

A Corpus-based Study of Persian Noun and Adjective Homographs to help Correct POS Tagging

Elham Alayiaboobar

PhD of General Linguistics; Iranian Research Institute for Information Science and Technology (IranDoc);
Email: Alayi@irandoc.ac.ir

Received: 14, Oct. 2017 | Accepted: 25, Feb. 2018

Abstract: Present research studies morphological structure of nouns and adjectives. There are two main reasons for studying them in the process of making any POS tagger system for tagging nouns: 1. If the system faces an out of vocabulary word (OOV word), one way to identify its tag would be considering its morphological structure; 2. In Persian, lots of homographs are made due to Persian complex morphology. Studying morphological structure of nouns in order to distinguish them from adjectives seems to be necessary, since many adjectives, having the same orthographic forms of nouns, would be wrongly tagged as "noun" or vic e versa. After studying morphological structure of nouns and adjectives in present study, Persian writing system is studied. Then definition of homographs and the related classifications are presented. Finally, the study uses different famous Persian corpora (including Bijankhan, and syntactical dependency corpus (vabastegi ye nahvi) for searching for homographs (using search tools) and Data Center for Persian Language (Paygah e Dadegan) whose non-tagged file was available (the homographs are searched and tagged manually)) to make a list of homographs. The result of studying the mentioned list showed that the frequency of homographs, especially those which are made due to identical orthographic form of indefinite morpheme, adjective-maker morpheme and second person inflectional morpheme is high in Persian corpora which makes POS tagging difficult.

Keywords: POS Tagger System, Morphological Structure of Persian Nouns and Adjectives, Persian Writing System, Homographs

Iranian Journal of
**Information
Processing and
Management**

Iranian Research Institute
for Information Science and Technology
(IranDoc)

ISSN 2251-8223

eISSN 2251-8231

Indexed by SCOPUS, ISC, & LISTA

Vol. 34 | No. 2 | pp. 897-922

Winter 2019



بررسی پیکره-بنیاد هم‌نگاره‌های اسمی و صفتی فارسی جهت کمک به برچسب‌گذاری صحیح اجزای کلام

الهام علایی ابودر

دکتری زبان‌شناسی همگانی؛ استادیار؛
پژوهشگاه علوم و فناوری اطلاعات ایران (ایرانداک)؛
Alayi@irandoc.ac.ir



دریافت: ۱۳۹۶/۷/۲۲ پذیرش: ۱۳۹۶/۱۲/۰۶ مقاله برای اصلاح به مدت شش روز نزد پدیدآوران بوده است.

فصلنامه | علمی پژوهشی
پژوهشگاه علوم و فناوری اطلاعات ایران
(ایرانداک)
شاپا (چاپی) ۲۲۵۱-۸۲۲۳
شاپا (الکترونیکی) ۲۲۵۱-۸۲۳۱
نمایه در SCOPUS، ISC، LISTA و
jipm.irandoc.ac.ir
دوره ۳۴ | شماره ۲ | صص ۸۹۷-۹۲۲
زمستان ۱۳۹۷



چکیده: در تهیه سامانه‌های برچسب‌گذاری اجزای کلام در زبان فارسی، بررسی ساخت‌واژی اسم‌ها و صفت‌ها از دو نظر حائز اهمیت است: ۱. اگر در یک پیکره متنی فارسی کلمه‌ای در واژگان حضور نداشته باشد (کلمه خارج از واژگان)، نمی‌توان برچسب‌های مربوط به کلمه را بازیابی کرد. در این صورت، برچسب کلمه را تنها می‌توان با توجه به شکل کلمه (انواع پیشوندها و پسوندهایی که به کلمات متصل می‌شوند) یا بافتی که کلمه در آن ظاهر می‌شود، یا هر دو، حدس زد. ۲. زبان فارسی ظرفیت بالایی برای ساخت هم‌نگاره‌های جدید که از ساخت‌واژه فارسی نشأت می‌گیرند، نیز دارد. بنابراین، بررسی ساخت‌واژی اسم‌ها و صفت‌ها، به‌منظور تفکیک آن‌ها از هم ضروری به نظر می‌رسد، زیرا اکثر صفت‌ها در بافت‌های گوناگون، با صورت نوشتاری یکسان می‌توانند برچسب «اسم» بگیرند. در تحقیق حاضر ساخت‌واژه اسم‌ها و صفت‌ها در فارسی بررسی شده است. نظام نوشتاری زبان فارسی نیز مورد بررسی قرار گرفته تا از این رهگذر بتوان به شناسایی انواع هم‌نگاره‌ها در زبان فارسی پرداخت. سپس، انواع هم‌نگاره‌ها در زبان فارسی مورد مطالعه قرار گرفته و در نهایت، از طریق جست‌وجو به دو روش ماشینی و دستی، فهرست مبسوطی از هم‌نگاره‌ها از پیکره‌های «پیکره متنی زبان فارسی»، «پایگاه دادگان زبان فارسی» و «پیکره وابستگی نحوی زبان فارسی» تهیه شده است. بررسی کلی هم‌نگاره‌ها در پیکره‌های مورد مطالعه نشان می‌دهد که بیشتر هم‌نگاره‌ها، فراوانی بالایی در پیکره‌های متنی فارسی دارند و اکثر آن‌ها در اثر یکسان بودن نمود نوشتاری تک‌واژ یای نکره، یای اسم‌ساز، شناسه دوم شخص مفرد، یای صفت‌ساز و یای متصل به گروه اسمی ایجاد شده‌اند.

