

بازشناسی متون فارسی با استفاده از مدل زبانی n-gram و پالایش گرامری

پریسا شیروانی^۱، مهرداد وطن خواه خوزانی^۲ و خشایار یغمایی^۳
^۱ دانشکده مهندسی برق و کامپیوتر، دانشگاه سمنان، سمنان، ایران
^۲ دانشکده مهندسی کامپیوتر، دانشگاه شفیلد هالام، شفیلد، انگلستان

چکیده

بازشناسی متون، در سال‌های اخیر بسیار مورد توجه قرار گرفته است. ارائه الگوریتم‌های بازشناسی، برگرفته از ساختار گرامری و معنایی این زبان می‌تواند روش مؤثری در پردازش‌های دیگر مربوط به خط و زبان فارسی باشد. در این مقاله با استفاده از شاخه علمی پردازش زبان‌های طبیعی، یک الگوریتم سه مرحله‌ای به منظور بازشناسی متون فارسی بر مبنای بازشناسی جملات فارسی ارائه می‌شود. این روش شامل مراحل ترکیب زیرکلمات به منظور ساخت کلمات و سپس جملات بالقوه معنی‌دار و در نهایت استفاده از دو مدل زبانی و چند قاعده گرامری، به منظور تشخیص جمله صحیح بر اساس انطباق با گرامر زبان فارسی است. آزمایش‌های متعدد نشان می‌دهد که دقت روش ارائه شده برای مرحله ساخت کلمات و سپس جملات بالقوه معنی‌دار ۹۸ درصد و برای تشخیص جمله صحیح با استفاده از مدل زبانی باگرام ۸۵ درصد و برای مدل زبانی تراگرام ۸۸ درصد است.

واژگان کلیدی: بازشناسی متن، فارسی، مدل‌سازی زبان فارسی، پردازش زبان‌های طبیعی.

۱- مقدمه

بازشناسی متن^۱ که یکی از محوری‌ترین شاخه‌های بازشناسی الگو^۲ است، با پیشرفت روزافزون رایانه، تبدیل به یکی از موضوعات مهم شد. در واقع به دلیل ویژگی‌هایی از قبیل انجام سریع‌تر و دقیق‌تر کارها و خستگی‌ناپذیر بودن رایانه‌ها، انجام امور به وسیله آنها ترجیح داده می‌شود. بنابراین می‌توان یکی از مهم‌ترین و اصلی‌ترین دلایل گسترش بی‌شمار تحقیقات در حوزه بازشناسی متون را همین تمایل به انجام کارها با رایانه دانست.

تحقیقات در زمینه بازشناسی متون چایی فارسی و عربی به طور تقریبی از اوایل دهه ۱۹۸۰ آغاز شد (Amin et al., 1980, Badie et al., 1980).

وجود برخی ویژگی‌های خاص در نگارش این متون مانند وجود نقاط، علائم، تنوع قلم‌ها و هم‌پوشانی حروف با یکدیگر، بازشناسی متون فارسی و عربی را با پیچیدگی‌ها و دشواری‌هایی همراه کرده است.

تاکنون کارهای انجام شده در زمینه بازشناسی متون بر بازشناسی حروف و کلمات متمرکز بوده‌اند (Al-Muhtaseb et al., 2008, El-Abed et al., 2008, Jacobs et al., 2005, Sarfaraz et al., 2003, Khosravi et al., 2008, Ebrahimi et al., 2008, رضوی و همکاران ۱۳۸۴ و ۱۳۸۳) (Ebrahimi et al., 2011). Bahmani 2010). این بازشناسی‌ها به طور معمول با سه رویکرد کلی صورت می‌گیرند (ابراهیمی، ۱۳۸۴): در رویکرد اول جداسازی و بازشناسی دو مرحله مستقل از یکدیگر هستند. در رویکرد دوم کلمه به عنوان یک الگوی واحد در نظر گرفته و بازشناسی کلمه بدون جداسازی حروفش انجام می‌شود. در رویکرد سوم یا رویکرد ترکیبی اطلاعات مربوط به جداسازی حروف تشکیل‌دهنده کلمه و اطلاعات

¹ Text Recognition

² Pattern Recognition

مربوط به شکل کلی کلمه در قالب یک سیستم ترکیبی جاسازی-بازشناسی به کار گرفته می‌شود. در این مقاله از مدل‌سازی آماری زبان^۱ در حوزه پردازش زبان‌های طبیعی^۲ و همچنین چند قاعده گرامری به منظور بازشناسی جملات فارسی استفاده شده است.

در دهه‌های اخیر با توجه به کاربردهای گسترده مدل‌های زبان، تحقیقات زیادی برای مدل‌سازی زبان‌های پرکاربرد جهانی و به خصوص زبان انگلیسی انجام شده است. مدل‌های آماری زبان انگلیسی از سال ۱۹۸۰ در سیستم‌های واقعی به کار رفته‌اند (Rosenfeld, 2000). اما برای مدل‌سازی زبان فارسی هنوز تحقیقات گسترده و قابل اعتنایی صورت نگرفته است. پیش‌بینی واژگان فارسی، یکی از تحقیقاتی است که با استفاده از مدل‌سازی آماری زبان فارسی انجام شده است (Ghayoomi et al., 2004). در این تحقیق که با وارد شدن حروف اولیه یک واژه، سیستم واژه‌هایی را که با آن حرف آغاز می‌شود در پنجره‌ای فهرست می‌کند و به کاربر پیشنهاد می‌دهد. با وارد شدن حروف بیشتر، پیشنهادها محدودتر می‌شود تا واژه مورد نظر کاربر از میان آنها یافت شود.

