

Keywords Extraction from Persian Thesis Using Statistical Features and Bayesian Classification

Behzad Hejazi¹  & Jalal A. Nasiri^{2*} 

Vol. 12, No. 6, Tome 66
pp. 339-367
January & February 2022

Received: 23 October 2019
Received in revised form: 22 December 2019
Accepted: 13 January 2020

Abstract

Keyword extraction aims to extract words that are able to represent the corpus meaning. Keyword extraction has a crucial role in information retrieval, recommendation systems and corpora classification. In Persian language, keyword extraction is known as hard task due to Persian's inherent complication. In this research work, we aim to address keyword extraction with a combination of statistical and Machine Learning as a novel approach to this problem. First the required preprocessing is applied to the corpora. Then three statistical methods and Bayesian classifier was utilized to the corpora to extract the keywords pattern. Also, a post processing methods was used to decrease the number of True Positive outputs. It should be pointed out that the built model can extract up to 20 keywords and they will be compared with keywords in the corresponding corpus. The evaluation results indicate that the proposed method, could extract keywords from scientific corpora (Specifically Thesis and Dissertations) with a good accuracy.

Keywords: Extraction, Bayesian Classification, statistical features, preprocessing, post-processing

-
1. M.Sc Student of Artificial Intelligence, Islamic Azad University North Tehran Branch, Tehran, Iran; ORCID: <https://orcid.org/0000-0001-6772-8715>
 2. Corresponding author: Department of Computational Linguistics, Information Science Research Department, Iranian Research Institute for Information Science and Technology (IRANDOC), Tehran, Iran; Email: j.nasiri@irandoc.ac.ir, ORCID: <https://orcid.org/0000-0003-1821-5037>

1. Introduction

Automated keyword extraction is the process of identifying document terms and phrases that can appropriately represent the subject of our writing. With the proliferation of digital documents today, extracting keywords manually can be impractical. Many applications such as auto-indexing, summarization, auto-classification, and text filtering can benefit from this process since the keywords provide a compact display of the text. Automated keyword generation can be broadly classified into two categories: keyword allocation and keyword extraction.

In keyword allocation, a set of potential keywords is selected from a set of controlled vocabularies, while keyword extraction examines the words in the text. Keyword extraction methods can be broadly classified into four groups: statistical approaches, linguistic approaches, machine learning approaches, and hybrid approaches.

2. Literature Review

working on Persian words is a big challenge for the paucity of sufficient research. The inadequacy of text pre-processing programs has made it more complex than the Latin language. Also, the presence of large dimensions of input data is one of the challenges that has always arisen in such researches and this problem becomes more apparent due to the variety of Persian written forms (Gandomkar, 2017, p. 233:256). In Moin Maedi's article (2015, p. 34:42) A method for extracting keywords in Persian language is presented. This article extracts keywords from each text separately and without seeing another text as training data.

In the article by Mohammad Razaghnoori (2017, P. 16:27) using the Word2Vec method and the TIF-IDIF frequency, they created a question and answer system in Persian, which is a new work due to the use of Word2Vec in Persian. However, with size reduction techniques and Word2Vec, this 72% success rate can be enhanced in the future.

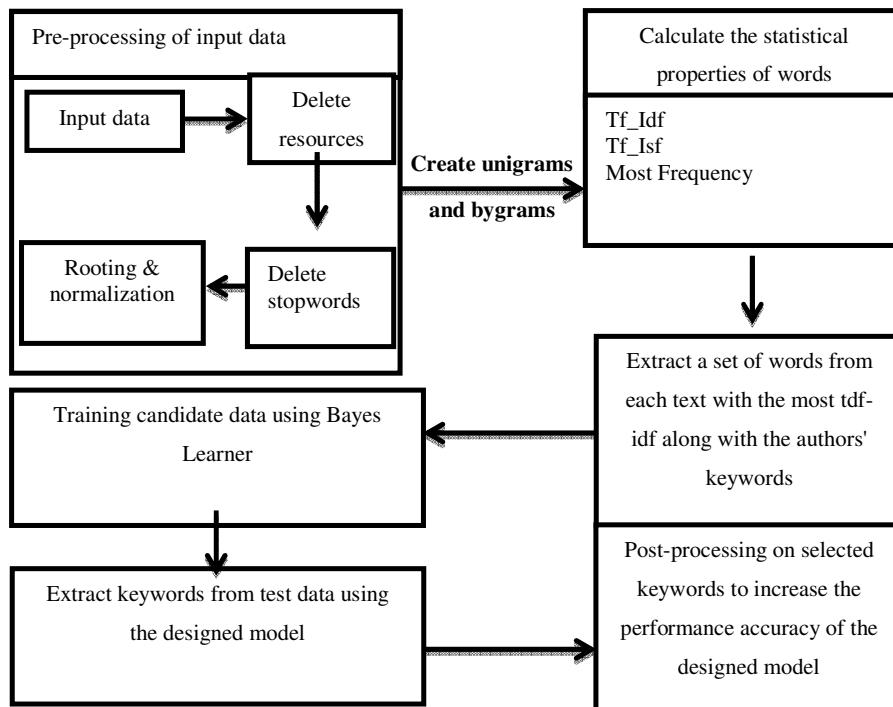
3. Methodology

Accordingly, the current paper examines the integration of statistical keyword extraction methods with the Naive Bayes Classifier. Initially, we integrated input texts which are dissertations in Persian by using preprocessing (deletion of stop words, etymology, etc.) methods. Then, using the available statistical features, each word has been given a certain weight. Then, the valuable words of each text were selected and the proposed model was taught using the selected category, then the selected words were processed by the trained model, and at the end, the words extracted from the final model were evaluated using the keywords suggested by the authors themselves. Figure 1 depicts all the steps performed.

4. Results

Literature review shows that this is the first time that these combinations are used to extract Persian keywords, so that unlike other studies, each text is as a sample for category input and words as its properties, however, in this paper the words of each text input are categorized and words are extracted using statistical methods that are considered as features. The choice of keywords by the authors has always been a personal decision and people may not make a single decision to choose a set of words for a single text.

Figure 1
Proposed research framework for keyword extraction



The current paper attempts to create a model and program with a new approach, due to the small number of input documents, which to extract keywords without dependence on the orientation of dissertations and the meaning of their words and only by using statistical features of words in each text. According to Tables 1 and 2, the developed model is able to extract a maximum of 20 keywords from each dissertation with an overall accuracy of 98.1%, in best condition which that is the use of a maximum frequency feature. The keywords written in each dissertation with 84% and 98% accuracy, correspond to one-word and two-word expressions, respectively.

Table 1*Evaluation criteria for Bayesian outputs in different states of statistical Features*

Statistical Features	Accuracy	Recall	F1-Score	Precision
Tf_Idf, Most Frequent, Tf_Isf	97.2%	0.98	0.98	0.98
Most Frequent	98.1%	0.982	0.99	0.99
Tf_Idf, Tf_Isf	99.8%	0.91	0.94	0.99

Table 2*Evaluation of post-processing test data for outputs that have been categorized by keyword.*

Step	Statistical Features	Recall	F1-Score	Precision	Number of words	Number of keywords that selected by writers
Uni-Grams	Most Frequent	0.84	0.323	0.2	210	42
By-Grams	Most Frequent	0.98	0.888	0.8	158	34

استخراج کلیدواژگان پایان‌نامه فارسی با استفاده از ویژگی آماری و دسته‌بند بیز

بهزاد حجازی^۱، جلال‌الدین نصیری^{۲*}

- دانشجوی کارشناسی ارشد هوش مصنوعی، دانشگاه آزاد اسلامی واحد تهران شمال، تهران، ایران.
- استادیار دانشکده علوم ریاضی، دانشگاه فردوسی، مشهد، ایران.

