





مرکز آموزشهای الکترونیکی

پایان نامه کارشناسی ارشد مهندسی هوش مصنوعی

شناسایی نویسنده با ویژگی‌های آماری یک سند متنی

نگارنده: سمانه وزیریان

استاد راهنما

دکتر مرتضی زاهدی

استاد مشاور

دکتر حسن پور

شهریور ۱۳۹۵



دفتر مدیریت تحصیلات تکمیلی
فرم شماره (۶)

باسمه تعالی

شماره: ۲۲۹۸
تاریخ: ۹۵/۷/۱۰
ویرایش:

فرم صورت جلسه دفاع از پایان نامه تحصیلی دوره کارشناسی ارشد

با تأییدات خداوند متعال و با استعانت از حضرت ولی عصر (عج) نتیجه ارزیابی جلسه دفاع از پایان نامه کارشناسی ارشد خانم سمانه وزیربان به شماره دانشجویی ۹۲۱۴۸۵۴ رشته کامپیوترگرایش... هوش مصنوعی تحت عنوان: شناسایی نویسنده با ویژگی‌های آماری یک سته متنی که در تاریخ ۹۵/۶/۱۵ با حضور هیأت محترم داوران در دانشگاه شاهرود برگزار گردید به شرح ذیل اعلام می‌گردد:

قبول (با درجه): عالی امتیاز ۱۹ دفاع مجدد مردود

۱- عالی (۲۰ - ۱۹)

۲- بسیار خوب (۱۸/۹۹ - ۱۸)

۳- خوب (۱۷/۹۹ - ۱۶)

۴- قابل قبول (۱۵/۹۹ - ۱۴)

۵- نمره کمتر از ۱۴ غیر قابل قبول

عضو هیأت داوران	نام و نام خانوادگی	مرتبه علمی	اسماء
۱- اسنادارهما	دکتر مرتضی راجدی	استاد یار	
۲- استاد مستاور	دکتر حسن نور	استاد یار	
۳- نماینده شورای تحصیلات تکمیلی	دکتر علی بنی	استاد یار	
۴- استاد ممتحن	دکتر هدی منابخی	استاد یار	
۵- استاد ممتحن	دکتر علی بویار	استاد یار	

رئیس دانشکده علوم تجربیات و یادگار امضاء



از پدر و مادر عزیزم که در همه مراحل زندگی در کنار من و بعنوان حامی حضور داشتند، سپاسگزارم و

از استاد عزیزم جناب آقای مرتضی زاهدی که با نهایت صبر به اینجانب در پیشرفت و به سرانجام

رسیدن این پایان نامه کمک نمودند، نهایت تشکر را دارم.

تعهد نامه

اینجانب **سमानه وزیریان** دانشجوی دوره کارشناسی ارشد رشته مهندسی هوش مصنوعی دانشکده مرکز آموزش های الکترونیکی دانشگاه صنعتی شاهرود نویسنده پایان نامه شناسایی نویسنده با ویژگی های آماری یک سند متنی تحت راهنمایی **دکتر مرتضی زاهدی** متعهد می شوم.

- تحقیقات در این پایان نامه توسط اینجانب انجام شده است و از صحت و اصالت برخوردار است.
- در استفاده از نتایج پژوهش های محققان دیگر به مرجع مورد استفاده استناد شده است.
- مطالب مندرج در پایان نامه تاکنون توسط خود یا فرد دیگری برای دریافت هیچ نوع مدرک یا امتیازی در هیچ جا ارائه نشده است.
- کلیه حقوق معنوی این اثر متعلق به دانشگاه صنعتی شاهرود می باشد و مقالات مستخرج با نام « دانشگاه صنعتی شاهرود » و یا « Shahrood University of Technology » به چاپ خواهد رسید.
- حقوق معنوی تمام افرادی که در به دست آمدن نتایج اصلی پایان نامه تأثیرگذار بوده اند در مقالات مستخرج از پایان نامه رعایت می گردد.
- در کلیه مراحل انجام این پایان نامه، در مواردی که از موجود زنده (یا بافت های آنها) استفاده شده است ضوابط و اصول اخلاقی رعایت شده است.
- در کلیه مراحل انجام این پایان نامه، در مواردی که به حوزه اطلاعات شخصی افراد دسترسی یافته یا استفاده شده است اصل رازداری، ضوابط و اصول اخلاق انسانی رعایت شده است.

تاریخ

امضای دانشجو

مالکیت نتایج و حق نشر

- کلیه حقوق معنوی این اثر و محصولات آن (مقالات مستخرج، کتاب، برنامه های رایانه ای، نرم افزارها و تجهیزات ساخته شده است) متعلق به دانشگاه صنعتی شاهرود می باشد. این مطلب باید به نحو مقتضی در تولیدات علمی مربوطه ذکر شود.
- استفاده از اطلاعات و نتایج موجود در پایان نامه بدون ذکر مرجع مجاز نمی باشد.

چکیده

شناسایی نویسنده یکی از زیر شاخه‌های پردازش متن است و هدف آن مشخص کردن نویسنده یک متن با نویسنده ناشناس است که با عنوان متن دیده نشده از آن یاد می‌شود. انتخاب نویسنده برای متن دیده نشده، از میان مجموعه‌ای از نویسنده‌گان با عنوان نویسنده‌گان کاندید، انجام می‌گیرد. برای طراحی سیستم شناسایی نویسنده نیاز به ویژگی‌های مشخص کننده سبک، و انتخاب روش برای طبقه‌بندی و یا تخصیص نویسنده مناسب برای متن دیده نشده است. معمولاً انتخاب ویژگی از میان داده‌های آموزشی موجود برای هر نویسنده کاندید صورت می‌گیرد.

با افزایش روند تولید متن در زبان‌های مختلف نیاز به توسعه سیستم‌های تشخیص نویسنده که مستقل از زبان نگارش متن عمل نمایند ضروری به نظر می‌رسد. در این پایان‌نامه برای حل مسئله تشخیص نویسنده از n -گرام کارکترها و n -گرام کلمات به عنوان ویژگی‌ها با مزیت استخراج آسان و مستقل از زبان، استفاده شده است. و به عنوان یک روش آماری و احتمالاتی مدل‌سازی زبانی به عنوان روش تخصیص استفاده شده است. برای بهبود نتایج ترکیب 1-گرام کلمات با 2-گرام کلمات همراه با مقدار IDF به عنوان وزن دهی احتمال n -گرام‌ها با عنوان مدل‌سازی زبانی تغییر یافته پیشنهاد شده است. همچنین عوامل موثر بر نتیجه ارزیابی مانند کاهش داده آموزشی و افزایش تعداد نویسنده کاندید و متعادل بودن یا نا متعادل بودن داده آموزشی در حل مسئله تشخیص نویسنده بررسی شده است.

برای ارزیابی از چهار پایگاه داده استفاده شده است. سه پایگاه داده در زبان فارسی و پایگاه داده چهارم در زبان انگلیسی است. روش مدل‌سازی تغییر یافته در تمام آزمایشات انجام شده بهبود را نسبت به n -گرام کاراکترها و n -گرام کلمات نشان می‌دهد. بهترین نتیجه با رسیدن به دقت ۱۰٪ با مدل‌سازی تغییر یافته در پایگاه داده فارسی بهبود خوبی را در استفاده از n -گرام‌ها نشان می‌دهد.

کلمات کلیدی: شناسایی نویسنده، تخصیص نویسنده، انگرام، مدل‌سازی زبانی، پایگاه داده

WMPR-AA2016-B، پایگاه داده WMPR-AA2016-A.

فهرست مطالب

صفحه	عنوان
۱	۱- فصل اول: معرفی و کلیات مسئله
۲	۱-۱ مقدمه
۲	۲-۱ ضرورت حل مسئله شناسایی نویسنده
۵	۳-۱ بیان مسئله
۵	۴-۱ ساختار پایان نامه
۷	۲- فصل دوم: مروری بر کارهای پیشین
۸	۱-۲ مقدمه
۸	۲-۲ ویژگی‌های سبکی
۱۰	۱-۲-۲ ویژگی‌های لغوی
۱۲	۲-۲-۲ ویژگی‌های کاراکتری
۱۳	۳-۲-۲ ویژگی‌های نحوی
۱۴	۴-۲-۲ ویژگی‌های معنایی
۱۵	۵-۲-۲ ویژگی‌های وابسته به کاربرد
۱۵	۳-۲ روش‌های کمی
۱۵	۱-۳-۲ تشکیل داده‌های آموزشی
۱۷	۲-۳-۲ روش‌های تفاضلی
۲۰	۳-۳-۲ روش‌های تولیدی
۲۲	۴-۲ کارهای انجام شده بر روی پایگاه داده‌های فارسی
۲۵	۵-۲ جمع بندی
۲۷	۳- فصل سوم: مباحث نظری و روش پیشنهادی
۲۸	۱-۳ مقدمه
۲۸	۲-۳ پایگاه داده
۲۸	۱-۲-۳ پایگاه داده WMPR-AA2016-A
۲۹	۲-۲-۳ پایگاه داده WMPR-AA2016-B
۳۰	۳-۲-۳ پایگاه داده R40
۳۱	۴-۲-۳ پایگاه داده RCV
۳۳	۳-۳ شناسایی نویسنده
۳۳	۱-۳-۳ پیش پردازش
۳۴	۲-۳-۳ انتخاب ویژگی
۳۵	۳-۳-۳ روش تخصیص
۳۶	۴-۳-۳ مدل سازی زبانی در شناسایی نویسنده
۳۹	۴-۳ روش پیشنهادی
۳۹	۱-۴-۳ استفاده از ضریب IDF به جای حذف کلمات پرتکرار
۴۱	۲-۴-۳ ترکیب ۱-گرام و ۲-گرام کلمات
۴۳	۵-۳ خصوصیات پایگاه داده
۴۶	۶-۳ جمع بندی

۴۷	فصل چهارم: نتایج
۴۸	۱-۴ مقدمه
۴۸	۱-۴ آزمایشات در پایگاه داده WMPR-AA2016-A
۴۹	۱-۱-۴ آزمایش با مدل سازی زبانی و n-گرام کلمات
۵۰	۲-۱-۴ آزمایش با مدل سازی زبانی تغییر یافته
۵۰	۳-۱-۴ آزمایش با مدل سازی و n-گرام کاراکترها
۵۱	۴-۱-۴ بررسی نتایج
۵۲	۲-۴ آزمایشات در پایگاه داده WMPR-AA2016-B
۵۳	۱-۲-۴ آزمایش با مدل سازی زبانی و n-گرام کلمات
۵۴	۲-۲-۴ آزمایش با مدل سازی زبانی تغییر یافته
۵۴	۱-۲-۴ آزمایش با مدل سازی زبانی و n-گرام کاراکترها
۵۵	۲-۴-۴ بررسی نتایج
۵۶	۳-۴ آزمایشات در پایگاه داده R40
۵۷	۱-۳-۴ آزمایش با مدل سازی زبانی و n-گرام کلمات
۵۸	۲-۳-۴ آزمایش با مدل سازی زبانی تغییر یافته
۵۸	۳-۳-۴ آزمایش با مدل سازی زبانی و n-گرام کاراکترها
۵۹	۴-۳-۴ بررسی نتایج
۶۰	۴-۴ زیر مجموعه با شش نویسنده از پایگاه داده RCV
۶۱	۱-۴-۴ آزمایش با مدل سازی زبانی تغییر یافته و مدل سازی زبانی با n-گرام کلمات
۶۲	۲-۴-۴ بررسی نتایج
۶۲	۵-۴ آزمایشات در پایگاه داده RCV
۶۳	۱-۵-۴ آزمایش با مدل سازی زبانی تغییر یافته و مدل سازی زبانی با n-گرام کلمات
۶۳	۲-۵-۴ آزمایش با مدل سازی زبانی و n-گرام کاراکترها
۶۴	۳-۵-۴ بررسی نتایج
۶۵	۶-۴ بررسی خصوصیات پایگاه داده
۶۵	۱-۶-۴ تعداد نویسنده در مجموعه نویسنده‌گان کاندید
۶۶	۲-۶-۴ کاهش داده آموزشی
۷۰	۳-۶-۴ متعادل کردن داده آموزشی
۷۳	۵- فصل پنجم: نتیجه گیری و پیشنهادات
۷۴	۱-۵ جمع بندی
۷۹	۲-۵ فعالیتها
۷۹	۳-۵ پیشنهادات

فهرست شکل‌ها

- شکل ۱-۱: مراحل کار در حل مسئله تشخیص نویسنده ۴
- شکل ۱-۲: طبقه بندی ویژگی‌ها ۹
- شکل ۲-۲: داده آموزشی مبتنی بر پروفایل ۱۶
- شکل ۳-۲: داده آموزشی مبتنی بر نمونه ۱۷
- شکل ۱-۳: توزیع طول متن در پایگاه داده WMPR-AA2016-A ۲۹
- شکل ۲-۳: تولید n -گرامها در یک متن ۳۵
- شکل ۳-۳: توزیع داده‌های آموزشی در بین نویسنده گان کاندید ۴۶
- شکل ۱-۴: تشکیل داده آموزشی و مجموعه داده آزمایشی در پایگاه داده WMPR-AA2016-A ۴۹
- شکل ۲-۴: تشکیل داده آموزشی و مجموعه داده آزمایشی در پایگاه داده R40 ۵۷
- شکل ۳-۴: اثر افزایش تعداد نویسنده گان در میزان دقت ۶۶
- شکل ۴-۴: اثر کاهش داده آموزشی در میزان دقت در پایگاه داده WMPR-AA2016-B ۶۸
- شکل ۵-۴: اثر کاهش داده آموزشی در میزان دقت WMPR-AA2016-A ۶۹
- شکل ۶-۴: تغییرات دقت در مدل سازی در دو حالت متعادل و غیرمتعادل در پایگاه داده WMPR-AA2016 ۷۱
- شکل ۷-۴: تغییرات دقت در مدل سازی در دو حالت متعادل و غیرمتعادل در پایگاه داده WMPR-AA2016-B ۷۱

فهرست جدول‌ها

جدول ۱-۲: لیستی از کارهای پیشین.....	۲۴
جدول ۱-۳: جزئیات پایگاه داده فارسی WMPR-AA2016-A.....	۲۹
جدول ۲-۳: جزئیات پایگاه داده فارسی WMPR-AA2016-B.....	۳۰
جدول ۳-۳: جزئیات پایگاه داده فارسی R40.....	۳۱
جدول ۴-۳: جزئیات پایگاه داده انگلیسی RCV.....	۳۲
جدول ۱-۴: نتایج ارزیابی روش مدلسازی زبانی با n-گرام کلمات در پایگاه داده WMPR-AA2016-A.....	۴۹
جدول ۲-۴: نتایج ارزیابی روش مدل سازی زبانی تغییر یافته در پایگاه داده WMPR-AA2016-A.....	۵۰
جدول ۳-۴: نتایج ارزیابی روش مدل سازی زبانی با n-گرام کاراکترها در پایگاه داده WMPR-AA2016-A.....	۵۱
جدول ۴-۴: رتبه بندی نتایج آزمایشات با مدل سازی زبانی در پایگاه داده WMPR-AA2016-A.....	۵۲
جدول ۵-۴: نتایج ارزیابی روش مدل سازی زبانی با n-گرام کلمات در پایگاه داده WMPR-AA2016-B.....	۵۳
جدول ۶-۴: نتایج ارزیابی روش مدل سازی زبانی تغییر یافته در پایگاه داده WMPR-AA2016-B.....	۵۴
جدول ۷-۴: نتایج ارزیابی روش مدل سازی زبانی با n-گرام کاراکترها در پایگاه داده WMPR-AA2016-B.....	۵۵
جدول ۸-۴: رتبه بندی نتایج آزمایشات با مدل سازی زبانی در پایگاه داده WMPR-AA2016-B.....	۵۶
جدول ۹-۴: نتایج ارزیابی روش مدل سازی زبانی با n-گرام کلمات در پایگاه داده R40.....	۵۷
جدول ۱۰-۴: نتایج ارزیابی روش مدل سازی زبانی تغییر یافته در پایگاه داده R40.....	۵۸
جدول ۱۱-۴: نتایج ارزیابی روش مدل سازی زبانی و n-گرام کاراکترها در پایگاه داده R40.....	۵۹
جدول ۱۲-۴: رتبه بندی نتایج آزمایشات با مدل سازی زبانی در پایگاه داده R40.....	۶۰
جدول ۱۳-۴: نتایج ارزیابی روش مدل سازی زبانی با n-گرام کلمات و مدلسازی زبانی تغییر یافته در زیر مجموعه ۶ عضوی از پایگاه داده RCV.....	۶۱
جدول ۱۴-۴: رتبه بندی نتایج آزمایشات با مدل سازی زبانی در پایگاه داده RCV.....	۶۲
جدول ۱۵-۴: نتایج ارزیابی روش مدل سازی زبانی با n-گرام کلمات و مدل سازی زبانی تغییر یافته در پایگاه داده RCV.....	۶۳
جدول ۱۶-۴: نتایج ارزیابی روش مدل سازی زبانی با n-گرام کاراکترها در پایگاه داده RCV.....	۶۳
جدول ۱۷-۴: رتبه بندی نتایج آزمایشات در پایگاه داده RCV.....	۶۴
جدول ۱۸-۴: دقت در نتایج ارزیابی روش مدل سازی زبانی با کاهش داده آموزشی در پایگاه داده WMPR-AA2016-B.....	۶۸
جدول ۱۹-۴: دقت در نتایج ارزیابی روش مدل سازی با کاهش داده آموزشی در پایگاه داده WMPR-AA2016-A.....	۶۹
شکل ۴-۶: تغییرات دقت در مدل سازی در دو حالت متعادل و غیر متعادل در پایگاه داده WMPR-AA2016-A.....	۷۱
جدول ۱-۵: میانگین دقت در تمامی آزمایشات انجام شده در چهار پایگاه داده.....	۷۷
جدول ۲-۵: مقایسه نتایج مدل سازی زبانی تغییر یافته با پژوهش پیشین.....	۷۸
جدول ۳-۵: مقایسه نتایج مدل سازی زبانی تغییر یافته با پژوهش پیشین.....	۷۸

۱- فصل اول: معرفی و کلیات مسئله

۱-۱ مقدمه

در این پایان نامه به بررسی مسئله شناسایی نویسنده پرداخته شده است. شناسایی نویسنده که به نوعی طبقه بندی متن است در دسته پردازش متن قرار می‌گیرد که خود زیر مجموعه مجموعه بزرگتر پردازش زبان طبیعی است.

شناسایی نویسنده عمل انتساب محتمل‌ترین نویسنده از میان مجموعه نویسنده‌گان شناخته شده^۱، به متنی است که نویسنده آن نامشخص^۲ است. در تعریف دیگر از (Juola, et al., 2006) عمل استنتاج محاسباتی^۳ نویسنده یک متن دیده نشده، با استفاده از اطلاعات آماری که در متن وجود دارد را تخصیص نویسنده^۴ می‌نامند. در طول این پایان‌نامه همواره منظور از شناسایی نویسنده، شناسایی با استفاده از روشهای آماری و محاسباتی است. این فرایند با نام‌های Authorship Attribution, Authorship Identification, Stylometric, non-traditional Authorship Attribution, در منابع مختلف نام برده شده است

۱-۲ ضرورت حل مسئله شناسایی نویسنده

مسئله تشخیص نویسنده از زمان‌های دور اهمیت داشته و مورد بررسی پژوهشگران زمان خود بوده است اما در دهه‌های اخیر (از اواخر ۱۹۹۰) با فراگیر شدن اینترنت، استفاده از صفحات وب و استفاده از محیط‌های چت و همچنین پیشرفتی که در زمینه‌های پردازش زبان طبیعی و الگوریتم‌های یادگیری ماشین انجام شده است توجه و پیشرفت‌های بیشتری در زمینه حل مسئله تشخیص نویسنده صورت گرفته است. اما متأسفانه این تحقیقات در زمینه زبان فارسی چشم‌گیر نبوده است.

^۱ known Authors

^۲ Unknown Author

^۳ Computational Inferring

^۴ Authorship Attribution

پیشرفت پردازش زبان‌های طبیعی و الگوریتم‌های یادگیری ماشین بر حل مسئله شناسایی نویسنده همان‌طور که در (Stamatatos, 2009) به آن اشاره شده است می‌تواند به صورت زیر تاثیر گزار باشد:

- پردازش زبان‌های طبیعی: با گسترش ابزارهای پردازش زبان امکان استفاده از انواع جدید از ویژگی‌های مانند ویژگی‌های نحوی برای آشکار کردن سبک نویسنده فراهم شده است.

- الگوریتم‌های یادگیری ماشین: با قدرتمندتر شدن این الگوریتم‌ها، استفاده آنها به عنوان روش تخصیص و یا طبقه بندی متن نیز با اثر مثبت همراه بوده است.

از طرف دیگر مسئله تشخیص نویسنده به خاطر کاربردهای گسترده آن در حوزه‌های مختلف، روز به روز از اهمیت بیشتری برخوردار می‌شود. از جمله کاربردهای آن می‌توان به موارد زیر اشاره کرد:

- تجزیه و تحلیل نوشته‌ها به صورت قانونی^۱: تجزیه و تحلیل نوشته‌ها به صورت قانونی به منظور پیدا کردن نویسنده در تحقیقات قانونی و بررسی‌های مجرمانه مانند حملات سایبری صورت می‌گیرد.

- تجارت الکترونیکی^۲: بررسی ادعاهای مالکیت نرم‌افزارهای تجاری با بررسی کد برنامه‌ها .

- شناسایی نویسنده متن‌های ادبی (Bozkurt, et al., 2007).

- شناسایی تروریسم با استفاده از شناسایی نویسنده در پیام‌ها

- تشخیص نویسنده ایمیل و متن‌های الکترونیکی (Argamon, et al., 2003).

- تشخیص سرقت ادبی (Alzahrani, et al., 2012).

- و ...

مسئله تشخیص نویسنده، بسته به نحوه تخصیص و هدف از حل مسئله، به چند زیر دسته تقسیم می‌شود؛ در تقسیم بندی که در (Joula, 2008) انجام شده است به سه گروه تقسیم شده است:

^۱ Forensic Analysis

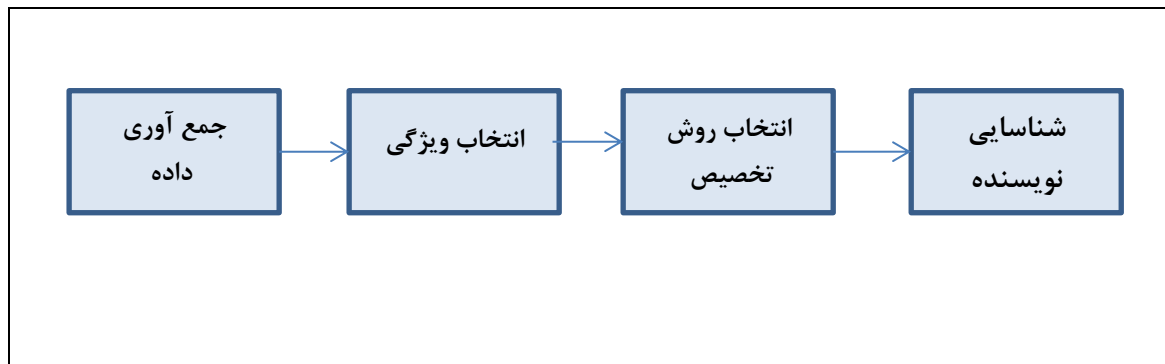
^۲ Electronic Commerce

- مجموعه بسته^۱: انتخاب نویسنده متن از روی مجموعه بسته‌ای از نویسندگان کاندید که برای آنها داده آموزشی وجود دارد.

- مجموعه باز^۲: در این حالت نویسنده متن دیده نشده می‌تواند یکی از نویسندگان مجموعه بسته نویسندگان کاندید با داده آموزشی باشد یا هیچکدام از آنها نباشد.

- پروفایل^۳: در این دسته تمرکز در مشخص کردن ویژگی‌های فردی نویسنده مانند جنسیت است و یا ویژگی متن مانند آیا نوشته به زبان مادری نویسنده نگارش شده است یا متن توسط یک نویسنده نوشته شده است یا چند نویسنده و...

به طور کلی می‌توان گفت در شناسایی نویسنده با دو مسئله روبرو هستیم، یکی ویژگی‌ها که مشخص کننده سبک نوشتاری نویسنده هستند و دیگری روش‌های کمی^۴ که ویژگی‌ها را به منظور انجام انتساب یا تخصیص به کار می‌برند. یک نمای کلی از مراحل کار در حل مسئله شناسایی شکل ۱-۱ نمایش داده شده است.



شکل ۱-۱: مراحل کار در حل مسئله تشخیص نویسنده

^۱ Close- set Authorship Attribution

^۲ Open- set Authorship Attribution

^۳ Profiling

^۴ Quantitative Methods

۱-۳ بیان مسئله

همان‌طور که عنوان شد عمل شناسایی نویسنده در مجموعه بسته، سعی بر شناسایی نویسنده و یا تخصیص نویسنده در متن مورد مناقشه از میان مجموعه نویسندگان شناخته شده دارد.

در این تعریف منظور از متن مورد مناقشه یا متن دیده نشده یا متن مورد سوال یا متن با نویسنده نامشخص، متنی است که نویسنده آن مشخص نیست. منظور از نویسندگان شناخته شده یا مجموعه نویسندگان کاندید، مجموعه‌ای از نویسندگان است که برای آنها داده آموزشی وجود دارد یا به عبارتی شناخته شده هستند و منظور از حل مسئله شناسایی نویسنده تخصیص یکی از نویسندگان مجموعه کاندید به متن دیده نشده است. به این منظور به طور معمول از یک روش تخصیص و گروهی از ویژگی‌ها، اعم از لغوی، نحوی، معنایی، و... یا ترکیب ویژگی‌ها با یکدیگر استفاده می‌شود.

با افزایش روند تولید متن در زبان‌های مختلف نیاز به توسعه سیستم‌های تشخیص نویسنده که مستقل از زبان نگارش متن عمل نمایند ضروری به نظر می‌رسد. در این پایان‌نامه برای حل مسئله تشخیص نویسنده از n -گرام کارکترها و n -گرام کلمات به عنوان ویژگی‌ها با مزیت استخراج آسان و مستقل از زبان، استفاده شده است. همچنین به عنوان یک روش آماری، روش احتمالاتی مدل‌سازی زبانی به عنوان روش تخصیص استفاده شده است. برای بهبود نتایج ترکیب 1-گرام کلمات با 2-گرام کلمات و استفاده همراه با مقدار IDF به عنوان وزن دهی احتمال n -گرام‌ها با عنوان مدل سازی زبانی تغییر یافته پیشنهاد شده است.

