

تشخیص خود کار جنسیت نویسنده در متون فارسی

مهدی مرادی و محمد بحرانی

مرکز زبان‌ها و زبان‌شناسی، دانشگاه صنعتی شریف، تهران، ایران

چکیده

با گسترش روزافزون استفاده از اینترنت، شاهد تبادل هزاران گیگابایت اطلاعات متنی در فضای مجازی هستیم. فضای سایبر این امکان را به اشخاص می‌دهد که هویت حقیقی خود را مخفی کنند و با هویت ساختگی جدیدی وارد آن شوند. از این رو اهمیت حفظ امنیت این فضا، کنترل بر محتوای تولیدشده توسط کاربران و شناسایی مشخصات تولیدکنندگان محتوا هر روز پررنگ‌تر می‌شود. موضوع مورد بررسی در این پژوهش که مربوط به حوزه شناسایی نویسنده است، شناسایی خودکار جنسیت نویسنده متن فارسی است. به منظور شناسایی جنسیت، با توجه به مطالعات زبان‌روان‌شناختی صورت گرفته، ۴۸ ویژگی روان‌شناختی و سبک‌شناسی تعریف شد. دو پیکره متنی جهت آموزش طبقه‌بندها تهیه و سپس برای شناسایی جنسیت، سه الگوریتم یادگیری ماشینی مختلف (ماشین بردار پشتیبان، بیز ساده و درخت تصمیم) طراحی شد. نتایج اعتبارسنجی متقابل ده تایی نشان داد که بیشترین دقت مربوط به طبقه‌بند درخت تصمیم با دقت ۷۳/۸٪ است.

واژگان کلیدی: تشخیص جنسیت، شناسایی نویسنده، متن کاوی.

۱- مقدمه

عصر ارتباطات و گسترش استفاده از نامه‌های الکترونیکی، پیام‌های فوری (IM)، اسناد، وبلاگ‌ها، مقالات خبری، صفحات خانگی، تالارهای گفتگو و شبکه‌های اجتماعی، متن را به اصلی‌ترین وسیله ارائه و انتقال اطلاعات تبدیل کرده است.

گم‌نامی^۱ به‌عنوان یکی بارزترین ویژگی‌های جامعه مجازی، (ژنگ و همکاران، ۲۰۰۶) از یک‌سو امکان بیان آزادانه عقاید و نظرات و از سوی دیگر امکان جرایم اینترنتی را فراهم می‌سازد. بی‌چهرگی و بی‌نامی افراد در اینترنت نه تنها پیگرد قانونی جرم و مجرم را در فضای رایانه‌ای مشکل می‌سازد، بلکه زمینه گسترش بیش‌تر آن را نیز فراهم می‌کند. مجرمان اینترنتی برای ناشناس ماندن، از سرورهای گمنام‌کننده استفاده می‌کنند و یا این‌که با ارائه اطلاعات نادرست در مورد جنسیت، مکان، سن و ملیت خود دست به اقدامات خراب‌کارانه می‌زنند. از این رو استفاده از روش‌های

جدید جهت ردیابی هویت بزه‌کاران ضروری به‌نظر می‌رسد. علاوه بر این تشخیص جنسیت نویسنده می‌تواند کاربردهای تجاری متعددی دیگری نیز داشته باشد. برای مثال با مشخص شدن جنسیت نویسنده می‌توان فهمید که زنان یا مردان بیش‌تر درباره چه موضوعات یا محصولات صحبت می‌کنند و از چه خدمات یا کالاهایی بیزارند. دسترسی به این اطلاعات برای تبلیغات هدفمند بسیار حائز اهمیت است (آرجون و همکاران، ۲۰۱۰). در حوزه ترجمه ماشینی^۲، وقتی قرار است ترجمه بین دو زبانی که نظام ضمیرهایشان متفاوت است، انجام شود مشخص بودن جنسیت ضمیر، دقت ترجمه را افزایش خواهد داد (برای مثال ترجمه جمله "من تنها هستم" بسته به جنسیت نویسنده، در ایتالیایی دو جمله "Sono solo" برای مرد و "Sono sola" برای زن خواهد بود). کاربرد شناسایی جنسیت از دیدگاه علوم اجتماعی و تعیین خصوصیات جمعیت‌شناسی محتوای تولیدشده در اینترنت نیز می‌تواند کاربردی باشد. مطالعات

²Machine translation

¹Anonymity

زبان‌شناختی بر روی تغییرات زبان زنان و مردان در بازه‌های زمانی یا در حوزه‌های مختلف از کاربردهای دیگر شناسایی جنسیت نویسنده خواهد بود.

مطالعات روان‌شناختی نشان می‌دهد که کلمات مورد استفاده هر فرد بیان‌گر ویژگی‌های شخصیتی و روانی آن فرد است (پنباکر و همکاران، ۱۹۹۹)، (استیرمن و همکاران، ۲۰۰۱). نتایج پژوهش‌های روان‌شناسان نشان می‌دهد که ویژگی‌های سبک‌شناسی^۱ هر شخص (مانند کاربرد ضمائر، حروف ربط و یا حروف اضافه) با رفتارها و احساساتش در ارتباط است (پنباکر، ۲۰۰۳). پژوهش‌گران از این ویژگی‌های زبان‌شناختی به‌عنوان ابزاری قانونی^۲ جهت شناسایی نویسنده، جنسیت، سن و یا سطح تحصیلات او استفاده می‌کنند.

۲- مروری بر ادبیات

تشخیص جنسیت نویسنده متن، در حوزه کارهایی قرار می‌گیرد که مربوط به شناسایی نویسنده^۳ (دول و همکاران، ۲۰۰۱) یا تصدیق نویسنده^۴ است (کوپل و همکاران، ۲۰۰۹). در فرایند شناسایی نویسنده، پژوهش‌گر به دنبال یافتن هویت واقعی نویسنده‌ای است که هویتش را مخفی نگاه داشته است.

نخستین پژوهش‌های کمی مربوط به شناسایی نویسنده در سال ۱۸۷۷ توسط مندن هال هواشناس اهل اوهایو صورت گرفت، که پیشنهاد کرد که از توزیع طول کلمات^۵ به‌عنوان ویژگی وابسته به نویسنده استفاده شود (مندنهال، ۱۸۸۷). وی در سال ۱۹۰۱ از این تکنیک برای شناسایی نویسنده اشعار مورد تردید شکسپیر-بیکون استفاده کرد (مندنهال، ۱۹۰۱). از پیش‌گامان رویکرد آماری به مسئله سبک‌شناسی در عصر رایانه، موستلر و والاس بودند که در سال ۱۹۶۴ از لغات دستوری^۶ و تحلیل بیزی^۷ (طبقه‌بند بیز ساده) برای شناسایی هویت نویسندگان مقالات فدرالیست^۸ استفاده کردند (موستلر و والاس، ۱۹۶۴).