بررسی یک مدل آماری زبان بر اساس دسته‌های منطقی دستوری زبان فارسی برای استفاده از بازشناسی گفتار پیوسته، نوعی مدل‌سازی جدید زبان فارسی است. این مدل زبان بر اساس دسته‌های منطقی N-Gram با طول متغیر کار می‌کند و در آن به جای یافتن الگوهای آماری مربوط به دنباله‌های کلمات، روابط بین دسته‌های منطقی از کلمات، مورد بررسی قرار می‌گیرند. این تحقیق در سه مرحله خوشه‌بندی کلمات، به دست آوردن مدل آماری زبان با استفاده از ساختار درختی و اعمال مدل زبان به سیستم بازشناسی انجام شده است.

تقطیع و برچسب‌دهی نحوی - معنایی داده‌های نوشتاری یکی از فعالیت‌های اصلی در طراحی و ساخت هر دادگان زبانی برای استخراج مدل زبانی است (بی‌جن‌خان محمود، ۱۳۸۱). برای استخراج مدل زبان فارسی یک بسته نرم‌افزاری نوشته شده، که در چارچوب فرآیند مارکوف صفر تا سه مرحله‌ای، توزیع احتمال مشروط کلمات فارسی را در چهار حالت به دست می‌دهد.

با استفاده از یک مدل آماری زبان می‌توان احتمال کلمات بعدی را پیش‌بینی کرد. همچنین مدل‌های زبانی دیگری از جمله مدل‌های توانی^۳ (Rosenfeld et al., 2001) و گرامرهای مستقل از متن^۴ (Corazza et al., 1993) برای مدل‌سازی زبان‌های طبیعی پیشنهاد شده‌اند؛ اما در عمل مدل‌های آماری N-gram (Rosenfeld, 2000) به‌علت سادگی پیاده‌سازی ترجیح داده می‌شوند.

در الگوریتم پیشنهادی از مدل‌سازی زبان فارسی در جهت پالایش گرامری جملات آن استفاده شده است. برای این منظور ابتدا زیرکلمات از متن مورد پردازش استخراج شده و از ترکیب آنها کلمات و سپس جملات بالقوه معنی‌دار ساخته می‌شوند. در نهایت از مدل زبانی 2-gram یا بایگرام (مدل پنهان مارکوف مرتبه اول) و 3-gram یا تراگرام (مدل پنهان مارکوف مرتبه دوم) و همچنین حذف برخی ترکیبات دستوری غیرممکن به منظور تشخیص جمله صحیح از میان مجموعه‌ای از جملات بالقوه معنی‌دار استفاده شده است. در بخش دوم به ساخت کلمات بامعنی و جملات بالقوه معنی‌دار از ترکیب زیرکلمات پرداخته می‌شود. بخش سوم به تشخیص جمله صحیح با استفاده از مدل‌های زبانی و همچنین حذف ترکیبات دستوری غیرممکن می‌پردازد. در بخش چهارم نتایج ارائه شده و بخش پنجم به بحث و تحلیل پیرامون آن‌ها اختصاص می‌یابد. در نهایت نتیجه‌گیری در بخش ششم ارائه خواهد شد.

۲- ساخت کلمات بامعنی و جملات بالقوه

معنی‌دار

کلمات را در زبان فارسی می‌توان ترکیبی از زیرکلمات دانست. هر بخش از دنباله حروف که از بخش‌های قبل و بعد بتواند جدا در نظر گرفته شود، یک زیرکلمه است. زیرکلمه از یک حرف یا ترکیبی از حروف به هم پیوسته تشکیل شده است. به‌عنوان مثال کلمه "مدرسه" دارای سه زیرکلمه "مد"، "ر"، "سه" است. درحالی که کلمه "شبنم" تنها شامل یک زیر کلمه چهارحرفی "شبنم" است. از این رو در این تحقیق سعی بر این بوده است که کلمات و به دنبال آن جملات از ترکیب زیرکلمات ایجاد و شناسایی شوند. برای این منظور ابتدا بانک اطلاعاتی از کلمات پرسامد در زبان فارسی تهیه شد (حسنی، ۱۳۸۴).

¹ Statistical Language Models

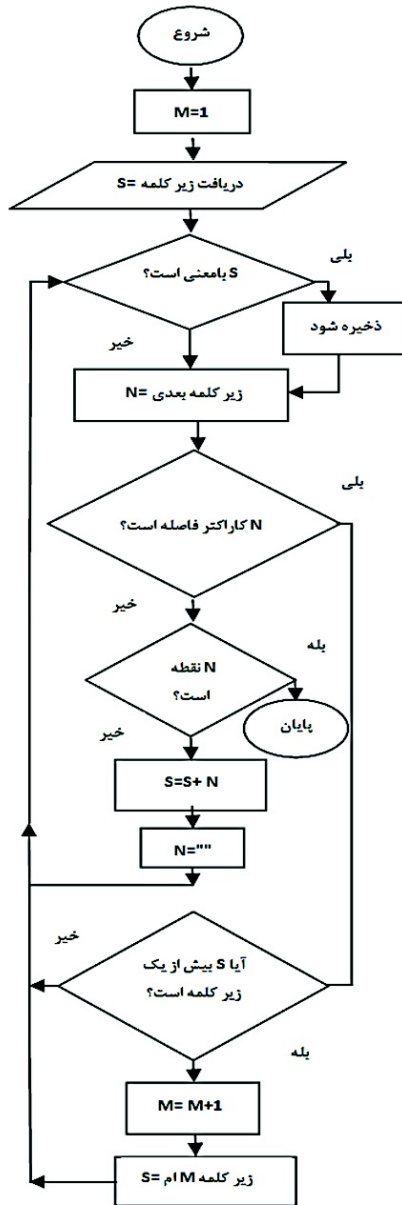
² Natural Language Processing

³ Exponential Models

⁴ Context Free Grammars

رفت" از زیر کلمات استخراج می‌شوند. سپس این کلمات در کنار هم قرار می‌گیرند و سه جمله بالقوه معنی‌دار زیر را می‌سازند:

- پزشک با عجله به داروخانه رفت.
- پزشک با عجله به دارو خانه رفت.
- پزشک با عجله به دار و خانه رفت.