در این مطالعه از چهار پایگاه داده، استفاده شده است که دو پایگاه داده آن توسط نگارنده جمع‌آوری شده است. پایگاه داده اول مجموعه رباعیات شش شاعر در زبان فارسی استفاده شده است. پایگاه داده دوم از مجموعه غزلیات هفت شاعر در زبان فارسی تهیه شده است. پایگاه داده سوم مجموعه متن‌های ادبی جمع‌آوری شده از چهل نویسنده فارسی زبان است و پایگاه داده چهارم در زبان انگلیسی و دارای پنجاه نویسنده کاندید است. همچنین در پایان عوامل موثر بر نتیجه ارزیابی مانند کاهش داده

آموزشی و افزایش تعداد نویسندگان و متعادل بودن یا نامتعادل بودن داده آموزشی در حل مسئله تشخیص نویسندگان بررسی شده است.

۴-۱ ساختار پایان نامه

در فصل اول تعریف و مقدمه‌ای از شناسایی نویسندگان و کاربردهای آن گفته شده است. در فصل دوم به دسته بندی ویژگی‌ها و متدهای تخصیص استفاده شده پرداخته شده و نمونه پژوهش‌های گذشته مرور شده است. در فصل سوم به معرفی چهار پایگاه داده، خصوصیات پایگاه داده، روش و ویژگی انتخاب شده در این پایان نامه پرداخته شده است. در فصل چهارم نتایج آزمایشات و ارزیابی سیستم گزارش شده است و در فصل پنجم جمع بندی و نتیجه گیری انجام شده است.

۲- فصل دوم: مروری بر کارهای پیشین

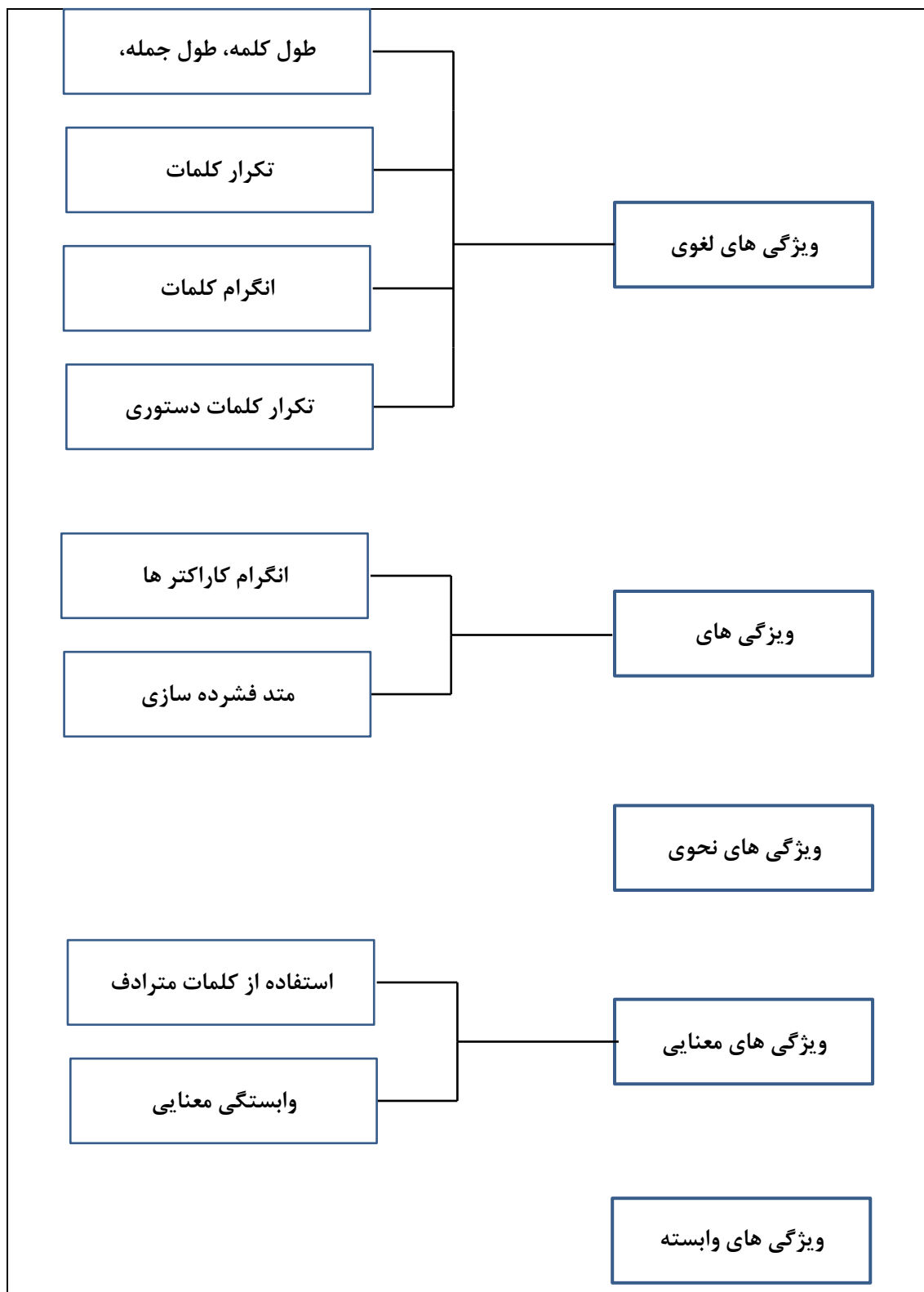
۱-۲ مقدمه

در این فصل همراه با مروری بر نمونه پژوهش‌های پیشین در حل مسئله نویسنده، دسته بندی ویژگی‌ها و روش‌های تخصیص به کار گرفته شده در این زمینه، نیز انجام شده است. فصل ابتدا با توضیح ویژگی‌های سبکی در بخش ۲-۲ شروع می‌شود. در بخش ۳-۲ پس از تعریفی در رابطه با نحوه تشکیل داده‌های آموزشی در روش‌های مختلف به متدهای تخصیص پرداخته شده است. و در انتها در بخش ۴-۲ مروری بر کارهای انجام شده بر روی پایگاه داده با زبان فارسی انجام شده است.

۲-۲ ویژگی‌های سبکی

ویژگی‌های سبکی^۱ ویژگی‌های معرفی کننده سبک نویسنده به صورت کمی هستند. در ادامه به طبقه بندی انواع ویژگی‌ها پرداخته شده است و همچنین به نمونه کارهای انجام شده در این راستا، اشاره شده است. بر مبنای مطالعات انجام شده در (Stamatatos, 2009) ویژگی‌ها را می‌توان به پنج دسته اصلی تقسیم بندی کرد. این دسته بندی در شکل ۱-۲ نمایش داده شده است.

^۱ Stylometric Features



شکل ۲-۱: طبقه بندی ویژگی‌ها

۱-۲-۲ ویژگی‌های لغوی

در دسته بندی ویژگی‌های لغوی^۱ متن به صورت رشته‌ای از توالی توکن‌ها در نظر گرفته شده است. هر توکن در آن می‌تواند یک کلمه یا یک عدد و یا یک علامت نگارشی باشد و حتی در تقسیم بندی بزرگتر می‌تواند به صورت جمله تقسیم بندی شود. با این دید به متن، از جمله ویژگی که در این گروه قرار می‌گیرند طول کلمه، طول جمله، غنی بودن از لحاظ لغت، (که در ساده‌ترین روش از نسبت بین تعداد لغات بدون تکرار بر کل لغات متن به دست می‌آید)، است. در اولین کارهای انجام شده در زمینه شناسایی نویسنده تنها از یک متغیر آماری^۲ مانند میانگین طول کلمات یا میانگین طول جملات استفاده شده است (Holmes, 1998).

ویژگی بعدی در این گروه، تکرار کلمات^۳ و n -گرام کلمات^۴ است. مزیت ویژگی‌ها در این دسته سادگی محاسبه آنها است؛ آنها برای محاسبه نیاز به ابزارهای پردازش زبان طبیعی پیشرفته‌ای ندارند و در اغلب موارد با ابزارهای ساده‌ای مثل جداکننده توکن^۵ بر اساس فاصله به دست می‌آیند. این مزیت آنها را تبدیل به ویژگی مستقل از زبان می‌کند. یعنی در زبان‌های مختلف ابزار محاسبه آن یکسان است (Keselj, et al., 2003). البته در این میان استثناهایی نیز وجود دارد مثل زبان چینی که در آنها از فاصله برای جدا کردن کلمات استفاده نمی‌شود و برای جدا کردن کلمات به ابزارهای پیچیده‌تری نیاز است. از نمونه کارها با این ویژگی می‌توان به کار (Stamatatos, 2006) اشاره کرد که در آن از هزار کلمه پر تکرار به عنوان ویژگی استفاده شده است و همچنین (Howedi & Mohd, 2014) که از n -گرام کلمه استفاده کرده است و یا (Luyckx & Daelemans, 2011) 3 -گرام کلمات را به کار برده است و کار (Mikros & Perifanos, 2013) که تلفیق n -گرام کاراکتر با n -گرام کلمه را به عنوان ویژگی انتخاب کرده است همچنین در کار (Rappoport & Koppel, 2013) از الگوی توالی کلمات با

¹ Lexical Features

² Univariate Statistic

³ Word Frequency

⁴ Word n -gram

⁵ Tokenizer

شروع و پایان با کلمات پرتکرار به عنوان ویژگی تعریف شده است. به این صورت که توالی کلمات منطبق بر الگوی یاد شده از داده‌های آموزشی استخراج می‌شود و داده‌های آزمایشی براساس تطبیق با ویژگی‌های استخراج شده دسته‌بندی می‌گردد.

ویژگی دیگر که می‌توان آن را نوعی ویژگی تکرار کلمات نیز در نظر گرفت استفاده از کلمات دستوری^۱ است. استفاده از کلمات دستوری از این جهت مورد توجه قرار گرفته است که به محتوای متن و موضوع بستگی ندارند و از این رو گمان می‌رود که استفاده از آنها توسط نویسندگان، خارج از اراده است و در نتیجه می‌تواند نمایش دهنده سبک نوشتاری نویسنده باشند (Stamatatos, 2009). به عنوان یکی از شناخته شده‌ترین و اولین کارها که با استفاده از روش‌های آماری به حل مسئله نویسنده پرداخته است کار Mosteller و Wallace (Mosteller, et al., 1964) از این دسته از ویژگی‌ها استفاده کرده است. آنها سعی داشته‌اند تا با استفاده از کلمات دستوری، دوازده عدد از مقالات فدرالی را که نویسنده آنها مشخص نبوده است را به سه نویسنده (Alexander Hamilton, James و John Jay و Adison) تخصیص دهند. بعد از آن در کارهای اولیه بعدی از چند متغیره آماری^۲ استفاده شده است از جمله (Burrows & F, 1992) که در آن از هفتاد و پنج کلمات دستوری پرتکرار به عنوان متغیر آماری انتخاب شده است و سپس با استفاده از PCA متغیرها کاهش داده شده است. از دیگر کارها در این گروه، می‌توان به کار (Boukhaled & Ganascia, 2015) اشاره کرد که در آن با مقایسه استفاده از تکرار کلمات دستوری در مقابل ویژگی‌های استخراج شده با روش قوانین پی‌درپی^۳، استفاده از کلمات دستوری را کارا تر عنوان کرده است.

¹ Function Word

² Multivariate Statistic

³ Sequential Rule

۲-۲-۲ ویژگی‌های کاراکتری

در دسته ویژگی‌های کاراکتری^۱ متن به صورت رشته‌ای از توالی کاراکترها در نظر گرفته می‌شود. از جمله ویژگی‌هایی که در این گروه قرار می‌گیرند نوع کاراکتر (عددی، حرفی و...)، تکرار کاراکترها، n-گرام در سطح کاراکتر^۲ و متدهای فشرده سازی است. با وجود کارهای زیادی که با استفاده از کاراکتر n-گرام‌ها انجام شده است و نتایج خوبی که به دست آمده است مفید بودن آنها نشان داده شده است. مزیت استفاده از کاراکتر n-گرام در این است که اختلافات جزئی در سبک نوشتاری را مشخص می‌کند و تحمل پذیری خوبی در برابر نویز در متن (متن نویزی متنی است که دارای خطاهای گرامری است و علائم نگارشی در آن درست استفاده نشده است) دارد. اما مشکلی که وجود دارد در تعیین مقدار n به صورت کارا است. یک مقدار بزرگ برای n باعث در بر گرفتن ویژگی‌های موضوعی علاوه بر ویژگی‌های متنی می‌شود و از طرفی انتخاب مقدار کوچک n در برگیرنده ویژگی‌های متنی به طور کافی نخواهد بود. بهترین انتخاب برای مقدار n وابسته به زبان و متد استفاده شده است ((Stamatatos, 2009), (Keselj, et al., 2003)).

از نمونه کارهای انجام شده در این زمینه، کار (Howedi & Mohd, 2014) است که در آن از کاراکتر n-گرام‌ها ی ۱، ۲ و ۳ به صورت جداگانه استفاده شده است و یا (Stamatatos, 2006) که از کنار هم قرار دادن کاراکتر n-گرام‌ها با مقدار n، ۳، ۴ و ۵ با تعداد تکرارهای مشخص در یک بردار ویژگی، استفاده شده است. و همچنین کار (Escalante, 2011) که از هیستوگرام n-گرام کاراکترها استفاده کرده است.

¹ Character Features

² Character n-gram

۲-۲-۳ ویژگی‌های نحوی

در گروه ویژگی‌های نحوی^۱ از ویژگی‌های ساختاری و نحوی متن، مانند نقش کلمات به عنوان ویژگی استفاده می‌گردد. علت قرار گرفتن این دسته از ویژگی‌های متنی در گروه ویژگی‌های سبکی، این عقیده است که نویسندگان الگوهای نحوی را به صورت ناخودآگاه استفاده می‌کنند و بنابراین در تمامی نوشته‌های خود الگوهای مشابه‌ای را دنبال خواهند کرد (Stamatatos, 2009).

یکی از ساده‌ترین حالات آن استفاده از یک برچسب زن گفتار^۲ است. از نمونه کاربرد آن می‌توان از کار (Diederich, et al., 2003) و (Zhao & Zobel, 2007) نام برد، در کار اول با استفاده از یک برچسب زن گفتار ساختار کلمات در سه گروه اسم، فعل، صفت و کلمات دستوری مشخص شده است و سپس بردار ویژگی متشکل از سه زیر بردار تکرار برچسب ها، تکرار ۲-گرام برچسب ها و تکرار کلمات با طول متفاوت ایجاد شده است. در کار دوم برچسب گفتار و ۲-گرام برچسب ها از جمله ویژگی‌های به کار رفته است.

حالت بعدی مربوط به استفاده از ابزارهای پردازش متن برای تشخیص جمله و عبارات اسمی و فعلی است. مانند کار (Luyckx & Daelemans, 2005) که در آن عبارات اسمی یکی از نه گروه ویژگی است که مورد استفاده قرار داده است. دامنه استفاده از این ویژگی‌ها تا استفاده از میزان خطای نحوی، وابستگی نحوی توکن‌ها و بسیاری موارد دیگر به عنوان ویژگی ادامه پیدا کرده است.

مشکلی که در رابطه با این دسته از ویژگی‌ها وجود دارد، این است که برای استخراج این گروه از ویژگی‌ها نیاز به ابزارهای پردازش متن قوی و دقیق است و به این ترتیب مسئله استخراج این ویژگی‌ها از زبانی به زبان دیگر متفاوت می‌شود و به عبارت دیگر استخراج آنها را وابسته به زبان کرده است. علاوه بر آن به خاطر خطای حاصل از استفاده از ابزارهای پردازش و تجزیه کننده متن خطای استفاده از ابزارها به سیستم تحمیل می‌گردد (Stamatatos, 2009).

¹ Syntactic Features

² Part Of Speech Tagger

۲-۲-۴ ویژگی‌های معنایی

در ویژگی‌های معنایی^۱ توجه روی معنا و نقش معنایی کلمات است که در سطح کلمات عبارات و یا جملات وجود دارد. برای روشن شدن نقش ویژگی معنایی، جمله "درب به علی لگد زد" را در نظر بگیرید این جمله از لحاظ گرامری و نوشتاری بدون خطا است اما چیزی که در این جمله اشتباه است معنا است. این که کلمه "درب" با وجود معنایی که دارد نمی‌تواند در این جایگاه نحوی قرار گیرد. می‌توان از تجزیه‌گر معنایی در مشخص کردن ویژگی‌هایی مانند لغات مترادف، متضاد، لغات با چند معنا و وابستگی معنایی^۲ استفاده کرد. اما همچنان که گفته شد برای استخراج این ویژگی‌ها نیاز به ابزارهای پارس و تحلیل معنایی است. در اکثر کارها از ترکیب این ویژگی‌ها با ویژگی‌های نحوی و لغوی استفاده شده است. به عنوان نمونه در کار (Gamon, 2004) از ویژگی‌های معنایی استخراج شده از گراف وابستگی معنایی، در ترکیب با دو ویژگی نحوی و کلمات دستوری استفاده شده است. از دیگر نمونه‌های ویژگی معنایی استفاده از تکنیک نمایه سازی معنایی نهفته^۳ است. در این تکنیک ابتدا ماتریس لفظ-سند^۴ تشکیل می‌گردد (در واقع هر سطر از این ماتریس بردار ویژگی هر سند می-باشد) در مرحله بعد با متد^۵ SVD، فاکتورگیری روی ماتریس انجام می‌شود سپس با استفاده از مقادیر ویژه به دست آمده ماتریسی معادل ماتریس اول ساخته می‌شود و در نهایت با استفاده از متدهای اندازه‌گیری فاصله برداری، شباهت‌سندها در فضای جدید سنجیده می‌شود. ایده اصلی در این تکنیک بر اساس نگاشت بردار ویژگی‌ها به فضای برداری با تعداد بعد کمتر است به طوری که ویژگی‌ها در بعد جدید آشکارکننده ارتباطات معنایی مورد نظر باشند. از کارهای انجام شده در این زمینه می‌توان به کارهای (Soboroff, et al., 1997) و (Satyam, et al., 2014) اشاره کرد که از کار کتر n -گرام‌ها برای تشکیل ماتریس لفظ-سند استفاده کرده است.

¹ Semantic Features

² Semantic Dependencies

³ Latent Semantic Indexing

⁴ Term-document

⁵ Singular Value Decomposition

۵-۲-۲ ویژگی‌های وابسته به کاربرد^۱

بسته به کاربرد مسئله تشخیص نویسنده در حوزه‌های مختلف مانند فروم‌ها و یا پیام‌های الکترونیکی می‌توان از ویژگی‌های مختص آن حوزه استفاده کرد. مانند کار (Abbasi & Chen, 2005) که در آن از ویژگی‌هایی از قبیل سائز و رنگ فونت استفاده شده است.

۳-۲ روش‌های کمی

پس از بررسی ویژگی‌های مختلف استفاده شده در حل مسئله تخصیص نویسنده در این بخش به دسته بندی و بررسی متدهای تخصیص در این حوزه پرداخته خواهد شد. اما قبل از آن تعریف روش‌های تشکیل و به کارگیری داده‌های آموزشی در روش‌های مختلف تخصیص پرداخته شده است.

۱-۳-۲ تشکیل داده‌های آموزشی

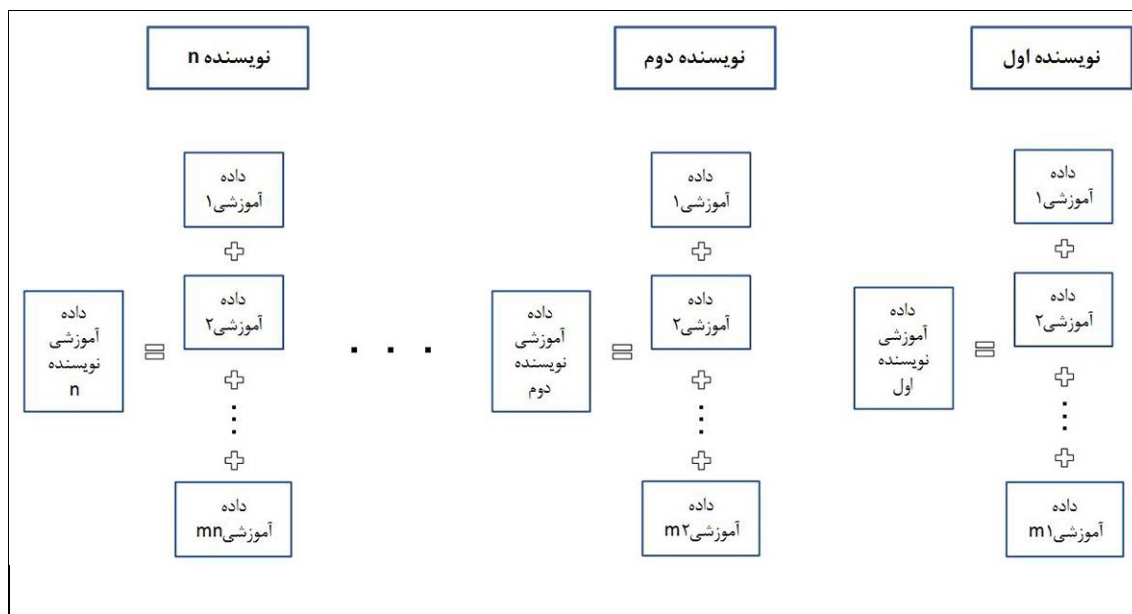
به طور کلی در روش‌های مختلف تخصیص، داده آموزشی نویسنده‌گان کاندید به دو صورت تشکیل و استفاده می‌شود:

- مبتنی بر پروفایل^۲

تمام متن‌های آموزشی هر نویسنده در یک فایل الحاق می‌گردد به صورتی که برای هر نویسنده در مجموعه نویسنده‌گان کاندید یک فایل متنی با عنوان فایل آموزشی وجود دارد و ویژگی‌های سبکی از فایل الحاق شده استخراج خواهند شد (شکل ۲-۲).

¹ Spplication Specific

² Profile-Based



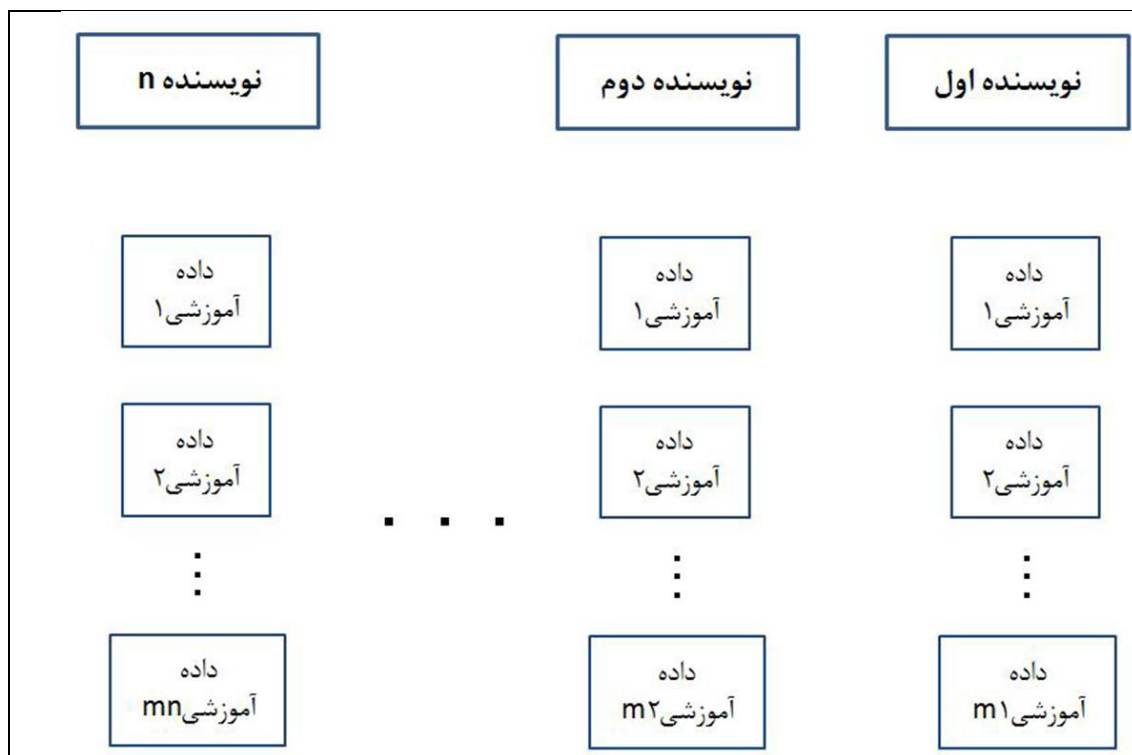
شکل ۲-۲: داده آموزشی مبتنی بر پروفایل

▪ مبتنی بر نمونه^۱

هر یک از فایل‌های متنی موجود برای هر نویسنده به صورت جداگانه در تخصیص شرکت می‌کنند. در نتیجه برای هر نویسنده در مجموعه کاندید می‌تواند چند فایل آموزشی وجود داشته باشد (شکل ۲-۳). در ادامه به معرفی متدهای تخصیص که در حل مسئله شناسایی نویسنده به کار رفته است می‌-

پردازیم.

^۱ Instance-Based



شکل ۲-۳: داده آموزشی مبتنی بر نمونه

۲-۳-۲ روش‌های تفاضلی^۱

در این دسته از روش‌ها برای طبقه بندی کلاس‌ها (در این جا متن‌ها) با استفاده از احتمال شرطی $P(c|x)$ به طور مستقیم نمونه x را به کلاس c نگاشت می‌کنند. (Jordan, 2002) تشکیل بردار آموزشی در این روش‌ها به هر دو صورت مبتنی بر نمونه و مبتنی بر پروفایل صورت گرفته است.

