

روش‌های یادگیری ماشین در بررسی ویژگی‌های زبان شعری در اشعار شاعران دفاع مقدس^۱

(مطالعه موردی: اشعار دو شاعر دفاع مقدس؛ قیصر امین پور و محمدرضا عبدالملکیان)

کامیار جوانمردی^{۲*}؛ منوچهر اکبری^۳

تاریخ دریافت: ۱۳۹۷/۰۲/۲۸ تاریخ پذیرش: ۱۳۹۷/۰۵/۱۶

چکیده: شناسایی سبک نویسنده و پردازش زبان طبیعی از اهمیت ویژه‌ای برخوردار است و پژوهش در این زمینه به دو صورت کیفی و کمی انجام می‌شود. از آنجایی که شعر و ادبیات همواره یکی از افتخارات تاریخی زبان فارسی به‌شمار می‌آید، شناسایی سبک نویسندگان و شاعران به‌صورت گسترده، بدون اعمال نظرات انسانی و به شیوه کمی، امری ضروری است. در این پژوهش کاربرد روش‌های آماری برای شناسایی سبک نویسنده مورد بررسی قرار می‌گیرد؛ به‌همین دلیل، ویژگی‌های واژگانی، حرفی و نحوی از متن‌های منتخب استخراج می‌شود. هدف اصلی مقاله، استخراج ویژگی‌های صوری متن و دسته‌بندی اشعار مربوط به دو شاعر حوزه دفاع مقدس (قیصر امین پور و محمدرضا عبدالملکیان) است. به‌این منظور، از دو دسته‌بند K نزدیک‌ترین همسایه و بیز ساده جهت انجام عمل دسته‌بندی و انتساب داده‌ها استفاده شد. بررسی هر کدام از دسته‌بندها با استفاده از معیارهای ارزیابی، انجام شد. نتایج ارزیابی‌ها روی سه نوع ویژگی نشان داد، ویژگی‌های واژگانی بدون حذف واژه‌های دستوری در دسته‌بند بیز ساده با ۹۲ درصد معیار F از بالاترین دقت در میان انواع ویژگی‌ها برخوردار است. این نتیجه، نشان‌دهنده کارایی قابل ملاحظه این نوع ویژگی در شناسایی سبک نویسنده است.

واژگان اصلی: شناسایی نویسنده، پردازش متن، سبک شناسی، دسته‌بندی خودکار متن، شعر دفاع مقدس.

۱. این مقاله از طرح پژوهشی با همین عنوان بدست آمده که در پژوهشگاه علوم و معارف دفاع مقدس انجام شده است.

۲. کارشناسی ارشد دانشگاه صنعتی شریف، نویسنده مسئول (javanmardi69@yahoo.com).

۳. استاد زبان و ادبیات فارسی دانشگاه تهران.

مقدمه

شناسایی نویسنده^۱ و ویژگی‌های زبانی متن یکی از مسائل مهم در پردازش زبان طبیعی^۲ است. با استفاده از تجزیه و تحلیل رایانشی^۳، سبک ادبی نویسنده، روش‌های مختلف دسته‌بندی متن^۴ و یا نوشته نویسنده مشخص می‌شود. شناسایی نویسنده یکی از مسائل چالش‌برانگیز حوزه ادبیات و هم‌چنین، پردازش متن است. شناسایی نویسنده، در واقع نشان‌دادن برخی از ویژگی‌های خاص نویسنده در متن و نوشته است. با این عمل می‌توان نوشته‌های نویسنده را از دیگران متمایز نمود و یا نوشته‌هایی که نویسنده‌ای ناشناس دارد، به نویسنده‌ای مشخص منسوب کرد.

شناسایی نویسنده در ادبیات و سبک‌شناسی همواره با نظرات کارشناسان و صاحب‌نظران در حوزه ادبیات انجام می‌گیرد. این‌گونه تحلیل‌ها باعث ایجاد اختلاف میان صاحب‌نظران شده است. برای نمونه، می‌توان به نظرات متفاوت نسبت به برخی از اشعار خیام، حافظ و شاعران دیگر از ادبیات فارسی و هم‌چنین، نوشته‌های شکسپیر و بیکن از ادبیات انگلیسی اشاره کرد.

ایجاد تنوع در مطالعات سبک‌شناسی و ورود به عرصه پهن‌آور ادبیات با کمک علم آمار و ارقام، وسوسه‌ای جذاب و لذت‌بخش به نظر می‌رسد. نگاشت احساسات شاعرانه، مضامین و مفاهیم پیچیده ادبی به ارقام، اعداد و بررسی عینی‌تر آن‌ها و ساخت پل ارتباطی میان این دو حوزه از اهداف اصلی این پژوهش است. به‌ویژه در عرصه زبان فارسی که به نظر مطالعات اندکی در آن انجام شده است. شعر معاصر به دلیل قابل فهم بودن، قرابت فکری و زبانی از لحاظ واژگان و ساختار جمله به زبان معیار کنونی و شعر نو، انتخاب مناسبی به نظر می‌رسد و اهداف پژوهش را با وضوح بالاتری نشان می‌داد.

شناسایی نویسنده در حوزه ادبیات، هوش مصنوعی و حقوق قضایی اهمیت ویژه‌ای دارد و در شاخه بازیابی اطلاعات^۵ و زبان‌شناسی رایانشی^۶ قرار می‌گیرد. شناسایی سبک و سیاق نویسنده متون ادبی، رسمی و حتی نامه‌های تهدیدآمیز با استفاده از ویژگی‌های کمی، بسیار کاربردی است.

1. Author identification
2. Natural language processing
3. Computational analysis
4. Text classification
5. Information retrieval
6. Computational linguistics

«هر چه قدر حرف نو زدن در ادبیات مشکل است، ورود کامپیوتر به این عرصه از آن هم مشکل‌تر است. با وجود پیشرفت‌های فراوان در تحقیقات کامپیوتری (به‌طوری‌که، محققان در برخی زبان‌ها همانند انگلیسی و اسپانیایی موفق به سرایش شعر به وسیله کامپیوتر شده‌اند)، هنوز اساتید ادبیات و ادیبان ما حاضر به قبول توانایی‌های این ابزار نیستند. از طرفی اساتید کامپیوتر اطلاع چندانی با مشکلات ادبیات ندارند و حاضر نیستند کار نوینی انجام دهند. باتوجه‌به اینکه تاکنون پژوهش‌های اندکی در این زمینه انجام شده، ورود به این عرصه امری ضروری به‌نظرمی‌رسد» (مجیری و مینایی، ۱۳۸۷: ۲).

مسئله اصلی این پژوهش، کاربرد شیوه‌های آماری روش‌های یادگیری ماشین در بررسی سطح زبانی شعر است. شناسایی سبک دو شاعر دفاع مقدس قیصر امین‌پور و محمدرضا عبدالملکیان با استفاده از روش‌های آماری به عنوان مطالعه موردی جهت نشان‌دادن این مسئله انتخاب شده است. انتخاب این دو شاعر به دلیل ویژگی‌های زبانی و شعری آن‌ها، برخورداری از سبک ویژه و شاخص بودن اشعار آنان در زمینه شعر جنگ و به‌ویژه شعر دفاع مقدس بوده است.

