



Razi University



Linguistics Society of Iran

Automatic Recognition of Authors Identity in Persian based on Systemic Functional Grammar

Fatemeh Soltanzadeh¹✉, Azadeh Mirzaei², Mohammad Bahrani³, and Shahram Modarres Khiabani⁴

1. Corresponding Author, Ph.D. in Linguistics, Department of Linguistics, Faculty of Persian Literature and Foreign Languages, Allameh Tabataba'i University, Tehran, Iran. E-mail: fatemeh.slt@gmail.com
2. Associate Professor of Linguistics, Department of Linguistics, Faculty of Persian Literature and Foreign Languages, Allameh Tabataba'i University, Tehran, Iran. E-mail: azadeh.mirzaei@atu.ac.ir
3. Assistant Professor, Department of Computer, Faculty of Statistics, Mathematics and Computer, Allameh Tabataba'i University, Tehran, Iran. E-mail: bahrani@atu.ac.ir
4. Assistant Professor, Department of English Language and Translation, Islamic Azad University, Karaj, Iran. E-mail: shmodarress@yahoo.com

Article Info

Article type:

Research Article

Article history:

Received: 19 Jul 2023

Received in revised form: 20 Sep 2023

Accepted: 23 Sep 2023

Published online: 22 Sep 2024

Keywords:

author identification, forensic linguistics, systemic functional grammar, function words, conjunctive adjunct, mood adjunct, comment adjunct.

ABSTRACT

Automated author identification is one of the important fields in forensic linguistics. In this study, the effectiveness of systemic functional grammar (Halliday and Matthiessen, 2014) features in Persian authorship attribution was compared with that of function words. First, a corpus composed of documents written by seven contemporary Iranian authors was collected. Second, a list of function words was extracted from the corpus. Moreover, conjunction, modality and comment adjunct system networks were applied to form a lexicon using linguistics resources. Then, the relative frequency of function words in addition to systemic functional features were calculated in each document. Multilayer perceptron classifier, a type of neural network, was used for learning phase which resulted in a desirable accuracy in evaluation phase. The results of the study showed that using function words method is superior to systemic functional approach alone in Persian author identification, however, simultaneous use of the two methods increases the effectiveness in comparison to each alone.

Cite this article: Soltanzadeh, F., Mirzaei, A., Bahrani, M., & Modarres Khiabani, Sh. (2024). Automatic Recognition of Authors Identity in Persian based on Systemic Functional Grammar. *Research in Western Iranian Languages and Dialects*, 12 (3), 59-84. <http://doi.org/10.22126/jlw.2023.9391.1716> (in Persian).



© The Author(s).

DOI: <https://doi.org/10.22126/jlw.2023.9391.1716>

Publisher: Razi University

Introduction

Recently, automated author identification has become a key focus for forensic linguistics. Author identification involves determining the writer of a text from a set of potential authors. The text in question could be a threatening letter, an email, a literary work, or a scientific article or book. The basis for author identification rests on the idea that different authors may write about the same topic using overlapping, yet distinct, lexico-grammatical units—an issue referred to as idiolect (Coulthard, 2004).

The first significant attempt to identify writing styles was Mendenhall's study of Shakespeare's plays (1887). The play *Henry VIII* is widely recognized as a collaborative work, not solely authored by William Shakespeare. Plechac (2020) investigated the use of accent or stress to identify the contributions of other authors to the play.

In Persian, several studies have been conducted to determine authorship (Farahmandpour et al., 2013; Arefi et al., 2021). These studies utilized repetitive features, such as lexical richness, frequency of syntactic groups, collocations, and the relative frequency of punctuation marks, to detect writing styles. Measuring the frequency of function words is one valid method for author identification. Function words, which have limited meanings, indicate the functional relationships between components of a sentence. Golshaie (2019) and Dabagh (2007) applied the frequency of Persian function words to identify authors. This study aims to compare the efficacy of function word frequency with systemic functional grammar methods in automatically identifying writing styles.

Theoretical Framework

Systemic Functional Grammar (SFG) is a component of the social semiotic approach to language known as systemic functional linguistics (Halliday & Matthiessen, 2014). SFG conceptualizes language as a network of systems, or interrelated sets of options for creating meaning. Since the 1960s, SFG has been applied in various contexts within computational linguistics (Matthiessen & Bateman, 1991; Teich, 1995). In SFG, the clause is considered the fundamental unit of language, and it is analyzed through three perspectives, defined as the ideational, interpersonal, and textual metafunctions. The ideational function is further divided into the experiential and logical aspects (Halliday & Matthiessen, 2014).

This study employs three system networks: conjunction, modality, and comment. These networks correspond to three types of adjuncts: conjunctive adjuncts, mood adjuncts, and comment adjuncts, respectively. In the systemic environment of conjunction, conjunctions function as conjunctive adjuncts within the clause structure. They establish relationships where one segment of text elaborates on, extends, or enhances another segment (Halliday & Matthiessen, 2014).

Modal adjuncts express the speaker's or writer's judgment or attitude toward the content of the message. There are two types of modal adjuncts: (i) mood adjuncts and (ii) comment adjuncts. Mood adjuncts and comment adjuncts are categorized within the modality and comment adjunct system networks, respectively. Modality encompasses intermediate degrees between positive and negative poles, defining the region of uncertainty between 'yes' and 'no.' The modality system allows writers to qualify events or entities in terms of their probability, typicality, obligation, or inclination (Halliday & Matthiessen, 2014). The comment adjunct system provides a means for the writer to comment on the status of a message concerning the textual and interactive context of the discourse (Argamon et al., 2007). Comments can target either the ideational content of the proposition or the interpersonal aspects of the speech function (Halliday & Matthiessen, 2014).

Method

A corpus was compiled from the works of seven contemporary Persian writers: Hoshang Golshiri, Bozorg Alavi, Ahmad Mahmoud, Mahmoud Dolatabadi, Nader Ebrahimi, Jalal Al-e Ahmad, and

Gholamhossein Saedi, totaling 2,069,243 words. From this corpus, a list of 197 function words was extracted using the Term Frequency-Inverse Document Frequency (TF-IDF) technique. Conjunction, modality, and comment adjunct system networks were then used to create a lexicon.

An author identification system was designed using machine learning techniques. The system tokenized the texts, extracted instances of lexical units specified in the lexicon, and computed the relative frequencies of semantic attribute values for each text, resulting in an overall "feature vector" that described each text. This approach was inspired by the method introduced by Argamon et al. (2007). For the learning phase, a multilayer perceptron classifier, a type of neural network, was utilized.

Results

To evaluate the system, the collected corpus was divided into five segments, and a 5-fold cross-validation method was applied. The 5-fold cross-validation demonstrated a satisfactory accuracy when focusing exclusively on function words. The combined use of function words and SFG methods achieved an accuracy of 74.47% for Persian author identification. Subsequent feature selection identified the most effective features for the machine learning phase. The results indicated that the relative frequency of function words outperformed SFG-based attributes in terms of effectiveness.

Discussion and Conclusions

The evaluation phase revealed that the function words-based method outperformed the systemic functional grammar (SFG) approach in identifying authors. However, the simultaneous use of both methods improved effectiveness compared to using either method alone. The superior performance of the function words-based method may be attributed to the high frequency of function words and the author's unconscious control over their use.

Among the SFG-based features, the combination of top features—namely conjunctive, mood, and comment adjuncts—produced higher accuracy than any single system network alone. Additionally, the results from feature selection indicated that features derived from the modality system network were more effective than those from the conjunction and comment adjunct system networks for Persian author identification.

Overall, while the function words-based method proved to be highly effective on its own, integrating it with SFG-based methods provided a more comprehensive approach, enhancing the accuracy of author identification.



پروہشگاہ علوم انسانی و مطالعات فرہنگی
پرتال جامع علوم انسانی

تشخیص خودکار هویت نویسنده متن در زبان فارسی براساس دستور نقش‌گرای نظام‌مند

فاطمه سلطان‌زاده^۱ | آزاده میرزایی^۲ | محمد بحرانی^۳ | شهرام مدرس خیابانی^۴

۱. نویسنده مسئول، دکتری زبان‌شناسی، گروه زبان‌شناسی، دانشکده ادبیات فارسی و زبان‌های خارجه، دانشگاه علامه طباطبایی، تهران، ایران. رایانامه: fatemeh.slt@gmail.com

۲. دانشیار گروه زبان‌شناسی، دانشکده ادبیات فارسی و زبان‌های خارجه، دانشگاه علامه طباطبایی، تهران، ایران. رایانامه: azadeh.mirzaei@atu.ac.ir

۳. استادیار گروه رایانه، دانشکده آمار، ریاضی و رایانه، دانشگاه علامه طباطبایی، تهران، ایران. رایانامه: bahrani@atu.ac.ir

۴. استادیار گروه آموزش مترجمی زبان انگلیسی، واحد کرج، دانشگاه آزاد اسلامی، کرج، ایران. رایانامه: shmodarress@yahoo.com

چکیده

اطلاعات مقاله

تشخیص خودکار هویت نویسنده متن یکی از مسائل مهم زبان‌شناسی حقوقی تلقی می‌شود. در پژوهش حاضر تلاش می‌شود کارایی ویژگی‌های مبتنی بر مفاهیم دستور نقش‌گرای نظام‌مند هالیدی (هالیدی و متیسن، ۲۰۱۴) با کارایی واژه‌های دستوری در تشخیص هویت نویسنده مقایسه شود. به این منظور، در ابتدا، پیکره‌ای از آثار هفت نویسنده معاصر ایرانی گردآوری شد. در مرحله دوم، از واژه‌های دستوری استخراج‌شده از پیکره فهرستی تهیه شد؛ به علاوه، یک مجموعه واژگان براساس شبکه نظام حروف ربط، شبکه نظام افزوده وجه و شبکه نظام افزوده نگرشی با استفاده از منابع زبانی تهیه شد. سپس بسامد نسبی واژه‌های دستوری و ویژگی‌های مبتنی بر دستور نقش‌گرای نظام‌مند در هر متن محاسبه شد. طبقه‌بند پرسپترون چند لایه، نوعی شبکه عصبی، برای مرحله آموزش سامانه به کار گرفته شد و به دقت مطلوبی در مرحله ارزیابی متجر شد. بررسی نتایج ارزیابی سامانه نشان داد که روش محاسبه بسامد واژه‌های دستوری نسبت به روش مبتنی بر دستور نقش‌گرای نظام‌مند در تشخیص هویت نویسنده متون فارسی برتری دارد؛ باوجوداین، هنگامی که ویژگی‌های دستور نقش‌گرای نظام‌مند هالیدی در کنار ویژگی بسامد واژه‌های دستوری به کار روند، کارایی سامانه نسبت به حالتی که تنها از ویژگی بسامد واژه‌های دستوری استفاده شود، ارتقا می‌یابد.

نوع مقاله: پژوهشی

تاریخ دریافت: ۱۴۰۲/۴/۲۸

تاریخ بازنگری: ۱۴۰۲/۶/۲۹

تاریخ پذیرش: ۱۴۰۲/۷/۱

تاریخ انتشار: ۱۴۰۳/۷/۱

کلیدواژه‌ها:

تشخیص هویت نویسنده،

زبان‌شناسی حقوقی،

دستور نقش‌گرای نظام‌مند، واژه‌های

دستوری،

افزوده ربطی،

افزوده وجه،

افزوده نگرشی.

استناد: سلطان‌زاده، فاطمه؛ میرزایی، آزاده؛ بحرانی، محمد؛ مدرس خیابانی، شهرام (۱۴۰۳). تشخیص خودکار هویت نویسنده متن در زبان فارسی براساس دستور نقش‌گرای نظام‌مند. *مطالعات زبان‌ها و گویش‌های غرب ایران*، ۱۲ (۳)، ۸۴-۵۹. <https://doi.org/10.22126/jlw.2023.9391.1716>

۱- مقدمه

در جهان امروز باتوجه به پیشرفت فناوری و استفاده گسترده از شبکه جهانی وب^۱، این امکان برای کاربران فراهم است که در فضای مجازی هویت واقعی خود را پنهان سازند، با هویت جعلی خود را معرفی کنند، رایانامه‌هایی بدون نام نویسنده واقعی بفرستند و حتی به سرقت ادبی یا علمی دست بزنند. از این رو، لازم است سامانه‌هایی طراحی شوند که امنیت فضای مجازی را تأمین کنند، سرقت ادبی و علمی را کشف نمایند و به صورت خودکار هویت واقعی نویسندگان متون را در این فضا شناسایی کنند. به این علت، در دهه‌های اخیر تشخیص خودکار هویت نویسنده متن^۲ به یکی از مسائل مهم در زبان‌شناسی حقوقی^۳ تبدیل شده است.

تشخیص خودکار هویت نویسنده متن بر این فرض استوار است که هر فرد یک گویش فردی^۴ یا سبک خاص و منحصر به فرد در گفتار یا نوشتار دارد. در واقع، هنگامی که مردم از زبان استفاده می‌کنند، ساخت‌های خاصی از واژگان-دستور^۵ را برمی‌گزینند که دیگر افراد ممکن است از آن‌ها کمتر استفاده کنند. در مرحله بعد هر فرد این ساخت‌ها را به شکلی متفاوت از دیگر افراد ترکیب می‌کند تا پیام خود را منتقل کند و این به معنای منحصر به فرد بودن گویش فردی و سبک گفتار یا نوشتار افراد است (کلنارد^۶، ۲۰۰۴).

در زبان فارسی درباره تشخیص هویت نویسنده متن، پژوهش‌هایی مانند فرهمندیور و دیگران (۱۳۹۱) و عارفی و دیگران (۱۴۰۰) انجام شده است. در این پژوهش‌ها بیشتر جنبه‌های فنی و مهندسی کار مدنظر بوده و به جنبه‌های زبانی و نظری کمتر توجه شده است. در پژوهش‌های یادشده عمدتاً از ویژگی‌های تکراری (غنا و واژگانی^۷، بسامد گروه‌های نحوی^۸، باهم‌آیی کلمات^۹، بسامد نسبی علائم نگارشی^{۱۰}) برای طراحی سامانه‌های تشخیص خودکار هویت نویسنده متن استفاده شده است و به طراحی ویژگی‌های زبانی بر مبنای یک نظریه زبانی خاص پرداخته نشده است. در زبان فارسی همچنین دو پژوهش دباغ (۲۰۰۷) و گلشائی (۱۳۹۸) در خصوص تشخیص خودکار هویت نویسنده متن انجام شده است که ناظر بر جنبه‌های زبانی هستند. در این پژوهش‌ها از روش محاسبه بسامد واژه‌های دستوری^{۱۱} در تشخیص هویت نویسنده متن استفاده شده است. محاسبه بسامد واژه‌های دستوری یک روش مهم و معتبر در تشخیص هویت نویسنده است. از آنجاکه نویسنده در به کارگیری واژه‌های دستوری، کنترل خودآگاه ندارد و این واژه‌ها در متن به تعداد زیاد به کار می‌روند، آن‌ها می‌توانند ویژگی‌های منحصر به فردی در تشخیص خودکار نویسنده تلقی شوند (سگارا^{۱۲} و دیگران، ۲۰۱۵؛ گلشائی، ۱۳۹۸).

