

# Automatic Detection of the Boundary between Metadata and Body in Persian Theses using BA\_SVM

## Mohadese Rahnama

M.Sc. Student in Computer Engineering; Alzahra University;  
Tehran, Iran Email: m.rahnama@student.alzahra.ac.ir

## Seyed Mohammad Hossein Hasheminejad\*

PhD in Computer Engineering; Assistant Professor; Alzahra  
University; Tehran, Iran Email: smh.hasheminejad@alzahra.ac.ir

## Jalal A Nasiri

PhD in Computer Engineering; Assistant Professor;  
Iranian Research Institute for Information Science and Technology  
(IranDoc); Tehran, Iran Email: j.nasiri@irandoc.ac.ir

Received: 09, May 2020 Accepted: 02, Nov. 2020

**Abstract:** Metadata extraction facilitates the process of indexing and improves information retrieval. Also automation of this process increases efficiency more than manual extraction. The example of the thesis metadata are names of students, professors, title, field, degree, abstract, keywords, etc. In this paper the aim is automatic boundary detection of metadata from the main body in Persian theses. Therefore, 250 theses collected from IRANDOC system. Features were extracted from paragraphs of each thesis then paragraphs were classified using support vector machine into 2 classes: metadata and body. In this study, Bat algorithm is used to set the parameter of SVM. The result reveals that the proposed method predicts type of paragraphs with 96.6 percent accuracy.

**Keywords:** Metadata Extraction, Information Extraction, Support Vector Machine (SVM), Metaheuristic Algorithm, Bat Algorithm (BA)

Iranian Journal of  
Information  
Processing and  
Management

Iranian Research Institute  
for Information Science and Technology  
(IranDoc)

ISSN 2251-8223

eISSN 2251-8231

Indexed by SCOPUS, ISC, & LISTA

Vol. 36 | No. 4 | pp. 1159-1180

Summer 2021



\* Corresponding Author

# استخراج هوشمند مرز فراداده و متن

## در پایان‌نامه‌های فارسی با رویکرد BA\_SVM

محدثه رهنما

دانشجوی کارشناسی ارشد مهندسی کامپیوتر؛  
دانشگاه الزهرا (س)؛ تهران، ایران؛  
m.rahnama@student.alzahra.ac.ir

سید محمدحسین هاشمی‌نژاد

دکتری مهندسی کامپیوتر؛ استادیار؛  
دانشگاه الزهرا (س)؛ تهران، ایران؛  
پدیدآور رابط smh.hasheminejad@alzahra.ac.ir

جلال‌الدین نصیری

دکتری مهندسی کامپیوتر؛ استادیار؛ پژوهشگاه علوم  
و فناوری اطلاعات ایران (ایرانداک)؛ تهران، ایران؛  
j.nasiri@irandoc.ac.ir



مقاله برای اصلاح به مدت ۳۵ روز نزد پدیدآوران بوده است.

پدیش: ۱۳۹۹/۰۸/۱۲

دریافت: ۱۳۹۹/۰۲/۲۰

نشریه علمی | رتبه بین‌المللی  
پژوهشگاه علوم و فناوری اطلاعات ایران  
(ایرانداک)

شاپا (چاپی) ۸۲۲۳-۲۲۵۱

شاپا (الکترونیکی) ۸۲۳۱-۲۲۵۱

نمایه در SCOPUS، LISTA، ISC، و

jipm.irandoc.ac.ir

دوره ۳۶ | شماره ۴ | صص ۱۱۵۹-۱۱۸۰

تایستان ۱۴۰۰



**چکیده:** استخراج فراداده باعث تسهیل در فرایند نمایه‌سازی و بهبود در بازبایی اطلاعات است. از سوی دیگر، خودکارسازی این فرایند سبب افزایش کارایی نسبت به استخراج دستی فراداده‌هاست. نام دانشجو، نام اساتید، عنوان، رشته و مقطع تحصیلی، چکیده، و کلمات کلیدی نمونه‌ای از فراداده‌های پایان‌نامه است. هدف در این مقاله شناسایی خودکار مرز فراداده و بدنه اصلی در پایان‌نامه‌های فارسی است. بدین منظور، ۲۵۰ پایان‌نامه ثبت‌شده در سامانه «ایرانداک» جمع‌آوری شده است. ویژگی‌های مد نظر از هر پاراگراف استخراج شده و سپس، پاراگراف‌های پایان‌نامه با روش ماشین بردار پشتیبان به دو کلاس فراداده و بدنه طبقه‌بندی شد. در این پژوهش برای تنظیم پارامترهای الگوریتم ماشین بردار پشتیبان، الگوریتم فرامکاشفه‌ای خفشا به کار گرفته شده است. نتایج نشان می‌دهد که روش پیشنهادی با دقت ۹۶/۶ درصد نوع پاراگراف را تشخیص می‌دهد.

**کلیدواژه‌ها:** استخراج فراداده، استخراج اطلاعات، ماشین بردار پشتیبان، الگوریتم فرامکاشفه‌ای، الگوریتم خفشا

## ۱. مقدمه

با گسترش نوشتارهای علمی، دسترسی سریع و ساده به فراداده‌های مقالات و پایان‌نامه‌ها بیش از پیش ضرورت می‌یابد. فرایند استخراج اطلاعات از متون علمی به صورت دستی و توسط انسان، بسیار زمان‌بر و پرهزینه خواهد بود. از این رو، خودکارسازی این فرایند می‌تواند در تسریع شاخص‌گذاری نوشتارهای علمی مؤثر باشد. در زمینه استخراج فراداده از مقالات علمی و پایان‌نامه‌ها پژوهش‌هایی به‌ویژه در زبان انگلیسی صورت گرفته، اما در زبان فارسی تحقیقات کمتری انجام شده است.

فراداده‌ها از اطلاعات بسیار باارزش پایان‌نامه‌ها و رساله‌ها هستند. اولین و مهم‌ترین گام در استخراج خودکار فراداده‌ها تشخیص مرز شروع بدنه پایان‌نامه است. از این رو، تشخیص خودکار مرز بین فراداده و متن اصلی پایان‌نامه در سامانه‌های ذخیره‌سازی موضوعی و بازبایی اطلاعات نقش مهمی ایفا می‌کند. در بسیاری از سامانه‌های ثبت و اشاعه متون علمی ورود اطلاعات فراداده توسط عامل انسانی انجام می‌پذیرد.

هدف این مقاله تشخیص خودکار مرز فراداده‌ها و متن اصلی پایان‌نامه‌های فارسی است. بدین ترتیب، فرایند استخراج فراداده‌ها ساده‌تر صورت گرفته و نیازی به پردازش تمامی متن پایان‌نامه نیست. به عبارت دیگر، استخراج هوشمند پاراگرافی که ابتدای متن اصلی پایان‌نامه یا رساله باشد، هدف این پژوهش است. در شکل ۱، نمونه‌هایی از مرز فراداده و بدنه در پایان‌نامه‌های فارسی نشان داده شده است.

پرسش‌های اصلی این پژوهش به شرح زیر است که در انتهای پژوهش پاسخ داده می‌شود:

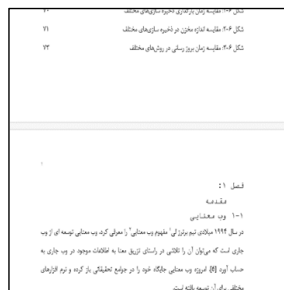
- ◇ آیا ترکیبی از ویژگی‌های مکاشفه، لغوی، هندسی و ویژگی‌های مرتبط با توالی پاراگراف‌ها می‌تواند دقت شناسایی هوشمند فراداده را افزایش دهد؟
- ◇ آیا استفاده از الگوریتم فراابتکاری خفاش<sup>۲</sup> می‌تواند در دستیابی به پارامترهای بهینه الگوریتم ماشین بردار پشتیبان مؤثر باشد؟



(پ)



(ب)



(الف)

شکل ۱. نمونه‌هایی از مرز فراداده و بدنه در پایان‌نامه. در شکل‌های (الف)، (ب)، و (پ) به ترتیب عبارات‌های «فصل ۱»، «فصل اول»، و «کلیات تحقیق» شروع متن اصلی و بدنه پایان‌نامه هستند.