تاریخ پذیرش: ۱۳۹۸/۱۰/۲۲

تاریخ دریافت: ۱۳۹۸/۰۸/۱

چکیده

هدف اصلی استخراج کلمات کلیدی انتخاب مجموعه‌ای از لغات در متن است که می‌تواند موضوع اصلی متن را بازگو کند. استخراج کلیدواژگان در بازیابی اطلاعات، سیستم‌های پیشنهاده‌نده متنی و دسته‌بندی متون، نقش مهم را ایفا می‌کند. در زبان فارسی با توجه به پیچیدگی ذاتی زبان فارسی استخراج کلیدواژگان به مراتب دشوارتر شده است. در این پژوهش سعی شده است با رویکرد نوین ترکیبی آماری و یادگیری ماشین به استخراج کلیدواژگان پرداخته شود. ابتدا با توجه به ساختار زبان فارسی پیش‌پردازهای لازم برای حذف کلمات و عالم نگارشی صورت می‌گیرد. سپس با استفاده از سه نوع ویژگی آماری و دسته‌بند بیز سیستم به صورت خودکار الگوی کلمات کلیدی با کلمات عادی را آموزش می‌بیند. همچنین پس‌پردازش کارا برای کم کردن کلمات مثبت کاذب در چارچوب پیشنهادی طراحی شده است. گفتنی است که مدل ساخته شده قادر به شناسایی تعداد حداقل ۲۰ کلیدواژه در هر پایان‌نامه است و این کلمات با کلیدواژگان نوشته شده در هر متن مناسبی توانسته است کلمات کلیدی ارزیابی‌های متنوع نشان می‌دهد روش پیشنهادی با دقت مناسبی توانسته است کلمات کلیدی نوشتارهای فارسی علمی (پایان‌نامه و رساله) را استخراج کند.

واژه‌های کلیدی: استخراج کلیدواژگان، دسته‌بند بیز، ویژگی‌های آماری، پیش‌پردازش، پس‌پردازش.

۱. مقدمه

استخراج خودکار کلمات کلیدی فرایند شناسایی اصطلاحات و عبارات سندی است که به‌طور مناسب می‌تواند نماینده موضوع نوشتار ما باشد. امروزه با گسترش تعداد اسناد دیجیتال، استخراج کلمات کلیدی به‌صورت دستی می‌تواند یک موضوع غیرعملی باشد. از زمانی که کلیدواژگان نمایشی فشرده از متن را فراهم کرده‌اند، بسیاری از کاربردها مانند نمایه‌سازی خودکار، خلاصه‌سازی، طبقه‌بندی خودکار و فیلترینگ متن می‌توانند از این فرایند سود ببرند. تولید کلمات کلیدی خودکار را می‌توان به‌طور گسترده به دو دسته تقسیم کرد: تخصیص کلیدواژه و استخراج کلمات کلیدی. در تخصیص کلمه کلیدی، مجموعه‌ای از کلمات کلیدی محتمل از مجموعه‌ای واژگان کنترل شده انتخاب می‌شوند، درحالی که در استخراج کلیدواژه، کلمات موجود در نوشتار موربدبررسی قرار می‌گیرند. روش‌های استخراج کلیدواژه می‌توانند به‌طور عمده به چهار گروه تقسیم شوند: رویکردهای آماری، رویکردهای زبانی، رویکردهای یادگیری ماشین و رویکردهای ترکیبی.

دسته‌بندی متن بخش مهمی در داده‌کاوی است که متنی را به یک یا چند کلاس یا دسته‌های از پیش تعیین شده اختصاص می‌دهد. چندین شکل از مجموعه‌های متنی مانند مقالات خبری، کتابخانه‌های دیجیتالی و صفحات وب متابع مهم اطلاعات هستند. از این‌رو، طبقه‌بندی متن یک جهت تحقیقاتی مهم در علوم کتابداری، اطلاع‌رسانی و علوم کامپیوتر است. بسیاری از کاربردهای داده‌کاوی را می‌توان به‌منزله یک مشکل طبقه‌بندی متن مدل‌سازی کرد. این کاربردها فیلترینگ اخبار، سازماندهی اسناد، بازیابی، تحلیل احساسات، فیلترینگ هرزنامه‌ها و استخراج کلیدواژگان است.

به‌دلیل نبود پژوهش‌های کافی، کار بر روی کلمات فارسی چالشی بزرگ است. کافی نبودن برنامه‌های مرتبط با پیش‌پردازش متن این امر را نسبت به زبان لاتین پیچیده‌تر کرده است. همچنین، وجود ابعاد بزرگ داده‌های ورودی یکی از چالش‌هایی است که همواره در این گونه پژوهش‌ها به وجود آمده است و به‌دلیل تنوع در صورت‌های نگارشی زبان فارسی (گندمکار، ۱۳۹۶، صص. ۲۲۲-۲۵۶) این مشکل، بیشتر نمایان می‌شود.

باتوجه به مسائل یادشده، این مقاله به بررسی ادغام روش‌های استخراج آماری کلمات کلیدی با دسته‌بند بیز ساده می‌پردازد. در ابتدا با استفاده از پیش‌پردازش (حذف

ایستوازدها، ریشه‌یابی و ...) به یکپارچه‌سازی متن‌های ورودی که پایان‌نامه‌هایی به زبان فارسی است، پرداخته شده است. سپس با استفاده از ویژگی‌های آماری موجود، به هر کلمه وزن مشخصی اعطا شده است. در مرحله بعد کلمات بالرزش هر متن انتخاب و با استفاده از دسته‌بند انتخابی به آموزش مدل مطرح شده پرداخته شده است. سپس پس‌پردازش کلمات انتخاب‌شده توسط مدل آموزش‌دیده صورت می‌پذیرد. در انتهای کار با بهره‌گیری از کلیدواژگان مطرح شده توسط خود نویسنده‌گان به ارزیابی کلمات استخراج شده از مدل نهایی پرداخته شده است.

باتوجه به مطالعات انجام شده، این اولین بار است که برای استخراج کلیدواژگان فارسی از این ترکیب‌ها استفاده می‌شود، به طوری که برخلاف آنکه همانند سایر پژوهش‌های صورت‌گرفته، هر متن به منزله یک نمونه برای ورودی دسته‌بند و کلمات به منزله ویژگی‌های آن باشند، در این پژوهش کلمات هر متن ورودی‌های دسته‌بند هستند و با بهره‌گیری از روش‌های آماری که به منزله ویژگی‌ها در نظر گرفته شده‌اند، به استخراج کلمات پرداخته می‌شود. در این پژوهش فرض بر این است که با بهره‌گیری از اطلاعات آماری داده‌های ورودی، مدلی هوشمند برای استخراج کلیدواژگان با دقت بالا ساخته شده است.

ادامه مقاله بدین شکل سازمان‌دهی شده است که در بخش دوم به بررسی پژوهش‌های انجام‌شده در زبان لاتین و فارسی پرداخته شده است. در بخش سوم به شرح چارچوب پیشنهادی پرداخته و قسمت‌های مختلف برنامه تعریف شده است. در بخش چهارم نتایج عملی مورد بررسی قرار می‌گیرد. بدین شکل که نحوه جمع‌آوری داده‌ها و عملیاتی که به روی آن‌ها صورت می‌پذیرد به طور کامل شرح داده می‌شود و نتایج آن موردن تحلیل قرار می‌گیرد. در بخش پنجم نتیجه‌گیری مقاله موردن بحث قرار گرفته و در بخش آخر به ایده‌های آینده برای انجام پژوهش‌های مختلف در این مبحث پرداخته شده است.