شکل-۱: روندنمای ساخت کلمات بامعنی

هرچند تعداد این جملات بالقوه معنی‌دار در بسیاری از موارد محدود است؛ اما برای شناسایی درست جملات

با ترکیب زیر کلمات، کلمات معنی‌دار ممکن به‌دست می‌آیند. به‌عنوان مثال جمله "پدر من آمد" از زیر کلمات "پد"، "ر"، "من"، "ا"، "مد" تشکیل شده است. از ترکیب ترتیبی (از راست به چپ) این زیر کلمات، تنها می‌توان کلمات معنی‌دار "پدر"، "من"، "آمد" را استخراج کرد و به همین دلیل ترکیب جمله به‌راحتی قابل شناسایی است.

این در حالی است که برخی از جملات فارسی را به‌دلیل امکان ترکیب زیر کلمات در ساخت کلمات با معنی متفاوت، می‌توان به صورت‌های مختلف خواند. اگرچه تمامی این ترکیب‌ها دارای معنی نیستند (این همان اشتباهی است که در کودکان و افراد نوسواد در خواندن متن‌ها پیش می‌آید). به‌طور مثال، زیر کلمات جمله "مادر بزرگ آمد"، عبارتند از "ما"، "د"، "ر"، "ب"، "ز"، "ر"، "ا"، "گ"، "ا"، "ا"، "مد". اکنون اگر از ترکیب زیر کلمات (از راست به چپ) استفاده کنیم، چهار جمله بالقوه معنی‌دار زیر قابل استخراج می‌باشند:

- ما در بز رگ آمد.
- ما در بزرگ آمد.
- مادر بز رگ آمد.
- مادربزرگ آمد.

روندنمای قسمت ساخت کلمات بامعنی از ترکیب زیر کلمات در شکل (۱) نشان داده شده است. همان‌طور که در این شکل مشاهده می‌کنید، ابتدا زیر کلمه اول دریافت می‌شود. سپس اگر این زیر کلمه به تنهایی معنی‌دار باشد، ذخیره و زیر کلمه بعدی از ورودی دریافت می‌شود. همچنین اگر زیر کلمه دریافتی اول بامعنی نباشد، زیر کلمه بعدی از ورودی دریافت می‌شود. اگر زیر کلمه دریافت‌شده، نقطه باشد به معنای رسیدن به انتهای جمله و پایان زیر کلمات جمله است. اگر نویسه فاصله باشد به معنای پایان زیر کلمات مربوط به یک کلمه است؛ زیرا بین زیر کلمات در یک کلمه فارسی حداکثر نیم‌فاصله وجود دارد و فاصله کامل بین کلمات ظاهر می‌شود. از این‌رو زیر کلمه دریافتی دوم در صورتی که نقطه یا نویسه فاصله نباشد، در کنار زیر کلمه اول قرار می‌گیرد و معنی‌دار بودن ترکیب ساخته شده بررسی می‌شود. در این روندنما M به‌عنوان یک اندیس برای انتخاب زیر کلمات در نظر گرفته شده است. به‌عنوان مثال دیگر زیر کلمات "پز"، "شک"، "با"، "عجله"، "به"، "د"، "ا"، "ر"، "و"، "خا"، "نه"، "ر"، "فت" را در نظر بگیرید. کلمات "پزشک، با، عجله، به، داروخانه، دارو خانه، دار و خانه،

ضروری است که با استفاده از قواعد دستوری و معنایی فارسی، جمله یا جملات معنی‌دار از این ترکیبات بالقوه استخراج شوند.

۳- تشخیص جمله صحیح

به‌منظور تشخیص جمله صحیح از میان مجموعه جملات بالقوه معنی‌دار ساخته شده، از دو مدل زبانی بایگرام و تریگرام و همچنین حذف ترکیبات دستوری غیرممکن استفاده کردیم.

۳-۱- اعمال مدل‌های زبانی

در این مرحله از دو مدل زبانی بایگرام و تریگرام (Srihar et al., 2007) که برای پالایش گرامری زبان فارسی تغییر یافته‌اند، استفاده می‌کنیم. دادگان آموزشی مورد استفاده در این تحقیق، متون فارسی استاندارد موجود در نشریات و مجلات هستند. درحقیقت، این داده‌ها از طیف وسیعی از متون فارسی در زمینه‌ها و موضوعات مختلف جمع‌آوری شده‌اند. برای تهیه داده‌های آموزشی با حجم زیاد تنها راه عملی استفاده از متون موجود در شبکه اینترنت است. این آموزش در مدل زبانی بایگرام به‌صورت دوتایی و در مدل زبانی تریگرام به‌صورت سه‌تایی انجام می‌شود.

به این ترتیب که در مدل زبانی اول (بایگرام) متون آموزشی به‌صورت دو کلمه دو کلمه، بررسی می‌شوند و تعداد تکرار برای هر دو نوع متوالی کلمه در متون آموزش داده شده بر حسب درصد، در یک جدول به نام "فراوانی توالی دوتایی مقولات واژگانی" ذخیره می‌شود؛ درحالی‌که در مدل زبانی دوم (تریگرام) متون آموزشی به‌صورت سه‌کلمه سه‌کلمه، مورد بررسی قرار می‌گیرند. در این حالت تعداد تکرار برای هر سه نوع متوالی کلمه در جدول "فراوانی توالی سه‌تایی مقولات واژگانی" نشان داده می‌شود.