دباغ (۲۰۰۷) بر اساس واژه‌های دستوری پربسامد و با روش‌های آماری به جداسازی سبک نگارش نظامی گنجوی/شهریار و عبدالحسین زرین کوب/سیمین دانشور پرداخته است. نتایج این پژوهش نشان می‌دهد که به کمک واژه‌های دستوری می‌توان سبک نویسنده‌ها در نظم و نثر فارسی را مشخص کرد. گلشائی (۱۳۹۸) باتکیه بر گویش فردی و با استفاده از واژه‌های دستوری زبان فارسی، تشخیص هویت نویسنده متن را برای پنج نویسنده معاصر بررسی کرده است. نتایج این پژوهش نشان می‌دهد که نویسنده‌های مختلف واژه‌های دستوری را به طور مشابه به کار نمی‌گیرند. در واقع، اگرچه همه نویسندگان، برخی واژه‌های دستوری پربسامد را به کار می‌برند، اولویت نویسندگان در به کارگیری آن‌ها متفاوت است.

پژوهشگران دو پژوهش نام‌برده (دباغ، ۲۰۰۷؛ گلشائی، ۱۳۹۸) اثبات می‌کنند که واژه‌های دستوری در نثر فارسی توانایی جداسازی سبک نویسندگان را دارند؛ اما آن‌ها برای واژه‌های دستوری، دسته‌بندی معنایی^{۱۳} ارائه نمی‌کنند. دباغ (۲۰۰۷) در انتخاب واژه‌های دستوری صرفاً به معیار پربسامد بودن واژه‌های دستوری بسنده کرده است. گلشائی (۱۳۹۸) نیز حداقل کلمات مورد نیاز برای جداسازی

1. World Wide Web
2. author identification
3. forensic linguistics
4. idelect
5. lexicogrammar
6. M. Coulthard
7. lexical richness
8. frequency of syntactic groups
9. collocation
10. relative frequency of punctuation marks
11. function words
12. S. Segarra
13. semantic classification

سبک نویسندگان براساس واژه‌های دستوری را بررسی کرده است. وی توالی‌های یک تا سه‌واژه‌ای (تک‌نگاشتی^۱، دونگاشتی^۲، سه‌نگاشتی^۳) واژه‌های دستوری را در تشخیص هویت نویسنده متن ارزیابی کرده و درنهایت، واژه‌های تک‌نگاشتی یا همان واژه‌های دستوری تک‌واژه‌ای را کاراترین نوع واژه‌های دستوری در نظر گرفته است. علاوه‌براین، در پژوهش وی، گستره مفهوم واژه‌های دستوری، اندکی محدود شده و ضمائر و تمامی افعال (ربطی و غیرربطی) حذف شده است.

باتوجه‌به نکته‌های یادشده در بررسی نقش واژه‌های دستوری زبان فارسی در جداسازی سبک نویسندگان، نیاز است در انتخاب واژه‌های دستوری به مؤلفه‌های معنایی بیشتر توجه شود و برای آن‌ها دسته‌های معنایی ارائه شود. درنهایت، می‌توان انواع مختلف واژه‌های دستوری را آزمون و بهترین ترکیب واژه‌های دستوری را از این بین انتخاب کرد. در پژوهش حاضر تلاش می‌شود که این خلأ، در حوزه تشخیص هویت نویسنده متون فارسی پر شود و برای واژه‌های دستوری، دسته‌بندی معنایی ارائه شود. علاوه‌براین، کارایی سامانه تشخیص هویت نویسنده با واژه‌های دستوری، مبنای مقایسه قرار گیرد و یک نظریه معتبر معنایی در تشخیص سبک نویسنده زبان فارسی آزموده شود. همچنین به نظر می‌رسد طراحی ویژگی‌های زبانی جدید که مبتنی بر یک نظریه زبانی معتبر باشد، در کنار واژه‌های دستوری، در ارتقای کارایی سامانه‌های تشخیص هویت نویسنده مؤثر باشد. بنابراین، تشخیص خودکار نویسنده متن برای زبان فارسی، مبتنی بر یک نظریه زبانی معتبر با استفاده از ویژگی‌های زبانی جدید، یک ضرورت برای زبان فارسی در جهان پیچیده امروز محسوب می‌شود. در این زمینه، در پژوهش حاضر تلاش بر آن است تا به کمک ویژگی‌های زبانی مبتنی بر دستور نقش‌گرای نظام‌مند^۴ هالیدی^۵ (هالیدی و متیسن^۶، ۲۰۱۴)، هویت نویسنده متن به صورت خودکار تشخیص داده شود و تأثیر ویژگی‌های دستور نقش‌گرا در مقایسه با ویژگی بسامد واژه‌های دستوری سنجیده شود.

پرسش‌های پژوهش حاضر عبارت‌اند از: ۱. کدام‌یک از ویژگی‌های زبانی مبتنی بر دستور نقش‌گرای نظام‌مند هالیدی (هالیدی و متیسن، ۲۰۱۴) بر تشخیص سبک نویسندگان فارسی زبان تأثیرگذار است؟ ۲. اگر ویژگی‌های زبانی مبتنی بر دستور نقش‌گرای نظام‌مند هالیدی (هالیدی و متیسن، ۲۰۱۴) در کنار ویژگی بسامد واژه‌های دستوری به کار روند، بر کارایی سامانه نسبت به حالتی که تنها از ویژگی بسامد واژه‌های دستوری استفاده شود، چه تأثیری دارد؟

برای تشخیص خودکار نویسنده متن از ویژگی‌های زبانی در سطوح مختلف همچون ویژگی‌های آوایی^۷، واژگانی، نحوی، معنایی و ترکیب آن‌ها استفاده می‌شود. پس از تعریف ویژگی‌ها، اسناد موجود در پیکره که هر کدام به یک نویسنده خاص در یک مجموعه بسته تعلق دارند، با بهره‌گیری از ویژگی‌های زبانی تعریف‌شده به روش یادگیری ماشین دسته‌بندی و اسناد مشابه به یک نویسنده واحد منسوب می‌گردد.

نخستین تلاش برای سنجیدن سبک نگارش به سده نوزدهم میلادی و بررسی نمایشنامه‌های شکسپیر^۸ باز می‌گردد (مندنهل^۹، ۱۸۸۷). این پرسش که آیا شکسپیر تمام آثارش را خود به تنهایی نوشته است، سال‌ها برای پژوهشگران بحث‌برانگیز بوده است. برای مثال، برخی منتقدان ادبی معتقد بودند که شکسپیر نمایشنامه هنری هشتم^{۱۰} را به کمک فرد دیگری نوشته است. پژوهشگران رایانه و هوش مصنوعی این ادعا را سال‌ها بررسی کرده و درنهایت، به این نتیجه رسیده‌اند که وی با همکاری فردی به نام فلچر^{۱۱} نمایشنامه هنری هشتم را نوشته است (پلچاک^{۱۲}، ۲۰۲۱). این کار باتکیه بر ویژگی‌های آوایی و الگوهای غالب آهنگین^{۱۳} (توزیع هجاهای تکیه‌دار و بدون تکیه در یک سطر) انجام شده است.

از شمارش کلمات پرتکرار، میانگین طول واژه‌ها در متن و تعداد تکرار کلمات منحصر به فرد و غنای واژگانی نیز به‌منزله ویژگی‌های

1. unigram
2. bigram
3. trigram
4. systemic functional grammar (SFG)
5. M. A. K. Halliday
6. C. M. M. Matthiessen
7. phonetic
8. W. Shakespeare
9. T. C. Mandenhal
10. Henry VIII
11. J. Fletcher
12. P. Plecháč
13. rhythmic

واژگانی برای تشخیص خودکار نویسنده استفاده می‌شود (استاماتاتوس^۱، ۲۰۰۹). در برخی پژوهش‌ها از ترکیب چند نوع ویژگی بهره گرفته شده است. برای مثال، در پژوهش ویراسینگه^۲ و دیگران (۲۰۲۱)، ویژگی بسامد واژه‌های دستوری، بسامد گروه‌های نحوی و غنای واژگانی به کار گرفته شده است. نجفی و تاوان (۲۰۲۲) که از بسامد گروه‌های نحوی، موجودیت‌های نامدار^۳ و باهم‌آیی کلمات در کنار یادگیری عمیق^۴ بهره می‌گیرند به کارایی مطلوبی در تعیین هویت نویسنده در زبان انگلیسی دست می‌یابند. مارتینز^۵ و دیگران (۲۰۲۲) نیز ویژگی بسامد گروه‌های نحوی را با روش شبکه عصبی گرافی^۶ می‌آزمایند و به دقت مناسبی دست می‌یابند.

دو پژوهش آرگامون^۷ و دیگران (۲۰۰۷) و آرگامون و کوپل^۸ (۲۰۱۳) نیز برای تشخیص خودکار نویسنده متن در زبان انگلیسی انجام شده است؛ در این پژوهش‌ها از مجموعه‌ای از ویژگی‌ها که براساس دستور نقش‌گرای نظام‌مند هالیدی (هالیدی و متیسن، ۲۰۱۴) نظام یافته‌اند، بهره گرفته شده است. در پژوهش آرگامون و دیگران (۲۰۰۷)، از دستور نقش‌گرای نظام‌مند در تشخیص هویت، جنسیت^۹، رده سنی و ملیت نویسنده متن استفاده شده و کارایی این دستور در مقایسه با بسامد واژه‌های دستوری سنجیده شده است. در پژوهش آرگامون و کوپل (۲۰۱۳)، از این دستور برای تشخیص جنسیت، رده سنی، ملیت و ویژگی‌های شخصیتی^{۱۰} نویسنده استفاده شده و کارایی این دستور در مقایسه با واژه‌های محتوایی^{۱۱} مقایسه شده است. به این منظور، یک واژگان خاص طراحی شده است که کلمات و عباراتی را دربرمی‌گیرد که از شبکه واژگان و اصطلاح‌نامه‌های^{۱۲} برخط گردآوری شده است. هر مدخل واژگان ویژگی‌هایی دارد که برگرفته از دستور نقش‌گرای نظام‌مند است. نتایج پژوهش آرگامون و دیگران (۲۰۰۷) نشان می‌دهد که این ویژگی‌ها در ترکیب با ویژگی‌های مبتنی بر واژه‌های دستوری بیشترین کارایی را دارند. پژوهش آرگامون و کوپل (۲۰۱۳) نیز بیانگر آن است که در تشخیص جنسیت، رده سنی و ملیت نویسنده، ترکیب واژه‌های محتوایی و ویژگی‌های مبتنی بر دستور نقش‌گرای نظام‌مند بالاترین کارایی را داشته است. همچنین روش دستور نقش‌گرای نظام‌مند در قیاس با مبنای مقایسه خود (واژه‌های محتوایی) کارایی پایین‌تری در تشخیص جنسیت، رده سنی و ملیت نویسنده متن در زبان انگلیسی داشته است.

۲- مبانی نظری

۲-۱ دستور نقش‌گرای نظام‌مند

در پژوهش حاضر، از مفاهیم دستور نقش‌گرای نظام‌مند هالیدی (هالیدی و متیسن، ۲۰۱۴) برای تشخیص خودکار هویت نویسنده متن استفاده می‌شود. این دستور، زبان‌ها را به‌منزله نظامی از انتخاب‌ها برای بیان معنا الگوسازی و همه انتخاب‌های ساختاری و واژگانی را براساس نقش‌های معنایی آن‌ها بازنمایی می‌کند. این نظریه در پردازش زبان طبیعی^{۱۳} و زبان‌شناسی رایانشی از دهه ۱۹۶۰ میلادی به کار گرفته شده است؛ به‌ویژه اینکه در تولید متن^{۱۴} بسیار کاربرد دارد (تیج^{۱۵}، ۱۹۹۵؛ متیسن و بتمن^{۱۶}، ۱۹۹۱). دستور نقش‌گرای نظام‌مند هالیدی مجموعه‌ای از محدودیت‌ها را در بیان معنا تعریف می‌کند. این دستور شبکه‌ای از گزینه‌های محتمل است. این شبکه از انتخاب‌ها، شبکه نظام^{۱۷} نامیده می‌شود. با داشتن این شبکه می‌توان برای توصیف سبک یک متن از ویژگی‌های آماری بهره برد؛ به این شکل که بسامد شرطی هر گره^{۱۸} را در این شبکه با داشتن گره والد^{۱۹} آن محاسبه کرد و آن را معیاری برای توصیف سبک متن به کار

1. E. Stamatos
2. J. Weerasinghe
3. named entities
4. deep learning
5. J. A. Martinez-Galicia
6. graph neural networks
7. S. Argamon
8. M. Koppel
9. gender
10. personality
11. content words
12. thesaurus
13. natural language processing
14. text generation
15. E. Teich
16. J. A. Bateman
17. system network
18. node
19. parent node

برد. این تفسیر بر این اصل استوار است که نویسنده دارای سبک خاص و منحصر به فرد، برای بیان معنا گزینیه‌های خاصی را بر دیگر گزینه‌ها ترجیح می‌دهد (آرگامون و دیگران، ۲۰۰۷). بنابراین، در پژوهش حاضر، از دستور نقش‌گرای نظام‌مند هالیدی برای چارچوب نظری پژوهش و از روش پیشنهادی در پژوهش آرگامون و دیگران (۲۰۰۷) برای تعریف ویژگی‌ها در تشخیص هویت نویسنده استفاده می‌شود.