۱-۲-۳-۲ روش‌ها بر اساس فاصله

در روش‌های بر اساس فاصله^۲ از یکی از معیارهای اندازه‌گیری فاصله برداری به عنوان روش تخصیص استفاده می‌گردد. از نمونه کارهایی که در این دسته قرار می‌گیرد می‌توان به کار (Stamatatos, et al., 2000) اشاره کرد که در آن با استخراج بیست و دو متغیر و تشکیل بردار ویژگی برای متن دیده نشده و اندازه‌گیری فاصله مالهونوبی^۳ از مرکز هر گروه از داده‌های آموزشی، متن دیده نشده به یکی از

¹ Discriminative Methods

² Distance Method

³ Mahalonobi

نویسنده‌گان کاندید تخصیص داده می‌شود. به عنوان نمونه دیگر در کار (Chaski, 2001) با استفاده از آزمون مربع کای¹ متن‌های دیده نشده در بین چهار نویسنده کاندید دسته‌بندی شده است. در (Keselj, et al., 2003) برداری از تکرار کارکتر n-گرام‌هایی که بیشترین تکرار را دارند برای هر نویسنده کاندید با استفاده از داده‌های آموزشی تشکیل می‌شود و همین‌طور برداری مشابه برای هر یک از متن‌های دیده نشده آزمایشی تشکیل می‌شود سپس با اندازه‌گیری یک معیار فاصله بر اساس مجموع مربعات، میزان شباهت بین دو بردار سنجیده می‌شود و متن دیده نشده به نویسنده‌ای که با آن کمترین فاصله یا بیشترین شباهت را دارد تخصیص می‌یابد. مشکل این روش این‌طور بیان شده است که اگر داده آموزشی مربوط به یک نویسنده کوچک باشد باعث می‌شود که فاصله بین بردار آموزشی نویسنده با متن کوچک تر در مقایسه با سایر بردارهای آموزشی کمتر شود.

در کار (Frantzeskou, et al., 2006) مشابه کار قبل بردار ویژگی بر اساس تکرار کارکتر n-گرام‌هایی که بیشترین تکرار را دارند برای داده‌های آموزشی و آزمایشی تشکیل شده است اما این بار به عنوان معیار فاصله از اندازه اشتراک بردار داده آموزشی و داده آزمایشی استفاده شده است و در نهایت متن دیده نشده به نویسنده‌ای که بزرگترین معیار فاصله را با آن دارد تخصیص می‌یابد.

در (Stamatatos, 2007) ابتدا یک متن از الحاق تمام متن‌های آموزشی مربوط به تمام نویسنده‌گان کاندید ایجاد شده است و برداری از تکرار کارکتر n-گرام‌ها برای متن الحاقی تشکیل می‌شود و بردار داده نرمال خوانده می‌شود؛ سپس اختلاف بین تکرار n-گرام‌ها در متن دیده نشده با بردار نرمال محاسبه شده است و به عنوان وزن به معیار فاصله در کار (Keselj, et al., 2003) اضافه می‌گردد؛ در نهایت متن دیده نشده به نویسنده‌ای که با آن کمترین فاصله یا بیشترین شباهت را دارد تخصیص می‌یابد. در (Savoy, 2012) با در نظر گرفتن یک توزیع دو جمله‌ای روی کل کلمات برای هر کلمه وزنی را بر اساس امتیاز² محاسبه شده است سپس به عنوان معیار فاصله بین متن دیده نشده و متن

¹ Chi-squared test

² Z score

آموزشی، میانگین اختلاف امتیاز z برای کلمات موجود در متن دیده نشده و کلمه متناظر آن در داده آموزشی محاسبه می گردد. در نهایت متن دیده نشده به نویسنده‌ای که با داده آموزشی مربوط به آن کمترین فاصله را دارد تخصیص می‌یابد.

۲-۲-۳-۲ روش‌های فشرده سازی

از نمونه کارهای انجام شده در این زمینه می‌توان به کار (Marton, et al., 2005) اشاره کرد. این روش را می‌توان به صورت مراحل زیر توصیف کرد:

۱- داده آموزشی به صورت مبتنی بر پروفایل برای هر یک از نویسندگان کاندید تشکیل می‌گردد.

۲- به وسیله یک روش فشرده‌سازی فایل آموزشی فشرده می‌گردد.

۳- در این مرحله فایل متن دیده نشده به فایل داده آموزشی هر نویسنده کاندید اضافه می‌گردد. و متد فشرده سازی بر فایل حاصل اعمال می‌گردد.

۴- اختلاف بیتی بین دو فایل فشرده ایجاد شده در مرحله دو و مرحله سه محاسبه می‌شود و نویسنده با کمترین اختلاف، نویسنده برگزیده برای متن دیده نشده خواهد بود.

۲-۲-۳-۲ روش‌های یادگیری ماشین

از جمله رهیافت‌ها در دسته روش‌های یادگیری ماشین^۱ می‌توان به موارد زیر اشاره کرد:

بردار ماشین : ۲۰۰۱ (Vel, et al., 2001)، ۲۰۰۳ (Diederich, et al., 2003)، ۲۰۰۳ (Koppel & Schler, 2003)، ۲۰۰۴ (Gamon, 2004)، ۲۰۰۴ (BEKKERMAN & ALLAN, 2004)، ۲۰۰۵ (Abbasi & Chen, 2005)، ۲۰۰۶ (Zheng, et al., 2006)، ۲۰۰۷ (Türkoğlu, et al., 2007)، ۲۰۰۷ (Pavelec, et al., 2007)، ۲۰۰۸ (Stamatatos, 2008)، ۲۰۱۱ (Silva, et al., 2011)، ۲۰۱۱ (Kourtis, et al., 2011)، ۲۰۱۳ (Sidorov, et al., 2013)، ۲۰۱۵ (Nagaprasad, et al., 2015).

¹ Machine Learning Methods

درخت تصمیم^۱: ۲۰۰۴ (Ramya, et al., 2004)، ۲۰۰۵ (Zhao & Zobel, 2005)، ۲۰۰۵ (Uzuner ۲۰۰۵، (Zhang, et al., 2006)، ۲۰۰۵ & Katz, 2005)، ۲۰۰۶ (Abdalla, et al., 2013)، ۲۰۱۳ (Koppe, et al., 2009) شبکه عصبی^۲: ۱۹۹۴ (KJELL, 1994)، ۱۹۹۶ (Tweedie, et al., 1996)، ۲۰۰۴ (Ramya, et al., 2004)، ۲۰۰۶ (Zheng, et al., 2006)، ۲۰۰۸ (Tearle, et al., 2008)، ۲۰۱۲ (Jamak, et al., 2012)، ۲۰۱۲ (Selman, et al., 2012).

۲-۳-۳ روشهای تولیدی^۳

در روشهای تولیدی از احتمال توام کلاس c و متغیر x ($P(c, x)$) استفاده می‌شود و احتمال شرطی $P(c|x)$ با استفاده از قانون بیز محاسبه می‌گردد (Jordan, 2002). و داده‌های آموزشی عموماً به صورت مبتنی بر پروفایل استفاده می‌شوند.

۲-۳-۳-۱ مدل احتمالاتی بیز

در این مدل از طبقه بند بیز استفاده شده است. از نمونه کارهای انجام شده در این زمینه می‌توان به کار (Peng & Schuurmans, 2003) اشاره کرد که در آن با استفاده از تکنیک‌های هموار کردن مدل-سازی زبانی و در نظر گرفتن وابستگی بین n -گرام‌ها بیز ساده را بهبود داده است. در کار (Altheneyan, 2014) از بیز روی بانک اطلاعاتی به زبان عربی استفاده شده است و از چهار مدل بیز متفاوت در ویژگی‌ها استفاده شده است. در کار (Boutwell, 2014) از کارکتر n -گرام‌های دو، سه، چهار...شش و بیز ساده استفاده شده است و در بانک اطلاعاتی جمع آوری شده از پیام‌های کوتاه به عنوان یک پایگاه داده متن کوتاه اعمال شده است. همچنین در این کار با الحاق چند متن

¹ Decision Tree

² Neural Network

³ Generative Method

کوتاه در یک متن و ایجاد متن بزرگتر توانسته با فراهم کردن اطلاعات زبانی بیشتر به میزان دقت بالاتری برسد.

۲-۳-۳-۲ مدل سازی زبانی

از آنجا که مدل سازی^۱ زبانی روش انتخابی روی این تحقیق است در بخش سوم به صورت کامل توضیح داده خواهد شد. در این قسمت تنها به ذکر چند نمونه از کارهای انجام شده در این زمینه می پردازیم. در اکثر کارها داده آموزشی به صورت پروفایل استفاده شده است و داده آموزشی هر نویسنده به صورت یک مدل زبانی در نظر گرفته شده است و احتمال هر متن با نویسنده ناشناس بر اساس مدل سازی زبانی اندازه گیری می شود و سپس با معیارهایی مانند پرپلکسیتی^۲ یا KLD با یکدیگر مقایسه می شوند و نویسنده با بهترین مقدار به عنوان نویسنده متن ناشناس انتخاب می گردد. به عنوان یکی از اولین کارهای انجام شده در این زمینه کار (Peng & Schuurmans, 2003) است. در این کار از کاراکتر n-گرامها به عنوان ویژگی استفاده شده است و از مدل سازی زبانی برای طبقه بندی استفاده شده است. بهترین نتیجه برای کاراکترگرامها با مقدار n ۲ و ۳ به دست آمده است. برای نمایش دقت روش ارائه شده از چهار پایگاه داده در سه زبان انگلیسی، گریک و چینی استفاده شده است و بهترین نتیجه را در پایگاه داده گریک^۳ گزارش شده است. در کار (Zhao, et al., 2006) از مدل سازی زبانی ساده و کلمات دستوری، تگ های دستوری و جداکننده ها هر یک به صورت جداگانه و همچنین به صورت تلفیق هر سه به عنوان ویژگی استفاده شده است و با مقایسه مقدار معیار KLD^۴ محاسبه شده برای متن دیده نشده با هر یک از نویسنده گان مجموعه آموزشی، نویسنده با کمترین مقدار را به عنوان نویسنده متن ناشناس انتخاب کرده است. برای سنجش میزان دقت از سه بانک اطلاعاتی در زبان انگلیسی استفاده شده است. در کار (Azarbondy, et al., 2015) چالش اصلی

¹ Language Modeling

² Perplexity

³ Greek

⁴ Kullback-Leibler Divergence

تشخیص نویسنده در متن‌های کوتاه با در نظر گرفتن تغییرات سبک نوشتاری هر نویسنده در گذر زمان است. در این مدل دوره زمانی هر نویسنده کاندید به چندین دوره با طول ثابت تقسیم شده است و سپس با استفاده از مدل سازی زبانی در سطح کاراکتر همراه با ضربی از تابع تغییرات سبک نوشتاری در طول دوره زمانی نویسنده، احتمال تولید متن دیده نشده توسط هر نویسنده محاسبه می‌گردد و در نهایت نویسنده کاندید با بالاترین احتمال به عنوان نویسنده متن دیده نشده انتخاب می‌گردد.

در بعضی از کارها از معیارهای هموارسازی در مدل سازی زبانی برای بهبود روش‌های احتمالاتی مانند بیز استفاده شده است. در این زمینه می‌توان به کار (Peng & Schuurmans, 2003) اشاره کرد که در آن با استفاده از تکنیک‌های هموار کردن مدل سازی زبانی و در نظر گرفتن وابستگی بین n -گرام‌ها بیز ساده را بهبود داده است. در دسته‌ای از کارها از مقادیر وابسته به مدل سازی زبانی به عنوان یک ویژگی استفاده شده است مانند کار پیلا در مقاله (Pillay & Solorio, 2010) که از یک مدل دو مرحله‌ای ترکیبی از دو روش یادگیری با ناظر و بدون ناظر برای تشخیص نویسنده روی پست‌های یک فروم وب استفاده کرده است. در این مدل مقدار پرپلکسیتی مدل زبانی در سطح سه کلمه به عنوان یک ویژگی در بردار ویژگی استفاده شده است.

از آنجا که استفاده از n -گرام کاراکتر و n -گرام کلمه همراه با مدل سازی زبانی راه حل به کار گرفته شده در این پژوهش است در جدول ۱-۲ تعدادی از کارهای پیشین که به نوعی از هر یک از آنها استفاده شده است لیست شده است.

۱-۲ کارهای انجام شده بر روی پایگاه‌داده‌های فارسی

در ادامه به عنوان آخرین بخش از این فصل مروری بر کارهای انجام شده بر روی پایگاه‌داده با زبان فارسی در زمینه مورد تحقیق در این پایان‌نامه یعنی شناسایی نویسنده می‌پردازیم.

کار (شاهمیری، 1385) در این کار پایگاه داده فارسی با چهار نویسنده کاندید تشکیل شده است. ۵۰ ویژگی گوناگون استخراج شده از اشعار ایرانی در سه دسته فیزیکی، مفهومی و آوایی دسته‌بندی شده است. برای تخصیص از روش‌های درخت تصمیم و شبکه عصبی مصنوعی استفاده شده است. دقت گزارش شده در درخت تصمیم تا ۹۴٪ و با شبکه عصبی تا ۹۵٫۹٪ درستی اعلام شده، در صورتی که دقت انسان تنها ۳۴٫۴٪ بوده است.

در کار (فرهمندپور، 1390) از ویژگی‌های واژگانی، نحوی، معنایی و وابسته به کاربرد استفاده شده است. الگوریتم رقابت کشورهای استعماری نیز جزء ایده‌ها و روش‌های جدید عنوان شده است که به علت سرعت همگرایی بالا، به آن پرداخته شده است. همچنین ایشان الگوریتم‌های یادگیری ماشین مانند دلتا، کا نزدیک‌ترین همسایه، درخت تصمیم، شبکه‌های عصبی، ترکیب الگوریتم ژنتیک و کا نزدیک‌ترین همسایه، ترکیب الگوریتم رقابت استعماری و کا نزدیک‌ترین همسایه و LDA را در دو پایگاه داده جمع‌آوری شده با هم مقایسه کرده است. روش دسته‌بندی LDA با بهترین دقت و ویژگی‌های نحوی، با بیشترین کارایی گزارش شده است. دو پایگاه داده ایجاد شده یکی در یک موضوع خاص از ۲۰ دانشجو با ۲۰۰۹ کلمه تهیه شده است و از دید نویسنده به علت کوتاه بودن و محاوره‌ای بودن متون باعث نادقیق شدن پارامترهای اندازه‌گیری شده است. در پایگاه دوم از مجموعه کتاب‌های هشت نویسنده معاصر استفاده شده است.

در کار (آذین، 1392) از پایگاه داده تشکیل شده از چهار شاعر شعر نو استفاده شده است. روش‌های تخصیص مورد استفاده کا نزدیک‌ترین همسایه، ماشین بردار پشتیبان و بیز ساده است. ویژگی‌های مورد استفاده در سه دسته ویژگی‌های واژگانی، حرفی و نحوی و ادغام این سه گروه ویژگی بوده است و بهترین گزارش در گروه ویژگی‌های نحوی بوده است.

و در نهایت کار رضانی و همکاران (Ramezani, et al., 2013) که به بررسی ویژگی‌های مختلف در پایگاه داده R40 (در بخش ۳-۲-۳ معرفی شده است) با استفاده بردار ماشین پرداخته است.

جدول ۱-۲: لیستی از کارهای پیشین. کارهایی که در آنها از n-گرام کاراکترها یا n-گرام کلمات و یا مدل سازی زبانی استفاده شده است

عنوان	ویژگی	تخصیص
۱	N-Gram Feature Selection for Authorship Identification ۲۰۰۶	استفاده از کاراکتر n-گرام ها با مقدار n ۳، ۴ و ۵ و با تعداد تکرار مشخص در یک بردار ویژگی
۲	Local histograms of character n-grams for authorship attribution	کاراکتر n-گرام
۳	Text Classification for Authorship Attribution Using Naive Bayes Classifier with Limited Training Data	۱- کاراکتر n-گرام ها با مقدار n ۱، ۲ و ۳ به صورت جداگانه استفاده کرده است. ۲- کلمه n-گرام های با مقدار n ۱، ۲ و ۳
۵	Searching with Style: Authorship Attribution in Classic Literature	۱- کلمات دستوری ۲- نقش های گرامری ۳- گرام نقش های گرامری ۴- ترکیب هر سه
۶	Using Function Words for Authorship Attribution: Bag of-Words vs. Sequential Rules	مقایسه استفاده از تکرار کلمات دستوری در مقابل ویژگی های استخراج شده با روش قوانین پی در پی
۷	Visualizing Document Authorship Using N-grams and Latent Semantic Indexing	تشکیل ماتریس لفظ-سند با استفاده از کاراکتر n-گرام ها و SVD
۷	Augmenting Naive Bayes Classifiers with Statistical Language Models	استفاده از n-گرام ها در سطح کاراکتر و کلمه
۹	Authorship Attribution of Short Messages Using Multimodal Features.	استفاده از n-گرام ها در سطح کاراکتر با مقدار n از ۲ تا ۶
۱۰	Time-aware authorship attribution for short text streams	کاراکتر n-گرام
۱۱	Combining naive Bayes and n-gram language models for text classification	کاراکتر n-گرام

بردار ماشین	<p>۱- سه بردار n-گرام کاراکترها + n-گرام کلمات + الگوی استخراج شده از توالی با شروع و پایان با کلمات پرتکرار و توالی از کلمات در بین آنها.</p> <p>۲- دو بردار n-گرام کاراکترها + الگوی استخراج شده</p> <p>۳- n-گرام کاراکترها + n-گرام کلمات</p> <p>۴- n-گرام کاراکترها</p>	Authorship Attribution of Micro-Messages	۱۲
<p>TIMBL (روشی بر اساس کا نزدیک‌ترین همسایه)</p>	<p>۱- n-گرام کلمات</p> <p>۲- n-گرام کلمات دستوری</p> <p>۳- n-گرام کاراکترها</p> <p>۴- n-گرام ریشه کلمات</p> <p>۵- n-گرام نقش دستوری کلمات</p>	The effect of author set size and data size in authorship attribution	۱۳
<p>متدی با ایده جداسازی کلاس ها با حاشیه بین صفحات کلاس</p>	تلفیق n -گرام کلمات با n -گرام کاراکترها	Authorship Attribution in Greek Tweets Using Author's Multilevel N-Gram Profiles	۱۴

۲-۲ جمع بندی

در این فصل به دسته بندی ویژگی‌ها در دسته‌های لغوی، کاراکتری، نحوی، معنایی و وابسته به کاربرد پرداخته شد و به چند نمونه از کارهای انجام شده اشاره شد. مزیت ویژگی‌ها در دسته لغوی سادگی و بی‌نیازی از ابزارهای پردازش زبان طبیعی عنوان شد و این که این مزیت آنها را تبدیل به ویژگی مستقل از زبان می‌کند همچنین وجود کارهای زیادی که با استفاده از کاراکتر n -گرام‌ها انجام شده است و نتایج خوبی گزارش شده نشان از مفید بودن آنها دارد. مزیت آنها در مستقل از زبان بودن، مشخص کردن اختلافات جزئی در سبک و تحمل پذیری خوب در برابر نویز شمرده شد. از طرفی مشکل تعیین مقدار n به صورت کارا در استفاده از آنها مطرح شد. ایده ویژگی‌های ساختاری و نحوی در استفاده ناخودآگاه نویسندگان از الگوهای نحوی بیان شد. همچنین گفته شد در ویژگی‌های معنایی توجه روی معنا و نقش معنایی کلمات است. و در نهایت از ویژگی‌های بسته به کاربرد به

عنوان استفاده از ویژگی‌های مختص حوزه‌های مختلف مانند سائز و رنگ فونت در فرم‌ها و یا پیام‌های الکترونیکی یاد شد.

در بخش ۲-۳ ابتدا به معرفی انواع روش‌های تشکیل داده آموزشی در دو دسته مبتنی بر نمونه و مبتنی بر پروفایل پرداخته شده است و سپس به دسته بندی متدهای تخصیص که در پژوهش‌های پیشین به کار گرفته شده، در دسته: متدها بر اساس فاصله، روش‌های فشرده سازی و روش‌های یادگیری ماشین مانند شبکه‌های عصبی، بردار ماشین و درخت تصمیم و همچنین روش‌های احتمالاتی مانند بیز و مدل سازی زبانی پرداخته شده است. و در انتها در بخش ۲-۴ مروری بر کارهای پیشین در پایگاه داده با زبان فارسی صورت گرفته است.

۳- فصل سوم: مباحث نظری و روش پیشنهادی

۱-۳ مقدمه

در این فصل به توصیف و تجزیه تحلیل روش‌های استفاده شده برای حل مسئله شناسایی نویسنده در این پایان‌نامه پرداخته شده است. ابتدا با توضیح در رابطه با پایگاه داده شروع شده است سپس به چگونگی حل مسئله تخصیص نویسنده پرداخته شده است و در ادامه ویژگی‌های مورد استفاده و روش تخصیص توضیح داده شده است. در انتها به معرفی خصوصیات از پایگاه داده که در حل مسئله تخصیص نویسنده تاثیر گزار است، پرداخته شده است.

۲-۳ پایگاه داده

از مشخصات یک پایگاه داده خوب آن چنان که در (Stamatatos, 2009) آمده است استفاده از داده‌هایی با نوع و موضوع یکسان است و همچنین استفاده از مجموعه نویسنده‌گانی با سطح تحصیلات، سن، ملیت و جنسیت یکسان است. با رعایت این موارد می‌توان اطمینان حاصل کرد که این عوامل تاثیری در طبقه‌بندی متن ندارند و طبقه‌بندی صرفاً بر اساس نویسنده‌گی صورت می‌گیرد. با این توضیحات سعی شده است تا در دو پایگاه داده WMPR-AA2016-A و WMPR-AA2016-B که توسط نگارنده جمع‌آوری شده از این اصول پیروی کنند. از طرف دیگر وجود پایگاه داده‌هایی با زبان مختلف و اندازه داده‌های آموزشی متفاوت و اندازه داده آزمایشی متفاوت شرایط را برای بررسی روش‌های بررسی شده در شرایط متفاوت فراهم کرده است. بحث بیشتر در این زمینه در فصل نتایج انجام خواهد شد.

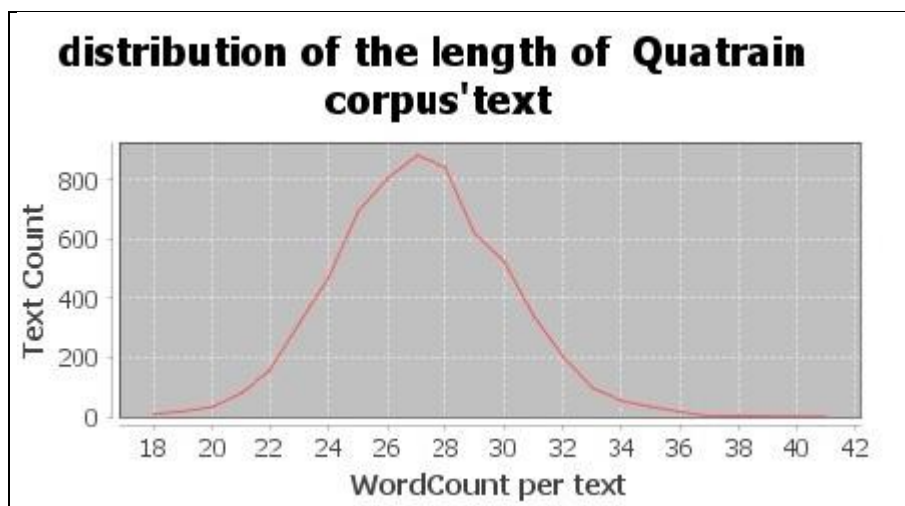
۱-۲-۳ پایگاه داده WMPR-AA2016-A

این پایگاه داده توسط نگارنده تهیه شده است و از مجموعه اشعار رباعی شش شاعر نامدار، ابوسعید ابوالخیر، عطار نیشابوری، انوری، خاقانی، مولوی، و ثنائی ایجاد شده است. در این پایگاه داده هر یک

رباعی به عنوان یک داده آزمایشی در نظر گرفته شده است از این رو با اندازه میانگین تعداد کلمه، ۲۷,۴ و انحراف از میانگین ۲,۶ این پایگاه داده در دسته پایگاه داده‌های متن کوتاه محسوب شده است. در شکل ۱-۳ توزیع طول متن در متن‌های آزمایشی در پایگاه داده WMPR-AA2016-A نشان داده شده است. نحوه چگونگی تخصیص داده آموزشی و داده آزمایشی در بخش ۱-۴ توضیح داده شده است. جزئیات مربوط به اطلاعات پایگاه داده از قبیل تعداد کلمه در داده‌های آموزشی و داده‌های آزمایشی و کل اندازه پایگاه داده در جدول ۱-۳ نمایش داده شده است.

جدول ۱-۳: جزئیات پایگاه داده فارسی WMPR-AA2016-A

پایگاه داده		نویسنده	(کلمات) کل	(کلمات) داده آموزشی	داده تست
WMPR-AA2016-A	A0	ابوسعید ابوالخیر	19352	15507	3854
	A1	انوری	12097	9684	2413
	A2	عطار نیشابوری	62248	49808	12440
	A3	خاقانی	9712	7797	1915
	A4	مولوی	53375	42700	10657
	A5	ثنائی	11711	9389	2322



شکل ۱-۳ توزیع طول متن در پایگاه داده WMPR-AA2016-A

۲-۲-۳ پایگاه داده WMPR-AA2016-B

این پایگاه داده توسط نگارنده تهیه شده است و از مجموعه غزلیات هفت شاعر بزرگ فارسی زبان، سعدی، ثنائی، عطار نیشابوری، اوحدی، ساوجی، مولوی و کرمانی تشکیل شده است. در این پایگاه

داده هر یک غزل به عنوان داده آزمایشی در نظر گرفته شده است. چگونگی تشکیل داده آموزشی و داده آزمایشی در بخش ۴-۲ توضیح داده شده است. جزئیات مربوط به اطلاعات پایگاه داده از قبیل تعداد کلمه در داده‌های آموزشی و داده‌های آزمایشی و کل اندازه پایگاه داده در جدول ۳-۲ نمایش داده شده است.

جدول ۳-۲: جزئیات پایگاه داده فارسی WMPR-AA2016-B

داده آموزشی (کلمات)	داده آزمایشی (کلمات)	کل (کلمات)	نویسنده		پایگاه داده WMPR- AA2016-B
70847	15826	88558	سعدی	B0	
42116	9797	52512	ثنائی	B1	
100264	23179	125265	عطار نیشابوری	B2	
92816	20563	115893	اوحدی	B3	
39555	8759	49417	ساوجی	B4	
98405	21841	122841	کرمانی	B5	
200240	50083	475558	مولوی	B6	

۳-۲-۳ پایگاه داده R40

این پایگاه داده توسط رضانی و همکارانشان تهیه شده است (Ramezani, et al., 2013) و از مجموعه نوشته‌های ۴۰ نویسنده نامی فارسی زبان مانند جمله جلال آل احمد، بزرگ علوی صادق چوبک و... در زمینه متن‌های ادبی و با سبک نگارشی نثر تهیه شده است. لذا این پایگاه داده در دسته پایگاه داده با تعداد نویسنده کاندید زیاد در نظر گرفته شده است. برای هر نویسنده پنج متن با اندازه بین ۸۰۰ تا ۱۰۰۰ کلمه در نظر گرفته شده است. جزئیات مربوط به اطلاعات این پایگاه داده در منبع عنوان شده، آورده نشده است لذا در این جا سعی شده تا اطلاعاتی از قبیل نام نویسنده‌گان و تعداد کلمات در

متن‌های هر نویسنده جمع آوری شود. اطلاعات جمع آوری شده در جدول ۳-۳ نمایش داده شده است.