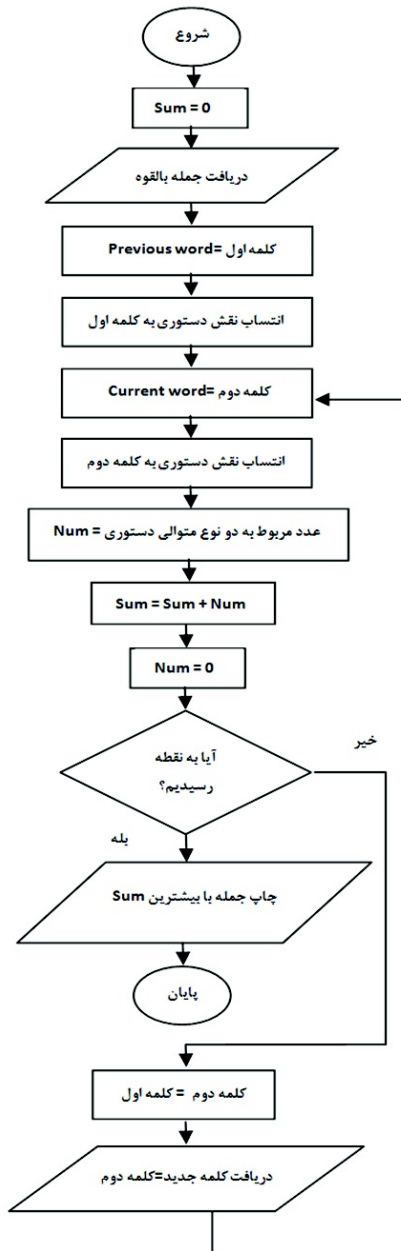
به‌عنوان نمونه جدول "فراوانی توالی دوتایی مقولات واژگانی" برای مدل زبانی بایگرام در شکل (۲) نشان داده شده است. همچنین روندنمای اعمال مدل زبانی بایگرام در شکل (۳) آورده شده است. همان‌طور که در این شکل مشاهده می‌کنید کلمات موجود در جمله به‌صورت دوتایی مورد بررسی قرار می‌گیرند به این ترتیب که ابتدا دو کلمه اول جمله در نظر گرفته شده و نقش دستوری آنها مشخص می‌شود؛ سپس عدد مربوط به ترکیب دستوری از جدول "فراوانی توالی دوتایی مقولات واژگانی" خوانده می‌شود. این

روند برای کلمات بعدی نیز تکرار می‌شود تا زمانی که به نویسه نقطه برسیم. درنهایت درصدهای خوانده شده از جدول "فراوانی توالی دوتایی مقولات واژگانی" مربوط به ترکیبات گرامری دوتایی واژگان جمله با هم جمع می‌شوند و درصد حاصل از اعمال مدل زبانی به جمله بالقوه معنی‌دار مورد بررسی، بدست می‌آید؛ سپس از بین تمام جملات بالقوه معنی‌دار که به هر کدام درصدی براساس اعمال مدل زبانی داده شده است، جمله با بیشترین مقدار درصد به‌عنوان جمله صحیح تشخیص داده خواهد شد.

به‌عنوان مثال جمله "من بادام دارم" را در نظر بگیرید. ابتدا دو کلمه اول این جمله یعنی "من" و "بادام" بررسی می‌شوند و نقش دستوری هر یک مشخص می‌شود. "من" به عنوان ضمیر و "بادام" به عنوان اسم در نظر گرفته می‌شوند یعنی دارای ترکیب "ضمیر+اسم" هستند. سپس عدد مربوط به این ترکیب دستوری از جدول "فراوانی توالی دوتایی مقولات واژگانی" خوانده می‌شود. درواقع درصدهای موجود در جداول بایگرام و تریگرام، اطلاعات آماری هستند که از متون آموزش داده شده، استخراج شده‌اند. در مرحله بعد، نقش دستوری کلمه بعدی (کلمه "دارم" در مثال مذکور) مشخص می‌شود و درصد مربوط به ترکیب گرامری این کلمه با کلمه قبل از آن با استفاده از جدول "فراوانی توالی دوتایی مقولات واژگانی" تعیین می‌گردد. در اینجا "دارم" به‌عنوان فعل و کلمه قبل از آن یعنی "بادام" به‌عنوان اسم، ترکیب دستوری "اسم+فعل" را تشکیل می‌دهند. درنهایت درصدهای مربوط به ترکیب‌های گرامری دوتایی کلمات جمله با هم جمع و عدد حاصل از اعمال مدل زبانی بایگرام به جمله مشخص می‌شود. در مثال بیان‌شده از جمع دو درصد مربوط به ترکیب‌های گرامری "ضمیر+اسم" و "اسم+فعل"، درصد مربوط به مدل زبانی بایگرام جمله تعیین می‌شود.