مهم‌ترین مرحله انجام این کار و اصلی‌ترین گام برای آغاز تحلیل آماری، سبک‌سنجی گزینش و تعریف ویژگی‌ها است. این ویژگی‌ها در شناسایی سبک نویسنده یا شاعر، شرایطی را نشان خواهد داد که قابلیت کمی‌سازی و تحلیل آماری را فراهم می‌کند. شاعرانی که در این پژوهش مورد بررسی قرار گرفته‌اند، آثارشان مشخص و در دسترس است. بخشی از این اشعار پس از جمع‌آوری، به‌عنوان اسناد آموزش^۱ و بخشی دیگر برای سنجش دقت سامانه نهایی این پژوهش استفاده می‌شوند. داده‌هایی که در بخش آموزش قرار ندارند، اسناد آزمون^۲ یا آزمایش نامیده می‌شوند. درنهایت در این پژوهش با استفاده از ویژگی‌های واژگانی، نحوی، حرفی و بهره‌گیری از دو دسته‌بند K نزدیک‌ترین همسایه^۳ و بیز ساده^۴ سبک دو شاعر حوزه دفاع مقدس، قیصر امین‌پور و محمدرضا عبدالملکیان از یکدیگر متمایز خواهد شد.

1. Training document
2. Test document
3. K-nearest neighbors
4. Naïve Bayes

در طراحی و پیاده‌سازی تمام بخش‌ها از جمله دسته‌بندها و برچسب‌زن نقش نحوی کلمات، از ابزارهای آماده استفاده نشده است. همچنین، در تمامی مراحل پژوهش، زبان برنامه‌نویسی پایتون^۱ برای طراحی و پیاده‌سازی برنامه‌ها کاربرد داشته است. این امر باعث پیاده‌سازی سطوح مختلف پردازش زبانی مستقل می‌شود تا در پژوهش‌های آینده، قابلیت بهره‌برداری داشته و آمارهای پیش‌بینی‌نشده، امکان محاسبه و بررسی بیشتری داشته باشند.

پیشینه پژوهش

نخستین پژوهش در زمینه شناسایی مؤلف به سال ۱۷۱۱ میلادی بازمی‌گردد که روی متن کتاب مقدس انجام شده است. کشیش آلمانی به نام اچ. بی. ویترا^۲ در سال ۱۷۱۱ میلادی به این نکته اشاره کرد، نام‌های متفاوتی از خداوند در کتب پنج‌گانه وجود دارد که نویسندگان متعددی آن‌ها را نوشته‌اند.

پژوهش‌های شناسایی نویسنده متن با استفاده از روش‌های نوین آماری و رایانشی، در اواخر قرن نوزدهم آغاز شده است (Olsson, 2004: 3). «تلاش‌های اولیه جهت سنجش سبک نوشتاری نویسنده به قرن نوزدهم بازمی‌گردد. این مطالعات توسط مندهال در سال ۱۸۸۷ روی نمایش‌نامه‌های شکسپیر انجام شد. سپس، مطالعات آماری توسط زیف در سال ۱۹۳۲ و یول در سال‌های میانی قرن بیستم انجام شد. مطالعات دقیق‌تر توسط ماستلر و والاس در سال ۱۹۶۴ و روی پایگاه داده فدرالیست پیپر^۳ انجام شد که بی‌شک، یکی از قدرتمندترین و مؤثرترین تحقیقات در تعیین هویت نویسنده بوده است. این پژوهش‌ها، آغاز مطالعه غیرستی تشخیص هویت نویسنده محسوب می‌شود» (فرهمندپور و همکاران، ۱۳۹۱: ۲۷).

در مطالعه‌ای (Abbasi & Chen, 2005: 68) محتوای سایت‌ها و پیام‌های گروه‌های تروریستی برای شناسایی خودکار، بررسی شده است. در این تحقیق از تحلیل پیام‌های آنلاین موجود در وب‌سایت‌ها و انجمن‌های گفتگوی^۴ عربی و انگلیسی و چارچوب‌های موجود برای

1. Python programming language
2. H. B. Witter
3. The Federalist Paper
4. Web forum

شناسایی نویسنده استفاده شده است. در این پژوهش ویژگی‌هایی از جمله طول واژه، حذف و اضافه، علائم سجاوندی، ریشه‌ واژگان، ساختار واژه و محتوا بررسی شده است.

پژوهش دیگر (Baayen et al., 2002: 2) به زبان هلندی انجام گرفته است. در این پژوهش از هشت دانشجوی هلندی رشته ادبیات خواسته شده تا هر یک، سه اثر در سبک‌های تخیلی، تحلیلی و توصیفی با طول هزار واژه بنویسند. در نهایت، نه اثر از هر نویسنده و در کل از ۷۲ یادداشت به‌عنوان پیکره پژوهش استفاده شده است. دقت این کار در صورتی که ۴۲ واژه دستوری و ۸ علامت سجاوندی لحاظ شده باشد، ۸۳/۵ درصد است و اگر ۵۰ واژه دستوری و ۸ علامت سجاوندی در نظر گرفته شود، دقت به ۸۸/۱ درصد خواهد رسید.

در پژوهش دیگری (Patton & Can, 2004: 468) چهار داستان از یک نویسنده بررسی شده است. این تحلیل با استفاده از ویژگی‌های پرکاربردترین واژه‌ها، شمارش هجاها، اطلاعات نوع کلمه و برچسب نحوی، طول جملات برحسب واژه‌ها، طول واژه در متن و طول واژه در واژگان، انجام گرفته است. در این بررسی، هر داستان در قالب پنج‌هزار واژه قرار گرفته و هر یک از ویژگی‌های مذکور با توجه به پیکره به‌دست آمده، شمارش شده است. براساس نتایج، دو رمان اول از همدیگر متمایز و دو رمان آخر درهم تنیده شده‌اند. بهترین ویژگی‌هایی که در این پژوهش براساس جداسازی و تفکیک میان متون به‌دست آمد، پرکاربردترین واژه‌ها و طول جمله است. این دو ویژگی باعث دسته‌بندی متون به ترتیب با دقت ۸۷ و ۸۱ درصد شد.

در پژوهش‌هایی، روش‌های آماری و نحوی در ترکیب با یکدیگر استفاده شده‌اند. در پژوهشی (Mechti et al., 2015: 2) روش‌های آماری و سبکی در الگوریتم‌های یادگیری ماشین ترکیب شده است. مهدی بازرگان، از نخستین افرادی است که به ویژگی‌های آماری در متون فارسی توجه کرد. وی با استفاده از «جمله‌نگار» که نمایش هندسی از تعداد واژه‌های تشکیل‌دهنده جملات و تعداد جملات است به بررسی دیباچه گلستان پرداخت. نتایج این پژوهش نشان می‌دهد، شاخص‌های اصلی جمله‌نگار در بخش‌های دیگر گلستان کم‌وبیش همانند دیباچه است؛ اما، جمله‌نگار مقدمه گلستان که نوشته محمدعلی فروغی است، چهره متفاوتی دارد (استاجی، ۱۳۸۷: ۱۷).