کلیدواژه‌ها: برچسب‌گذاری اجزای کلام، ساخت‌واژه اسم‌ها و صفت‌ها، نظام نوشتاری، هم‌نگاره‌ها

۱. مقدمه

مفهوم برچسب‌گذاری اجزای کلام^۱، مفهومی برگرفته از زبان‌شناسی است که به طبقه‌نحوی کلمات مربوط می‌شود. به این صورت که ویژگی‌های ترکیبی یک کلمه در بافت/ ساخت نحوی مشخص می‌شود. در زبان‌شناسی پیکره‌ای، برچسب‌گذاری اجزای کلام، در واقع، عمل انتساب برچسب به کلمات تشکیل‌دهنده یک متن یا یک پیکره است. این برچسب‌گذاری بر اساس مقوله آن کلمه در متن (مانند «اسم»، «فعل»، «قید»، «صفت»، و غیره) صورت می‌پذیرد. کلمات در جایگاه‌های مختلف می‌توانند برچسب‌های واژگانی متفاوت داشته باشند. به عنوان مثال، کلمه «آسمانی» در جمله «فردا آسمانی صاف در انتظار شهروندان عزیز خواهد بود»، دارای برچسب «اسم، نکره»، ولی در جمله «او انسانی آسمانی بود»، دارای برچسب «صفت» است. همچنین، کلمه «مردی» در جمله «در راه مدرسه مردی را دیدم»، دارای برچسب «اسم، نکره»، ولی در جمله «مثلاً تو مردی؟»، دارای برچسب «اسم + فعل اسنادی» است. بنابراین، درصد بالایی از کلمات از نقطه‌نظر برچسب‌واژگانی، مبهم هستند و می‌توانند فارغ از بافت، بیش از یک برچسب‌واژگانی داشته باشند.

برچسب‌گذاری اجزای واژگانی کلام، عملی کاربردی در بسیاری از حوزه‌های پیشرفته‌تر پردازش زبان طبیعی^۲ از جمله ترجمه ماشینی، خطایاب، تبدیل متن به گفتار، بازیابی اطلاعات، موتورهای جست‌وجو و کمک به مدل‌های آماری است (Megerdoomian 2004). همچنین، در طراحی سیستم نمایه‌ساز ماشینی، یکی از بخش‌ها، طراحی زیرسیستم‌هاست که خود بخش‌های گوناگونی دارد. یکی از این بخش‌ها، زیرسیستم تحلیل واژگانی است. این زیرسیستم، متن را به واژه‌ها تفکیک می‌کند و ماهیت هر کلمه را تشخیص می‌دهد و تشخیص نوع واژه و شناسایی فعل‌ها، الفاظ و اصطلاح‌ها را دربردارد. بنابراین، سیستم برچسب‌دهی خودکار، ماهیت مقوله کلمات را مشخص می‌کند تا بتوان در مراحل بعدی از این اطلاعات در جهت استخراج کلیدواژه‌ها یا هر نوع بازیابی اطلاعات از متن استفاده کرد. تاکنون از مدل‌ها و روش‌های بسیاری برای برچسب‌گذاری

1. part of speech (POS) tagging

2. natural language processing

در زبان‌های مختلف استفاده شده است. این روش‌ها و مدل‌ها را می‌توان به دو دسته کلی تقسیم‌بندی کرد: دسته اول، رویکردهای آماری است که از به کار بردن روش‌های احتمالاتی برای پیکره‌های برچسب خورده^۱ بهره می‌جویند و دسته دوم، رویکردهای غیر آماری و مبتنی بر دانش^۲ / مبتنی بر قاعده^۳ است که بر مبنای یادگیری ماشینی و دانش بشری استوار هستند. البته، رویکرد سومی هم وجود دارد که در واقع، تلفیقی از دو رویکرد آماری و مبتنی بر قانون است که برچسب‌گذاری تلفیقی یا دووجهی^۴ خوانده می‌شود (Megerdoozian 2004).

در زبان‌هایی که ساخت واژه تصریفی و اشتقاقی آن‌ها پیچیده است، ظرفیت بالایی برای ساخت کلمات با اشکال جدید وجود دارد، زیرا در این زبان‌ها، تکواژها به روش‌های گوناگون به یکدیگر متصل می‌شوند و کلمات جدیدی را تولید می‌کنند که امکان دارد این کلمات حتی در پیکره‌های بزرگ وجود نداشته باشند و یا فراوانی آن‌ها بسیار اندک باشد. تکواژهای تصریفی معمولاً مفاهیم دستوری را به کلمه اضافه می‌کنند و بنابراین، مقوله نحوی را تغییر نمی‌دهند. اما تکواژهای اشتقاقی کلمات جدید می‌سازند و دارای مفهومی خاص هستند یا مقوله نحوی کلمه را تغییر می‌دهند. بررسی ساخت واژگی کلمات از دو نظر حائز اهمیت است: ۱. شناسایی برچسب کلماتی که در واژگان وجود ندارند و برچسب نحوی چنین کلماتی را تنها می‌توان با توجه به شکل کلمه یا بافتی که کلمه در آن ظاهر می‌شود، یا هر دو، حدس زد. ۲. کمک به رفع ابهام از برچسب نحوی هم‌نگاره‌ها^۵ در فارسی. بسیاری از وندها، اشتقاقی و تصریفی، نمود نوشتاری مشابهی دارند و در اتصال به ستاک، هم‌نگاره‌های گوناگون می‌سازند. به عنوان مثال، نمود نوشتاری وند تصریفی نکره ساز /i/ و وند اشتقاقی صفت ساز /i/ یکسان است («ی») و کلماتی مانند: «آسمانی» و «اسلامی»، فارغ از بافت، می‌توانند هم برچسب «اسم» داشته باشند و هم «صفت». گویشوران فارسی به دلیل آشنایی کامل با قواعد ساخت واژگی زبان خود، هنگام مواجهه با چنین کلماتی که حاوی تکواژهای تصریفی و اشتقاقی هستند یا برای اولین بار با آن‌ها در متون گوناگون فارسی مواجه می‌شوند، در فهم و پردازش آن‌ها معمولاً دچار مشکل نمی‌شوند. اما این مسئله راجع به سامانه‌های پردازش رایانه‌ای متن، صحیح نیست.