جدول ۳-۳: جزئیات پایگاه داده فارسی R40

کل (کلمات)	نویسنده		کل (کلمات)	نویسنده	
۵۵۹۳	ناشناس	R20	۶۱۱۴	احمد میر علایی	R0
۴۹۳۹	غلام حسین ساعدی	R21	۵۲۱۰	محمد علی کاتوزیان	R1
۳۶۵۲	فروغ فرخزاد	R22	۳۶۰۵	امیر حسین چهل تن	R2
۶۱۴۱	مهرداد انتظاری	R23	۶۷۱۰	بزرگ علوی	R3
۸۲۵۳	گلی ترقی	R24	۱۳۶۰۵	سراسر حادثه	R4
۵۳۰۳	نغمه کیانی راد	R25	۱۹۶۶۴	بهنام دیانی	R5
۵۷۴۶	ناشناس	R26	۶۴۴۶	مژده کوهی	R6
۶۳۷۷	محمد علی	R27	۶۲۸۶	ناشناس	R7
۶۰۰۸	ناشناس	R28	۲۱۲۰۹	جلال ال احمد	R8
۴۱۴۳	جمال زاده	R29	۷۴۲۹	زهرا دلگرمی	R9
۵۰۶۶	ناشناس	R30	۵۲۲۵	رضا جولایی	R10
۵۳۱۶	مهرداد انتظاری	R31	۳۴۰۷	رضا قاسمی	R11
۷۶۹۶	میهن بهرامی	R32	۵۹۹۹	ناشناس	R12
۵۱۳۹	ناشناس	R33	۴۷۶۰	زویا پیرزاد	R13
۳۲۱۷	جعفر مدرس صادقی	R34	۹۱۶۱	شهریار مدنی پور	R14
۸۲۰۴	فریده شجاعی	R35	۸۱۲۵	صادق چوبک	R15
۷۲۰۴	الهه موذنی	R36	۳۱۷۲۰	صادق هدایت	R16
۹۲۰۳	نیلا	R37	۷۲۹۰	صمد بهرنگی	R17
۷۳۱۹	پروانه شیخو	R38	۳۸۹۹	علی خدایی	R18
۹۱۳۷	محنا عزیزی	R39	۱۰۳۷۷	ناشناس	R19

پایگاه داده
R40

۳-۲-۴ پایگاه داده RCV

پایگاه داده RCV (Lewis, et al., 2004) در زبان انگلیسی تهیه شده است. و با پنجاه نویسنده در دسته پایگاه داده با تعداد نویسنده زیاد در نظر گرفته شده است. متن‌ها در این پایگاه داده به صورت

نثر است و جنسیت نویسندگان به صورت مرد و زن می‌باشد. برای هر نویسنده ۵۰ متن آموزشی و ۵۰ متن آزمایشی وجود دارد. جزئیات مربوط به اطلاعاتی از قبیل تعداد کلمه در داده‌های آموزشی، داده‌های آزمایشی و کل اندازه پایگاه داده در جدول ۳-۴ نمایش داده شده است.

جدول ۳-۴: جزئیات پایگاه داده انگلیسی RCV

داده آموزشی (کلمات)	داده آزمایشی (کلمات)	کل (کلمات)	نویسنده	
24386	22499	46885	AaronPressman	C0
25884	22579	48463	AlanCrosby	C1
27875	29083	56958	AlexanderSmith	C2
26216	25517	51733	BenjaminKangLim	C3
25487	23336	48823	BernardHickey	C4
22815	24504	47319	BradDorfman	C5
29249	26210	55459	DarrenSchuettler	C6
26595	29024	55619	DavidLawder	C7
24506	21906	46412	EdnaFernandes	C8
26395	26572	52967	EricAuchard	C9
27483	27228	54711	FumikoFujisaki	C10
25493	25340	50833	GrahamEarnshaw	C11
25602	26961	52563	HeatherScofield	C12
28152	29620	57772	JaneMacartney	C13
21677	22485	44162	JanLopatka	C14
16358	17031	33389	JimGilchrist	C15
24014	23568	47582	JoeOrtiz	C16
24984	25294	50278	JohnMastrini	C17
25912	22347	48259	JonathanBirt	C18
25879	26416	52295	JoWinterbottom	C19
26248	21389	47637	KarlPenhaul	C20
22431	24715	47146	KeithWeir	C21
22872	23250	46122	KevinDrawbaugh	C22
26636	26612	53248	KevinMorrison	C23
31640	26219	57859	KirstinRidley	C24
24024	22307	46331	KourosKarimkhany	C25
16419	17121	33540	LydiaZajc	C26
29335	28920	58255	LynneO'Donnell	C27
30706	27521	58227	LynnleyBrowning	C28
26557	26034	52591	MarcelMichelson	C29
20714	22392	43106	MarkBendeich	C30
24285	25771	50056	MartinWolk	C31
26468	25439	51907	MatthewBunce	C32
22680	22071	44751	MichaelConnor	C33

پایگاه داده
RCV

25442	25556	50998	MureDickie	C34
25950	28587	54537	NickLouth	C35
23124	24580	47704	PatriciaCommins	C36
30723	29082	59805	PeterHumphrey	C37
27067	27199	54266	PierreTran	C38
24850	24026	48876	RobinSidel	C39
27252	27000	54252	RogerFillion	C40
27344	27728	55072	SamuelPerry	C41
28523	29007	57530	SarahDavison	C42
26455	24463	50918	ScottHillis	C43
24727	25228	49955	SimonCowell	C44
25500	22892	48392	TanEeLyn	C45
29171	28753	57924	TheresePoletti	C46
22391	24172	46563	TimFarrand	C47
25046	25812	50858	ToddNissen	C48
24421	23147	47568	WilliamKazer	C49

۳-۳ شناسایی نویسنده

همان‌طور که در فصل اول توضیح داده شد در فرایند شناسایی نویسنده هدف تخصیص یک نویسنده از میان مجموعه نویسنده‌گان کاندید (شناخته شده)، به متنی است که نویسنده آن مشخص نیست. همچنین نیازمندی‌های این فرایند انتخاب ویژگی‌های سبکی مناسب و انتخاب روش تخصیص کارآمد به منظور انتساب نویسنده درست به متن دیده نشده معرفی شد. از این جهت در ادامه پس از توضیح مرحله پیش پردازش، ویژگی‌ها و روش تخصیص استفاده شده در این پایان‌نامه شرح داده خواهد شد.

۱-۳-۳ پیش‌پردازش

در مرحله پیش‌پردازش تمامی علامت‌های خاص، اعداد و حروف الفبای خارجی از تمامی بانک‌های اطلاعاتی حذف شده و با فضای خالی جایگزین شده‌اند. برای جداسازی جملات در پایگاه داده به ترتیب در ابتدا و انتهای هر جمله از تگ <s> و </s> استفاده شده است. تا در مرحله استخراج ویژگی‌ها ابتدا و انتهای هر جمله مشخص باشد. در پایگاه داده با سبک نوشتاری نظم هر مصرع به

عنوان یک جمله در نظر گرفته شده است و در نثر علامت نقطه و علامت سوال به عنوان مشخص کننده جمله در نظر گرفته شده است.

۳-۳-۲ انتخاب ویژگی

در این پایان نامه n -گرام کاراکتر و n -گرام کلمه به عنوان ویژگی انتخاب شده است. علت این انتخاب همان طور که در فصل قبل توضیح داده شد این است که n -گرام کارکترها و کلمات در دسته‌ای از ویژگی‌ها قرار دارند که برای استخراج احتیاج به ابزارهای پردازش زبان ندارند. از این رو استخراج این ویژگی‌ها هزینه زمان اجرا و خطای ناشی از خطای ابزار پردازش زبان را ندارند. مورد پر اهمیت‌تر استقلال این ویژگی‌ها از زبان است. از این جهت که استفاده از ابزارهای پردازش زبان وابسته به زبان است و در هر زبانی ابزارهای پردازش زبان خود را نیاز دارد. به عنوان نمونه استخراج قوانین نحوی در زبان فارسی با زبان انگلیسی متفاوت است. در ادامه روش استخراج n -گرام‌ها شرح داده شده است.

۳-۳-۲-۱-۲-۳-۳-۳ گرام‌ها

در یک تعریف کلی می‌توان گفت، n -گرام دنباله‌ای از n توکن پشت سر هم در یک رشته از توکن‌ها می‌باشد. توکن‌ها بنا بر استفاده از هر جنسی اعم از کاراکتر، کلمه، ریشه کلمات، نقش دستوری کلمات و ... می‌توانند باشند. استخراج n -گرام‌ها را به صورت مراحل زیر می‌توان توضیح داد:

- برای استخراج n -گرام‌ها متن به صورت توالی از توکن‌ها یا به عبارت دیگر توالی از کاراکترها یا کلمات یا ریشه کلمات یا نقش دستوری کلمات و... در می‌آید.
- مرحله بعد انتخاب طول n -گرام یا همان مقدار n است.
- و بعد از آن یک پنجره با طول n در نظر گرفته می‌شود.
- به اندازه $n-1$ کاراکتر فضای خالی برای پوشش تمام توکن‌های متن به انتهای متن اضافه می‌گردد.

▪ حرکت پنجره از اولین توکن متن شروع شده و جابجایی آن توکن به توکن تا آخرین توکن

متن ادامه می‌یابد. در هر بار حرکت پنجره آنچه در داخل پنجره قرار گرفته است، n -گرام با

اندازه n می‌باشد. (شکل 3-2)

تعداد n -گرام‌های تولید شده مستقل از مقدار n برابر با تعداد توکن‌ها در متن بدون در نظر گرفتن $n-1$

توکن فاصله اضافه شده، است.

متن : یک صبح روز یکشنبه ماه تیر هوای شهر برلین تیره و خفه کننده بود.

۱- تولید توالی از توکن‌ها در شکل کاراکترها
 ی،ک،-،ص،ب،ح،-،ر،و،ز،-،ی،ک،ش،ن،ب،ه،-،م،ا،ه،-،ت،ی،ر،-،ه،و،ا،ی،-،ش،ه،ر،-،ب،ر،ل،ی،ن،-،
 ت،ی،ر،ه،و،خ،ف،ه،ک،ن،د،ه،-،ب،و،د،

۲- تولید انگرام با اندازه مشخص
 حرکت پنجره با اندازه $n=2$

ی،ک،-،ص،ب،ح،-،ر،و،ز،-،ی،ک،ش،ن،ب،ه،-،م،ا،ه،-،ت،ی،ر،-،ه،و،ا،ی،-،ش،ه،ر،-،
 ب،ر،ل،ی،ن،-،ت،ی،ر،ه،و،خ،ف،ه،ک،ن،د،ه،-،ب،و،د،

اندازه	انگرام‌های تولید شده
۲-گرام	یک،ک-؛ ص-؛ صب؛ بح؛ ح-؛ ر-؛ رو؛ وز؛ ز-؛ ی-؛ یک؛ کش؛ شن؛ نب؛ به؛ ه-؛ م-؛ ما؛ اه؛ ه-؛ - ت؛ تی؛ یر؛ ر-؛ ه-؛ هو؛ و....
۳-گرام	یک-؛ ک-ص-؛ صب-؛ صبح؛ بح-؛ ح-ر-؛ رو-؛ روز؛ وز-؛ ز-ی-؛ یک-؛ یکش؛ کشن؛ شنب؛ نبه؛ به-؛ -ما؛ ماه؛ اه-؛ ه-ت-؛ تی؛ تیر؛ و....

شکل ۳-۲: تولید n -گرام‌ها در یک متن. در این مثال فاصله به عنوان یک توکن در نظر گرفته شده است و با علامت "-" مشخص شده است. علامت کاما جدا کننده n -گرام‌ها و کاراکترها است.

۳-۳-۳ روش تخصیص

متد تخصیص مورد استفاده در این پایان‌نامه مدل‌سازی زبانی است. علت انتخاب این روش قرار گرفتن

آن در دسته روش‌های آماری احتمالاتی است. همچنین به جهت قرار گرفتن آن در دسته روش-

های پویا و استفاده از احتمال پیشین باعث باز شدن راهی برای بسط احتمال پیشین و بهبود نتایج در

مدل به دست آمده در کارهای آینده می‌شود. علاوه بر آن با بررسی‌های انجام شده، این روش تاکنون روی پایگاه داده با زبان فارسی برای حل مسئله تشخیص نویسنده استفاده نشده است. در ادامه به توضیح این روش پرداخته شده است.

۳-۳-۴ مدل سازی زبانی در شناسایی نویسنده

در مدل سازی زبانی هدف پیدا کردن احتمال رخداد توالی از کلمات در یک مجموعه از جملات زبان است. به عبارت دیگر سنجیدن اینکه با توجه به داده‌های آموزشی که از یک زبان در دست است با چه فراوانی، توالی از کلمات می‌تواند توسط مدل زبانی داده‌های آموزشی تولید شود.

استفاده اولیه از مدل سازی زبانی در پردازش سیگنال بوده است و دامنه استفاده از آن تا حوزه‌های بایوانفورماتیک، یادگیری ماشین و بازیابی اطلاعات نیز گسترش پیدا کرده است. در ترجمه ماشینی آماری، یکی از اجزا اصلی سیستم ترجمه، استفاده از مدل سازی زبانی به این منظور است که ترجمه یک عبارت چقدر در زبان مورد هدف، عبارتی کاربردی و متعارفی است (Koehn, 2010). در بازیابی اطلاعات ایده اصلی در استفاده از مدل سازی زبانی این است که یک سند در صورتی کاندید خوبی برای تطبیق پیدا کردن با کوئری مورد جستجو است که سند با احتمال بالایی بتواند کوئری را تولید کند. (D.Manning, et al., 2009)

پایه بیشتر تحقیقات و راه کارهای ارائه شده برای حل مسئله شناسایی نویسنده بر این اساس است که هر نویسنده از الگوهای نوشتاری مخصوص به خود در نوشتار استفاده می‌نماید که با عنوان "اثر انگشت مولف" از آن یاد می‌شود و توضیح آن هم می‌تواند این باشد که هر شخصی زبان را خودش و با تجربیات و برداشت‌های شخصی خودش فرا می‌گیرد (Joula, 2008). از این مسئله می‌توان استفاده کرد و استفاده از مدل سازی زبانی در حل مسئله شناسایی نویسنده را این گونه توضیح داد که با جمع‌آوری مجموعه نوشته‌های هر نویسنده در مجموعه نویسنده‌گان کاندید و الحاق تمام متن‌های جمع‌آوری شده برای هر نویسنده در یک متن بزرگتر، برای هر نویسنده کاندید الگویی از زبان آن

نویسنده ایجاد شود. سپس با به کارگیری مدل سازی زبانی به دسته بندی متن دیده نشده با نویسنده نامشخص پرداخته شود.

▪ محاسبه مدل سازی زبانی

همان طور که گفته شد در مدل سازی زبانی هدف پیدا کردن احتمال یک کلمه بر اساس $n-1$ کلمه قبل از آن و یا احتمال پیدا کردن یک جمله است. ارتباط این دو احتمال به این صورت است که اگر فرض شود جمله U از کلمات $u_1 \dots u_n$ تشکیل شده است آنگاه طبق قانون زنجیره ای داریم :

$$P(U = u_1, \dots, u_n) = P(u_1) P(u_2 | u_1) P(u_3 | u_1, u_2) \dots P(u_n | u_1, \dots, u_{n-1}) \quad 1-3$$

رابطه ۱-۳ ارتباط بین احتمال جملات با احتمال کلمات را نشان می دهد. از محاسبه احتمال کلمه به شرط وقوع کلمات قبلی به احتمال جمله رسیده است. برای محاسبه هر یک از احتمالات شرطی در رابطه ۱-۳ داریم :

$$P(u_n | u_{n-1}) = \frac{\text{Count}(u_1, \dots, u_{n-1} u_n)}{\text{Count}(u_1, \dots, u_{n-1})} \quad 2-3$$

در رابطه ۲-۳ صورت از شمارش تعداد عبارات حاصل از کنار هم قرار دادن کلمات u_1 تا u_n و مخرج از شمارش تعداد عبارات حاصل از کنار هم قرار دادن کلمات u_1 تا u_{n-1} در متن مدل، محاسبه می گردد. مشکلی که در استفاده از رابطه ۱-۳ و قانون زنجیره ای وجود دارد این است که اگر تمامی ارتباط بین کلمات در نظر گرفته شود در محاسبه هر یک از احتمالات شرطی در رابطه ۲-۳ احتیاج به یک مجموعه متن بسیار بزرگ است در غیر این صورت بیشتر احتمالات صفر خواهد شد و در نتیجه آن احتمال جمله U نیز صفر خواهد شد. از این رو در مدل سازی زبانی با استفاده از قانون مارکوف ارتباط بین کلمات محدود می شود و این محدودیت در مدل سازی زبانی ساده با در نظر نگرفتن ارتباط بین کلمات یا مستقل فرض کردن کلمات از یکدیگر منتهی می شود. در نتیجه آن رابطه ۱-۳ به صورت زیر در می آید (Manning & Schitze, 2000)(Jurafsky & Martin, 2006).

$$P(U = u_1, \dots, u_n) = P(u_1) P(u_2) \dots P(u_n) \quad 3-3$$

▪ شناسایی نویسنده با مدل سازی زبانی ساده

اگر فرض شود مجموعه $A = \{a_1, a_2, \dots, a_m\}$ ، مجموعه نویسندگان کاندید هستند هدف پیدا کردن مناسبترین نویسنده برای متن دیده نشده U (متن با نویسنده نامشخص) است. با توضیحاتی داده شده به راحتی می توان متوجه شد که مدل سازی زبانی در گروه روش های بر مبنای پروفایل قرار دارد (تعریف روش های بر مبنای پروفایل در بخش ۳-۲-۱ آورده شده است). بنابراین به منظور تشکیل پروفایل هر نویسنده a_i ، در مرحله اول تمامی متن های آموزشی مربوط به نویسنده a_i در مجموعه A ، در یک متن الحاق شده و متن حاصل را d_i می نامیم. حاصل مرحله اول تشکیل مجموعه آموزشی $D = \{d_1, d_2, \dots, d_m\}$ خواهد بود. که در آن هر d_i پروفایل یا داده آموزشی یا داده مدل نویسنده a_i را تشکیل می دهد. هدف پیدا کردن اندیس نویسنده i است به طوری که:

$$d_i = \operatorname{argmax}\{P(d_i | U) \mid d_i \in D\} \quad 4-3$$

در آن d_i متن آموزشی مربوط به نویسنده a_i و U متن مورد سوال (متن با نویسنده نامشخص) است. طبق قانون بیز رابطه ۳-۴ می تواند به صورت زیر نوشته شود:

$$d_i = \operatorname{argmax}\{P(d_i)P(U | d_i)/P(U)\} \mid d_i \in D \quad 5-3$$

در رابطه ۳-۵ از آنجا که $P(U)$ یا احتمال evidence مستقل از اندیس i است و برای تمامی داده های آموزشی مقدار یکسانی دارد می تواند از رابطه حذف شود و $P(d_i)$ یا احتمال پیشین برای تمامی داده های آموزشی یکسان در نظر گرفته شده در نتیجه رابطه ۳-۵ به صورت رابطه زیر در می آید:

$$d_i = \operatorname{argmax}\{P(U | d_i)\} \mid d_i \in D \quad 6-3$$

در رابطه ۳-۶، $P(U | d_i)$ یا احتمال بیشترین شباهت، احتمال تولید متن مورد سوال U است که توسط مدل زبانی تشکیل شده روی داده آموزشی d_i تولید شده است. برای محاسبه احتمال $P(U | d_i)$ طبق

قانون زنجیره ای در رابطه ۱-۳ و قانون مارکوف در رابطه ۲-۳، احتمال هر جمله در متن U از حاصل- ضرب احتمال هر یک از کلمات u_1 تا u_n در مدل زبانی داده آموزشی d_i به دست می‌آید. در متن‌های بزرگ ضرب احتمال کلمات در یکدیگر ممکن است خطای میز شناور^۱ را ایجاد نماید. برای جلوگیری از این خطا از لگاریتم استفاده می‌شود و در نهایت خواهیم داشت:

$$i = \operatorname{argmax}_{d_i \in D} \left\{ \sum_{u_j \in U} \log P(u_j | LM_{d_i}) = \sum_{u_j \in U} \log \frac{t_{u_j, d_i}}{L_{d_i}} \right\} \quad ۷-۳$$

که در این رابطه u_i کلمات در متن مورد سوال U است. LM_{d_i} مدل زبانی ساخته شده روی داده آموزشی d_i است. و $\frac{t_{u_j, d_i}}{L_{d_i}}$ احتمال هر کلمه u_j روی مدل زبانی ساخته شده در داده آموزشی d_i است که از تقسیم تعداد کلمات u_j در متن آموزشی d_i تقسیم بر کل کلمات متن آموزشی d_i به دست می‌آید.

۴-۳ روش پیشنهادی

روش پیشنهادی با ایجاد دو تغییر در مدل سازی زبانی ساده ایجاد شده است و از این پس از آن با نام مدل سازی زبانی تغییر یافته یاد شده است. در ادامه به توضیح تغییرات داده شده می‌پردازیم.

۱-۴-۳ استفاده از ضریب IDF به جای حذف کلمات پرتکرار

در استفاده از کلمات به عنوان ویژگی یکی از مشکلات، رخداد کلماتی است که به طور عام و مشترک استفاده می‌شوند و با نام کلمات پرتکرار^۲ شناخته شده‌اند. این کلمات به خاطر مشترک بودن و با تکرار بالا در بین بیشتر نوشته‌ها ارزش جداکنندگی زیادی به عنوان یک ویژگی ندارند و در اکثر روش‌ها به علت نداشتن اثر منفی در روش تخصیص، حذف می‌شوند. روش حذف آنها استفاده از

^۱ Floating Point Under Low

^۲ Stop Words

لیستی از پرتکرارترین کلمات در میان متون پیکره و انتخاب و حذف این کلمات به عنوان کلمات پرتکرار است. مشکلاتی که در استفاده از لیست برای حذف کلمات پرتکرار وجود دارد این است که:

▪ تضمینی بر اینکه آیا هیچ یک از کلمات قرار گرفته در لیست در دسته کلمات با قدرت جداکنندگی نیستند، وجود ندارد.

▪ آیا تمام کلمات پرتکرار در یک لیست ایستا وجود دارند؟

▪ استفاده از لیست ایستا روش را به زبان وابسته می‌کند.

یکی از راهکارها، که به نظر می‌رسد بتواند تا حدودی مشکلات مطرح شده را حل نماید استفاده از یک روش وزن‌دهی به کلمات است. تا با استفاده از آن، کلمات با قدرت جداکنندگی بیشتر وزن بیشتری نسبت به کلمات پرتکرار به دست آورند و اثر بیشتری در تعیین نویسنده داشته باشند. یکی از این روش‌های وزن‌دهی استفاده از ضریب IDF^1 است.

$$IDF = \frac{N}{DF_u} \quad ۸-۳$$

در رابطه ۸-۳، N تعداد متن‌های آموزشی و DF_u تعداد متن‌های آموزشی که کلمه u در آنها تکرار شده است. در مدل‌سازی زبانی تغییر یافته به جای استفاده از لیست کلمات پرتکرار، و حذف کلمات لیست از داده‌های آموزشی، هر یک از کلماتی که به عنوان یک ویژگی انتخاب می‌شوند در مقدار ضریب IDF ضرب می‌گردد. و به این ترتیب به جای حذف کلماتی که ممکن است دربرگیرنده خاصیت سبکی نویسنده باشند وزن آنها تغییر داده می‌شود. به عنوان نمونه اگر احتمال رخداد دو کلمه متفاوت u_1 و u_2 در مدل زبانی نویسنده a_1 با یکدیگر برابر باشند و فرض شود u_1 در دسته کلمات پرتکرار قرار دارد که توسط اکثر نویسندگان کاندید استفاده شده است و u_2 در دسته کلماتی قرار داشته باشد که در متن‌های نویسندگان کمتری شرکت دارد، در این صورت مقدار ضریب IDF برای u_2 مقدار عددی بزرگتری نسبت به ضریب IDF برای u_1 خواهد بود و نتیجه ضرب مقادیر IDF به

¹ Inversing Document Frequency

دست آمده در احتمالات برابر u_1 و u_2 مقدار عددی بزرگتری برای u_2 است. به این ترتیب به جای حذف کلمه u_1 که ممکن است دربرگیرنده خاصیت سبکی نویسنده باشد وزن آن تغییر داده می‌شود.

۳-۴-۲ ترکیب ۱-گرام و ۲-گرام کلمات

همانطور که در قسمت قبل شرح داده شد مدل سازی در سطح دو کلمه به دنبال ارتباطات دوتایی کلمات است و مدل سازی در سطح سه کلمه به دنبال ارتباطات سه تایی بین کلمات است. از آنجا که با بالا رفتن مقدار n در n -گرام کلمات احتمال وقوع آن در یک متن کمتر می‌شود در مدل سازی در سطح سه کلمه و بالاتر نیاز به داده آموزشی با تعداد کلمه زیاد است تا بتوان برای متن دیده نشده به مقدار احتمال معنا داری دست یافت در غیر این صورت برای اکثر n تایی‌ها، در متن دیده نشده مقدار احتمال صفر به دست می‌آید. در مدل سازی در سطح تک کلمه، کلمات به تنهایی و بدون در نظر گرفتن ارتباط بین کلمات و در اصطلاح به صورت کیسه‌ای از کلمات^۱، به عنوان ویژگی استفاده شده است. به نظر می‌رسد استفاده از این ویژگی همیشه نمی‌تواند کافی باشد. به عنوان مثال فرض کنیم دو نویسنده کاندید با دو مجموعه آموزشی داریم و مدل سازی زبانی روی هر یک از این داده‌های آموزشی با نام‌های LM_1 و LM_2 پیاده‌سازی شده است. هدف تخصیص نویسنده درست به متن دیده نشده U با محتوای "خانه دوست کجاست؟ در فلق بود که پرسید سوار"، است. اگر تعداد تکرار کلمات در متن U در دو LM_1 و LM_2 یکسان فرض شود در این صورت در تخصیص متن U با دو مدل زبانی LM_1 و LM_2 ، هر دو امتیاز یکسانی را دریافت خواهند کرد. اما با فرض اینکه که در LM_1 علاوه بر تک کلمات موجود در متن U ، ترکیبات دوتایی مانند "خانه دوست" نیز وجود دارد و این ترکیبات در LM_2 وجود ندارد و با گسترش نگاه از سطح تک کلمه به دو کلمه در این صورت در رقابت بین دو مدل بانی LM_1 مدل مناسب‌تری به نظر خواهد رسید. زیرا با وجود کلمات دوتایی مشترک بین متن U و LM_1 ، U به مدل زبانی LM_1 نزدیک‌تر است.