۲. مرور ادبیات

۱-۲. پژوهش‌های انجام‌شده در حوزه زبان لاتین

در تازه‌ترین پژوهش‌هایی که در این حوزه انجام شده می‌توان به آثارنو دی^۱ (2018, p.92, 105) اشاره کرد که در آن باتوجه به تجزیه و تحلیل احساسات به ارزیابی عملکرد

محصولات و خدمات، از محتویات تولیدشده توسط کاربر می‌پردازد. رویکردهای تجزیه و تحلیل احساسات مبتنی بر زنجیرهای نحوی، زمانی که اطلاعات آموزشی کافی نیست بر مبنای یادگیری ترجیح داده می‌شوند. واژگان موجود فقط دارای یگانگی و نمرة احساسات خود هستند. مشاهده می‌شود که در استخراج عبارات چندکلمه‌ای، کلمه‌های احساس شده به وسیلهٔ ترکیبی از تک‌کلمه‌ای‌ها^۳ با تقویت‌کننده‌ها نتایج بهبودیافته را نشان می‌دهد. چنین واژگان عاطفی در دسترس عموم نیست. این مقاله روشی را برای ایجاد عباراتی به نام جملات چندبخشی^۴ ارائه می‌دهد. رویکرد مبتنی بر قاعدة پیشنهادی، نمرات احساسات چندبخشی را از یک مجموعهٔ تصادفی شامل بررسی محصول و رتبه‌بندی عددی مربوطه در مقیاس پنج‌گانه استخراج می‌کند. به علاوهٔ روش‌هایی که گفته شده، قوانین اضافی دیگری نیز وجود دارند که می‌توانند به همراه روش‌های گفته شده به کار گرفته شوند و کیفیت عبارات کاندید را بهبود بخشدند. چنین ایده‌هایی در مقالهٔ رافائل^۵ (2010, p.190, 193) به کار گرفته شده است. در روش آن‌ها اگر یک عبارت در ۴۰۰ کلمهٔ اول متن ظاهر نشود آنگاه به منزلهٔ عبارت کاندید چهت پیش‌بینی کلیدواژگان درنظر گرفته نمی‌شود. ایدهٔ پشت شرط آن‌ها این است که اگر یک عبارت در ابتدای متن به کار نزود پس احتمال کمی وجود دارد که عبارت کلیدی باشد. یک رویکرد مشابه توسط نیومن^۶ (2010, p.150, 153) به کار گرفته شده که محدودهٔ استخراج عبارات کاندید را به ۲۰۰۰ کلمهٔ اول هر متنی محدود می‌کند. در مقالهٔ لانگ^۷ (2010, p.21, 26) نیز رویکردی به کار گرفته شده که در آن، انتخاب عبارات کاندید را به بخش‌های چکیده، مقدمه، کارهای مرتبط، عناوین موجود در متن و اولین خط هر پاراگراف برای ادامهٔ متن محدود می‌کند. این روش کیفیت نتایجشان را بهبود بخشیده است.

لیتواک^۸ (2008, p.17, 24) به بررسی کارایی روش‌های استخراج کلمات کلیدی متمرکز و بدون نظارت بر استخراج خلاصهٔ اسناد پرداختند. برای نشان دادن اسناد متنی، نمایندگی نحو براساس گراف استفاده شده است.

در مطالعهٔ دیگری هوان^۹ (2006, p.275, 284) یک الگوریتم استخراج کلمات کلیدی خودکار ارائه داد که می‌توان آن را در هر دو روش یادگیری بدون نظارت و تحت نظارت استفاده کرده است. الگوریتم ارائه شده هر سند متن را به منزلهٔ شبکه‌ای معنایی مدل می‌کند. عبارات کلیدی براساس داینامیک ساختاری شبکهٔ معنایی استخراج می‌شوند.

ترنی^۹ (2002) روش استخراج کلمات کلیدی را بهمنزله یک مسئله یادگیری ماشین نظارت شده ارائه کرد. در این روش از ۹ ویژگی برای امتیاز دادن به عبارات کاندید استفاده شده است. سپس عبارات کلیدی براساس ویژگی‌ها از میان کاندیدهای استخراج می‌شوند. ویتن^{۱۰} (1999) یک الگوریتم استخراج کلمات کلیدی ساده و کارآمد (KEA) را ارائه داد که از الگوریتم بیز ساده برای استخراج کلمات کلیدی مبتنی بر دامنه استفاده می‌کند. در این روش، عبارات کلیدی ممکن باید با روش‌های لغوی تعیین شود و عبارات کلیدی خوب توسط الگوریتم یادگیری ماشین بهدست آید. کمر^{۱۱} (2005, p.657, 669) اثربخشی روش‌های خودکار استخراج کلیدواژه‌ها و روش‌های یادگیری را در مقالات علمی به زبان انگلیسی مورد بررسی قرار دادند. روش‌های استخراج کلمات کلیدی با روش‌های مختلف یادگیری ماشین ارزیابی می‌شوند و نتایج تجربی نشان می‌دهد که الگوریتم C4.5 دارای بالاترین نرخ پیش‌بینی برای دامنه است. در تحقیقات عملی صورت گرفته در آنان^{۱۲} (2016, p.232, 247) که برای استخراج کلمات کلیدی بر روی کتابخانه دیجیتال ACM صورت گرفته است از ترکیب روش‌های مختلف و مقایسه آن‌ها با هم پرداخته است که بهترین نتایج مربوط به ترکیب الگوریتم‌های دسته‌بندی‌کننده^{۱۳} و جنگل تصادفی^{۱۴} است.

۲-۲. پژوهش‌های انجام‌شده در حوزه زبان فارسی

همچنین مقالاتی در این حوزه در زبان فارسی به چاپ رسیده است که درادامه به اختصار توضیحاتی درموردشان داده می‌شود. در مقاله معین ماعدی^{۱۵} (2015, p.34, 42) روشی برای استخراج کلمات کلیدی در زبان فارسی ارائه شده است که این مقاله به استخراج کلیدواژگان از هر متن به صورت جداگانه و بدون دیدن متن دیگری به عنوان داده‌های آموزش، می‌پردازد. روش پیشنهادی مقاله با ترکیب ویژگی‌های آماری، بردار موقعیت تکرار را برای هر کلمه و قوانین ساده زبان ارائه داده است. روش پیشنهادی در سه مرحله ارائه شده است که پیش‌پردازش متن، وزن‌دهی به کلمات، ایجاد بردار موقعیت تکرار کلمه در متن و محاسبه شباهت و استفاده از اسناد تشکیل شده برای این مطالعه در زبان فارسی موردارزیابی قرار گرفته است. از این روش نتیجه می‌شود که زبان فارسی را می‌توان بهمنزله دستاوردهای جدید و امیدوارکننده برای زبان‌های دیگر درنظر گرفت. با وجود این، در

بررسی نتایج حاصل از نتایج به دست آمده می توان مشکلات زیر را در بررسی نتایج به دست آورد: ۱. انتخاب کلمات توصیفی به جای کلمات کلیدی توسط نویسنده، ۲. استفاده از اختصار انگلیسی به جای کلمه انتخاب شده فارسی. مانند کلمه «بردارهای ماشین پشتیبان» در کلمات کلیدی، اما استفاده از کلمه «SVM» در متن، ۳. استفاده کوچک در متن یا به منزله اشاره مختصر به ویژه کلمات مهم در مقالات، معمولاً در ابتدای متن مانند عنوان، انتزاعی یا مقدمه بیشتر برای حضور وجود دارد و اغلب نویسندها به ضمایر یا مرجع آنها اشاره دارند، ۴. وجود کلمات و عبارات مربوط به سبک نوشتاری نویسنده، ۵. ایجاد ابهام در مرحله تبدیل فایل PDF به متن، ۶. احترام نگذاشتن به زبان فارسی.

در مقاله محمد رzagوری^{۱۰} (2017, p.16, 27) با استفاده از روش ورد تو وک^۷ و بسامد تی اف - ایدیاف سیستم پرسش و پاسخ را در زبان فارسی به وجود آورده که با توجه به استفاده از ورد تو وک در زبان فارسی کاری تازه بوده است، گرچه با تکنیکهای کاهش ابعاد و ورد تو وک می توان این ۷۲ درصد میزان موقفيت را در آينده ارتقا داد.