1. tagged corpora

2. knowledge-based

3. rule-based

4. hybrid tagging

5. homographs

این سامانه‌ها، به‌خصوص سامانه‌های برجسب‌گذاری، به‌دلیل عدم اشراف به قواعد ساخت‌واژی زبان، در برخورد با کلمات دارای پیچیدگی‌های ساخت‌واژی، توان محدودی دارند. همچنین، ویژگی‌های خاص دستوری و نگارشی زبان و خط فارسی، دشواری‌هایی را در ذخیره و بازیابی اطلاعات در محیط رایانه‌ای نیز پدید آورده است. رسم‌الخط فارسی نیز از یک‌سو به‌علت اختلاف نظر پدیدآورندگان متون و از سوی دیگر، پیچیدگی‌های ذاتی خود، به‌هنگام ذخیره، جست‌وجو و بازیابی اطلاعات، چالش‌های متعددی را برای طراحان و نمایه‌سازان پایگاه‌ها، کاربران و پدیدآورندگان منابع به‌وجود آورده است (آخشیک و فتاحی ۱۳۹۱). این پیچیدگی‌ها مشکلات بسیاری را در مسیر برجسب‌گذاری رایانه‌ای اجزای واژگانی کلام ایجاد می‌کنند. برخی از این مشکلات شامل موارد زیر است:

الف. پیچیدگی ساخت‌واژه زبان فارسی و ارتباط آن با خط فارسی: اگر چند وند در یک کلمه ظاهر شوند، همه آن‌ها معمولاً به کلمه می‌چسبند (Megerdooian 2000). نشانه‌های جمع، کسره اضافه، نشانه نکره، ضمایر ملکی و غیره می‌توانند به کلمه متصل شوند؛ مانند: «کتابهایم» که شامل «کتاب + ها (علامت جمع) + ی» (نمود نوشتاری همخوان میانجی /j/ + م (ضمیر ملکی اول شخص مفرد)) است. این ویژگی باعث اشکال متفاوتی از کلمات با ریشه یکسان می‌گردد و این کلمات در سامانه‌های محاسباتی، متفاوت از یکدیگر فرض می‌شوند. در مورد افعال می‌توان گفت از نقطه‌نظر ساخت‌واژی، فعل شامل بن فعل و وندهای تصریفی است. افعال با توجه به شخص و شمار صرف می‌شوند و بنابراین، اشکال متفاوتی از آن‌ها به‌وجود می‌آید. شکل یکسان برخی از تکواژها نیز باعث ابهام در متون فارسی می‌شود. به‌عنوان مثال، پسوند «-ی» در کلمه «مردی»، فارغ از بافت، می‌تواند هم نشانه نکره در نظر گرفته شود و هم شناسه دوم شخص مفرد در یک فعل اسنادی («مرد» + «ی») (به معنی مرد هستی)). علاوه بر این، در فارسی معمولاً واژه‌های کوتاه نمود نوشتاری ندارند و این مسئله نیز باعث ابهام در تحلیل می‌گردد؛ مانند: کلمه «مردم» که می‌تواند به‌صورت‌های /mardam/ و /mordam/ و /mardom/ تلفظ شود. این ابهام، ابهام‌هم‌نگاره خواننده می‌شود.

ب. عبارت‌های چند کلمه‌ای: حضور افعال مرکب و کلمات مرکب اسمی مشکلی دیگر در پردازش زبانی محسوب می‌شود. به‌عنوان مثال، در کلمات مرکب اسمی که هسته آغاز

هستند، یعنی هسته گروه اسمی در ابتدای عبارت ظاهر شده است، وند جمع به قسمت اول متصل می‌شود. مانند: ماشین لباس شویی + وند جمع = ماشین‌های لباس شویی.

ج. مرز کلمات: فاصله نیز یکی از عوامل ابهام در متون فارسی است. در متون گوناگون، فاصله ماهیت اختیاری دارد. این مسئله در برخی موارد باعث می‌شود که کلمات مجزا یک قطعه^۱ در نظر گرفته شوند (مانند: «رفتند مردم»). همچنین، وندهای تصریفی (مانند: «-تر»، «-ترین»، «-ها»، «-می») را می‌توان به سه صورت به ستاک مربوط کرد: ۱) به ستاک چسباند، ۲) به صورت مجزا (غیرچسبان) با یک فاصله از ستاک در آخر کلمه / قطعه آورد، و ۳) به صورت غیرچسبان ولی با استفاده از نویسه نیم‌فاصله به ستاک مربوط کرد. به عنوان مثال، تکواژ جمع «-ها» در اسم‌ها می‌تواند به چند شکل ظاهر شود. در مورد کلمه «کتاب»، سه شکل «کتابها»، «کتاب‌ها» و «کتاب‌ها» برای حالت جمع در خط فارسی وجود دارد. در مورد افعال نیز می‌توان به پیشوند تصریفی زمان استمراری (یعنی «-می») اشاره کرد. مانند، حال استمراری اول شخص مفرد بن مضارع «رو» که هنگام اضافه‌شدن وند تصریفی «-می» می‌تواند به سه شکل «می‌روم»، «می‌روم» و «می‌روم» نوشته شود. این مورد نیز عامل تأثیرگذار دیگری در برجسب‌گذاری اجزای واژگانی کلام در زبان فارسی محسوب می‌شود، زیرا چند شکل نوشتاری متفاوت از یک کلمه عملاً به عنوان کلمات متفاوت تفسیر می‌شوند. هر سامانه برجسب‌گذاری باید توانایی تشخیص این شکل‌ها را داشته باشد تا برجسب دستوری صحیح به کلمه بدهد (Megerdooonian 2004). لازم به ذکر است که «فرهنگستان زبان و ادب فارسی» در باب پیوسته‌نویسی یا جدانویسی ترکیبات در زبان فارسی سه فرض را متصور است: ۱. تدوین قواعدی برای جدانویسی همه کلمات مرکب و تعیین موارد استثنا، ۲. تدوین قواعدی برای پیوسته‌نویسی همه کلمات مرکب و تعیین موارد استثنا، و ۳. تدوین قواعدی برای جدانویسی الزامی بعضی

1. token

متن، رشته‌ای از قطعات نوشتاری است و قطعه، نمود نوشتاری یک تکواژ یا دنباله‌ای از تکواژهاست که معمولاً از راست یا چپ، یا هر دو، از طریق فاصله یا علائم نقطه‌گذاری از دیگر قطعات تفکیک می‌شود. به عنوان مثال، کلمه‌ای مانند «می‌روم» یک قطعه نوشتاری است که از سه واحد زبانی تشکیل شده است: «می»، «رو» و «م». دو واحد زبانی اول با نیم‌فاصله از هم متمایز شده‌اند و چون حرف واو («و») به حرف بعد از خود نمی‌چسبد، استفاده از نیم‌فاصله بین واحد دوم و سوم ضرورتی ندارد (علایی و بی‌جن‌خان ۱۳۹۲).

