



Automatic Persian Multi-Text Summarization Techniques based on Meta-Heuristic Algorithms

F. Ahangari¹
S. Karbasi²
M. Yaghoubi³

Received: 7, Jan. 2019
Accepted: 12, May 2019

doi: 10.30484/nastinfo.2019.2330

Purpose: The main objective of this study is to present a pattern for standard summarization of Persian texts with the approach of converting the problem to optimization problem by compatible meta-heuristic algorithms.

Methodology: In this research, standard multi-text "Pasokh" collection, which contains 50 different types of news from the most popular news agencies in Iran, each containing 20 documents, as well as 5 summaries of abstractive and 5 extractive, used for evaluation. First, the preprocessing performed on the input texts and the initial summary generated with TF-ISF benchmark, readability and consistency criteria of the sentences, similarity to the title, position of the sentence in the text, and the length of the sentence. With respect to each of these criteria, weighting function assigned to extracted sentences and a similarity matrix created. Then, output of the extraction system processed by Genetic algorithm and Cuckoo search algorithm for the final summary. Eventually, the output obtained from the previous step analyzed using the Rouge evaluation tools and the comparison with the human abstracts.

Findings: The average of all values obtained in Rouge evaluation tools for calculation the overlapping of common samples of human summaries and machine summaries by Cuckoo search algorithm were higher than the values obtained by Genetic algorithm as well as Ijaz online summarizer system.

Meanwhile, among the eight criteria, the longest common sub-sentence with a value of 0.33 and the number of common words in the text with 0.40 were better than the rest.

Conclusion: The results of the comparison of two algorithms indicate that the Cuckoo search algorithm is better in the entire criteria. On the other hand, comparing the results suggests that the average time calculated for summarizing by the proposed system is also less.

Keywords: Automatic text summarization, Extractive summarization, Meta-Heuristic algorithms, Genetic algorithm, Cuckoo search algorithm, Rouge evaluation tools

¹ MA of Computer Science, Golestan University, Gorgan, fatemehahangari.fa@gmail.com

² Assistant Professor, Computer Science, Golestan University, Gorgan (Corresponding author), s.karbasi@gu.ac.ir

³ Assistant Professor, Computer Science, Golestan University, Gorgan, m.yaghoubi@gu.ac.ir



تکنیک‌های خلاصه‌سازی چندسندی خودکار متون فارسی مبتنی بر الگوریتم‌های فرااکتشافی

فاطمه آهنگری¹
سهیلا کرباسی²
مهدی یعقوبی³

چکیده:

هدف: ارائه الگوی خلاصه‌سازی استاندارد متون فارسی با رویکرد تبدیل مسئله خلاصه‌سازی به مسئله بهینه‌سازی توسط الگوریتم‌های فرااکتشافی سازگار. **روش‌شناسی:** در این پژوهش از اسناد استاندارد پیکره چندسندی «پاسخ» که شامل ۵۰ موضوع مختلف از انواع گونه‌های خبری از خبرگزاری‌های پربیننده ایران، برای ارزیابی استفاده شده است. هر موضوع حاوی ۲۰ سند و همچنین ۵ خلاصه چکیده‌ای و ۵ خلاصه استخراجی است. ابتدا عملیات پیش‌پردازش روی متون ورودی انجام و خلاصه‌های اولیه تولید شدند. این کار به کمک معیار TF-ISF، معیارهای خوانایی و انسجام جملات، ویژگی شباهت با عنوان، ویژگی موقعیت جمله در متن، و ویژگی طول جمله انجام شد. با توجه به هر یک از این معیارها، وزنی به هر یک از جملات خلاصه اختصاص داده و ماتریس شباهت ایجاد شد. سپس، خروجی سیستم استخراج توسط دو الگوریتم فرااکتشافی ژنتیک و جستجوی فاخته برای رسیدن به خلاصه نهایی پردازش شد. در نهایت، خروجی به دست آمده از مرحله قبل به کمک ابزار ارزیابی Rouge و مقایسه با خلاصه‌های انسانی تحلیل شدند.

یافته‌ها: میانگین همه مقادیر به دست آمده از ابزار ارزیابی Rouge در محاسبه میزان هم‌پوشانی نمونه‌های مشترک خلاصه‌های انسانی و خلاصه ماشینی توسط الگوریتم جستجوی فاخته بیشتر از مقادیر به دست آمده توسط الگوریتم ژنتیک و همچنین سامانه خلاصه‌ساز برخط ایجاز بودند. از میان هشت معیار موجود در این ابزار، دو معیار ارزیابی طولانی‌ترین زیررشته مشترک با مقدار ۰.۳۳ و تعداد لغات مشابه در متن با مقدار ۰.۴۰ نتایج بهتری نسبت به بقیه معیارها داشتند.

نتیجه‌گیری: نتایج حاصل از مقایسه دو الگوریتم به کاررفته، حاکی از عملکرد بهتر الگوریتم جستجوی فاخته در هر یک از معیارهای ابزار Rouge است. از طرفی مقایسه زمانی نتایج نشان می‌دهد که میانگین زمانی محاسبه شده برای خلاصه‌سازی توسط سیستم پیشنهادی با الگوریتم جستجوی فاخته کمتر است.

کلیدواژه‌ها: خلاصه‌سازی خودکار متن، خلاصه استخراجی، الگوریتم‌های فرااکتشافی، الگوریتم ژنتیک، الگوریتم جستجوی فاخته، ابزار ارزیابی Rouge

¹ کارشناس ارشد کامپیوتر، دانشکده فنی و مهندسی، دانشگاه گلستان، گرگان Fatemehahangari.fa@gmail.com

² استادیار گروه کامپیوتر، دانشکده فنی و مهندسی، دانشگاه گلستان، گرگان (نویسنده مسئول) S.karbasi@gu.ac.ir

³ استادیار گروه کامپیوتر، دانشکده فنی و مهندسی، دانشگاه گلستان، گرگان M.yaghoubi@gu.ac.ir

امروزه با توجه به افزایش حجم مستندات متنی، توسعه منابع مبتنی بر وب و نیازهای اساسی به نگهداری، دسته‌بندی، بازیابی و پردازش آنها، توجه به پردازش زبان‌های طبیعی و بهره‌گیری از ابزارهایی نظیر خلاصه‌سازهای خودکار و مترجم‌های ماشینی، بیش از پیش احساس می‌شود. فرایند فشرده‌سازی متن به‌نحوی که ویژگی‌ها و نکات اصلی سند اولیه حفظ شود، خلاصه‌سازی نام دارد که به استفاده از منابع بیشتر با سرعت بالاتر و در نتیجه به‌دست‌آوردن اطلاعات مفیدتر منجر می‌شود. به‌طور کلی، خلاصه‌سازی یکی از زیر مجموعه مشکلات پیچیده پردازش زبان طبیعی است که هنوز زمان زیادی باقی است تا به‌حدی برسیم که ماشین همانند انسان به‌طور کامل مفهوم متن را استخراج کند.

یک خلاصه خوب و به‌اصطلاح باکیفیت خلاصه‌ای است که خوانا باشد و جملاتش درباره موضوعی یکسان و مرتبط بحث کند. همچنین جمله‌ها باید با عنوان سند ارتباط داشته باشند و اطلاعاتی درباره عنوان ارائه دهند. بنابر تعریف ارائه‌شده در استاندارد ایزو 215¹، خلاصه یک بازگویی مختصر از سند است. خلاصه‌سازی خودکار سند، یعنی تولید یک نسخه مختصرتر از سند اصلی توسط یک سیستم کامپیوتری به‌نحوی که ویژگی‌ها و نکته‌های اصلی سند اولیه حفظ شود (هوی²، 2005).

خلاصه‌سازی را می‌توان به‌عنوان یک روش برای نمایش بخش‌های اصلی یک سند و یا اطلاع‌رسانی سریع با پوشش تمام اطلاعات متن اصلی بیان کرد. در هر صورت، مهم‌ترین هدف خلاصه‌سازی ارائه تصویری کلان از محتوای متن است که باعث کاهش زمان خواندن متن اصلی می‌شود. در این راستا، ابزارهای خلاصه‌سازی متن برای تشخیص عناوین و موضوعات کلیدی متون استفاده می‌شوند که این امور نیز به‌نوبه خود در بیشتر کاربردهای سیستم‌های ماشینی ضروری هستند (طالب‌علی و ریاحی، 1394؛ عرب احمدی، 1397). از مهم‌ترین مزایای خلاصه‌سازی خودکار نسبت به خلاصه انسانی می‌توان به قابل کنترل بودن اندازه خلاصه، قابل پیش‌بینی بودن محتوای آن، و قابلیت شناسایی متن خلاصه در متن اصلی اشاره کرد. به‌طور کلی، روش‌های خلاصه‌سازی متن را می‌توان به خلاصه استخراجی و چکیده‌ای طبقه‌بندی کرد (هوی، 2005). با این حال، روش طبقه‌بندی دیگری متشکل از خلاصه‌سازی براساس تعداد اسناد ورودی وجود دارد که شامل خلاصه چندسندی و تک‌سندی است. طبقه‌بندی براساس هدف خلاصه‌سازی نیز نوع دیگری است که شامل روش‌های پرس‌وجوی عمومی است (گوپتا³، 2010).

¹ ISO 215

² Hovy

³ Gupta

روش خلاصه‌سازی استخراجی شامل انتخاب جمله‌ها و پاراگراف‌های مهم از سند اصلی و ترکیب آنها به شکل کوتاه‌تر است که در این مطالعه استفاده شده است. روش خلاصه‌سازی چکیده‌ای شامل درک متن اصلی با استفاده از روش‌های زبانی برای بررسی و تفسیر متن و سپس پیدا کردن مفاهیم و عبارت‌های جدید به بهترین نحو است و در نهایت متنی کوتاه‌تر تولید می‌کند که مهم‌ترین اطلاعات سند اصلی را بیان می‌کند.

