X واژه کلیدی این‌گونه است که به صورت تصادفی از هر قسمت تعدادی واژه کلیدی را که از قبل توسط خبره برچسب‌گذاری شده است، انتخاب می‌کنیم. اولویت انتخاب با واژگان کلیدی است که در قسمت‌های بیشتری از متن برچسب خورده وجود داشته باشند. حال با استفاده از مدل word2Vec که به‌ازای سند مربوطه ساخته شده است و این X واژه کلیدی کاندید شده، نزدیک‌ترین واژگان را به این واژه‌ها می‌یابیم. این فهرست جدید از واژگان را فهرست نزدیک‌ترین همسایگی‌ها می‌نامیم. این فهرست را بر اساس امتیاز واژگان که معرف میزان شباهت معنایی و موضوعی آن‌ها با واژگان نامزدشده اولیه هستند به صورت نزولی مرتب می‌کنیم. حال از بین نزدیک‌ترین همسایگی‌ها، پنج درصد از واژگان ابتدای فهرست را به‌عنوان واژگان کلیدی انتخاب، سپس دقت، بازخوانی و معیار F را محاسبه و این کار را به‌ازای پنج درصد واژگان ابتدای فهرست، ده بار تکرار و درنهایت به‌ازای دقت، بازخوانی و معیار F، میانگین این ده مرحله را ثبت می‌کنیم. این فرایند را به‌ازای انتخاب درصدهای مختلفی از فهرست نزدیک‌ترین همسایگی‌ها، تا جایی ادامه می‌دهیم که معیار F به بالاترین مقدار خود برسد. در شکل (۶) به‌خوبی مراحل کار برای انتخاب واژگان کلیدی نامزد و ارزیابی آنها نشان داده شده است، همچنین در نمودار (۱) و جدول (۱) به‌خوبی نتایج کار قابل مشاهده است.

۸- نتایج و بحث

در این مقاله با استفاده از مدل Skip-gram برای آموزش شبکه‌ی عصبی و الگوریتم Word2Vec به‌منظور ساخت

^۱ Cross-validation

بردار ویژگی‌ها، توانستیم بردار ویژگی‌های تمام واژگان موجود در یک متن را استخراج کنیم، با این تفاوت که در روش به‌کار گرفته‌شده نسبت به سایر روش‌ها از ماتریس نهایی می‌توان رابطه‌ی معنایی و همبستگی موضوعی بین لغات را محاسبه کرد. از این ویژگی استفاده کردیم و با بهره‌گیری از یک فهرست اولیه از واژگان کلیدی توانستیم واژگان مشابه با آن‌ها را در قالب فهرست نزدیک‌ترین همسایگی‌ها استخراج کنیم. فهرست به‌دست‌آمده را به صورت نزولی مرتب و از ابتدای فهرست، درصدهای مختلفی از واژگان را انتخاب و به‌ازای هر درصد، ۱۰ مرتبه فرایند آموزش شبکه‌ی عصبی و ساخت ماتریس همبستگی و استخراج فهرست نزدیک‌ترین همسایگی‌ها را تکرار و درنهایت میانگین دقت، فراخوانی و معیار F را محاسبه می‌کنیم. این کار را تا جایی ادامه می‌دهیم که به بهترین نتایج در ارزیابی دست یابیم، نتایج نشان می‌دهند که به‌ازای انتخاب حداکثر چهار درصد واژگان از ابتدای فهرست نزدیک‌ترین همسایگی‌ها، نتایج مورد قبولی به دست می‌آید. نتایج آزمایش‌ها نشان می‌دهد که دقت روش پیشنهادی ۷۸ درصد خواهد بود.

نتایج به‌دست‌آمده در روش پیشنهادی با [27] مقایسه و در جدول (۲) نشان داده شده است.