از کلمات مرکب و پیوسته‌نویسی بعضی دیگر و دادن اختیار در خصوص سایر کلمات به نویسندگان. «فرهنگستان» در تدوین و تصویب دستور خط فارسی، فرض سوم را برگزیده و تنها موارد الزامی جدانویسی و پیوسته‌نویسی را مشخص کرده است (دستور خط فارسی ۱۳۸۸ نقل در آخشیک و فتاحی ۱۳۹۱).

د. قطعه‌های پیچیده^۱: منظور از قطعه‌های پیچیده، شکل‌های چندجزئی هستند که در واقع، از چسبیدن مقولات دستوری (مانند: «به»، «را»، «که» و «آن») به ستاک به‌وجود می‌آیند که این تکواژهای وابسته/مقولات دستوری چسبیده به ستاک، تکواژ شبه کلمه^۲ خوانده می‌شوند. مانند: «بدین صورت»، «بعقیده»، «اینکار».

ح. قواعد آوایی و واجی: در فارسی، شکل وندهای متصل به ستاک بر اساس کاراکتر پایانی ستاک تغییر می‌کند. به‌عنوان مثال، اگر اسم جاندار به همخوان ختم شود، تکواژ جمع «ان-»/an- به ستاک متصل می‌شود. مانند: «زن + ان = زنان»؛ و اگر به واکه ختم شود، همخوان غلت /z/ با نمود نوشتاری <ی> بین ستاک و پسوند قرار می‌گیرد. مانند: گدا+ان = گدایان.

و. مرز عبارت‌ها: مسائل گوناگونی باعث ابهام در ساختار عبارت‌های اسمی در فارسی می‌شود. تکواژهای آشکار بسیار اندکی در فارسی برای نشان‌دادن مرزهای عبارات اسمی وجود دارد. اغلب هیچ جزئی در خط فارسی برای متصل کردن اجزای یک عبارت اسمی وجود ندارد، زیرا تکواژ اضافه^۳، اغلب واکه کوتاهی است که در خط ظاهر نمی‌شود. علاوه بر این، از آنجا که ترتیب پایه و غالب اجزای جمله در فارسی، فاعل-مفعول-فعل^۴ است، فقدان تکواژهای آشکار نشان‌دهنده مرزها، تشخیص محل به پایان رسیدن فاعل و شروع مفعول را سخت می‌کند و دشواری‌هایی را در پردازش متن پدید می‌آورد (Megerdooian 2004).

از میان مسائل مطرح‌شده فوق، مسائلی که از ساخت‌واژه فارسی ناشی می‌شوند از اهمیت بیشتری برخوردار هستند، زیرا ساخت‌واژه علاوه بر شکل کلمه، برچسب کلمه را نیز تحت تأثیر قرار می‌دهد. در زبان فارسی، تکواژهای تصریفی و اشتقاقی بسیاری وجود دارد که با اتصال به کلمات باعث تغییر برچسب کلمات در پیکره‌های زبانی و

1. complex tokens

2. word-like morpheme

3. ezafe morpheme

4. subject-object-verb (SOV)

همچنین، کاهش فراوانی کلمات و در نهایت، کاهش کارایی سامانه‌های برچسب‌گذاری آماری در زبان فارسی می‌شود. بنابراین، سامانه‌های برچسب‌گذاری باید حاوی تحلیلگر ساخت‌واژی نیز باشند (محسنی ۱۳۸۷).

در تهیه سامانه‌های برچسب‌گذاری در زبان فارسی (نمونه فارسی: سامانه «هضم» *sobhe/hazm*)، بررسی ساخت‌واژی اسم‌ها و صفت‌ها از دو نظر حائز اهمیت است:

۱. اگر کلمه‌ای قبلاً در پیکره متنی ظاهر نشده باشد، نمی‌توان از پیکره، اطلاعات دقیقی راجع به آن کلمه به دست آورد و حتی از توزیع کلمات در پیکره نیز نمی‌توان استفاده کرد، زیرا توزیع کلمات ناشناخته کاملاً متفاوت از کلمات شناخته شده در متن است. بنابراین، اگر در یک پیکره متنی فارسی، کلمه‌ای در واژگان حضور نداشته باشد (کلمه خارج از واژگان)^۱ نمی‌توان برچسب‌هایی را که کلمه به آن‌ها منتسب می‌شود، بازبانی کرد. در این صورت، برچسب کلمه را تنها می‌توان با توجه به شکل کلمه یا بافتی که کلمه در آن ظاهر می‌شود، یا هر دو، حدس زد. منظور از شکل کلمه، ویژگی‌هایی مانند انواع پیشوندها و پسوندهایی است که به کلمات متصل می‌شوند. طبق گزارشات، روش‌هایی که بر استفاده از چنین اطلاعاتی مبتنی هستند، حاوی دقت بالاتری در زبان انگلیسی هستند (Brill 1994 و weischedel 1993 نقل در محسنی ۱۳۸۷). به عنوان مثال، یکی از این روش‌ها، روش برچسب‌گذار زیراکس^۲ است که در آن از قوانینی بر اساس بخش‌های پایانی کلمات، یعنی حروف انتهایی کلمات، برای انتساب مجموعه‌ای از برچسب‌ها به کلمات خارج از واژگان استفاده می‌کند (Cutting et al. 1992 نقل در محسنی ۱۳۸۷). بنابراین، در وهله اول، بررسی ساخت‌واژی اسم‌ها و صفت‌ها، فارغ از بافت ضروری به نظر می‌رسد تا از این رهگذر بتوان اسم‌ها و صفت‌ها را در صورت دارا بودن شاخص‌های صرفی و اشتقاقی، شناسایی و سپس برچسب‌دهی کرد.