روش‌های متعددی برای استخراج متن وجود دارد که TF-IDF¹ یکی از اصلی‌ترین آنهاست (لدنوا، ژلبوخ، و هرناندز²، 2008؛ هرناندز و لدنوا³، 2009). در این روش، وزن‌دهی جملات براساس تکرار کلمه‌ها و جمله‌ها انجام می‌شود. اصطلاح «تکرار جمله» به تعدادی از جمله‌های سند اشاره دارد که شامل کلمه‌های پرتکرار هستند.

استفاده از خلاصه‌سازی متن براساس خوشه‌بندی جمله‌ها نیز از روش‌های معمول است (زنگ و لی⁴، 2009). در این روش، هر خوشه نشان‌دهنده یک موضوع است. شباهت‌ها میان جمله‌ها براساس مجموعه‌ای از پارامترها بررسی می‌شوند و سپس عبارت‌های مشابه در یک خوشه قرار می‌گیرند. در هر خوشه، به جمله‌هایی که شباهت بیشتری به عنوان خوشه دارند، امتیازهای بیشتری اختصاص می‌یابد و در نهایت، می‌توانند برای خلاصه انتخاب شوند (مشکی، 1388). مزیت اصلی این روش، این است که موضوع هر متن به‌سادگی مشخص می‌شود؛ اما از آنجایی که انتخاب تعداد خوشه‌ها موضوع مهمی است و ممکن است خیلی زیاد یا خیلی کم باشد، بر نتایج خلاصه می‌تواند تأثیرگذار باشد. به عبارت دیگر، انتخاب تعداد مطلوب خوشه‌ها کار دشواری است که مهم‌ترین ضعف این روش به حساب می‌آید (زنگ و لی، 2009).

در بعضی از روش‌های خلاصه‌سازی استخراجی، مجموعه داده‌هایی که توسط کاربر برچسب‌گذاری شده‌اند به‌عنوان ابزار خلاصه‌سازی استفاده می‌شود. به عبارت دیگر، فرض می‌کنیم مجموعه‌ای از متن ورودی و متن خلاصه آنها را داریم که جمله‌ها به‌وسیله یک نشانه‌گذار خاص به چند جزء شکسته می‌شوند. هر جزء با مجموعه‌ای از ویژگی‌های از پیش تعریف‌شده (مکان، تعداد تکرار، و کلمه‌های عنوان متن) نشان داده می‌شود. سپس یک روش یادگیری نظارت‌شده به خلاصه‌ساز آموزش می‌دهد تا بخش‌های مهم را براساس بردار ویژگی استخراج کند. بعضی از این روش‌ها عبارت‌اند از: درخت تصمیم، تئوری بیزین، شبکه‌های عصبی، و منطق فازی (سانمالی، سلیم، و بینواهلان⁵، 2009؛ سونگ، چوی، پارک، و دینگ⁶، 2011). کاهش دقت و سرعت عملیات در

¹ Term Frequency-Inverse Document Frequency

² Ledeneva, Gelbukh, & Hernández

³ Hernandez & Ledeneva

⁴ Zhang & Li

⁵ Suanmali, Salim, & Binwahlan

⁶ Song, Choi, Park, & Ding

مجموعه اسناد بزرگ، از معایب اصلی این روش‌هاست. به‌عنوان نمونه، مشکل اصلی روش مبتنی بر منطق فازی کاهش دقت خلاصه‌سازی در صورت عدم تعریف دقیق قوانین است (سونگ و همکاران، 2011).

در مورد پژوهش‌های مرتبط در حوزه خلاصه‌سازی متون فارسی می‌توان به سامانه خلاصه‌ساز FarsiSum اشاره کرد که به‌صورت برخط دسترس‌پذیر است¹ (هاسل و مزدک²، 2004). این سامانه، نسخه تغییر یافته سامانه خلاصه‌ساز متون سوئدی به‌نام SweSum برای پوشش زبان فارسی است. البته این خلاصه‌ساز فقط براساس ویژگی‌های آماری عمل می‌کند و خصوصیات زبان‌شناسی متن و چالش‌های خاص زبان فارسی را در نظر نمی‌گیرد. در پژوهش کریمی و شمس‌فرد (1385) یک روش خلاصه‌سازی تک‌سندی پیشنهاد شده است که بر مبنای گزینش جمله‌ها کار می‌کند. همچنین، محتوای خلاصه می‌تواند کلی یا براساس پرس‌وجوی کاربر باشد. ایده به‌کاررفته در گزینش جمله‌ها در این خلاصه‌ساز، ترکیبی از دو روش زنجیره لغوی و نظریه گراف است.

در پژوهش مشکئی (1388) پس از بررسی موضوعات و چالش‌های مربوط به پردازش متون فارسی، روشی مبتنی بر خوشه‌بندی برای خلاصه‌سازی چندسندی متون فارسی پیشنهاد شده است که در آن از الگوریتم خوشه‌بندی Kmeans استفاده شده است. در این روش، پس از پیش‌پردازش متن، ابتدا جمله‌ها خوشه‌بندی می‌شوند و سپس به‌ازای هر خوشه، جمله‌ای گزینش می‌شود که بیشترین ارتباط را با سایر جمله‌ها دارد. جهت ارزیابی روش نیز از قضاوت انسانی استفاده شده است. بدین صورت که برای خلاصه هر یک از 10 مجموعه موجود، سه داور انسانی رأی خود را به‌صورت خوب، متوسط و ضعیف ارائه می‌دهند.

در پژوهش اخوان، شمس‌فرد، و عرفانی جورابچی (1387)، یک روش خلاصه‌سازی گزینشی پیشنهاد شده است که قابلیت به‌کارگیری در دو حالت تک‌سندی و چندسندی را دارد. در این روش، از معیارهایی مانند واژه‌های مهم، عبارت‌های اشاره، واژه‌های عنوان، نقل‌قول برای امتیازدهی به جمله‌ها بهره‌گیری شده است. همچنین، برای جلوگیری از افزونگی چنانچه دو جمله شباهتی بیش از یک مقدار آستانه داشته باشند، جمله‌ای که امتیاز کمتر دارد را نادیده می‌گیرد. از دیگر ویژگی این روش می‌توان به قابلیت دریافت درخواست از کاربر اشاره کرد.

در پژوهش هنرپیشه، قاسم ثانی، و میرروشندل³ (2008)، یک روش برای خلاصه‌سازی چندسندی پیشنهاد شده است. مبنای این روش، استفاده از اتصال میانگین خوشه‌ها در خوشه‌بندی سلسله‌مراتبی جمله‌ها و به‌کارگیری روش «تجزیه مقادیر منفرد»⁴ برای تعیین اهمیت جمله‌هاست. در این پژوهش، از دو منبع زبانی ساده، یکی برای تجزیه متن به واژه‌ها و دیگری برای تعیین فراوانی واژه‌ها در اسناد استفاده شده است.

¹ <http://swesum.nada.kth.se/index-eng.html>

² Hassel & Mazdak

³ Honarpisheh, Ghassem-Sani, & Mirroshandel

⁴ Singular Value Decomposition (SVD)

استفاده از «الگوریتم‌های فرااکتشافی»¹، تکنیک دیگری در خلاصه‌سازی و انتخاب جمله‌های باارزش است. پژوهشگران الگوریتم‌های اکتشافی و فرااکتشافی را به‌عنوان فصل مشترک تکنیک‌هایی مانند هوش مصنوعی، یادگیری ماشین، جستجوی عملیات و دیگر تکنیک‌های مهندسی می‌دانند که استفاده از این الگوریتم‌ها برای رسیدن به خلاصه‌ای باکیفیت می‌تواند بسیار مفید باشد (یو و ژن²، 2010). الگوریتم‌های ژنتیک (فتاح و رن³، 2009؛ قزوینیان، حسن‌آبادی، و حلاوتی⁴، 2008)، «بهینه‌سازی ازدحام ذرات»⁵ (فونگ و اوکسلی⁶، 2011)، و «بهینه‌سازی کلونی مورچه‌ها»⁷ (مارتین، دوبیکر، و هاسن⁸، 2007) از شناخته‌شده‌ترین روش‌های فرااکتشافی هستند که نتایج مناسبی در خلاصه‌سازی اسناد متنی نشان داده‌اند. الگوریتم‌های فرااکتشافی به الگوریتم‌هایی گفته می‌شود که برپایه رفتار طبیعی بعضی از گونه‌های خاص جانداران به‌وجود می‌آیند (مریخ بیات، 1393). این الگوریتم‌ها در بسیاری از زمینه‌ها، برای حل مسائل مربوط به بهینه‌سازی استفاده می‌شوند که توسط روش‌های معمولی قابل حل نیستند. پیشوند یونانی «متا» که در نام آن استفاده شده است نشان می‌دهد این الگوریتم‌ها سطح بالایی از هوشمندی را فراهم می‌کنند. هدف اصلی این روش‌ها، جستجوی مؤثر فضای جواب است و برخلاف الگوریتم‌های اکتشافی، وابسته به نوع خاصی از مسئله نیستند (ری و وارشنی⁹، 2015). رفتارهای شگفت‌انگیز فاخته، ایده «الگوریتم جستجوی فاخته»¹⁰ (یانگ و دب¹¹، 2014) را به‌وجود آورده است. فاخته‌ها تاکتیک تولیدمثل خاصی دارند به این صورت که تخم‌های بارور شده خود را در لانه گونه‌های دیگری قرار می‌دهند تا والدین جایگزین به‌طور ناخواسته تخم فاخته را در کنار تخم‌های خودشان پرورش دهند (گاول، شارما، و بدی¹²، 2011). فاخته به‌طور غریزی لانه‌ای را انتخاب می‌کند که تخم‌هایش زودتر از تخم‌های پرنده میزبان، به جوجه تبدیل شود. هنگامی که اولین جوجه فاخته از تخم بیرون می‌آید، اولین اقدام غریزی تخریب تخم‌های میزبان و خارج کردن آنها از لانه آغاز می‌شود. این کار سهم مواد غذایی جوجه فاخته را افزایش می‌دهد (یانگ و دب، 2009).