در مدل زبانی تریگرام، ابتدا سه کلمه اول جمله در نظر گرفته شده و نقش دستوری‌شان مشخص می‌شود؛ سپس کلمه اول حذف شده و سه کلمه بعدی بررسی می‌شوند. به‌عنوان مثال جمله "پدر به خانه آمد" را در نظر بگیرید. در ابتدا نقش دستوری سه کلمه اول یعنی کلمات "پدر"، "به"، "خانه" به صورت "اسم+حرف اضافه+اسم" تعیین می‌شود. سپس کلمه اول یعنی پدر حذف شده و نقش دستوری سه کلمه بعدی "به"، "خانه"، "آمد" به‌صورت "حرف اضافه+اسم+فعل" مشخص می‌شود. تا اینجا دو ترکیب دستوری سه‌تایی "اسم+حرف اضافه+

فعل + فعل"، "فعل + فعل + اسم"، "فعل + فعل + ضمیر"، "فعل + فعل + حرف"، "فعل + فعل + قید"، "فعل + فعل + صفت" در مدل زبانی تریگرام از نظر گرامر زبان فارسی صحیح نبوده و حذف می‌شوند. به‌عنوان مثال جمله "در جریان زندگی او خود را به سرنوشت سپرد" را در نظر بگیرید که از آن دو جمله بالقوه معنی‌دار زیر استخراج می‌شود:



(شکل-۳): روندنمای تشخیص جمله صحیح با استفاده از مدل زبانی بایگرام

اسم" و "حرف اضافه + اسم+فعل" حاصل شده است که با استفاده از جدول "فراوانی توالی سه تایی مقولات واژگانی" درصدهای مربوط به هر ترکیب مشخص شده و این مقادیر با هم جمع می‌شوند. این مجموع، درصد مربوط به اعمال مدل زبانی تریگرام به جمله مذکور را می‌دهد.

۳-۲- حذف ترکیبات دستوری غیرممکن

به‌منظور بهبود عملکرد روش پیشنهادی و افزایش دقت علاوه بر استفاده از مدل‌های زبانی بیان شده، از حذف برخی از ترکیبات غیرممکن دستوری نیز به‌منظور پالایش جملات استفاده کردیم. به این صورت که ابتدا مدل زبانی به جملات بالقوه معنی‌دار اعمال می‌شود و درصد متناسب با هر جمله با توجه به نقش دستوری کلمات آن، با استفاده از جدول گرامر رایج مشخص می‌شود.

برج کلمات گرامر اشتباه آموزش نمایش آماري دوتایی تست دوتایی نمایش آماري سه تایی تست سه تایی

اطلاعات آماري ثبت شده

کلمه اول	کلمه دوم	درصد
اسم	اسم	۱۰
اسم	فعل	۷۰
اسم	صفت	۲۰
اسم	قید	۳۵
اسم	حرف اضافه	۵۳
اسم	حرف ربط	۱۰
حرف اضافه	اسم	۵۱
حرف اضافه	قید	۱۵
حرف اضافه	حرف اضافه	۴
حرف ربط	اسم	۴
حرف ربط	فعل	۳

(شکل-۲): فراوانی توالی دوتایی مقولات واژگانی

سیس جملاتی که دارای ترکیبات غیرممکن هستند، حذف می‌شوند و درصدهای مربوط به جملات باقی‌مانده با هم مقایسه شده و جمله با بیشترین درصد به‌عنوان جمله صحیح، تشخیص داده می‌شود. ترکیبات غیرممکن دستوری به‌کار برده شده در الگوریتم پیشنهادی، به‌منظور پالایش جملات عبارتند از:

ترکیب غیرممکن اول: در گرامر فارسی دو فعل پشت سر هم قرار نمی‌گیرند. بنابراین جملات دارای ترکیب "فعل + اسم + فعل + فعل"، "صفت + فعل + فعل"، "فعل + فعل + فعل" و "ضمیر + فعل + فعل"، "قید + فعل + فعل"، "حرف + فعل

البته از نظر صرفی کلمه مشترک بین صفت و قید و همچنین بین ضمیر و حرف ربط وجود دارد که در برخی از موقعیت‌ها پس از آنها در جمله نشانهٔ مفعول نیز می‌تواند قرار گیرد. به‌عنوان مثال کلمهٔ "خوب" بین صفت و قید و کلمهٔ "که" بین ضمیر و حرف ربط مشترک است.

۴- نتایج تجربی

در این قسمت نتایج حاصل از اعمال الگوریتم مربوط به بخش چهارم (یعنی تشخیص جملهٔ صحیح) بر روی تعدادی از جملات که از یک متن بیست صفحه‌ای استخراج شده‌اند، ارائه می‌شود. این نتایج در جداول (۱ و ۲) نشان داده شده‌اند. در این جداول تعداد چهارصد و بیست جمله مورد بررسی قرار گرفته‌اند. که از این تعداد، چهل و یک جمله به دو صورت، بالقوه معنی‌دار و یک جمله به سه صورت، بالقوه معنی‌دار هستند. همان‌طور که انتظار می‌رود در متون فارسی تعداد جملاتی که دارای حالات بالقوه معنی‌دار هستند کم می‌باشد. در جدول (۲) نتایج حاصل از اعمال مدل‌های زبانی بایگرام و تراگرام به جملات انتخابی آورده شده است. همان‌طور که در جدول مشخص است، تعداد جملات اشتباه تشخیص داده شده توسط مدل‌های زبانی اعمال شده بسیار کم است. مدل تراگرام عملکرد بهتری نسبت به مدل زبانی بایگرام داشته است؛ زیرا تعداد جملات اشتباه تشخیص داده شده توسط مدل تراگرام کمتر از مدل بایگرام است. با توجه به اعمال مدل‌های زبانی بر جملات مختلف و جداول (۱) و (۲) کارایی مدل زبانی بایگرام ۸۵ درصد و کارایی مدل زبانی تراگرام ۸۸ درصد محاسبه شده است. به‌عنوان مثال نمونه‌ای از این جملات در زیر بررسی شده است:

جملهٔ ورودی: با مداد نرم سرعت نوشتن خود را زیاد کنید.
از جملهٔ فوق، چهار جملهٔ بالقوه معنی‌دار زیر را می‌توان استخراج کرد:

- با مداد نرم سرعت نوشتن خود را زیاد کنید.
- بامداد نرم سرعت نوشتن خود را زیاد کنید.
- با مداد نرم سرعت نوشتن خود را زیاد کنید.
- بامداد نرم سرعت نوشتن خود را زیاد کنید.