زبان در دستور نقش‌گرای نظام‌مند به صورت نظامی از معناها در نظر گرفته می‌شود که با صورت^۱ همراه می‌شود. واحد مطالعه در دستور نقش‌گرای نظام‌مند بند^۲ است. از لحاظ معنایی، بند سه لایه مستقل دارد که باهم در تعامل هستند. این لایه‌ها فرانش^۳ نامیده می‌شوند و عبارت‌اند از: ۱. فرانش تجربی^۴ که تجربه نویسنده از هستی را بازنمایی می‌کند؛ ۲. فرانش بینافردی^۵ که به نحوه تعامل میان شرکت‌کنندگان در یک رویداد گفتگویی می‌پردازد؛ ۳. فرانش متنی^۶ که انسجام و پیوستگی میان مطالب را در یک متن سامان می‌دهد (هالیدی و متیسن، ۲۰۱۴).

در پژوهش پیش‌رو، افزوده‌ها^۷ در دو فرانش بینافردی و متنی مدنظر هستند. در این پژوهش، افزوده ربطی^۸ که در فرانش متنی مطرح است و افزوده وجهی^۹ (افزوده وجه^{۱۰} و افزوده نگرشی^{۱۱}) که در فرانش بینافردی به کار می‌رود، برای تشخیص هویت نویسنده استفاده می‌شود. در این زمینه، سه نوع شبکه نظام معرفی می‌شود: شبکه نظام حروف ربط^{۱۲}، شبکه نظام وجه‌نمایی^{۱۳} و شبکه نظام افزوده نگرشی. در ابتدا، درباره فرانش متنی و شبکه نظام حروف ربط بحث می‌شود. سپس به فرانش بینافردی و دو نوع شبکه نظام وجه‌نمایی و شبکه نظام افزوده نگرشی پرداخته می‌شود.

۲-۱-۱ فرانش متنی

در فرانش متنی نحوه چینش سازه‌ها با عنوان ساخت آغازگری-پایان‌بخشی^{۱۴} بررسی می‌شود. آغازگر به منزله نقطه عزیمت پیام، نخستین سازه بند است به شرط آنکه مشارک^{۱۵}، فرایند^{۱۶} یا افزوده حاشیه‌ای باشد (هالیدی و متیسن، ۲۰۱۴). بقیه جمله منهای آغازگر، پایان‌بخش است.

افزوده‌های ربطی (حروف ربط) در فرانش متنی مطرح هستند و بین دو بند یا دو جمله پیوند منطقی برقرار می‌کنند. جایگاه بی‌نشان آن‌ها در ابتدای بند یا پس از آغازگر است. این افزوده‌ها که با انسجام متن مرتبط هستند، در دستور نقش‌گرای نظام‌مند در شبکه نظام حروف ربط سازمان‌بندی می‌شوند. انواع مختلف حروف ربط برای پیوند دادن بندها به کار گرفته می‌شوند. حروف ربط نشان می‌دهند که چگونه یک بند براساس بافت پیشین بسط و توسعه می‌یابد. در این شبکه نظام، افزوده ربطی (حرف ربط) به سه دسته کلی افزوده تشریحی^{۱۷}، گسترشی^{۱۸} و تفصیلی^{۱۹} تقسیم می‌شود (هالیدی و متیسن، ۲۰۱۴) که در شکل (۱) مشاهده می‌شود.

در افزوده تشریحی، افزوده‌ها خود به دو دسته کلی تقسیم می‌شوند: بدلی^{۲۰} و واضح‌سازی^{۲۱}. در افزوده بدلی، یک عنصر دوباره معرفی یا

1. form
2. clause
3. metafunction
4. experiential metafunction
5. interpersonal metafunction
6. textual metafunction
7. adjunct
8. conjunctive adjunct
9. modal adjunct
10. mood adjunct
11. comment adjunct
12. conjunction
13. modality
14. theme-rheme
15. participant
16. process
17. elaborating
18. extending
19. enhancing
20. appositive
21. clarification

بیان می‌شود که به صورت توضیحی^۱ یا به صورت نمونه‌سازی^۲ است. در افزوده واضح‌سازی نیز آن عنصر به شکل واضح‌تری بیان می‌شود.^۳ مثال‌های (۱) و (۲) به ترتیب نمونه‌ای از «افزوده تشریحی، بدلی، توضیحی» و «افزوده تشریحی، بدلی، نمونه‌سازی» هستند.

۱. آنچه به تو می‌گویم یقین و حتمی است مثل این است که خودم از دهن صاحب‌کار شنیده باشم (علوی، ۱۳۸۶).

۲. ما مسلمانیم ارباب؛ مثلاً خود من تا حالا حتی یک بار هم نشده نمازم قضا بشود (گلشیری، ۱۳۷۰).

افزوده گسترشی خود به سه دسته کلی تقسیم می‌شود: ۱. افزوده افزایشی^۴: چیزی را به معنا می‌افزاید؛ ۲. افزوده تغییری^۵: به‌نوعی معنا را تغییر می‌دهد؛ ۳. افزوده تبیینی^۶: برخلاف آن چیزی را بیان می‌کند. در مثال (۳) و (۴) می‌توانید به ترتیب نمونه‌ای از «افزوده گسترشی، افزایشی، مثبت^۷» و «افزوده گسترشی، تبیینی» را مشاهده کنید.

۳. بچه‌ها هورا کشیدند و کف زدند (آل احمد، ۱۳۴۶).

۴. شب‌ها اگر از فرسنگ‌ها راه به جهت تپه‌ماهورهای شبانکاره نگاه می‌کردی، نور نارنجی پایداری می‌دید که بیش‌وکم می‌شد؛ اما مرگ نداشت (ابراهیمی، ۱۳۹۹).

افزوده تفصیلی که پیوستگی^۸ را از طریق مقایسه بین عناصر یا بندها به وجود می‌آورد خود به چهار دسته کلی تقسیم می‌شود: مضمون^۹، شیوه^{۱۰}، فضایی-زمانی^{۱۱}، سببی-شرطی^{۱۲}. مثال‌های (۵) و (۶) نمونه‌هایی از «افزوده تفصیلی، شیوه^{۱۳}، طریق» هستند.

۵. یکی دو مرتبه که مردم ده بیچاره می‌شدند، کدخدا را پیش خان همسایه می‌فرستادند و از او کمک می‌گرفتند و بدین طریق دهکده‌ای به تصرف خانی درمی‌آمد (علوی، ۱۳۹۹).

۶. راه افتاد طرف در و پیش از اینکه خارج شود ایستاد و بعد با صدای محکمی پرسید: «برویانف را چطور کشتی؟» (ساعدی، ۱۳۹۷).

مثال (۷) نمونه‌ای از «افزوده تفصیلی، فضایی-زمانی، ساده، هم‌زمانی» است که به لحاظ توالی زمانی، به دو رویداد هم‌زمان اشاره دارد.

۷. خودت عا‌جز نکن ننه؛ حالا کی می‌خواد بفروشه؟ (محمود، ۱۳۵۳).

افزوده سببی-شرطی نیز از دو نوع کلی سببی و شرطی تشکیل شده است. افزوده سببی خود به دو نوع عام^{۱۴} و مشخص^{۱۵} تقسیم می‌شود. در نوع «مشخص» افزوده به نتیجه^{۱۶}، علت^{۱۷} و قصد^{۱۸} اختصاص می‌یابد. در مثال (۸) و (۹) می‌توانید به ترتیب نمونه‌ای از «افزوده سببی-شرطی، سببی، مشخص، علت» و «افزوده سببی-شرطی، شرطی، منفی^{۱۹}» را مشاهده کنید.

۸. وکیل رعیت وفادارترین رعیت خویش را به هوسی به زنجیر کشید و کشان به بازار آورد و به دار مجازاتی که هیچ حقش نبود آویخت، فقط به‌خاطر آنکه خود را در برابر اجانب کوچک‌تر از او می‌دید (ابراهیمی، ۱۳۹۹).

۹. با این حرف‌ها نمی‌خواهم ناراحت کنم؛ اما تو باید بدانی، باید بفهمی، حتماً کریستین برایت نگفته اگر نه کار به اینجاها نمی‌کشید (گلشیری، ۱۳۵۰).

همان‌گونه که پیش‌تر اشاره شد در شکل (۱) شبکه نظام حروف ربط هالیدی و متیسن (۲۰۱۴) مشاهده می‌شود که سه دسته کلی

1. expository
2. exemplifying

۳. شایان ذکر است که افزوده «واضح‌سازی» نیز خود به هفت زیربخش تقسیم می‌شود؛ به دلیل جزئیات زیاد، در این پژوهش این تقسیم‌بندی مورد توجه نبوده است و در شکل (۱) نمایش داده نشده است. همین مسئله درباره افزوده «مرکب، فضایی-زمانی، تفصیلی» نیز صادق است.

4. additive
5. varying
6. adversative

۷. افزوده افزایشی به دو صورت مثبت و منفی است.

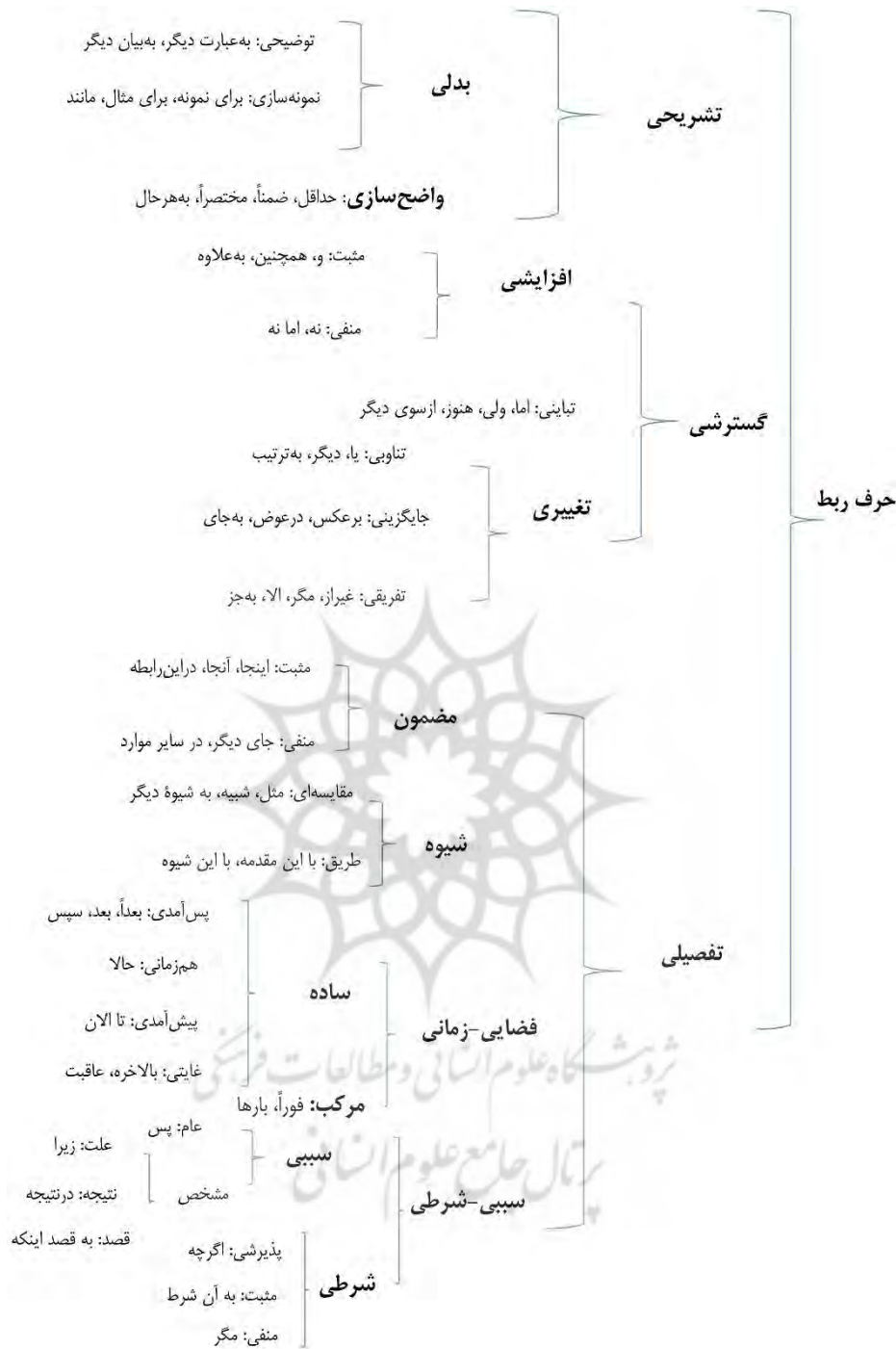
8. cohesion
9. matter
10. manner
11. spatio-temporal
12. causal-conditional

۱۳. افزوده شیوه، پیوستگی را به وسیله مقایسه یا ارجاع بین عناصر ایجاد می‌کند و به دو گروه افزوده مقایسه‌ای و طریق تقسیم می‌شوند.

14. general
15. specific
16. result
17. reason
18. purpose

۱۹. افزوده‌هایی همچون «مگر»، «اگر نه» و «در غیر این صورت» از نوع «افزوده سببی-شرطی، شرطی، منفی» هستند.

افزوده تشریحی، گسترشی و تفصیلی را شامل می‌شود.



شکل (۱). شبکه نظام حروف ربط (هالیدی و متیسن، ۲۰۱۴: ۶۱۲؛ جعفری، ۱۳۸۸: ۱۳۷-۱۳۹)

۲-۱-۲ فرانتش بینافردی

در فرانتش بینافردی هر بند نوعی تبادل است. تعامل زبانی در این فرانتش یک رابطه دادوستد کالا یا خدمات یا اطلاعات است. در واقع، در این فرانتش پیام از دو نوع گزاره^۱ یا پیشنهاد^۲ است. در صورت دادن یا خواستن اطلاعات، پیام از نوع گزاره و در صورت عرضه یا درخواست کالا یا خدمات^۳ پیام از نوع پیشنهاد است (هالیدی و متیسن، ۲۰۱۴). گوینده در فرانتش بینافردی ابزارهایی در دست دارد که

1. proposition
2. proposal
3. goods and services

با آن‌ها می‌تواند نظر خود را در پیام منعکس کند، میزان قطعیت^۱ گزاره را تخمین بزند یا سبب شود که احتمال اجرایی شدن پیشنهاد افزایش یابد. میزان قطعیت هر گزاره با صورت‌بندی‌هایی که نماینده مفهوم احتمال^۲ هستند (حتماً، احتمالاً، شاید، به احتمال زیاد، احتمال می‌رود و...)، همچنین با قیود تکرار (همیشه، اغلب، گاهی و...) مشخص می‌شود. میزان قطعیت یا به بیان بهتر، ضمانت اجرایی پیشنهاد با صورت‌بندی‌هایی منتقل می‌شود که نماینده مفهوم اجبار یا التزام^۳ (اجبار، پیشنهاد و اجازه) و تمایل^۴ (تصمیم، توانایی و تمایل) هستند (میرزایی، ۱۳۹۷).