الگوریتم جستجوی فاخته (CSA) پارامترهای طراحی و محدودیت‌های مختلفی را در نظر می‌گیرد:

1. هر یک از فاخته‌ها لانه‌ای را به‌طور تصادفی انتخاب می‌کنند و تخم خود را در آن قرار می‌دهند.

¹ Metaheuristic algorithms

² Yu & Gen

³ Fattah & Ren

⁴ Qazvinian, Hassanabadi, & Halavati

⁵ Particle Swarm Optimization (PSO)

⁶ Foong & Oxley

⁷ Ant Colony Optimization (ACO)

⁸ Martens, De Backer, & Haesen

⁹ Rai & Varshney

¹⁰ Cuckoo search algorithm

¹¹ Yang & Deb

¹² Goel, Sharma, & Bedi

2. بهترین و باکیفیت‌ترین لانه، لانه‌ای است که از تخم محافظت می‌کند و آن را به نسل بعدی منتقل می‌کند.

3. تعداد لانه‌های میزبان همیشه ثابت است و تخمی که فاخته در لانه میزبان قرار می‌دهد به احتمال $P_a \in [0,1]$ توسط پرنده میزبان شناسایی می‌شود. در چنین شرایطی، پرنده میزبان می‌تواند تخم را دور بیندازد یا لانه را ترک کند و یک لانه جدید ایجاد کند. برای حل این مشکل، می‌توان لانه‌ها را با احتمال p_a با n لانه (راه‌حل‌های تصادفی) جدید جایگزین کرد (یانگ و دب، 2014). تکنیک بهینه‌سازی جستجوی فاخته در زمینه خلاصه‌سازی متون غیرفارسی استفاده شده و نتایج حاصل موفقیت‌آمیز بوده است. همچنین دقت و موفقیت آن بسیار بهتر از سایر روش‌ها مانند PSO تأیید شده است (روتري و بالابانتاری¹، 2017؛ میرشجاعی و معصومی²، 2015). به‌طور کلی، چالش اصلی روش‌های مبتنی بر استخراج در متون بزرگ که شامل تعداد زیادی جمله هستند، انتخاب جمله‌ها و رتبه‌بندی آنهاست. روش پیشنهادی این پژوهش، یک روش خلاصه‌سازی متون فارسی مبتنی بر الگوریتم‌های فرااکتشافی سازگار با مسئله خلاصه‌سازی است. پژوهش‌ها در زمینه خلاصه‌سازی متن نشان می‌دهد این مبحث می‌تواند از علوم گوناگون همچون هوش مصنوعی، یادگیری ماشین، تکنیک‌های جستجوی عملیات، و بهینه‌سازی بهره‌مند شود. منظور از مسائل یادگیری این است که ما از سیستم کامپیوتر انتظار داریم برای حل مسائل بتواند تا حدی هوش انسانی را تقلید کند و از تکنیک‌های «هوش مصنوعی»³ استفاده کند. یکی از اهداف اصلی در مسائل خلاصه‌سازی، داشتن خلاصه‌ای نزدیک به خلاصه انسانی است. بنابراین بهره‌مندی از تکنیک‌های یادگیری و هوش مصنوعی از ملزومات اصلی سیستم خلاصه‌ساز به‌شمار می‌رود. از طرفی، مسائل بهینه‌سازی و مسائل یادگیری نیز با یکدیگر در ارتباطند. گاهی تابعی که نیازمند بهینه‌شدن است بسیار پیچیده است و نمی‌توان ارزش هدف را برای هر راه‌حل آن محاسبه کرد. در چنین شرایطی الگوریتم‌های یادگیری برای محاسبه تقریبی تابع تناسب⁴ استفاده می‌شوند.

روش‌شناسی

به‌منظور پیاده‌سازی روش پیشنهادی و بررسی نتایج از پیکره پاسخ⁵ به‌عنوان پایگاه داده استفاده شده است. این پیکره شامل 1500 سند است که می‌توان از آنها به‌عنوان ورودی به الگوریتم‌های پیشنهادی استفاده کرد. پیکره

¹ Rautray & Balabantaray

² Mirshojaei & Masoomi

³ Artificial intelligence

⁴ Fitness function

⁵ <http://dadegan.ir/catalog/pasokh>

پاسخ اولین پیکره متنی استاندارد برای ارزیابی خلاصه‌سازی تک‌سندی و چندسندی در زبان فارسی است (بهمدی مقدس، کاهانی، طوسی، پورمعصومی، و استیری¹، 2013).

این پیکره در دو بخش تک‌سندی و چندسندی سازماندهی شده است که در این پژوهش از اسناد بخش چندسندی این پیکره استفاده شده است. پیکره چندسندی شامل 50 موضوع است که هر موضوع حاوی 20 سند بوده و همچنین هر موضوع شامل 5 خلاصه چکیده‌ای و 5 خلاصه استخراجی است. در تولید این مجموعه سعی شده است تمامی استانداردهای لازم برای تولید یک پیکره خلاصه‌سازی رعایت شود که مشخصات کامل آن در جدول 1 آمده است.

جدول 1. مشخصات اجزای مجموعه چندسندی در پیکره پاسخ

تعداد	عنوان
۵۰	موضوعات
۲۰	اسناد مربوط به هر موضوع
۷	خبرگزاری‌ها
۵	خلاصه‌های استخراجی برای هر سند
۵	خلاصه‌های چکیده‌ای برای هر سند
۱۰٪	نرخ فشردگی
۱۲	میانگین تعداد جملات اسناد

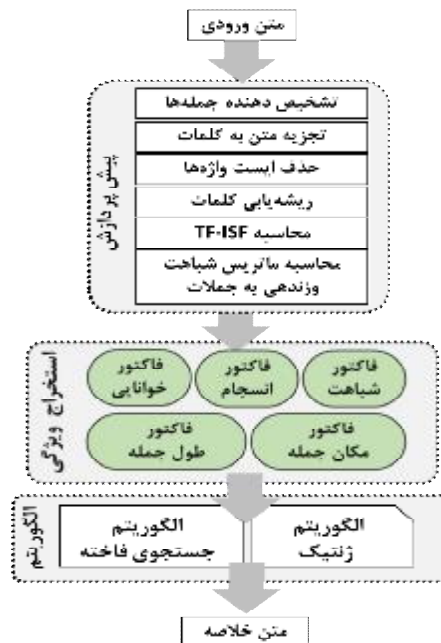
معماری کلی فرایند خلاصه‌سازی استخراجی پیشنهادی از دو مرحله پیش‌پردازش و پردازش تشکیل می‌شود. در مرحله پیش‌پردازش، پایان جمله‌ها توسط نرم‌افزار و به‌طور خودکار مشخص می‌شوند، ایست‌واژه‌ها^۲ و کلمه‌های بدون معنی حذف می‌شوند و واژه‌ها به‌صورت زنجیره‌ای در کنار هم قرار می‌گیرند. ارتباط جمله‌ها با موضوع اصلی (عنوان) در مرحله پردازش شناسایی می‌شوند. ارزش جمله‌های متن براساس معیارهایی که در ادامه بیان شده‌اند ارزیابی می‌شوند و سپس، وزنی به هر یک از آنها اختصاص داده می‌شود.

مرحله پردازش شامل استخراج ویژگی‌ها و تولید خلاصه به‌کمک الگوریتم‌هاست. در فرایند استخراج ویژگی، به‌کمک پنج معیار در نظر گرفته شده برای این پژوهش که عبارت‌اند از: شباهت با عنوان، خوانایی، انسجام، ارزش طولی، و ارزش مکانی (در ادامه هر کدام به‌طور کامل معرفی می‌شوند)، خلاصه اولیه تولید می‌شود. سپس به‌کمک دو الگوریتم فرااکتشافی ژنتیک و فاخته، خلاصه‌های نهایی تولید شده و مورد مقایسه (میرشجاعی و معصومی، 2015؛ رحیمی‌راد، 1393؛ پورمعصومی، کاهانی، طوسی، استیری، و قائمی، 1393) قرار می‌گیرند.

¹ Behmadi Moghaddas, Kahani, Toosi, Pourmasoumi, & Estiri

² Stop words

همان‌طور که ذکر شد روش خلاصه‌سازی استخراجی شامل انتخاب جمله‌های مهم و پاراگراف‌ها از سند اصلی و ترکیب آنهاست. اهمیت جمله‌ها براساس ویژگی‌های آماری و زبانی جمله‌ها تعیین می‌شود. مهم‌ترین مزایای این روش عبارت از سادگی، سرعت بالا در فرایند خلاصه‌شدن و در کل، کاهش زمان مطالعه کاربران است. با این حال، این روش معایبی نیز دارد. به‌طور مثال، اگر طول جمله‌ها بیش از حد کوتاه یا بلند باشد یا اطلاعات مرتبط و مهم بین جمله‌های دیگر پخش شده باشند، روش استخراج نمی‌تواند آنها را شناسایی کند. با توجه به خصوصیات زبان فارسی، به‌کارگیری مناسب معیارهای ارزش‌گذاری جملات بسیار مهم است. معیارهای تعریف‌شده و به‌کارگرفته‌شده در این مرحله از پژوهش شامل مجموعه کاملی از معیارهای ارزش‌گذاری جملات در متون فارسی محسوب می‌شود که در تعداد تکرار زیاد به‌صورت تجربی آزمون شدند تا به‌شکل مناسب با خصوصیات جملات فارسی مقاداردهی و منطبق شوند.