در مقاله بهنام ثابتی^{۱۱} (2018) به شرح میراثتکس^{۱۲} پرداخته اند که می شود گفت بزرگترین مجموعه اسناد در حوزه زبان فارسی است که دارای بیش از ۲/۴ میلیون سند و دارای بیش از ۱/۴ میلیارد کلمه است، اگرچه هنوز به نسبت زبان انگلیسی دارای مقدار کمتری است. در مقاله مرتضی اخوت^{۱۳} (2010, p.94,101) برای برچسب‌گذاری متنون فارسی نحوی از روش مارکوف^{۱۴} استفاده کردند. به طوری که برای برچسب‌گذاری کلمات فارسی ابتدا متن فارسی را به انگلیسی ترجمه می کردند، سپس متن انگلیسی را برچسب‌گذاری می کردند. ایراد این روش این است که در هین ترجمه یک سری کلمات نامشخص باقی می ماندند که دقت برای این کلمات پایین است. اما برای کلماتی که درست ترجمه شده اند که بخش عمده ای از متن هستند دارای دقت ۹۸ درصد است و به طور ميانگين با احتساب کلمات درست و غلط تشخيص داده شده دارای دقت ۹۷ درصد است.

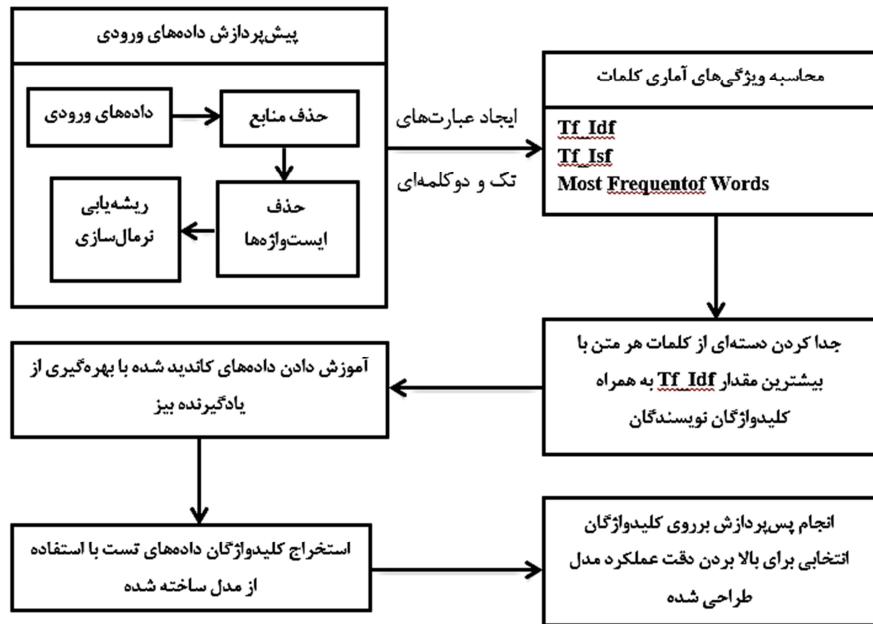
محمد^{۱۵} (2011, p.112,116) با استفاده از متن های روزنامه همشهری به استخراج کلمات کلیدی يك‌كلمه‌اي و دو‌كلمه‌اي پرداخته است. اما جواب‌هایي که استخراج شده مورد ارزیابی واقع نشده است. حال که بررسی پژوهش‌های پیشین انجام شده است، نوآوری‌های انجام شده در این پژوهش موردنحلیل قرار گرفته است. این پژوهش سه نوآوری بارز دارد که

بدین شرح است:

- ۱) استخراج کلیدواژه‌ها بر مبنای ویژگی‌های آماری هر کلمه به جای ویژگی‌های هر سند،
- ۲) کاربردی‌سازی الگوریتم‌های استخراج کلیدواژه لاتین برای متون علمی فارسی با توجه به پیچیدگی‌های زبان علمی فارسی،
- ۳) طراحی پس‌پردازش مناسب جهت کارا کردن چارچوب پیشنهادی با توجه به ویژگی‌های زبان فارسی.

۳. چارچوب پیشنهادی

در این قسمت راهکار پیشنهادی برای استخراج کلیدواژگان از متون علمی توضیح داده شده است. نمودار مراحل پیشنهادی در شکل ۱ به تصویر کشیده شده است. همان‌گونه که از شکل مشخص است، ابتدا داده‌های ورودی دریافت شده‌اند، سپس پیش‌پردازش‌های لازم بر روی داده‌ها انجام شد. در گام بعدی عبارت‌های تک‌کلمه‌ای و دوکلمه‌ای ایجاد می‌شوند و ویژگی‌های آماری آن‌ها به دست می‌آید. سپس انتخاب کلمات برای داده‌های آموزش انجام شده است. در مرحله بعد به وسیله کلمات انتخابی و یادگیرنده بیز به استخراج کلمات از داده‌های تست پرداخته شده است. سپس پردازش‌های لازم بر روی کلمات انتخابی مدل یادگیرنده پرداخته می‌شود تا بدین وسیله دقت عملکرد برنامه بهبود یابد و کلمات کلیدی برنامه استخراج شوند. در گام آخر، با بهره‌گیری از معیارهای ارزیابی استاندارد تعریف شده که در ادامه مقاله به شرح آن‌ها نیز پرداخته شد، سعی می‌شود دقت عملکرد مدل ساخته شده نمایش داده شود.

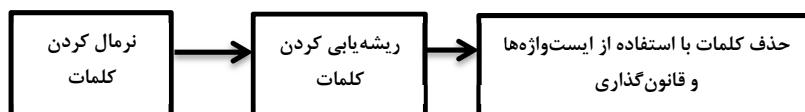


شکل ۱: چارچوب پیشنهادی پژوهش برای استخراج کلیدواژگان

Figure 1: Proposed research framework for keyword extraction

۱-۳. پیش‌پردازش متن

پیش‌پردازش متن به مفهوم یکدست کردن متن، شناسایی و ازبین بردن کلماتی در متن است که به منزله واژه‌ای مهم در متن تلقی نمی‌شوند و عموماً این کلمات در تمامی اسناد موجود است و باعث تمایز نوشه‌ها با یکدیگر نمی‌شوند.



شکل ۲: نمودار مرحله‌های پیش‌پردازش داده‌ها

Figure 2: Chart of data preprocessing steps

برای پیش‌پردازش متن می‌شود از ریشه‌یابی و نرمال کردن کلمات استفاده کرد. این امر باعث می‌شود تا کلمات تبدیل به ریشه‌شان شوند و تمامی فاصله‌ها و نیم‌فاصله‌ها استاندارد شوند. برای انجام این کار از کتابخانه هضم در زبان برنامه‌نویسی پایتون استفاده شده است. این کتابخانه توسط گروه سبجه ساخته شده و به صورت رایگان در اختیار کاربران قرار گرفته است. نمونه‌ای از تغییرات کلمات، قبل و بعد از استفاده از این برنامه در جدول ۱ آمده است.

جدول ۱: مقایسه کلمات قبل و بعد از ریشه‌یابی توسط برنامه هضم

Table 1: Compare words before and after rooting by Hazm program

نمونه اولیه کلمات	کلمات نرمال و ریشه‌یابی شده
خصوصی‌سازی	خصوص‌ساز
سازمانی	سازمان
تکنیک داده‌کاوی	تکنیک داده‌کاو
دانش‌آموزان	دانش‌آموز

یکی از ایرادات این قسمت آن است که این برنامه کلماتی مانند «استخوان» را نیز تبدیل به «استخو» می‌کند، اما به دلیل زیاد بودن واژگان و حجم زیاد ابعاد به‌اجبار باید از این برنامه استفاده کرد.

ایست‌واژه‌ها کلماتی مانند حروف اضافه و افعال هستند که در پیش‌پردازش از متن ما حذف می‌شوند. نمونه‌ای از این لغات که در این پژوهش مورد استفاده قرار گرفته است در جدول ۲ آمده است.