در این فرانش محتوای یک بند به دو عنصر وجه^۵ و مانده^۶ تقسیم می‌شود. وجه آن بخش از بند است که از دو قسمت اصلی تشکیل شده است: ۱. فاعل^۷ که یک گروه اسمی است؛ ۲. عنصر زمان‌دار^۸ که بخشی از گروه فعلی است و بیانگر زمان، نمود^۹، وجه، جهت^{۱۰} و قطبیت^{۱۱} است. عنصر بعدی در سازمان‌بندی نقش‌ها در معنای بینافردي، مانده است که خود از سه عنصر نقشی مختلف تشکیل شده است: الف) محمول^{۱۲}؛ ب) متمم^{۱۳}؛ ج) افزوده. محمول در تمام بندها با گروه فعلی تحقق می‌یابد. متمم آن قسمت از مانده است که فاعل نیست؛ ولی توانایی فاعل شدن را دارد و معمولاً به صورت گروه اسمی ظاهر می‌شود. افزوده، عنصری از مانده است که توانایی فاعل شدن را ندارد و به طور مشخص به صورت یک گروه قیدی یا حرف اضافه‌ای ظاهر می‌شود (جعفری، ۱۳۸۸: ۱۳۳).

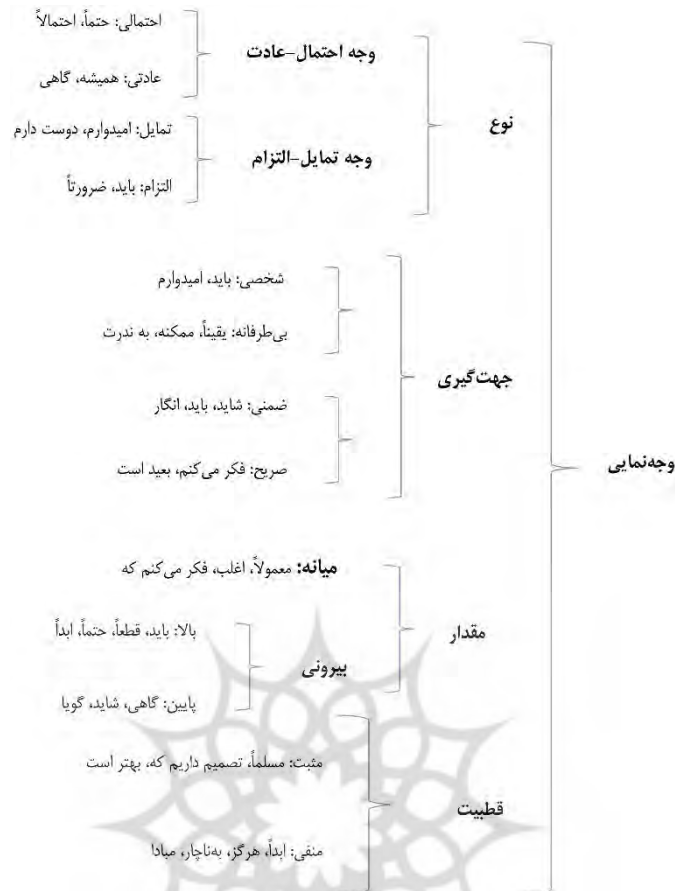
افزوده وجهی در فرانش بینافردي مطرح است و خود به دو دسته تقسیم می‌شود: الف) افزوده‌های وجه که مفاهیمی چون احتمال، اجبار^{۱۴} (التزام)، تمایل^{۱۵}، عادت^{۱۶}، شدت^{۱۷} و زمان را بیان می‌کنند؛ ب) افزوده‌های نگرشی که تفسیر، قضاوت و اظهار نظر در مورد پیام و گزاره‌ای را که بیان می‌شود، ارائه می‌دهد. افزوده‌های وجه در شبکه نظام وجه‌نمایی جای می‌گیرند. آن‌ها نویسنده را قادر می‌سازند که حوادث و موجودیت‌ها را در متن براساس احتمال، عادت، تمایل و اجبار ارزش‌گذاری کند. این شبکه نظام از چهار زیرنظام نوع^{۱۸}، جهت‌گیری^{۱۹}، مقدار^{۲۰} و قطبیت تشکیل شده است که در شکل (۲) سلسله‌مراتب آن نشان داده شده است.

زیرنظام نخست، «نوع» است که از دو نوع کلی وجه‌سازی^{۲۱} (وجه احتمال-عادت) و تعدیل‌سازی^{۲۲} (وجه تمایل-التزام) تشکیل شده است. وجه‌سازی در ارتباط با گزاره موضوعیت می‌یابد. گزاره‌ها یا تأیید می‌شوند یا انکار. نقش نظام وجهیت تحلیل قلمرو عدم قطعیت مابین آری یا خیر است (حسین حمه و دیگران، ۱۴۰۰؛ هالیدی و متیسن، ۲۰۱۴). در وجه احتمال-عادت، از یک سو، احتمال وقوع گزاره مطرح است (قطعاً، احتمالاً، شاید، بعیده، غیرممکنه) و از سوی دیگر، میزان تکرار وقوع گزاره (همیشه، معمولاً، گاهی، هرگز). در مورد اول، وجه احتمال-عادت از نوع احتمالی و دومی از نوع عادت است. تعدیل‌سازی در ارتباط با پیشنهاد مطرح می‌شود. همان‌طور که پیش‌تر اشاره شد، میزان قطعیت پیشنهاد با صورت‌بندی‌هایی که نماینده مفهوم اجبار یا التزام (اجبار، پیشنهاد و اجازه) و تمایل (تصمیم، توانایی و تمایل) هستند، منتقل می‌شود. پس در تعدیل‌سازی نیز دو نوع وجهیت مطرح است: وجه التزامی (مانند باید، ضرورتاً، قهراً و...) و وجه تمایلی (مانند می‌خواهم، دوست داریم، تصمیم داریم و...).

همان‌طور که پیش‌تر گفته شد شبکه افزوده وجه از چهار زیرنظام نوع، جهت‌گیری، مقدار و قطبیت تشکیل شده است که در شکل

1. certainty
2. probability
3. obligation
4. inclination or readiness
5. mood element
6. residue
7. subject
8. finite
9. aspect
10. voice
11. polarity
12. predicator
13. complement
14. obligation
15. desire
16. usuality
17. intensity
18. type
19. orientation
20. value
21. modalization
22. modulation

(۲) می‌توانید سلسله‌مراتب آن را مشاهده کنید.



شکل (۲). شبکه نظام افزوده وجه (هالیدی و متیسن، ۲۰۱۴: ۱۸۲)

دومین زیرنظام، «جهت‌گیری» است. این زیرنظام، دو نوع تقابل را نشان می‌دهد: تقابل شخصی^۱ و بی‌طرفانه^۲ و تقابل صریح^۳ و ضمنی^۴. برای درک بهتر، این دو نوع تقابل در جدول (۱) با مثال نشان داده شده است.

جدول (۱). تحلیل نمونه‌ها براساس زیرنظام جهت‌گیری (هالیدی و متیسن، ۲۰۱۴: ۱۸۱)

بی‌طرفانه	شخصی	جهت‌گیری
قطعاً	باید	ضمنی
مسلم است	مطمئنم	صریح

در مثال‌های «باید» و «مطمئنم» که از نوع شخصی هستند، نویسنده منبع نظر را ذکر می‌کند؛ درحالی‌که در «قطعاً» و «مسلم است»، منبع آن مشخص نیست و حالت بی‌طرفانه دارند. درباره «باید»، ذکر منبع نظر به‌صورت ضمنی و در «مطمئنم» به‌صورت صریح بیان شده است (هالیدی و متیسن، ۲۰۱۴: ۱۸۱).

سومین زیرنظام، «مقدار» است که سه سطح «بالا»، «پایین» و «میانهم» دارد. چهارمین زیرنظام وجه‌نمایی، «قطبیت» است. قطبیت تقابل میان مثبت و منفی است. قطبیت را می‌توان به‌نوعی دو سر پیوستار و جهت در نظر گرفت (میرزایی، ۱۴۰۰).

در ادامه، مثال‌هایی از افزوده‌های وجه در جملات پیکره پژوهش حاضر ذکر شده است. در جدول (۲)، افزوده‌های وجه در این مثال‌ها بنابر شبکه نظام وجه‌نمایی تحلیل شده‌اند.

1. subjective
2. objective
3. explicit
4. implicit

۱۰. فولاد گفت: «یعنی ممکنه که تو تب و هذیان از خونه زده باشه بیرون و تو نخلستانا، گوشه‌ای افتاده باشه و سرما خشکش کرده باشه؟» (محمود، ۱۳۸۱).

۱۱. با این حرف‌ها نمی‌خواهم ناراحت کنم؛ اما تو باید بدانی، باید بفهمی، حتماً کریستین برایت نگفته اگر نه کار به اینجاها نمی‌کشید (گلشیری، ۱۳۵۰).

۱۲. عمه بزرگ گفت: «خسرو خان، از یک شازده بعید است که بادبادک پسر باغبان را بردارد (گلشیری، ۱۴۰۰).

۱۳. این بهت و درنگ گهگاهی از صبح با شیرو بود (دولت‌آبادی، ۱۴۰۱).

۱۴. من تصمیم دارم هر شبه بیایم پیش تو، به شرط اینکه بازهم از این چیزها بنویسی (ساعدی، ۱۳۷۷).

جدول (۲). تحلیل نمونه‌ها بر اساس شبکه نظام وجه‌نمایی

مثال	نوع (سطح اول)	نوع (سطح دوم)	جهت‌گیری	مقدار	قطبیت
ممکنه	وجه‌سازی	احتمالی	بی‌طرفانه	میانه	مثبت
حتماً	وجه‌سازی	احتمالی	بی‌طرفانه	بالا	مثبت
بعید است	وجه‌سازی	احتمالی	بی‌طرفانه	پایین	منفی
گهگاهی	وجه‌سازی	عادی	بی‌طرفانه	پایین	مثبت
نمی‌خواهم	تعدیل‌سازی	تمایل	شخصی	بالا	منفی
باید	تعدیل‌سازی	التزام	شخصی	بالا	مثبت
تصمیم دارم	تعدیل‌سازی	تمایل	شخصی	بالا	مثبت

شبکه نظام دیگر، شبکه افزوده نگرشی است. این افزوده‌ها نظر نویسنده یا گوینده را بیان می‌کنند و هدف آن‌ها محتوای گزاره یا بخشی از نقش گفتار است؛ از این رو، به دو گروه گزاره‌ای^۱ و نقش گفتار^۲ تقسیم می‌شوند. در افزوده نگرشی از نوع گزاره‌ای، گوینده یا نویسنده نظر خود را درباره کل گزاره یا بخشی از آن (فاعل) بیان می‌کند. پس افزوده نگرشی، گزاره‌ای به دو نوع بندی و فاعلی تقسیم می‌شود. افزوده نگرشی از نوع نقش گفتار نیز به دو نوع فاقد شرایط^۳ و واجد شرایط^۴ تقسیم می‌شود (هالیدی و متیسن، ۲۰۱۴: ۱۹۰). در شکل (۳)، شبکه نظام افزوده نگرشی مشاهده می‌شود.



شکل (۳). شبکه نظام افزوده نگرشی (هالیدی و متیسن، ۲۰۱۴)

در شکل (۴) و (۵)، افزوده بندی و فاعلی و در شکل (۶) افزوده نقش گفتار نمایش داده شده است. در مثال (۱۵) و (۱۶) می‌توانید به ترتیب نمونه‌ای از «افزوده گزاره‌ای، بندی، صلاحیتی^۵، فرض، شایعه^۶» و «افزوده گزاره‌ای، بندی،

1. propositional
2. speechfunctional
3. unqualified
4. qualified
5. qualificative
6. hearsay

مطلوبیت، غیرمطلوب را مشاهده کنید.

۱۵. بر فرض هم که خواهرم با این عنکبوت مشغولیتی پیدا کرده باشد، تازه به من چه؟ عنکبوت، عنکبوت است دیگر (آل احمد، ۱۳۵۰).

۱۶. اگر خود شما دختر داشتید، به همچو آدمی شوهرش می‌دادید؟ گفت: «متأسفانه من دختر ندارم» (آل احمد، ۱۳۵۰).

شکل (۴) شبکه نظام افزوده گزاره‌ای بندی هالیدی و متیسن (۲۰۱۴) را نشان می‌دهد. در این شکل می‌توانید سلسله‌مراتب افزوده گزاره‌ای بندی به کاررفته در مثال‌های (۱۵) و (۱۶) را مشاهده کنید.



شکل (۴). شبکه نظام افزوده گزاره‌ای بندی (هالیدی و متیسن، ۲۰۱۴: ۱۹۰)

در مثال (۱۷)، نمونه‌ای از «افزوده گزاره‌ای، فاعلی، اخلاقی، منفی» ذکر شده است.

۱۷. یک کتابخانه نفیس تخصصی در باب معماری داشت؛ اما حتی اشتباهاً، یک دیوان حافظ و مولوی در کنار آن مجموعه نفیس، جای نگرفته بود (ابراهیمی، ۱۳۷۴).

شکل (۵) شبکه نظام افزوده گزاره‌ای فاعلی هالیدی و متیسن (۲۰۱۴) را نشان می‌دهد. در این شکل می‌توانید سلسله‌مراتب افزوده گزاره‌ای فاعلی به کاررفته در مثال (۱۷) را مشاهده کنید.



شکل (۵). شبکه نظام افزوده گزاره‌ای فاعلی (هالیدی و متیسن، ۲۰۱۴: ۱۹۰)

مثال‌های (۱۸) و (۱۹) به ترتیب نمونه‌هایی از «افزوده نقش گفتار، فاقد شرایط، حقیقی» و «افزوده نقش گفتار، واجد شرایط، اعتبار، عام» هستند.

۱۸. شاید حقیقتاً هم اسد نمی‌ترسید و حس ترس سامون بود که او را وامی‌داشت بپندارد برادرش هم می‌ترسد (دولت‌آبادی، ۱۳۹۵).

۱۹. به‌طور کلی، تو کم‌وکسر نداری؟ (ساعدی، ۱۳۷۷).