۲. زبان فارسی ظرفیت بالایی برای ساخت هم‌نگاره‌های جدید را نیز که از ساخت‌واژه فارسی نشأت می‌گیرند، دارد. بسیاری از وندها، اشتقاقی و تصریفی، نمود نوشتاری مشابهی دارند و در اتصال به ستاک، هم‌نگاره‌های گوناگون می‌سازند. به عنوان مثال، نمود نوشتاری وند تصریفی نکره‌ساز /i/ و وند اشتقاقی صفت‌ساز /i/ یکسان است («-ی») و کلماتی مانند: «آسمانی» و «اسلامی»، فارغ از بافت، می‌توانند هم برچسب

1. out of vocabulary word (OOV word)

2. xerox

«اسم» داشته باشند و هم «صفت». بنابراین، بررسی ساخت‌واژی اسم‌ها و صفت‌ها، به‌منظور تفکیک آن‌ها از هم ضروری به نظر می‌رسد، زیرا اکثر صفت‌ها در بافت‌های گوناگون، با صورت نوشتاری یکسان، می‌توانند برچسب «اسم» بگیرند.

در پژوهش حاضر ابتدا به پیشینه پژوهش و بررسی مطالعات انجام‌شده در زمینه ساخت‌واژه فارسی و سامانه‌های برچسب‌دهی پرداخته می‌شود. سپس، خلاصه‌ای از بررسی ساخت‌واژه اسم‌ها و صفت‌ها در فارسی ارائه خواهد شد و به دنبال آن، بخش بعد، مطالعه هم‌نگاره‌ها در فارسی است و در نهایت، روش جمع‌آوری داده‌ها و معرفی پیکره‌ها و پایگاه‌های دادگان و نتیجه‌گیری ارائه خواهد شد.

۲. پیشینه پژوهش

از آنجا که در این تحقیق، هدف از بررسی ساخت‌واژی اسم‌ها و به دنبال آن، صفت‌ها، کمک به شناسایی و برچسب‌گذاری هم‌نگاره‌های اسمی و صفتی در پیکره‌های علمی است، بررسی مطالعات و تحقیقاتی که در زمینه برچسب‌گذاری اجزای واژگانی کلام صورت گرفته است، ضروری به نظر می‌رسد. بنابراین، در این بخش علاوه بر ذکر چند نمونه از تحقیقاتی که در زمینه ساخت‌واژه اسم‌ها و بعضاً، صفت‌ها، صورت گرفته است، به ذکر نمونه‌هایی از تحقیقاتی که در زمینه برچسب‌گذاری اجزای واژگانی کلام انجام شده نیز پرداخته خواهد شد.

۱-۲. پیشینه مطالعاتی که در زمینه برچسب‌گذاری اجزای واژگانی کلام انجام شده است

«هرست» روشی دقیق و ساده برای ابهام‌زدایی از هم‌نگاره‌های اسمی معرفی می‌کند که از پیکره متنی بزرگی در این راستا بهره می‌گیرد. الگوریتم معرفی شده، بافت اطراف اسم هدف را چک می‌کند و معنایی را انتخاب می‌کند که بیشترین شاهد برای انتخاب آن موجود باشد. این شواهد شامل مجموعه‌ای از ویژگی‌های واژگانی، نحوی، و نوشتاری است (Hearst 19991). «مریالدو» «مدل مارکوفی پنهان»^۱ را معرفی می‌کند که در آن دو کلمه قبل به‌عنوان سابقه^۲ در نظر گرفته می‌شود. در این مدل، دادن هر برچسبی بستگی به برچسب دو کلمه قبل دارد که به‌صورت زیر نمایش داده می‌شود:

1. Hidden Markov Model (HMM)

2. history

$$p(t_i | w_{i-1}, t_{i-1}, w_{i-2}, t_{i-2}, \dots, w_1, t_1) = p(t_i | t_{i-1}, t_{i-2})$$

(Merialdo 1994). «عاصی و عبدالحسینی» سامانهٔ برچسب‌دهی نحوی را به‌عنوان پروژه‌ای در «پژوهشگاه علوم انسانی و مطالعات فرهنگی» معرفی می‌کنند که سامانهٔ مذکور، خود، در خدمت تهیهٔ پیکره‌ای فارسی به نام «پایگاه دادگان زبان فارسی»^۱ قرار خواهد گرفت و می‌توان گفت که اولین تلاش برای تهیهٔ سامانه‌ای جهت برچسب‌دهی نحوی زبان فارسی است. آن‌ها به بررسی روشی می‌پردازند که Schuetze (1995) برای برچسب‌دهی نحوی در متون انگلیسی به کار برده است و هدف آن‌ها بررسی به کارگیری همان روش در زبان فارسی است. آن‌ها معتقدند که با بهره‌گیری از محاسبات دقیق‌تر می‌توان همان روش را در فارسی به کار برد (Assi and Abdolhosseini 2000). «هو» و همکاران از طریق اضافه کردن اطلاعات نحوی در کنار اطلاعات ساخت‌واژی مربوط به هر کلمه در داخل شبکهٔ مربوطه به بهبود «لینگویستیکا»^۲ کمک می‌کنند. آن‌ها علاوه بر درج فهرست ستاک‌ها و وندها، فهرست دیگری که مربوط به برچسب‌های نحوی کلمات است در شبکه وارد کرده‌اند (Hu et al. 2005). «الحاج» نحوهٔ برچسب‌دهی کلمات موجود در «قرآن کریم» را معرفی می‌کند که خود بخشی از سامانه‌های تجزیه و تحلیل رایانه‌ای «قرآن کریم» است (Elhaj 2009). در واقع، سامانه‌ای که وی معرفی می‌کند، تلفیقی است از تجزیه و تحلیلی ساخت‌واژی و «مدل مارکوفی پنهان» که بر اساس ساختار جملات عربی است. «فرودل» برچسب‌دهی با استفاده از شبکه‌های عصبی را معرفی و برای این منظور از دو روش استفاده می‌کند: ۱. استفاده از شبکه‌های پرسپترون چندلایهٔ تکرار شونده که محتمل‌ترین برچسب نحوی برای کلمه را در متن به‌طور احتمالاتی در نظر می‌گیرد. ۲. در روش دوم کلمات به بردارهایی از مشخصه‌ها در فضای چندبعدی تبدیل می‌شوند. وی معتقد است روش معرفی شده ۹۴ درصد برچسب‌دهی نحوی درست از متن ارائه می‌دهد (Frodl 2013). «راسخ و فخر احمد» سعی دارند ضمن بررسی الگوریتم‌های موجود در رفع ابهام معنایی کلمه، روش‌هایی بر اساس بررسی مفهومی و ساختاری کلمات در جهت رفع ابهام کلمات در متون ارائه نمایند و به این ترتیب، راهکاری در جهت رفع مشکلات ترجمهٔ ماشینی و بهبود کیفیت آن ارائه می‌کنند (۲۰۱۴). «محسنی» با در نظر گرفتن مسائل و مشکلاتی که در مسیر برچسب‌گذاری