جدول ۲: نمونه‌ای از کلماتی که به منزله ایست‌واژه به‌کار گرفته شده‌اند

Table 2: An example of words used as a StopWord

از	به	و	آن	همچنین	را	شد	تا	می	سپس	نمودار	جدول	=
این	ولی	آن	همچنین	را	شد	تا	می	سپس	نمودار	جدول	+	تا

در این قسمت بعضی مواقع برای یکپارچگی متن قوانینی توسط خود پژوهشگران علاوه‌بر حذف ایست‌واژه‌ها نیز وضع می‌شود که در این پژوهش نیز شرط‌هایی برای واژگان

گذاشته شده، که شامل این موارد است: ۱) حذف اعداد، ۲) حذف علائم نگارشی، ۳) حذف منابع و صفحات ابتدایی متون که شامل تقدیر و تشکر و تعهدنامه‌ها هستند. یکی از موارد مشکل‌ساز برای انجام این کار در زبان فارسی به‌طور نمونه آن است که زبان‌های برنامه‌نویسی به‌طور مستقیم قادر به شناسایی و حذف اعداد به زبان فارسی نیستند و باید ابتدا آن‌ها را به زبان انگلیسی برگرداند و سپس به حذف آن‌ها اقدام کرد.

۲-۳. استخراج عبارات تک‌کلمه‌ای و دوکلمه‌ای

اولین اقدامی که در شناسایی عبارات کاندید درنظر گرفته می‌شود آن است که آیا کلمات تک استخراج شوند یا دنباله‌ای از کلمات. استخراج عبارات چندکلمه‌ای یا همان N-gram یکی از روش‌های رایج در استخراج عبارات کلیدی است. مفهوم و هدف N-gram استخراج گروه‌هایی از n کلمه‌ای پشت‌سر هم در یک جمله است. متغیر n تعداد کلمات را تعیین می‌کند. در این پژوهش باتوجه به آنکه اکثریت کلیدواژگان مطرح شده توسط نویسنده‌گان یک‌کلمه‌ای و دوکلمه‌ای هستند، عبارات تک‌کلمه‌ای و دوکلمه‌ای استخراج شده‌اند.

۳-۳. وزن‌دهی به کلمات

در این بخش با استفاده از ویژگی‌های آماری به کلمات مربوط به هر سند وزن مشخصی تخصیص داده می‌شود. در این روش برای تعیین کلمات کلیدی از روش‌های آماری فرکانس کلمه و معیار TF-IDF و معیار TF-ISF استفاده شده است. روش استخراج کلیدی آماری می‌تواند مستقل از دامنه باشد و نیازی به اطلاعات آموزشی نیز نداشته باشد (Beliga, 2015, vol.1/p.20).

۱-۳-۳. بسامد کلمه - معکوس بسامد سند^{۲۳}

این روش فرکانس کلمه در یک سند را ارائه می‌دهد. معیار بهتر این است که فرکانس کلمه در سند با مقایسه ارتباط آن با اسناد دیگر به‌دست آورده شود. این معیار بسامد کلمه - معکوس بسامد سند یا به اختصار تی.اف - آی.دی.اف نام دارد که میزان خاص بودن کلمه در متن موردنظر را نشان می‌دهد. برای محاسبه آن از فرمول زیر استفاده می‌شود (Lott,

.(2012

$$TFIDF(t, d, n, N) = TF(t, d) \times IDF(n, N) \quad (1)$$

تی - اف از رابطه زیر و از تقسیم فرکانس کلمه بر تعداد کل کلمات سند به دست می آید.

ای - دی - اف نیز از رابطه زیر محاسبه می شود که در آن، N تعداد اسنادی از پیکره است که دارای کلمه موردنظر است و n تعداد کل اسناد پیکره است.

$$TF(t, d) = \frac{m_{td}}{\sum_k m_{kd}} \quad (2)$$

$$IDF(n, N) = \log \frac{(N + 1)}{(n + 1)} \quad (3)$$

با استفاده از این معیار کلماتی از یک سند که در اسناد دیگر کمتر استفاده شده اند امتیاز بالاتری دارند.

۲-۳-۳. استخراج کلمات کلیدی براساس تکرار^{۲۴}

شایع ترین روش استخراج کلمه کلیدی (MF) است، مؤلفه های متدالول در متن را به منزله کلمات کلیدی مشخص می کند. برای نشان دادن اسناد متن، ماتریس اصطلاحی استفاده می شود. در این بازنمایی، مقدار رخداد اصطلاح برای یک جمله خاص ارزش یک یا صفر را به ترتیب مربوط به وقوع یا عدم رخداد اصطلاح در آن جمله می گیرد. براساس این نمایندگی، نمره فرکانس برای هر اصطلاح با محاسبه تعداد وقایع اصطلاح در ماتریس محاسبه می شود. برای هر اصطلاح (t_k)، نمره فرکانس به صورت زیر تعریف می شود (Rossi, 2014, p.17, 37).

$$MF(t_k) = \sum_{s_l \in S} occurrence_{t_k, s_l} \quad (4)$$

۳-۳-۳. بسامد کلمه - معکوس بسامد جمله^{۲۵}

بسامد کلمه - معکوس بسامد جمله اندازه ای آماری است که سازگاری با «TF-IDF» را به جملات متن می دهد. در اصطلاح TF-ISF، استخراج کلمات کلیدی براساس فرکانس هر جمله از متن سند به منزله یک بردار از وزن درنظر گرفته شده است (Neto, 2000, p.41, 55). در اندازه گیری TF-ISF، نرخ t_k از طریق فراوانی معکوس جمله S به وسیله فراوانی محاسبه

می شود (Fiori, 2014)

$$TF-ISF(S) = \sum_{t \in S} freq(t_k) \times isf(t_k) \quad (5)$$

$$isf(t_k) = 1 - \frac{\log(n(t_k))}{\log(n)} \quad (6)$$

جایی که n نشان دهنده تعداد جملات در مجموعه استناد است. در این پژوهش هر کلمه یک جمله در نظر گرفته شده است.

۴-۳. استفاده از یادگیرنده بیز ساده

استخراج عبارات کلیدی می تواند به منزله فرایند یادگیری ماشین در نظر گرفته شود. این گونه روش ها در کل به صورت زیر عمل می کنند: در ابتدا یک مجموعه از سند های آموزشی برای سیستم فراهم می شود که کلمات کلیدی به صورت دستی برای هر کدام از سند های متنی این مجموعه انتخاب شده است. سپس سیستم با استفاده از این مجموعه آموزش داده می شود و بعد از آن سیستم این توانایی را دارد که برای استناد جدید، کلمات کلیدی انتخاب کند. در روش های یادگیری ماشین یک الگوریتم یادگیری، مانند ماشین های بردار پشتیبانی، بیز ساده، درخت تصمیم گیری، برای ساخت یک مدل طبقه بندی استفاده می شود. در ساخت مدل، یک مجموعه آموزش استناد با برچسب استفاده می شود و مدل دارای یک مجموعه آزمون از استناد معتبر است.