در پژوهشی شاهمیری و بروجردی (۱۳۸۵) سه ویژگی فیزیکی، آوایی و مفهومی را از اشعار خیام، حافظ، فردوسی و مولوی استخراج کرده‌اند. سپس، با استفاده از شبکه عصبی مصنوعی و درخت تصمیم‌گیری دسته‌بندی و شناسایی انجام داده‌اند. در مطالعه دیگری (آذین، ۱۳۹۲) اشعار نیما یوشیج، مهدی اخوان‌ثالث، احمد شاملو و سهراب سپهری با استفاده از روش‌های یادگیری ماشین، بررسی شده است. این روش‌ها در تحلیل مستقل بوده و از معیارهای انسانی تأثیر نمی‌گیرند.

چارچوب نظری

شناسایی سبک نویسنده از مسائل چالش‌برانگیز حوزه ادبیات و پردازش متن است. «واژه سبک، عربی و از ریشه «سَبَك» به معنی «ذوب کردن و در قالب ریختن» است. سابقه کاربرد اصطلاح سبک در زبان عربی به سال‌های سوم هجری قمری در آثار جاحظ و ابوهلال عسکری می‌رسد. در ادبیات فارسی و عربی واژه‌های طرز، روش، سیاق و نمط مترادف سبک است. در زبان فارسی اصطلاح سبک‌شناسی را نخستین بار محمدتقی بهار برای کتاب خود «سبک‌شناسی نثر» استفاده کرد. در زبان‌های لاتین «استیل»^۱ به معنی سبک است که قلمی نوک تیز بوده و با آن حروف را روی صفحه‌های سخت می‌تراشیدند. بعدها این واژه به‌طور مجازی به روش خاص شکل دادن به حروف در نوشتار تعبیر شده است. هر اثر ادبی همانند هر انسانی دارای مشخصاتی باشد؛ از این رو، در پاره‌ای از صفات و خصال با آثار دیگر شریک است» (سمیعی گیلانی، ۱۳۸۶: ۵۵).

بررسی و مطالعه سبک، از موضوع‌های نوین در حوزه زبان‌شناسی است. در سبک‌شناسی، دانش گونه‌های کاربردی زبان به عنوان شیوه عمل تعریف می‌شود؛ اما، واژه سبک در کاربرد جاری زبان، بر معنا و مفاهیم متعددی همانند سبک چیدمان منزل و سبک ورزشی دلالت دارد. همچنین در مواردی که هنرمند به درجه‌ای از تسلط دست می‌یابد، اذعان می‌شود که وی در هنر خود صاحب سبک است (سمیعی گیلانی، ۱۳۸۴: ۸۷).

تاکنون عموم پژوهش‌ها در حوزه سبک‌شناسی به شیوه کیفی انجام شده است. در این روش، پژوهشگر با در نظر گرفتن دانش و درک خود از اثر مورد مطالعه و مکاتب ادبی، به

نتیجه‌گیری می‌رسد. «در این شیوه نتایج به دست آمده معمولاً موردی و غیر قابل تعمیم است؛ زیرا برخلاف آزمایش‌هایی که در رشته‌های تجربی انجام می‌شود، این مطالعات نمی‌تواند مجدداً توسط افراد و شرایط گوناگون آزمایش شود. از این رو دانش زبان‌شناسی که تمایل بسیاری به شیوه‌های علمی و تجربی دارد، به دنبال شیوه‌های کمی و قابل تکرار بوده است. این امر منجر به ظهور حوزه مطالعات سبک‌شناسی و تشخیص نویسنده شده است» (یاحقی و ایزانلو، ۱۳۸۵: ۱۵۳).

در این پژوهش، شناسایی نویسنده در قالب یک مسئله دسته‌بندی بررسی شده است. برای انجام بررسی‌ها، فضای مسئله به دسته‌هایی به تعداد نویسندگانی که در داده‌های آموزشی هستند، تقسیم می‌شود. روش‌های مختلفی برای دسته‌بندی اطلاعات وجود دارد که در این پژوهش به دو دسته‌بندی اشاره می‌شود.

دسته‌بند بیز ساده^۱

در زبان طبیعی، ارجاع سند D به مجموعه‌های از پیش تعریف شده $C = \{c_1, c_2, \dots, c_n\}$ اطلاق می‌شود. به طور معمول، ساختار آموزشی با نظارت^۲ برای آموزش دسته‌بند استفاده می‌شود. الگوریتم آموزشی، مجموعه‌ای از N نمونه برچسب خورده آموزشی $\{(d_i, c_j) : i, j = 1, \dots, N\}$ فراهم می‌کند که با استفاده از آن یک تابع دسته‌بندی ایجاد می‌شود و اسناد را به دسته‌های نگاشت می‌دهد. d_i سند i و c_j دسته j است که d_i به آن نگاشت می‌شود. در این رابطه D و C متغیرهای تصادفی برای ارزش‌های عددی سند و دسته هستند (آذین، ۱۳۹۲: ۱۳). رایج‌ترین الگوریتم آموزش برای دسته‌بندی متون براساس قانون بیز (رابطه شماره ۱) است.

رابطه شماره ۱:

$$P(C = c | D = d) = \frac{P(C = c) \cdot P(D = d | C = c)}{P(D = d)}$$

دسته‌بند K نزدیک‌ترین همسایه

1. Naive Bayes Classifiers
2. Supervised learning framework

این دسته‌بند به دلیل سادگی، کاربرد بسیاری در طبقه‌بندی متون دارد. مبنای کار این الگوریتم، مقایسه متن آزمایش شده با متون آموزش داده شده و در نهایت تعیین میزان شباهت میان آن‌ها است. متون آموزشی با n ویژگی به عنوان یک نقطه در فضای n بعدی نمایش داده می‌شود. هنگامی که متن ناشناخته‌ای الگوریتم دریافت می‌کند، فضای الگو برای یافتن متون آموزشی که مشابه متن ناشناخته باشد، جستجو می‌کند (سرایی و شاهقلیان، ۱۳۸۹).

«در این شیوه، نمونه‌ها در یک فضای n -بعدی به صورت نقاطی متمایز قرار می‌گیرند و دسته‌بندی با استفاده از اندازه‌گیری و مقایسه فاصله اقلیدسی^۱ نقاط مختلف روی بردار ویژگی‌ها انجام می‌گیرد. در مدل‌های برداری^۲ که از روش K نزدیک‌ترین همسایه استفاده می‌شود، دو مقیاس اندازه‌گیری فاصله، یعنی فاصله اقلیدسی و کسینوسی^۳ کاربرد دارد» (آذین، ۱۳۹۲: ۱۳).

معیارهای ارزیابی

دسته‌بندی متون با استفاده از مقایسه داده‌های پیش‌بینی شده و مقادیر واقعی، ارزیابی و سنجش می‌شود. جدول شماره ۱ ماتریس دسته‌بندی ساده‌ای را برای مسئله دسته‌بندی نشان می‌دهد.