1. Farsi Linguistic Data Base (FLDB)

2. Linguistica

اجزای کلام در فارسی وجود دارد، ابتدا طرحی کلی برای برچسب‌گذاری خودکار با دقت بالا در زبان فارسی پیشنهاد می‌کند. سپس، تحلیل ساخت‌وازی و استفاده از آن را برای پوشش دادن تعداد زیادی از برچسب‌های پیکره با حفظ دقت بالا در برچسب‌گذاری کلمات مورد بررسی دقیق‌تر قرار داده و تأثیر وجود یک تحلیلگر ساخت‌وازی در سطح تصریف را بر برچسب‌گذاری اجزای واژگانی کلام در زبان فارسی بررسی می‌کند. وی معتقد است که نتایج به‌دست آمده نشان از کارایی بسیار مناسب این روش پیشنهادی در برچسب‌گذاری دارد (۱۳۸۸).

۲-۲. نمونه‌هایی از مطالعاتی که در زمینه ساخت‌واژه اسم‌ها و بعضاً، صفت‌ها در فارسی صورت گرفته است

«صادقی» طی ۱۲ مقاله، به شیوه‌ها و امکانات واژه‌سازی در زبان فارسی معاصر می‌پردازد. وی این گفتارها را با بحث از اشتقاق در زبان فارسی معاصر آغاز می‌کند و در صدر آن‌ها از امکانات واژه‌سازی از افعال مرکب، که در سال‌های اخیر به‌عنوان عامل عقیم و غیرزایا بودن زبان فارسی مورد حمله بعضی از محققان قرار گرفته است، گفت‌وگو می‌کند. وی همچنین، طبقه‌بندی از انواع پسوندهای اسم‌ساز و صفت‌ساز را ارائه می‌دهد (۱۳۷۰-۱۳۷۲). «کشانی» ابتدا به تعریف پسوند و نقش آن در زبان فارسی می‌پردازد، و سپس، پسوند و پسوندواره‌ها و انواع آن‌ها را در زبان فارسی به تفصیل شرح می‌دهد (نوعاً پسوندهای اشتقاقی که در ساختن انواع اسم، انواع صفت و قید به کار می‌روند به دقت بررسی می‌شود) و در نهایت، حرکت پسوندها را مورد مطالعه قرار می‌دهد (۱۳۷۱). «عاصی» موضوعات اصلی در پردازش زبان طبیعی را معرفی می‌کند که مشتمل است بر حوزه خط (که خود شامل غلط‌یابی املائی و بازشناسی خودکار متن^۱ است) و حوزه زبان (که خود مستلزم بهره‌گیری از حوزه‌های آواشناسی، ساخت‌واژه، نحو و معناشناسی است) (۱۳۸۳). «لازار»^۲ با رویکردی زبانشناختی و نگاهی موشکافانه، نکته‌سنج و بی‌طرفانه و با پرداختن به مسائلی که کمتر به آن‌ها توجه شده است، خصوصیات آوائی، صرفی و نحوی و نیز چگونگی شکل‌گیری واژه‌ها را در زبان فارسی معاصر تبیین می‌کند و گاه نیز به کشف ویژگی‌هایی در عناصر سازنده زبان و روابط میان آن‌ها دست می‌یابد

1. Optical Character Recognition (OCR)

2. Lazar

(۱۹۵۷) که طرح آن‌ها در نوشته‌های دستورنویسان فارسی‌زبان پیش از او سابقه نداشته است. نمونه‌هایی که از زبان فارسی معاصر در این دستور زبان ارائه شده، بر اساس آثار نویسندگانی چون «صادق هدایت»، «صادق چوبک»، «علی دشتی»، «جمالزاده» و دیگران و نیز نوشته‌های روزنامه‌ها و مجلات و سرانجام شنیده‌های مؤلف هنگام اقامتش در ایران بوده است.

بررسی پژوهش‌های انجام‌شده در زمینه برجسب‌گذاری اجزای کلام نشان می‌دهد که مسئله برجسب‌دهی در پردازش متون حائز اهمیت است. برخی از پژوهش‌های اشاره‌شده در بخش ۲-۱، استفاده از ابزارهایی مانند «مدل مارکوفی پنهان»، تلفیقی از تجزیه و تحلیلی ساخت‌واژی و «مدل مارکوفی پنهان» و شبکه‌های پرسپترون چندلایه تکرار شونده را پیشنهاد می‌کنند و برخی دیگر روش‌هایی را معرفی می‌کنند که مستلزم بهره‌گرفتن از ابزارهای تحلیل‌کننده ساخت‌واژی و نحوی است. در این راستا و به‌منظور ارائه بهینه چنین سیستم‌های برجسب‌دهی به اجزای کلام، ناگزیریم به بررسی ساخت‌واژی اجزای کلام در فارسی پردازیم. بنابراین، در بخش ۲-۲، نمونه‌هایی از مطالعات انجام‌شده در این زمینه ارائه شده است.