۴-۴. بیز ساده

الگوریتم بیز یک الگوریتم یادگیری آماری است که بر اساس قضیه بیز است. این الگوریتم در زمینه معرفی، استفاده و یادگیری دانش احتمالی واضح است. بر اساس طبقه بندی فرض استقلال شرطی آن است که هزینه های محاسباتی موردنیاز را ساده می کند. از این رو، الگوریتم می تواند به خوبی قیاس کند و به راحتی می تواند در تعدادی از حوزه ها، از جمله مشتق متنا استفاده شود. این الگوریتم دقت بالا و سرعت را در مجموعه داده های بزرگ به ارمغان می آورد و نتایج مشابهی را با الگوریتم های طبقه بندی دیگر مانند تصمیم گیری درختان و شبکه های عصبی ارائه می دهد.

$$P(c|x) = \frac{P(x|c)P(c)}{P(x)}$$

کلاس احتمال پیشین همسایگی
 پیش‌بینی احتمال پیشین

(۷)

$$P(c|x) = P(x_1|c) \times P(x_2|c) \times \dots \times P(x_n|c) \times P(c) \quad (۸)$$

۵-۳. معیارهای ارزیابی

برای ارزیابی عملکرد پیش‌بینی شده در روش‌های استخراج کلمات کلیدی آماری، الگوریتم‌های طبقه‌بندی و روش‌های ترکیبی از دقت طبقه‌بندی، اندازه‌گیری F و منطقه تحتمنحنی (AUC) به منزله معیارهای ارزیابی استفاده می‌شود. دقت طبقه‌بندی (ACC) یکی از معیارهای به‌کاررفته در بررسی طبقه‌بندی‌هاست. این نسبت، نسبت مثبت واقعی و منفی واقعی در کل تعداد موارد است که توسط معادله نشان داده شده است:

$$ACC = \frac{TN + TP}{TP + FP + FN + TN} \quad (۹)$$

جایی که TN، TP و FN نشان‌دهنده تعداد منفی واقعی، تعداد مثبت واقعی، تعداد مثبت کاذب و تعداد منفی‌های کاذب است. دقت (PRE) درصد مثبت واقعی نسبت به مجموع مثبت واقعی و مثبت کاذب است که توسط معادله زیر ارائه شده است:

$$PRE = \frac{TP}{TP + FP} \quad (۱۰)$$

در (REC)، نسبت مثبت‌های واقعی نسبت به مثبت‌های واقعی و منفی‌های کاذب است که توسط معادله زیر داده می‌شود:

$$REC = \frac{TP}{TP + FN} \quad (۱۱)$$

مقدار F مقداری بین (۰ و ۱) است. این میانگین هارمونیک دقیق و فراخوانی است که توسط معادله زیر تعیین می‌شود:

$$F_measure = \frac{2 \times PRE \times REC}{PRE + REC} \quad (12)$$

محدوده تحت منحنی (AUC) یکی دیگر از شاخص‌های معنول برای ارزیابی طبقه‌بندی‌ها است. این برابر با احتمالی است که یک طبقه‌بندی نمونه مثبت تصادفی انتخاب شده را بالاتر از یک انتخاب منفی انتخاب می‌کند. مقدار آن از ۰ به ۱ طول می‌کشد. مقادیر بالاتر AUC نشان‌دهنده عملکرد بهتر الگوریتم‌های طبقه‌بندی است.

۴. نتایج تجربی

۱-۴. عملکرد و تحلیل تجربی

تمامی داده‌ها که شامل ۱۴۱ فایل پایان‌نامه با فرمت ورد هستند، از مجموعه ایران‌دک گردآوری شده است. این پایان‌نامه‌ها مربوط به رشته‌های علوم انسانی، فنی و مهندسی، هنر، علوم پایه و پزشکی هستند که به ترتیب هرکدام از آن‌ها دارای ۲۳، ۴۸، ۲۳ و ۲۵ پایان‌نامه از مجموع داده‌های ما را سهیم‌اند. گفتنی است که تعداد کلیدواژگان مطرح شده در جدول ۳ شامل کلیدواژگان بیش از دو کلمه‌ای نیز است، ولی در این پژوهش فقط آن دسته از عباراتی که تک‌کلمه‌ای و دو‌کلمه‌ای هستند در نظر گرفته شده‌اند.

جدول ۳: اطلاعات مربوط به داده‌های ورودی

Table 3: Input data information

گرایش	تعداد استناد	تعداد کلمات	تعداد کلیدواژگان
انسانی	۱۶	۳۰۷۱۶۲	۷۰
فنی	۴۸	۷۷۴۹۸۱	۲۳۹
هنر	۲۳	۵۳۶۹۰۳	۱۱۶
پایه	۲۱	۴۲۲۰۰۳	۱۰۰
پزشکی	۲۵	۲۹۶۹۰۹	۱۰۹
جمع کل	۱۴۱	۲۲۸۳۹۸۰	۶۸۴

ابتدا داده‌های ورودی مورد پیش‌پردازش قرار گرفته‌اند. در گام بعدی هر متن به صورت

لیستی از عبارات تک‌کلمه‌ای و دوکلمه‌ای درنظر گرفته شده‌اند. سپس ویژگی‌های آماری گفته شده برای هر کدام از این عبارات محاسبه می‌شود و ماتریسی شامل ۵ ستون ایجاد می‌شود، که سه ستون آن ویژگی‌های آماری گفته شده، یک ستون مربوط به عبارت‌های ورودی و ستون آخر که نشان هر کلمه است. اگر این کلمه توسط نویسنده پایان‌نامه کلیدواژه باشد، نشان ۱ و در غیر این صورت نشان صفر را می‌گیرد. نمونه‌ای از این ماتریس برای عبارت‌های تک‌کلمه‌ای در جدول ۴ آمده است.

جدول ۴: نمونه‌ای از ماتریس داده‌های ورودی تک‌عبارتی

Table 4: An example of a single-phrase (Unigrams) input data matrix

كلمات	Tf-idf	Tf-isf	Mf	Label
فرهنگ	۰.۰۲۴۵۷۷	۰.۰۰۸۶۱۴	۵۱	۰
استقرار	۰.۰۰۵۲۶۸	۰.۰۰۲۰۲۳	۱	۰
پرسشنامه	۰.۰۰۳۶۶۸	۰.۰۰۱۰۸۲	۳	۰
نوآور	۰.۰۰۱۹۸۷	۰.۰۰۰۵۱۳	۱	۱

در این مرحله برای کم کردن اختلاف زیاد تعداد کلیدواژه‌ها و غیرکلیدواژه‌ها، باتوجه به بیشترین مقدار ویژگی TF-IDF، تعداد ۵۰ عبارات تک‌کلمه‌ای و ۵۰ عبارات دوکلمه‌ای از هر متن به صورت جداگانه انتخاب می‌شوند و به همراه تمام کلیدواژگان مطرح شده توسط نویسنگان جدا می‌شوند. این کلمات انتخاب شده همراه دو ویژگی دیگر آن‌ها در ماتریسی قرار داده می‌شوند. پس در این حالت از هر متن حدود ۱۰۰ الی ۱۱۰ کلمه در ماتریس ساخته شده موجود است. در گام بعدی به طور تصادفی کلمات مربوط به یک متن را از ماتریس ورودی جدا کردیم و این ماتریس را به منزله داده‌های آموزش به یادگیرنده بیز داده می‌شود و تمامی کلمات متنی که جدا شده بودند به منزله داده‌های تست درنظر گرفته می‌شوند. این مرحله را ۲۱ بار انجام می‌پذیرد و هر بار کلمات یک متن را به طور کامل به منزله داده‌های تست درنظر گرفته می‌شود. در انتها باتوجه به معیارهای ارزیابی گفته شده به تحلیل نتایج پرداخته می‌شود. همان‌طور که در جدول ۵ دیده می‌شود، این آزمایشات با ترکیب‌های مختلف ویژگی‌های آماری مطرح شده نیز مقایسه شده است.

جدول ۵: میزان معیارهای ارزیابی برای خروجی دسته‌بند بیز در حالت‌های مختلف ویژگی‌های آماری

Table 5: Evaluation criteria for Bayesian outputs in different states of statistical Features

ویژگی‌های آماری	Accuracy	Recall	F1-Score	Precision
Tf_Idf, Most Frequent, Tf_Isf	۹۷.۲%	۰.۹۸	۰.۹۸	۰.۹۸
Most Frequent	۹۸.۱%	۰.۹۸۲	۰.۹۹	۰.۹۹
Tf_Idf, Tf_Isf	۹۹.۸%	۰.۹۱	۰.۹۴	۰.۹۹

باتوجه به نتایج به دست آمده، هنگامی که ویژگی بیشترین فرکانس را از ماتریس داده‌ها حذف می‌شود یادگیرنده بیز قادر به تشخیص داده‌ها با نشان ۱ نیست. بیشترین دقت نیز برای جدا کردن کلیدواژگان از متن نیز مربوط به وقتی است که بیشترین فرکانس به تنها یی مورد بررسی قرار می‌گیرد.