شکل (۶) شبکه نظام افزوده نقش گفتار هالیدی و متیسن (۲۰۱۴) است.



شکل (۶). شبکه نظام افزوده نقش گفتار (هالیدی و متیسن، ۲۰۱۴: ۱۹۱؛ جعفری، ۱۳۸۸: ۱۴۲)

۲-۲ واژه‌های دستوری

بنابر تعریف ردفورد^۲ (۲۰۰۴)، واژه‌های دستوری، اطلاعاتی درباره نقش‌های دستوری متناسب به برخی مفاهیم (مانند شخص، شمار، جنس و حالت) دارند و معنای قائم به ذات ندارند. با این توضیح، حرف تعریف^۳، کمی‌نما^۴، ضمیر، فعل کمکی، علامت مصدر و متمم‌نما^۵ جزء مقوله‌های دستوری به حساب می‌آیند. در این پژوهش، در تعریف واژه‌های دستوری دو ملاک در نظر گرفته شده است. اگر واژه موردنظر مسئول انتقال مفهومی دستوری باشد یا اگر خودش به‌صورت مستقل معنی نداشته باشد، آنگاه آن واژه، دستوری است. به این ترتیب، حروف اضافه که حوزه معنایی‌شان در هم‌نشینی با واژه‌های دیگر مشخص می‌شود، دستوری محسوب می‌شوند (میرزایی و صفری، ۱۳۹۴: ۲۵۴-۲۵۵). در مجموع، در پژوهش حاضر، حرف اضافه، کمی‌نما، ضمیر، فعل کمکی، همکرد فعل (فعل سبک^۶)، متمم‌نمای «که»، حرف تعریف «یه» و نقش‌نمای مفعولی «را» کلمات دستوری به شمار می‌آیند.

۳- پیکره پژوهش

در پژوهش حاضر، به‌منظور طراحی سامانه‌ای برای تشخیص خودکار هویت نویسنده متون فارسی، پیکره‌ای از آثار هفت نویسنده معاصر ایرانی به نام‌های هوشنگ گلشیری، بزرگ علوی، احمد محمود، محمود دولت‌آبادی، نادر ابراهیمی، جلال آل‌احمد و غلامحسین ساعدی گردآوری و هنجارسازی^۷ شد. منبع به‌کاررفته برای گردآوری این داده، پایگاه داده‌های زبان فارسی (عاصی، ۱۹۹۷) و پیکره

1. validity
2. A. Radford
3. article
4. quantifier
5. complementizer
6. light verb
7. normalization

متنی زبان فارسی (بی‌جن‌خان و دیگران، ۲۰۱۱) بوده است. علاوه‌براین، آثار مختلف هفت نویسنده در شبکه جهانی وب جست‌وجو شد و کتاب‌های آن‌ها که در قالب پی‌دی‌اف^۱ بودند، با نرم افزار ایبو^۲ به قالب متنی^۳ درآمدند. درنهایت، کل پیکره هنجارسازی و یکدست شد. در جدول (۳)، آماره‌های این پیکره نشان داده شده است. شایان ذکر است از آنجاکه در بیشتر این آثار از زبان «محویره» استفاده شده است و ابزارهای پیش‌پردازشگر^۴ زبان فارسی درباره زبان محاوره کارایی نامطلوبی دارند یا دسترسی به آن‌ها محدود است، پیکره به‌صورت خام و فاقد برچسب در نظر گرفته شده است. همچنین روش پیشنهادی در این پژوهش مستقل از ابزارهای پیش‌پردازشگر است.

جدول (۳). آماره‌های پیکره داستانی

نام نویسنده	تعداد آثار	مجموع کل واژه‌ها
نادر ابراهیمی	۸	۲۶۱۸۶۶
احمد محمود	۹	۳۶۱۵۶۵
جلال آل احمد	۹	۲۵۴۲۷۶
بزرگ علوی	۷	۳۲۶۸۰۲
هوشنگ گلشیری	۱۲	۲۴۷۷۹۱
محمود دولت آبادی	۸	۳۴۶۱۳۲
غلامحسین ساعدی	۶	۳۲۳۳۵۸
مجموع	۵۹	۲۱۲۱۷۹۰

۴- روش پژوهش

در پژوهش پیش‌رو، برای تشخیص خودکار هویت نویسنده متن از رویکرد نقش‌گرایی نظام‌مند هالیدی (هالیدی و متیسن، ۲۰۱۴) و روش پیشنهادی در پژوهش آرگامون و دیگران (۲۰۰۷) استفاده شده است. به این منظور، ابتدا فهرستی شامل واژه‌های دستوری، یک مجموعه واژگان برای افزوده‌های ربطی، وجه و نگرشی براساس رویکرد نقش‌گرایی نظام‌مند هالیدی و پیکره‌ای از آثار هفت نویسنده معاصر گردآوری شد (ارجاع به بخش ۳). سپس الگوریتمی^۵ طراحی شد که کلماتی از متون پیکره را که در فهرست واژه‌های دستوری یا مجموعه واژگان تعریف شده‌اند، استخراج و بسامد نسبی شاخص‌های تعریف‌شده برای هر متن را محاسبه می‌کرد. بسامد نسبی شاخص‌ها در این پژوهش شامل بسامد نسبی هریک از واژه‌های دستوری و بسامد نسبی هر گره در شبکه نظام‌های مبتنی بر دستور نقش‌گرایی هالیدی (هالیدی و متیسن، ۲۰۱۴) است. در مرحله بعد، مقدار هریک از شاخص‌ها در هر متن به‌منزله معیاری برای تمایز گذاشتن میان متون مختلف در نظر گرفته شد. به این شکل که الگوریتم یادگیری ماشینی، بردار ویژگی‌ها^۶ را که شامل هریک از شاخص‌های تعریف‌شده برای متون است، دریافت می‌کرد و باتوجه به تفاوت میان مقادیر این شاخص‌ها بین نویسندگان تمایز قائل می‌شد. در ادامه، همه این مراحل به تفصیل شرح داده می‌شود.

۱-۴ استخراج کلمات از متن

واژگان، کلمات و عباراتی را شامل می‌شود که به‌صورت نیمه‌خودکار^۷ از پیکره یا شبکه واژگان گردآوری شده است. هر مدخل واژگان شماری از ویژگی‌های معنایی براساس دستور نقش‌گرایی نظام‌مند را دربردارد. برای گردآوری ویژگی‌های معنایی واژگان از سه منبع پیکره گفتمانی زبان فارسی (میرزایی و صفری، ۲۰۱۸)، فارس‌نت (شمس‌فرد و دیگران، ۲۰۱۰) و پیکره گردآوری‌شده در پژوهشی در حوزه دستور نقش‌گرایی نظام‌مند (جعفری، ۱۳۸۸) بهره گرفته شده است. پیکره گفتمان زبان فارسی جزء معدود پیکره‌های جهان است که در سطح گفتمان به‌صورت دستی برچسب‌گذاری شده است. پیکره گفتمان زبان فارسی در دو سطح درون‌جمله‌ای (بینابندی) و درون‌متنی (بیناجمله‌ای) تهیه شده است. کلیات شیوه برچسب‌گذاری این پیکره که در مجموع حجمی بالغ بر یک میلیون کلمه دارد، براساس

1. PDF
2. <https://www.eboo.ir/>
3. text (txt)
4. preprocessor
5. algorithm
6. feature vector

۷. منظور از نیمه‌خودکار این است که در ابتدا فهرست کلمات به‌صورت خودکار گردآوری می‌شود؛ سپس زبان‌شناس به‌صورت دستی آن را دسته‌بندی و برچسب‌گذاری می‌کند.

استاندارد درخت‌بانک گفتمانی پن^۱ است (پراساد^۲ و دیگران، ۲۰۰۸). این شیوه‌نامه به‌لحاظ نظری در بحث روابط منطقی از رویکرد نقش‌گرایی هالیدی و متیسن (۲۰۱۴) در فرانش منطقی^۳ پیروی می‌کند. علاوه‌براین، در پیکره گفتمان برای سه مفهوم وجه، قطبیت و قطعیت برچسب‌گذاری صورت گرفته است (شمس‌فرد و بی‌جن‌خان، ۱۴۰۱: ۷۳). فارس‌نت نیز بزرگ‌ترین شبکه واژگان^۴ زبان فارسی است. جعفری (۱۳۸۸)، افزوده‌ها در زبان فارسی را براساس رویکردهای نقشی و صوری بررسی و در این زمینه پیکره‌های گردآوری کرده است. در ادامه، نحوه گردآوری واژگان برای انواع افزوده‌ها و فهرست واژه‌های دستوری بیان شده است.

همان‌گونه که پیش‌تر اشاره شد، افزوده‌های ربطی در فرانش منطقی مطرح هستند و بین دو بند یا دو جمله پیوند منطقی برقرار می‌کنند. برای تهیه فهرست افزوده‌های ربطی، کل حروف ربط موجود در پیکره گفتمان استخراج شد. در پیکره گفتمان به پیروی از درخت‌بانک گفتمانی پن برچسب‌های نماینده معنای منطقی در چهار گروه کلی شامل زمانی^۵، وابستگی^۶، مقایسه^۷ و گسترش^۸ دسته‌بندی می‌شوند. در این دسته‌بندی هر گره به دو زیربخش «نوع» و «زیرنوع» تقسیم‌بندی می‌شود (شمس‌فرد و بی‌جن‌خان، ۱۴۰۱: ۸۰). از آنجاکه هریک از حروف ربط می‌تواند در پیکره بیش از یک برچسب داشته باشد، فقط آن حروفی انتخاب شد که در بیش از ۶۰ درصد موارد، یک نوع برچسب معین دارد. به این ترتیب، ۳۷۲ حرف ربط انتخاب شد. پس از استخراج این فهرست، برچسب هریک از حروف ربط به‌صورت دستی بررسی شد و براساس شبکه نظام ربط نمایش داده‌شده در شکل (۱) بازبینی و دوباره دسته‌بندی شد.

برای تهیه فهرست افزوده‌های وجه و نگرشی از پیکره گفتمان (میرزایی و صفری، ۲۰۱۸)، فارس‌نت (شمس‌فرد و دیگران، ۲۰۱۰) و پیکره پژوهش جعفری (۱۳۸۸) استفاده شد. برای تهیه فهرست افزوده‌های وجه، فهرست کلمات دارای سه نوع برچسب «وجه»، «قطبیت» و «قطعیت» از پیکره گفتمان استخراج شد و براساس شبکه نظام وجه‌نمایی که در شکل (۲) مشاهده می‌شود، دسته‌بندی شد. علاوه‌براین، از پیکره پژوهش جعفری (۱۳۸۸) برای تهیه فهرست افزوده‌های وجه و نگرشی استفاده شد. برای افزوده‌های نگرشی مستخرج از پیکره جعفری، فهرست کلمات مترادف با آن‌ها با کمک فارس‌نت یافته و به مجموعه واژگان افزوده شد. از آنجاکه در برخی آثار نویسندگان پیکره معاصر از زبان محاوره استفاده شده است، بعضی عبارتها که صورت محاوره‌ای دارند؛ مانند «معلومه»، «ممکنه» نیز برچسب‌گذاری و به مجموعه واژگان افزوده شد. همچنین صورتهای منفی مانند «معلوم نیست»؛ درمقابل «معلوم است» نیز در مجموعه واژگان گنجانده شد. علاوه‌براین، صورتهای تصریفی افعالی مانند «فکر کردن»، «حس کردن»، «خواستن» و «دوست‌داشتن» نیز در مجموعه واژگان قرار گرفت. در مجموع، ۳۲۲ افزوده وجه و ۲۰۱ افزوده نگرشی برچسب‌گذاری شد.

از آنجاکه فهرست واژه‌های دستوری و محتوایی به حوزه تخصصی یا گونه^۹ هر متن وابسته است (میرزایی و صفری، ۱۳۹۴)، برای تهیه فهرست واژه‌های دستوری، از پیکره گردآوری شده در پژوهش حاضر استفاده شد. به این منظور، از روش وزن‌دهی فرکانس کلمه-معکوس فرکانس سند^{۱۰} استفاده شد^{۱۱}. براساس معیار فرکانس کلمه-معکوس فرکانس سند، ۱۹۷ واژه‌ای که مقادیر کوچک‌تری داشتند، واژه دستوری در نظر گرفته شدند. تا اینجا نحوه گردآوری واژگان شرح داده شد. در ادامه، نحوه عملکرد سامانه تشریح می‌شود.

در ابتدا، کلماتی از متن که در مجموعه واژگان تعریف شده‌اند، استخراج می‌شوند؛ سپس براساس برچسب آن واژه یا عبارت در مجموعه واژگان، برچسب‌گذاری می‌شوند. برای نمونه، حرف ربط «برای مثال»، در مجموعه واژگان براساس نظام حروف ربط به ترتیب سلسله‌مراتب دارای برچسب «تشریحی»، «بدلی»، «نمونه‌سازی» است. سامانه، تمام موارد وقوع عبارت «برای مثال» را در متن کشف و براساس برچسب واژگان برچسب‌گذاری می‌کند. برای استخراج کلمات از متن سعی بر آن بود که اگر یک کلمه یا عبارت زیرمجموعه عبارت دیگری در مجموعه واژگان باشد (مانند «که» و «به‌محض آنکه»)، برچسب عبارت بزرگ‌تر در نظر گرفته شود.

1. Penn discourse treebank

2. R. Prasad

3. logical metafunction

4. WordNet

5. temporal

6. contingency

7. comparison

8. expansion

9. genre

10. tf-idf

۱۱. در این روش میزان تکرار یک کلمه در یک مستند درمقابل تعداد تکرار آن در مجموعه کلیه مستندات سنجیده می‌شود. در این روش به هر کلمه در یک متن عددی تخصیص می‌دهد. این معیار برای کلماتی که تقریباً در تمام اسناد پیکره یافت شوند، دارای مقادیر کمینه است. این کلمات همان واژه‌های دستوری تلقی می‌شوند (شوتر و دیگران، ۲۰۰۸).

۲-۴ محاسبه بسامد نسبی شاخص‌ها

پس از استخراج کلمات از متن، بسامد نسبی مقادیر شاخص‌های معنایی برای هر متن محاسبه می‌شود. خروجی این مرحله بردار ویژگی‌های معنایی است که توصیفگر متن است. براساس ویژگی‌های مبتنی بر دستور نقش‌گرای هالیدی (هالیدی و متیسن، ۲۰۱۴) و واژه‌های دستوری هر متن با بردار ویژگی‌های عددی تعریف می‌شود.