ابتدا نقش دستوری کلمات هر جمله مشخص می‌شود و درصد متناسب با ترکیب‌های گرامری دوتایی آن با استفاده از جدول "فراوانی توالی مقولات واژگانی" تعیین می‌شود. در مرحلهٔ بعد به‌دلیل این که هیچ‌کدام از جملات بالا دارای

در جریان زندگی او خود را به سرنوشت سپرد.

در جریان زندگی او خود را به سر نوشت سپرد.

در جملهٔ دوم به‌علت این که دو فعل "نوشت" و "سپرد" پشت سر هم قرار گرفته‌اند، یعنی جملهٔ دارای ترکیب "فعل + فعل" در مدل زبانی بایگرام و ترکیب "اسم + فعل + فعل" در مدل زبانی تراگرام است از نظر قواعد گرامری زبان فارسی صحیح نبوده و حذف می‌شود. از این رو جملهٔ اول از میان دو جملهٔ بالقوه معنی‌دار مذکور، به‌عنوان جملهٔ صحیح تشخیص داده می‌شود.

ترکیب غیرممکن دوم: حروف اضافه قبل از فعل قرار نمی‌گیرند؛ یعنی ترکیب "حرف اضافه + فعل" در جملات فارسی وجود ندارد. از این رو جملات دارای ترکیب "حرف اضافه + فعل" در مدل زبانی بایگرام و جملات دارای ترکیبات "اسم + حرف اضافه + فعل"، "صفت + حرف اضافه + فعل"، "فعل + حرف اضافه + فعل"، "ضمیر + حرف اضافه + فعل"، "قید + حرف اضافه + فعل"، "حرف اضافه + حرف اضافه + فعل"، "حرف ربط + حرف اضافه + فعل"، "حرف اضافه + فعل + اسم"، "حرف اضافه + فعل + صفت"، "حرف اضافه + فعل + ضمیر"، "حرف اضافه + فعل + قید"، "حرف اضافه + فعل + حرف اضافه"، "حرف اضافه + فعل + حرف" در مدل زبانی تراگرام صحیح نبوده و حذف می‌شوند. به‌عنوان مثال از جمله "اهداف جدیدی برای زندگی خود برگزیدم" دو جملهٔ بالقوه معنی‌دار زیر استخراج می‌شود:

اهداف جدیدی برای زندگی خود برگزیدم.

اهداف جدیدی برای زندگی خود برگزیدم.

جملهٔ دوم به‌علت این که حرف اضافه "بر" قبل از فعل "گزیدم" قرار گرفته است یعنی دارای ترکیب "حرف اضافه + فعل" در مدل زبانی بایگرام و ترکیب "ضمیر + حرف اضافه + فعل" در مدل زبانی تراگرام است صحیح نبوده و حذف می‌شود. از این رو جملهٔ بالقوه معنی‌دار اول به‌عنوان جملهٔ صحیح در نظر گرفته می‌شود.

ترکیب غیرممکن سوم: در زبان فارسی، نشانهٔ مفعول (را) بعد از اسم یا صفت قرار می‌گیرد. از این رو ترکیبات نحوی "قید + را"، "حرف اضافه + را"، "حرف ربط + را" و همچنین "را + را" در گرامر فارسی وجود ندارند و نادرست هستند. از این رو ترکیبات بالا از نظر مدل زبانی بایگرام و به‌طور مشابه ترکیب سه‌تایی ترکیبات دوتایی بالا از نظر مدل زبانی تراگرام صحیح نبوده و حذف می‌شوند.

مشکلاتی از این قبیل بایستی از قواعد معنایی زبان فارسی استفاده کرد؛ زیرا از نظر قواعد دستوری (استفاده از جداول "فراوانی توالی مقولات واژگانی" و حذف ترکیبات غیرممکن)، هر دو جمله صحیح تشخیص داده می‌شوند.

۵- بحث و تحلیل

با توجه به نتایج حاصل از جدول (۲) مشاهده می‌شود اشتباهات، بیشتر در مواردی است که جمله کلمه‌ای داشته که شامل یک حرف و یک اسم معنی‌دار باشد؛ مانند: "باهوش، بامزه، برنامه و ...". به این دلیل که با توجه به جدول گرامر رایج بایگرام (ترایگرام) احتمال ترکیب دوتایی (سه تایی) "حرف اضافه + اسم" در جملات فارسی زیاد است. بنابراین این کلمات به صورت ترکیب "حرف اضافه + اسم" صحیح تشخیص داده می‌شوند. به عنوان مثال کلمه "برچسب" به صورت دو کلمه مجزا "بر" به عنوان حرف اضافه و "چسب" به عنوان اسم تشخیص داده می‌شود. همچنین حضور کلماتی که از دو اسم معنی‌دار تشکیل شده باشند مانند: "کارخانه، آشپزخانه، داروخانه و ..." احتمال تشخیص بیش از یک جمله به عنوان جمله صحیح را افزایش می‌دهند. بیش تر در این گونه موارد، دو جمله به عنوان جملات صحیح تشخیص داده می‌شوند.

نتیجه دیگری که از جداول (۱ و ۲) می‌توان گرفت این است که نتایج حاصل از مدل ترایگرام بهتر از نتایج حاصل از مدل بایگرام است. در واقع به دلیل افزایش تعداد کلمات مورد بررسی در مدل ترایگرام دقت آن نسبت به مدل بایگرام افزایش یافته است. همچنین می‌توان گفت در مدل ترایگرام از تعداد جملات اشتباه تشخیص داده شده مدل بایگرام کاسته شده و به تعداد جملات به طور تقریبی صحیح تشخیص داده شده، افزوده شده است.