با گسترش رایانه، استفاده از ویژگی‌های سبک‌شناسی برای شناسایی نویسنده بیش از پیش مورد توجه پژوهش‌گران قرار گرفت، به‌طوری که تاکنون بیش از هزار ویژگی سبک‌شناسی مانند ویژگی‌های واژه‌بنیاد، نویسه‌بنیاد^۹ (به‌عنوان مثال (یول، ۱۹۹۴)) و نیز ویژگی‌های مبتنی بر کلمات دستوری (به‌عنوان مثال موستلر و والاس) موستلر و والاس (۱۹۸۴)) و علایم نگارشی^{۱۰} (باین و همکاران، ۲۰۰۲) برای این منظور تعریف شده است. هر چند پژوهش‌های بسیار زیادی بر روی زبان انگلیسی صورت گرفته و نتایج قابل قبولی نیز به‌دست آمده است (زاجاری و همکاران، ۲۰۱۲) ولی با این وجود کارهای محاسباتی صورت گرفته بر روی شناسایی نویسنده زبان فارسی بسیار محدود و پراکنده بوده است (زینب فرهمند پور و همکاران، اسفند ۹۰). تاکنون روش‌های مختلفی برای شناسایی نویسنده استفاده شده است. پژوهش‌های نخستین در این حوزه بیش‌تر معطوف به استفاده از هیستوگرام توزیع طول کلمات (مندنهال، ۱۸۸۷)، طبقه‌بندهای بیزی (موستلر و والاس، ۱۹۸۴)، تحلیل مؤلفه‌های اصلی^{۱۱} (بوروز، ۱۹۸۴)، تحلیل خوشه^{۱۲} (هولمز، ۱۹۹۲) و غیره اشاره کرد. پیشرفت‌های سخت‌افزاری و افزایش قدرت پردازش رایانه‌ها، پژوهش‌های این حوزه را به سوی استفاده از تکنیک یادگیری ماشین مانند درخت تصمیم^{۱۳} (آپت و همکاران، ۱۹۹۸) شبکه‌های عصبی^{۱۴} (تویدی و همکاران، ۱۹۹۶)، ماشین‌بردار پشتیبان^{۱۵} سوق داد (دیدریچ و همکاران، ۲۰۰۰).

مسئله مورد کنکاش در این پژوهش بررسی تمایز یا عدم تمایز زبان زنان و مردان است، اینکه جنسیت نمودی در سبک زبانی افراد دارد یا نه و کدام ویژگی‌های زبانی نقش پررنگ‌تری در تمایز زبان این دو جنس دارند. هر چند این تفاوت به قدری نیست که ارتباط بین زنان و مردان را مختل کند، ولی همان‌طور که نتایج پژوهش‌ها نشان می‌دهد، زنان و مردان هر کدام گونه متفاوتی از زبان را به کار می‌برند. طی چند دهه گذشته، پژوهش‌گران از زوایای مختلفی، مسئله

تشویق مردم این ایالت به امضای قانون اساسی پیشنهادی برای ایالات متحده در کنوانسیون ۱۷۸۷ نوشته شدند“ بپولیوس“ نام مستعار نویسندگان این مقالات بود. (ویکی پدیا)

⁹Character-based

¹⁰Punctuation

¹¹Principle component analysis

¹²Cluster analysis

¹³Decision tree

¹⁴Neural networks

¹⁵Support vector machine

¹ Linguistic style

² Forensic tools

³ Authorship identification

⁴ Authorship verification

⁵ Word-length distribution

⁶ Function words

⁷ Bayesian analysis

^۸فدرالیست مجموعه‌ای از ۸۵ مقاله است که

در سال‌های ۱۷۸۷ و ۱۷۸۸ میلادی در دو نشریه در شهر نیویورک برای

موضوع را در مقاله‌ای فهرست کرده است (داوری اردکانی، ۱۳۸۷). در این بین، از جمله پژوهش‌هایی که موضوع را نه از دیدگاه روان‌شناختی و اجتماعی بلکه از نگاه زبان‌شناختی (بررسی نظام ساختاری زبان از نظر انعکاس موضوع جنسیت) مورد کنکاش قرار داده‌اند، می‌توان به بررسی جنسیت در واژگان (فارسیان، ۱۳۷۸)، نگاه مقابله‌ای به زبان فارسی و انگلیسی (چاکانی، ۲۰۰۰) و جنسیت و شمول معنایی (شجاع رضوی، ۱۳۸۶) اشاره کرد. ما در این پژوهش با تعریف ۴۸ ویژگی مختلف، توانستیم با دقت به‌نسبته بالایی، جنسیت نویسنده را تشخیص دهیم.

۳- تبیین مسئله

مسئله تشخیص جنسیت نویسنده، نوعی مسئله طبقه‌بندی دوطبقه^۸ است. که هدف در این مسئله قراردادن متن بی‌نام e در یکی از طبقه‌های {زن، مرد} است (رابطه ۱).

$$e \in \begin{cases} \text{اگر نویسنده } e \text{ زن باشد کلاس } +1 \\ \text{اگر نویسنده } e \text{ مرد باشد کلاس } -1 \end{cases} \quad (1)$$

برای بررسی فرض شماره ۱، باید ویژگی‌های تعریف شوند که در بیشتر متون نوشته‌شده توسط زنان یا متون نوشته‌شده توسط مردان به‌طورنسبی ثابت باشند. پس از مشخص شدن فضای ویژگی‌ها، هر متن e را می‌توان به‌صورت برداری متشکل از d بعد که در آن d تعداد کل ویژگی‌ها است، نشان داد. با داشتن مجموعه‌ای از متون از قبل دسته‌بندی‌شده، طبقه‌بند آموزش داده شده و سپس از آن برای دسته‌بندی متون بدون برچسب استفاده می‌شود. بیان مسئله به زبان ریاضی عبارت است از آموزش طبقه بنی $y=f(x)$ ، از روی دادگان آموزش $D=(x_1, y_1), (x_2, y_2), \dots, (x_N, y_N)$ آن $x_i=[x_{i1}, x_{i2}, \dots, x_{id}]$ برداری متشکل از d بعد است. اگر $Y=\{y_i, i=1, 2, \dots, N\}$ نشان‌دهنده مجموعه برچسب باشد، $y_i \in \{+1, -1\}$ نشان‌دهنده طبقه ۱- (مرد) و طبقه ۱+ (زن) و N تعداد نمونه در مجموعه داده است. در کل شناسایی جنسیت نویسنده را می‌توان به چهار مرحله تقسیم کرد شکل (۱).