## ۲. پیشینه پژوهش

استخراج اطلاعات<sup>۱</sup>، فرایند تعیین و شناسایی اصطلاحات از پیش تعریف‌شده در حوزه‌ای خاص، بدون در نظر گرفتن اطلاعات غیر مرتبط است (Piskorski & Yangarber 2013). در این فرایند، موجودیت‌ها، خصیصه‌ها و ارتباط آن‌ها با یکدیگر به دست می‌آید. هدف از استخراج اطلاعات، دستیابی به اطلاعات ساخت یافته<sup>۲</sup> از اسناد بدون ساختار<sup>۳</sup> یا نیمه‌ساخت یافته<sup>۴</sup> است (Adrian et al. 2017). فراداده، داده‌ای است که اطلاعاتی درباره داده دیگر می‌دهد و در انواع مختلفی دسته‌بندی می‌شوند؛ از جمله: توصیفی (مانند عنوان و نویسنده)، ساختاری (مثل تعداد صفحات) و اجرایی (مانند مالک) (Meng et al. 2018). طبق مقاله مروری «ناسار، جعفری و مالیک»، به‌طور کلی، دو نوع اطلاعات از مقالات علمی قابل استخراج است. فراداده‌ها یکی از این اطلاعات است. برای نمونه، می‌توان به‌عنوان مقاله، نویسندگان، محل انتشار، تاریخ انتشار، کلمات کلیدی و ... اشاره کرد. تحلیل صفحات آغازین مقالات و تجزیه رشته‌های مرجع<sup>۵</sup> دو رویکرد متداول برای دستیابی به فراداده‌ها هستند. اطلاعات بالقوه دیگری مانند موضوع پژوهش، مسئله، الگوریتم‌ها، نحوه پاسخ به سؤالات پژوهش، مجموعه داده‌های به کار رفته، ابزارها، محدودیت‌های حوزه مطالعه، روش ارزیابی، نتایج و پیشنهاد برای پژوهش‌های آینده نیز در مقالات علمی و دانشگاهی ظاهر می‌شوند که دستیابی به آن‌ها می‌تواند مفید باشد (Nasar, Jaffry, & Malik 2018).

1. information extraction (IE)  
4. semi-structured documents

2. structured information  
5. references strings

3. unstructured documents

بیشتر پژوهش‌های انجام‌شده در این زمینه بر استخراج فراداده‌ها و سایر اطلاعات ذکر شده از مقالات -به‌ویژه مقالات لاتین- تمرکز کرده‌اند. برای نمونه، هدف مقاله «منگ و همکاران»، به‌دست آوردن اطلاعاتی درباره موضوع، روش‌های استفاده‌شده و داده‌های به‌دست‌آمده در مقاله ورودی است (Meng et al. 2018). در پژوهشی دیگر، «صفدر و همکاران» به شناسایی شبه‌کدها و تحلیل منابع موجود درباره شبه‌کد و الگوریتم پیشنهادی در اسناد مقاله توجه کرده‌اند (Safder et al. 2020).

نمونه‌هایی از پژوهش‌های موجود در زمینه استخراج فراداده‌ها از متون علمی عبارت‌اند از: «پنگ و مک کالوم» در مطالعه خود سعی در تحلیل بخش‌های اول مقالات دارند و همچنین، رشته مراجع مقاله را به‌منظور استخراج فراداده‌های مورد نظر پردازش می‌کنند (Peng & McCallum 2006). استخراج فراداده‌های توصیفی از مقالات نشریه‌های علمی و تجزیه فهرست منابع آن‌ها از اهداف مقاله Tkaczyk et al. (2015) است. «بوخرز، آمهور و ستاب» بر استخراج بخش منابع و دستیابی به اطلاعات آن متمرکز شده‌اند (Boukhers, Ambhore & Staab 2019). همچنین، هدف پژوهش Rizvi, Dengel & Ahmed (2019) تشخیص خودکار بخش فهرست منابع در مقالات پژوهشی است.

هدف پژوهش حاضر، استخراج فراداده‌های توصیفی صفحات آغازین پایان‌نامه است. برخی مطالعات انجام‌شده در این زمینه عبارت‌اند از:

«تن‌سازان و مهدوی» در پژوهش خود سعی دارند فراداده‌ها را از سرآیند مقاله و بخش فهرست منابع فارسی و انگلیسی اسناد و مقالات علمی فارسی استخراج کنند. در این روش، اولین گام شناسایی بخش‌های آغازین و فهرست منابع متن مقاله ورودی است. فراداده‌های هر مقاله از ابتدای متن آن شروع شده و به ابتدای مقدمه ختم می‌شود. بنابراین، برای شناسایی مرز شروع بدنه اصلی مقاله، در ساده‌ترین حالت، یافتن کلمه کلیدی «مقدمه» نشان‌دهنده شروع بدنه اصلی است (۱۳۹۶). البته، به‌دلیل تنوع قالب‌های ساختاری مقالات، این روش به‌تنهایی کفایت نمی‌کند. به‌عنوان نمونه ممکن است پیش از شروع بدنه اصلی مقاله، چکیده انگلیسی قرار داشته باشد. در پژوهش حاضر این بخش جزو سرآیند مقاله در نظر گرفته نشده و راه‌حلی برای استثناها ارائه شده است. شناسایی فهرست منابع، تنها بر اساس ویژگی‌های متنی انجام شده است و استخراج

فراداده از سرآیند مقاله و فهرست منابع فارسی و انگلیسی با مدل آماری «میدان تصادفی شرطی»<sup>۱</sup> انجام گرفته است.

در مطالعه «پنگ و مک کالموم» مقاله‌ای با داده‌های متنی بررسی شده و فراداده‌های سرآیند و اطلاعات بخش منابع مورد توجه قرار گرفته است. به منظور استخراج اطلاعات از سرآیند مقالات، داده‌های استخراج شده از سرآیند شامل کلمات ابتدای پایان‌نامه تا اولین بخش مقاله (به‌طور معمول مقدمه) و یا کل کلمات صفحه اول بوده‌اند. برای استخراج فراداده از بخش منابع نیز از رشته متنی منابع استفاده شده است. در این روش، افزون بر ویژگی‌های لغوی، ویژگی‌های مربوط به طرح‌بندی و هجی کلمات نیز در نظر گرفته شده و روش میدان تصادفی شرطی برای استخراج فراداده‌ها استفاده شده است (Peng & McCallum 2006).

هدف مقاله «کرن» و همکاران استخراج فراداده از مقالات در قالب «پی‌دی‌اف»<sup>۲</sup> است. پژوهشگران، مرحله پیش‌پردازش را در دو گام انجام داده‌اند: (۱) بلوک‌های متنی به‌دست آمده از هر سند را در گروه فراداده‌های مورد نظر طبقه‌بندی می‌کنند (بلوک‌های متنی که فراداده نیستند، در گروه «سایر» دسته‌بندی می‌شوند)؛ و (۲) هر بلوک متنی از چندین «توکن»<sup>۳</sup> تشکیل شده است و در گام دوم طبقه‌بندی در سطح «توکن»‌ها انجام می‌شود. بدین ترتیب که هر یک از کلمات در چندین گروه طبقه‌بندی می‌شود. به‌عنوان نمونه، بلوک متنی که در گام اول برچسب «نویسنده» گرفته، در گام بعدی به هر کلمه بلوک یکی از برچسب‌های نام، نام خانوادگی، ایمیل، وابستگی و ... تخصیص می‌یابد. در انتهای فرایند، فراداده‌های مورد نظر استخراج می‌شود (Kern et al. 2012).