بسامد نسبی هر گره در نظام‌های مبتنی بر دستور نقش‌گرای هالیدی^۱، بسامد انتخاب گره O1 برحسب انتخاب O2 است که گره والد آن در شبکه نظام محسوب می‌شود. درحقیقت، گره O2 دارای برچسب درشت^۲ و O1 دارای برچسب ریز^۳ است و زیرمجموعه آن برچسب تلقی می‌شود. بسامد نسبی O1 با داشتن گره O2 به‌منزله گره والد به این شکل تعریف می‌شود:

$$RF_d(O_1 | O_2) = \frac{N_d(O_1, O_2)}{N_d(O_2)} \quad (\text{فرمول ۱})$$

(آرگامون و دیگران، ۲۰۰۷)

در این فرمول، $N_d(O_1, O_2)$ تعداد رخداد هم‌زمان O1 و O2 است و $N_d(O_2)$ تعداد رخداد O2 در هر متن است و d به معنای متن یا همان سند است (آرگامون و دیگران، ۲۰۰۷: ۸). به‌بیان ساده‌تر، چه نسبت از تعداد رخداد کلمات در متن که دارای برچسب O2 هستند، دارای برچسب ریزتر O1 نیز هستند.

برای نمونه، در شبکه نظام حروف ربط، گره‌های «تشریحی»، «گسترشی» و «تفصیلی» زیرمجموعه گره حرف ربط محسوب می‌شوند. بسامد نسبی گره «تشریحی» با داشتن گره «حرف ربط» به این معناست که چه نسبت از حروف ربط به‌کارگرفته‌شده در یک متن، از نوع تشریحی هستند. درباره شبکه نظام وجه‌نمایی و افزوده نگرشی نیز بسامد نسبی براساس فرمول (۱) محاسبه می‌شود. برای مثال، در شبکه نظام وجه‌نمایی، گره‌های «بالا»، «پایین» و «میان» زیرمجموعه گره «مقدار» هستند. بسامد نسبی گره «بالا» با داشتن گره «مقدار» به این معناست که چه نسبت از افزوده‌های وجه در یک متن، دارای مقدار «بالا» هستند.

علاوه بر ویژگی بسامد نسبی، صرفاً برای نظام وجه‌نمایی، ویژگی‌هایی براساس ترکیب گره‌ها تعریف می‌شود. برای هر جفت گره در زیرنظام‌های مختلف وجه‌نمایی (مانند «نوع» و «مقدار»)، بسامد نسبی واژه‌هایی که با نوع دو گره برچسب خورده‌اند، براساس ترکیب گره‌های والد آن‌ها تعریف می‌شود (آرگامون و دیگران، ۲۰۰۷). برای مثال، بسامد نسبی جفت گره «وجه احتمال-عادت» و «بالا» با داشتن جفت گره «نوع» و «مقدار» بیان می‌کند که چه نسبت از کلماتی که دارای «وجه احتمال-عادت» هستند، در زیرنظام «مقدار» دارای مقدار «بالا» نیز هستند؛ (مانند «حتماً» و «دائماً»).

برای مقایسه کارایی ویژگی‌های مبتنی بر دستور نقش‌گرای هالیدی و ویژگی‌های دستوری، ویژگی بسامد واژه‌های دستوری تعریف می‌شود. به این منظور، واژه‌های دستوری در هر متن شناخته می‌شوند و بسامد نسبی هریک از آن‌ها به‌منزله یک ویژگی محاسبه می‌گردد. برای محاسبه بسامد نسبی واژه‌های دستوری از فرمول (۲) استفاده شده است:

$$\frac{\text{count}(w)}{\sum_{w \in FW} \text{count}(w')} \quad (\text{فرمول ۲})$$

(آرگامون و دیگران، ۲۰۰۷)

در این فرمول، W هریک از واژه‌های دستوری است که بسامد آن نسبت به بسامد کل واژه‌های دستوری در یک متن خاص محاسبه می‌شود (آرگامون و دیگران، ۲۰۰۷: ۱۳). در صورت کسر فرمول (۲)، تعداد رخداد یک واژه دستوری خاص و در مخرج آن تعداد رخداد کل واژه‌های دستوری در متن قرار دارد.

۳-۴ یادگیری ماشینی

برای مرحله یادگیری ماشینی از روش یادگیری با ناظر^۴ استفاده می‌شود. در این روش به یک سامانه، مجموعه‌ای از جفت‌های ورودی-خروجی ارائه می‌شود و سامانه می‌کوشد تابعی از ورودی به خروجی را فراگیرد. یادگیری باناظر به تعدادی داده ورودی برای آموزش

۱. هر شبکه نظام در هر سطح شامل چند گره است. برای مثال، شبکه نظام حروف ربط در سطح اول شامل سه گره، افزوده تشریحی، گسترشی و تفصیلی است. در سطح بعد، گره تشریحی شامل دو گره بدلی و واضح‌سازی است. درحقیقت، گره تشریحی والد دو گره بدلی و واضح‌سازی است. در اینجا گره تشریحی نسبت به دو گره بدلی و واضح‌سازی دارای برچسب درشت و دو گره بدلی و واضح‌سازی دارای برچسب ریزتر هستند.

2. coarse-grained
3. fine-grained
4. supervised learning

سامانه نیاز دارد. همان‌طور که پیش‌تر اشاره شد، خروجی مرحله دوم بردار ویژگی‌های معنایی است که توصیفگر متن است. در مورد تشخیص خودکار نویسنده متن، ورودی‌ها بردار ویژگی‌های زبانی متون هستند و خروجی همان نام نویسنده متن است. در پژوهش حاضر، برای دسته‌بندی نویسندگان، از روش یادگیری عمیق^۱ استفاده شده است. رویکردهای مبتنی بر یادگیری عمیق در پژوهش‌های نوین در حوزه پردازش زبان طبیعی از جمله تشخیص خودکار هویت نویسنده بسیار مورد توجه هستند (اوجندو^۲ و دیگران، ۲۰۲۰). یکی از انواع طبقه‌بندی‌های روش یادگیری عمیق، طبقه‌بند پرسپترون چند لایه^۳ است که از آن در این پژوهش بهره گرفته شده است. برای استفاده از این طبقه‌بند از یک ابزار^۴ یادگیری ماشینی به زبان پایتون^۵ استفاده شده است. الگوریتم یادگیری بردار ویژگی‌های معنایی هر متن را به‌منزله ورودی دریافت می‌کند. بردار ویژگی‌های معنایی می‌تواند انواع ویژگی‌های مبتنی بر دستور نقش‌گرای هالیدی (شبکه نظام حروف ربط، وجه‌نمایی و افزوده نگرشی) یا ترکیب انواع مختلف آن‌ها را دربرگیرد. از سوی دیگر، این بردار می‌تواند ویژگی بسامد نسبی واژه‌های دستوری را شامل شود. در حالت دیگر، در انتخاب نوع بردار ویژگی می‌توان تمام ویژگی‌های مبتنی بر دستور نقش‌گرای هالیدی و ویژگی بسامد نسبی واژه‌های دستوری را باهم آمیخت. در نهایت، الگوریتم یادگیری نام نویسنده متن را به‌منزله خروجی برمی‌گرداند. در ادامه، کارایی این الگوریتم در حالت‌های مختلف بردار ویژگی ارزیابی و گزارش می‌شود.

۵- یافته‌ها

برای ارزیابی کیفیت سامانه طراحی شده در تشخیص خودکار هویت نویسنده متن از روش ارزیابی متقاطع K تایی^۶ استفاده شده است. اگر کل داده‌ها را به‌طور تصادفی به K زیرنمونه تقسیم کنیم، می‌توانیم در هر مرحله از فرایند، تعداد K-1 از این لایه‌ها را به‌منزله مجموعه داده آموزشی و یکی را به‌منزله مجموعه داده آزمون در نظر بگیریم. مشخص است که تعداد تکرارهای فرایند برابر با K خواهد بود (رفائیل زاده و دیگران، ۲۰۰۹). در واقع، یادگیری سامانه با کمک داده آموزشی صورت می‌گیرد و با استفاده از داده آزمون ارزیابی می‌شود. در این پژوهش، $K=5$ در نظر گرفته شده است. نکته درخور توجه در ارزیابی این است که نه تنها داده آموزشی و آزمون با یکدیگر هم‌پوشانی ندارند، بلکه برای ارزیابی سامانه تشخیص سبک نگارش یک نویسنده، نباید یک داستان خاص بین داده آموزشی و آزمون تقسیم شود. این نکته‌ای است که در برخی پژوهش‌ها نادیده گرفته شده است؛ به این علت، در ارزیابی سامانه، کل آثار هر نویسنده به پنج قسمت تقسیم شد؛ چهار قسمت برای آموزش و یک قسمت برای آزمون. این تقسیم‌بندی به گونه‌ای است که داستان‌های موجود در بخش آموزش از داستان‌های موجود قسمت آزمون متمایز هستند. برای محاسبه کارایی سامانه از معیار دقت^۷ استفاده شد. دقت در تشخیص خودکار نویسنده عبارت است از نسبت نتایج درست به دست‌آمده از سامانه به کل تعداد نمونه‌های موجود در داده آزمون (درونه و شریعتی، ۲۰۱۴). از آنجاکه فرایند ارزیابی پنج بار تکرار شد، در نهایت، از نتایج به دست‌آمده در پنج مرحله برای کل سامانه و همین‌طور برای هر نویسنده میانگین وزن دار^۸ گرفته شد.^۹

الگوی طراحی شده در چند حالت آزموده شد: حالت نخست، بردار ویژگی شامل ویژگی‌های تعریف شده بر اساس دستور نقش‌گرای نظام‌مند است؛ حالت دوم، بردار ویژگی از نوع ویژگی‌های واژه‌های دستوری است؛ حالت سوم، ترکیب حالت نخست و دوم است. در هر حالت، دقت سامانه به روش ارزیابی متقاطع K تایی سنجیده شد. در جدول (۴) نتایج ارزیابی سامانه به‌طور خلاصه گزارش شده است. همان‌گونه که مشاهده می‌شود، بهترین نتایج حاصل از ترکیب روش مبتنی بر دستور نقش‌گرای نظام‌مند و بسامد واژه‌های دستوری است. همچنین دقت روش مبتنی بر واژه‌های دستوری از روش دستور نقش‌گرا بسیار بیشتر است. در شکل (۷)، نتایج ارزیابی سامانه به‌طور کامل و برای هر نویسنده به‌طور مجزا نمایش داده شده است. علاوه بر این، میزان کارایی تمام ویژگی‌های تعریف شده در سامانه

1. deep learning
2. A. Uchendu
3. multilayer perceptron
4. <https://scikit-learn.org/>
5. python
6. K-Fold Cross Validation
7. accuracy
8. weighted mean

۹. در محاسبه میانگین وزن‌دار یک مجموعه عامل نابرابر، برای هر یک از عامل‌ها، وزن یا ارزش معینی در نظر گرفته می‌شود؛ سپس آن عامل در وزن معین ضرب می‌گردد. آنگاه جمع این ارقام به دست‌آمده بر مجموع وزن‌ها تقسیم می‌شود.

بررسی و کاراترین ویژگی‌ها در عملکرد سامانه انتخاب شد. برای این کار از آزمون F تحلیل واریانس یک‌طرفه^۱ استفاده شد.^۲ در این پژوهش مقادیر بردار ویژگی (کل ویژگی‌های مبتنی بر دستور نقش‌گرای نظام‌مند و ویژگی واژه‌های دستوری)، براساس نتایج به‌دست‌آمده از آزمون F تحلیل واریانس یک‌طرفه مرتب و ۱۵۰ ویژگی برتر انتخاب شد. براساس این معیار، بسامد نسبی واژه‌های دستوری «آن»، «رو» و «شد» به‌ترتیب دارای بالاترین امتیاز درمیان کل ویژگی‌ها هستند.

جدول (۴). نتایج ارزیابی کلی سامانه براساس معیار دقت

نام مجموعه و ویژگی	کلمات دستوری	کلمات دستوری، افزوده وجه و نگرشی	کلمات دستوری و افزوده ربطی	کل دستوری نقش‌گرا	کلمات دستوری و کل دستوری نقش‌گرا
میانگین وزن‌دار	۶۸.۹۸	۷۲.۵۳	۷۲.۰۵	۴۶	۷۴.۲۱

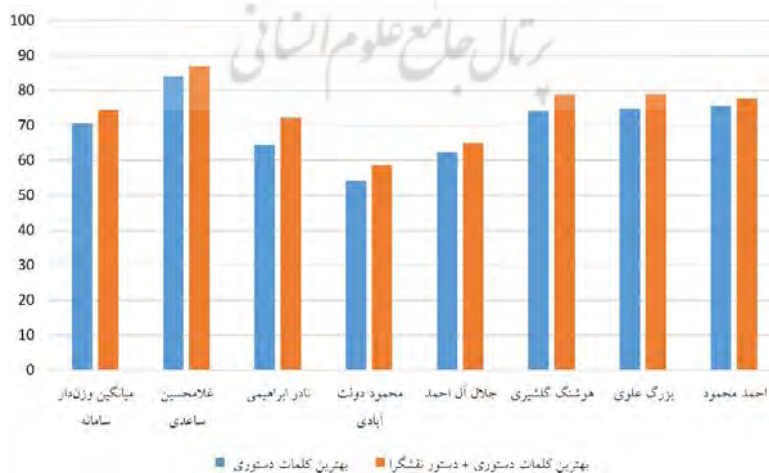


شکل (۷). نمودار مقایسه‌ای ارزیابی سامانه براساس معیار دقت

پس از انتخاب بهترین ویژگی‌ها کارایی سامانه در دو حالت سنجیده شد: ۱. کلمات دستوری؛ ۲. کلمات دستوری (بهترین ویژگی‌ها) و ویژگی‌های مبتنی بر دستور نقش‌گرا. این نتایج در جدول (۵) گزارش شده است. همان‌گونه که در شکل (۸) مشاهده می‌شود، بهترین نتایج را ترکیب حالت بهترین کلمات دستوری و دستور نقش‌گرا به دست می‌دهد.

جدول (۵). نتایج ارزیابی سامانه براساس معیار دقت (بهترین ویژگی‌ها)

نام مجموعه و ویژگی	میانگین وزن‌دار
بهترین کلمات دستوری	۷۰/۴۷
بهترین کلمات دستوری و دستور نقش‌گرا	۷۴/۴۷



شکل (۸). نمودار مقایسه‌ای ارزیابی سامانه براساس معیار دقت (بهترین ویژگی‌ها)

1. ANOVA

۲. آزمون تحلیل واریانس یک‌طرفه برای آزمون مقایسه میانگین یک متغیر کمی درمیان بیش از دو گروه مستقل استفاده می‌شود.