۱. تهیه پیکره دارای برچسب جنسیت

^۸Binary Classification problem

تفاوت بین زنان و مردان و نیز مسئله جنسیت و زبان را مورد پژوهش قرار داده‌اند (تالبوت، ۱۹۹۸). این نوع پژوهش‌ها که در آغاز بیش‌تر معطوف به بررسی تأثیر متغیر جنسیت در رفتار کلامی افراد در سطح آوایی و نیز شیوه تعامل بودند (ووداک و همکاران، ۲۰۰۷). با انتشار مقاله روبین لیکاف ("جایگاه زنان در زبان") در ۱۹۷۵ بدل به نقطه عطفی در مطالعات زبان و جنسیت شد. لیکاف در این مقاله با برشمردن ویژگی‌های واژگانی، نحوی و کاربردشناختی^۱، مانند استفاده از واژه‌های خاص، پوچواژه‌ها^۲، ضمیمه‌های سؤالی^۳، عبارات غیر صریح^۴ و قواعد از نوع تصحیح افراطی^۵ سعی در شناسایی سبک نگارش زنان داشت (لیکاف، ۱۹۷۵). تالبوت در کتاب خود می‌گوید: الگوهای زبانی در گفتار زنان و مردان، نشان‌گر طبقه اجتماعی آنان است. در محیط کار، زنان طبقه پایین‌تر هنگام گفت‌گو با رؤسایشان از زبان مؤدبانه‌تری استفاده می‌کنند (تالبوت، ۱۹۹۸). طبق نتایج پژوهش مولاک و همکاران، در گفتگوهای دو نفره، جملات پرسشی بیش‌تر در کلام زنان و جملات امری بیش‌تر در کلام مردان دیده می‌شود (مولاک، ۱۹۹۸). باید توجه کرد که پژوهش‌گران بین مفهوم جنس^۶ و جنسیت^۷ تفاوت قائلند. جنس به ابعاد بیولوژیکی مردانگی و زنانگی فرد اشاره دارد؛ درحالی‌که جنسیت به صفات و ویژگی‌های اجتماعی دو جنس اطلاق می‌شود (اونگر، ۱۹۷۹). جنس هر فرد از بدو تولدش مشخص است؛ حال آن‌که جنسیت را جامعه تعیین می‌کند؛ که به ویژگی‌های شخصی و روانی شخص دلالت دارد. به عبارت دیگر هر زنی رفتار زنانه ندارد و هر مردی نیز لزوماً مردانه رفتار نمی‌کند. مسئله مورد بررسی این پژوهش تفاوت‌های زبانی وابسته به جنسیت است و به تفاوت‌های مرتبط با جنس پرداخته خواهد شد.

۲-۱- جنسیت و زبان فارسی

هر چند کارهایی بر روی زبان‌های دیگر انجام شده، ولی در مورد شناسایی خودکار جنسیت نویسنده در زبان فارسی کار چندان مهمی صورت نگرفته است. سروری برکارهای (غیر محاسباتی) انجام‌شده بر روی زبان فارسی نشان از انجام پژوهش‌هایی بر روی گونه نوشتاری و گفتاری فارسی دارد. نگار داوری اردکانی مقالات و پایان‌نامه‌های مرتبط با این

^۱ Pragmatic

^۲ Expletive

^۳ Tag questions

^۴ Hedges

^۵ Hypercorrect

^۶ Sex

^۷ Gender

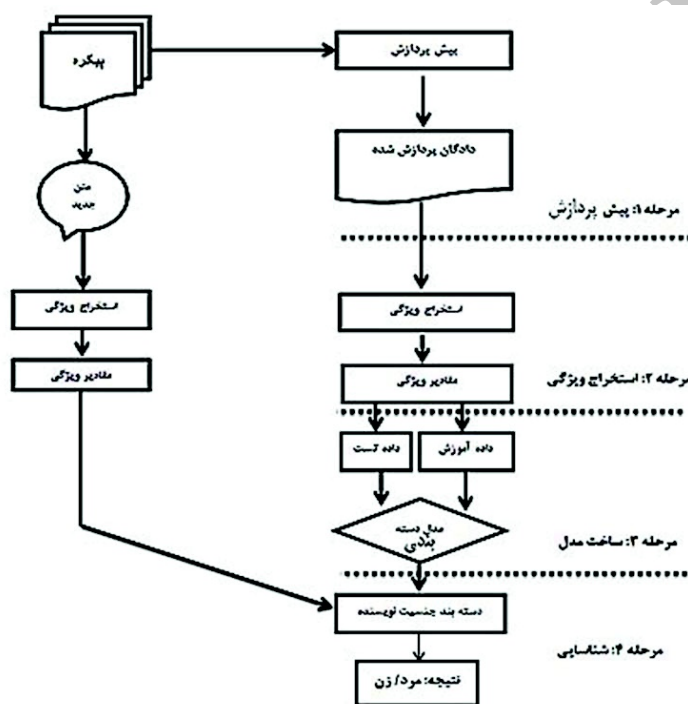
که متون موجود در این پیکره فاقد برجسب مشخص کننده جنسیت بودند، با توجه به اسم نویسنده داستان یا رمان، به صورت دستی تمامی متون موجود در بخش کتاب‌های داستانی این پیکره به دو گروه زن و مرد تفکیک شد. بخش دیگری از داده‌ها نیز از رمان‌ها، داستان‌های بلند و کوتاهی که به صورت پی‌دی‌اف در اینترنت موجود بودند تهیه گشت. برای حذف تأثیر عامل جنسیت مترجم بر روی متن اصلی اثر، تمامی داستان‌ها و رمان‌های غیر فارسی از مجموعه مورد نظر حذف شد. پس از تهیه پیکره به منظور حذف عامل تعداد متن بر روی استخراج ویژگی‌ها، در هر دو دسته زن و مرد تعداد مساوی ۲۵۴ داستان قرار داده شد.

۲. شناسایی ویژگی‌های زبانی متمایز کننده جنسیت
۳. استخراج خودکار مقادیر هر ویژگی از هر متن
۴. ساخت مدل دسته‌بندی برای شناسایی جنسیت نویسنده در متون هدف

۴- تهیه پیکره آموزشی (دارای برجسب جنسیت)

۱-۴ - پیکره بی جن خان (متون ادبی)

بخشی از داده‌های مورد استفاده در این پژوهش از پیکره متنی بی جن خان (بی جن خان، ۱۳۸۶) به دست آمد. از آنجا



(شکل-۱): فرآیند شناسایی جنسیت نویسنده

بود، ولی از آنجا که کاربران ایرانی به طور معمول جملات دیگران را توییت کرده یا در دیوار خود می نویسند، تصمیم گرفته شد که از نظرات کاربران در سایت هلو کیش استفاده شود (سایت هلوکیش، ۲۰۱۲). برای خزش^۷ در این سایت از برنامه بیوتیفول سوپ^۸ که به صورت کتابخانه‌ای در محیط

۲-۴ - پیکره متون گفتاری

برای تهیه این پیکره، از رویکرد وب برای پیکره^۱ استفاده شد. مهم ترین چالش، شناسایی سایتی فارسی بود که نخست محتوای متنی آن توسط خود کاربر نوشته و دوم نام نویسنده هر متن نیز مشخص شده باشد. در ابتدا تصمیم بر استفاده از توییت^۲ های سایت تویتر^۳ یا استاتوس^۴ های سایت فیس بوک^۵

⁴ Statues= وضعیت

⁵ Www.facebook.com

⁶ Wall

⁷ Crawling

⁸ Beautiful Soup

¹ Web for corpus

² Tweet

³ Www.tweeter.com

تعداد نویسه‌های ویژه (# و %) و غیره می‌شود (مریم، ۱۹۹۶). ویژگی‌های مبتنی بر واژه شامل یازده معیار آماری مانند غنای واژگانی، معیار یول (K) و معیار سیمپسون (D) هستند. پژوهش‌های چهار دهه گذشته نشان می‌دهد که شرایط جسمی و روانی افراد در واژه‌گزینی آن‌ها تأثیر دارد (گوتسچالک، ۱۹۹۶). تحلیل‌های انجام گرفته بر روی متون نشان داده است که تعداد کلمات مثبت (مانند عشق، خوب) در نوشته‌های افراد خوش‌قلم بیش‌تر به چشم می‌خورد؛ این افراد از کلمات دارای بار معنایی منفی (مانند زشت، کریه) در حد متعارف استفاده می‌کنند و کلمات شناختی (مانند دانستن، علت) را فراوان مورد استفاده قرار می‌دهند (پنباکر و همکاران، ۲۰۰۷). ویژگی‌های نحوی، سبک نگارشی نویسنده را در سطح جمله نشان می‌دهند. این ویژگی‌ها شامل تعداد علائم نگارشی (مانند ویرگول، نقطه، علامت تعجب و غیره) است. به‌علت تفاوت زنان و مردان در استفاده از علائم نگارشی، از این ویژگی‌ها می‌توان برای تشخیص جنسیت نویسنده استفاده کرد. به‌عنوان مثال طبق پژوهش‌ها مولاک، بسامد علامت سؤال در متون زنان نسبت به مردان بیش‌تر است (مولاک و همکاران، ۱۹۹۰).