«کان، لانگ و انگوین» و «کونگ» و همکاران در مطالعات خود با تحلیل ساختار منطقی سند<sup>۴</sup>، فراداده‌های مورد نظر را استخراج کرده‌اند (Kan, Luong & Nguyen 2010; Cuong et al. 2015). روش ارائه‌شده در «کان، لانگ و نوین» به هر بخش از مقاله بر اساس نقش آن‌ها در سند، برچسبی از فراداده‌ها تخصیص می‌دهند (Kan, Luong & Nguyen 2010). «دو» و همکاران در مطالعه خود به نحوه استخراج نام نویسندگان، راه‌های دسترسی به آن‌ها (مانند ایمیل، شماره تلفن، محل خدمت و ...) و تطبیق آن‌ها با نام هر نویسنده

1. conditional random field (CRF)

2. PDF

3. token

4. logical structure analysis

پرداخته‌اند (Do et al. 2013). در پژوهش حاضر، پیش از آغاز فاز شناسایی فراداده‌ها، نسخه «پی‌دی‌اف»، با هدف حفظ ویژگی‌های مکانی و هندسی محتوا به «یکس‌ام‌ال»<sup>۱</sup> تبدیل شده و به هر خط بر اساس نقش آن در متن مقاله، برجسی به روش Kan, Luong (2010) & Nguyen اختصاص داده‌اند. سپس، خطوطی که یکی از دو برجسب «نویسنده» و «وابستگی» را داشته باشند، وارد سیستم اصلی می‌شوند. سیستم شامل دو مرحله است: در مرحله نخست، نام نویسنده (یا نویسندگان) و اطلاعات وابستگی‌ها شناسایی شده، و در مرحله بعد، وابستگی‌های هر نویسنده با نام او تطبیق داده می‌شود. بخش شناسایی با روش «یادگیری ماشین با نظارت» و روش «میدان تصادفی شرطی» انجام شده و بخش تطبیق و یافتن ارتباطات نیز با «روش ماشین بردار پشتیبان» صورت می‌گیرد.

«سوزا، موریرا و هیوسر» در مطالعه خود، با توجه به اینکه فراداده‌های یک مقاله در صفحه اول آن ظاهر می‌شود، فقط به تحلیل صفحه اول اکتفا کرده‌اند. پردازش در دو مرحله تعریف شده است: نخست، بخش‌های مختلف صفحه اول، یعنی سرعنوان، عنوان، اطلاعات نویسنده، بدنه و پاورقی مشخص شده، در مرحله بعد، فراداده‌های موجود در هر بخش استخراج می‌شوند (Souza, Moreira & Heuser 2014).

Tkaczyk et al. (2015) برای استخراج فراداده از صفحات ابتدایی و فهرست منابع مقالات علمی، سیستمی با ورودی قالب «پی‌دی‌اف» ارائه کردند<sup>۲</sup>. از این رو، نخست، محتوای متنی و ویژگی‌های تصویری آن‌ها استخراج می‌شود. این سیستم خود به سه بخش استخراج ساختار پایه<sup>۳</sup>، استخراج فراداده، و استخراج فهرست مقالات ارجاعی<sup>۴</sup> تقسیم می‌شود. نخست، ساختار سلسله‌مراتبی سند به دست می‌آید. سپس، نواحی مقاله به چهار بخش فراداده، بدنه، مراجع و سایر طبقه‌بندی می‌شود؛ در ادامه، ناحیه‌هایی که برجسب فراداده گرفته‌اند، به انواع فراداده‌های از قبیل تعریف شده طبقه‌بندی می‌شود. هر دو طبقه‌بندی با روش «ماشین بردار پشتیبان» انجام شده است (Cortes & Vapnik 1995). در بخش استخراج فهرست منابع، رشته متنی مراجع تجزیه شده و فراداده‌های منابع مقالات مشخص می‌شود (Tkaczyk et al. 2015).

تطبیق گراف ساختار<sup>۵</sup> سند نیز می‌تواند در یافتن موجودیت‌های مورد نظر مؤثر

1. XML

2. <http://cermine.ceon.pl/index.html>

3. basic structure extraction

4. bibliography extraction

5. structure graph matching

باشد. برای نمونه، «کولی و بلید» در مطالعه خود روشی ارائه کرده‌اند که افزون بر موجودیت‌های مقالات علمی، اسناد سازمانی و فاکتور خرید هم بررسی می‌شود (Kooli & Belaid 2016). در این روش، برچسب‌زنی موجودیت‌ها پس از استخراج متن از تصویر ورودی با استفاده از عبارات منظم و لغت‌نامه انجام می‌شود. در واقع، این برچسب‌زنی تنها بر اساس متن انجام شده و ویژگی‌های تصویری و ساختاری در مراحل بعدی مورد استفاده قرار می‌گیرد. در مرحله بعد، ساختارهای محلی در هر بخش از تصویر سند با یک گراف مدل می‌شود و با مدل‌هایی که سیستم از قبل یادگرفته (یادگیری به روش بدون نظارت یا خوشه‌بندی) تطبیق داده شده، نزدیک‌ترین مدل یافت می‌شود. از این مدل برای اصلاح خطاهای احتمالی و به‌روزرسانی استفاده می‌شود. در آخرین مرحله نیز موجودیت‌ها شناسایی می‌شوند.

«لیو» و همکاران در مقاله خود با به‌کارگیری شبکه‌های عمیق<sup>۱</sup> و ترکیب ویژگی‌های متنی و تصویری، فراداده‌ها را از سرآیند فایل «پی‌دی‌اف» مقالات استخراج کرده‌اند. در این روش با استفاده از شبکه عصبی پیچشی<sup>۲</sup>، ویژگی‌های متنی و تصویری به‌دست آمده و خروجی به شبکه حافظه طولانی کوتاه‌مدت<sup>۳</sup> فرستاده می‌شود (Liu et al. 2017).

هدف مقاله «فریز» و همکاران، ارائه ابزاری برای تبدیل فایل «پی‌دی‌اف» به «ایکس‌ام‌ال»، دستیابی به محتوای متنی و بخش‌های مختلف مقاله، مانند فراداده‌ها، بخش‌ها، زیربخش‌ها و مراجع است (Ferréset al. 2018). «کیو و ژو» مطالعه دیگری در زمینه استخراج فراداده‌ها از کتاب‌های دیجیتال انجام داده‌اند که در آن دسترسی به اطلاعاتی مانند عنوان کتاب، نویسنده، ناشر، تاریخ انتشار، شماره استاندارد بین‌المللی کتاب<sup>۴</sup> (شابک) و ... بررسی شده است (Qiu & Zhou 2019). اگرچه مقاله ذکر شده بر اسناد علمی تمرکز ندارد، اما از این جهت که به‌دنبال فراداده‌های توصیفی کتاب‌های دیجیتال است، با مسئله‌ای مشابه روبه‌روست.

در پژوهش حاضر، برای تنظیم پارامترهای ماشین بردار پشتیبان، از الگوریتم خفاش به‌عنوان یک الگوریتم فرامکاشفه‌ای<sup>۵</sup>، بهره گرفته شده است؛

1. deep neural networks      2. convolutional neural network (CNN)      3. long short-term memory (LSTM)  
4. ISBN      5. metahuristic



## ۲-۱. الگوریتم خفاش

مفهوم الگوریتم خفاش را نخستین بار (2010) Yang مطرح کرد. این الگوریتم با الهام از نحوه جهت‌یابی خفاش با پژواک صدا طراحی شده و برای بهینه‌سازی مسایل پیوسته مناسب است. برای سادگی موضوع، قواعد زیر در نظر گرفته شده است:

◇ همه خفاش‌ها از طریق ارسال صدا فاصله را درک کرده و تفاوت بین طعمه و موانع را تشخیص می‌دهند؛

◇ خفاش‌ها به صورت تصادفی و با سرعت  $v_i$  در مکان  $x_i$  پرواز می‌کنند و با فرکانس ثابت  $f_{min}$ ، طول موج  $\lambda$  و بلندی صدای  $A_0$  متغیر به دنبال طعمه می‌گردند. خفاش‌ها طول موج (یا فرکانس) پالس صدا و نرخ پالس انتشار  $r \in [0,1]$  را بر اساس میزان نزدیکی به هدف خود به صورت خودکار تنظیم می‌کنند؛

◇ بلندی صدا نیز می‌تواند تغییر کند. در این الگوریتم فرض می‌شود که بلندی صدا از مقدار  $A_0$ ، که عددی مثبت است، تا مقدار ثابت و کمینه  $A_{min}$  تغییر می‌کند؛

◇ همچنین، برای سادگی فرض می‌شود که از ردیابی صدا برای محاسبه میزان تأخیر و نقشه‌برداری سه‌بعدی استفاده نمی‌شود.