شکل 1. معماری کلی روش ارائه‌شده

در این پژوهش به‌جای معیار وزندهی کلاسیک TF-IDF (بایزایتز و ریبرونتو¹، 2011) که مبنای وزندهی براساس محاسبه تکرار کلمات است، از معیار TF-ISF استفاده شده است. در معیار TF-ISF محاسبه وزن در

¹ Baeza-Yates & Ribeiro-Neto

واحد جمله انجام می‌شود (فزونیان و همکاران، 2008) که نحوه محاسبه آن به صورت رابطه‌های 1 و 2 می‌باشد.

$$tf_{i,j} = \frac{freq_{t,j}}{\max(freq_{t,j})} \quad (1)$$

$$isf_i = \log \frac{N}{n_i} \quad (2)$$

در رابطه 1، نشان‌دهنده تعداد تکرار کلمه i در سند j است که به علت کاهش تأثیر اندازه اسناد بزرگتر، نرمال‌سازی می‌شود. $\max(freq_{t,j})$ تعداد تکرار کلمه‌ای است که از سایر کلمات بیشتر تکرار شده است. isf_i معکوس تعداد جمله‌هایی است که شامل کلمه i هستند. N تعداد کل جمله‌های سند و n_i تعداد جمله‌های حاوی کلمه i است. در ادامه، وزن مربوط به کلمه i در سند j مطابق رابطه 3 محاسبه می‌شود.

$$w_{i,j} = tf_{i,j} \times isf_i \quad (3)$$

در گراف وزن‌دهی یک یال میان هر جمله با عنوان متن ایجاد می‌شود. یعنی عنوان به صورت یک پرسش یا درخواست مطرح می‌شود و پاسخ، ارتباط (یال‌های) میان جمله‌ها با عنوان را نشان می‌دهد. بنابراین وزن آن طبق رابطه 4 محاسبه می‌شود (بایزایتز و ریبرونتو، 2011).

$$w_{l,q} = \left(0.5 + \frac{0.5 \times freq_{l,q}}{\max_t freq_{t,q}} \right) isf_l \quad (4)$$

که در آن $freq_{l,q}$ شباهت هر جمله با عنوان تعریف می‌شود و به کمک معیار کسینوسی طبق رابطه 5 به سادگی محاسبه می‌شود (بایزایتز و ریبرونتو، 2011).

$$sim(s_j, q) = \frac{\sum_{i=1}^t w_{l,j} \times w_{l,q}}{\sqrt{\sum_{i=1}^t w_{l,j}^2} \times \sqrt{\sum_{i=1}^t w_{l,q}^2}} \quad (5)$$

به صورت مشابه، شباهت بین دو جمله نیز بر اساس رابطه 6 محاسبه می‌شود.

$$sim(s_m, s_n) = \frac{\sum_{i=1}^t w_{l,m} \times w_{l,n}}{\sqrt{\sum_{i=1}^t w_{l,m}^2} \times \sqrt{\sum_{i=1}^t w_{l,n}^2}} \quad (6)$$

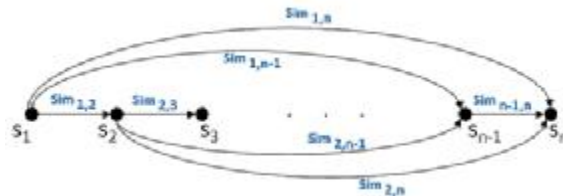
در پژوهش حاضر، برای نمایش اسناد و متون از نوع خاصی از گراف استفاده شده است که به آن «گراف جهت‌دار بدون دور»¹ گفته می‌شود. برای نمایش یک سند با استفاده از این نوع گراف، هر جمله یک گره از گراف را تشکیل می‌دهد. ارتباط گره‌ها به کمک یال‌هایی است که مقدار شباهت گره‌ها را به کمک وزنی که به یال اختصاص داده می‌شود نشان می‌دهد. چگونگی اختصاص وزن به یال‌ها در رابطه 7 نشان داده شده است.

$$\forall (s_i, s_j) \in E, W(s_i, s_j) = sim(s_i, s_j) \quad (7)$$

در واقع، وزن هر یال که دو گره را به هم متصل می‌کند، شباهت میان جمله‌های مرتبط با هم را نشان می‌دهد.

¹ Directed Acyclic Graph (DAG)

همانطور که در شکل 2 دیده می‌شود بین جمله‌هایی که به لحاظ زمانی ترتیب دارند یک یال قرار می‌گیرد. به‌طور مثال، اگر s_i و s_j دو جمله از یک سند باشند، که در آن s_i به لحاظ زمانی قبل از s_j است، نمایش گراف می‌تواند شامل یال (s_i, s_j) باشد ولی نمی‌تواند یال (s_j, s_i) را نمایش دهد (میترا، سینقال، و بوکلی، 1997).



شکل 2. نمایش چگونگی وزندهی به گراف جهت‌دار بدون دور

در ادامه، برای دستیابی به ویژگی‌های توصیفی جمله‌ها از تعدادی معیارهای قابل اندازه‌گیری استفاده کرده‌ایم که برای همه فاکتورهایی که در ادامه تعریف می‌شوند فرض می‌کنیم طول خلاصه ثابت است و این طول ثابت با S نمایش داده می‌شود. تعداد جمله‌ها در سند اصلی نیز N فرض می‌شود.

الف) فاکتور شباهت با عنوان TRF^1

یک متن استخراج شده شامل جمله‌هایی است که درباره عنوان متن اولیه صحبت می‌کند (سیلا، ناسیمتو، پاپا، فریتاز و کاستنر²، 2004). به کمک ماتریس تشابه می‌توان فرمول فاکتور شباهت با عنوان را محاسبه کرد. یک روش ساده برای محاسبه این فاکتور این است که از میانگین شباهت جمله‌ها در خلاصه، تقسیم بر حداکثر مقدار میانگین به عنوان مبنا استفاده کنیم. برای محاسبه TRF ابتدا TR طبق رابطه 8 محاسبه می‌شود:

$$TR_s = \frac{\sum_{s_j \in \text{summary}} \text{stm}(s_j, q)}{S} \quad (8)$$

به کمک TR می‌توان TRF را به صورت رابطه 9 محاسبه کرد:

$$TRF_s = \frac{TR}{\max_{\text{summary}}(TR)} \quad (9)$$

که مقدار ماکسیمم در فرمول بالا از میان همه خلاصه‌های ممکن با طول S به دست می‌آید. برای یافتن \max باید میانگین بیشترین مقدارهای S که برای تشابه جمله‌ها با عنوان به دست می‌آید را محاسبه کرد.

ب) فاکتور انسجام³ (CF)

یک خلاصه با کیفیت شامل جمله‌هایی است که همدیگر را همراهی می‌کنند. با توجه به این فرض، نیازمند محاسبه شباهت جمله‌ها به صورت دو به دو هستیم. این شباهت به کمک ماتریس شباهت محاسبه می‌شود که در

¹ Topic Relation factor

² Silla, Nascimento, Pappa, Freitas, & Kaestner

³ Cohesion factor

بخش وزن‌دهی گراف توضیح داده شد (میترا و همکاران، 1997). دو فرض اولیه درباره گراف نمایش جمله‌ها بیان می‌شود:

1. شباهت گره با خودش صفر است.

2. شباهت برای گره‌هایی تعریف می‌شود که از نظر زمانی پشت سر هم قرار گرفته باشند.

برای محاسبه فاکتور انسجام نیازمند میانگین وزن‌های اختصاص داده‌شده به یال‌ها هستیم. در این صورت، مقدار ایده‌آل برای فاکتور انسجام با تقسیم میانگین همه وزن‌های اختصاص داده‌شده به یال‌ها بر حداکثر مقدار میانگین برای همه خلاصه‌های ممکن، به دست می‌آید.

بنابراین برای محاسبه CF، ابتدا یک زیرگراف از گراف اصلی سیستم به دست می‌آوریم. سپس C_S که میانگین شباهت میان همه جمله‌های خلاصه است، به صورت رابطه 10 محاسبه می‌شود:

$$C_S = \frac{\sum_{s_i, s_j \in \text{summary}} W(s_i, s_j)}{N_S} \quad (10)$$

N_S در رابطه بالا تعداد کل یال‌ها در زیرگراف تشکیل شده است. CF باید نشان‌دهنده میزان نزدیکی همه جمله‌ها با هم باشد که به صورت رابطه 11 محاسبه می‌شود (میرشجاعی و معصومی، 2015):

$$CF_S = \frac{\lg(C_S \times 9 + 1)}{\lg(M \times 9 + 1)} \quad (11)$$

M در رابطه بالا نشان‌دهنده بیشترین وزن جمله‌ها در گراف یا همان حداکثر شباهت میان جمله‌هاست. این رابطه براساس تنظیمات تجربی به دست آمده است. در واقع، تابع لگاریتم در مواقعی که میانگین بسیار کوچک‌تر از ماکسیمم است از مقداردهی بسیار پایین به CF جلوگیری می‌کند. اگر بیشتر جمله‌های استخراج‌شده درباره موضوع واحدی صحبت کنند، CF افزایش می‌یابد و از طرف دیگر، اگر جمله‌ها مربوط به موضوعاتی دور از هم باشند، CF به صفر متمایل می‌شود.