به‌کمک ویژگی مبتنی بر ساختار، نویسنده چیدمان نوشته خود را سازماندهی می‌کند. افراد عادت‌های متفاوتی در سازماندهی متن خود دارند. ویژگی‌های مانند میانگین طول جمله، تعداد خطوط، تعداد خطوط خالی و غیره جزء ویژگی‌های ساختاری یک متن شمرده می‌شوند.

لغات دستوری، لغاتی هستند که معنای واژگانی کمی دارند، ولی کاربرد دستوری زیادی داشته و برای نشان دادن روابط دستوری مابین کلمات استفاده می‌شود. از این ویژگی‌ها می‌توان به ضمیر، اصوات و حروف اضافه اشاره کرد.

نشانه‌های روانی‌زبانی، مجموعه کلماتی هستند که بیان‌کننده ویژگی‌های شخصیتی و روانی شخص (مانند بسامد استفاده از صفات مثبت یا منفی، رنگ‌ها، کلمات رکیک) است.

۵-۱- استخراج خودکار و بازنمایی ویژگی‌ها

برای استخراج ویژگی‌های روانی‌زبانی جدول (۳) مانند تعداد صفات و تعداد قیدها و ضمیر، از برجسبزن اجزای کلام با مدل مارکوف مخفی مرتبه دو، که از روی پیکره برجسب‌خورده بی‌جن‌خان (بی‌جن‌خان، ۱۳۸۶) آموزش دیده شده بود استفاده شد. فهرست کلمات مربوط به ضمیر

برنامه‌نویسی پایتون وارد می‌شود و قابلیت‌های بالایی در پارس صفحات HTML/XML دارد استفاده شد. برای استخراج نظرات مرتبط با نظردهندگان زن و مرد، ابتدا فهرستی تفکیک‌شده از اسامی فارسی زن و مرد تهیه شد؛ سپس نام نویسنده نظر با این فهرست اسامی تطبیق داده شده و نظرات بر حسب این فهرست تفکیک‌شده، در دو دسته زن و مرد قرار گرفت.

(جدول-۱): پیکره داستان و رمان فارسی

میانگین تعداد کلمه در هر سند	تعداد سند	
۱۷۷۸۰	۲۵۴	نویسنده مرد
۲۸۶۹۰	۲۵۴	نویسنده زن
۲۳۳۳۵	۵۰۸	کل

(جدول-۲): پیکره گفتاری نظرات

میانگین تعداد کلمه در هر سند	تعداد سند	
۸۵	۱۱۸۹	نویسنده مرد
۸۸	۱۱۸۹	نویسنده زن
۸۶/۵	۲۳۷۸	کل

از آن‌جا که بسیاری از متون فارسی با صفحه کلید عربی و یا غیراستاندارد نوشته شده و می‌شوند، احتمال وجود نویسه‌های غیر فارسی (برای مثال عربی یا اردو) در متون وجود دارد که در صورت عدم هنجارسازی این نویسه‌ها، نتایج پردازش‌های انجام‌شده بر روی متون قابل اطمینان نخواهند بود. قبل از انجام پژوهش، تمامی نویسه‌های غیرفارسی به معادل فارسی تبدیل شدند. جداول (۱) و (۲) مشخصات این دو پیکره را نشان می‌دهد.

۵- انتخاب مجموعه ویژگی

با توجه به ادبیات پژوهش و بر پایه پژوهش‌های روان‌شناختی، پنج مجموعه ویژگی مرتبط با جنسیت شامل: (۱) نویسه‌بنیاد (۲) واژه‌بنیاد (۳) نحوی (۴) ساختاربنیاد و (۵) کلمات دستوری برای هر سند محاسبه شده است. جدول شماره ۳ فهرست ۴۸ ویژگی استخراج‌شده را نشان می‌دهد. ویژگی‌های مبتنی بر نویسه شامل شش ویژگی سبک‌شناسی هستند که به‌طور گسترده در مسائل تشخیص نویسنده به کار می‌روند و شامل ویژگی‌هایی مانند تعداد نویسه خالی،

¹Import

۶-۱- بیز ساده به عنوان الگوریتم دسته‌بندی

طبقه‌بندهای بیز ساده، سرعت محاسباتی بالایی داشته و با داده‌های آموزشی کم نیز دقت خوبی نشان می‌دهند. برای فرموله کردن تشخیص جنسیت نویسنده با استفاده از مدل استاندارد بردارند به صورت زیر عمل می‌کنیم: هر متن را به صورت $D = (F_1, F_2, \dots, F_n)$ در نظر می‌گیریم، که در آن F_i ویژگی مورد نظر (مانند طول کلمه، تعداد صفت و غیره) است با داشتن این بردار، هدف ما تشخیص جنسیت نویسنده آن یعنی G است:

$$G \in \{\text{male, female}\} \quad (۳)$$

جنسیت G از طریق بیشینه کردن احتمال $P(G|D)$ به دست می‌آید. طبق قانون بیز داریم:

$$P(G|D) = \frac{P(G)P(D|G)}{P(D)} \quad (۴)$$

$$P(D|G) = P(F_1, F_2, \dots, F_n|G) \quad (۵)$$

با توجه به فرض بیز ساده مبنی بر استقلال متغیرها داریم:

$$P(F_1, F_2, \dots, F_n|G) = P(F_1|G)P(F_2|G) \dots P(F_n|G), \quad (۶)$$

و فرمول به شکل زیر تغییر می‌کند:

$$\text{gender} = \underset{G}{\text{argmax}} P(G) \prod_{i=1}^n p(F_i|G) \quad (۷)$$

از آنجا که داده آموزشی ما برای دو جنس حجم یکسانی داشت، احتمال پیشین $P(G)$ برای هر دو یکسان خواهد بود.