به‌طور کلی، محدوده فرکانس  $f$  در بازه  $[f_{min}, f_{max}]$  است و با طول موج در بازه  $[\lambda_{min}, \lambda_{max}]$  متناسب است.

از آن‌جا که حاصل ضرب  $\lambda f$  مقداری ثابت و برابر سرعت است، در پیاده‌سازی الگوریتم فقط مقدار فرکانس تغییر می‌یابد.

در شبیه‌سازی، خفاش‌های مصنوعی هر کدام موقعیت و سرعتی دارند که در یک فضای  $d$  بعدی حرکت می‌کنند و در گام  $t$  موقعیت  $x_i^t$  و سرعت  $v_i^t$  آن‌ها به‌روزرسانی می‌شود (رابطه ۱، رابطه ۲، و رابطه ۳).

$$f_i = f_{min} + (f_{max} - f_{min})\beta \quad \text{رابطه ۱}$$

$$v_i^t = v_i^{t-1} + (x_i^{t-1} - x_i^*)f_i \quad \text{رابطه ۲}$$

$$x_i^t = x_i^{t-1} + v_i^t \quad \text{رابطه ۳}$$

پارامتر  $\beta \in [0,1]$  عددی تصادفی از توزیع یکنواخت است.  $x_i^*$  بهترین مکان سراسری است که از مقایسه مکان همه خفاش‌ها (در هر تکرار) به‌دست می‌آید.

هنگام جست‌وجوی محلی، پس از یافتن بهترین راه‌حل از راه‌حل‌های فعلی، راه‌حل

جدیدی به روش مسیر تصادفی برای هر خفاش تولید می‌شود (رابطه ۴).

$$x_{new} = x_{old} + \varepsilon A^t \quad \text{رابطه ۴}$$

پارامتر  $\varepsilon \in [-1, 1]$  عددی تصادفی است.  $A^t = \langle A_i^t \rangle$  میانگین بلندی صدای همه خفاش‌ها در تکرار  $t$  است.

از طرفی، هنگامی که خفاش طعمه را پیدا می‌کند، بلندی صدای خود را کاهش و نرخ پالس را افزایش می‌دهد (رابطه ۵).

$$A_i^{t+1} = \alpha A_i^t, \quad r_i^{t+1} = r_i^0 (1 - \exp(-\gamma t)) \quad \text{رابطه ۵}$$

$\alpha$  و  $\gamma$  دو عدد ثابت هستند. در واقع،  $\alpha$  فاکتور خنک‌سازی است و برای هر  $0 < \alpha < 1$  و  $\gamma > 0$  اگر  $t \rightarrow \infty$  داریم:

$$A_i^t \rightarrow 0, \quad r_i^t \rightarrow r_i^0$$

برای سادگی پیاده‌سازی این الگوریتم، حالت  $\alpha = \gamma = 0.9$  در نظر گرفته شده است.

## ۲-۲. ترکیب ماشین بردار پشتیبان و الگوریتم خفاش

«ثروت، حسینیان و النقی» در مقاله خود برای تنظیم پارامترهای ماشین بردار پشتیبان از الگوریتم خفاش بهره گرفته‌اند (Tharwat, Hassanien & Elnaghi 2017). این بهینه‌سازی به صورت دو آزمایش روی ۹ مجموعه داده دانشگاه کالیفرنیا<sup>۱</sup> انجام شده است. در آزمایش اول، هدف، بهینه‌سازی پارامتر خطا  $C$  و پارامتر هسته چند جمله‌ای<sup>۲</sup>  $d$  بوده و در آزمایش دوم، پارامتر خطا  $C$  و پارامتر هسته تابع پایه شعاعی<sup>۳</sup>  $\sigma$  بهینه‌سازی شده است. روش آن‌ها به این صورت است که ابتدا داده‌ها نرمال شده (رابطه ۶) و به روش اعتبارسنجی متقابل  $k$ -fold و  $k=10$  به دو دسته داده‌های آموزشی و داده‌های آزمایشی تقسیم می‌شوند. سپس، الگوریتم خفاش اجرا شده و برای وضعیت هر خفاش، پارامترهای به دست آمده برای آموزش داده‌ها با ماشین بردار پشتیبان اعمال می‌شود. پس از محاسبه نرخ خطا (رابطه ۷) بهترین وضعیت ( $X^*$ ) خفاش‌ها به دست می‌آید. در این جا، تابع برازش<sup>۴</sup> همان نرخ خطاست،  $N_e$  تعداد نمونه‌هایی است که برچسب اشتباه خورده‌اند، و  $N$  تعداد کل داده‌های آزمایش است. این الگوریتم تا برآورده شدن معیار خاتمه ادامه می‌یابد (شکل ۲).

1. <https://archive.ics.uci.edu/ml/datasets.php>

2. polynomial kernel

3. radial basis function (RBF)

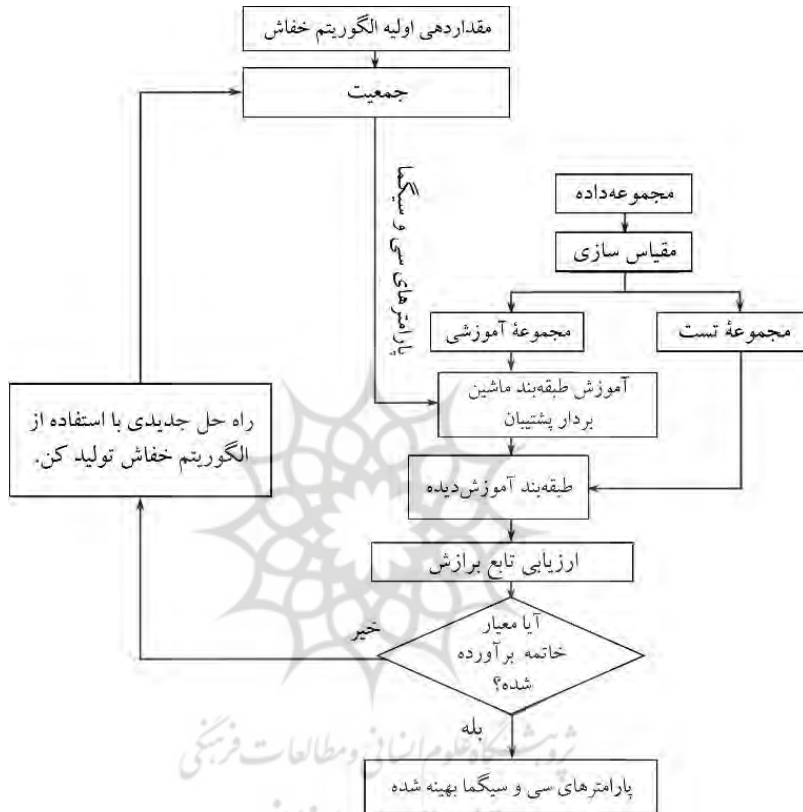
4. fitness function

$$v' = \frac{v - \min}{\max - \min}$$

رابطه ۶

$$\text{Minimize: } F = \frac{N_e}{N}$$

رابطه ۷



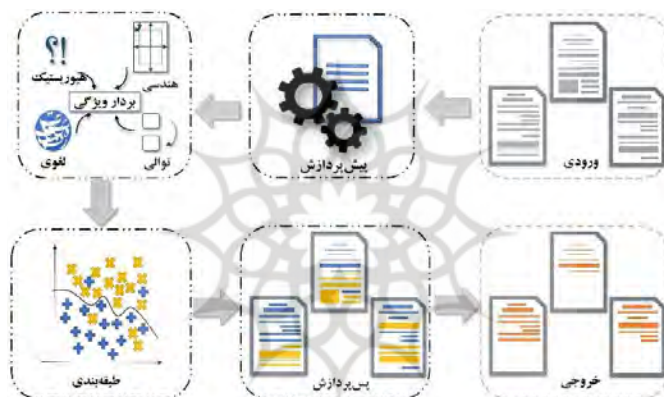
شکل ۲. نمودار گردش کار بهینه‌سازی ماشین بردار پشتیبان با الگوریتم خفاش (Tharwat, Hassanien & Elnaghi 2017)

### ۳. روش پژوهش

در این بخش درباره روش پیشنهادی برای تشخیص مرز بدنه اصلی در اسناد پایان‌نامه‌های فارسی توضیح داده می‌شود.