جدول شماره ۱: دسته‌بندی دو حالت

| | | | |
|---|--------|------------------------|------------------------|
| مقدار خروجی سامانه دسته پیش‌بینی شده ^۴ | | | |
| درست (مثبت) | | نادرست (منفی) | |
| مقدار واقعی ^۵ | درست | مثبت صحیح ^۶ | منفی غلط ^۷ |
| | نادرست | مثبت غلط ^۸ | منفی صحیح ^۹ |

1. Euclidean distance
2. Vector models
3. Cosine similarity
4. Predicted Class
5. False Negative
6. True Positive
7. True Class
8. True Negative
9. False Positive

در جدول شماره ۱، مقادیر موجود در خانه‌های مثبت و منفی مبین دقت و پیش‌بینی صحیح سامانه است. همچنین، مقادیر موجود در خانه‌های مثبت و منفی غلط، مبین خطای دسته‌بند و پیش‌بینی اشتباه است. معیارهای دقت، بازخوانی و معیار F براساس جدول شماره ۲ به‌دست می‌آیند.

جدول شماره ۲: معیارهای ارزیابی

| فرمول | معیار |
|---|----------|
| $\frac{\text{مثبت صحیح}}{\text{مثبت غلط} + \text{مثبت صحیح}}$ | دقت |
| $\frac{\text{مثبت صحیح}}{\text{منفی غلط} + \text{مثبت صحیح}}$ | بازخوانی |
| $\frac{\text{دقت} \times \text{بازخوانی} \times 2}{\text{دقت} + \text{بازخوانی}}$ | معیار F |

داده‌ها و مباحث تفصیلی تحقیق

در این پژوهش، ویژگی‌های صوری در سه سطح واژگانی، نحوی و حرفی دو شاعر به صورت کمی بررسی شده است؛ از این رو، هر شعر به‌عنوان یک سند و مجموعه اسناد به‌عنوان پیکره^۱ سامانه پنداشته شده است. برای محاسبه این ویژگی‌ها از زبان برنامه‌نویسی پایتون و ابزار پردازش زبان طبیعی^۲ استفاده شده است. ویژگی‌های کمی محاسبه شده، به‌عنوان ویژگی‌های دسته‌بندها در امر دسته‌بندی متون لحاظ شده است. همچنین، دو دسته‌بند K نزدیکترین همسایه و بیز ساده کاربرد داشته است که این دسته‌بندها عمل انتساب متون را به شاعران موردنظر انجام می‌دهند.

پیش‌از آنکه پیکره‌ای برای انجام پردازش استفاده شود، باید ابتدا روی داده‌ها پیش‌پردازش انجام گیرد تا صورت‌های غیراستاندارد موجود در پیکره به شکل استاندارد تبدیل شوند؛ زیرا،

1. Corpus
2. Natural language toolkit

چنانچه حروف، نشانه‌های نگارشی و کلمات فارسی به شکل یکسانی نوشته نشوند، متون مورداستفاده توسط سامانه‌های رایانه‌ای قابل تحلیل نخواهند بود.

جدول شماره ۳: تعداد کل اشعار پیکره شاعران

| شاعر | تعداد اشعار |
|-------------|-------------|
| امین پور | ۱۲۲ |
| عبدالملکیان | ۳۷۶ |

پیکره کاربردی این پژوهش، متن اشعار قیصر امین پور و محمدرضا عبدالملکیان است که تعداد آن‌ها در جدول شماره ۳ نشان داده شد. تمام اشعار به صورت دستی تایپ و هر شعر در قالب یک سند و مجموع اشعار به عنوان کل پیکره برای استخراج ویژگی‌ها و دسته‌بندی به کار گرفته شد. کتاب‌های شعر قیصر امین پور که در این مطالعه استفاده شد، کتاب‌های از آینه‌های ناگهان دفتر اول، از آینه‌های ناگهان دفتر دوم، از تنفس صبح، بی‌بال پریدن و دستور زبان عشق است. همچنین، از آثار محمدرضا عبدالملکیان دفتر شعرهای ساده با تو حرف می‌زنم، هوای حوصله ابری‌ست، دنیا به کبوترانش پشت کرده است و پل خواب انتخاب شد. هر کدام از کتاب‌های شعر این دو شاعر شامل چندین شعر و نوشته است. در جدول شماره ۳، تعداد اشعار و در جدول شماره ۴ تعداد کل کلمات اشعار هر شاعر نمایش داده شده است.

جدول شماره ۴: تعداد کل کلمات پیکره اشعار شاعران

| شاعر | تعداد کل کلمات |
|-------------|----------------|
| امین پور | ۱۴۵۱۳ |
| عبدالملکیان | ۱۸۲۵۵ |

ویژگی‌های واژگانی

بردار ویژگی‌های واژگانی که شامل تعداد کلمه در هر سند یا شعر است، به دست می‌آید. در واقع، هزار کلمه پراهمیت از لحاظ معیار $tf-idf$ به عنوان بردار ویژگی‌های واژگانی در نظر گرفته می‌شود.

الگوی وزن‌دهی بر مبنای بسامد کلمات و خواص تکرار کلمات در مجموعه اسناد تعیین می‌شود. وزن‌های محلی، تابعی از تکرار هر کلمه در یک سند است. همچنین، کلمه‌ای که بیشترین تکرار را داشته باشد، در زمره کلمه‌های مهم سند است (کامیار، ۱۳۹۰). پنج کلمه از پراهمیت‌ترین کلمه‌ها برای بررسی واژگانی اشعار در جدول شماره ۵ نمایش داده شده است.

جدول شماره ۵: واژه‌های پر کاربرد پیکره اشعار بر اساس معیار $tf-idf$

| معیار $tf-idf$ | کلمه | ترتیب اهمیت |
|----------------|------|-------------|
| ۰/۰۰۴۴۳۹۹۶ | با | ۱ |
| ۰/۰۰۴۳۷۱۴۵ | حالا | ۲ |
| ۰/۰۰۴۳۴۱۲۰ | من | ۳ |
| ۰/۰۰۴۲۶۵۶۲ | نه | ۴ |
| ۰/۰۰۴۲۲۲۷۷ | هم | ۵ |

یکی از روش‌های ساده برای نشان‌دادن متون، استفاده از بردار وقوع کلمات است. در بسیاری از پژوهش‌های شناسایی نویسنده، برای مطالعه و شناخت سبک وی از ویژگی‌های واژگانی استفاده می‌شود. به این شیوه در حوزه بازبایی اطلاعات و دسته‌بندی متون مبتنی بر موضوع^۱، کیسه لغات^۲ می‌گویند. در این روش، متن به عنوان مجموعه‌ای از کلمات در نظر گرفته می‌شود که هر کلمه تعداد وقوع یا بسامدی دارد و مستقل از اطلاعات بافتی^۴ است.

1. Term frequency-inverse document frequency
2. Topic-based text classification
3. Bag of words
4. Contextual information

باین‌حال، تمایزهایی میان شناسایی سبک نویسنده با روش‌های بازیابی اطلاعات و دسته‌بندی متون مبتنی بر موضوع وجود دارد. کلمه‌هایی که تکرار بسیاری دارند (مانند حروف تعریف، حروف اضافه، ضمایر و حروف ربط)، از ویژگی‌های مفید در شناسایی سبک نویسنده به‌شمار می‌آیند؛ اما، این نوع کلمات در روش‌های دسته‌بندی متون براساس موضوع، معمولاً کلمات زائد در نظر گرفته می‌شوند و در همان ابتدای کار به‌عنوان یکی از مراحل پیش‌پردازش از متن حذف می‌شوند. این نوع کلمات را ایست‌واژه^۱ می‌نامند و از آنجایی که فاقد اطلاعات معنایی هستند، در محاسبات لحاظ نمی‌شوند.