۶- نتیجه‌گیری و کارهای آتی

در این مقاله یک روش جدید برای بازشناسی متون فارسی ارائه شد. الگوریتم پیشنهادی در مرحله اول از یکی از روش‌های برچسب‌زنی اجزای متصل به منظور جداسازی زیرکلمات استفاده می‌کند. در مرحله دوم زیرکلمات استخراج شده کنار هم قرار گرفتند و همه کلمات معنی‌دار و سپس جملات بالقوه بامعنی تشکیل شدند. در آخرین مرحله از مدل‌های زبانی بایگرام و ترایگرام و همچنین از قواعد گرامری زبان فارسی به منظور تشخیص جمله صحیح از میان

ترکیبات دستوری غیرممکن نیستند، از این رو هیچ یک از آنها حذف نمی‌شوند؛ سپس، درصدهای مربوط به این چهار جمله با یکدیگر مقایسه می‌شوند و مدل زبانی بایگرام جملات اول و سوم را صحیح تشخیص می‌دهد. درحالی‌که مدل زبانی ترایگرام تنها جمله اول را صحیح تشخیص می‌دهد.

البته غالب جملات موجود در متون استاندارد تنها به یک صورت خوانده می‌شوند و دارای چندین جمله بالقوه معنی‌دار نیستند. به عنوان مثال دیگر، جمله زیر را در نظر بگیرید:

جمله ورودی: شکارچیان شورا را برای فراری دادن شیر از جنگل تشکیل دادند.
از جمله فوق، چهار جمله بالقوه معنی‌دار زیر استخراج می‌شوند:

شکارچیان شورا را برای فراری دادن شیر از جنگل تشکیل دادند.

شکارچیان شو را برای فراری دادن شیر از جنگل تشکیل دادند.

شکارچیان شورا را برای فراری دادن شیراز جنگل تشکیل دادند.

شکارچیان شو را برای فراری دادن شیراز جنگل تشکیل دادند.

جملات دوم و چهارم به علت ترکیب دستوری غیرممکن "را + حذف می‌شوند؛ سپس درصدهای مربوط جملات باقی‌مانده یعنی جملات اول و سوم با یکدیگر مقایسه می‌شوند و هر دو مدل زبانی، جمله اول را به عنوان جمله صحیح تشخیص می‌دهند.

مثال سوم را به صورت زیر در نظر بگیرید:

جمله ورودی: دود کارخانه هوا را آلوده می‌کند.

دو جمله بالقوه معنی‌دار زیر را می‌توان از آن استخراج کرد:

دود کارخانه هوا را آلوده می‌کند.

دود کارخانه هوا را آلوده می‌کند.

بعد از تعیین نقش دستوری کلمات در جملات بالقوه معنی‌دار فوق، عدد مربوط به هر یک از این جملات مشخص می‌شود؛ سپس از نظر ترکیبات دستوری غیرممکن مورد بررسی قرار می‌گیرند و به دلیل این که هیچ کدام دارای ترکیبات دستوری غیرممکن نیستند، حذف نمی‌شوند. مدل زبانی بایگرام و همچنین مدل زبانی ترایگرام هر دو جمله را به عنوان جملات صحیح تشخیص می‌دهند. برای رفع

همکاران، ۱۳۸۶)، به نظر می‌رسد با استفاده از ساختارهای دقیق‌تر گرامر و همچنین ارتباط معنایی بین کلمات و کل جمله، بتوان خطاهای باقی‌مانده را رفع کرد.

مجموعه‌ای از جملات استفاده شد. نتایج نشان داد که با اعمال الگوریتم پیشنهادی می‌توان غالب جملات فارسی را بازشناسی کرد. در عین حال، با توجه به استفاده گسترده از کسره (در ارتباط بین صفت-موصوف و مضاف-مضاف الیه) که در شکل ظاهری کلمات مشخص نمی‌شود (عیسی پورو

(جدول-۱): آمار جملات مورد بررسی

تعداد کل جملات	تعداد جملات با دو جمله بالقوه معنی‌دار	تعداد جملات با سه جمله بالقوه معنی‌دار	تعداد جملات با بیش از سه جمله بالقوه معنی‌دار
۴۲۰	۴۱	۱	۱

(جدول-۲): نتایج حاصل از اعمال الگوریتم پیشنهادی

مدل‌های زبانی	تعداد جملات با چندین جمله بالقوه معنی‌دار	تعداد جملات صحیح تشخیص داده شده	تعداد جملات به‌طور تقریبی صحیح تشخیص داده شده	تعداد جملات اشتباه تشخیص داده شده
بایگرام	۴۳	۳۴	۶	۳
تراگرام	۴۳	۳۵	۷	۱

منابع

عیسی پور، شهریار. همایون پور، محمد مهدی. بی‌جن خان، محمود. "شناسایی محل کسره اضافه در زبان فارسی با استفاده از گرامر مستقل از متن احتمالاتی"، سیزدهمین کنفرانس ملی انجمن کامپیوتر ایران، جزیره کیش، خلیج فارس، ایران ۱۹ الی ۲۱ اسفند ۱۳۸۶.

Adab, M., M.B., "Simultaneous segmentation and recognition of Farsi/Latin printed texts with MLP". Inter-national Joint Conference on Neural Networks., 2002, pp.1534-1539.

AL-Muhtaseb, H., "Recognition of off-line printed Arabic text using hidden Markov Model". Signal Processing., 2008, pp.2902-2912.

Amin, A., Kaced, A., Haton, J.P., Mohr, R., " Hand-written Arabic character recognition by the IRAC system". Proceedings of Fifth International Conference on Pattern Recognition., Miami Beach, FL, USA, 1980, pp.729-731.