ج) فاکتور خوانایی¹ (RF)

دستیابی به مفهوم خوانایی در سیستم استخراج کاری دشوار است. یک سند خوانا، سندی است که جملاتش ارتباط زیادی با جمله‌های بعد از خود داشته باشند. جمله اول با جمله دوم ارتباط و شباهت زیادی دارد، همین‌طور برای جمله‌های دوم و سوم و الی آخر. درحقیقت، یک خلاصه خوانا، همان‌طور که در ابتدا تعریف شد از جمله‌هایی تشکیل می‌شود که یک زنجیره واضح و مرتبط از جمله‌ها را تشکیل می‌دهد (میرشجاعی و معصومی، 2015). محاسبه خوانایی خلاصه s با طول S، که با R_S نمایش داده می‌شود از رابطه 12 به دست می‌آید.

$$R_S = \sum_{0 \leq t \leq S} W(s_t + s_{t+1}) \quad (12)$$

¹ Readability factor

بنابراین، فاکتور خوانایی برای خلاصه S، به صورت رابطه 13 محاسبه می‌شود.

$$RF_S = \frac{R_S}{\max_{\forall i} R_i} \quad (13)$$

دوباره این فرض یادآوری می‌شود که طول جمله‌های خلاصه ثابت در نظر گرفته شده است. زمانی که هدف یافتن خواناترین خلاصه است، یافتن خلاصه‌ای با طول S، برابر است با یافتن مسیری با طول S، که حداکثر وزن را در گراف سند دارد.

(د) فاکتور طول جمله¹ (LF)

جمله‌های بسیار طولانی برای حضور در خلاصه مناسب نیستند زیرا خلاصه را طولانی می‌کنند. از سوی دیگر، جمله‌های خیلی کوتاه هم اغلب بار معنایی زیادی ندارند. جمله‌هایی مناسب‌تر هستند که طول (تعداد کلمه‌ها) آنها به میانگین طول جمله‌های متن نزدیک‌تر باشد. بر این اساس LF برای یک جمله i می‌تواند به صورت رابطه 14 محاسبه شود (شاگری، تقویان، و بهبودی، 1390).

$$L_i = W_{avg} - |W_{avg} - W_i| \quad (14)$$

که در آن W_{avg} میانگین طول جمله‌های متن و W_i طول جمله i است.

برای محاسبه مقادیر ویژگی طولی جمله از رابطه 15 استفاده شده است. این تعریف با ترکیب مفاهیم معیار طول جمله و مفاهیم نرمال‌سازی حاصل شده است.

$$LF_i = \frac{1}{\sqrt{2\pi}\sigma} e^{\frac{-(L_i - \mu)^2}{2\sigma^2}} \quad (15)$$

در این رابطه LF ویژگی طولی است که هر چه این مقدار به حد متوسط نزدیک‌تر باشد مقدار آن بیشتر است (پورمعصومی و همکاران، 1393). σ انحراف معیار مقادیر طول جمله‌ها و μ هم میانگین مقادیر طول جمله‌هاست.

(ه) فاکتور مکان جمله² (PF)

یک جمله هر چه به ابتدای متن نزدیک‌تر باشد، مهم‌تر و برای حضور در خلاصه مناسب‌تر است. به خصوص در بسیاری از کاربردها باید به جمله اول شانس ممتازی برای حضور در خلاصه داده شود (شاگری و همکاران، 1390). بنابراین امتیاز PF برای جمله i به صورت رابطه 16 در نظر گرفته شده است:

$$PF = \begin{cases} N + \alpha & \text{اگر جمله اول است} \\ N - P_i & \text{در غیر این صورت} \end{cases} \quad (16)$$

¹ Length factor

² Position factor

که در آن N طول متن اصلی (تعداد کل جمله‌ها) و P_i موقعیت جمله i و α امتیاز فوق‌العاده برای جمله اول است. در نهایت، برای محاسبه ارزش مکانی در سیستم پیشنهادی از رابطه 17 استفاده می‌شود که به هر جمله از سند، عددی به‌عنوان ارزش مکانی آن جمله اختصاص می‌دهد.

$$PF_i = \frac{1}{1 + \frac{1}{\sqrt{2\pi}\sigma} + e^{-\frac{(p_i - \mu)^2}{2\sigma^2}}} \quad (17)$$

در این رابطه موقعیت جمله در متن با p_i نشان داده شده و در نهایت عدد حاصل برابر با مقدار ارزش مکانی جمله در سند است. درباره معیار مکان جمله، با بررسی های انجام‌شده بر خلاصه‌های انسانی نیز مشاهده شده است که جمله‌های ابتدایی و انتهایی متون خبری از ارزش بیشتری نسبت به سایر جمله‌ها برخوردار هستند.

با توجه به تأثیر استفاده از الگوریتم‌های فرااکتشافی در بهبود خلاصه‌سازی متن، اکثر آنها بر متون غیرفارسی و به‌طور بر کلکسیون‌های تک‌سندی آزمایش شده‌اند (روتیری و بالاباتاری، 2017ب). در استفاده از الگوریتم ژنتیک برای خلاصه‌سازی متن، اگر برداری را در نظر بگیریم، که i مین موجودیت از بردار مقدار یک بگیرد، به این معناست که جمله s_i در سند خلاصه حضور خواهد داشت، و اگر مقدار صفر بگیرد، خلاصه شامل s_i نخواهد بود. فرض کنید تعداد جمله‌های متن خلاصه، S در نظر گرفته شود و $S < N$ ، که N تعداد کل جمله‌ها در سند اصلی است. تعدادی از عناصر بردار، به طول S به‌صورت تصادفی مقدار یک می‌گیرند و $N-S$ عنصر بردار مقدار صفر می‌گیرند. این بردارها، جمعیت اولیه الگوریتم ژنتیک را تشکیل می‌دهند و فرزندان به‌وسیله عملگرهای تقاطع و جهش تولید می‌شوند.

تابع تناسب برای انتخاب بهترین بخش از والدین برای تولید فرزندان جدید به‌کار برده می‌شود. این فرایند ادامه می‌یابد تا زمانی که خلاصه مناسب‌تر برای تولید وجود نداشته باشد. ممکن است خلاصه‌ای نیارمند خوانایی بالا باشد درحالی‌که برای خلاصه‌ای دیگر، ارتباط زیاد با عنوان اهمیت داشته باشد و خوانایی مدنظر نباشد. تابع تناسب، نیازمند یک طراحی مناسب است که با وزن‌دهی معیارهای استفاده‌شده حاصل می‌شود. با توجه به پنج فاکتور معرفی‌شده در بخش قبل، تابع تناسب مطابق رابطه 18 تعریف می‌شود.

$$F = \frac{\alpha \times TRF + \beta \times CF + \gamma \times RF + \delta \times LF + \epsilon \times PF}{\alpha + \beta + \gamma + \delta + \epsilon} \quad (18)$$

α ، β ، γ ، δ و ϵ ضرایبی هستند که براساس کاربرد تعریف می‌شوند. واضح است با انتخاب مقادیر بزرگ برای α نسبت به بقیه ضرایب، در خلاصه نهایی معیار ارتباط با عنوان وزن بیشتری نسبت به معیارهای انسجام، خوانایی و ارزش‌های طولی و مکانی خواهد داشت. از سوی دیگر، اگر ضرایب β و γ مقادیر بالایی داشته باشند، خلاصه منسجم‌تر و خواناتر خواهد بود.

گاهی سند اولیه عنوان ندارد، در چنین شرایطی مقدار α صفر خواهد بود. زمانی که کاربر قصد دارد خلاصه نهایی را به عنوان یک سند مستقل استفاده کند، طبیعتاً معیار انسجام برای خلاصه ارزشمند تلقی می شود. بنابراین ضریب β باید مقدار بیشتری نسبت به ضرایب دیگر داشته باشد. روشی مبتنی بر یادگیری به نام «اعتبارسنجی متقابل»¹ وجود دارد که برای مقداردهی به ضرایب به کار می رود. در این روش اعدادی متفاوت به هر کدام از ضرایب اختصاص داده می شود و هر کدام از حالتها بررسی و بهترین مقدار برای هر یک از ضرایب مشخص می شود. یادآوری این نکته ضروری است که در نهایت تابع تناسب مقداری نرمال شده ارائه می دهد. بنابراین مقادیر اختصاصی به هر یک از ضرایب نیز باید بین صفر و یک باشند.

در عمل ترکیب و جهش، جستجو و انتخاب کرموزومها به صورت تصادفی انجام می شود (کلامی، 2015). در نهایت بعد از اتمام دو عمل ترکیب و جهش از میان جمعیت جدید مجدداً عمل انتخاب² انجام می شود. این عملیات در تعداد تکرارهای زیاد اعمال می شوند تا زمانی که بهترین فرزند انتخاب شود، این همان خروجی الگوریتم ژنتیک است. در روش پیشنهادی از مقادیر $MP=0.3$ ³ و $CP=0.8$ ⁴ و برای عملگرهای جهش و ترکیب استفاده شده است که 2000 بار تکرار شده اند.

برای استفاده از الگوریتم فاخته برای استخراج جمله های مهم، مجموعه ای از پارامترها در ابتدا مقداردهی می شوند و تطبیق پارامترها با سیستم استخراج که پیش تر به طور کامل توضیح داده شد انجام می شود. براساس نحوه کارکرد الگوریتم جستجوی فاخته، ابتدا پارامترهای تعداد پرندگان و تعداد تکرار متناسب با سیستم پیشنهادی مقادیر ثابت می گیرند. سپس تعداد جمله های اولیه، جمله های خلاصه و ماتریس شباهت حاصل از سیستم استخراج به عنوان پارامترهای ورودی الگوریتم جستجوی فاخته در نظر گرفته می شوند. سپس جمله های باارزش بالاتر انتخاب می شوند و به صورت متن خلاصه نمایش داده می شوند.