(جدول ۲): مقایسه نتایج حاصل از روش پیشنهادی با

حدآستانه ۵ و ۱۰ درصد

(Table-2): Comparison of the results of the proposed method with thresholds of 5% and 10%

معیار ارزیابی	روش پیشنهادی		روش [27]	
	٪۵	٪۱۰	٪۵	٪۱۰
دقت	۸۵/۳۵	۸۴/۷۸	۴۹/۲۳	۲۶/۱۹
فراخوانی	۲۶/۲۵	۳۰/۵۰	۶۰/۱۴	۶۷/۹۰
معیار F	۴۰/۰۵	۴۴/۷۹	۵۳/۹۲	۳۷/۸۰

۹- نتیجه‌گیری و کارهای آینده

در این مقاله سعی بر آن شد تا روشی جدید و مستقل از زبان برای بهبود واژگان کلیدی استخراج‌شده از متون ارائه شود که علاوه بر در نظر گرفتن فراوانی واژگان، از روابط و همبستگی بین واژگان نیز برای بهبود کیفیت استفاده می‌کند، همان چیزی که در روش‌های قبلی که در این مقاله و در بخش کارهای پیشین به آن‌ها اشاره شده است، مورد توجه قرار نگرفته است.

به‌عنوان پیشنهاد جهت پروژه‌های آینده می‌توان موارد زیر را مطرح کرد:

- [11] S. Rose, D. Engel, N. Cramer, and W. Cowley, "Automatic keyword extraction from individual documents," *Text mining: applications and theory*, vol. 1, pp. 1-20, 2010.
- [12] J. Wang, H. Peng, and J.-s. Hu, "Automatic keyphrases extraction from document using neural network," in *Advances in Machine Learning and Cybernetics: Springer*, 2006, pp. 633-641.
- [13] A. Ahmadi and T. Hosseinkhah, "Extract keywords from a text using neural networks," presented at the 10th International Conference on Industrial Engineering, Tehran, Iran Industrial Engineering Association, Amirkabir University of Technology, 2013.
- [14] S. De Deyne, S. Verheyen, and G. Storms, "Structure and organization of the mental lexicon: A network approach derived from syntactic dependency relations and word associations," in *Towards a theoretical framework for analyzing complex linguistic networks: Springer*, 2016, pp. 47-79.
- [15] E. L. Lin and G. L. Murphy, "Thematic relations in adults' concepts," *Journal of experimental psychology: General*, vol. 130, no. 1, p. 3, 2001.
- [16] F. Liu, D. Pennell, F. Liu, and Y. Liu, "Unsupervised approaches for automatic keyword extraction using meeting transcripts," in *Proceedings of human language technologies: The 2009 annual conference of the North American chapter of the association for computational linguistics*, 2009, pp. 620-628.
- [17] X. Ao, X. Yu, D. Liu, and H. Tian, "News keywords extraction algorithm based on TextRank and classified TF-IDF," in *2020 International Wireless Communications and Mobile Computing (IWCMC), 2020: IEEE*, pp. 1364-1369.
- [18] F. Liu, X. Huang, and W. Huang, "Comparing Machine Learning Algorithms to Predict Topic Keywords of Student Comments," in *International Conference on Cooperative Design, Visualization and Engineering*, 2020: Springer, pp. 178-183.
- [19] R. Campos, V. Mangaravite, A. Pasquali, A. Jorge, C. Nunes, and A. Jatowt, "YAKE! Keyword extraction from single documents using multiple local features," *Information Sciences*, vol. 509, pp. 257-289, 2020.
- [20] J. R. Thomas, S. K. Bharti, and K. S. Babu,

۱. استفاده از روش‌های رتبه‌بندی واژگان در فهرست نزدیک‌ترین همسایگی‌ها؛
۲. استفاده از ویژگی هم‌رخدادی واژگان برای شناسایی عبارات کلیدی؛
۳. استفاده از ضریب وابستگی بین واژگان با توجه به ترادف و تضاد بین آنها.