¹ Bag Of Words

با توجه به توضیحات داده شده به نظر می‌رسد که تلفیق دو ویژگی ۱-گرام کلمه و ۲-گرام کلمه، و استفاده آنها به عنوان ویژگی در مدل سازی زبانی می‌تواند منجر به آشکار شدن ویژگی‌های سبکی بیشتری در مسئله تشخیص نویسنده گردد. در مدل سازی زبانی تغییر یافته احتمال ۲-گرام کلمات با احتمال ۲-گرام کلمات تلفیق شده است. با این کار به مجموعه داده‌های آموزشی که علاوه بر کلمات متن دیده نشده دارای کلمات در سطح ۲-گرام از متن دیده نشده هستند امتیاز بیشتری داده شده است. علاوه بر این آن‌چنان‌که در قسمت ۳-۴-۱ در استفاده از IDF به جای حذف کلمات پرتکرار توضیح داده شد برای حذف دوتایی‌های پر تکرار مانند دو تایی "من و" که در اکثر متن‌ها دیده می‌شود از لگاریتم ضرب IDF در سطح دو گرام کلمه استفاده شده است. در نهایت فرمول مدل سازی زبانی تغییر یافته از فرمول مدل سازی زبانی ساده در سطح کلمه در رابطه ۳-۷ به شکل زیر تغییر پیدا خواهد کرد:

$$d_i = \operatorname{argmax}_{d_i \in D} \left\{ \prod_{u_j \in U} \operatorname{IDF} P(u_j | LM_{d_i}) \operatorname{IDF} P(u_j | u_{j-1}, LM_{d_i}) \right\} \quad 9-3$$

که در این رابطه $P(u_j | LM_{d_i})$ احتمال تک کلمات در مدل زبانی ساخته شده روی داده آموزشی d_i مربوط به نویسنده a_i کاندید است و برای محاسبه مانند رابطه ۳-۷ از $\frac{t_{u_j, d_i}}{L_{d_i}}$ که احتمال هر کلمه u_j روی مدل زبانی ساخته شده در داده آموزشی d_i است و از تقسیم تعداد کلمات u_j در متن آموزشی d_i تقسیم بر کل کلمات متن آموزشی d_i به دست می‌آید، استفاده شده است. همچنین $P(u_j | u_{j-1}, LM_{d_i})$ احتمال ۲-گرام کلمات در مدل زبانی ساخته شده روی داده آموزشی d_i مربوط به نویسنده a_i است و برای محاسبه آن مطابق با رابطه ۳-۲ از تقسیم تعداد دوتایی‌های u_j u_{j-1} در متن آموزشی d_i تقسیم بر تعداد کلمات u_j در متن آموزشی d_i به دست می‌آید. در نهایت برای تبدیل ضرب به جمع از لگاریتم استفاده شده است.

۳-۵ خصوصیات پایگاه داده

در (Stamatatos, 2009) خصوصیات از پایگاه داده را که در میزان کارایی روش‌های حل مسئله نویسنده موثر هستند به صورت زیر معرفی کرده است:

- مقدار داده آموزشی

مقدار داده آموزشی به این موضوع که چه مقدار داده آموزشی برای حل مسئله نویسنده در نظر گرفته شده است اشاره دارد و اینکه کارایی به دست آمده وابسته به چه مقدار داده آموزشی است و تغییر مقدار داده آموزشی و کاهش آن چه تاثیری در نتایج دارد. این مسئله از آن جا اهمیت بیشتری پیدا می‌کند که در بعضی از مسائل وابسته به تشخیص نویسنده مانند تشخیص نویسنده در متن‌های مجرمانه معمولاً داده آموزشی زیادی در دست نیست. البته مقدار داده آموزشی در نظر گرفته شده می‌تواند به طول داده آزمایشی نیز بستگی داشته باشد؛ زیرا زمانی که داده آزمایشی دارای متن کوتاه است برای استخراج ویژگی‌های جداکننده احتیاج به داده آموزشی بیشتری خواهد بود (Luyckx & Daelemans, 2011).

- تعداد نویسنده در مجموعه نویسنده‌گان کاندید

یکی از چالش‌ها در حل مسئله تخصیص نویسنده تعداد نویسنده‌گان در مجموعه نویسنده‌گان کاندید است. هر چه میزان تغییرات دقت در مدل با افزایش تعداد نویسنده‌گان کمتر باشد یا به عبارت دیگر روش استفاده شده، در تعداد نویسنده زیاد کارایی خود را از دست ندهد، نشان‌دهنده پایداری بیشتر روش مورد استفاده است (Luyckx & Daelemans, 2011).

از جمله کارهایی که در این زمینه بررسی انجام داده‌اند می‌توان به کار (Argamon, et al., 2003) اشاره کرد که در نتایج آن گزارش شده است با افزایش تعداد نویسنده‌گان از ۲ به ۲۰ نویسنده میزان دقت تا ۴۰٪ کاهش پیدا کرده است. در کار (Grieve, 2007) با اعمال ویژگی‌های متفاوت و گزارش نتایج برای تعداد نویسنده کاندید در بازه ۲ تا ۴۰ نویسنده کمترین میزان تغییرات در تقریباً بازه ۹۶ تا ۶۵ درصد و بیشترین میزان تغییرات با به کارگیری ویژگی‌هایی مانند میانگین اندازه جملات در بازه

۷۰ تا ۶ درصد قرار دارد. در کار (Abbasi & Chen, 2008) با استفاده از تکنیک Writeprints که بر مبنای تبدیل karhunen-koeve می‌باشد دقت ۹۲ تا ۸۳ و ۹۶ تا ۹۱ را برای تعداد نویسنده ۲۰ تا ۱۰۰ نمایش داده است.

- متعادل بودن داده آموزشی

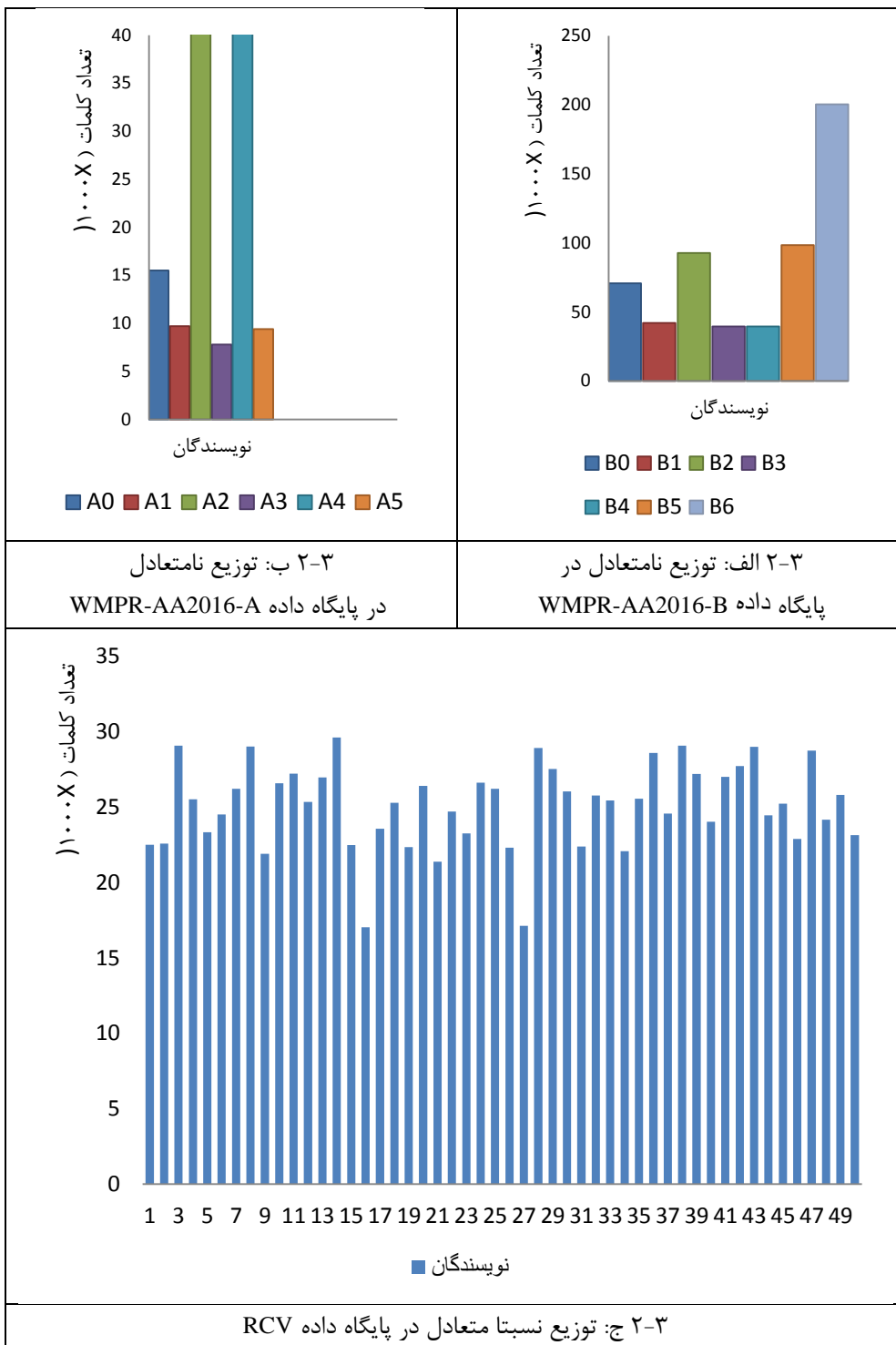
توزیع داده‌های آموزشی به صورت یکسان یا متفاوت در بین نویسنده‌گان کاندید با نام داده آموزشی متعادل و داده آموزشی نامتعادل نامیده شده است.

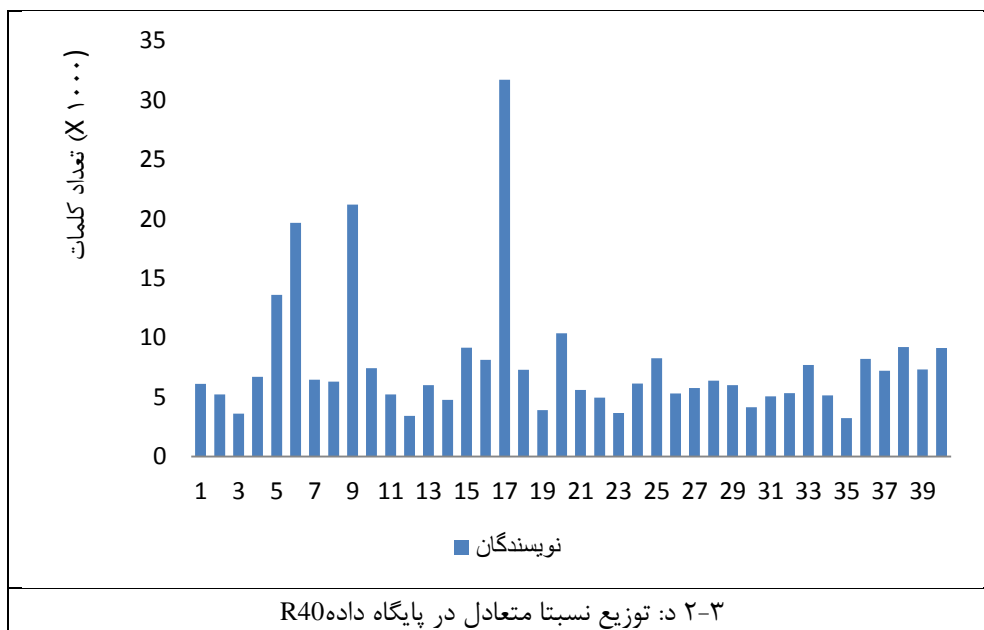
در بیشتر روش‌ها احتمال پیشین در بین تمام نویسنده‌گان کاندید یکسان فرض شده است اما در حالتی که توزیع داده‌ها در بین داده‌های آموزشی نویسنده‌گان کاندید یکسان نباشد، باعث از بین رفتن این برابری احتمال پیشین می‌شود. (Alexander Yun-chung Liu, 2004). راه حل ارائه شده برای حل این مسئله استفاده از اضافه کردن آموزشی^۱ و یا کم کردن داده آموزشی^۲ است و راه حل‌های متفاوتی برای این منظور ارائه شده است. در ساده‌ترین حالت برای اضافه کردن داده آموزشی، داده‌های آموزشی موجود تکرار می‌شوند و برای کم کردن داده آموزشی، از داده‌های آموزشی به صورت تصادفی کم می‌شود (Stamatatos, 2008) (Alexander Yun-chung Liu, 2004).

در ادامه توزیع داده‌های آموزشی در بین نویسندگان کاندید چهار پایگاه داده استفاده شده در این پایان‌نامه نمایش داده شده است. در بین این پایگاه داده‌ها، پایگاه داده WMPR-AA2016-A و WMPR-AA2016-B دو پایگاه داده نامتعادل و پایگاه داده RCV نسبتاً متعادل و پایگاه داده R40 توسط رضانی و همکارانشان متعادل معرفی شده است.

¹ Up Sampling

² Down Sampling





شکل ۳-۳: توزیع داده‌های آموزشی در بین نویسندگان کاندید. در چهار پایگاه داده WMPR-AA2016-A، WMPR-AA2016-B، RCV و R40

۳-۶ جمع بندی

در این فصل به معرفی چهار پایگاه داده استفاده شده در این پایان نامه پرداخته شد. سه پایگاه داده در زبان فارسی و پایگاه داده چهارم در زبان انگلیسی است. پایگاه داده‌ها در سبک نگارش و اندازه داده داده آموزشی، آزمایشی و تعداد نویسنده کاندید متفاوت است. پس از آن در بخش ۳-۳ مراحل مورد نیاز برای حل مسئله شناسایی نویسنده اعم از پیش پردازش و انتخاب ویژگی و روش تخصص عنوان شده است و n-گرام کاراکترها و کلمات و مدل سازی زبانی به عنوان ویژگی و روش تخصیص انتخابی شرح داده شده است. در بخش ۳-۴ یک روش پیشنهادی با ترکیب ۱-گرام کلمات با ۲-گرام کلمات و استفاده از مقدار IDF به عنوان وزن دهی احتمال n-گرام‌ها بر اساس میزان تکرار در داده‌های آموزشی هر نویسنده، معرفی شده است.

و در نهایت در بخش ۳-۵ به معرفی خصوصیتی از پایگاه داده که در نتایج ارزیابی در حل مسئله تخصیص نویسنده موثر هستند پرداخته شده است. و از آنها با عنوان مقدار داده آموزشی، تعداد نویسندگان کاندید و توزیع داده‌های آموزشی بین نویسندگان کاندید نام برده شده است.

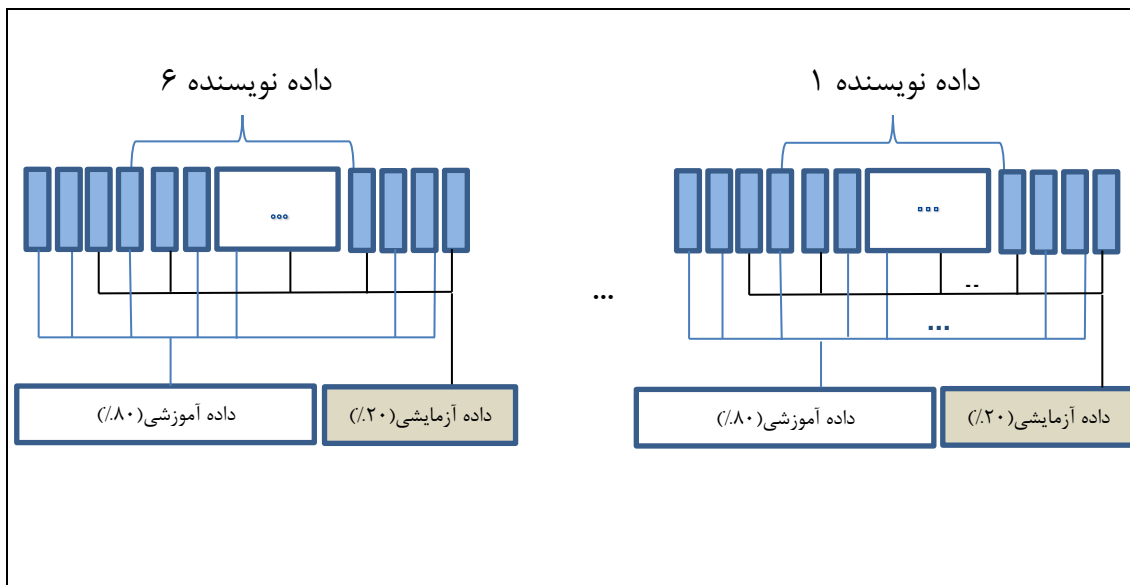
۴- فصل چهارم : نتایج

۱-۴ مقدمه

در فصل قبل به معرفی روش مدل سازی زبانی با استفاده از n-گرام کلمات و n-گرام کاراکترها پرداخته شد. همچنین ترکیب 1-گرام تک کلمات با 2-گرام کلمات و استفاده از مقدار IDF به عنوان وزن دهی احتمال n-گرامها (بر اساس میزان تکرار در داده‌های آموزشی هر نویسنده) در مدل سازی زبانی با نام مدل سازی تغییر یافته، پیشنهاد شد. در این فصل به ارزیابی حل مسئله نویسنده با روش مدل سازی زبانی ساده با استفاده از n-گرام کلمات با مقدار $n=2, 3$ و n-گرام کاراکترها با تغییر مقدار n از $n=2$ تا $n=6$ همچنین مدل سازی زبانی تغییر یافته در چهار پایگاه داده معرفی شده در بخش ۳-۲ پرداخته شده است. سپس به بررسی اثر خصوصیات پایگاه داده که در بخش ۳-۵ توضیح داده شد، پرداخته خواهد شد.

۱-۴ آزمایشات در پایگاه داده WMPR-AA2016-A

در آزمایشات صورت گرفته در این پایگاه داده شش نویسنده در رقابت تعیین نویسنده ناشناس برای داده‌های آزمایشی شرکت دارند. برای تشکیل داده‌های آزمایشی و داده‌های تست از مجموعه کل متون موجود برای هر نویسنده ۲۰٪ داده‌ها برای داده آزمایشی و ۸۰٪ درصد داده‌ها برای داده آموزشی انتخاب شده است. انتخاب داده‌ها برای قرارگیری در دو مجموعه آزمایشی و آموزشی به صورت تصادفی انجام شده است. شکل ۱-۴ نمایی کلی از این تقسیم بندی را نمایش می‌دهد.



شکل ۴-۱: تشکیل داده آموزشی و مجموعه داده آزمایشی در پایگاه داده WMPR-AA2016-A

۴-۱-۱ آزمایش با مدل سازی زبانی و n-گرام کلمات

نتایج حاصل از اجرای مدل سازی زبانی ساده با n-گرام کلمات با مقدار ۳ و ۲، ۱ و n در جدول ۴-۱ نمایش داده شده است. برای روشن تر شدن ارزیابی، علاوه بر میانگین دقت، نتایج ارزیابی به تفکیک هر نویسنده نیز محاسبه شده است.

جدول ۴-۱: نتایج ارزیابی روش مدل سازی زبانی با n-گرام کلمات در پایگاه داده WMPR-AA2016-A

<i>F-measure</i>	<i>Recall</i>	<i>Precision</i>	نویسنده	
41.19	38.73	44.0	A0	مدل سازی زبانی با ۱-گرام کلمات
43.26	51.72	37.19	A1	
77.15	79.78	74.69	A2	
39.26	47.05	33.68	A3	
62.26	53.26	74.91	A4	
34.65	43.20	28.92	A5	
60.27				میانگین دقت
30.43	29.57	31.34	A0	مدل سازی زبانی با ۲-گرام کلمات
26.22	27.58	25.0	A1	
74.54	81.09	68.97	A2	
32.16	33.82	30.66	A3	
62.74	56.28	70.88	A4	
15.38	14.81	16.0	A5	
56.				میانگین دقت

22.48	39.71	15.68	A0	مدل سازی زبانی با ۳-گرام کلمات
17.10	14.94	20.0	A1	
64.45	64.17	64.74	A2	
18.34	14.70	24.39	A3	
43.58	35.42	56.62	A4	
14.86	13.58	16.41	A5	
42.52				میانگین دقت

۴-۱-۲ آزمایش با مدل سازی زبانی تغییر یافته

نتایج ارزیابی مدل سازی زبانی تغییر یافته در جدول ۳-۲ نمایش داده شده است. نتایج روی این پایگاه داده یک دقت ۵,۴۴ درصدی را نسبت به بهترین نتیجه در روش مدل سازی زبانی با n-گرام کلمات، را نشان می‌دهد. نکته قابل توجه در این است که با وجود این که به علت ناکافی بودن داده‌های آموزشی مدل سازی زبانی ساده در سطح 2-گرام کلمه موفق نبوده است اما تلفیق 1-گرام کلمات با 2-گرام باعث بهبود نتایج شده است.

جدول ۴-۲: نتایج ارزیابی روش مدل سازی زبانی تغییر یافته در پایگاه داده WMPR-AA2016-A

مدل سازی زبانی تغییر یافته			نویسنده
<i>F-measure</i>	<i>Recall</i>	<i>Precision</i>	
41.40	38.73	46.49	A0
42.04	42.52	41.57	A1
80.27	88.13	73.71	A2
45.80	44.11	47.61	A3
70.21	66.33	74.57	A4
32.43	29.62	35.82	A5
65.71%			میانگین دقت

۴-۱-۳ آزمایش با مدل سازی و n-گرام کاراکترها

این آزمایش با استفاده از مدل سازی زبانی و n-گرام کاراکترها انجام شده است و آزمایش با افزایش مقدار n از ۱ تا ۶ تکرار شده است. میزان دقت در n=3 به بالاترین دقت در این سری از آزمایش‌ها

رسیده است و بعد از آن با افزایش مقدار n ، میزان دقت کاهش پیدا کرده است. نتایج در جدول ۳-۴ نمایش داده شده است.

جدول ۳-۴: نتایج ارزیابی روش مدل سازی زبانی با n -گرام کاراکترها در پایگاه داده WMPR-AA2016-A

میانگین دقت در مدل سازی زبانی ساده و n -گرام کاراکترها	n
50.69	2
62.95	3
60.27	4
60.27	5
59.13	6

۴-۱-۴ بررسی نتایج

در میان آزمایش‌های انجام شده در پایگاه داده WMPR-AA2016-A، بهترین نتیجه در مدل سازی زبانی تغییر یافته و کمترین نتیجه با مدل سازی زبانی در ۳-گرام کلمات حاصل شده است. نتایج حاصل از آزمایش‌ها در این پایگاه داده نشان دهنده آن است که دقت مدل سازی زبانی ساده در ۱-گرام کلمات بهتر از مدل سازی زبانی ساده در سطح ۲-گرام کلمات و حتی ۳-گرام کلمات است. این می تواند به این علت باشد که با توجه به طول کوتاه داده‌های آزمایشی در این پایگاه داده داده‌های آموزشی، برای فراهم کردن اطلاعات زبانی در سطح ۲-گرام و ۳-گرام به اندازه کافی فراهم نبوده است. در مدل سازی زبانی تغییر یافته نتایج از ۱-گرام کلمات نیز بهتر شده است. علت این بهبود می تواند استفاده همزمان از دو ویژگی ۱-گرام کلمات ۲-گرام کلمات همراه با ضریب IDF باشد که به نوعی توانسته است که کاستی‌های مربوط به کمبود اطلاعات زبانی را در هریک از این ویژگی‌ها به تنهایی جبران کند. در مدل سازی در سطح کاراکتر نتایج به دست آمده با تغییر مقدار n تغییر می-کند به طوری که در $n=3$ بهترین نتیجه به دست آمده است و این نتیجه از نتایج مدل سازی با ۱-گرام کلمات، ۲-گرام و ۳-گرام کلمات بهتر بوده است و در رتبه بندی آزمایش‌ها در این پایگاه داده در رتبه دوم قرار گرفته است. این نتیجه می‌تواند تایید کننده قدرت بالای کاراکتر n -گرام‌ها در ثبت اختلافات

جزئی سبکی باشد. در جدول ۴-۴ نتایج حاصل از آزمایشات در پایگاه داده WMPR-AA2016-A به ترتیب میانگین دقت، رتبه بندی شده است.

جدول ۴-۴: رتبه بندی نتایج آزمایشات با مدل سازی زبانی در پایگاه داده WMPR-AA2016-A

رتبه	روش	دقت
۱	مدل سازی زبانی تغییر یافته	۶۵,۷۱
۲	مدل سازی زبانی با ۳- گرام کاراکترها	62.95
3	مدل سازی زبانی با ۱- گرام کلمات	60.27%
4	مدل سازی زبانی ۴و۵- گرام کاراکترها	60.27
5	مدل سازی زبانی با ۶- گرام کاراکترها	59.13
6	مدل سازی زبانی با ۲-گرام کلمات	56
7	مدل سازی زبانی با ۲- گرام کاراکترها	50.69
8	مدل سازی زبانی ساده با ۳-گرام کلمات	42.52

۲-۴ : آزمایشات در پایگاه داده WMPR-AA2016-B

در آزمایشات صورت گرفته در این پایگاه داده هفت نویسنده در رقابت تعیین نویسنده ناشناس برای داده‌های آزمایشی شرکت دارند. برای تشکیل داده‌های آموزشی و داده‌های آزمایشی از مجموعه کل متون موجود برای هر نویسنده ۲۰٪ داده‌ها برای داده آزمایشی و ۸۰٪ داده‌ها برای داده آموزشی انتخاب شده است. انتخاب داده‌ها برای قرار گیری در دو مجموعه آزمایشی و آموزشی به صورت تصادفی انجام شده است. شکل ۴-۱ نمایی از این تقسیم بندی را نمایش می‌دهد.