۶-۲- الگوریتم ADtree (درخت تصمیم تناوبی)

در الگوریتم درخت تصمیم‌گیری، هر گره درخت، نمایان‌گر یک ویژگی خاص است. در هر نقطه از پیمایش این درخت، تصمیم‌گیری براساس مقدار این ویژگی خاص در نمونه دسته بندی نشده صورت می‌گیرد. در کل، درخت تصمیم‌گیری روشی متداول است که حوزه کاربرد گسترده‌ای دارد (صفویان و لاندگرب، ۱۹۹۱). درخت تصمیم ADtree، نوع وزن دار و تعمیم‌یافته‌ای از درخت تصمیم‌گیری است که از دو نوع گره تشکیل شده: گره تصمیم‌گیری و گره پیش‌بینی. ترکیب امتیاز گره‌های پیش‌بینی پیمایش شده و گره‌های

فاعلی، پرسشی، حرف ربط، حرف ربط گروهی، صوت و حروف اضافه از واژگان زایای فارسی به دست آمد (اسلامی و همکاران، ۱۳۸۳). برای شمارش صفات مثبت و منفی در متون، به صفات موجود در واژگان زایای زبان فارسی برچسب مثبت (۱۸۴۰ صفت) و منفی (۱۵۳۶) داده شد و سپس این صفات در متن شمارش شد. برای استخراج رنگ‌ها از ویژگی سلسله‌مراتبی^۱ صفحات ویکی‌پدیا فارسی استفاده شد (برای مثال در صفحه "رده رنگ‌ها" می‌توان نام رنگ‌های مختلف زبان فارسی را مشاهده کرد). برای استخراج کلمات مربوط به شک و تردید و قطعیت، فهرست واژگان مربوط به این دو حوزه در زبان انگلیسی از سایت‌های مختلف گردآوری و معادل فارسی آن‌ها نوشته و در متون شمارش شدند. فهرست کلمات رکبیک با جستجو در اینترنت جمع‌آوری شد (به ضمایم مراجعه شود).

استخراج‌کننده ویژگی برای هر مستند متنی، یک بردار ۴۸ بعدی جهت بازنمایی مقادیر تولید می‌کند. از آنجا که اعداد مربوط به ویژگی‌ها از روش‌های مختلفی به دست آمده بود، این اعداد ممکن بود از صفر تا بیست‌هزار تغییر کنند. برای مثال مقدار مربوط به تعداد کلمه برای یک مستند ممکن است ۱۴۵۴۰ باشد؛ درحالی‌که میانگین تعداد نویسه در کلمه برای همان مستند ۳.۴ باشد. از این‌رو، برای اطمینان از این‌که در فرایند دسته‌بندی با تمامی مقادیر به صورت یکسان برخورد خواهد شد، تمامی مقادیر را به کمک روش هنجارسازی کمینه-بیشینه^۲ نرمال کردیم تا در حد واسط صفر تا یک قرار گیرند:

$$\text{Normalized } x_{ij} = \frac{x_{ij} - \min(x_j)}{\max(x_j) - \min(x_j)} \quad (۲)$$

که در این فرمول: x_{ij} ویژگی i ام در نمونه j ام و $\min(x_j)$ و $\max(x_j)$ کمترین و بیشترین مقادیر ویژگی i ام هستند.

۶-۳ روش دسته‌بندی

برای دسته‌بندی از سه دسته‌بند بیز ساده^۳، درخت تصمیم ADtree و بردار ماشین پشتیبان^۴ که به صورت گسترده در مسائل دسته‌بندی متون کاربرد دارند، استفاده شد. برای دسته‌بندی از بسته یادگیری ماشین وکا استفاده شد (مارک هال و همکاران، ۲۰۰۹).

^۱ Hierarchical

^۲ max-min normalization method

^۳ Naïve Bayes

^۴ Support vector machine

ادامه (جدول-۳): مجموعه ویژگی‌ها و توضیحات

ویژگی	تعداد
تعداد کلمات دستوری	F_{30}
تعداد ضمیر فاعلی / N	F_{31}
تعداد ضمیر پرسشی / N	F_{32}
حرف ربط / N	F_{33}
حرف ربط گروهی / N	F_{34}
صوت / N	F_{35}
حرف اضافه / N	

ادامه (جدول-۳): نشانه‌های روانی-زبانی

ویژگی	تعداد
صفات مثبت	F_{36}
صفات منفی	F_{37}
رنگ‌ها	F_{38}
کلمات رکیک	F_{39}
تعداد صفات	F_{40}
تعداد قید	F_{41}
تعداد ضمائر	F_{42}
نسبت ضمائر ۳ ویژگی	F_{43-46}
شک و تردید	F_{47}
قطعیت	F_{48}

۳-۶- ماشین بردار پشتیبان

این روش یکی از پرکاربردترین الگوریتم‌ها در حوزه شناسایی نویسنده است (دیدریچ و همکاران، ۲۰۰۳). اساس این روش، یافتن بهترین ابرصفحه^۸ دارای بیشترین فاصله است، که بتواند داده‌های دو دسته را از هم تفکیک کند. این الگوریتم قادر است بدون بیش‌برازش شدن، تعداد زیادی ویژگی را در قالب یک مسئله دسته‌بندی خطی یا غیر خطی به کار گیرد. اگر ورودی این الگوریتم را بتوان به صورت خطی جدا کرد، SVM فاصله بین این دو طبقه را با جستجوی ابرصفحه بهینه جدا کننده، بیشینه می‌کند. بردار وزنی بهینه ω^* را با حل مسئله بهینه‌سازی زیر می‌توان یافت:

$$\text{Minimize } j(\omega) = \frac{1}{2} \|\omega\|^2 \quad (A)$$

$$\text{به شرطی که } y_i (\omega \cdot x_i - b) \geq 1$$

در مواردی که داده‌ها به صورت غیر خطی قابل جداسازی نیستند، SVM با استفاده از یک متغیر بی‌اثر ξ و اجازه قرار گرفتن نمونه‌های آموزشی در ناحیه‌ای مابین دو ابرصفحه‌ای که از نقطه پشتیبانی دو طبقه می‌گذرد، اقدام به ساختن یک حاشیه نرم می‌کند (کورتکس و واپنیک،

⁸ Hyperplane
⁹ Soft margin

انتهایی، تعیین کننده دسته هر نمونه است (ماسون و همکاران، ۱۹۹۱).

جدول-۳: مجموعه ویژگی‌ها و توضیحات

ویژگی	توضیح ویژگی
F_1	تعداد کل نویسه‌ها (C)
F_2	تعداد کل حروف الفبا(الف-ب)/C
F_3	تعداد کل اعداد/C
F_4	تعداد نویسه فاصله/C
F_5	تعداد نویسه تب/C
F_6	تعداد نویسه‌های ویژه/C
F_7	تعداد کل کلمات (N)
F_8	میانگین تعداد نویسه در کلمه
F_9	غناي واژگانی (کل لغات یکتا/N)
F_{10}	کلمات طولانی (بزرگ‌تر از ۲ نویسه)/N
F_{11}	کلمات کوتاه (کوچک‌تر از ۲ نویسه)/N
F_{12}	هپکس لگومنا ^۱ (کلمات ۱ تکراره)/N
F_{13}	هپکس دیسلگومنا ^۲ (کلمات ۲ تکراره)/N
F_{14}	معیار K یول ^۳
F_{15}	معیار D سیمپسون ^۴
F_{16}	معیار S سیشل ^۵
F_{17}	معیار R هونور ^۶
F_{18}	تعداد کاما C(/,)
F_{19}	تعداد نقطه C(/.)
F_{20}	تعداد دونقطه C(/:)
F_{21}	تعداد سمیکلون C(/@)
F_{22}	تعداد علامت سوال C(/?)
F_{23}	تعداد علامت تعجب C(/!)
F_{24}	تعداد علامت تعجب ۳ تایی C(/!!!)
F_{25}	تعداد کل خط
F_{26}	تعداد کل جملات (S)
F_{27}	تعداد کلمه در هر جمله
F_{28}	تعداد خطوط خالی / تعداد کل خطوط
F_{29}	میانگین طول خطوط غیر خالی