کلمات دستوری یا همان ایست‌واژه‌ها به‌وفور و به‌صورت ناخودآگاه، فارغ از موضوع و بافت‌های گوناگون استفاده می‌شوند؛ بنابراین، این نوع کلمات در موضوع‌های گوناگون می‌توانند عامل مهمی در شناسایی سبک نویسنده باشند (Stamatatos, 2009: 541).

پس از تعیین ویژگی‌ها و تشکیل بردار برای هر یک، با استفاده از دو دسته‌بند به‌انتساب اسناد آزمایش به شاعران پرداخته می‌شود. برای آزمایش دسته‌بند بیز ساده، ۱۵ درصد داده‌ها به‌عنوان داده آزمون استفاده می‌شود و مابقی به‌صورت داده آموزش در دسته‌بند مورد استفاده قرار می‌گیرد. جدول شماره ۶ نمایان‌گر تعداد اسناد (اشعار)، داده‌های بخش‌های آموزش و آزمون برای هر دسته (شاعر) است.

جدول شماره ۶: تعداد اشعار انتخاب شده برای بخش آموزش و آزمون

| داده‌های آموزش | داده‌های آزمون | |
|----------------|----------------|-------------|
| ۱۰۴ | ۱۸ | امین‌پور |
| ۳۳۰ | ۴۶ | عبدالملکیان |
| ۴۳۴ | ۶۴ | مجموع |

پس از اجرای مرحله آموزش، داده‌های آزمون در دسته‌بند قرار می‌گیرند. درنهایت، دسته‌بند بیز ساده با ویژگی‌های واژگانی آموزش داده می‌شود. سپس، نتایج ارزیابی این روش با

1. Stop word

معیارهای دقت، بازخوانی و معیار F (جدول شماره ۷) و همچنین، نتایج ماتریس درهم‌ریختگی^۱ (جدول شماره ۸) اندازه‌گیری می‌شود.

به‌منظور بهبود هر الگوریتمی، نیاز است تا این مدل بررسی شود و نقاطی که دچار خطا است، شناسایی شود. تحلیل خطا در دسته‌بندیها از طریق ماتریس درهم‌ریختگی یا جدول وابستگی^۲ انجام می‌گیرد. ماتریس درهم‌ریختگی برای دسته‌بندی با N دسته، یک ماتریس N در N است. در این ماتریکس، درایه x, y شامل تعداد دفعات داده‌ای است که در دسته x قرار دارد و توسط مدل، به‌عنوان داده‌ای از دسته y قرار می‌گیرد.

جدول شماره ۷: نتایج ارزیابی دسته‌بند بیز ساده بدون حذف ایست‌واژه‌ها روی ویژگی‌های واژگانی

| معیار F | دقت | بازخوانی | |
|---------|-------|----------|-------------|
| ۸۹/۴۷ | ۸۵/۰۰ | ۹۴/۴۴ | امین پور |
| ۹۵/۵۵ | ۹۷/۷۲ | ۹۳/۴۷ | عبدالملکیان |
| ۹۲/۵۱ | ۹۱/۳۶ | ۹۳/۹۶ | میانگین |

جدول شماره ۸: نتایج ماتریس درهم‌ریختگی دسته‌بند بیز ساده بدون حذف ایست‌واژه‌ها روی ویژگی‌های واژگانی

| | | |
|----------|-------------|-------------|
| امین پور | عبدالملکیان | |
| ۳ | ۴۳ | عبدالملکیان |
| ۱۶ | ۲ | امین پور |

در مثال دسته‌بند بیز ساده، ویژگی‌های واژه‌ها، جداول و اعداد درون درایه‌های قطر اصلی ماتریس بدین معنا است که اسناد توسط دسته‌بند درست پیش‌بینی شده‌اند. برای نمونه، از میان ۴۶ سند عبدالملکیان برای آزمایش در دسته‌بندی، ۴۳ مورد به‌درستی تشخیص داده شده‌اند.

1. Confusion Matrix
2. Contingency Table

همچنین، سه نمونه دیگر به عنوان اشعار امین پور دسته‌بندی شده‌اند. تحلیل خطا به این صورت، بخش بسیار مهمی از انواع کاربرد زبان‌شناسی رایانشی است و می‌تواند مبین مشکلاتی در داده‌های آموزش باشد. مهم‌تر از همه، تحلیل خطا در توسعه انواع دانش و الگوریتم‌ها برای استفاده در حل مسائل مفید است (Jurafsky & Martin, 2006).

در گام بعدی، تمام ایست‌واژه‌ها از میان اشعار حذف شده‌اند و بار دیگر دسته‌بندی بیز ساده روی داده‌ها، آزمایش شد. برخی از ایست‌واژه‌ها در جدول شماره ۹ نشان داده شده است. سپس، نتایج ارزیابی دسته‌بند و ماتریس درهم‌ریختگی بعد از حذف ایست‌واژه‌ها در جداول شماره ۱۰ و ۱۱ درج شد.

جدول شماره ۹: برخی از انواع ایست‌واژه‌های زبان فارسی

| ضمیر | حرف اضافه | حرف ربط | فعل |
|------|-----------|---------|-------|
| من | از | و | است |
| ما | برای | اما | بود |
| تو | به | که | شد |
| همین | را | پس | هستند |

جدول شماره ۱۰: ارزیابی دسته‌بند بیز ساده پس از حذف ایست‌واژه‌ها روی ویژگی‌های واژگانی

| معیار F | دقت | بازخوانی | |
|---------|-------|----------|-------------|
| ۷۷/۲۷ | ۶۵/۳۸ | ۹۴/۴۴ | امین پور |
| ۸۸/۰۹ | ۹۷/۳۶ | ۸۰/۴۳ | عبدالملکیان |
| ۸۲/۶۸ | ۸۱/۳۷ | ۸۷/۴۳ | میانگین |

جدول شماره ۱۱: ماتریس درهم‌ریختگی دسته‌بند بیز ساده پس از حذف ایست‌واژه‌ها روی ویژگی‌های واژگانی

| | | |
|----------|-------------|-------------|
| امین پور | عبدالملکیان | |
| ۹ | ۳۷ | عبدالملکیان |
| ۱۷ | ۱ | امین پور |

نتایج به‌دست آمده، پایین‌بودن دقت پیش‌بینی دسته‌بند را نشان می‌دهد. این امر مؤثر بودن عدم حذف ایست‌واژه‌ها و واژه‌های دستوری در دسته‌بندی شناسایی سبک نویسنده را ثابت می‌کند؛ از این رو، حذف این واژه‌ها تأثیر منفی در دسته‌بندی اسناد دارد.

سپس، دسته‌بند K نزدیک‌ترین همسایه برای دسته‌بندی اسناد استفاده شده است. آزمایش‌های مختلف با مقادیر متفاوت K روی اسناد آموزش انجام گرفت. در این مرحله، جهت انتخاب K، مقادیر متفاوتی آزمایش شد که در نهایت مقدار آن هشت در نظر گرفته شد. نتایج ارزیابی این دسته‌بند با به‌کارگیری فاصله کسینوسی و اقلیدسی به ترتیب در جداول شماره ۱۲ و ۱۳ نشان داده شده است.