۲-۳. بررسی ساخت‌واژی اسم‌ها و صفت‌ها در فارسی

زبان فارسی دارای شمار نسبتاً زیادی از پسوندها برای اشتقاق در مقوله اسم است. اسم در فارسی با نشانه‌های تصریفی خاصی می‌تواند همراه شود. در زبان فارسی معاصر، بسیاری از این پسوندها چندان زایا نیستند. شکل ۱، حاصل جمع‌بندی نشانه‌های تصریفی و اشتقاقی مربوط به اسم و مستخرج از مجموعه ۱۲ مقاله «دکتر علی اشرف صادقی» تحت عنوان «شیوه‌ها و امکانات واژه‌سازی در زبان فارسی معاصر ۱ تا ۱۲: ۱۳۷۰-۱۳۷۲» (در بحث اشتقاق) است؛ همچنین در بخش تصریف از کتاب دستور زبان فارسی معاصر «لازار» (۱۹۵۷) و مقاله مشخصه‌های تصریفی در زبان فارسی امروز «قطره» (۱۳۸۶) استفاده شده است. این فهرست تنها پسوندهایی را عرضه می‌دارد که مهم‌تر و زایاتر شمرده می‌شوند؛ اگرچه برخی از پسوندهایی که در فارسی معاصر چندان زایا نیستند نیز بررسی شده است.

نیز در مورد آن‌ها به کار برد. هم‌نگاره‌ها در زبان‌های مختلف یکی از مهم‌ترین لایه‌های ابهام را در سطح متن ایجاد می‌کنند و منجر به چالش‌هایی در پردازش متن می‌شوند. در برخی زبان‌ها، مانند زبان فارسی، به دلیل ساختار نوشتاری خاص آن، تعداد هم‌نگاره‌ها نسبت به دیگر زبان‌ها بیشتر است. تعداد زیادی از هم‌نگاره‌ها از ساختار زبان ناشی می‌شوند، به عبارت دیگر، هم‌نگاره‌ها در هر زبان به آن زبان خاص وابسته‌اند و هرچه ساختار و اثر اشتقاقی و تصریفی زبان پیچیده‌تر باشد، می‌توان انتظار داشت هم‌نگاره‌های بیشتری وجود داشته باشد و در نتیجه، ابهام‌زدایی از آن‌ها نیز مشکل‌تر است. با این اوصاف می‌توان دریافت که بازنمایی هم‌نگاره‌ها نیز تا حد زیادی وابسته به ساختار زبان باشد (محسنی ۱۳۸۷). در این بخش به بررسی نظام نوشتاری زبان فارسی پرداخته می‌شود تا از این رهگذر بتوان به شناسایی انواع هم‌نگاره‌ها در زبان فارسی پرداخت. سپس، انواع هم‌نگاره‌ها از جمله هم‌نگاره‌های اسمی در زبان فارسی مورد مطالعه قرار خواهد گرفت.

۳-۱. برخی از مسائل مربوط به نظام نوشتاری زبان فارسی که در پردازش متن حائز اهمیت است

در نظام نوشتاری زبان فارسی می‌توان گفت اطلاعات واجی، که از طریق صورت نوشتاری کلمات به دست می‌آید، کامل نیست، زیرا واژه‌های کوتاه معمولاً نمود نوشتاری ندارند (Sproat 2000, 163)، نقل در علایی و بی‌جن‌خان (۱۳۹۲). تنها تعداد انگشت‌شماری نگاره وجود دارد که برای نشان دادن واژه‌های کوتاه در خط به کار می‌روند؛ از جمله نگاره «ه»، که برای نشان دادن واژه کوتاه /e/ یا /a/ به کار می‌رود؛ مانند «به» و «نه»؛ و حرف «و» که می‌تواند نمود نوشتاری واژه کوتاه /o/ باشد؛ مانند «تو». فقدان واژه‌های کوتاه در نظام نوشتاری فارسی باعث ایجاد ابهام می‌شود. به عنوان مثال، هم‌نگاره‌هایی مانند «مرد» /mard/ و «مرد» /mord/ در اثر فقدان واژه‌های کوتاه در نظام نوشتاری فارسی ایجاد شده‌اند (علایی و بی‌جن‌خان ۱۳۹۲). همچنین، استفاده نکردن از برخی از علائم زیروبری در خط فارسی نیز می‌تواند بعضاً هم‌نگاره‌هایی را ایجاد کند. برخی از این علائم شامل فتحه، کسره، ضمه، سکون، تنوین، مد، تشدید و همزه است.

۳-۲. طبقه‌بندی هم‌نگاره‌ها در زبان فارسی

هم‌نگارگی در زبان فارسی ناشی از بازنمایی واجی و صرفی عناصر زبانی در خط فارسی است. به این صورت که رابطه‌ای چند-به-چند و در برخی موارد، غیرنظام‌مند میان عناصر واجی و صرفی زبان فارسی و تظاهر نوشتاری آن‌ها در خط فارسی دیده می‌شود (بی‌جن‌خان و مرادزاده ۱۳۸۳، نقل در محسنی ۱۳۸۷). هم‌نگاره‌ها در زبان فارسی می‌توانند واژگانی باشند. منظور از هم‌نگاره‌های واژگانی، هم‌نگاره‌هایی هستند که در فرهنگ‌های لغات به صورت مدخل‌های مجزا ذکر شده‌اند و حضورشان وابسته به بافت نحوی و اضافه‌شدن تکواژهای اشتقاقی و تصریفی است. هم‌نگاره‌های موجود در هر دو دسته، می‌توانند هم‌آوا (مانند: «دوش» و «دوش»، «شانه»، «شانه» و «شیر» و «شیر») یا غیرهم‌آوا (مانند: «کند» /kond/، «کند» /kanad/، و «کند» /konad/، «ملک» /malak/، «ملک» /melk/ و «ملک» /molk/ و «ملک» /malek/، «حسن» /hasan/ و «حسن» /hosn/ باشند. هم‌نگاره‌های فارسی را با توجه به علل پیدایش آن‌ها، می‌توان در طبقات مختلفی قرار داد که در زیر به برخی از آن‌ها اشاره می‌شود:

۳-۲-۱. هم‌نگاره‌هایی که در اثر عدم نمایش علائم زیروزبری در خط فارسی ایجاد شده‌اند؛ مانند: «مرد» /mard/ و «مرد» /mord/، «کند» /kanad/، «کند» /kond/ و «کند» /konad/.