۴-۲-۱ آزمایش با مدل سازی زبانی و n-گرام کلمات

نتایج حاصل از اجرای مدل سازی زبانی ساده با n-گرام کلمات با مقدار ۳ و ۱، ۲=n در جدول ۴-۵ نمایش داده شده است. برای روشن تر شدن ارزیابی، علاوه بر میانگین دقت، نتایج ارزیابی به تفکیک هر نویسنده نیز محاسبه شده است.

جدول ۴-۵: نتایج ارزیابی روش مدل سازی زبانی با n-گرام کلمات در پایگاه داده WMPR-AA2016-B

<i>F-measure</i>	<i>Recall</i>	<i>Precision</i>	نویسنده	مدل سازی زبانی با ۱-گرام کلمات
88.25	86.50	90.08	B0	
81.48	84.61	78.57	B1	
90.53	87.30	94.01	B2	
84.63	84.39	84.88	B3	
72.51	82.66	64.58	B4	
90.90	88.70	93.22	B5	
94.65	91.93	97.55	B6	
87.95			میانگین دقت	
82.88	86.50	79.56	B0	مدل سازی زبانی با ۲-گرام کلمات
67.11	54.94	86.20	B1	
84.57	90.24	79.56	B2	
79.27	76.30	82.5	B3	
72.30	62.66	85.45	B4	
88.83	94.08	84.13	B5	
93.04	94.52	91.62	B6	
۸۵,۱۱			میانگین دقت	
71.48	۷۹,۲۶	۶۵,۶۵	B0	مدل سازی زبانی با ۳-گرام کلمات
43.93	31.86	70.73	B1	
71.82	79.26	65.65	B2	
64.09	62.42	65.85	B3	
42.85	32.0	64.86	B4	
81.44	84.94	78.21	B5	
84.65	89.04	80.67	B6	
73.32			میانگین دقت	

۲-۲-۴ آزمایش با مدل سازی زبانی تغییر یافته

مشابه پایگاه داده WMPR-AA2016-A در آزمایشات در پایگاه داده WMPR-AA2016-B نیز، روش مدل سازی زبانی تغییر یافته یک بهبود ۳,۷ درصدی را نسبت به استفاده n-گرام کلمات نشان می‌دهد. نتایج ارزیابی در جدول ۴-۶ نمایش داده شده است.

جدول ۴-۶: نتایج ارزیابی روش مدل سازی زبانی تغییر یافته در پایگاه داده WMPR-AA2016-B

مدل سازی زبانی تغییر یافته			نویسنده
<i>F-measure</i>	<i>Recall</i>	<i>Precision</i>	
90.24	88.05	92.05	B0
83.33	76.92	90.90	B1
88.11	92.68	83.97	B2
88.51	91.32	85.86	B3
81.69	77.33	86.56	B4
95.23	96.77	93.75	B5
95.63	94.81	96.48	B6
91.56			میانگین دقت

۱-۲-۴ آزمایش با مدل سازی زبانی و n-گرام کاراکترها

این سری از آزمایش‌ها با استفاده از مدل سازی زبانی و n-گرام کاراکترها انجام شده است و آزمایش با افزایش مقدار n از ۱ تا ۶ تکرار شده است. میزان دقت در n=3 به بالاترین دقت در این سری از آزمایش‌ها رسیده است و بعد از آن با افزایش مقدار n، میزان دقت کاهش پیدا کرده است. نتایج در جدول ۴-۷ نمایش داده شده است.

جدول ۴-۷: نتایج ارزیابی روش مدل سازی زبانی با n-گرام کاراکترها در پایگاه داده WMPR-AA2016-B

n	میانگین دقت در مدل سازی زبانی ساده و n-گرام کاراکترها
2	71.25
3	85.88
4	76.16
5	65.54
6	63.33

۴-۲-۲ بررسی نتایج

در میان آزمایش‌های انجام شده در پایگاه داده WMPR-AA2016-B بهترین نتیجه در مدل سازی زبانی تغییر یافته و کمترین در مدل سازی زبانی با ۶-گرام کاراکترها به دست آمده است. نتایج حاصل از اجرای مدل سازی زبانی ساده با n-گرام کلمات مشابه پایگاه داده WMPR-AA2016-A نمایانگر بهتر بودن نتایج در مدل سازی زبانی با ۱-گرام کلمات نسبت به مدل سازی زبانی با ۲-گرام و ۳-گرام کلمات است. در مدل سازی با n-گرام کاراکترها نتایج به دست آمده با تغییر مقدار n تغییر می‌کند به طوری که در n=6 کمترین میزان دقت و در n=3 بهترین نتیجه به دست آمده است. نکته قابل توجه حاصل شدن بهترین نتیجه در هر دو پایگاه داده WMPR-AA2016-A و WMPR-AA2016-B با n=3 در گروه آزمایشات با n-گرام کاراکترها است. این می‌تواند تا حدودی نشان دهنده وابستگی مقدار n به زبان و سبک نوشتاری باشد.

به طور کلی در تمام آزمایش‌ها در این پایگاه داده نتایج نسبت به پایگاه داده WMPR-AA2016-A رشدی بالاتر از ۲۰ درصد را داشته است. این می‌تواند به علت دو عامل بیشتر شدن داده‌های آموزشی و همچنین بزرگتر شدن داده‌های آزمایشی به معنای تعداد کلمات در هر متن داده آزمایشی باشد. در جدول ۴-۸ نتایج حاصل از آزمایشات انجام شده در پایگاه داده WMPR-AA2016-B به ترتیب میانگین دقت، رتبه بندی شده است.

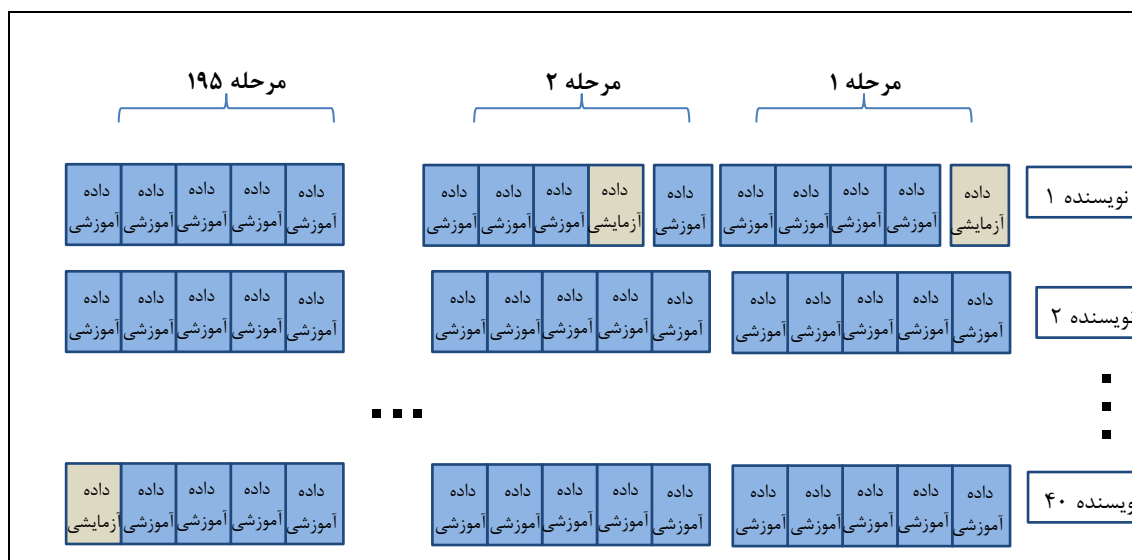
جدول ۴-۸: رتبه بندی نتایج آزمایشات با مدل سازی زبانی در پایگاه داده WMPR-AA2016-B

رتبه	روش	دقت
۱	مدل سازی زبانی تغییر یافته	91.56
2	مدل سازی زبانی با ۱-گرام کلمات	87.11
۳	مدل سازی زبانی با ۳-گرام کاراکترها	85.88
۴	مدل سازی زبانی با ۲-گرام کلمات	85.11
۵	مدل سازی زبانی با ۴-گرام کاراکترها	76.16
6	مدل سازی زبانی با ۳-گرام کلمات	73.32
7	مدل سازی زبانی با ۲-گرام کاراکترها	71.25
8	مدل سازی زبانی با ۵-گرام کاراکترها	65.54
9	مدل سازی زبانی با ۶-گرام کاراکترها	63.33

۴-۳ آزمایشات در پایگاه داده R40

این پایگاه داده که توضیح آن در بخش ۳-۲-۳ داده شد دارای ۴۰ نویسنده است و برای هر نویسنده بین ۴ تا پنج متن با طول کلمه‌ای بین ۸۰۰ تا ۹۰۰ کلمه موجود است. و در بین سه پایگاه داده در زبان فارسی جزء پایگاه داده با تعداد نویسنده کاندید زیاد قرار می‌گیرد. به علت کم بودن متن‌های موجود برای هر نویسنده از روش اعتبار سنجی متقابل با کنار گذاشتن یک متن^۱ استفاده شده است به این صورت که برای هر نویسنده یکی از متن‌های موجود به عنوان داده آزمایشی در نظر گرفته شده است و متن‌های باقی مانده برای تشکیل داده آموزشی برای هر نویسنده الحاق می‌گردد. به این ترتیب برای هر نویسنده بین چهار تا پنج مرحله وجود دارد و این مرحله به تعداد نویسنده‌گان تکرار می‌شود در مجموع در این پایگاه داده با ۴۰ نویسنده، ۱۹۵ مرحله تشکیل شده است. میانگین میزان دقت در هر مرحله به عنوان دقت روش عنوان شده است. شکل ۴-۲ تشکیل داده آموزشی و مجموعه داده آزمایشی در پایگاه داده R40 را نمایش می‌دهد.

^۱ Leave-one-out Cross Validation



شکل ۴-۲: تشکیل داده آموزشی و مجموعه داده آزمایشی در پایگاه داده R40

۴-۳-۱ آزمایش با مدل سازی زبانی و n-گرام کلمات

نتایج حاصل از اجرای مدل سازی زبانی ساده با n-گرام کلمات با مقدار ۳ و $n=2,1$ در جدول ۴-۹ نمایش داده شده است.

جدول ۴-۹: نتایج ارزیابی روش مدل سازی زبانی با n-گرام کلمات در پایگاه داده R40

میانگین دقت	نویسنده	مدل سازی زبانی با
100	R0-R38	۱-گرام کلمات
۸۰	R39	
۹۹,۴	میانگین دقت	
100	R0-R ^۸	مدل سازی زبانی با ۲-گرام کلمات
۶۰	R9	
100	R10	
80	R11	
50	R12	
100	R13-R39	
97.43	میانگین دقت	
100	R0-R1	مدل سازی زبانی با ۳-گرام کلمات
80	R2	
100	R3-R4	
80	R5	
100	R6	
80	R7	

100	B6-R10	
80	R11	
50	R12	
80	R13	
100	R14-R15	
80	R16	
100	R17-R26	
80	R27	
100	R28-R29	
75	R30	
100	R31	
80	R32	
100	R33-R37	
40	R38	
100	R39	
92.82	میانگین دقت	

۲-۳-۴ آزمایش با مدل سازی زبانی تغییر یافته

نتایج حاصل از مدل سازی زبانی تغییر یافته در جدول ۴-۱۰ نمایش داده شده است.

جدول ۴-۱۰: نتایج ارزیابی روش مدل سازی زبانی تغییر یافته در پایگاه داده R40

نویسنده	مدل سازی زبانی تغییر یافته
R0-R39	%۱۰۰
میانگین دقت	%۱۰۰

۳-۳-۴ آزمایش با مدل سازی و n-گرام کاراکترها

در این آزمایش مدل سازی زبانی با n-گرام کاراکترها انجام شده است و آزمایش با افزایش مقدار n از ۱ تا ۷ تکرار شده است. میزان دقت در n=2 به بالاترین دقت در این سری از آزمایشها رسیده است و بعد از آن با افزایش مقدار n، میزان دقت کاهش پیدا کرده است. نتایج در جدول 4-11 نمایش داده شده است.

جدول ۴-۱۱: نتایج ارزیابی روش مدل سازی زبانی و n-گرام کاراکترها در پایگاه داده R40

میانگین دقت	n
97.43	2
93.33	3
93.84	4
95.38	5
97.30	6
96.41	7

۴-۳-۴ بررسی نتایج

همانند دو پایگاه داده قبل، در این پایگاه داده نیز بهترین نتیجه در مدل سازی زبانی تغییر یافته و کمترین در مدل سازی زبانی در سطح کاراکتر با $n=3$ به دست آمده است.

نتایج در این پایگاه داده بهترین نتیجه در بین تمام نتایج به دست آمده در تمام پایگاه داده‌های استفاده شده را نمایش می‌دهد. که می‌تواند به علت طول داده آزمایشی مناسب و توزیع نسبتاً متعادل داده‌های آموزشی در بین نویسندگان کاندید باشد. همچنین نسبت به پژوهش پیشین (Ramezani, et al., 2013) در این پایگاه داده روش مدل سازی زبانی تغییر یافته نتیجه را بهبود داده است. در جدول ۵-۲ مقایسه بین مدل سازی زبانی تغییر یافته و کار (Ramezani, et al., 2013) نشان داده شده است.

نکته قابل توجه در آزمایش‌های n-گرام کاراکترها، حاصل شدن بهترین نتیجه در $n=2$ است. مقدار n با وجود ثابت بودن زبان اما تفاوت در سبک نوشتاری و همچنین تعداد نویسندگان کاندید تغییر کرده است. رتبه بندی نتایج آزمایشات به ترتیب میانگین دقت در جدول ۴-۱۲ نمایش داده شده است.

جدول ۴-۱۲: رتبه بندی نتایج آزمایشات با مدل سازی زبانی در پایگاه داده R40.

رتبه	نام روش	بهترین دقت
۱	مدل سازی زبانی تغییر یافته	100
2	مدل سازی زبانی با ۱-گرام کلمات	۹۹,۴
۳	مدل سازی زبانی با ۲-گرام کلمات	97.43
۴	مدل سازی زبانی با ۳-گرام کلمات	92.82
۵	مدل سازی زبانی با ۲-گرام کاراکترها	97.43
۶	مدل سازی زبانی با ۶-گرام کاراکترها	۹۷,۳۰
۷	مدل سازی زبانی با ۷-گرام کاراکترها	96.41
۸	مدل سازی زبانی با ۵-گرام کاراکترها	95.38
۹	مدل سازی زبانی با ۴-گرام کاراکترها	93.84
9	مدل سازی زبانی با ۳-گرام کاراکترها	93.33

۴-۴ زیر مجموعه با شش نویسنده از پایگاه داده RCV

همان طور که در بخش ۳-۵ در رابطه با خصوصیات پایگاه داده توضیح داده شد تعداد نویسندگان نویسنده یکی از عوامل موثر در مسئله تخصیص نویسنده است. از این جهت برای مقایسه بهتر نتایج حاصل از این پایگاه داده در زبان انگلیسی با نتایج دو پایگاه داده WMPR-AA2016-A و WMPR-AA2016-B در زبان فارسی، آزمایشی روی یک مجموعه با شش نویسنده از زیر مجموعه نویسندگان پایگاه داده RCV انجام شده است. برای شرکت تمامی نویسندگان در نتایج این ارزیابی، ارزیابی با ۳۰ تکرار آزمایش در زیر مجموعه‌های شش عضوی از مجموعه نویسندگان ۵۰ عضوی پایگاه داده RCV انجام شده است. سعی شده است انتخابها به گونه‌ای انجام شود که تمام نویسندگان در ارزیابی شرکت داشته باشند. میانگین این مجموعه سی تایی به عنوان دقت روش در نظر گرفته شده است.

۴-۴-۱ آزمایش با مدل سازی زبانی تغییر یافته و مدل سازی زبانی با n-گرام کلمات

نتایج حاصل از اجرای مدل سازی زبانی ساده با n-گرام کلمات با مقدار ۳ و ۱،۲=n و همچنین مدل سازی زبانی تغییر یافته در جدول ۴-۱۳ نمایش داده شده است.

جدول ۴-۱۳: نتایج ارزیابی روش مدل سازی زبانی با n-گرام کلمات و مدل سازی زبانی تغییر یافته در زیر مجموعه ۶ عضوی از پایگاه داده RCV

نویسنده	مدل سازی زبانی با ۱-گرام کلمات	مدل سازی زبانی با ۲-گرام کلمات	مدل سازی زبانی با ۳-گرام کلمات	مدل سازی زبانی تغییر یافته
میانگین دقت				
گروه ۱	۹۶	۹۶,۶۶	۹۶	۹۷
گروه ۲	۹۵	۹۶,۳۳	۹۵	۹۶
گروه ۳	۸۸,۶۶	۸۷,۳۳	۸۸	۸۸,۳۳
گروه ۴	۹۸,۶۶	۹۹	۹۷,۳۳	۹۹
گروه ۵	۹۶	۹۷,۳۳	۹۶,۳۳	۹۸
گروه ۶	۹۶,۳۳	۹۵,۳۳	۹۶	۹۶,۳۳
گروه ۷	۸۳	۸۵,۳۳	۸۸,۳۳	۸۶,۳۳
گروه ۸	۹۰,۶۶	۹۰,۳۳	۹۲,۶۶	۹۱,۶۶
گروه ۹	۹۶	۹۹,۳۳	۹۷,۶۶	۹۸,۶۶
گروه ۱۰	۹۴	۹۴	۹۵,۶۶	۹۵,۳۳
گروه ۱۱	۹۰,۶۶	۹۱,۶۶	۹۱	۹۱,۶۶
گروه ۱۲	۹۴,۳۳	۹۳,۶۶	۹۴	۹۵,۳۳
گروه ۱۳	۹۰,۳۳	۹۳	۹۴	۹۳,۳۳
گروه ۱۴	۹۱	۸۹	۸۹,۶۶	۹۰
گروه ۱۵	۹۵,۳۳	۹۶	۹۴,۳۳	۹۶
گروه ۱۶	۹۷,۳۳	۹۷,۶۶	۹۷	۹۸,۶۶
گروه ۱۷	۹۳,۳۳	۹۴,۳۳	۹۴,۶۶	۹۴,۳۳
گروه ۱۸	۹۲,۶۶	۹۴	۹۱	۹۴,۶۶
گروه ۱۹	۸۱,۳۳	۸۵,۶۶	۸۶,۶۶	۸۵,۳۳
گروه ۲۰	۹۷	۹۴,۶۶	۹۶	۹۶
گروه ۲۱	۹۱,۳۳	۹۰,۶۶	۹۱,۶۶	۹۱
گروه ۲۲	۹۲	۹۲,۳۳	۹۲	۹۲,۶۶
گروه ۲۳	۹۰	۸۹,۶۶	۸۶,۳۳	۹۱
گروه ۲۴	۷۷,۶۶	۸۲,۳۳	۸۴	۸۱,۳۳
گروه ۲۵	۹۲,۶۶	۹۴,۶۶	۹۵,۶۶	۹۴,۳۳

۹۵,۳۳	۹۶,۳۳	۹۴	۹۵	گروه ۲۶
۸۹	۹۱	۸۹,۶۶	۸۸	گروه ۲۷
۹۰	۹۰,۳۳	۹۰	۹۰,۶۶	گروه ۲۸
۹۲	۹۲,۳۳	۹۲,۶۶	۸۸,۳۳	گروه ۲۹
۹۸	۹۷,۳۳	۹۷,۳۳	۹۷,۶۶	گروه ۳۰
93.22	92.94	92.8	92.03	میانگین

۴-۴-۲ بررسی نتایج

نکته قابل توجه در نتایج این سری از آزمایش‌ها در این است که با وجود بهتر شدن نتیجه مدل سازی زبانی در ۳-گرام کلمات نسبت به ۲-گرام و ۱-گرام کلمات، بهترین نتیجه در مدل سازی زبانی تغییر یافته به دست آمده است. البته اختلاف بهبود ایجاد شده کمتر از بهبود ایجاد شده در دو پایگاه داده WMPR-AA2016-A و WMPR-AA2016-B است. نتایج در این پایگاه داده از دو پایگاه داده WMPR-AA2016-A و WMPR-AA2016-B بهتر است. اما باید توجه داشت طول داده آزمایشی در این پایگاه داده بزرگتر است و همچنین پایگاه داده متعادل تری است. رتبه‌بندی نتایج آزمایشات بر اساس میانگین دقت به دست آمده در جدول ۴-۱۴ نمایش داده شده است.

جدول ۴-۱۴: رتبه بندی نتایج آزمایشات با مدل سازی زبانی در پایگاه داده RCV

رتبه	نام روش	بهترین دقت
۱	مدل سازی زبانی تغییر یافته	۹۳,۳۴
۲	مدل سازی زبانی با ۳-گرام کلمات	۹۲,۹۴
۴	مدل سازی زبانی با ۲-گرام کلمات	۹۲,۸۰
۵	مدل سازی زبانی با ۳۱-گرام کلمات	92.03

۴-۵ آزمایشات در پایگاه داده RCV

این پایگاه داده در زبان انگلیسی بوده و با ۵۰ نویسنده در دسته پایگاه داده با نویسنده زیاد محسوب می‌گردد. برای داده آزمایشی و آموزشی برای هر نویسنده ۵۰ متن وجود دارد. در روش مدل سازی زبانی همان‌طور که عنوان شد از آنجا که یک روش بر پایه پروفایل است، برای تشکیل داده آموزش

برای هر نویسنده تمامی متن‌های آموزشی هر نویسنده الحاق می‌گردد و یک متن آموزشی برای هر نویسنده تشکیل می‌شود.

۴-۵-۱ آزمایش با مدل سازی زبانی تغییر یافته و مدل سازی زبانی با n-گرام کلمات

نتایج حاصل از مدل سازی زبانی با ۱-گرام ، ۲-گرام و ۳-گرام کلمات و همچنین مدل سازی زبانی تغییر یافته در جدول ۴-۱۵ نمایش داده شده است.

جدول ۴-۱۵: نتایج ارزیابی روش مدل سازی زبانی با n-گرام کلمات و مدل سازی زبانی تغییر یافته در پایگاه داده RCV

نویسنده	مدل سازی زبانی با ۱-گرام کلمات	مدل سازی زبانی با ۲-گرام کلمات	مدل سازی زبانی با ۳-گرام کلمات	مدل سازی زبانی تغییر یافته
میانگین دقت				
C0-c49	67.48	70.	69.92	70.48

۴-۵-۲ آزمایش با مدل سازی زبانی و n-گرام کاراکترها

در این آزمایش مدل سازی زبانی در سطح کاراکتر انجام شده است. آزمایش با افزایش مقدار n از ۲ تا n=6 تکرار شده است. میزان دقت در n=3 به بالاترین دقت در این روش رسیده است و بعد از آن با افزایش مقدار n ، میزان دقت کاهش پیدا کرده است. نتایج در جدول ۴-۱۶ نمایش داده شده است.

جدول ۴-۱۶: نتایج ارزیابی روش مدل سازی زبانی با n-گرام کاراکترها در پایگاه داده RCV

n	آزمایش با مدل سازی و n-گرام کاراکترها
2	61.40
۳	63.96
4	63
5	63.24

۳-۵-۴ بررسی نتایج

در این پایگاه داده بهترین نتیجه برای مدل سازی تغییر یافته به دست آمده است. اگرچه این نتیجه با نتیجه مدل سازی در سطح ۲-گرام کلمه اختلاف زیادی ندارد (48٪). به طور کلی نتایج به دست آمده در n-گرام کلمات بهتر از حاصل n-گرام کاراکترها است. در مقایسه دو پایگاه داده R40 و پایگاه داده RCV به عنوان دو پایگاه داده با تعداد نویسنده زیاد و در دو زبان متفاوت، پایگاه داده R40 با بهترین نتیجه 100٪، نتیجه بسیار بهتری را نسبت به بهترین نتیجه 70.48 در RCV، نمایش می دهد. البته عواملی مانند تعداد نویسنده گان بیشتر در پایگاه داده RCV نسبت به پایگاه داده R40، می تواند در اختلاف نتایج موثر باشد. رتبه بندی نتایج آزمایشات بر اساس میانگین دقت در جدول 4-17 نمایش داده شده است. در پژوهش (Stamatatos, 2006) از این پایگاه داده استفاده شده است در جدول 5-3 مقایسه بین نتایج مدل سازی زبانی تغییر یافته با این پژوهش نمایش داده شده است. در این مقایسه مدل سازی زبانی تغییر یافته با دقت 70.48، در رتبه ششم قرار گرفته است که با دقت در رتبه اول، به اندازه 3.56 اختلاف دارد.

جدول 4-17: رتبه بندی نتایج آزمایشات در پایگاه داده RCV

رتبه	روش	بهترین دقت
1	مدل سازی زبانی تغییر یافته	70.48
2	مدل سازی زبانی ساده با ۲-گرام کلمات	70
3	مدل سازی زبانی ساده با ۳-گرام کلمات	69.92
4	مدل سازی زبانی ساده با ۱-گرام کلمات	67.48
5	مدل سازی زبانی ساده با ۳-گرام کاراکترها	63.96
6	مدل سازی زبانی ساده با ۲-گرام کاراکترها	61.40
7	مدل سازی زبانی ساده با ۵-گرام کاراکترها	63.24
8	مدل سازی زبانی ساده با ۴-گرام کاراکترها	63

۴-۶ بررسی خصوصیات پایگاه داده

در این قسمت به بررسی خصوصیات مختلف یک پایگاه داده که در بخش ۳-۵ توضیح داده شد و تاثیری که روی نتایج می‌تواند داشته باشد می‌پردازیم. ویژگی‌های مورد بحث تعداد نویسنده‌گان، اندازه داده آموزشی و متعادل بودن یا نامتعادل بودن داده‌های آموزشی در پایگاه داده است.