Badie, K., Shimura, M., "Machine recognition of Arabic cursive scripts". Proceedings of International Workshop Pattern Recognition in Practice., Amsterdam, Netherlands, 1980, pp. 315-323.

ابراهیمی، افشین. "استفاده از شکل کلی زیرکلمات چایی در بازیابی تصویر مستندات و بازشناسی متون فارسی"، رساله دکتری بخش مهندسی برق، دانشگاه تربیت مدرس، تهران، ایران، ۱۳۸۴.

حسنی، حمید. پژوهشگر گروه فرهنگ‌نویسی فرهنگستان زبان و ادب فارسی. "واژه‌های پرکاربرد فارسی امروز"، نشر کانون فرهنگی، ۱۳۸۴.

رضوی، سید محمد. کبیر، احسان الله. "بازشناسی برخط حروف مجزای فارسی با شبکه عصبی"، مجموعه مقالات سومین کنفرانس ماشین بینایی و پردازش تصویر ایران ص ۸۳-۸۹، دانشگاه تهران، اسفند ماه ۱۳۸۳.

رضوی، سید محمد. کبیر، احسان الله. "بازشناسی برخط زیرکلمات فارسی با استفاده از نقاط و علائم حروف"، دومین کنفرانس فناوری اطلاعات و دانش، دانشگاه صنعتی امیرکبیر، تهران، خرداد ماه ۱۳۸۴.



پریسا شیروانی دوره کارشناسی خود را در سال ۱۳۸۷ در رشته مهندسی برق (مخابرات) گذرانده و مدرک کارشناسی ارشدش را در سال ۱۳۹۰ در رشته مهندسی برق (مخابرات) از دانشگاه سمنان اخذ کرد. وی هم‌اکنون در دانشگاه علوم و فناوری سپاهان شهر اصفهان مشغول به فعالیت و زمینه‌های علمی مورد علاقه وی پردازش سیگنال به‌ویژه پردازش متن است.

نشانی رایانامه ایشان عبارت است از:

Shirvani.parisa@gmail.com



مهرداد وطن‌خواه خوزانی دوره کارشناسی خود را در سال ۱۳۸۵ در رشته مهندسی کامپیوتر (نرم‌افزار) در دانشگاه آزاد اسلامی واحد نجف‌آباد گذرانده و هم‌اکنون دانشجوی دوره کارشناسی ارشد در رشته مهندسی نرم‌افزار گرایش شبکه در دانشگاه شفیلد هالام انگلستان است. زمینه‌های علمی مورد علاقه وی شبکه‌های عصبی، طراحی و پیاده‌سازی نرم‌افزارهای یکپارچه و نهان‌نگاری اطلاعات است.

نشانی رایانامه ایشان عبارت است از:

Mehrdad.vatankhah@gmail.com



خشایار یغمائی: دوره کارشناسی خود را در سال ۱۳۶۲ در رشته مهندسی برق (الکترونیک) و دوره کارشناسی ارشد خود را در رشته مهندسی برق (مخابرات) در سال ۱۳۶۴ در دانشگاه تهران گذرانده است. همچنین مدرک دکترای خود را در سال ۱۳۷۶ در رشته مهندسی برق (مخابرات) از دانشگاه Surrey انگلستان اخذ کرد و به دانشکده مهندسی برق و کامپیوتر دانشگاه سمنان پیوست و در این دانشگاه مشغول به فعالیت است. زمینه‌های علمی مورد علاقه وی پردازش سیگنال، پردازش صوت و تصویر و پردازش متن است.

نشانی رایانامه ایشان عبارت است از:

khashayar.yaghmaie@gmail.com

Bahmani.Z and Alamdar.F, "Off-line Arabic/Farsi Handwritten Word Recognition Using RBF Neural Network and Genetic algorithm", *IEEE*, vol.2, no.3, pp.352-357, 2010.

Corazza, A., De Mori, R., Gretter, R., Satta, G., "Language modeling using stochastic context-free grammars". *Speech Communication*, Vol.13(1-2), 1993, pp. 163-170.

Ebrahimi, A., Kabir, E., "A pictorial dictionary for printed Farsi sub-words". *Pattern Recognition Letters*, 2008, pp. 656-663.

Ebrahimpour.R and Davoudi.R, "Decision Templates with Gradient based Features for Farsi Handwritten word Recognition", *International Journal of Hybrid Information Technology*, Vol.4, No.1, 2011.

El-Abed, H., Margner, V., "Arabic text recognition systems-state of the art and future trends". *International conference on Innovations in Information Technology*, 2008, pp. 692-696.

Jacobs, C., Simard, P., Rinker, Viola and J., "Text recognition of low-resolution document images". *Eight International Conference on Document Analysis and Recognition (ICDAR'05)*, Vol.2, 2005, pp. 695-699.

Khosravi, H., Kabir, E., "Farsi font recognition based on Sobel-Roberts features". *Pattern Recognition Letters*, 2008, pp.75-

Rosenfeld, R., "Two decades of statistical language modeling: where do we go from here?". *Proceedings of the IEEE*, Vol.88, 2000, pp.1270-1278.

Rosenfeld, R., Chen, S. F. and Zhu, X., "Whole-sentence exponential language models: a vehicle for linguistic-statistical integration". *Computer Speech & Language*, Vol.15 (1), 2001, pp.55-73.

Sarfaraz, M., Nawaz, N., Al-Khuraidly, A., "Offline Arabic text recognition system". *International Conference on Geometric Modeling and Graphics (GMAG'03)*, 2003, pp.30.

Srihar, R., Shetty, S., Srihar, S., "Use of language models in handwriting recognition". *Center of excellence for document analysis and recognition (CEDAR)*, 2007.