¹ Hapax legomena
² Hapax dislegomena
³ Yule's K measure
⁴ Simpson's D measure
⁵ Sichel's S measure
⁶ Honore's R measure
⁷ Structural features

۱۹۹۵). در این حالت مسئله بهینه‌سازی به صورت زیر تغییر می‌کند.

$$\text{Minimize } j^*(\omega) = \frac{1}{2} \|\omega^*\|^2 + G \sum_{i=1}^N \xi_i \quad (9)$$

به شرطی که $\xi_i \geq 0$ و $y_i(\omega^* \cdot x_i - b) \leq 1 - \xi_i$

وقتی از SVM برای حل مسائل غیرخطی استفاده می‌شود، حقه هسته^۱ کمک می‌کند تا فضای ویژگی x به یک فضای بعد بالای $\phi(x)$ نگاشت شود و ابرصفحه دارای بیش‌ترین فاصله را در فضای جدید جستجو کند. در این مقاله برای دسته‌بندی از هسته PUK^۲ استفاده شده است که هم مقاوم است و هم قدرت نگاشت خوبی دارد (اوستون و همکاران، ۲۰۰۶).

۷- نتایج به دست آمده از دسته‌بندی

سه الگوریتم بیز ساده، درخت تصمیم ADtree و SVM را بر روی دو پیکره اجرا کردیم. برای ارزیابی مدل‌ها از روش ارزیابی اعتبارسنجی متقابل^۳ ده‌تایی استفاده شد. دقت دسته بندی بر روی پیکره رمان و داستان به ترتیب ۰.۶۶/۳، ۰.۷۳/۶ و ۰.۷۱/۶ بودند؛ درحالی‌که دقت همین دسته‌بندیها بر روی پیکره متنی نظرات کاربران به ترتیب ۰.۵۳، ۰.۶۳ و ۰.۵۸.۵ بود شکل (۲). به منظور مقایسه نتایج حاضر با نتایج تشخیص جنسیت در زبان انگلیسی، در جدول (۴) نتایج چند نمونه پژوهش در زبان انگلیسی آورده شده است.

۸- ویژگی‌های متمایز کننده

به منظور بررسی اهمیت هر گروه از ویژگی‌های پیشنهادشده، دسته‌بندی را پنج مرتبه با هر کدام از این ویژگی‌ها و به وسیله Adtree تکرار کردیم. همان‌طور که نتایج نشان می‌دهد، تمامی ویژگی‌ها در میزان دقت طبقه‌بند دخیل بوده‌اند؛

ولی متمایزکننده‌ترین ویژگی‌ها به ترتیب ویژگی‌های مبتنی بر نویسه، مبتنی بر واژه و ساختاری و کم‌تأثیرترین ویژگی‌ها نیز ویژگی‌های دستوری و روانی‌زبانی بوده‌اند جدول (۵).

به منظور انتخاب ده ویژگی برتر و کاهش فضای ویژگی^۴، از الگوریتم بهره اطلاعات^۵ استفاده شد (نجیب، ۱۹۹۹). این

الگوریتم اهمیت هر ویژگی را با محاسبه بهره اطلاعاتی آن ویژگی و با توجه به طبقه مشخص می‌کند. با انجام دوباره عمل دسته‌بندی با SVM بر روی این ده ویژگی، مشاهده شد که دقت طبقه‌بند ۲٪ افت داشته است. جدول (۶) فهرست ده ویژگی نخست متمایزکننده را نشان می‌دهد.

(جدول-۴): بعضی از نتایج به دست آمده برای زبان انگلیسی

مرجع	داده آموزشی	دسته بند	دقت (%)
(جرمی و همکاران، ۲۰۱۰)	وضعیت‌های فیس بوک	بیز ساده	۶۷/۷
(لای، ۲۰۱۰)	کتاب و ویلاگ	پرسپترون	۵۹
(نا چانگ و همکاران، ۲۰۱۱)	ایمیل‌های انرون گروه خبری رویتر	بیز ساده	۶۹
		SVM	۸۲/۲۳
			۷۶/۷۵

(جدول-۵): دقت طبقه‌بندی با هر کدام از زیرویژگی‌ها

ویژگی	دقت (%)
مبتنی بر واژه	۶۹/۱
روانی-زبانی	۵۵/۹
مبتنی بر نویسه	۶۹/۳
نحوی	۵۸/۹
ساختاری	۶۳/۸
دستوری	۵۴/۴
کل	۷۳/۸

۹- نتیجه‌گیری و بحث

این پژوهش نشان داد که به کمک ویژگی‌های سبک‌شناسی و روان‌شناختی جنسیت نویسنده را می‌توان تشخیص داد. نتیجه طبقه‌بندی بر روی دو پیکره نشان داد که تمایز زبانی مربوط به جنسیت نویسندگان در پیکره مربوط به متون ادبی نسبت به پیکره مربوط به نظرات کاربران بیشتر است. دلیل این عدم تمایز هم ممکن است اقتضای محتوایی متون مربوط به نظرات باشد (از لحاظ حجم متن، عدم امکان بیان احساسات متفاوت، گونه زبانی دو پیکره). از بین ویژگی‌های مختلف طراحی شده، ویژگی‌های واژه‌بنیاد و نشانه‌های زبانی‌روانی تمایزدهنده‌ترین عوامل در دسته‌بندی متون بودند. به عنوان پیشنهاد برای کارهای بعدی، جهت رسیدن به دقت‌های بالاتر می‌توان ویژگی کیسه کلمات^۷ و یا دیگر ویژگی‌های متمایزکننده جنسیت را نیز به مجموعه ویژگی‌ها افزود. همچنین برای رسیدن به نتایج بهتر، از داده‌های آموزشی مربوط به چت‌های فارسی استفاده کرد. برای

^۶Status

^۷Bag of words

^۱Kernel trick

^۲Pearson VII Universal kernel

^۳Cross-validation

^۴Feature space reduction

^۵Information gain

ضمیمه الف:

معیار K یول:

$$Yules K = 10^4 \left(-\frac{1}{N} + \sum_{i=1}^V V_i \left(\frac{i}{N} \right)^2 \right)$$

معیار D سیمپسون:

$$Simpsons D = \sum_{i=1}^V V_i \frac{i}{N} \frac{i-1}{N-1}$$

معیار S سیشل:

$$Sichels S = \frac{\text{count of Hapax Dislegomena}}{V}$$

معیار R هونورس:

$$Honores R = \frac{100 \log_{10} N}{1 - \frac{\text{count of Hapax Legomena}}{V}}$$