۴-۶-۱ تعداد نویسنده در مجموعه نویسنده‌گان کاندید

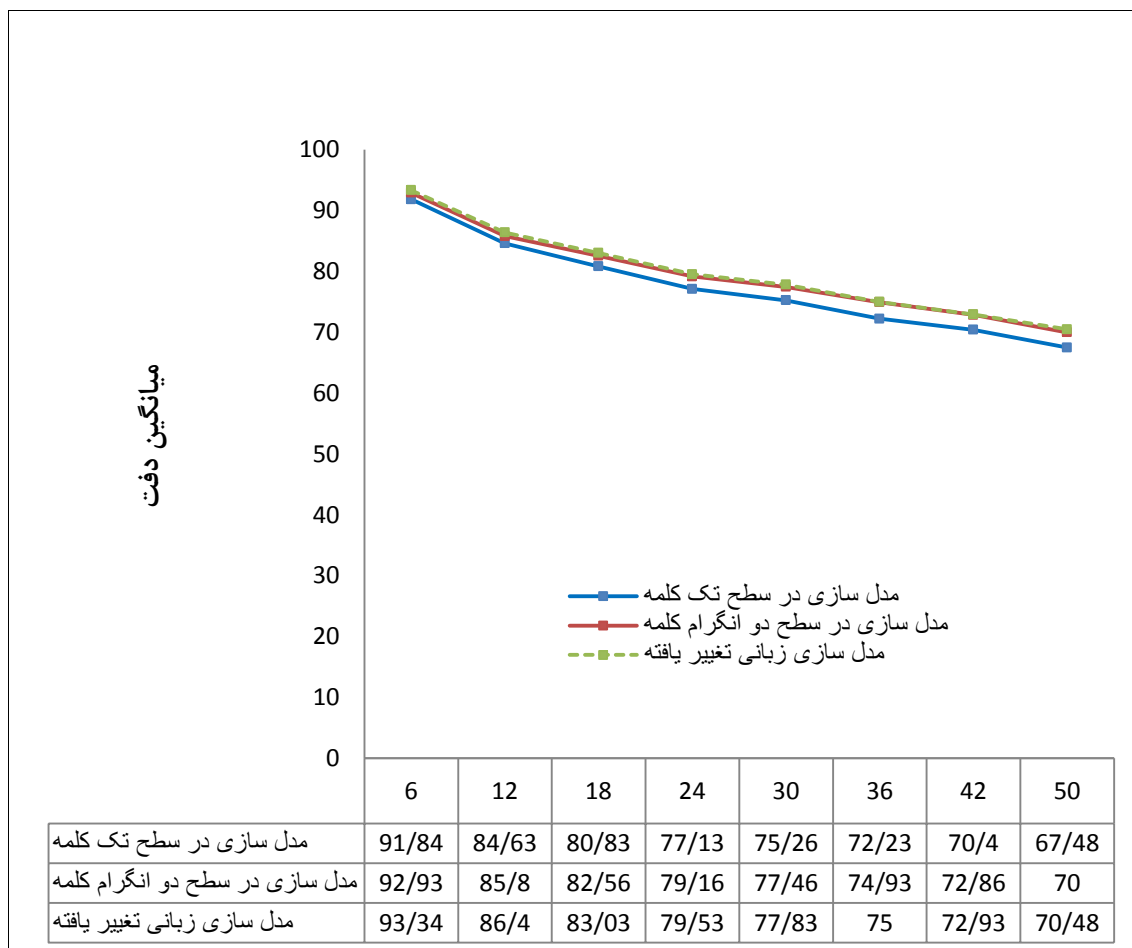
پایگاه داده RCV با ۵۰ نویسنده در دسته پایگاه داده با تعداد نویسنده زیاد قرار می‌گیرد به همین علت برای روشن شدن اثر تعداد نویسنده در میزان دقت در حل مسئله شناسایی نویسنده و بررسی این که ویژگی‌های آزمایش شده چقدر در برابر افزایش تعداد نویسنده مقاومت دارند و نتایج چگونه تغییر پیدا می‌کند؛ آزمایشی به صورتی که شرح داده خواهد شد ترتیب داده شده است.

۴-۶-۱-۱ آماده سازی زیر مجموعه‌ها

از مجموعه نویسنده‌گان پایگاه داده RCV به ترتیب زیر مجموعه‌هایی با تعداد نویسنده‌گان شش، دوازده، هجده، بیست و چهار، سی، سی و شش، چهل و دو و پنجاه عضوی انتخاب شده است. برای به دست آوردن میزان دقت مدل در هر یک از زیر مجموعه‌ها به صورت تصادفی سی زیر مجموعه متفاوت با تعداد نویسنده‌گان یکسان انتخاب شده است و میانگین دقت در سی زیر مجموعه به عنوان دقت زیر مجموعه عنوان شده است. به عنوان نمونه برای محاسبه میزان دقت زیر مجموعه شش تایی، سی زیر مجموعه ۶ تایی از مجموعه ۵۰ تایی نویسنده‌گان انتخاب شده است و میانگین دقت سی زیر مجموعه به عنوان دقت در مجموعه شش تایی در نظر گرفته شده است.

۴-۶-۱-۲ بررسی نتایج:

نتیجه افزایش تعداد نویسنده گان در شکل ۳-۴ نمایش داده شده است. با افزایش تعداد نویسنده گان، میانگین دقت در هر مدل کاهش پیدا می کند. در تمامی گروه ها میانگین دقت در روش مدل سازی زبانی تغییر یافته بالاتر از دو مدل دیگر است.



شکل ۴-۳: اثر افزایش تعداد نویسنده گان در میزان دقت

۴-۶-۲ کاهش داده آموزشی

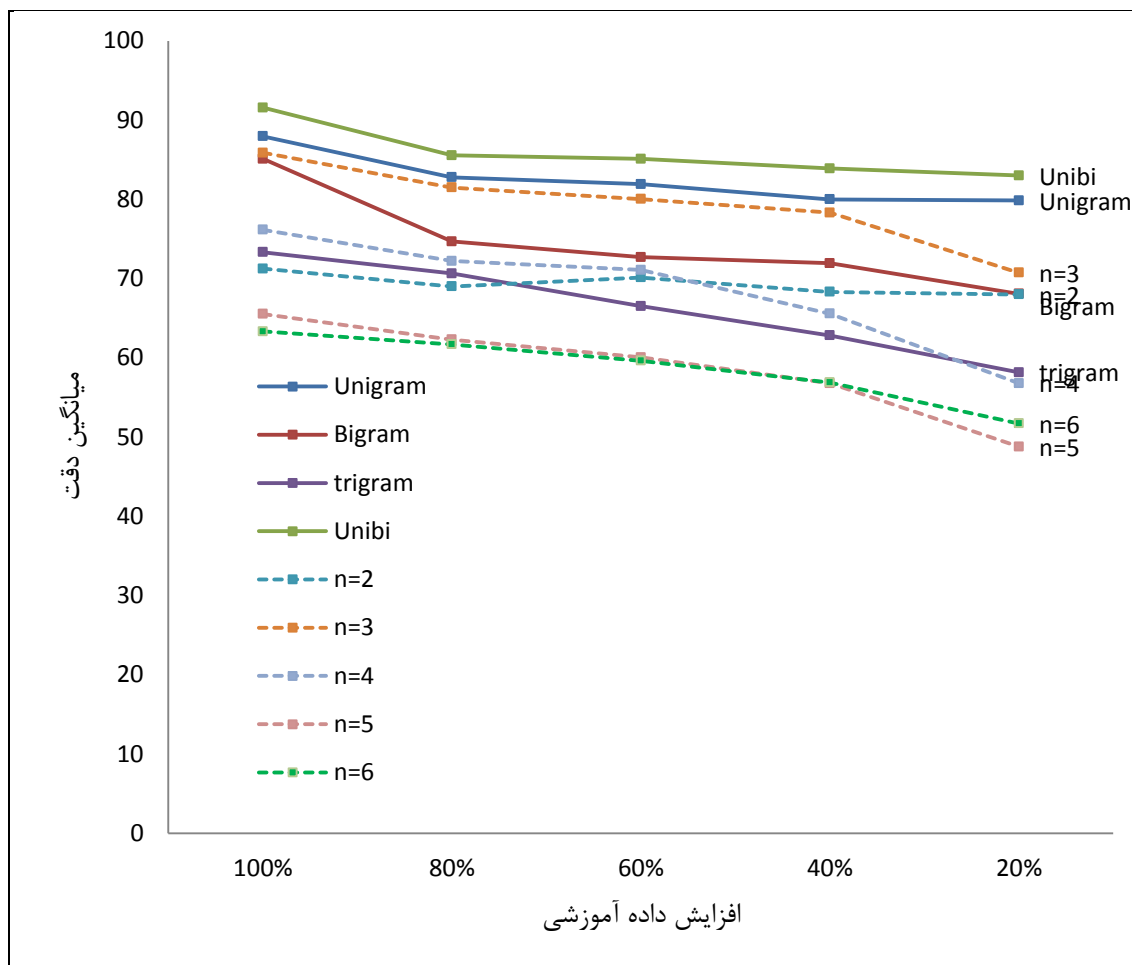
این آزمایش برای بررسی دقت مدل سازی با افزایش داده آموزشی در حل مسئله شناسایی نویسنده در نظر گرفته شده است. دو پایگاه داده WMPR-AA2016-A و WMPR-AA2016-B به خاطر داده آموزشی زیاد، فارسی بودن زبان و ویژگی متن کوتاه در پایگاه داده WMPR-AA2016-A انتخاب شده اند.

۴-۶-۲-۱ آماده سازی داده

داده‌های آموزشی در اندازه ۲۰، ۴۰، ۶۰، ۸۰ و ۱۰۰ درصد از داده آموزشی در دو پایگاه داده WMPR-AA2016-A و WMPR-AA2016-B که در قسمت ۴-۱ و ۴-۲ توضیح داده شده است، انتخاب شده‌اند. داده‌های آزمایشی داده‌های معرفی شده در قسمت ۴-۱ و ۴-۲ است و در تمامی تکرار آزمایش با افزایش داده آموزشی یکسان است.

۴-۶-۲-۲ بررسی نتایج

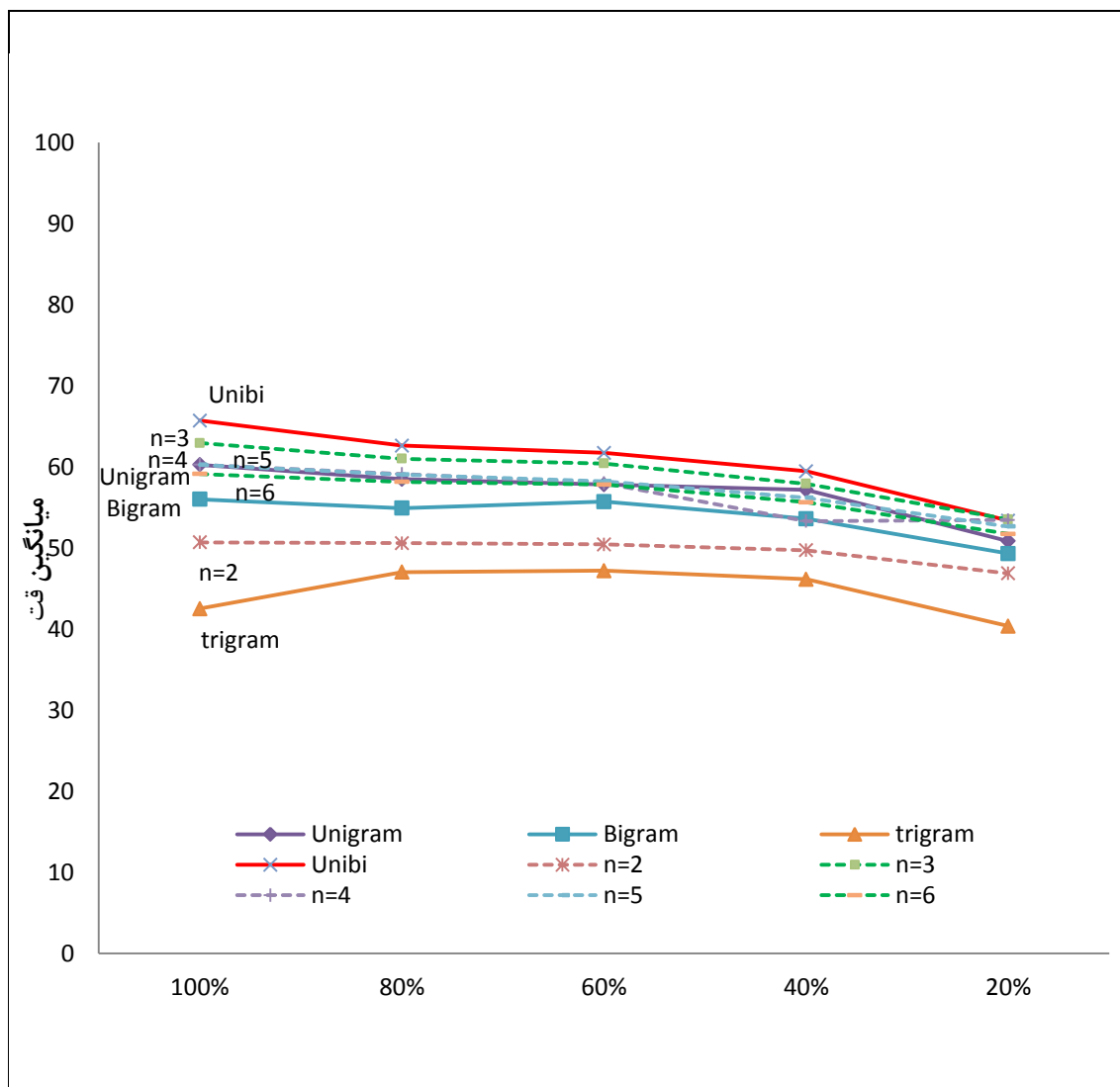
نتایج حاصل از کاهش داده آموزشی با مدل سازی‌ها در دو پایگاه داده WMPR-AA2016-B و WMPR-AA2016-A به ترتیب در شکل ۴-۴ و ۴-۵، نمایش داده شده است. در هر دو پایگاه داده در تمامی اندازه‌های داده آموزشی، مدل سازی تغییر یافته بهتر از سایر مدل سازی‌ها عمل می‌نماید و کاهش داده آموزشی در تمام مدل سازی‌ها، با کاهش میزان دقت همراه است. مقدار دقت در زمانی که فقط از ۲٪ داده‌های آموزشی استفاده می‌شود در پایگاه داده WMPR-AA2016-A از ۶۵٫۷۱ به ۵۳٫۳۷ کاهش پیدا کرده است و در پایگاه داده WMPR-AA2016-B از ۹۱٫۵۶ به ۸۳ کاهش پیدا کرده است.



شکل ۴-۴: اثر کاهش داده آموزشی در میزان دقت در پایگاه داده WMPR-AA2016-B

جدول ۴-۱۸: دقت در نتایج ارزیابی روش مدل سازی زبانی با کاهش داده آموزشی در پایگاه داده WMPR-AA2016-B

نسبت داده آموزشی	مدل سازی زبانی با ۱-گرام کلمات	مدل سازی زبانی با ۲-گرام کلمات	مدل سازی زبانی با ۳-گرام کلمات	مدل سازی زبانی با ۴-گرام کلمات	مدل سازی زبانی با ۵-گرام کلمات	مدل سازی زبانی با ۶-گرام کلمات	مدل سازی زبانی با ۷-گرام کلمات	مدل سازی زبانی با ۸-گرام کلمات	مدل سازی زبانی با ۹-گرام کلمات	مدل سازی زبانی با ۱۰-گرام کلمات
۲۰٪	79.86	68.07	58.17	83	68	70.74	56.79	48.79	51.72	
۴۰٪	80	71.94	62.82	83.9	68.33	78.31	65.57	56.79	56.88	
۶۰٪	81.92	72.71	66.52	85.11	70.13	80.03	71.08	60.06	59.63	
۸۰٪	82.78	74.69	۷۰,۶۵	85.55	69.01	81.49	72.22	62.3	61.7	
۱۰۰٪	87.95	85.11	73.32	91.56	71.25	85.88	76.16	65.54	63.33	



شکل ۴-۵: اثر کاهش داده آموزشی در میزان دقت WMPR-AA2016-A

جدول ۴-۱۹: دقت در نتایج ارزیابی روش مدل سازی با کاهش داده آموزشی در پایگاه داده WMPR-AA2016-A

نسبت داده آموزشی	مدل سازی زبانی با ۱-گرام کلمات	مدل سازی زبانی با ۲-گرام کلمات	مدل سازی زبانی با ۳-گرام کلمات	مدل سازی زبانی با تغییر یافته	مدل سازی زبانی با ۲-گرام کاراکترها	مدل سازی زبانی با ۳-گرام کاراکترها	مدل سازی زبانی با ۴-گرام کاراکترها	مدل سازی زبانی با ۵-گرام کاراکترها	مدل سازی زبانی با ۶-گرام کاراکترها
۲۰٪	50.85	49.3	40.37	53.37	46.87	53.53	53.45	52.64	51.74
۴۰٪	57.18	53.61	46.14	59.46	49.71	57.92	53.32	56.21	55.64
۶۰٪	57.83	55.72	47.19	61.73	50.44	60.43	58	58.24	57.83
۸۰٪	58.48	54.91	47.03	62.63	50.6	61	59.13	59.05	58.16
۱۰۰٪	60.27	56	42.52	65.71	50.69	62.95	60.27	60.27	59.13

۳-۶-۴ متعادل کردن داده آموزشی

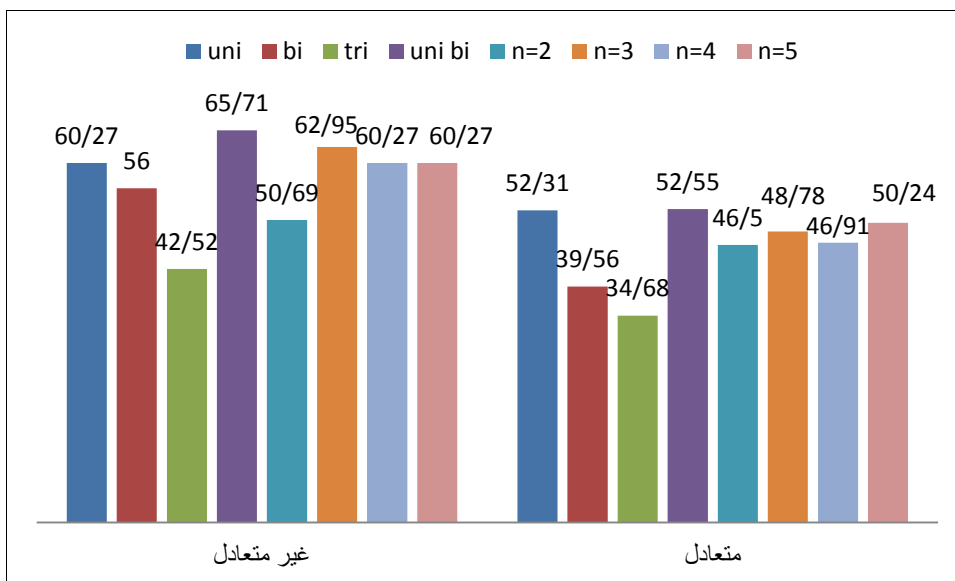
در این بخش برای بررسی اثر متعادل شدن داده آموزشی آزمایشی انجام شده است. و برای این منظور دو پایگاه داده کاملاً نامتعادل WMPR-AA2016-A و WMPR-AA2016-B در نظر گرفته شده است.

۱-۳-۶-۴ آماده سازی داده آموزشی

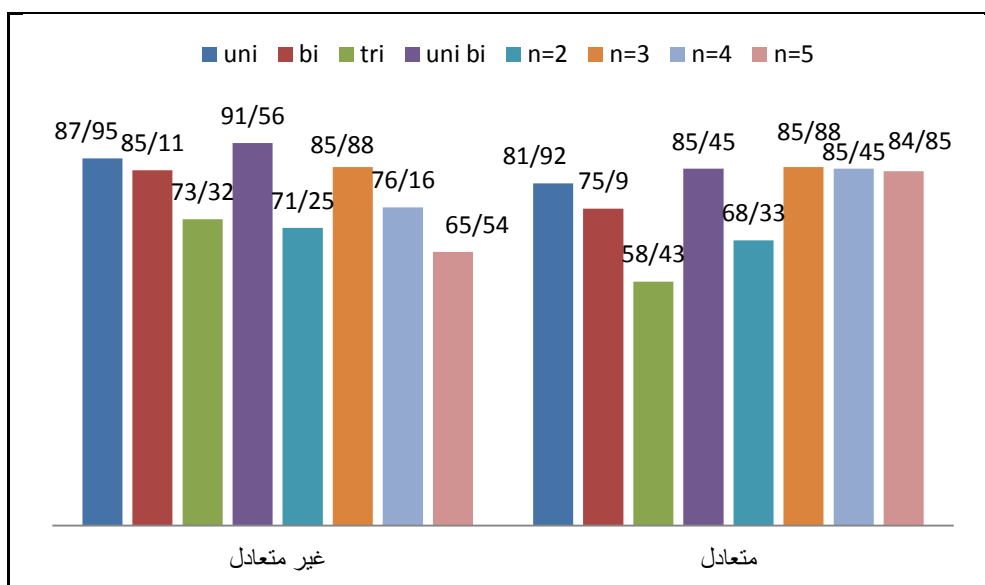
برای متعادل کردن از روش کم کردن نمونه‌ها به صورت تصادفی استفاده شده است. و به این منظور تا یکنواخت شدن توزیع داده‌های آموزشی در نویسندگان کاندید از داده‌های آموزشی به صورت تصادفی داده کم شده است.

۲-۳-۶-۴ بررسی نتایج

نتایج در جدول ۴-۶ و ۴-۷ برای دو پایگاه داده WMPR-AA2016-A و WMPR-AA2016-B نمایش داده شده است. در پایگاه داده WMPR-AA2016-A پس از متعادل سازی در هیچکدام از روش‌ها مدل سازی بهبودی ایجاد نشده است. این می‌تواند از چندین دلیل مانند کارا نبودن روش متعادل سازی انتخابی یا کافی نبودن داده‌های آموزشی در اثر کم شدن داده‌های آموزشی ناشی شده باشد. در پایگاه داده WMPR-AA2016-B بعد از متعادل سازی نتایج در مدل سازی در سطح کارا کتر بهبود داشته است. به طوری که در $n=3$ با اختلاف کمی، بهتر از مدل تغییر یافته عمل می‌کند. اما در سطح کلمه همچنان داده‌های غیر متعادل بهتر عمل می‌کنند. به طور کلی عملکرد داده‌های متعادل شده در پایگاه داده WMPR-AA2016-B بهتر از پایگاه داده WMPR-AA2016-A عمل کرده است که می‌تواند به علت داده آموزشی بالاتر و طول داده آزمایشی بیشتر باشد.



شکل ۴-۶: تغییرات دقت در مدل سازی در دو حالت متعادل و غیر متعادل در پایگاه داده WMPR-AA2016-A.



شکل ۴-۷: تغییرات دقت در مدل سازی در دو حالت متعادل و غیر متعادل در پایگاه داده WMPR-AA2016-B.

۵- فصل پنجم: نتیجه‌گیری و پیشنهادات

۵-۱ جمع بندی

در فصل اول در بیان مسئله تشخیص نویسنده، به عنوان تخصیص نویسنده برای یک متن دیده نشده بیان شد و ضرورت حل آن با توجه به کاربردهای مختلف آن در حوزه‌های مختلف از جمله تجارت الکترونیکی، شناسایی نویسنده در متن‌های ادبی، شناسایی تروریسم با استفاده از پیام‌ها، تشخیص نویسنده در ایمیل و متن‌های الکترونیکی، تشخیص سرقت ادبی و .. و همچنین گسترش استفاده و انتشار متن با گسترش استفاده از اینترنت و صفحات وب، بیان شد.

در فصل دوم به دسته بندی ویژگی‌ها در دسته‌های لغوی، کاراکتری، نحوی، معنایی و وابسته به کاربرد پرداخته شد. و به چند نمونه از کارهای انجام شده اشاره شد. مزیت ویژگی‌ها در دسته لغوی سادگی و بی‌نیازی به ابزارهای پردازش زبان طبیعی عنوان شد و این که این مزیت آنها را تبدیل به ویژگی مستقل از زبان می‌کند. همچنین بیان شد که وجود کارهای زیادی که با استفاده از کاراکتر n -گرام‌ها انجام شده است و نتایج خوبی گزارش شده نشان از مفید بودن آنها دارد. و مزیت آن در مستقل از زبان بودن، مشخص کردن اختلافات جزئی در سبک و تحمل پذیری خوب در برابر نویز بیان شد. از طرفی مشکل تعیین مقدار n به صورت کارا در استفاده از آنها مطرح شد. در رابطه با ایده استفاده از ویژگی‌های ساختاری گفته شد که نویسنده‌گان الگوهای نحوی را به صورت ناخودآگاه استفاده می‌کنند و بنابراین در تمامی نوشته‌های خود الگوهای مشابه‌ای را دنبال خواهند کرد. توجه در ویژگی‌های معنایی روی معنا و نقش معنایی کلمات عنوان شد و در نهایت ویژگی بسته به کاربرد استفاده از ویژگی‌های مختص حوزه‌های مختلف، مانند سبک و رنگ فونت در فروم‌ها و یا پیام‌های الکترونیکی تعریف شد. بعد از معرفی ویژگی‌ها، ابتدا به معرفی انواع روش‌های تشکیل داده آموزشی در دو دسته مبتنی بر نمونه و مبتنی بر پروفایل پرداخته شد و دسته بندی متدهای تخصیص در پژوهش‌های پیشین در دسته‌های متدها بر اساس فاصله، روش‌های فشرده سازی و روش‌های یادگیری ماشین مانند شبکه‌های عصبی، بردار ماشین و درخت تصمیم و همچنین روش‌های احتمالاتی مانند بیز و

مدل‌سازی زبانی انجام شد و در انتها مروری بر کارهای پیشین در پایگاه داده با زبان فارسی صورت گرفت.

در فصل سوم به معرفی چهار پایگاه داده استفاده شده در این پایان نامه پرداخته شد. سه پایگاه داده در زبان فارسی و پایگاه داده چهارم در زبان انگلیسی معرفی شد. و تفاوت آنها در سبک متفاوت نظم و نثر، اندازه داده آموزشی و آزمایشی و تعداد نویسنده کاندید بیان شد. پس از آن مراحل مورد نیاز برای حل مسئله شناسایی نویسنده پیش پردازش و انتخاب ویژگی و متد تخصص ذکر شد و n -گرام کاراکترها و کلمات به عنوان ویژگی و مدل سازی زبانی برای متد تخصیص شرح داده شده‌اند. پس از آن روشی با ترکیب ۱-گرام کلمات با ۲-گرام کلمات و استفاده از مقدار IDF به عنوان وزن دهی احتمال n -گرامها (بر اساس میزان تکرار در داده‌های آموزشی هر نویسنده) در مدل سازی زبانی با نام مدل سازی تغییر یافته، پیشنهاد شد. در نهایت در بخش آخر به معرفی خصوصیات از پایگاه داده که در نتایج ارزیابی در حل مسئله تخصیص نویسنده موثر هستند پرداخته شد. و از آنها به عنوان مقدار داده آموزشی، تعداد نویسندگان کاندید و توزیع داده های آموزشی بین نویسندگان کاندید نام برده شد.