جدول ۶: میزان ارزیابی داده‌های تست برای خروجی‌هایی که دسته‌بند آن‌ها را کلیدواژه تشخیص داده است.

Table 6: Evaluation of test data for outputs whose keywords have been categorized

تعداد کلیدواژگان انتخابی نویسنده‌گان	تعداد کلمات	Precision	F1-Score	Recall	ترتیب ویژگی‌های آماری
۷۶	۷۴۳۳	۰.۰۱	۰.۰۲	۰.۹۲	Tf_Idf, Most Frequent, Tf_Isf
۷۶	۴۶۱۸	۰.۰۲	۰.۰۳	۰.۹۵	Most Frequent

اما آنچه بیشتر از همه در این مبحث حائز اهمیت است، نتایج مربوط به داده‌هایی از تست است که یادگیرنده به منزله کلیدواژه ارائه داده است. باتوجه به آزمایشات انجام شده هنگامی که بیشترین فرکانس از ویژگی‌ها حذف می‌شود، یادگیرنده قادر به تشخیص هیچ‌کدام از کلیدواژگان نیست. باتوجه به جدول ۶، در سطر دوم که بیشترین فرکانس فقط استفاده شده است نتایج بهتری از سطر اول که تمامی ویژگی‌ها به بیز انتقال پیدا کرده است، نشان می‌دهد. در سطر دوم هم تعداد کلیدواژگان واقعی بیشتری تشخیص داده شده و Recall بیشتری دارد.

و از طرفی دیگر تعداد کلماتی که یک واقعی نیستند و فقط صرفاً شبیه به کلمات کلیدی نویسنده هستند حدود ۳۰۰۰ عبارت کمتر از حالت سطر اول است.

حال برای اینکه تعداد کلیدواژگان انتخابی دارای محدودیت شوند و میزان دقت مدل طراحی شده را بالاتر ببرود، ابتدا افعال و کلمات کمتر از سه حرف حذف می‌شوند. سپس اقدام به جستجوی کلمات انتخابی در بخش‌هایی از متن شده است که دارای اهمیت روان‌شناسنامه بیشتری هستند. با توجه به اینکه بخش‌هایی چکیده و مقدمه در هر متن باید نشان‌دهنده محتوای آن متن باشند، این راهکار در نظر گرفته شده است. نتایج پژوهش لانگ (2010, p.21, 26) نشان می‌دهد که عموماً کلمات مهم و اثرگذار هر سند حداقل یکبار در بخش‌های چکیده و مقدمه آمده‌اند. از آنجایی که انتخاب این بخش‌ها به صورت اتوماتیک و هوشمند دشوار است، ۱۰ درصد ابتدای پایان‌نامه برای اعمال پس‌پردازش انتخاب شده است. نتیجه به دست آمده این است که اگر کلمات در بخش چکیده و فصل اول پایان‌نامه‌ها وجود نداشته باشند از مجموع کلمات انتخابی حذف می‌شوند. در مرحله بعد هر کدام از عبارت‌های تک‌کلمه‌ای و دو کلمه‌ای مربوط به داده‌های تست را به صورت جداگانه با توجه به بیشترین مقدار فرکانس آن‌ها لیست می‌شوند و فقط ۱۰ اتای اول آن‌ها در نظر گرفته می‌شوند. البته این تعداد برای قسمت عبارت‌های دوکلمه‌ای ممکن است کمتر نیز باشد. پس در این صورت خروجی انتها بی مدل طراحی شده حدود ۲۰ عبارت تک‌کلمه‌ای و دوکلمه‌ای است که برای هریک از ۲۱ متن به منزله داده تست در نظر گرفته شده است. همان‌طور که در جدول ۷ مشاهده می‌شود، پس از انجام عملیات گفته شده تعداد کلیدواژگان انتخابی بسیار محدودتر می‌شوند.

جدول ۷: میزان ارزیابی داده‌های تست بعد از پس‌پردازش برای خروجی‌هایی که دسته‌بند آن‌ها را کلیدواژه تشخیص داده است.

Table 7: Evaluation of post-processing test data for outputs that have been categorized by keyword.

تعداد کلیدواژگان انتخابی نویسنده‌گان	تعداد کلمات	تعداد کلیدواژگان	تعداد کلمات	Precision	F1-Score	Recall	ویژگی آماری	مرحله	عبارت‌های تک‌کلمه‌ای
۴۲	۲۱۰	۰.۲۰	۰.۳۲۳	۰.۸۴	Most Frequent				عبارت‌های تک‌کلمه‌ای
۳۴	۱۵۸	۰.۸۰	۰.۸۸۸	۰.۹۸	Most Frequent				عبارت‌های دوکلمه‌ای

همان طور که از جدول ۷ مشخص است مجموع کلمات انتخابی بعد از پردازش گفته شده برابر با ۳۶۸ کلمه شده است، یعنی تعداد ۴۲۵ کلمه نسبت به بهترین حالت قبل از پردازش که در جدول ۶ آمده، کم شده است. درصد Precision داده های تست که در جدول ۷ آمده است مربوط به ۱۰ کلمه اول مجموع کلمات انتخابی است که با توجه به مقدار بیشترین فرکانسیان لیست شده اند. برای مثال داده های خروجی مربوط به یک پایان نامه که مربوط به گرایش انسانی است در جدول ۸ و ۹ آورده شده است. همان طور که مشاهده می کنید، در انتها تعداد ۱۰ عبارت تک کلمه ای و تعداد ۵ عبارت دو کلمه ای به منزله خروجی مدل طراحی شده موجود است. نویسنده این پایان نامه عبارت های {ساختار سازمانی، نوآوری، بروکراتیک، آموزش عالی} را به منزله کلیدواژگان درنظر گرفته است. همان گونه که مشاهده می شود، تمامی کلیدواژگان انتخابی نویسنده در خروجی مدل ساخته شده، آمده است. علاوه بر آن ها، ۱۱ کلمه دیگر نیز به منزله کلیدواژه تشخیص داده شده اند که همگی آن ها به غیر از دو کلمه {توجه، همبستگی} نشان دهنده مفهوم متن موردنظر است.

جدول ۸: نمونه ای از خروجی برنامه برای کلمات تک عبارتی

Table 8: An example of a program output for unigram words

کلمات	نشان نویسنده	نشان تشخیص داده شده توسط مدل	بیشترین فرکانس
نوآور	۷۸	۱	۱
سازمان	۷۶	۱	.
توجه	۵۷	۱	.
کارکن	۴۸	۱	.
سیستم	۴۸	۱	.
دانشگاه	۴۳	۱	.
آموزش	۴۰	۱	.
بروکراتیک	۳۷	۱	۱
همبستگی	۲۵	۱	.
انسان	۲۴	۱	.

جدول ۹: نمونه‌ای از خروجی برنامه برای کلمات دو عبارت

Table 9: An example of a program output for bygram words

کلمات	بیشترین فرکانس	نمایش نویسنده	نمایش تشخیص داده شده توسط مدل
خلاف نوآور	۶۱	۱	.
ساختار سازمان	۴۳	۱	۱
قوانين مقررات	۳۹	۱	.
سلسه‌مراتب	۳۳	۱	.
آموزش عالی	۲۶	۱	۱

۴-۲. ابزار مورد استفاده

در این پژوهش از زبان برنامه‌نویسی رایگان پایتون^{۲۷} استفاده شده است. برای پیش‌پردازش متن‌های ورودی نیز همان‌طور که بیان شد از برنامه رایگان هضم^{۲۸} استفاده شده است.