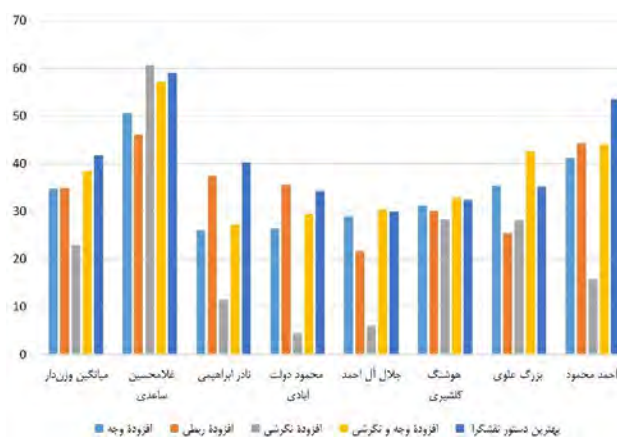
۶- بحث و نتیجه‌گیری

در پژوهش حاضر، برای تعیین هویت نویسنده متن یک سامانه خودکار به روش یادگیری عمیق معرفی شد. برای طراحی سامانه از ویژگی بسامد واژه‌های دستوری و مفاهیم دستور نقش‌گرای نظام‌مند استفاده شد. سپس عملکرد سامانه در سه حالت ویژگی‌های مبتنی بر واژه‌های دستوری، ویژگی‌های مبتنی بر دستور نقش‌گرای نظام‌مند و ترکیب این دو نوع ویژگی سنجیده و تحلیل شد و بهترین ویژگی‌ها انتخاب شدند. نتایج ارزیابی سامانه نشان می‌دهد روش محاسبه بسامد واژه‌های دستوری نسبت به روش مبتنی بر دستور نقش‌گرای نظام‌مند (شبکه نظام ربط، وجه‌نمایی و افزوده نگرشی) در تشخیص هویت نویسنده متن برتری دارد؛ اما در صورتی که ویژگی‌های زبانی مبتنی بر دستور نقش‌گرای نظام‌مند هالیدی در کنار ویژگی بسامد واژه‌های دستوری به کار روند، کارایی سامانه نسبت به حالتی که فقط از ویژگی بسامد واژه‌های دستوری استفاده شود، افزایش می‌یابد.

به نظر می‌رسد برای این نتیجه که ویژگی بسامد نسبی واژه‌های دستوری کارایی بالاتری نسبت به ویژگی‌های مبتنی بر دستور نقش‌گرا در تفکیک سبک نویسندگان داشته، دلایل متفاوتی وجود داشته است. دلیل نخست می‌تواند به کنترل ناپذیری خودآگاه نویسنده در استفاده از واژه‌های دستوری مربوط باشد. نویسنده متن واژه‌های دستوری را به‌طور ناخودآگاه به کار می‌گیرد و برخلاف واژه‌های محتوایی کنترل خودآگاه بر به‌کاربردن یا به‌کارنبردن آن‌ها ندارد (گلشایی، ۱۳۹۸). از آنجاکه تعیین مرز دقیق بین واژه‌های محتوایی و دستوری عملاً امکان‌پذیر نیست، می‌توان در خصوص واژه‌های دستوری یک پیوستار تعریف کرد. در این پیوستار هر چه به سمت دستوری بودن واژه پیش روییم، کنترل خودآگاه در استفاده از واژه‌ها کاهش می‌یابد. در این پژوهش، براساس آزمون F تحلیل واریانس یک‌طرفه، ضمایر «آن»، «این»، «آن‌ها»، «من» و «خود» ویژگی‌های بسیار مناسبی برای تمایز گذاشتن میان سبک نویسندگان تلقی می‌شوند. نتایج حاصل از پژوهش حاضر در مورد کارترین ویژگی‌ها در تشخیص هویت نویسنده گواهی بر این مدعاست که ضمایر در این پیوستار به سمت دستوری بودن تمایل دارند و نویسنده کنترل خودآگاه بر به‌کارگیری آن‌ها ندارد. بدیهی است که این پیوستار خود می‌تواند موضوع یک پژوهش جدید باشد.

دومین دلیل برای اینکه در تشخیص هویت نویسنده ویژگی‌های مبتنی بر دستور نقش‌گرا به‌طور مستقل مؤثرتر از واژه‌های دستوری عمل نکرده‌اند، می‌تواند این باشد که استفاده از واژه‌های دستوری به‌جای یکدیگر، تفاوت سبکی بیشتری ایجاد می‌کند. برای مثال، تفاوت سبک در متونی که از واژه‌های «من» و «این‌جانب» استفاده می‌کنند، درخور توجه است؛ این درحالی است که برای مثال، «مسلماً» و «مسلم است» که در زیرنظام جهت‌گیری از شبکه نظام وجه‌نمایی اولی «ضمنی» و دومی «صریح» تلقی می‌شوند، تفاوت سبکی چندانی ایجاد نمی‌کنند. این مطلب درباره «باید» و «نباید» که در زیرنظام قطبیت با یکدیگر تفاوت دارند، نیز صادق است.

دلیل دیگر برای کارایی بیشتر واژه‌های دستوری نسبت به ویژگی‌های مبتنی بر دستور نقش‌گرا، می‌تواند اختلاف درصد حضور واژه‌های دستوری و کلمات مرتبط با سه نوع ویژگی دستور نقش‌گرا در پیکره باشد. به‌جز افزوده ربطی «و» از شبکه نظام حروف ربط که فراوانی بسیار بالایی در پیکره دارد، درصد حضور واژه‌های دستوری بیشتر از واژه‌های مرتبط با ویژگی‌های مبتنی بر دستور نقش‌گرا است. برای بررسی بهتر ویژگی‌های مبتنی بر دستور نقش‌گرای هالیدی، هریک از این سه شبکه نظام حروف ربط، وجه‌نمایی و افزوده نگرشی به‌صورت مجزا و در ترکیب با یکدیگر ارزیابی شدند و نتایج آن‌ها با یکدیگر و با حالت بهترین ویژگی‌های سه نوع نظام مقایسه شدند که در شکل (۹) مشاهده می‌شود.



شکل (۹). نمودار مقایسه‌ای کارایی سامانه با انواع ویژگی‌های دستور نقش‌گرا

همان‌گونه که در نمودار (۹) مشاهده می‌شود، ترکیب بهترین ویژگی‌های مرتبط با هر سه نوع افزوده وجه، نگرشی و ربطی بالاترین دقت را در میان دیگر حالت‌ها دارد. پس‌از آن، افزوده وجهی (افزوده نگرشی) قرار دارد. گواه دیگر در برتری افزوده وجهی نسبت به افزوده ربطی، درصد حضور ویژگی‌های مرتبط با افزوده وجهی در میان ویژگی‌های منتخب است (۳۹/۰۷ درصد). این درحالی است که درصد حضور ویژگی‌های مرتبط با افزوده ربطی در این فهرست ۱۱/۲۵ درصد است.

تحلیل دقیق‌تر بهترین ویژگی‌های افزوده وجهی نشان می‌دهد که ویژگی‌های مرتبط با زیرنظام‌های «مقدار» و «نوع» از شبکه نظام افزوده وجهی به دو صورت مستقل و ترکیب با سایر زیرنظام‌ها بسیار کارا هستند. زیرنظام‌های «قطبیت» و «جهت‌گیری» در ترکیب با سایر زیرنظام‌ها به بهبود عملکرد سامانه کمک شایانی می‌کنند. درباره زیرنظام وجه‌سازی (وجه احتمال-عادت) علاوه بر اینکه کل زیرنظام و زیرنظام «وجه احتمالی» در کارایی سامانه بسیار مؤثر هستند، زیرنظام «وجه عادت» بیشترین تأثیر را دارد. در زیرنظام تعدیل‌سازی (وجه تمایل-التزام) نیز علاوه بر اثرگذاری کل زیرنظام و زیرنظام «وجه التزامی»، زیرنظام «وجه تمایلی» بالاترین کارایی را دارد. در شبکه نظام افزوده نگرشی، زیرنظام گزاره‌ای به‌ویژه از نوع «اخلاقی» و «فرض» بیشترین میزان اثربخشی را دارد.

در نظام شبکه حروف ربط، بسامد افزوده‌های ربطی هریک از سه زیرنظام «تشریحی»، «گسترشی» و «تفصیلی» (به‌طور مجزا) به کل افزوده‌های ربطی ویژگی‌هایی ارزشمند به شمار می‌روند. برپایه ویژگی‌های منتخب، زیرنظام «گسترشی» نتایج مطلوب را به دست می‌دهد. در این زیرنظام، بسامد نسبی افزوده‌هایی از نوع «افزایشی» و «تباینی» امتیاز بالایی دارند. در زیرنظام «تفصیلی»، زیرنظام «سببی-شرطی» بسیار مؤثر است. در این زیرنظام، زیرنظام «مشخص» به‌ویژه حروف ربط از نوع «علت» و پس‌از آن، «نتیجه» تأثیرگذار هستند. در زیرنظام «تشریحی»، افزوده‌های ربطی «بدلی» به‌ویژه از نوع «نمونه‌سازی» کارآمد هستند.

نتایج کلی این پژوهش همسو با نتایج مقاله مرجع آرگامون و دیگران (۲۰۰۷) است. در پژوهش آرگومان و دیگران (۲۰۰۷)، روش بسامد نسبی واژه‌های دستوری نسبت به روش دستور نقش‌گرا (شبکه نظام حروف ربط، شبکه نظام وجه‌نمایی و شبکه نظام افزوده نگرشی) برتری دارد و ترکیب روش واژه‌های دستوری و دستور نقش‌گرا بالاترین کارایی را دارد. در پژوهش آرگومان و دیگران (۲۰۰۷) روش دستور نقش‌گرا به افزایش کارایی سامانه کمک کرده است که با نتایج کلی پژوهش حاضر همسو است.

یکی از تفاوت‌های پژوهش حاضر با پژوهش آرگومان و دیگران (۲۰۰۷) این است که در پژوهش آرگومان و دیگران (۲۰۰۷)، تعداد واژه‌های دستوری ۶۷۵ واژه در نظر گرفته شده است؛ درحالی که در پژوهش حاضر تعداد کل واژه‌های دستوری ۱۹۷ واژه است و پس از مرحله یافتن بهترین ویژگی‌ها تعداد ۷۴ واژه از واژه‌های دستوری انتخاب شده است. آرگومان و دیگران (۲۰۰۷) با استفاده از الگوریتم طبقه‌بندی ماشین بردار پشتیبان^۱ و با ۶۷۵ واژه دستوری، تنها با ویژگی بسامد نسبی واژه‌های دستوری به دقتی حدود ۸۵ درصد در تشخیص هویت نویسنده برای زبان انگلیسی دست یافته‌اند. در این پژوهش، برای مقایسه نتایج کار با مقاله مرجع از ۶۶۴ واژه دستوری و الگوریتم ماشین بردار پشتیبان استفاده شد و دقتی معادل ۸۲.۷۵ درصد برای زبان فارسی به دست آمد که به نتایج پژوهش مرجع نزدیک است. در این مقایسه مشخص شد فرایند یادگیری دچار بیش‌برازش^۲ شده است^۳ و این مورد را می‌توان از نقاط ضعف پژوهش مرجع یاد کرد که البته در آن مقاله نیز چند بار به آن اشاره شده است. به این دلیل، در پژوهش حاضر برای جلوگیری از بیش‌برازش تعداد واژه‌های دستوری بسیار کمتر از مقاله مرجع و برابر با ۱۹۷ کلمه در نظر گرفته شده است. جدا از اینکه تعداد واژه‌های دستوری چقدر است، نتایج حاصل از افزودن ویژگی‌های مبتنی بر دستور نقش‌گرا به سامانه همسو با پژوهش آرگومان و دیگران (۲۰۰۷) است. در مقاله مرجع کارایی سامانه با افزودن ویژگی‌های مبتنی بر دستور نقش‌گرا، حدود ۵ درصد بیشتر از حالتی است که صرفاً از واژه‌های دستوری استفاده شود. این افزایش کارایی برابر با افزایش کارایی پژوهش حاضر در حالت قرارگرفتن ویژگی‌های مبتنی بر دستور نقش‌گرا در کنار واژه‌های دستوری است (ارجاع به جدول ۴). تفاوت دیگر نتایج مقاله مرجع با پژوهش حاضر این است که ویژگی‌های مبتنی بر دستور نقش‌گرای هالییدی (به‌طور مستقل از واژه‌های دستوری) برای زبان فارسی به اندازه زبان انگلیسی در تشخیص هویت نویسنده کارایی نداشته‌اند. علاوه بر تفاوت‌های دو زبان فارسی و انگلیسی، به نظر می‌رسد تفاوت در حجم واژگان تعریف‌شده برای شبکه نظام‌های دستور نقش‌گرا از عوامل این تفاوت باشد.

1. SVM

2. overfitting

۳. بیش‌برازش به این معناست که الگوریتم داده‌های آموزشی را بسیار خوب یاد گرفته است؛ ولی توانایی پیش‌بینی خوب داده‌های جدید را ندارد.

شبکه‌های مختلف دستور نقش‌گرا (شبکه‌ی نظام حروف ربط، شبکه‌ی نظام وجه‌نمایی و شبکه‌ی نظام افزوده‌ی نگرشی) به‌طور مستقل نتایج درخوردنی توجیهی را در تعیین هویت نویسنده به دست نمی‌دهند. درواقع، نتایج مطلوب، نتایجی است که از ترکیب ویژگی‌های منتخب از هر سه نظام و ویژگی‌های ویژه‌ی دستوری به دست می‌آید. گرچه رویکرد نقش‌گرای نظام‌مند هالیدی به‌ویژه در فرانش بینافردی درباره‌ی وجه، دسته‌بندی دقیقی ارائه می‌کند، شبکه‌ی نظام‌های این دستور به‌صورت مستقل از واژه‌های دستوری عملکرد مطلوبی در تشخیص هویت نویسنده متن ندارند و فقط در کنار واژه‌های دستوری به افزایش کارایی سامانه یاری می‌رسانند.