مراحل کلی اجرای الگوریتم CSA به شرح زیر است:

1. مقداردهی پارامترهای الگوریتم؛
2. اختصاص تصادفی جمله ها به لانه ها (مقداردهی تصادفی nest ها)؛
3. ارزیابی لانه ها به کمک عملکرد تابع هزینه⁵؛
4. بروزرسانی موقعیت لانه ها؛

¹ Cross validation

² Selection

³ Mutation percentage

⁴ Crossover percentage

⁵ Cost function

5. اگر شرط پایان حلقه (یافتن جمله‌های باکیفیت مدنظر) حاصل شود، الگوریتم پایان می‌یابد. در غیر این صورت، به مرحله 3 برمی‌گردد (میرشجاعی و معصومی، 2015).

این الگوریتم را می‌توان به موارد پیچیده‌تر توسعه داد که در آن هر لانه تخم‌های متعددی دارد که مجموعه‌ای از راه‌حل‌ها را نشان می‌دهد. در این پژوهش، برای هر لانه تنها یک تخم در نظر گرفته می‌شود. تعریف تابع هزینه در الگوریتم CSA همان تعریف تابع تناسب در الگوریتم ژنتیک است. در این جا نیز از ترکیب معیارهای ارتباط با عنوان، انسجام، خوانایی، ارزش طولی، و ارزش مکانی برای محاسبه تابع هزینه استفاده می‌شود. با این تفاوت که تابع هزینه در اینجا به وسیله تخم‌ها محاسبه می‌شود.

بروزرسانی لانه‌ها، به کمک معادله‌ای معروف به نام Levy flight انجام می‌شود. این معادله موقعیت $nest$ را براساس تابع هزینه بروز می‌کند و برای اعداد حقیقی تعریف شده است ولی سیستم انتخاب به صورت دودویی است (در سیستم خلاصه‌ساز پیشنهادی، یک جمله برای خلاصه یا انتخاب می‌شود یا نمی‌شود که این موضوع بیان‌کننده حالت دودویی است). Levy flight اساساً یک راه رفتن تصادفی¹ را فراهم می‌کند که طول گام‌ها² نیز به صورت تصادفی است. معادله Levy flight در رابطه 19 نشان داده شده است.

$$X_t^{(t+1)} = X_t^{(t)} + \alpha \otimes Levy(\lambda) \quad (19)$$

که در آن $\alpha > 0$ است و نشان‌دهنده طول گام است که متناسب با مقیاس مسئله تعریف می‌شود. به طور کلی راه رفتن تصادفی یک زنجیره مارکوف است که موقعیت بعدی تنها به مکان فعلی و احتمال انتقال بستگی دارد. \ddot{A} هم به معنای ضرب ورودی است. Levy هم به صورت رابطه 20 محاسبه می‌شود.

$$Levy \sim \mu = t^{-\lambda}, \quad (1 < \lambda \leq 3) \quad (20)$$

که میانگین و واریانس بی‌نهایت دارد. بعضی از راه‌حل‌های جدید که توسط Levy به دست می‌آیند جایگزین راه‌حل‌های موجود می‌شوند و برخی که مقدار تابع هزینه مناسبی ندارند دور انداخته می‌شوند. این عمل، سرعت جستجوی محلی را افزایش می‌دهد. در جدول 2، مقادیر پارامترهای کنترلی دو الگوریتم ژنتیک و جستجوی فاخته نشان داده شده‌اند.

جدول 2. پارامترهای کنترلی دو الگوریتم ژنتیک (GA) و جستجوی فاخته (CSA)

CSA		GA	
۰.۲۵	P_a	۰.۳	MP
۱	α	۰.۸	CP
۱.۵	λ	۰.۱N	S

¹ Random walk

² Step size

به منظور ارزیابی سیستم خلاصه‌ساز پیشنهادی، آزمایشگاه فناوری وب دانشگاه فردوسی مشهد نسخه فارسی ابزار Rouge را در اختیار قرار داده است (پورمعصومی و همکاران، 1393). ابزار Rouge میزان هم‌پوشانی نمونه‌های مشترک خلاصه‌های انسانی و خلاصه ماشینی را محاسبه می‌کند. البته این نمونه براساس معیارهای مختلف می‌تواند تولید شود. امتیازات بالاتر نشان می‌دهد که کیفیت خلاصه تولیدشده بهتر و قابل قبول‌تر است. معیار Rouge از چند معیار مبتنی بر n-gram تشکیل شده است که به کمک آن ارزیابی را انجام می‌دهد. n-gram به معنی تعداد n تایی‌های مشترک بین خلاصه‌های انسانی و خلاصه ماشینی است. نرم‌افزار Rouge شامل هشت معیار مبتنی بر n-gram است که هر یک بر ویژگی خاصی دلالت دارد. در این پژوهش $n=2$ در نظر گرفته شده است. در حال حاضر، جدیدترین سامانه خلاصه‌ساز فارسی در دسترس (برخط)، سامانه خلاصه‌ساز ایجاز¹ است که به همین دلیل از این سامانه برای مقایسه نتایج حاصل از سیستم پیشنهادی استفاده شده است.

برای تولید خلاصه توسط سامانه ایجاز، 20 سند موجود در 50 موضوع مختلف از پیکره چندسندی تبدیل به 50 سند واحد شدند، چرا که در حال حاضر سامانه ایجاز تنها به صورت تک‌سندی کار می‌کند. اگرچه سامانه ایجاز قابلیت خلاصه‌سازی متون با حجم بالا را دارد، اما برای اطمینان بیشتر بعد از یکپارچه‌کردن اسناد پیکره چندسندی، پیش‌پردازش اولیه روی متون ورودی انجام شد تا خلاصه‌های حاصل از کیفیت مناسبی برخوردار باشند. سپس خلاصه‌های برگردانده‌شده توسط سامانه ایجاز ذخیره شدند و مشابه سیستم پیشنهادی به کمک خلاصه‌های انسانی موجود در ابزار ارزیابی Rouge مقایسه شدند.

برای ارزیابی، 15 سند (D91A01 - D91A15) به عنوان نمونه از پیکره پاسخ به عنوان ورودی به سیستم پیشنهادی و سامانه خلاصه‌ساز ایجاز داده شد و در نهایت 10 درصد از مجموع جملات هر سند به عنوان خلاصه نهایی تولید شدند. شایان ذکر است که هر یک از این 15 سند از تجمیع 20 سند حول یک موضوع خاص حاصل شده‌اند که تعداد انتخاب‌شده برای ارزیابی دقت سیستم خلاصه‌سازی براساس پژوهش پورمعصومی و همکاران (1393) مناسب است.

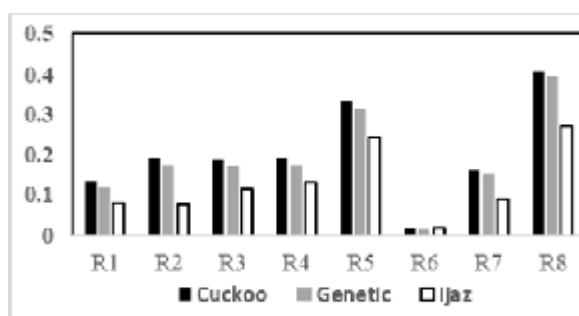
یافته‌ها

جدول 3 مقادیر معیارهای محاسبه‌شده به کمک ابزار Rouge را در متون خلاصه به دست‌آمده نشان می‌دهد که این مقادیر در شکل 3 نیز با یکدیگر مقایسه شده‌اند.

جدول 3. میانگین مقادیر به دست‌آمده از ابزار Rouge در روش پیشنهادی و مقادیر حاصل از سامانه ایجاز در شرایط مشابه

¹ ijaz.um.ac.ir

Cuckoo	Genetic	Ijaz	معیارهای نرم افزار Rouge
۰.۱۳۳۱۷۳	۰.۱۱۵۷۱۳	۰.۰۷۷۵۹	ارزیابی با بررسی n گرام‌های مشابه در کل متن (R1)
۰.۱۹۲۰۹۳	۰.۱۷۰۸۲۷	۰.۰۷۷۳۷	ارزیابی با بررسی n گرام‌های مشابه در جمله‌ها (R2)
۰.۱۸۷۲۴	۰.۱۷۵۴۹۳	۰.۱۱۳۸۸	ارزیابی با بررسی ویژه n گرام‌های مشابه در کل متن (R3)
۰.۱۹۲۰۹۳	۰.۱۷۰۸۲۷	۰.۱۲۷۹۸	ارزیابی با بررسی ویژه n گرام‌های مشابه در جمله‌ها (R4)
۰.۳۲۹۸۶۷	۰.۳۱۱۱۹۳	۰.۲۴۴۳۵	ارزیابی با بررسی طولانی‌ترین زیر رشته مشترک (R5)
۰.۰۱۵۸۱۳	۰.۰۱۵۳۳۷	۰.۰۱۸۶۳	ارزیابی با بررسی طولانی‌ترین زیر رشته مشترک وزن دار (R6)
۰.۱۵۸۸۵۳	۰.۱۴۸۹۸	۰.۰۸۹۳۶	ارزیابی با بررسی ۲ گرام‌های مشابه با فاصله آزاد (R7)
۰.۴۰۶۶۴	۰.۳۹۴۵۳۳	۰.۲۷۲۴۱	ارزیابی با بررسی تعداد لغت‌های مشابه در متن (R8)



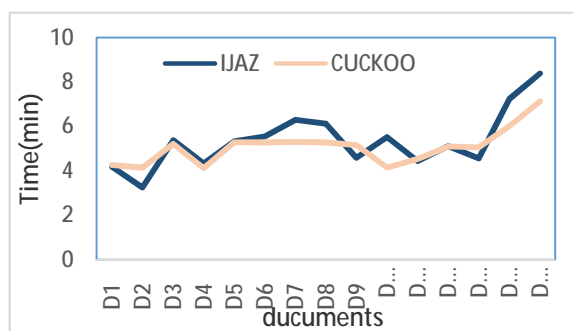
شکل 3. عملکرد الگوریتم‌های سیستم پیشنهادی و مقایسه آنها با سامانه خلاصه‌ساز ایجاز

مطابق شکل 3، عملکرد الگوریتم جستجوی فاخته در هر یک از معیارهای نرم‌افزار Rouge بهتر از الگوریتم ژنتیک است که به معنای شباهت بیشتر خلاصه‌های به‌دست‌آمده به خلاصه‌های انسانی است. ضمناً از میان 8 معیار موجود در نرم‌افزار Rouge، معیار ارزیابی با بررسی طولانی‌ترین زیررشته مشترک (R5) و ارزیابی با بررسی تعداد لغات مشابه در متن (R8) عملکرد بهتری نسبت به بقیه معیارها دارند.