V: تعداد کلمات یکتا

V_i : تعداد کلمات یکتا که i مرتبه رخ داده است

N: تعداد کل کلمات

Hapax Dislegomena: کلماتی که فقط دوبار تکرار شده‌اند

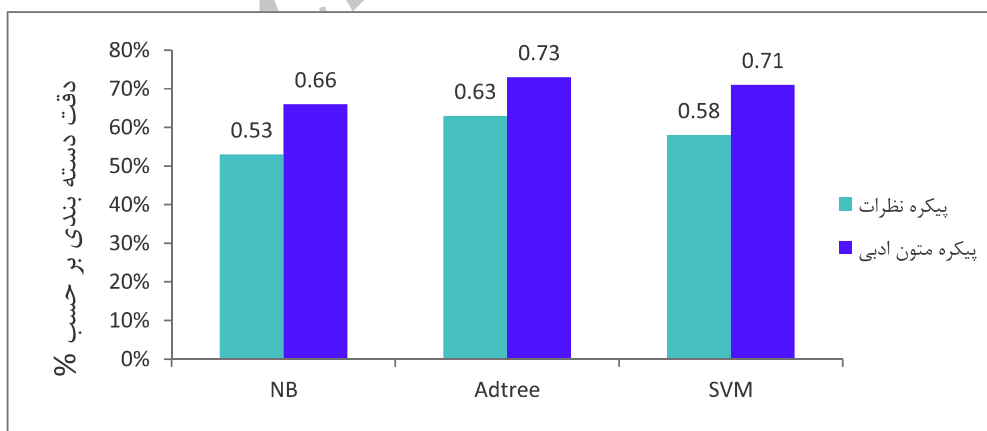
Legomena Hapax: کلماتی که فقط یکبار تکرار شده‌اند

پیشرفت در حوزه شناسایی نویسنده برای زبان فارسی یکی از اولویت‌ها تعریف معیارهای ارزیابی و مقایسه روش‌های مختلف با استفاده از پیکره‌های محک‌زنی^۱ است. تهیه پیکره‌های مختلف زبانی با دامنه‌های (زمانی، موضوعی و غیره) مشخص در پیشرفت کار بسیار تأثیرگذار خواهد بود. دستاوردهای این مقاله عبارت‌اند از:

- ارائه نخستین دادگان دارای برچسب جنسیت برای زبان فارسی (دو پیکره نوشتاری و گفتاری).
- تهیه واژگان دسته‌بندی‌شده زبان فارسی (فهرست کلمات رکیک، رنگ‌ها، کلمات بیان‌کننده تردید و یقین، حروف ربط و غیره) و ارائه فضای ویژگی به کمک این واژگان.
- تهیه فهرست کاملی از اسامی مرد و زن فارسی
- درنهایت ارائه طبقه‌بندی‌هایی جهت طبقه‌بندی متون فارسی از نظر جنسیت نویسنده

(جدول - ۶): ده ویژگی برتر محاسبه‌شده با بهره اطلاعاتی

۱	تعداد نویسه
۲	تعداد کلمات
۳	غنای واژگانی
۴	کلمات دوتکراره
۵	معیار یول
۶	کلمات کوچک‌تر از دو نویسه
۷	تعداد کل حروف الفبا
۸	کلمات یک تکراره
۹	ضمیم متصل اول شخص/ماقی ضمایر
۱۰	تعداد خطوط



(شکل - ۲): الگوریتم دسته‌بندی

^۱Benchmark

ضمیمه ب (لیست کلمات):

قطعیت							
پیوسته	همیشه	همواره	مطمئنا	یقینا	مسلمنا	بدون شک	قطعا
صد در صد	هیچ	حاشا	هیچ وقت	هیچگاه	ابدا	هرگز	همه وقت

تردید							
گمان	حدس	یحتمل	ممکن	گویا	ممکن	احتمالا	شاید
به نظر من	محتمل	احتمالای	فرض	ظن			

حروف ربط							
که	هر چند	هم	هكذا	گرچه	چون	چه	آنچنانکه
زیرا	یعنی	یا	ولو	ولیکن	ولی	و	تا
							اما
							اگر
							لیکن
							لکن

اصوات							
براوو	بالله	به‌به	بادا	آخه	آمین	عجبا	احسنت
زکی	زهاره	یا	خوشا	وای	والسفا	وا	اوخ
جانمی	هورا	هی‌هات	حیف	هان	ای	اوا	ان‌شالله
زرت	دالی	کریمنا	پیشت	مرده‌باد	مرسی	مرحبا	لبیک
کریمنا	دردا	سک سک	آفرین	کاش	دریغا	تبارک الله	آهای

پرسشی							
چی	چرا	چه	چندمین	چند	ایا	چطور	چگونه
کجایی	کجائی	کجا	کدامیک	کدامین	کدام	کی	کی
							آهای
							که

رنگ							
خاکستری	بنفش	برگ سنجدی	ارغوانی	آلبالویی	آبی	پوست پیازی	پسته‌ای
سرخابی	سبز	زیتونی	زرد	دوغی	دودی	دریایی	خرمایی
فیروزه‌ای	عدسی	طلایی	صورتی	شبدیز	سیاه	سفید	سرمه‌ای
قرمز	یاسی	نقره‌ای	نخودی	نارنجی	موشی	مغز	ماشی
							فیلی

۷مین کنفرانس ملی سالانه انجمن کامپیوتر ایران، اسفند ۹۰.

سعیده شجاع رضوی. جنسیت و شمول معنایی، فصل‌نامه پازند. ۱۳۸۶. ۱۰: سال سوم.

محرم اسلامی، مسعود شریفی آتشگاه، صدیقه علیزاده لمجیری، و طاهره زندی. واژگان زبانی فارسی،

۱۰- مراجع

داوری اردکانی، نگار؛ عیار، عطیه. کنکاشی در پژوهش‌های زبان‌شناسی جنسیت، مجله مطالعات راهبردی زنان، ۱۳۸۷، ۴۲-۱۶۲.

درگاه ملی آمار. ۱۳۹۱. <http://www.amar.org.ir>.

زینب فرهمند پور و همکاران. طراحی و پیاده‌سازی یک سامانه هوشمند تشخیص هویت نویسنده‌ی فارسی زبان،

۱۵۵, p.-91-120.

Internet World Usage Statistics, ۲۰۱۲, <http://www.internetworldstats.com/stats.htm>.

Jeremy K, Zach G, David K, Machine Learning and Feature Based Approaches to Gender Classification of Facebook Statuses .۲۰۱۰ .http://thekeesh.com/cs224n/final_writ.

Koppel, M., Schler, J., and Argamon, S. Computational methods in authorship attribution, Journal of American Society Information Science Technology .۶۰, ۱, ۲۰۰۹ .p9-26.

Lai Chao-Yue Author Gender Analysis, Applied Natural Language Processing. ۲۰۱۰ .

Lakoff, R. Language and woman's place . New York : Harper and Row, 1975.

Mason Yoav Freund and Llew The Alternating Decision Tree Algorithm, the 16th International Conference on Machine Learning .۱۹۹۹ .p.124-133.

Mendenhall T. C. A mechanical solution of a literary problem, Popular Science Monthly .

Mendenhall T. C. The characteristic curves of composition, Science . ۹, ۲۱۴ .۱۸۸۷ .p.237-248.

Merriam T. Marlowe's hand in Edward III revisited, Literary and Linguistic .۱۱, ۱, ۱۹۹۶ .p .

Miller Z., Dickinson B, Wei H. Gender Prediction on Twitter Using Stream Algorithms with N-Gram Character Features, International Journal of Intelligence Science . ۲, ۲۰۱۲ .P.143-148.

Mosteller F. and Wallace D. L. Inference and Disputed Authorship: The Federalist, Addison-Wesley, 1964.