جدول شماره ۱۲: دسته‌بند K نزدیک‌ترین همسایه با استفاده از فاصله کسینوسی روی ویژگی‌های واژگانی

| معیار F | دقت | بازخوانی | |
|---------|-------|----------|-------------|
| ۷۸/۱۶ | ۶۶/۰۱ | ۹۵/۷۷ | امین پور |
| ۹۴/۵۰ | ۹۹/۰۹ | ۹۰/۳۳ | عبدالملکیان |
| ۸۶/۳۳ | ۸۲/۵۵ | ۹۳/۰۵ | میانگین |

جدول شماره ۱۳: دسته‌بند K نزدیکترین همسایه با استفاده از فاصله اقلیدسی روی ویژگی‌های واژگانی

| معیار F | دقت | بازخوانی | |
|---------|-------|----------|-------------|
| ۶۹/۷۶ | ۶۰/۰۰ | ۸۳/۳۳ | امین پور |
| ۸۴/۷۰ | ۹۲/۳۰ | ۷۸/۲۶ | عبدالملکیان |
| ۷۷/۲۳ | ۷۶/۱۵ | ۸۰/۷۹ | میانگین |

همان‌گونه که در جداول فوق مشاهده می‌شود، نتایج دسته‌بند K نزدیک‌ترین همسایه با استفاده از دو فاصله اقلیدسی و کسینوسی تفاوت نسبتاً قابل توجهی در دسته‌بندی نشان می‌دهند؛ با این حال، نتایج با داشتن ویژگی‌های واژگانی قابل قبول است. همچنین، استفاده از فاصله کسینوسی، حدود ۱۰ درصد بر کارایی این دسته‌بند افزوده است.

ویژگی‌های نحوی

در این مرحله، برداری با صد ویژگی تشکیل شد که تمامی آن‌ها، دونگاشت‌های نحوی متداول در هر سند را نشان می‌دهند. در جدول شماره ۱۴، پنج دونگاشت رایج به همراه بسامدشان نمایش داده شده است.

جدول شماره ۱۴: پنج دونگاشت پرتکرار از نقش اجزای کلام در پیکره اشعار

| دونغاشت نقش اجزای کلام | وقوع در پیکره | |
|------------------------|---------------|---|
| N_N | ۶۲۱۷ | ۱ |
| P_N | ۳۵۰۸ | ۲ |
| PUNC_N | ۲۸۷۶ | ۳ |
| N_PUNC | ۲۷۰۹ | ۴ |
| CON_N | ۲۵۱۱ | ۵ |

در استخراج دونگاشت‌های نحوی، دقت برچسب‌زنی از اهمیت بالایی برخوردار است. به طوری که، هرچه نتایج به دست آمده از برچسب‌زن صحت بیشتری داشته باشد، کارایی ویژگی‌های نحوی در سامانه افزایش می‌یابد. در این پژوهش برای برچسب‌زنی از پیکره دو میلیون بی‌جن‌خان به عنوان داده آموزش استفاده شده است. این برچسب‌زنی نحوی با استفاده از مدل مخفی مارکوف^۱ طراحی شده است. دقت این برچسب‌زنی با معیار F، ۸۹ درصد است. در پیکره بی‌جن‌خان علائم و برچسب‌هایی استفاده شده‌اند که در جدول شماره ۱۵ نمایش داده شده است (بی‌جن‌خان، ۱۳۸۱).

جدول شماره ۱۵: علائم و برچسب‌های نحوی نقش اجزای کلام

| برچسب فارسی | برچسب انگلیسی | |
|----------------|---------------|---|
| اسم | N | ۱ |
| اسم جمع | N_PL | ۲ |
| اسم مفرد خاص | N_SING_NPR | ۳ |
| حرف اضافه | P | ۴ |
| گروه حرف اضافه | PP | ۵ |

پس از استخراج ویژگی‌ها و تشکیل بردار، انتساب اسناد آزمایش به شاعران با استفاده از دسته‌بندها، مورد مطالعه قرار گرفت. در این بخش ویژگی‌های واژگانی به روش K نزدیک‌ترین همسایه، برای آزمایش ۱۵ درصد از داده‌های کلی جدا شد و با استفاده از بقیه اسناد به آموزش دسته‌بند پرداخته شد. نتایج ارزیابی دسته‌بندی به روش K نزدیک‌ترین همسایه، با داشتن بردار ویژگی‌های نحوی و کاربرد دو فاصله اقلیدسی و کسینوسی در جداول شماره ۱۶ و ۱۷ نمایش داده شده است.

1. Hidden Markov Model

جدول شماره ۱۶: نتایج ارزیابی دسته‌بند K نزدیک‌ترین همسایه با استفاده از فاصله اقلیدسی روی ویژگی‌های نحوی

| معیار F | دقت | بازخوانی | |
|---------|-------|----------|-------------|
| ۶۷/۵۶ | ۷۲/۸۱ | ۶۳/۰۲ | امین پور |
| ۸۸/۸۸ | ۸۶/۶۶ | ۹۱/۰۸ | عبدالملکیان |
| ۷۸/۱۹ | ۷۹/۷۴ | ۷۷/۰۵ | میانگین |

جدول شماره ۱۷: نتایج ارزیابی دسته‌بند K نزدیک‌ترین همسایه با استفاده از فاصله کسینوسی روی ویژگی‌های نحوی

| معیار F | دقت | بازخوانی | |
|---------|--------|----------|-------------|
| ۸۳/۷۲ | ۷۲/۰۰ | ۱۰۰/۰۰ | امین پور |
| ۹۱/۷۶ | ۱۰۰/۰۰ | ۸۴/۷۸ | عبدالملکیان |
| ۸۷/۷۴ | ۸۶/۰۰ | ۹۲/۳۹ | میانگین |

سپس، در مرحله بعد، بردار ویژگی‌های نحوی با استفاده از دسته‌بند بیز ساده انجام شد که نتایج ارزیابی این روش و ماتریس ابهام به ترتیب در جداول شماره ۱۸ و ۱۹ نمایش داده شده است. باتوجه به نتایج به دست آمده در این بخش، به نظر می‌رسد استفاده از ویژگی‌های نحوی در حوزه شناسایی نویسنده منجر به دقت مناسبی در دسته‌بندی و انتساب متن می‌شود.

جدول شماره ۱۸: نتایج ارزیابی دسته‌بند بیز ساده روی ویژگی‌های نحوی

| معیار F | دقت | بازخوانی | |
|---------|-------|----------|-------------|
| ۸۲/۹۲ | ۷۳/۹۱ | ۹۴/۴۴ | امین پور |
| ۹۱/۹۵ | ۹۷/۵۶ | ۸۶/۹۵ | عبدالملکیان |
| ۷۴/۴۴ | ۸۵/۷۳ | ۹۰/۷ | میانگین |

جدول شماره ۱۹: نتایج ماتریس درهم‌ریختگی دسته‌بند بیز ساده روی ویژگی‌های نحوی

| | | |
|----------|-------------|-------------|
| امین پور | عبدالملکیان | |
| ۶ | ۴۰ | عبدالملکیان |
| ۱۷ | ۱ | امین پور |

ویژگی‌های حرفی

برای استخراج ویژگی‌های حرفی، ابتدا پیکره به‌عنوان دنباله‌ای از رشته‌های حرفی در نظر گرفته شد و سپس، صد، دونگاشت حرفی متداول استخراج شد. در جدول شماره ۲۰، پنج دونگاشت پرکاربرد به‌همراه بسامدشان قابل مشاهده است.