هر پایان‌نامه قبل از شروع فصل‌های اصلی از بخش‌هایی چون عنوان فارسی (مطابق با جلد فارسی)، اصالت و مالکیت، صورت جلسه دفاع، تقدیم‌نامه، سپاسگزاری، چکیده فارسی و فهرست‌ها تشکیل می‌شود. هدف پژوهش حاضر، پیشنهاد روشی برای شناسایی

و جداسازی پاراگراف‌های مربوط به این بخش‌ها از بدنه اصلی و شروع فصل‌های اصلی پایان‌نامه است. در این روش، متن پایان‌نامه ورودی در پیش‌پردازش یکسان‌سازی می‌شود. هدف این کار، یکسان‌سازی یونی‌کد کاراکترهایی است که ممکن است چند یونی‌کد داشته باشند؛ مانند حرف «ی»، «و» و «ک». به دنبال آن ویژگی‌ها برای پاراگراف‌های فایل DOCX پایان‌نامه‌های ورودی محاسبه می‌شود. سپس، پاراگراف‌ها با کمک ماشین بردار پشتیبان در دو گروه بدنه و فراداده طبقه‌بندی می‌شوند. در نهایت، در فاز پس‌پردازش بررسی می‌شود که از کدام پاراگراف به بعد اکثریت پاراگراف‌ها برچسب بدنه دارند. به این ترتیب، بخش‌های آغازین پایان‌نامه به دست می‌آید و می‌توان فراداده‌های مورد نظر را از این بخش استخراج کرد (شکل ۳).



شکل ۳. مراحل روش پیشنهادی برای تشخیص مرز بدنه پایان‌نامه‌های فارسی

یکی از بخش‌های سیستم ارائه شده در مقاله «کازیک» و همکاران، استخراج محتوا با هدف تعیین نقش هر بخش از مقاله است (Tkaczyki et al. 2015). بدین ترتیب، بخش‌های مختلف مقاله ورودی، با الگوریتم ماشین بردار پشتیبان در چهار گروه فراداده (عنوان، نویسندگان، چکیده، کلمات کلیدی و ...)، منابع، بدنه و سایر (تشکر و قدردانی، شماره صفحه و ...) طبقه‌بندی می‌شوند. بدین منظور، ویژگی‌هایی از بخش‌های مختلف مقاله استخراج می‌شود. در روش پیشنهادی نیز ویژگی‌های به دست آمده از هر پاراگراف، از برخی ویژگی‌های نام‌برده در سیستم ارائه شده در مقاله «کازیک» اقتباس شده است.

ویژگی‌های استخراج شده به چهار دسته کلی تقسیم می‌شوند: ویژگی‌های مکاشفه‌ای، ویژگی‌های هندسی، ویژگی‌های مربوط به توالی<sup>۱</sup>، و ویژگی‌های لغوی. در مجموع، ۱۵ ویژگی برای هر پاراگراف محاسبه شده است (جدول ۱). ویژگی‌های ۱ تا ۱۲ از نوع ویژگی‌های مکاشفه‌ای هستند. ویژگی‌های ۱۳، ۱۴ و ۱۵ به ترتیب، ویژگی‌های مربوط به توالی پاراگراف‌ها، لغوی، و هندسی هستند. در ویژگی ۱۴ منظور از کلمات کلیدی رایج در پاراگراف‌های ابتدای بدنه، عباراتی مانند مقدمه، فصل اول، کلیات تحقیق و ... هستند.

تشخیص مرز پاراگراف‌های فراداده و بدنه با استفاده از الگوریتم ماشین بردار پشتیبان به یک مسئله طبقه‌بندی دو گروهی تبدیل شده است. فرایند آموزش و بهینه‌سازی مشابه مقاله (Tharwat, Hassanien & Elnaghi (2017) بوده و تابع پایه شعاعی به‌عنوان هسته الگوریتم ماشین بردار پشتیبان انتخاب شده است. نتایج به‌دست آمده در این مقاله نشان می‌دهد که ترکیب ماشین بردار پشتیبان با الگوریتم خفاش، نرخ خطای کمتری نسبت به سایر روش‌ها، از جمله الگوریتم پرندگان<sup>۲</sup> دارد؛ زیرا در الگوریتم پرندگان، تمامی پارامترها در ابتدای الگوریتم مقادیر ثابتی دارند، اما پارامترهای الگوریتم خفاش در طول تکرار می‌توانند تغییر یابند. به عبارت دیگر، الگوریتم پرندگان نوع خاصی از الگوریتم خفاش است. همچنین، نتایج آن‌ها نشان می‌دهد که بهینه‌سازی پارامترهای ماشین بردار پشتیبان با الگوریتم ژنتیک در مقایسه با الگوریتم خفاش، نرخ خطای بیشتری دارد. از این‌رو، در پژوهش حاضر، برای بهینه‌سازی پارامترها از الگوریتم خفاش به‌عنوان الگوریتم تکاملی استفاده شده است. پارامترهای اولیه الگوریتم خفاش بر اساس مقاله پایه مقاردهی شده (جدول ۲) و برای یافتن تعداد خفاش‌ها و تعداد تکرار، آزمایش‌های همین مقاله انجام شده است.

---

1. sequential

2. particle swarm optimization (PSO)

## جدول ۱. ویژگی‌های استخراج‌شده از پاراگراف‌ها

ویژگی‌ها	
۱. تعداد تکرار کاراکترهای 'آ' یا 'ا'	۹. تعداد تکرار کاراکترهای '؟' یا '؟' نسبت به کل کاراکترها
۲. تعداد تکرار کاراکترهای 'آ' یا 'ا' نسبت به کل کاراکترها	۱۰. تعداد تکرار علائم نگارشی نسبت به کل کاراکترها
۳. تعداد ارقام فارسی و انگلیسی	۱۱. تعداد کلمات پاراگراف
۴. تعداد ارقام فارسی و انگلیسی نسبت به کل کاراکترها	۱۲. آیا پاراگراف با عدد شروع شده است؟
۵. تعداد تکرار کاراکتر نقطه (.)	۱۳. شماره پاراگراف
۶. تعداد تکرار کاراکتر نقطه نسبت به کل کاراکترها	۱۴. تعداد کلمات کلیدی رایج در پاراگراف‌های ابتدای بدنه
۷. تعداد خطوط	۱۵. نسبت طول به عرض پاراگراف
۸. تعداد تکرار کاراکترهای '؟' یا '؟'	

## جدول ۲. مقداردهی اولیه پارامترهای الگوریتم خفاش

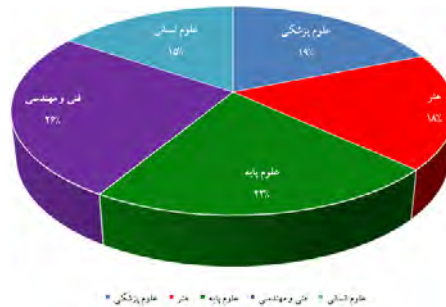
پارامتر	مقدار
فرکانس (مقادیر کمینه و بیشینه)	$f_{max} = 2, f_{min} = 0$
نرخ پالس	$r=0.5$
بلندی صدا	$A=0.5$

پس از دسته‌بندی پاراگراف‌ها در دو گروه فراداده و بدنه، ضروری است پاراگراف شروع بدنه اصلی پایان‌نامه شناسایی شود. به همین منظور، بررسی می‌شود که از کدام پاراگراف به بعد، اکثریت پاراگراف‌ها در گروه بدنه دسته‌بندی شده‌اند. بدین ترتیب، بخش‌های آغازین پایان‌نامه ورودی از متن اصلی تفکیک می‌شود.