۱۰- مراجع 10- References

- [1] F. Liu, X. Huang, W. Huang, and S. X. Duan, "Performance Evaluation of Keyword Extraction Methods and Visualization for Student Online Comments," *Symmetry*, vol. 12, no. 11, p. 1923, 2020.
- [2] H. Yan, Q. He, and W. Xie, "Crnn-Ctc Based Mandarin Keywords Spotting," in *ICASSP 2020-2020 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP), 2020: IEEE*, pp. 7489-7493.
- [3] Y. Zhang, M. Tuo, Q. Yin, L. Qi, X. Wang, and T. Liu, "Keywords extraction with deep neural network model," *Neurocomputing*, vol. 383, pp. 113-121, 2020.
- [4] M. Mohammadi and M. Analouyi, "Keyword extraction in Persian documents," presented at the 13th Conference of Iran Computer Association, Kish, Iran, 2007.
- [5] H. Veisi, N. Aflaki, and P. Parsafard, "Variance-based features for keyword extraction in Persian and English text documents," *Scientia Iranica*, vol. 27, no. 3, pp. 1301-1315, 2020.
- [6] C. Zhang, "Automatic keyword extraction from documents using conditional random fields," *Journal of Computational Information Systems*, vol. 4, no. 3, pp. 1169-1180, 2008.
- [7] X. Wan and J. Xiao, "Single Document Keyphrase Extraction Using Neighborhood Knowledge," in *AAAI*, 2008, vol. 8, pp. 855-860.
- [8] D. B. Bracewell, F. Ren, and S. Kuriowa, "Multilingual single document keyword extraction for information retrieval," in *2005 International Conference on Natural Language Processing and Knowledge Engineering, 2005: IEEE*, pp. 517-522.
- [9] P. D. Turney, "Learning to extract keyphrases from text," arXiv preprint cs/021, 2002, 2013.
- [10] Y. Matsuo and M. Ishizuka, "Keyword extraction from a single document using word co-occurrence statistical information," *International Journal on Artificial Intelligence Tools*, vol. 13, no. 01, pp. 157-169, 2004.



علی امیری جزه کارشناسی ارشد خود را در رشته هوش مصنوعی از دانشگاه امام حسین(ع) در سال ۱۳۹۹ اخذ کرد. حوزه مورد علاقه وی پردازش زبان طبیعی، تحلیل های زبانی و داده های حجیم و پایان نامه او در زمینه شناسایی موجودیت های اسمی با استفاده از یادگیری عمیق و رویکرد تقویتی است. نشانی رایانامه ایشان عبارت است از:

aamirij@ihu.ac.ir

"Automatic keyword extraction for text summarization in e-newspapers," in *Proceedings of the international conference on informatics and analytics*, 2016, pp. 1-8 .

- [21] D. M. Allen, "The relationship between variable selection and data augmentation and a method for prediction," *technometrics*, vol. 16, no. 1, pp. 125-127, 1974.
- [22] M. Stone, "Cross-validators choice and assessment of statistical predictions," *Journal of the Royal Statistical Society: Series B (Methodological)*, vol. 36, no. 2, pp. 111-133, 1974.
- [23] M. Stone, "An asymptotic equivalence of choice of model by cross-validation and Akaike's criterion," *Journal of the Royal Statistical Society: Series B (Methodological)*, vol. 39, no. 1, pp. 44-47, 1977.
- [24] T. Mikolov, K. Chen, G. Corrado, and J. Dean, "Efficient estimation of word representations in vector space," *arXiv preprint arXiv:1301.3781*, 2013.
- [25] C. Manning and R. Socher, "Natural language processing with deep learning," *Lecture Notes Stanford University School of Engineering*, 2017.
- [26] J. Hu, S. Li, Y. Yao, L. Yu, G. Yang, and J. Hu, "Patent keyword extraction algorithm based on distributed representation for patent classification," *Entropy*, vol. 20, no. 2, pp. 104, 2018.
- [27] H. Omid and S. Saeedeh, Sadidpour, "Automatic extraction of Persian short text keywords using word2vec," *Electronic and cyber defense*, vol. 8, 2, pp. 105-114, 2020.



محمد رضا حسینی آهنگر دکترای

خود را در رشته هوش مصنوعی و رباتیک از دانشگاه علم و صنعت در سال ۱۳۹۰ اخذ کرد. تخصص ایشان هوش مصنوعی است و در حال حاضر به عنوان

عضو هیأت علمی دانشکده مهندسی رایانه دانشگاه امام حسین(ع) به تدریس دروس هوش مصنوعی مشغول هستند. از سال ۱۳۹۱ ریاست دانشگاه جامع امام حسین(ع) به عهده ایشان است.

نشانی رایانامه ایشان عبارت است از:

mrhasani@ihu.ac.ir