جدول شماره ۲۰: پنج دونگاشت پرتکرار حرفی در پیکره اشعار

| بسامد | دونگاشت حرفی | |
|-------|--------------|---|
| ۳۷۰۵ | ـ ۵ | ۱ |
| ۳۵۸۹ | ـ ی | ۲ |
| ۲۸۱۴ | ـ ب | ۳ |
| ۲۸۰۲ | ـ ر | ۴ |
| ۲۷۴۳ | ـ ن | ۵ |

سپس، بسامد هر یک از این ویژگی‌ها در هر سند محاسبه و وارد بردار ویژگی‌های حرفی شد. نتایج دسته‌بندی به روش K نزدیک‌ترین همسایه با استفاده از دو معیار فاصله اقلیدسی و کسینوسی به ترتیب در جداول شماره ۲۱ و ۲۲ نشان داده شده است.

جدول شماره ۲۱: نتایج ارزیابی دسته‌بند K نزدیک‌ترین همسایه با استفاده از فاصله اقلیدسی روی ویژگی‌های حرفی

| معیار F | دقت | بازخوانی | |
|---------|-------|----------|-------------|
| ۷۳/۹۱ | ۶۰/۷۱ | ۹۴/۴۴ | امین پور |
| ۸۵/۳۶ | ۹۷/۲۲ | ۷۶/۰۸ | عبدالملکیان |
| ۷۹/۶۳ | ۷۸/۹۶ | ۸۵/۲۶ | میانگین |

جدول شماره ۲۲: نتایج ارزیابی دسته‌بند K نزدیک‌ترین همسایه با استفاده از فاصله کسینوسی روی ویژگی‌های حرفی

| معیار F | دقت | بازخوانی | |
|---------|-------|----------|-------------|
| ۷۳/۱۷ | ۶۵/۲۱ | ۸۳/۳۳ | امین پور |
| ۸۷/۳۵ | ۹۲/۶۸ | ۸۲/۶۰ | عبدالملکیان |
| ۸۰/۲۶ | ۷۸/۹۵ | ۸۲/۹۷ | میانگین |

استفاده از ویژگی‌های حرفی در مطالعه سبک شاعران امکان بررسی هجاها، وزن و قافیه که از عوامل مهم در متمایز کردن سبک آن‌هاست، در اختیار محققان قرار می‌دهد. روش پایانی در این بخش، استفاده از دسته‌بند بیز ساده است. نتایج ارزیابی و ماتریس ابهام به ترتیب در جداول شماره ۲۳ و ۲۴ قابل مشاهده است.

جدول شماره ۲۳: نتایج ارزیابی دسته‌بند بیز ساده روی ویژگی‌های حرفی

| معیار F | دقت | بازخوانی | |
|---------|-------|----------|-------------|
| ۷۶/۱۹ | ۶۶/۶۶ | ۸۸/۸۸ | امین پور |
| ۸۸/۳۷ | ۹۵/۰۰ | ۸۲/۶۰ | عبدالملکیان |
| ۸۲/۲۸ | ۸۰/۸۳ | ۸۵/۷۴ | میانگین |

جدول شماره ۲۴: نتایج ماتریس درهم‌ریختگی دسته‌بند بیز ساده روی ویژگی‌های حرفی

| | | |
|----------|-------------|-------------|
| امین پور | عبدالملکیان | |
| ۸ | ۳۸ | عبدالملکیان |
| ۱۶ | ۲ | امین پور |

باتوجه به نتایج ویژگی‌های واژگانی و نحوی، پایین بودن تعداد اسناد در نتایج دسته‌بندهای K نزدیک‌ترین همسایه و بیز ساده تأثیر قابل توجهی دارد. باتوجه به ماتریس ابهام این دو دسته‌بند، تعداد کم اشعار قیصر امین‌پور در پیکره باعث ایجاد دقت پایین‌تری در سامانه شده است.

بحث و نتیجه‌گیری

در این پژوهش سعی شد تا شیوه شناسایی خودکار اشعار دو شاعر دفاع مقدس یعنی قیصر امین‌پور و محمدرضا عبدالملکیان ارزیابی شود؛ به‌همین منظور، برخی از ویژگی‌های زبانی اشعار از طریق روش‌های آماری و الگوریتم‌های یادگیری ماشین، استخراج شد. همچنین، شناسایی ویژگی‌های زبانی شاعران با استفاده از تحلیل در سه سطح واژگانی، نحوی و حرفی انجام گرفت. ویژگی‌های واژگانی با استخراج هزار واژه پر کاربرد بر اساس معیار Tf-idf و ویژگی‌های نحوی و حرفی با استخراج صد دونگاشت پر کاربرد در پیکره اشعار به دست آمد.

نتایج ارزیابی دسته‌بند بیز ساده بدون حذف ایست‌واژه‌ها روی ویژگی‌های واژگانی (جدول شماره ۷) نشان می‌دهد، معیار F برای امین‌پور ۸۹ درصد است؛ بنابراین، سامانه طراحی شده، از میان صد شعر قیصر امین‌پور که در بخش آموزش مشاهده نکرده و برای نخستین مرتبه با آنان روبه‌رو شده است، ۸۹ شعر را به درستی متعلق به قیصر امین‌پور می‌داند. همچنین، سامانه یازده شعر از قیصر امین‌پور را به اشتباه در دسته اشعار عبدالملکیان قرار داده است.

نتایج ماتریس درهم‌ریختگی (جدول شماره ۸) نشان می‌دهد، از میان ۴۶ شعر عبدالملکیان که سامانه پژوهش برای نخستین مرتبه با آنان روبه‌رو شده، ۴۳ شعر را به‌عنوان اشعار عبدالملکیان دسته‌بندی کرده است. همچنین، سامانه به اشتباه ۳ شعر را به امین‌پور نسبت داده است. به‌این ترتیب، از میان ۱۸ شعر امین‌پور که در بخش آزمون در نظر گرفته شده است، سامانه ۲ شعر را به اشتباه از عبدالملکیان و ۱۶ شعر دیگر را به‌درستی از امین‌پور شناسایی کرده است.

باید اذعان داشت در چنین پژوهش‌هایی، کاستی‌هایی نسبت به تحلیل کارشناسان وجود دارد؛ اما، به این معنی نیست که این روش‌ها ناکارآمد هستند. در بررسی‌های آماری، روش‌های بسیاری برای یادگیری ماشینی با استفاده از الگوریتم‌های مختلف و اهداف مشخص طراحی شده است. هر یک از این روش‌ها، مزایا و کاستی‌های منحصر به فردی دارند. شاید این‌گونه روش‌ها ساده به نظر می‌رسند؛ اما، از آنان در تحلیل‌های انسانی استفاده می‌شود. ترکیب چندین روش و استفاده از آن‌ها در سطوح متفاوت زبانی می‌تواند در بررسی‌های ادبی مفید باشد.