در این مقاله زبان برنامه‌نویسی «پایتون»<sup>۱</sup> نسخه ۳/۶ برای پیاده‌سازی فرایند یادگیری و بهینه‌سازی به کار گرفته شده است. الگوریتم‌ها روی سیستمی با پردازنده «اینتل کور»<sup>۲</sup>، حافظه ۳۲ گیگابایت و سیستم عامل «اوبونتو»<sup>۴</sup> ۱۶/۰۴ اجرا شده است.

## ۴. تجزیه و تحلیل یافته‌ها

در این بخش دربارهٔ مجموعهٔ داده‌های به‌کار رفته، آزمایش‌های انجام‌شده برای یافتن پارامترها و نتایج به‌دست‌آمده از عملکرد طبقه‌بندی توضیح داده می‌شود.

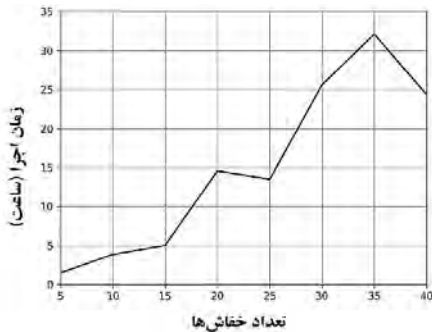


شکل ۴. نمودار دسته‌بندی پایان‌نامه‌ها

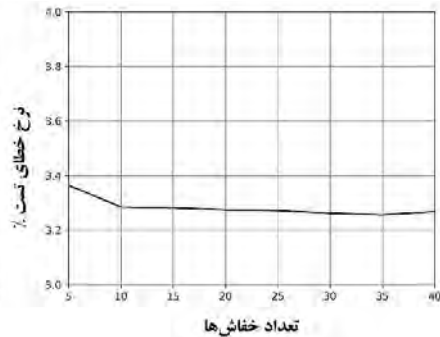
مجموعهٔ داده‌ها شامل ۲۵۰ پایان‌نامهٔ ثبت‌شده در «پژوهشگاه علوم و فناوری اطلاعات ایران (ایرانداک)»<sup>۱</sup> و در قالب DOCX هستند. داده‌ها به‌صورت تصادفی و از رشته‌های مختلف انتخاب شده‌اند. پایان‌نامه‌های انتخابی در پنج دستهٔ علوم انسانی، علوم پایه، فنی و مهندسی، هنر، و علوم پزشکی قرار می‌گیرند (شکل ۴). از فایل DOCX می‌توان اطلاعاتی چون محتوای پاراگراف‌ها، جداول، قالب و فونت، حاشیهٔ صفحات و ... را به‌دست آورد که در محاسبهٔ ویژگی‌ها مورد استفاده قرار می‌گیرند. از آنجا که جداکنندهٔ صفحات در فایل DOCX ظاهر نمی‌شود، در این پژوهش پاراگراف‌های پایان‌نامه‌ها مورد بررسی قرار گرفته‌اند. در مجموع، از ۲۵۰ پایان‌نامه، ۸۶۲۷۰ پاراگراف برای یادگیری به‌دست آمد. از بین آن‌ها تعداد ۴۰۴۱۸ پاراگراف‌های فراداده و تعداد ۴۵۸۵۲ پاراگراف‌های بدنهٔ اصلی هستند.

برای تعیین تعداد خفاش‌ها و تعداد تکرار آزمایش‌هایی انجام شد. در نخستین آزمایش میزان خطا به ازای مقادیر ۵ تا ۴۰ خفاش به‌دست آمد و زمان مورد نیاز برای اجرای الگوریتم محاسبه شد (نمودار ۱).

1. <https://irandoc.ac.ir/>



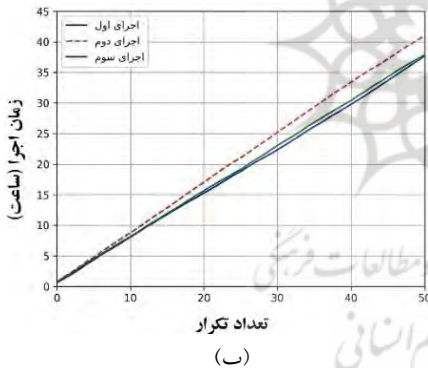
(ب)



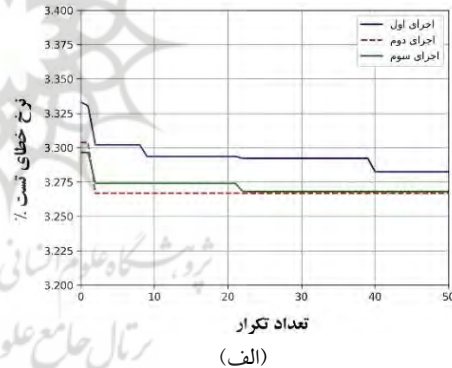
(الف)

نمودار ۱. آزمایش تأثیر تعداد خفاش بر عملکرد الگوریتم: (الف) نرخ خطای تست به ازای مقادیر مختلف خفاش‌ها؛ (ب) زمان اجرا به ازای مقادیر مختلف خفاش‌ها

مطابق نمودار ۱، با افزایش تعداد خفاش‌ها، زمان اجرا نیز افزایش می‌یابد. همچنین، برای بررسی تأثیر تعداد تکرار بر عملکرد الگوریتم، سه آزمایش انجام شده است (نمودار ۲).



(ب)

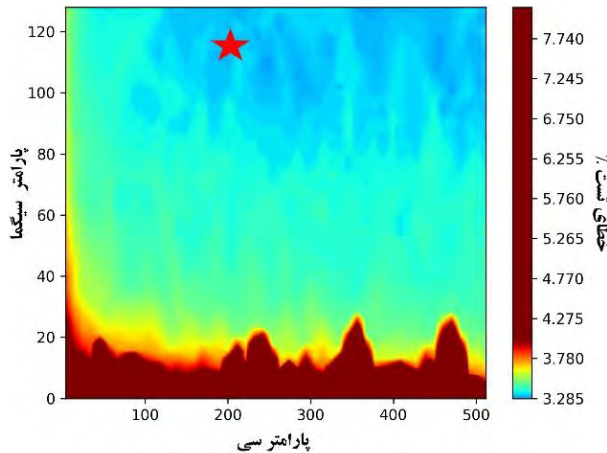


(الف)

نمودار ۲. آزمایش تأثیر تعداد خفاش بر عملکرد الگوریتم: (الف) نرخ خطای تست به تعداد تکرارهای مختلف؛ (ب) زمان اجرا به ازای افزایش تعداد تکرار

تابع نرخ خطا در این مسئله به ازای مقادیر مختلف پارامترهای C و  $\sigma$  رسم شده است (نمودار ۳). همان‌طور که مشاهده می‌شود، نرخ خطا دارای تعداد زیادی کمینه محلی است.





نمودار ۳. نمودار تراز نرخ خطا به ازای پارامترهای خطا (C) و سیگما. نقطه‌ای که با ستاره مشخص شده نقطه‌ی کمینه است که به کمک الگوریتم خفاش به دست آمده است

نتایج به دست آمده از الگوریتم خفاش با روش جست و جوی شبکه‌ای<sup>۱</sup> نیز مقایسه شده است (جدول ۳).