۵. نتیجه

انتخاب کلیدواژگان توسط نویسندهان همواره تصمیمی وابسته به شخص بوده است و ممکن است افراد مختلف برای انتخاب مجموعه‌ای از لغات برای یک متن تصمیمی واحد نگیرند. در این پژوهش سعی بر آن بوده است که با رویکری نوین، با توجه به تعداد کم اسناد ورودی، مدل و برنامه‌ای ساخته شود که بدون وابستگی به گرایش پایان‌نامه‌ها و معنای کلمات آن‌ها و صرفاً با استفاده از ویژگی‌های آماری کلمات در هر متن، به استخراج کلیدواژگان پرداخته شود. همان‌طور که در جدول ۵ و ۷ نیز نشان داده شده است، مدل طراحی شده در بهترین شرایط که مربوط به استفاده از یک ویژگی بیشترین فرکانس است، قادر است با دقت کلی ۹۸.۱ درصد، از هر پایان‌نامه تعداد حدکثر ۲۰ کلیدواژه را استخراج کند. کلیدواژگان نوشته شده در هر پایان‌نامه نیز با دقت‌های ۸۴ درصد و ۹۸ درصد که به ترتیب مربوط به عبارت‌های تک‌کلمه‌ای و دوکلمه‌ای هستند، دربر می‌گیرد. برای مثال نیز یکی از نتایج بدست آمده بر روی یکی از پایان‌نامه‌ها در جداول ۸ و ۹ نیز آورده و نشان داده شده است که مدل طراحی شده از این پایان‌نامه تعداد ۵ عبارت دوکلمه‌ای و ۱۰ عبارت تک‌کلمه‌ای را به منزله کلیدواژه تشخیص داده است که ۱۰۰ درصد کلیدواژگان نوشته شده در پایان‌نامه در آن آمده است.

۶. کارهای آینده

از این مدل طراحی شده سازمان هایی نظیر ایران داک، پایگاه مجلات تخصصی نور و خبررسانی های فارسی نظری همشهری برای نمایه سازی، استخراج کلیدواژگان و فیلترینگ می توانند بهره ببرند. یکی از مشکلات بزرگی که در این فرایند با آن مواجه هستیم، ریشه یابی و یکپارچه کردن متن است که این امر باعث می شود تا پیش پردازش مناسبی صورت نگیرد و کلمات زائد به شکل مناسبی حذف نشوند. این عوامل پیچیدگی ابعاد ماتریس ورودی برای دسته بند را بالا می برد و باعث پایین آمدن دقت می شود. در آینده سعی بر آن می شود به جای استفاده از برنامه هضم، از برنامه فارسی یار^{۲۹} استفاده شود و نتایج به دست آمده موردنرسی قرار داده شوند. همچنین، می توان به جای استفاده از دسته بند بیز از دیگر دسته بند ها استفاده کرد و نتایج به دست آمده را با این دسته بند مقایسه کرد.

۷. پی نوشت ها

1. Atanu Dey
2. unigrams
3. senti-n-gram
4. Rafeal
5. Neumann
6. Loung
7. Litvak
8. Huan
9. Turney
10. Witten
11. Hachoen Kerner
12. Onan
13. bagging
14. random forest
15. Moien Maadi
16. Mohammad Razzaghnoori
17. word2vec
18. Behnam Sabeti
19. miras text
20. Morteza Okhovvat
21. Hidden Markov
22. Moammad

23. Term Frequency Inverse Document Frequency (TF-IDF)
24. most frequent of words
25. term frequency inverse sentence frequency (TF-ISF)
26. naïve bayes
27. python
28. sobhe.ir
29. text-mining.ir

۸. منابع

- گندمکار، ر. (۱۳۹۶). توسعی معنایی در زبان فارسی؛ مطالعه موردی افعال. *جستارهای زبانی*، ۵(۵۲)، ۲۲۳-۲۵۶.

References

- Gandomkr, R. (2017). Semantic expansion in Persian; A case study of verbs. *Language Related Research*. 5(53), 233-256.[In Persian]
- Aytug Onan, Serdar Korukoglu, Hasan Bulut. (2016). *Ensemble of keyword extraction methods and classifiers in text classification*. *Expert Systems with Applications*. 232-247.
- Atanu Dey, Mamata Jenamani, Jitesh J.Thakkar. (2018).Senti-N-Gram: An n-gram lexicon for sentiment analysis. *Expert Systems with Applications*. 92–105.
- Rafea, S. El-Beltagy & A. (2010). Kp-miner: Participation in semeval-2” Proceedings of the 5th International Workshop. *ACL2010*. 190–193
- Neumann, K. E. and G. (2010).Dfki keywe: Ranking keyphrases extracted from scientific articles. In *Proceedings of the 5th International Workshop on Semantic Evaluation*. ACL 2010, 150–153.
- Luong, T. Nguyen and M. (2010).WINGNUS: Keyphrase extraction utilizing document logical structure. *Proceedings of the 5th International Workshop*. 21–26.
- Turney, P. D. (2002).Mining the Web for Lexical Knowledge to Improve Keyphrase Extraction: Learning from Labeled and Unlabeled Data P.D. *ArXiv*.
- Witten, I. H. , Paynter, G. W. , Frank, E. , Gutwin, C. (1999). KEA: Practical

automatic keyphrase extraction. In *Proceedings of the Fourth ACM Conference on Digital Libraries*.

- HaCohen-Kerner. (2005).Automatic extraction and learning of keyphrases from scientific articles. *Lecture Notes in Computer Science*. 657-669.
- Moien Maadi, Kazim Fouladi. (2015).Providing a method for extracting keywords in the Persian language. *International Academic Journal of Innovative Research*. 34-42.
- Mohammad Razzaghnoori, Hedieh Sajedi , Iman Khani Jazani. (2017).Question classification in Persian using word vectors and frequencies. *Journal Cognitive Systems Researc*. 16-27.
- Behnam Sabeti, Hossein Abedi Firouzjaei, Ali Janalizadeh Choobbasti. (2018).MirasText: An Automatically Generated Text Corpus for Persian. *Proceedings of the Eleventh International Conference on Language Resources and Evaluation (LREC 2018)*.
- Morteza Okhovvat, Behrouz Minaei Bidgoli. (2010). A hidden Markov model for Persian part-of-speech Tagging. 94–101.
- Sayyid Mohammad Hoseini Khozani, Hosein Bayat. (2011).Specialization of keyword extraction approach to Persian texts. *International Conference of Soft Computing and Pattern Recognition (SoCPaR)*. 112-116.
- Beliga, S. , Mestrovic, A. , & Martincic-Ipsic, S. (2015).An overview of graph-based keyword extraction methods and approaches. *Journal of Information and Organizational Sciences*. 39 (1), 1–20.
- Lott, B. 2012. *Survey of keyword extraction techniques*. *UNM Education*.
- Rossi, R. G. , Maracini, R. M. , & Rezende, S. O. (2014). Analysis of domain independent statistical keyword extraction methods for incremental clustering. *Learning and Nonlinear Models*. 17-37.
- Neto, L. J. , Santos, A . D. , Kaestner, C. A. (2000). Document cluster- ing and text summarization. In *Proceedings of the 4th International Conference on Practical*

Applications of Knowledge Discovery and Data Mining. 41-55.

- Fiori, A. (2014). Innovative document summarization techniques: Revolutionizing knowledge understanding. *Advances in Data Mining and Database Management*.
- Litvak, M. , & Last, M. (2008).Graph-based keyword extraction for single-document summarization. In *Proceedings of the Workshop on Multi-source Multilingual Information Extraction and Summarization*. 17-24.
- Huan, C. , Tian, Y. (2006).Keyphrase extraction us- ing semantic network structure analysis. In *Proceedings of the Sixth International Conference on Data Mining*. 275-284.