در فصل چهارم به ارزیابی حل مسئله نویسنده با روش مدل سازی زبانی ساده با استفاده از ۱-گرام ، ۲-گرام و ۳-گرام کلمات، n -گرامها کاراکترها با تغییر مقدار n و مدل سازی زبانی تغییر یافته در چهار پایگاه داده پرداخته شد. نکته قابل توجه در n -گرام کاراکترها کسب بهترین نتیجه در دو پایگاه داده با زبان و سبک نگارشی و تعداد نویسنده کاندید مشابه در مقدار $n=3$ است. جدا از نتایج در مدل سازی تغییر یافته می توان گفت در دو پایگاه داده با نویسنده زیاد n -گرام کلمات و در دو پایگاه داده دیگر n -گرام کاراکترها نتایج بهتری داشته است. نتایج بهتر n -گرام کاراکترها در دو پایگاه داده با طول داده آزمایشی کوتاه تر، می تواند تاییدکننده قدرت n -گرام کاراکترها در تشخیص اختلافات جزئی در سبک نوشتاری باشد. روش مدل سازی ترکیبی در تمام آزمایشات انجام شده در رتبه اول قرار گرفته است، اگرچه اختلاف آن با روش های مدل سازی بر اساس n -گرام کلمه در پایگاه RCV با تعداد نویسنده

زیاد، کم است اما در دو پایگاه داده با تعداد نویسنده کم، که یکی در دسته داده‌های متن کوتاه نیز قرار می‌گیرد، با بهبود خوبی نسبت به سایر روش‌ها در رتبه اول قرار گرفته است. در پژوهش (Stamatatos, 2006) از این پایگاه داده استفاده شده است در جدول ۵-۳ مقایسه بین نتایج مدل-سازي زبانی تغییر یافته با این پژوهش نمایش داده شده است. در این مقایسه مدل سازی زبانی تغییر یافته با دقت ۷۰,۴۸ در رتبه ششم قرار گرفته است که با دقت در رتبه اول، به اندازه ۳,۵۶ اختلاف دارد. همچنین در پایگاه داده R40 نتیجه با مدل سازی تغییر یافته نسبت به بهترین نتیجه در پژوهش پیشین (Ramezani, et al., 2013) انجام شده بهبود خوبی را نشان می‌دهد. نتایج در این پژوهش با مدل سازی زبانی تغییر یافته در جدول ۵-۲ نمایش داده شده است. در این مقایسه در رتبه اول قرار گرفته است.

در بررسی کاهش مقدار داده آموزشی مشخص شد که کاهش داده آموزشی در تمامی پایگاه داده‌ها باعث کاهش دقت در نتایج شده است. مقدار دقت در زمانی که فقط از ۲۰٪ داده‌های آموزشی استفاده می‌شود در پایگاه داده WMPR-AA2016-A از ۶۵,۷۱ به ۵۳,۳۷ کاهش پیدا کرده است و در پایگاه داده WMPR-AA2016-B از ۹۱,۵۶ به ۸۳ کاهش پیدا کرده است. از کاهش ۸ درصدی در پایگاه داده WMPR-AA2016-B و همچنین ثابت ماندن نسبی رتبه بندی روش‌ها شاید بتوان این نتیجه را گرفت که روش مدل سازی زبانی با n -گرام کاراکترها و n -گرام کلمات نسبت به کاهش داده آموزشی رفتار نسبتاً پایداری را داشته است. در پایگاه داده WMPR-AA2016-A پس از متعادل سازی در هیچکدام از روش‌ها مدل سازی بهبودی ایجاد نشده است. این می‌تواند از دلایلی مانند کارا نبودن روش متعادل سازی انتخابی یا کافی نبودن داده‌های آموزشی در اثر کم شدن داده‌های آموزشی ناشی شده باشد. در پایگاه داده WMPR-AA2016-B بعد از متعادل سازی، نتایج در مدل سازی با n -گرام کاراکترها بهتر شده است، به طوری که حتی در $n=3$ با اختلاف کم، بهتر از مدل تغییر یافته عمل می‌کند (قبل از متعادل سازی روش مدل تغییر یافته در رتبه بهترین نتیجه را داشته است) اما در سطح کلمه همچنان داده‌های غیر متعادل بهتر عمل می‌کنند. به طور کلی عملکرد داده‌های متعادل

شده در پایگاه داده WMPR-AA2016-B بهتر از پایگاه داده WMPR-AA2016-A عمل کرده است

که می‌تواند به علت داده آموزشی بالاتر و طول داده آزمایشی بیشتر باشد.

جدول ۵-۱: میانگین دقت در تمامی آزمایشات انجام شده در چهار پایگاه داده

پایگاه داده WMPR-AA2016-B			پایگاه داده WMPR-AA2016-A		
دقت	روش	رتبه	دقت	روش	رتبه
91.56	مدل سازی زبانی تغییر یافته	۱	۶۵,۷۱	مدل سازی زبانی تغییر یافته	۱
87.11	مدل سازی زبانی با ۱-گرام کلمات	2	62.95	مدل سازی زبانی با ۳-گرام کاراکترها	۲
85.88	مدل سازی زبانی با ۳-گرام کاراکترها	۳	60.27	مدل سازی زبانی با ۱-گرام کلمات	3
85.11	مدل سازی زبانی با ۲-گرام کلمات	۴	60.27	مدل سازی زبانی ۴-۵-گرام کاراکترها	4
76.16	مدل سازی زبانی با ۴-گرام کاراکترها	۵	59.13	مدل سازی زبانی با ۶-گرام کاراکترها	5
73.32	مدل سازی زبانی با ۳-گرام کلمات	6	56	مدل سازی زبانی با ۲-گرام کلمات	6
71.25	مدل سازی زبانی با ۲-گرام کاراکترها	7	50.69	مدل سازی زبانی با ۲-گرام کاراکترها	7
65.54	مدل سازی زبانی با ۵-گرام کاراکترها	8	42.52	مدل سازی زبانی ساده با ۳-گرام کلمات	8
63.33	مدل سازی زبانی در سطح کاراکتر با n=6	9			
پایگاه داده RCV			پایگاه داده R40		
70.48	مدل سازی زبانی تغییر یافته	۱	100	مدل سازی زبانی تغییر یافته	۱
70	مدل سازی زبانی ساده با ۲-گرام کلمات	2	۹۹,۴	مدل سازی زبانی با ۱-گرام کلمات	2
69.92	مدل سازی زبانی ساده با ۳-گرام کلمات	3	97.43	مدل سازی زبانی با ۲-گرام کلمات	۳
67.48	مدل سازی زبانی ساده با ۱-گرام کلمات	4	92.82	مدل سازی زبانی با ۳-گرام کلمات	۴
63.96	مدل سازی زبانی ساده با ۳-گرام کاراکترها	5	97.43	مدل سازی زبانی با ۲-گرام کاراکترها	۵
61.40	مدل سازی زبانی ساده با ۲-گرام کاراکترها	6	۹۷,۳۰	مدل سازی زبانی با ۶-گرام کاراکترها	۶
63.24	مدل سازی زبانی ساده با ۵-گرام کاراکترها	7	96.41	مدل سازی زبانی با ۷-گرام کاراکترها	۷
63	مدل سازی زبانی ساده با ۴-گرام کاراکترها	8	95.38	مدل سازی زبانی با ۵-گرام کاراکترها	۸
			93.84	مدل سازی زبانی با ۴-گرام کاراکترها	۹
پایگاه داده زیر مجموعه ۶ تایی RCV					
۹۳,۳۴	مدل سازی زبانی تغییر یافته	۱			
92.94	مدل سازی زبانی در سطح سه کلمه	2			
92.94	مدل سازی زبانی در سطح دو کلمه	4			
92.03	مدل سازی زبانی ساده در سطح تک کلمه	5			

جدول ۵-۲: مقایسه نتایج مدل سازی زبانی تغییر یافته با نتایج در (Ramezani, et al., 2013)

روش	ویژگی	دقت
مدل سازی زبانی تغییر یافته		۱۰۰
بردار ماشین	۲- گرم کاراکتر	۸۴
بردار ماشین	تکرار فعل‌ها	۷۶

جدول ۵-۳: مقایسه نتایج مدل سازی زبانی تغییر یافته با نتایج در (Stamatatos, 2006). در (Stamatatos, 2006) با کنار هم قرار دادن کارکتر n-گرام‌ها با مقدار 5 و 4 و n=3 عنوان ویژگی، استفاده شده است. آزمایش سه بار تکرار شده است و در هر بار تعدا ویژگی‌های استخراج شده افزایش پیدا کرده است. روش پیشنهاد شده در این کار با نام PM در جدول مشخص شده است.

روش	ویژگی	دقت
PM	۳-گرام+۴ گرام+۵ گرام=۸۱۷۸ ویژگی	74.04
IG ¹	۳-گرام+۴ گرام+۵ گرام=۸۱۷۸ ویژگی	۷۲,۵۶
PM	۳-گرام+۴ گرام+۵ گرام=۸۱۷۸ ویژگی	۷۲,۴۸
IG	۳-گرام+۴ گرام+۵ گرام=۴۶۹۱ ویژگی	72.16
PM	۳-گرام+۴ گرام+۵ گرام=۲۳۱۴ ویژگی	72
مدل سازی زبانی تغییر یافته		۷۰,۴۸
IG	۳-گرام+۴ گرام+۵ گرام=۲۳۱۴ ویژگی	۶۹,۴

¹ Informaion Gain

۲-۵ فعالیت‌ها

- ✓ گرد آوری دو مجموعه متن با شش و هفت شاعر از شاعران نامدار به زبان فارسی با نام پایگاه داده WMPR-AA2016-A و پایگاه داده WMPR-AA2016-B (فصل ۴)
- ✓ دسته بندی و بررسی در روش‌های انجام شده در حل مسئله تخصیص نویسنده (فصل ۲)
- ✓ پیشنهاد یک روش برای حل مسئله تخصیص نویسنده (فصل ۳)
- ✓ اعمال روش مدل سازی زبانی در سطح کلمه، ۲-گرام کلمه، ۳-گرام کلمه و در سطح کاراکتر با تغییر بازه متغیر n ، بر روی چهار پایگاه داده با سبک نوشتاری، زبان و تعداد نویسنده متفاوت. (فصل ۶)
- ✓ بررسی تاثیر خصوصیات پایگاه داده مانند، مقدار داده آموزشی و تعداد نویسنده کاندید و توزیع داده آموزشی در بین نویسنده‌گان کاندید.

۳-۵ پیشنهادات

- در این قسمت پیشنهادات برای کارهای آینده به صورت زیر ارائه می‌گردد.
- ایجاد پایگاه داده‌های استاندارد فارسی
- ایجاد پایگاه داده‌های استاندارد در زبان فارسی با تنوع در تعداد نویسنده کاندید، سبک نوشتاری و داده‌های آموزشی و آزمایشی با طول متفاوت باعث ارزیابی دقیق تر و نزدیک تر به مسائل واقعی خواهد شد.
- تمرکز در حل مساله نویسنده در حالت مجموعه باز و شناسایی ویژگی‌های فردی همان‌طور که در فصل اول توضیح داده شد، شناسایی نویسنده علاوه بر حالت مجموعه بسته در دو دسته مجموعه باز و شناسایی ویژگی‌های فردی نویسنده نیز انجام می‌شود. از این رو گسترش حل مسئله شناسایی نویسنده در این دو دسته در پایگاه داده با زبان فارسی می‌تواند باعث توسعه رهیافت-های تشخیص نویسنده در زبان فارسی شود. همچنین توسعه روش‌های شناسایی نویسنده با به

کارگیری روش‌های احتمالاتی و بهبود نتایج در این زمینه با بسط احتمال پیشین به نظر می‌رسد می‌تواند باعث بهبود نتایج گردد.

- تولید نرم افزار

در فصل اول کاربردهای متفاوت شناسایی نویسنده ذکر شد. توسعه تولید نرم افزارها در این زمینه در پایگاه داده‌های فارسی می‌تواند به کاربردی کردن رهیافت‌ها در زبان فارسی کمک نماید.

- Abbasi, A. & Chen, H., 2005. Applying authorship analysis to extremist-group Web forum messages. *Intelligent Systems, IEEE*, 20(5), pp. 67-75.
- Abbasi, A. & Chen, H., 2008. Writeprints: A stylometric approach to identity-level identification and similarity detection in cyberspace. *ACM Transactions on Information Systems (TOIS)*, 26(2), p 7.
- Abdallah, E. E. et al., 2013. Simplified features for email authorship identification. *International Journal of Security and Networks*, 8(2), pp. 72-81.
- Alexander Yun-chung Liu, B., 2004. *The Effect of Oversampling and Undersampling on Classifying Imbalanced Text Datasets*. Degree of Master of Science in Engineering ,The University of Texas at Austin.
- Altheneyan, A. S. a. M. M. E. B., 2014. Naive Bayes classifiers for authorship attribution of Arabic texts. *Journal of King Saud University-Computer and Information Sciences*, 26(4), pp. 473-484.
- Alzahrani, et al., 2012. Understanding plagiarism linguistic patterns, textual features, and detection methods. *Systems, Man, and Cybernetics, Part C: Applications and Reviews, IEEE Transactions on*, 42(2), pp. 133-149.
- Argamon, S., Saric, M. & Stein, a. S., 2003. Style Mining of Electronic Messages for Multiple Authorship Discrimination: First Results. *ACM SIGKDD*, Volume 19, pp. 475-480.
- Azarbonyad, H., Dehghani, M., Marx, M. & Kamps, J., 2015. *Time-aware authorship attribution for short text streams*. New York, NY, USA, Acm, pp. 727-730.
- BEKKERMAN, R. & ALLAN, J., 2004. *Using bigrams in text categorization*. Amherst, s.n., pp. 1-2.
- Boukhaled, M. A. & Ganascia, J.-G., 2015. Using Function Words for Authorship Attribution:Bag-Of-Words vs. Sequential Rules. *Natural Language Processing and Cognitive Science: Proceedings 2014*, p. 115.
- Boutwell, S. R., 2014. *Authorship attribution of short messages using multimodal features*, s.l: Master's thesis, Naval Postgraduate School.
- Bozkurt, I., Bilkent Univ., A., Baghoglu, O. & Uyar, E., 2007. Authorship attribution Performance of various features and classification methods. *Computer and information sciences*, pp. 1-5.
- Burrows & F, J., 1992. Not unless you ask nicely: The interpretative nexus between analysis and information. *Literary and Linguistic Computing*, 7(2), pp. 91-109.
- Chaski, C. E., 2001. Empirical evaluations of language-based author identification techniques. *Forensic Linguistics*, Volume 8, pp. 1-65.
- D.Manning, C., Raghavan, P. & Schütze, H., 2009. Language models for information retrieval. In: *Introduction to Information Retrieval*. s.l.:Cambridge University Press, p. 237.
- Diederich, J., Kindermann, J., Leopold, E. & Paass, G., 2003. Authorship Attribution with Support Vector. *Applied Intelligence* , 19(1), pp. 109-123.
- Diederich, J., Kindermann, J., Leopold, E. & Paass, G., 2003. machines, Authorship attribution with support vector. *Applied intelligence*, 19(1-2), pp. 109-123.
- Escalante, H. J. T. S. M. M.-y.-G., 2011. *Local histograms of character n-grams for authorship attribution*. s.l., s.n.

- Frantzeskou, G., Stamatatos, E., Gritzalis, S. & Katsikas, S., 2006. *Effective identification of source code authors using byte-level information*. ACM New York, s.n., pp. 893-896.
- Gamon, M., 2004. *Linguistic correlates of style: authorship classification with deep linguistic analysis features*. Stroudsburg, s.n., p. 611.
- Gamon, M., 2004. *Linguistic correlates of style: authorship classification with deep linguistic analysis features*. s.l., s.n.
- Grieve, J., 2007. Quantitative authorship attribution: An evaluation of techniques. *Literary and linguistic computing*, 22(3).
- Holmes, D. I., 1998. The evolution of stylometry in humanities scholarship. *Literary and linguistic computing*, 13(3), pp. 111-117.
- Howedi, F. & Mohd, M., 2014. Text Classification for Authorship Attribution Using Naive Bayes Classifier with Limited Training Data. *Computer Engineering and Intelligent Systems*, 5(4), pp. 48-58.
- Jamak, A., Savatić, A. & Can, M., 2012. Principal component analysis for authorship attribution. *Business Systems Research*, 3(2), pp. 49-56.
- Joula, P., 2008. Authorship Attribution. In: Boston: the essence of knowlege, pp. 238-239.
- Juola, P., Sofko, J. & Brennan, P., 2006. A Prototype for Authorship Attribution Studies. *Literary and Linguistic Computing*, 21(2), pp. 169-178.
- Jurafsky, D. & Martin, J. H., 2006. N-Grams. In: *An Introduction to Natural Language Processing, Computational Linguistics, and Speech Recognition*. s.l.:s.n.
- Keselj, V., Peng, F., Cercone, N. & Thomas, C., 2003. *N-gram-based author profiles for authorship attribution*. PACLING, s.n.
- KJELL, B., 1994. Authorship determination using letter pair frequency features with neural network classifiers. *Literary and Linguistic Computing*, 9(2), pp. 119-124.
- Koehn, P., 2010. Language Modeles. In: *Statistical Machine Translation*. s.l.:Cambridge University Press, p. 181.
- Koppel, M. & Schler, J., 2003. *Exploiting Stylistic Idiosyncrasies for Authorship Attribution*. s.l., s.n., p. 72.
- Koppe, M., Schler, J. & Argamon, S., 2009. Computational methods in authorship attribution. *Journal of the American Society for information Science and Technology*, 60(1), pp. 9-26.
- Kourtis, Stamatatos, I. a. & Efstathios, 2011. *Author identification using semi-supervised learning*. Amsterdam, The Netherlands, s.n.
- Lewis, D. D., Yang, Y., Rose, T. G. & Li, F., 2004. Rcv1: A new benchmark collection for text categorization research. *Journal of machine learning research*, Volume 5, pp. 361-97.
- Luyckx, K. & Daelemans, W., 2005. Shallow Text Analysis and Machine Learning for Authorship Attribution. *LOT Occasional Series*, Volume 4, pp. 149--160.
- Luyckx, K. & Daelemans, W., 2011. The effect of author set size and data size in authorship attribution. *Literary and linguistic Computing*, 26(1), pp. 35-55.
- Manning, C. D. & Schiitze, H., 2000. Words. In: *Foundations of Statistical Natural Language Processing*. London, England: The MIT Press, p. 191.
- Marton, Y., Wu, N. & Hellerstein, L., 2005. *On compression-based text classification*. s.l., Springer, pp. 300--314.
- Mikros, G. K. & Perifanos, K., 2013. *Authorship Attribution in Greek Tweets Using Author's Multilevel N-Gram Profiles*. s.l., s.n.

- Mosteller, Frederick, Wallace & David, 1964. *Inference and disputed authorship: The Federalist*, s.l.: s.n.
- Nagaprasad, S. et al., 2015. Empirical Evaluations Using Character and Word N-Grams on Authorship Attribution for Telugu Text. In: *Intelligent Computing and Applications*. s.l.:Springer, pp. 613--623.
- Pavelec, D., Justino, E. & Oliveira, L. S., 2007. Author Identification using Stylometric Features. *Inteligencia Artificial*, 11(36), pp. 59-66.
- Peng, F. & Schuurmans, D., 2003. *Combining naive Bayes and n-gram language models for text classification*. s.l., Springer.
- Pillay, S. R. & Solorio, T., 2010. *Authorship attribution of web forum posts*. Dallas, TX, IEEE.
- Ramezani, R., Sheydaei, N. & Kahani, M., 2013. *Evaluating the effects of textual features on authorship attribution accuracy*. s.l., IEEE.
- Ramyaa, He, C. & Rasheed, K., 2004. *Using machine learning techniques for stylometry*. s.l., s.n., pp. 897--903.
- Rappoport, R. S. O. T. A. & Koppel, M., 2013. Authorship attribution of micro-messages. In: s.l.:s.n.
- Satyam, A., Dawn, A. K. & Saha, S. K., 2014. *A Statistical Analysis Approach to Author Identification Using Latent Semantic Analysis*. s.l., s.n.
- Savoy, J., 2012. Authorship attribution based on specific vocabulary. *ACM Transactions on Information Systems*, 30(2), p. 12.
- Selman, S., Turan, K. & Kuşakçı, A. O., 2012. Distinction of The Authors of Texts Using Multilayered Feedforward Neural Networks. *SouthEast Europe Journal of Soft Computing*, 1(1), pp. 128-138.
- Sidorov, G. a. V. F. a. S., Gelbukh, E. a., Alexander & Chanona-Hern, 2013. Syntactic dependency-based n-grams as classification features. In: *Advances in Computational Intelligence*. San Luis Potosi, Mexico: Springer, pp. 1-11.
- Silva, R. S. et al., 2011. *twazn me!!! ;(' Automatic Authorship Analysis of Micro-Blogging Messages*. Alicante, Spain, s.n., pp. 161-168.
- Soboroff, I. M., Nicholas, C. K., Kukla, J. M. & Ebert, D. S., 1997. *Visualizing document authorship using n-grams and latent semantic indexing*. New York, NY, USA, s.n.
- Stamatatos, E., 2006. Authorship attribution based on feature set subsampling ensembles. *International Journal on Artificial Intelligence Tools*, 15(05), pp. 823-838.
- Stamatatos, E., 2007. *Author Identification Using Imbalanced and Limited Training Texts*. Regensburg, s.n., pp. 237 - 241.
- Stamatatos, E., 2008. Author identification: Using text sampling to handle the class imbalance problem. *Information Processing & Management*, 44(2), p. 790--799.
- Stamatatos, E., 2008. Author identification: Using text sampling to handle the class imbalance problem. *Information Processing & Management*, 44(2), pp. 790--799.
- Stamatatos, E., 2009. A survey of modern authorship attribution methods. *Journal of the American Society for information Science and Technology*, 60(3), pp. 538-556.
- Stamatatos, E., Fakotakis, N. & Kokkinakis, G., 2000. Automatic Text Categorization in Terms of Genre and Author. *Computational linguistics*, 26(4), pp. 471-495.
- Stamatatos, J. H. a. E., 2006. *N-gram feature selection for authorship identification*. s.l., Springer, pp. 77--86.

Tearle, M., Taylor, K. & Demuth, H., 2008. An algorithm for automated authorship attribution using neural networks. *Literary and linguistic computing*, 23(4), pp. 425-442.

Türkoğlu, F., Diri, B. & Amasyalı, M. F., 2007. *Author Attribution of Turkish Texts by Feature Mining*. Qingdao, China, s.n., pp. 1086-1093.

Tweedie, F. J., Singh, S. & Holmes, D. I., 1996. Neural network applications in stylometry: The Federalist Papers. *Computers and the Humanities*, 30(1), pp. 1-10.

Uzuner, Ö. & Katz, B., 2005. *A comparative study of language models for book and author recognition*. Jeju Island, Korea, Springer, pp. 969-980.

Vel, O. d., Anderson, A., Corney, M. & Mohay, G., 2001. Mining e-mail content for author identification forensics. *ACM Sigmod Record*, 30(4), pp. 55-64.

Zhao, Y. & Zobel, J., 2005. *Effective and Scalable Authorship Attribution Using Function Words*. Jeju Island, Korea, Springer, pp. 174-189.

Zhao, Y. & Zobel, J., 2007. *earching with style: Authorship attribution in classic literature*. s.l., s.n., pp. 59--68.

Zhao, Y., Zobel, J. & Vines, P., 2006. *Using relative entropy for authorship attribution*. s.l., Springer, pp. 92-105.

Zheng, R., Li, J., Chen, H. & Huang, Z., 2006. A framework for authorship identification of online messages: Writing-style features and classification techniques. *Journal of the American Society for Information Science and Technology*, 57(3), pp. 378-393.

آذین، ز. 1392، پایان نامه ارشد، شناسایی خودکار شاعران شعر نو با استفاده از ویژگی های زبانی، مهندسی کامپیوتر، دانشگاه صنعتی شریف.

امیرشهاب، ش. ۱۳۸۵. پایان نامه ارشد، ویژگی های نحوی، کارایی بیشتری را در بین سایر ویژگی ها دارد، مهندسی کامپیوتر و فناوری اطلاعات

فرهمندپور، زینب. 1390، پایان پایان نامه ارشد، طراحی و پیاده سازی یک سیستم هوشمند تشخیص هویت بر اساس سبک نوشتاری فارسی، مهندسی کامپیوتر، دانشگاه بوعلی سینا همدان.

Abstract

Authorship attribution (AA) or author identification refers to the problem of determining who has written a disputed text or unseen text. In the close class authorship attribution problem, the unseen text is assigned to any one of candidate authors set, that text sample as training data are available for them. Two main requirements of authorship attribution system are features and attribute method. Features are usually selected with training data.

With increasing text in different languages, seems to be an essential need for developing authorship attribution system which is language independent. Since the procedure of extracting character n-grams and word n-grams are language-independent and require no special tools, in this thesis they have been used as features. Also language modeling has been chosen as a statistical and probabilistic attribute method. We present an approach based on language modeling called modified language modeling. It aims to offer a solution for AA problem by combinations of both bigram words weighting and unigram words weighting. Moreover, the IDF value multiplied by related word probability has been used, instead of removing stop words and balancing word probability as weights, as well.

In order to evaluate the results, four corpora have been used. Two datasets of Persian poetry, one Persian prose dataset and the fourth is English prose dataset. The accuracy of AA is calculated by language modeling with character n-grams, language modeling with word n-grams and modified language modeling on the four datasets. In all databases, modified language modeling shows improvement. The best performance is obtained for modified language by Persian prose dataset, which is 100 percent.

Keywords: Authorship Attribution, Authorship Identification, Language Modeling, WMPR-AA2016-A corpora, WMPR-AA2016-B corpora



Department of Electronic Learning

M.Sc. Thesis in Artificial Intelligence Engineering

Authorship attribution with statistical modeling on text data

By: Samane Vazirian

Supervisor:

Dr. Morteza Zahedi

Advisor:

Dr. Hamid Hasanpour

September 2016