جدول ۳. مقایسه عملکرد الگوریتم خفاش و جست و جوی شبکه‌ای

BA-SVM	Grid search-SVM	درصد خطای تست
۳/۲۶۸	۳/۳۱۴	

برای ارزیابی فرایند یادگیری، داده‌های آموزشی به روش اعتبارسنجی متقابل k-fold و  $k=10$  جداسازی شده‌اند. با توجه به ویژگی‌های استخراج شده برای هر پاراگراف (جدول ۱)، ۸۰ درصد ویژگی‌های محاسبه شده از نوع ویژگی‌های مکاشفه‌ای هستند؛ از این رو، تأثیر سایر ویژگی‌ها و ترکیب آن‌ها با ویژگی‌های مکاشفه‌ای بر عملکرد طبقه‌بندی بررسی شده و نتایج در جدول ۴، خلاصه شده است. با توجه به یافته‌های به دست آمده، به نظر می‌رسد ویژگی‌های هندسی تأثیری بر بهبود عملکرد نداشته و مفید نیست.

1. grid search

#### جدول ۴. تأثیر انواع ویژگی‌ها بر عملکرد طبقه‌بندی پاراگراف‌ها

نوع ویژگی / معیارهای ارزیابی (میانگین)	دقت	صحت	بازخوانی	امتیاز اف ۱
مکاشفه‌ای	۷۶/۰	۸۷/۴	۷۶/۰	۷۹/۵
مکاشفه‌ای + هندسی	۷۶/۱	۸۷/۴	۷۶/۱	۷۹/۵
مکاشفه‌ای + لغوی	۸۷/۳	۸۳/۶	۸۳/۲	۸۳/۴
مکاشفه‌ای + توالی	۹۴/۴	۹۵/۵	۹۵/۲	۹۵/۳
مکاشفه‌ای + هندسی + لغوی	۸۳/۴	۸۳/۷	۸۳/۴	۸۳/۵
مکاشفه‌ای + توالی + لغوی	۹۶/۷	۹۵/۸	۹۵/۵	۹۵/۷
مکاشفه‌ای + هندسی + توالی	۹۵/۲	۹۵/۵	۹۵/۲	۹۵/۳
مکاشفه‌ای + هندسی + توالی + لغوی	۹۶/۶	۹۵/۸	۹۵/۴	۹۵/۶

برای ارزیابی عملکرد طبقه‌بندی، معیارهای ارزیابی دقت، صحت، بازخوانی<sup>۱</sup> و امتیاز اف ۱، برای گروه‌های فراداده و بدنه محاسبه شده است (جدول ۵). عملکرد روش پیشنهادی با بخش استخراج محتوا و طبقه‌بندی اولیه قسمت‌های مختلف مقاله در سیستم «سرمن»<sup>۲</sup> نیز مقایسه شده است (Tkaczyk 2015) (جدول ۶). سیستم «سرمن» در ارزیابی طبقه‌بندی الگوریتم ماشین بردار پشتیبان هسته‌های خطی، چندجمله‌ای، تابع پایه شعاعی و سیگموید<sup>۳</sup> را بررسی کرده و بهترین نتیجه در استفاده از هسته تابع پایه شعاعی به دست آمده است. در سیستم پیشنهادی نیز هسته تابع پایه شعاعی مورد استفاده قرار گرفته است.

#### جدول ۵. نتیجه طبقه‌بندی پاراگراف‌های پایان‌نامه

پاراگراف‌های فراداده	دقت	صحت	بازخوانی	امتیاز اف ۱
پاراگراف‌های فراداده	۹۲/۹	۹۴/۱	۹۲/۹	۹۳/۵
پاراگراف‌های بدنه	۹۷/۹	۹۷/۵	۹۷/۹	۹۷/۷
میانگین	۹۶/۶	۹۵/۸	۹۵/۴	۹۵/۶

1. recall

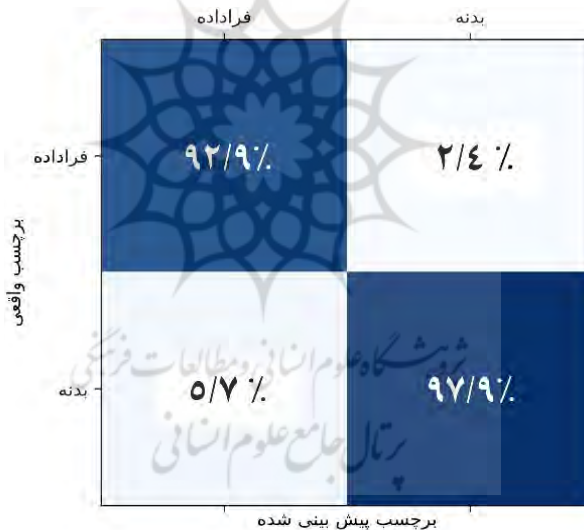
2. CERMINE

3. Sigmoid

جدول ۶. مقایسه عملکرد روش پیشنهادی و سیستم «سرمین» (Tkaczyk et al. 2015)

معیار/ روش	سرمین	روش پیشنهادی
میانگین امتیاز اف ۱ (%)	۹۳/۹	۹۵/۶

ماتریس درهم‌ریختگی دو گروه بدنه و فراداده نیز در شکل ۸، آمده است. همان‌طور که انتظار می‌رود، امکان رخداد دو خطا وجود دارد: نخست، تخصیص برچسب بدنه به پاراگراف‌های فراداده، و دوم تخصیص برچسب فراداده به پاراگراف‌های بدنه. خطای نوع اول به‌طور معمول، در پاراگراف‌های طولانی‌تر رخ می‌دهد. برای نمونه، پاراگراف‌های چکیده از دیگر بخش‌های آغازین پایان‌نامه دارای کاراکترهای بیشتری بود. از این رو، طولانی‌تر از سایر بخش‌هاست. در مقابل، خطای نوع دوم در تیرها یا پاراگراف‌های کوتاه بدنه اصلی پایان‌نامه امکان رخداد دارد.



شکل ۸. ماتریس درهم‌ریختگی

به‌منظور مقایسه عملکرد طبقه‌بندی در انواع پایان‌نامه‌ها، آزمایش دیگری انجام شده است (جدول ۷). همان‌طور که پیش‌تر اشاره شد، پایان‌نامه‌ها در پنج دسته علوم پایه، علوم انسانی، فنی و مهندسی، هنر، و علوم پزشکی قرار می‌گیرند. طبق نتایج به‌دست آمده، به نظر می‌رسد که طبقه‌بندی پاراگراف‌ها در پایان‌نامه‌های علوم انسانی، عملکرد بهتری نسبت به سایر رشته‌ها دارد.

### جدول ۷. مقایسه عملکرد طبقه‌بند در انواع پایان‌نامه‌ها

دسته‌بندی	میانگین دقت (درصد)	میانگین صحت (درصد)	میانگین بازخوانی (درصد)
علوم پایه	۹۷/۰	۹۶/۶	۹۶/۰
علوم انسانی	۹۷/۵	۹۷/۰	۹۶/۵
فنی و مهندسی	۹۶/۵	۹۵/۰	۹۴/۰
هنر	۹۷/۴	۹۶/۹	۹۶/۴
علوم پزشکی	۹۶/۱	۹۵/۸	۹۵/۱

### ۵. نتیجه‌گیری

در پژوهش حاضر، جداسازی خودکار مرز فراداده‌ها از بدنه پایان‌نامه‌های فارسی مورد بررسی قرار گرفته است. بدین منظور، مجموعه داده‌ای متشکل از پایان‌نامه‌های علوم پایه، علوم انسانی، فنی و مهندسی، هنر، و علوم پزشکی جمع‌آوری شده است. برای پاراگراف‌های هر پایان‌نامه، ویژگی‌های مکاشفه‌ای، لغوی، توالی و هندسی محاسبه می‌شود. بر اساس یافته‌های پژوهش به نظر می‌رسد که ویژگی هندسی تأثیری در بهبود عملکرد طبقه‌بند نداشته است. سایر انواع ویژگی‌های تعریف‌شده برای حل این مسئله کفایت می‌کنند. برای طبقه‌بندی پاراگراف‌های پایان‌نامه‌ها، از الگوریتم ماشین بردار پشتیبان استفاده شده و پارامترهای آن با الگوریتم فرامکاشفه‌ای خفاش، بهینه‌سازی شده است. پاراگراف‌ها در دو گروه فراداده و بدنه طبقه‌بندی شده و نتایج نشان می‌دهد که مدل یادگیری به‌دست آمده با دقت ۹۶/۶ درصد می‌تواند نوع پاراگراف را تشخیص دهد. افزون بر این، طبق نتایج به‌دست آمده، عملکرد طبقه‌بند در پایان‌نامه‌های علوم انسانی بهتر از سایر رشته‌های دیگر بوده است. به نظر می‌رسد که تنوع در ساختار پایان‌نامه‌ها و تفاوت نوشتارها در رشته‌های گوناگون موجب تفاوت در نتایج به‌دست آمده شده است. با توجه به یافته‌های به‌دست آمده و عملکرد خوب روش پیشنهادی، برای پژوهش‌های آینده می‌توان پاراگراف‌های فراداده را با هدف دستیابی به فراداده‌های از پیش تعریف‌شده پایان‌نامه‌ها تحلیل کرد.