معیارهای ارزیابی Rouge روی متون مختلف عملکردهای متفاوتی از خود نشان می‌دهند. علت نیز ماهیت و ویژگی‌های مختلف اسناد گوناگون است. به‌طور مثال، معیار ارزیابی با بررسی طولانی‌ترین زیررشته مشترک وزن دار (R6)، حالت وزن دار معیار ارزیابی با بررسی طولانی‌ترین زیررشته مشترک (R5) است و به‌نظر می‌رسد باید عملکرد بهتری نسبت به R5 نشان دهد. درحالی‌که این موضوع برای خلاصه‌سازی تک‌سندی صدق می‌کند ولی برای خلاصه‌سازی چندسندی معیار R6 کاهش زیادی داشته و این به‌علت اهمیت داشتن ترتیب کلمات زیر رشته است.

میانگین مقادیر به‌دست‌آمده توسط 8 معیار ارزیابی در نرم‌افزار Rouge توسط الگوریتم فاخته 0.20 و توسط الگوریتم ژنتیک 0.19 است. از طرفی در مقایسه الگوریتم‌های پیشنهادی با سامانه ایجاز، نتایج هر دو الگوریتم حاکی از عملکرد بهتر سیستم پیشنهادی نسبت به سامانه ایجاز است. میانگین مقادیر به‌دست‌آمده توسط

معیارهای ارزیابی در نرم افزار Rouge توسط سامانه ایجاز 0.13 است که تفاوت درخور توجه سامانه پیشنهادی نسبت به سامانه خلاصه ساز ایجاز، نشان دهنده عملکرد مناسب روش پیشنهادی است. ضمناً، نمودار مقایسه زمانی خلاصه ساز پیشنهادی توسط الگوریتم فاخته و خلاصه ساز ایجاز بر 15 متن چندسندی آزمایش شده و در شکل 4 نشان داده شده است. نتایج حاکی از آن است که میانگین زمانی محاسبه شده برای خلاصه سازی توسط سیستم پیشنهادی با الگوریتم فاخته 5.6 درصد کمتر از زمان مورد نیاز توسط سیستم ایجاز است که به طور یقین می توان انتظار داشت در صورت افزایش تعداد اسناد مورد خلاصه سازی، اختلاف زمان خلاصه سازی به شدت درخور توجه خواهد بود و این نیز تأییدی بر کارایی روش پیشنهادی است.



شکل 4. نمودار مقایسه زمانی روش پیشنهادی با خلاصه ساز ایجاز

نتیجه گیری

زمینه های کاربرد خلاصه سازی خودکار متن بسیار گسترده است که می توان به تولید خلاصه اخبار و انتقال آن از طریق سیستم هایی نظیر تلفن همراه، تسریع بازیابی اطلاعات (تولید سیستم های پاسخ گویی)، تولید سیستمی جهت استخراج اسناد مشابه، و کاربرد خلاصه سازی در وب معنایی اشاره کرد.

بیشتر فعالیت های اولیه مربوط به خلاصه سازی خودکار متن بر ساختار ظاهری متن مانند موقعیت جمله و عبارات اشاره، متمرکز بودند. سپس ایده استفاده از هوش مصنوعی بر مبنای استخراج دانش برای شناسایی موجودیت های مفهومی از متن و استخراج روابط بین موجودیت ها با مکانیزم های استخراج مطرح شد. مشکل اصلی کاربردهای هوش مصنوعی این است که الگوهای تعریف شده محدودیت هایی دارند و ممکن است به تحلیل کامل موجودیت های مفهومی منجر نشوند. نگرش های آماری در خلاصه سازی متن نیز ضعف هایی دارند که عبارت اند از: نیاز به دخالت انسان، ارجاعات مبهم، تفسیر محتویات غیرمتنی و مترادف ها. در سال های اخیر بحث خلاصه سازی مبتنی بر کاربر یا خلاصه سازی شخصی سازی شده مطرح است. ایده اصلی خلاصه سازی

شخصی‌سازی شده و یا مبتنی بر کاربر این است که کاربران مختلف با توجه به دانش و پیش‌زمینه اطلاعاتی که دارند، دیدگاه‌های متفاوتی روی اسناد یکسان دارند.

هدف اصلی این پژوهش ارائه روشی برای خلاصه‌سازی چندسندی متون فارسی بود که بخش عمده نوآوری‌های انجام‌شده در این پژوهش مربوط به پیشنهاد الگوریتم‌های فرااکتشافی سازگار برای انتخاب جمله‌های متن خلاصه است. در واقع، الگوریتم‌های فرااکتشافی، یکی از انواع الگوریتم‌های بهینه‌سازی تقریبی هستند که راهکارهای برون‌رفت از نقاط بهینه محلی و قابلیت کاربرد در طیف گسترده‌ای از مسائل را دارند. همین موضوع ایده استفاده از برخی الگوریتم‌های تقریبی را برای تبدیل مسئله خلاصه‌سازی به یک مسئله بهینه‌سازی و جستجوی یافتن بهترین روش برای حل این نوع مسائل بهینه‌سازی را تداعی کرد.

بر اساس نتایج به‌دست‌آمده، الگوریتم جستجوی فاخته خلاصه‌ای باکیفیت‌تر نسبت به الگوریتم ژنتیک فراهم می‌کند. میانگین مقادیر به‌دست‌آمده توسط معیارهای ارزیابی در نرم‌افزار Rouge توسط الگوریتم فاخته 0.20 و توسط الگوریتم ژنتیک 0.19 است. این نتیجه می‌تواند به‌علت این باشد که پارامترهای کمتری در الگوریتم جستجوی فاخته برای تنظیم‌شدن وجود دارند که همین موضوع پیچیدگی الگوریتم را نیز کاهش می‌دهد. درحقیقت، صرف‌نظر از اندازه جمعیت (N)، اساساً یک پارامتر P_a برای تنظیم‌شدن وجود دارد. علاوه بر این، نرخ همگرایی به پارامتر P_a خیلی زیاد نیست. این بدان معناست که مقداردهی مجدد مقادیر این پارامترها به حداقل می‌رسد. همچنین تنظیمات این پارامتر به‌گونه‌ای است که پارامترهای دیگر را تحت تاثیر رفتار خود قرار نمی‌دهد. از طرفی استفاده از معادله Levy flight و راه رفتن تصادفی در کد الگوریتم جستجوی فاخته مانع از گیرافتادن الگوریتم در نقاط بهینه محلی و خاتمه نیمه‌کاره آن می‌شود. این قابلیت می‌تواند علت دیگری برای عملکرد بهتر الگوریتم جستجوی فاخته نسبت به الگوریتم ژنتیک باشد.

همانطور که در شکل 3 مشخص است علاوه بر مقایسه عملکرد دو الگوریتم فراابتکاری ژنتیک و جستجوی فاخته در سیستم پیشنهادی، این سیستم با سامانه خلاصه‌ساز ایجاز به‌عنوان یک مرجع رایج و کاربردی، مقایسه و بررسی شده است. نتایج، نشان‌دهنده برتری خلاصه‌ساز پیشنهادی نسبت به سامانه ایجاز در هر یک از معیارهای ارزیابی است. علت اصلی این بهبود در نتایج را می‌توان استفاده از الگوریتم‌های فراابتکاری عنوان کرد؛ چراکه تفاوت اصلی سیستم پیشنهادی و سامانه ایجاز نیز در کاربرد همین الگوریتم‌هاست. این موضوع به‌خوبی موید موفقیت ایده اولیه سیستم پیشنهادی در این مطالعه است.

طبق پژوهش‌های انجام‌شده در زمینه الگوریتم‌های فرااکتشافی، الگوریتم بهینه‌سازی جستجوی فاخته برای بسیاری از مشکلات بهینه‌سازی عمومی بهتر از سایر الگوریتم‌ها عمل کرده است. این استراتژی بهینه‌سازی می‌تواند به‌طور بالقوه، برنامه‌های بهینه‌سازی چندمنظوره را با محدودیت‌های مختلف و حتی حل مشکلات

سخت NP مطالعه کند. در نهایت، این پژوهش می‌تواند شروعی برای ارائه تکنیک‌های مختلفی برای بهبود روش پیشنهادی محسوب شود. به‌عنوان پیشنهاد، می‌توان ویژگی‌های در نظر گرفته شده برای مرحله پیش‌پردازش را افزایش داد و با دخیل کردن ویژگی‌های معنایی و ادراکی با در نظر گرفتن ویژگی‌های زبان‌شناختی زبان فارسی، کیفیت خلاصه‌های تولید شده را افزایش داد و به‌طور خاص، می‌توان امیدوار بود که استفاده از الگوریتم‌های فراکتشافی ترکیبی موجب بهبود کارایی شوند. ضمناً به اعمال تکنیک‌هایی جهت کاهش زمان اجرای الگوریتم‌های فراکتشافی نیز باید توجه شود.