از کاستی پژوهش حاضر، می‌توان به این مسئله اشاره کرد که به‌جز شمار محدودی از افعال در فهرست افزوده‌ی وجهی و افعال کمکی، به مقوله‌ی فعل در طراحی سامانه‌ی هویت نویسنده توجه نشد. در پژوهش‌های آینده می‌توان حجم واژگان برای افزوده‌های ربطی و افزوده‌های وجهی به‌ویژه با مقوله‌ی فعل را گسترش داد و به این شکل به پژوهش‌های حوزه‌ی دستور نقش‌گرای نظام‌مند هالیدی (هالیدی و متیسن، ۲۰۱۴) کمک کرد. در پژوهش‌های آینده می‌توان الگوریتم استخراج سه نوع ویژگی دستوری نقش‌گرا و واژه‌های دستوری را درباره‌ی دیگر آثار داستانی نیز به کار برد و بسامد انواع ویژگی‌ها را به دست آورد و از آن در نقد ادبی و سبک‌شناسی برای اهداف زبانی استفاده کرد.

منابع

- آل احمد، جلال (۱۳۴۶). *نفرین زمین*. تهران: فردوس.
- آل احمد، جلال (۱۳۵۰). *پنج داستان*. تهران: فردوس.
- ابراهیمی، نادر (۱۳۷۴). *یک عاشقانه آرام*. تهران: روزبهان.
- ابراهیمی، نادر (۱۳۹۹). *بر جاده‌های آبی سرخ*. تهران: روزبهان.
- جعفری، آریتا (۱۳۸۸). بررسی افزوده‌ها در زبان فارسی: براساس رویکردهای نقشی و صوری. *دستور (ویژنامه نامه فرهنگستان)*، ۵(۱)، ۱۲۸-۱۵۵.
- حسین حمه، همزه؛ علی اکبری، نسرین؛ کریمی، یادگار (۱۴۰۰). بررسی وجه و وجهیت در کردی سورانی: تحلیلی نقش‌گرا. *مطالعات زبان‌ها و گویش‌های غرب ایران*، ۹(۴)، ۱-۲۳.
- دولت‌آبادی، محمود (۱۳۹۵). *روزگار سپری‌شده مردم سالخورده*. تهران: چشمه.
- دولت‌آبادی، محمود (۱۴۰۱). *کلیدر*. چاپ ۳۷. تهران: فرهنگ معاصر.
- ساعدی، غلامحسین (۱۳۷۷). *آشفته‌حالان بیدارینخت*. تهران: نگاه.
- ساعدی، غلامحسین (۱۳۹۷). *غریبه در شهر*. تهران: نگاه.
- شمس‌فرد، مهرنوش؛ بی‌جن‌خان، محمود (۱۴۰۱). *پردازش متن و گفتار فارسی: مروری بر مبانی نظری و آخرین یافته‌های پژوهشی*. تهران: سمت.
- عارفی، سمیه؛ بصیری، محمداحسان؛ روزمند، امید (۱۴۰۰). انتخاب ویژگی برای شناسایی نویسنده در متون کوتاه برخط فارسی. *فناوری اطلاعات و ارتباطات ایران*، ۱۳(۴۷-۴۸)، ۳۵-۵۷.
- علوی، بزرگ (۱۳۸۶). *ورق‌پاره‌های زنجان*. تهران: نگاه.
- علوی، بزرگ (۱۳۹۹). *گیله‌مرد*. تهران: نگاه.
- فرهمندپور، زینب؛ نیک‌مهر، هومن؛ منصوری‌زاده، محرم؛ طبیب‌زاده، امید (۱۳۹۱). یک سیستم نوین هوشمند تشخیص هویت نویسنده فارسی زبان براساس سبک نوشتاری. *محاسبات نرم*، ۱(۲)، ۲۶-۳۵.
- گلشنائی، رامین (۱۳۹۸). واژه‌های دستوری به‌مثابه نشانگرهای گویش فردی: رویکردی پیکره‌ای به شناسایی هویت نویسنده در زبان فارسی. *جستارهای زبانی*، ۱۰(۳)، ۳۱۷-۲۹۳.
- گلشیری، هوشنگ (۱۳۵۰). *کریستین و کید*. تهران: کتاب زمان.
- گلشیری، هوشنگ (۱۳۷۰). *در ولایت هوا*. استکهلم: عصر جدید.
- گلشیری، هوشنگ (۱۴۰۰). *شازده احتجاب*. چاپ ۱۸. تهران: نیلوفر.
- محمود، احمد (۱۳۵۳). *همسایه‌ها*. تهران: امیرکبیر.
- محمود، احمد (۱۳۸۱). *غریبه‌ها و پسرک بومی*. تهران: معین.
- میرزایی، آزاده (۱۳۹۷). بازتعریف مفاهیم بند پایه و بند پیرو براساس رویکرد نقش‌گرا. *زبان و زبان‌شناسی*، ۱۳(۲۶)، ۱۱۷-۱۳۳.
- میرزایی، آزاده (۱۴۰۰). رابطه قطبیت و وجهیت بندی در زبان فارسی. *مطالعات زبان‌ها و گویش‌های غرب ایران*، ۹(۱)، ۱۱۳-۱۳۵.

میرزایی، آزاده؛ صفری، پگاه (۱۳۹۴). ساخت واژه-متن‌های تخصصی و عمومی زبان فارسی براساس بسامدگیری واژه‌های نقشی و محتوایی. مجموعه مقالات نخستین همایش ملی زبان‌شناسی پیکرهای (صص. ۱۷۵-۱۹۱). تهران: نویسه پارسی.

References

- Alavi, B. (2007). *Scrap papers from prison*. Tehran: Negah. (In Persian)
- Alavi, B. (2020). *Gilemard*. Tehran: Negah. (In Persian)
- Al-e-Ahmad, J. (1967). *The cursing of the land*. Tehran: Ferdous. (In Persian)
- Al-e-Ahmad, J. (1971). *Five stories*. Tehran: Ferdous. (In Persian)
- Arefi, S., Basiri, M. E., & Roozmand, O. (2021). Feature selection for author identification of Persian online short texts. *Journal of Information and Communication Technology*, 13(47-48), 35-57. <https://dorl.net/dor/20.1001.1.27170414.1400.13.47.4.0> (In Persian)
- Argamon, S., & Koppel, M. (2013). A systemic functional approach to automated authorship analysis. *Journal of Law & Policy*, 21(2), 299-315.
- Argamon, S., Whitelaw, C., Chase, P., Hota, S. R., Garg, N., & Levitan, S. (2007). Stylistic text classification using functional lexical features. *Journal of the American Society for Information Science and Technology*, 58(6), 802-822. <https://doi.org/10.1002/asi.20553>.
- Assi, S. M. (1997). Farsi linguistic database (FLDB). *International journal of Lexicography*, 10(3), 5.
- Bijankhan, M., Sheykhzadegan, J., Bahrani, M., & Ghayoomi, M. (2011). Lessons from building a Persian written corpus: Peykare. *Language resources and evaluation*, 45(2), 143-164. <https://doi.org/10.1007/s10579-010-9132-x>.
- Coulthard, M. (2004). Author identification, idiolect, and linguistic uniqueness. *Applied Linguistics*, 25(4), 431-447. <https://doi.org/10.1093/applin/25.4.431>.
- Dabagh, R. M. (2007). Authorship attribution and statistical text analysis. *Advances in Methodology and Statistics*, 4(2), 149-163. <https://doi.org/10.51936/uvjx7198>.
- Daroonch, A. H., & Shariati, A. (2014). Metrics for evaluation of the author's writing styles: Who is the best? *Chaos: An Interdisciplinary Journal of Nonlinear Science*, 24(3). <https://doi.org/10.1063/1.4895468>.
- Dowlatabadi, M. (2016). *The elderly people's elapsed time*. Tehran: Cheshmeh. (In Persian)
- Dowlatabadi, M. (2022). *Kelidar* (37th print). Tehran: Farhange Moaser. (In Persian)
- Ebrahimi, N. (1995). *A quiet Romance*. Tehran: Rouzbahan. (In Persian)
- Ebrahimi, N. (2020). *On the blue and red paths*. Tehran: Rouzbahan. (In Persian)
- Farahmandpoor, Z., Nikmehr, H., Mansoorizade, M., & Tabibzadeh Ghamsary, O. (2013). A novel intelligent Persian authorship system based on writing style. *Soft Computing Journal*, 1(2), 26-35. <https://dorl.net/dor/20.1001.1.23223707.1391.1.2.60.9> (In Persian)
- Golshaie, R. (2019). Function words as idiolect markers: A corpus-based approach to authorship attribution in Farsi. *Language Related Research*, 10(3), 293-317. (In Persian)
- Golshiri, H. (1971). *Christine and kid*. Tehran: Ketab-e Zaman. (In Persian)
- Golshiri, H. (1991). *Dar velayat-e Hava*. Stockholm: Asr-e Jadid. (In Persian)
- Golshiri, H. (2021). *Prince Ehtejab* (18th print). Tehran: Niloufar. (In Persian)
- Halliday, M. A. K., & Matthiessen, C. M. M. (2014). *Halliday's introduction to functional Grammar* (4th ed.). Oxon: Routledge.
- Hussein Hama, H., Ali-akbari, N., & Karimi, Y. (2022). Mood and modality in Sorani Kurdish: A functional analysis. *Researches in Western Iranian Languages and Dialects*, 9(4), 1-23. <https://doi.org/10.22126/JLW.2021.6008.1504>. (In Persian)
- Jafari, A. (2008). An analysis of adjuncts in Persian: A syntacto-discoursal approach. *Dastoor*, 5(1), 128-155. (In Persian)
- Mahmoud, A. (1974). *The neighbors*. Tehran: Amirkabir. (In Persian)
- Mahmoud, A. (2002). *The strangers and the little native boy*. Tehran: Moin. (In Persian)
- Mandenthal, T. C. (1887). The characteristics curves of composition science. *Science*, 9(214s), 237-246.
- Martinez-Galicia, J. A., Embarcadero-Ruiz, D., Ríos-Orduña, A., & Gómez-Adorno, H. (2022). Graph-based siamese network for authorship verification. *CLEF 2022 Labs and Workshops, Notebook Papers*. Italy: Bologna.
- Matthiessen, C. M., & Bateman, J. A. (1991). *Text generation and systemic-functional linguistics: Experiences from English and Japanese*. London and New York: Frances Pinter Publishers.

- Mirzaei, A. (2017). Redefining the concepts of dependent and independent clauses according to a functional approach. *Language and Linguistics*, 13(26), 117-132. (In Persian)
- Mirzaei, A. (2018). *An introduction to corpus linguistics*. Tehran: Allameh Tabataba'i University Press. (In Persian)
- Mirzaei, A. (2021). The relationship between polarity and clausal modality in Persian. *Journal of Western Iranian Languages and Dialects*, 9(1), 113-135. <https://doi.org/10.22126/jlw.2020.5372.144>. (In Persian)
- Mirzaei, A., & Safari, P. (2014). Building specialized and general documents in Persian based on the frequency of function and content words. In *proceeding of The 1st National Conference on Corpus Linguistics* (PP. 175-192). Tehran: Nevisesh Parsi. (In Persian)
- Mirzaei, A., & Safari, P. (2018). Persian discourse treebank and coreference corpus. In *Proceedings of the eleventh international conference on language resources and evaluation (lrec 2018)*. Japan: Miazaky.
- Najafi, M., & Tavan, E. (2022). Text-to-text transformer in authorship verification via stylistic and semantical analysis. *CLEF 2022 Labs and Workshops, Notebook Papers*. Italy: Bologna.
- Plecháč, P. (2021). Relative contributions of Shakespeare and Fletcher in Henry VIII: An analysis based on most frequent words and most frequent rhythmic patterns. *Digital Scholarship in the Humanities*, 36(2), 430-438. <https://doi.org/10.1093/lc/fqaa032>.
- Prasad, R., Dinesh, N., Lee, A., Miltsakaki, E., Robaldo, L., Joshi, A., & Webber, B. (2008). The Penn discourse treebank 2.0. In *Proceedings of the Sixth International Language Resources and Evaluation (LREC'08)*. Marrakech, Morocco: European Language Resources Association.
- Radford, A. (2004). *Minimalist syntax: Exploring the structure of English*: Cambridge University Press.
- Refaeilzadeh, P., Tang, L., & Liu, H. (2009). Cross-validation. *Encyclopedia of database systems*, 5, 532-538. Boston: Springer. https://doi.org/10.1007/978-0-387-39940-9_565.
- Sa'edi, Gh. (1998). *Ashoftehalan-e Bidarbakht*. Tehran: Negah. (In Persian)
- Sa'edi, Gh. (2018). *Stranger in the town*. Tehran: Negah. (In Persian)
- Schütze, H., Manning, C. D., & Raghavan, P. (2008). *Introduction to information retrieval* (Vol. 39, pp. 234-265). Cambridge: Cambridge University Press.
- Segarra, S., Eisen, M., & Ribeiro, A. (2015). Authorship attribution through function word adjacency networks. *IEEE Transactions on Signal Processing*, 63(20), 5464-5478. <https://doi.org/10.1109/TSP.2015.2451111>.
- Shamsfard, M., & Bijankhan, M., (2022). *Text and speech processing for the Persian language: The state of the art and a brief review of the theoretical foundations*. Tehran: Samt. (In Persian)
- Shamsfard, M., Hesabi, A., Fadaei, H., Mansoori, N., Famian, A., Bagherbeigi, S., & Assi, S. M. (2010). Semi automatic development of farsnet: The persian wordnet. In *Proceedings of 5th global WordNet conference* (Vol. 29). India: Mumbai.
- Stamatatos, E. (2009). A survey of modern authorship attribution methods. *Journal of the American Society for information Science and Technology*, 60(3), 538-556. <https://doi.org/10.1002/asi.21001>
- Teich, E. (1995). *A proposal for dependency in systemic functional grammar—metasemiosis in computational systemic functional linguistics*. Ph.D. Doctoral dissertation, University of the Saarland and GMD/IPSI, Darmstadt, Germany.
- Uchendu, A., Le, T., Shu, K., & Lee, D. (2020). Authorship attribution for neural text generation. In *Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing* (pp. 8384-8395). (EMNLP 2020-2020 Conference on Empirical Methods in Natural Language Processing, Proceedings of the Conference). Association for Computational Linguistics (ACL).
- Weerasinghe, J., Singh, R., & Greenstadt, R. (2021). Feature vector difference based authorship verification for open-world settings. In *Proceedings of working of the evaluation Forum* (pp.2201-2207). Romania: Bucharest.