### قدردانی

بدین وسیله از حمایت آزمایشگاه متن کاوی و یادگیری ماشین «پژوهشگاه علوم و فناوری اطلاعات ایران» در انجام این پژوهش تشکر و قدردانی می‌شود.

## فهرست منابع

تن‌سازان، امیر، و محمدامین مهدوی. ۱۳۹۶. استخراج فراداده‌های متنی از مقاله‌های علمی به زبان فارسی با مدل آماری CRF پژوهش‌های نظری و کاربردی در علم اطلاعات و دانش‌شناسی ۷(۱): ۳۰۴-۳۲۱.

## References

- Adrian, W., N. Leone, M. Manna, & C. Marte. 2017. Document Layout Analysis for Semantic Information Extraction. In: Esposito F., Basili R., Ferilli S., Lisi F. (eds) *AI\*IA 2017 Advances in Artificial Intelligence*. *AI\*IA 2017. Lecture Notes in Computer Science*, vol 10640. Springer, Cham. [https://doi.org/10.1007/978-3-319-70169-1\\_20](https://doi.org/10.1007/978-3-319-70169-1_20)
- Boukhers, Z., S. Ambhore, & S. Staab., 2019. An end-to-end approach for extracting and segmenting high-variance references from pdf documents. *ACM/IEEE Joint Conference on Digital Libraries (JCDL)*. Champaign, IL, USA.
- Cortes, C., & V. Vapnik. 1995. Support-vector networks. *Machine Learning* volume. 20297–273 .:
- Cuong, N., M. Kumar, M.-Y. Kan, & W. Lee. 2015. Scholarly document information extraction using extensible features for efficient higher order semi-crfs. *JCDL '15: proceedings of the 15th acm/ieee-cs joint conference on digital libraries*. Knoxville, Tennessee, USA.
- Do, H., M. Chandrasekaran, P. Cho, & M.-Y. Kan. 2013. Extracting and matching authors and affiliations inscholarly documents. *JCDL '13: proceedings of the 13th acm/ieee-cs joint conference on digital libraries*. Indianapolis, Indiana, USA
- Ferrés, D., H. Saggion, F. Ronzano, & À. Bravo. 2018. PDFdigest: an adaptable layout-aware pdf-to-xml textual content extractor for scientific articles. 11<sup>th</sup> Language Resources and Evaluation Conference (LREC). Miyazaki, Japan
- Kan, M.-Y., M. Luong, & T. Nguyen. 2010. Logical structure recovery in scholarly articles with rich document features. *International Journal of Digital Library Systems* 1 (4): 1-23.
- Kern, R., K. Jack, M. Hristakeva, & M. Granitzer. 2012. TeamBeam — meta-data extraction from scientific literature. *D-Lib Magazine*.
- Kooli, N., & A. Belaid. 2016. Inexact graph matching for entity recognition in OCRed documents. 23rd International Conference on Pattern Recognition (ICPR). Cancun, Mexico. IEEE.
- Liu, R., L. Gao, D. An, Z. Jiang, & Z. Tang. 2017. Automatic document metadata extraction based on deep networks. *Natural Language Processing and Chinese Computing*. Lecture Notes in Computer Science, 10619305-317 .:
- Meng, B., L. Hou, E. Yang, & J. Li. 2018. Metadata extraction for scientific papers. *Chinese Computational Linguistics and Natural Language Processing Based on Naturally Annotated Big Data*. CCL 2018, NLP-NABD 2018. Lecture Notes in Computer Science, 11221: 111-122
- Nasar, Z., S. Jaffry, & M. Malik. 2018. Information extraction from scientific articles: a survey. *Scientometrics* 1171990–1931 (3) .
- Peng, F., & A. McCallum. 2006. Information extraction from research papers using conditional random fields. *Information Processing & Management* 42963-979 (4) .
- Piskorski, J., & R. Yangarber. 2013. Information extraction: past, present and future. In *Multi-source, Multilingual Information Extraction and Summarization*. Berlin: Springer.
- Qiu, S., & T. Zhou. 2019. A method of extracting metadata information in digital books. 10th International Conference on Information Technology in Medicine and Education (ITME). Qingdao, China.
- Rizvi, S., A. Dengel, & S. Ahmed 2019. DeepBiRD: An automatic bibliographic reference detection approach.

- Safder, I., S.-U. Hassan, A. Visvizi, T. Noraset, R. Nawaz, & S. Tuarob2020 .. Deep learning-based extraction of algorithmic metadata in full-text scholarly documents. *Information Processing & Management* 57102269 (6) .
- Souza, A., V. Moreira, & C. Heuser2014 .. ARCTIC: metadata extraction from scientific papers in pdf using two-layer CRF. DocEng '14: Proceedings of the 2014 ACM symposium on Document engineering. New York, NY, USA.
- Tharwat, A., A. Hassanien, & B. Elnaghi2017 .. A BA-based algorithm for parameter optimization of Support Vector Machine. *Pattern Recognition Letters* 9313-22 .:
- Tkaczyk, D., P. Szostek, M. Fedoryszak, P. Dendek, & Ł. Bolikowski2015 .. CERMINE: automatic extraction of structured metadata from scientific literature. *International Journal on Document Analysis and Recognition (IJ DAR)* 18: 317–335.
- Yang, X.-S. 2010. A new metaheuristic bat-inspired algorithm. *Nature Inspired Cooperative Strategies for Optimization (NICSO 2010)* 284: 65-74.

#### محدثه رهنما

متولد سال ۱۳۷۴، دانشجوی کارشناسی ارشد مهندسی کامپیوتر گرایش هوش مصنوعی در دانشگاه الزهرا (س) است. یادگیری ماشین، پردازش زبان‌های طبیعی و پردازش تصویر از جمله علایق پژوهشی وی است.



#### سید محمدحسین هاشمی نژاد

متولد سال ۱۳۶۴، دارای مدرک تحصیلی دکتری در رشته مهندسی نرم‌افزار از دانشگاه تربیت مدرس است. ایشان هم‌اکنون استادیار دانشگاه الزهرا (س)، گروه مهندسی کامپیوتر است. موضوعات معماری نرم‌افزار، یادگیری ماشین و الگوریتم‌های تکاملی علایق پژوهشی وی است.



#### جلال‌الدین نصیری

متولد سال ۱۳۶۲، دارای مدرک تحصیلی دکتری در رشته مهندسی کامپیوتر، گرایش نرم‌افزار از دانشگاه تربیت مدرس است. ایشان هم‌اکنون استادیار گروه زبان‌شناسی رایانشی پژوهشکده علوم اطلاعات و مدیر آزمایشگاه متن‌کاوی و یادگیری ماشین در پژوهشگاه علوم و فناوری اطلاعات ایران (ایرانداک) است.



پردازش زبان‌های طبیعی و یادگیری ماشین از جمله علایق پژوهشی وی است.