مآخذ

اخوان، تارا؛ شمس‌فرد، مهرنوش؛ و عرفانی جورابچی، مونا (1387). خلاصه‌ساز تک‌سندی و چندسندی متون فارسی: PARSUMIST. مقاله ارائه شده در چهاردهمین کنفرانس سالانه انجمن کامپیوتر ایران، تهران.

پورمعصومی، آصف؛ کاهانی، محسن؛ طوسی، سیداحمد؛ استیری، احمد، و قائمی، هادی (1393). ایجاز: یک سامانه عملیاتی برای خلاصه‌سازی تک‌سندی متون خبری فارسی. *پردازش‌های علائم و داده‌ها*, 17 (1), 33-48.

رحیمی‌راد، مژگان (1393). بهبود انتخاب ویژگی با الگوریتم‌های تکاملی بهینه‌سازی ازدحام ذرات و ژنتیک برای طبقه‌بندی متن. مقاله ارائه شده در نخستین سمپوزیوم ملی رباتیک و هوش مصنوعی، اهواز.

شاکری، حسین؛ تقویان، فاطمه؛ و بهبودی، فاطمه (1390). یک روش جدید خلاصه‌سازی متن فارسی مبتنی بر ویژگی‌های جملات. مقاله ارائه شده در دومین همایش فناوری اطلاعات، حال، آینده، مشهد.

طالب علی، لیلا؛ ریاحی، نوشین (1394). آبان. مروری بر روش‌های خلاصه‌سازی خودکار متون. مقاله ارائه شده در کنفرانس بین‌المللی پژوهش‌های کاربردی در فناوری اطلاعات، کامپیوتر و مخابرات، تربت حیدریه.

عرب احمدی، فاطمه زهرا (1397). بررسی تاثیر تکنیک‌های خلاصه‌سازی بر روی دسته‌بندی متون فارسی. پایان‌نامه کارشناسی ارشد، دانشگاه گلستان، گرگان.

کریمی، زهره؛ شمس‌فرد، مهرنوش (1385). سیستم خلاصه‌سازی خودکار متون فارسی. مقاله ارائه شده در دوازدهمین کنفرانس سالانه انجمن کامپیوتر، تهران.

مربخ بیات، فرشاد (1393). الگوریتم‌های بهینه‌سازی فراابتکاری. تهران: جهاد دانشگاهی.

مشکی، محسن (1388). خلاصه‌سازی گزینشی چندسندی متون فارسی. پایان‌نامه کارشناسی ارشد، دانشگاه علم و صنعت ایران، تهران.

- Baeza-Yates, R., & Ribeiro-Neto, B. (2011). *Modern Information Retrieval: the concepts and technology behind search*. New York; Toronto: Addison Wesley.
- Behmadi Moghaddas, B., Kahani, M., Toosi, S. A., Pourmasoumi, A., & Estiri, A. (2013). Pasokh: a standard corpus for the evaluation of Persian text summarizers. In *3rd International eConference on Computer and Knowledge Engineering*, October 31 - November 1, (pp. 471-475), IEEE.
- Fattah, M. A., & Ren, F. (2009). GA, MR, FFNN, PNN and GMM based models for automatic text summarization. *Computer Speech & Language*, 23 (1), 126-144.
- Foong, O.-M., & Oxley, A. (2011). A hybrid PSO model in extractive text summarizer. In *IEEE Symposium on Computers & Informatics*, March 20-23, (pp. 130-134). Piscataway, NJ: IEEE.

- Goel, S., Sharma, A., & Bedi, P. (2011). Cuckoo search clustering algorithm: a novel strategy of biomimicry. In World Congress on Information and Communication Technologies, December 11-14, (pp. 916-921). Piscataway: IEEE.
- Gupta, V. (2010). A survey of text summarization extractive techniques. *Journal of Emerging Technologies in web Intelligence*, 2 (3). 259-268.
- Hassel, M., & Mazdak, N. (2004). FarsiSum, a Persian Text Summarizer. *Proceedings of the 20th International Conference on Computational Linguistics*, August August 23-27, (pp. 82-84). East Stroudsburg, PA: Association for Computational Linguistics.
- Hernandez, R., & Ledeneva, Y. (2009). Word Sequence Models for Single Text Summarization. In *Proceedings of the Second International Conferences on Advances in Computer-Human Interactions*, February 1-7, (pp. 44-48), IEEE.
- Honarpisheh, M. A., Ghassem-Sani, G. R., & Mirroshandel, G. (2008). A multi-document multilingual automatic summarization system. In *Proceedings of the Third International Joint Conference on Natural Language Processing*, (pp. 733-738). Retrieved June 17, 2019, from <https://www.aclweb.org/anthology/I08-2101>
- Hovy, E. (2005). Text Summarization. In R. Mitkov (Ed.), *the Oxford Handbook of Computational Linguistics* (pp. 583-598). Oxford: Oxford University Press.
- ISO 215:1986. (1986). Documentation -- Presentation of contributions to periodicals and other serials. Retrieved June 27, 2019, from <https://www.iso.org/standard/4086.html>
- Kalami, S. (2015). Implementation of Binary Genetic Algorithm in MATLAB. Retrieved June 27, 2019, from <https://www.mathworks.com/matlabcentral/mlc-downloads/downloads/submissions/52856/versions/2/previews/YPEA101%20Genetic%20Algorithms/01%20Binary%20Genetic%20Algorithm/Crossover.m/index.html>
- Ledeneva, Y., Gelbukh, A., & Hernández, R. (2008). Terms derived from frequent sequences for extractive text summarization. In A. Gelbukh (Ed.), *Computational Linguistics and Intelligent Text Processing. Proceedings of the 9th international conference on Computational linguistics and intelligent text processing*, February 17-23, (pp. 593-604). Berlin, Heidelberg: Springer-Verlag.
- Martens, D., De Backer, M., & Haesen, R. (2007). Classification with ant colony optimization. *IEEE Transactions on Evolutionary Computation*, 11 (5), 651-665.
- Mirshojaei, H., & Masoomi, B. (2015). Text summarization using cuckoo search optimization algorithm. *Journal of Computer & Robotics*, 8 (2), 19-24.
- Mitra, M., Singhal, A., & Buckley, C., (1997). Automatic text summarization by paragraph extraction. Retrieved June 17, 2019, from <https://www.aclweb.org/anthology/W97-0707>
- Qazvinian, V., Hassanabadi, L. S., & Halavati, R. (2008). Summarising text with a genetic algorithm-based sentence extraction. *International Journal of Knowledge Management Studies*, 2 (4), 426-444.
- Rai, P., & Varshney, A. (2015). Comparative analysis of meta-heuristic algorithms based on their application areas. *International Journal of Innovative Research in Computer and Communication Engineering*, 3 (6), 5982-5988.
- Rautray, R., & Balabantaray, R. C. (2017a). An evolutionary framework for multi document summarization using Cuckoo search approach: MDSCSA. *Applied Computing and Informatics*, 14 (2), 134-144.
- Rautray, R., & Balabantaray, R. C. (2017b). Bio-inspired approaches for extractive document summarization: a comparative study. *Karbala International Journal of Modern Science*, 3 (3), 119-130.
- Silla, J., Nascimento, C., Pappa, G. L., Freitas, A. A., Kaestner, C. A. A. (2004). Automatic text summarization with genetic algorithm-based attribute selection. In C. Lemaître C., C. A. Reyes, & J. A. González (Eds.), *Advances in Artificial Intelligence – IBERAMIA 2004* (vol. 3315, pp. 305-314). Berlin, Heidelberg: Springer.
- Song, W., Choi, L. C., Park, S. C., & Ding, X. F. (2011). Fuzzy evolutionary optimization modeling and its applications to unsupervised categorization and extractive summarization. *Expert Systems with Applications*, 38 (8), 9112-9121.

- Suanmali, L., Salim, N., & Binwahlan, M. S. (2009). Fuzzy Logic Based Method for Improving Text Summarization. *International Journal of Computer Science and Information Security*, 2 (1). Retrieved June 12, 2019, from <https://pdfs.semanticscholar.org/2478/77f2f680fe8f81672c90dcfb9b7d2c94c388.pdf>
- Yang, X. S., & Deb, S. (2009). Cuckoo search via Lévy flights. In *World Congress on Nature & Biologically Inspired Computing*, December 9-11, (pp. 210-214). Retrieved June 17, 2019, from https://www.cs.tufts.edu/comp/150GA/homeworks/hw3/_reading7%20Cuckoo%20search.pdf
- Yang, X. S., & Deb, S. (2014). Cuckoo search: Recent advances and applications. *Neural Computing and Applications*, 24 (1) 169-174.
- Yu, X., & Gen, M. (2010). *Introduction to evolutionary algorithms*. London: Springer.
- Zhang, P., & Li, C., (2009). Automatic text summarization based on sentences clustering and extraction. In the 2nd IEEE International Conference on Computer Science and Information Technology, August 8-11, (pp. 167-170). IEEE

استناد به این مقاله:

آهنگری، فاطمه؛ کرباسی، سهیلا؛ و یعقوبی، مهدی (1398). تکنیک‌های خلاصه‌سازی چندسندی خودکار متون فارسی مبتنی بر الگوریتم‌های فرااکتشافی. *مطالعات ملی کتابداری و سازماندهی اطلاعات*، 30 (2)، 80-58.