در این پژوهش، اشعار دو شاعر دفاع مقدس مورد مطالعه قرار گرفته است. این شاعران در حوزه‌ای مشخص و فضایی یکسان اشعاری سروده‌اند. مطالعه این آثار با استفاده از ویژگی‌های واژگانی مهم، دونگاشت‌های نحوی و حرفی در سطح کل اشعار یا پیکره آموزش انجام گرفته است. بررسی نهایی این معیارها نشان می‌دهد، ویژگی‌های واژگانی در دسته‌بند بیز ساده از بالاترین دقت برخوردار است؛ از این رو، دسته‌بند بیز ساده با استفاده از این سطح از ویژگی‌ها قادر به شناسایی سبک شاعران است.

دسته‌بند بیز در تمام ویژگی‌ها، دقتی بالاتر از دسته‌بند K نزدیک‌ترین همسایه دارد. دسته‌بند بیز ساده علاوه بر سرعت بالا و سادگی، دقت بالایی دارد. علت این نتیجه، کم بودن پیکره مورد مطالعه است؛ زیرا، بیز ساده در پیکره‌های کم حجم، دقت مناسبی دارد.

نتایج به دست آمده، تأثیر عدم حذف ایست‌واژه‌ها و واژه‌های دستوری در دسته‌بندی شناسایی سبک نویسنده را ثابت می‌کند. حذف این‌گونه واژه‌ها بر دسته‌بندی مبتنی بر موضوع و شناسایی سبک نویسنده، تأثیر مثبتی و بر دسته‌بندی اسناد، تأثیر منفی دارد.

دسته‌بند K نزدیک‌ترین همسایه، یکی از روش‌های یادگیری ماشین در دسته‌بندی متون است. در این پژوهش این روش برای شناسایی سبک نویسنده بررسی شد. معیاری که برای محاسبه فاصله ویژگی‌ها روی بردارها در نظر گرفته می‌شود، اقلیدسی است؛ اما، در پردازش متن فاصله کسینوسی لحاظ می‌شود. در این پژوهش هر دو فاصله مورد مطالعه قرار گرفت و اختلاف آنان تعیین شد. به کارگیری فاصله کسینوسی دقت دسته‌بند را حدود ۱۰ درصد افزایش می‌دهد. این نتیجه نشان می‌دهد، فاصله کسینوسی در حالت‌ها و بخش‌های مختلف پردازش متن برای شناسایی نویسنده مفید است.

کتابنامه

منابع فارسی

- آذین، زهرا (۱۳۹۲). شناسایی خودکار شاعران شعر نو با استفاده از ویژگی‌های زبانی، پایان‌نامه کارشناسی ارشد زبان‌شناسی، تهران: دانشگاه صنعتی شریف، دانشکده زبان‌شناسی.
- استاجی، اعظم (۱۳۸۷). تشخیص مؤلف متون ادبی و قانونی، بحثی در زبان‌شناسی قانونی، نشریه زبان و زبان‌شناسی، دوره چهارم، شماره ۲، ۱۵-۳۲.
- بی‌جن‌خان، محمود (۱۳۸۱). طرح مدل‌سازی زبان فارسی مرحله دوم، آزمایشگاه گروه زبان شناسی دانشکده ادبیات و علوم انسانی دانشگاه تهران.
- سرایبی، محمدحسین و شاهقلیان، آذر (۱۳۸۹). کاوش متون فارسی بر مبنای روش طبقه‌بندی، نشریه انجمن کامپیوتر ایران، جلد هشتم، شماره ۱، ۱۳-۸.
- سمیعی گیلانی، احمد (۱۳۸۴). سبک، نشریه نامه فرهنگستان، شماره ۶، ۱۰۲-۸۶.
- سمیعی گیلانی، احمد (۱۳۸۶). مبانی سبک‌شناسی شعر، نشریه ادب پژوهش، شماره ۲، ۷۶-۴۹.
- شاهمیری، امیرشهاب و مطش بروجردی، محمدرضا (۱۳۸۶). تعیین شاعر به کمک روش‌های یادگیری ماشینی، مجموعه مقالات سومین کنفرانس بین‌المللی فناوری و دانش، مشهد: دانشگاه فردوسی مشهد.
- فرهمندپور، زینب؛ نیک‌مهر، هومن؛ منصوریزاده، محرم و طبیب‌زاده‌قمصری، امید (۱۳۹۱). یک سیستم نوین هوشمند تشخیص هویت نویسنده فارسی‌زبان براساس سبک نوشتاری، نشریه محاسبات نرم، دوره اول، شماره ۲، ۲۶-۳۵.
- کامیار، حسین (۱۳۹۰). روش جدید وزن‌دهی معنایی به کلمات در کاربردهای پردازش متن، پایان‌نامه کارشناسی ارشد مهندسی کامپیوتر، مشهد: دانشگاه فردوسی مشهد، دانشکده مهندسی.
- مجیری، محمدمهدی و مینایی، بهروز (۱۳۸۷). تشخیص وزن عروضی اشعار فارسی: کاربرد جدیدی از متن کاوی، دومین کنفرانس داده‌کاوی ایران، تهران: دانشگاه صنعتی امیرکبیر.

یاحقی، محمدجعفر و ایزانلو، علی (۱۳۸۵). سبک سنجی، نقد و بررسی شیوه آماری کیوسام در انتساب یک اثر، نشریه زبان و ادبیات فارسی دانشگاه خوارزمی، دوره چهاردهم، شماره ۵۳-۵۲، ۱۹۰-۱۵۱.

منابع انگلیسی

- Abbasi, A., & Chen, H. (2005). Applying authorship analysis to extremist group web forum messages, *IEEE Intelligent Systems*, 20 (5), 67-75.
- Baayen, H; Halteren, H. V; Neijt, A., & Tweedie, F. (2002). An experiment in authorship attribution. *JADT 2002: Sixth International Conference on Textual Data Statistical Analysis*, 29-37.
- Jurafsky, D., & Martin, J. H. (2006). *Speech and Language Processing: An introduction to natural language processing, computational linguistics, and speech recognition*, United State: Prentice Hall.
- Mechti, S; Jaoua, M; Faiz, R; Belguith, L. H., & Bsir, B. (2015). On the Empirical Evaluation of Author Identification Hybrid Method Notebook for PAN at CLEF 2015. *CLEF 2015 Evaluation Labs and Workshop*, France: Toulouse.
- Olsson, J. (2004). *Forensic linguistics, an introduction to language, crim and law*. London, New York: Continuum.
- Patton, J. M., & Can, F. (2004). A stylometric analysis of Yasar Kemal's Ince Memed tetralogy, *Computers and the Humanities*, 38 (4), 457-467.
- Salton, G. B. (1988). Term-weighting approaches in automatic text Retrieval, *Information Processing & Management*, 24 (5), 513-523.
- Stamatatos, E. (2009). A Survey of Modern Authorship Attribution Methods, *Journal of the American Society for information Science and Technology*, 60 (3), 538-556.