Mosteller F, Wallace DL. Applied Bayesian and Classical Inference: the case of the federalist papers, Springer series in statistic :Springer, 1984.

Mukherjee A., Liu B. Improving Gender Classification of Blog Authors, Empirical Methods in Natural Language Processing, ۲۰۱۰.p.۲۱۷-۲۰۷ .

Mulac A., Wiemann, J. M., Widenmann, S. J& ., Gibson, T. W. Male/female language differences and effects in same-sex and mixed-sex dyads: The gender-linked language effect, Communication Monographs .۵۵, ۱۹۸۸ .p.315-335.

Mulac A. Studley LB, Blau S. The gender-linked language effect in primary and secondary students' impromptu essays, Sex Roles .۲۳, ۱۹۹۰ .p-۹ .۱۰ .

مجموعه مقالات اولین کارگاه پژوهشی زبان فارسی و رایانه، تهران، ۱۳۸۳.

محمدرضا فارسیان. جنسیت در واژگان، دانشگاه تهران، ۱۳۷۸.

محمود بی‌جن‌خان. پیکره متنی زبان فارسی، پژوهشکده پردازش هوشمند علایم، ۱۳۸۶.

Apte C, Damerau F, Weiss SM, Apte C, Damerau F, Weiss S Text mining with decision trees and decision rules, The conference on automated learning and discovery, workshop 6: learning from text and the web, 1998.

Baayen H, van Halteren H, Neijt A, Tweedie F. An experiment in authorship attribution, The 6th international conference on the statistical analysis of textual data .St. Malo, France, 2002.

BijanKhan M. The role of the corpus in writing a grammar: An introduction to a software, IranianJ-ournal of Linguistics, 2004.

Burrows J. F, Word-Patterns and Story-Shapes: the Statistical Analysis of Narrative Style, Literary and Linguistic Computing, 1987, p. 61-70.

Chakani Sarvar Language and Gender: A Contrastive View on Farsi and English Discourse Islamic Azad University of Tabriz, 2000.

Cortes C, Vapnik V. Support-vector networks, Machine learning, 1995, p.۹۷-۲۷۳ .

De Vel, O., Anderson, A., Corney, M., and Mohay, G. Mining E-mail content for author identification forensics, SIGMOD Record, .۳۰ : ۴ ۲۰۰۱, p. 55-64.

Diederich J, Kinder Mann J, Leopold E, Paass G. Authorship attribution with support vector machine, Applied Intelligence ;19, .۲۰۰۰ p.۲۳-۱۰۹ .

Diederich, J., Kindermann, J., Leopold, E., and Paass, G. Authorship attribution with support vector machines, Applied Intelligence .۱۹, ۲۰۰۳ .p.۲۳-۱۰۹ .

Gottschalk LA Gleser GC. The measurement of psychological states through the content analysis of verbal behavior .Berkeley, University of California Press, 1969.

Hall, Mark, Eibe Frank, Geoffrey Holmes, Bernhard Pfahringer, Peter Reutemann, and Ian H. Witten. "The WEKA data mining software: an update." ACM SIGKDD explorations newsletter 11, no. 1 (2009): 10-18. Hello kish .۲۰۱۲ .www.hellokish.org.

Holmes D. I ,A stylometric analysis of formon scripture and related texts, Royal Statistical Society,



مهدی مرادی مدرک کارشناسی ارشد زبان‌شناسی رایانشی خود را در سال ۱۳۹۲ از دانشگاه صنعتی شریف دریافت و مدرک کارشناسی خود را نیز از دانشگاه زنجان اخذ کرده است. از ایشان چند مقاله در نشریات داخلی

و همایش‌های بین‌المللی و نیز یک فرهنگ تخصصی به چاپ رسیده است. وی هم‌اکنون مدرس دانشگاه آزاد واحد زنجان بوده و زمینه‌های پژوهشی مورد علاقه وی هستان‌شناسی، تشخیص نویسنده، متن‌کاوی و ساخت پیکره‌های زبانی است. نشانی رایانامه ایشان عبارت است از:

Mehdi.moradi.cl@gmail.com



محمد بحرانی درجه دکترای خود را در رشته هوش مصنوعی در سال ۱۳۸۹ از دانشگاه صنعتی شریف و درجه کارشناسی ارشد خود را از همان دانشگاه در سال ۱۳۸۲ دریافت

کرده است. وی هم‌اکنون عضو هیئت علمی گروه زبان‌شناسی رایانشی در مرکز زبان‌ها و زبان‌شناسی دانشگاه صنعتی شریف است. زمینه‌های پژوهشی مورد علاقه ایشان عبارت است از: زبان‌شناسی رایانشی، پردازش زبان طبیعی، پردازش و بازشناسی گفتار و مدل‌سازی زبانی.

نشانی رایانامه ایشان عبارت است از:

bahrani@sharif.edu

Na Cheng, R. Chandramouli, K.P. Subbalakshmi Author gender identification from text , digitalin vestigation .۲۰۱۱ .p .77-88.

Nacip F.A using information gain as feature weight , The 8th Turkish Symposium on Artificial Intelligence and Nural Networks IAINN .(۹۹'Istanbul, Turkey.۱۹۹۹)

Pennebaker J. W., Mehl, M. R& .,Niederhoffer, K. G. Psychological aspects of natural language use: Our words, our selves., Annual Review of Psychology , ۲۰۰۳

Pennebaker JW, Chung CK, Ireland M, Gonzales A, Booth RJ. The development and psychometric properties of LIWC2007 .Austin, Texas : LIWC Inc, 2007.

Pennebaker, J. W& .,King, L. A. Linguistic styles: Language use as an individualdifference .Personality and Social Psychology, 1999 .P. 6-14.

Safavian SR, Landgrebe D. A survey of decision tree classifier methodology , Cybernetics , ۳ .۱۹۹۱ . ۲۱p .660-674.

Stirman, S.W& .,Pennebaker, J.W. Word use in the poetry of suicidal and non-suicidal poets. , 63, . Psychosomatic Medicine, 2001 - .۶۳ - .p .517-522.

Symantec Internet Security Threat Report, Trends for 2010 .Security Response Publications.۲۰۱۱ , Talbot MM Language and gender: an introduction., Wiley-Blackwell, 1998.

Tweedie FJ, Singh S, Holmes DI. Neural network applications in stylometry: the federalist papers , Computers and the Humanities.۳۰ , ۱ .۱۹۹۶ .

Unger R. K Toward a redefinition of sex and gender , American Psychologist .۳۴ .۱۹۷۹ .p .1085-1094.

Üstün, Bülent, Willem J. Melssen, and Lutgarde MC Buydens. "Facilitating the application of Support Vector Regression by using a universal Pearson VII function based kernel." Chemometrics and Intelligent Laboratory Systems 81, no. 1 .2006. 29-40.

Wodak, Ruth & Benke, Gertrud "Gender as a Variation in Sociolinguistics " in The Handbook of Sociol-inguistics ,Cambridge : Blackwell, 2007.

Yule GU. The statistical study of literary vocabulary, Cambridge University Press, 1944.

Zheng R, Li J, Chen H, Huang Z . A framework for authorship identification of online messages: writing-style features and classification techniques .Journal of the American Society for Information Science and Technology, 2006.