۳-۲-۲. هم‌نگاره‌هایی که در اثر عدم تناظر یک-به-یک میان واج‌ها و نگاره‌ها در فارسی ایجاد شده‌اند؛ مانند: «رود» /rud/ و «رود» /ravad/، «شوم» /um/ و «شوم» /avam/.

۳-۲-۳. هم‌نگاره‌هایی که در اثر یکسانی تظاهر واجی و نوشتاری تکواژها به وجود می‌آیند. در زبان فارسی ممکن است چندین تکواژ متفاوت، نمود نوشتاری یکسانی داشته باشند که در زیر به برخی از آن‌ها اشاره شده است (بی‌جن‌خان و مرادزاده ۱۳۸۳، نقل در محسنی ۱۳۸۷):

۳-۲-۴. یکسان‌بودن نمود نوشتاری تکواژ یای نکره، یای اسم‌ساز (اسم مکان، اسمی که دال بر شغل یا محافظت و دارندگی است، اسم معنا یا اشیاء، تصغیر و تحجیب، اسم مصدر یا حاصل مصدر)، شناسهٔ دوم شخص مفرد، یای صفت‌ساز (صفت فاعلی و مفعولی، صفتی که دال بر نسبت است) و یای متصل به گروه اسمی؛ مانند:

الف. یای نکره: استادی که به موقع سر کلاس حاضر می‌شود.

ب. یای اسم‌ساز: وی در سال ۱۳۸۸ به درجه استادی رسید.

پ. یای شناسه دوم شخص مفرد: تو دیگر خود استادی.

ت. یای نسبت: مشکلات جوانی.

ث. یای متصل به گروه اسمی: کارگر ماهری که ما می‌شناسیم، ...

۳-۲-۵. یکسانی تکواژهای ضمیر متصل سوم شخص مفرد با تکواژ اسم‌ساز؛ مانند:

الف. ضمیر متصل سوم شخص مفرد: به رویش خندیدم. /ruja/

ب. تکواژ اسم‌ساز: رویش گل‌ها را بین. /ru'je/

هم‌نگاره‌ها را در دو سطح می‌توان بررسی کرد: سطح تکواژی و سطح کلمه. منظور از سطح تکواژی، تکواژهایی هستند که فارغ از پایه‌ای که به آن متصل می‌شود، با تکواژهای دیگر هم‌نگاره هستند؛ به‌عنوان مثال، تکواژ دال بر شمار در اسم‌ها، تکواژی که اسم مکان می‌سازد و تکواژی که صفت فاعلی و مفعولی می‌سازد، همه دارای نمود نوشتاری یکسان («-ی») هستند. برخی از مواردی را که هم‌نگارگی را در سطح تکواژ نشان می‌دهد، میتوان به‌صورت زیر دسته‌بندی کرد:

◇ پسوند «-ین» که هم شمار را در اسم‌ها نشان می‌دهد و هم می‌تواند در صفت‌ها، دال بر نسبت باشد.

◇ پسوند «-ه / -ه» که می‌تواند نشانه معرفه در گفتار، سازنده اسم اشیا و اسم معنا، تأنیث در اسم‌ها، دال بر دارندگی و اتصاف در صفت‌ها باشد یا صفت فاعلی و مفعولی بسازد.

◇ پسوند «-ش» که می‌تواند هم نشانه معرفه در گفتار باشد و هم اسم بسازد.

◇ پسوند «-گاه» که اسم مکان و قید زمان می‌سازد.

◇ پسوند «-ئیه» که اسم مکان می‌سازد یا دال بر تأنیث است.

◇ پسوند «-ک» که تصغیر و تحیب را در اسم‌ها نشان می‌دهد و یا اسم اشیا و معنا می‌سازد.

◇ پسوند «-ا» که هم می‌تواند دال بر تأنیث در اسم‌ها باشد و هم مصدر، اسم مصدر، صفت فاعلی و مفعولی بسازد.

◇ پسوند «-ئینه» که در صفت‌ها، دال بر نسبت است و در اسم‌ها، اسم معنا و اشیا

می‌سازد.

نوع دوم هم‌نگارگی، هم‌نگارگی در سطح کلمه است. به این صورت که، همان‌گونه که ذکر شد، کلماتی هم‌نگاره خوانده می‌شوند که صورت نوشتاری یکسان، اما معنا (و گاهی، تلفظ) متفاوتی دارند.

۴. روش پژوهش

پس از بررسی انواع هم‌نگاره‌ها در زبان فارسی به منظور شناخت آن‌ها در متون گوناگون، لازم است در پیکره‌های متنی گوناگون زبان فارسی به جست‌وجوی هم‌نگاره‌ها پرداخته شود و فهرستی از آن‌ها تهیه شود تا بتوان در تحقیقات بعدی، هم‌نگاره‌هایی را که فراوانی بالاتری در پیکره‌ها دارند، شناسایی کرده و اقداماتی در جهت رفع ابهام از آن‌ها در پیکره‌ها انجام داد. پیکره‌های استفاده‌شده برای این منظور، پیکره متنی زبان فارسی، پایگاه دادگان زبان فارسی و پیکره وابستگی نحوی زبان فارسی هستند.

در مورد برخی پیکره‌ها که فایل آن‌ها موجود است، باید از روش جست‌وجوی ماشینی جهت تهیه فهرست هم‌نگاره‌ها استفاده می‌شد. نگارنده در این بخش از یک متخصص برنامه‌نویس کمک گرفته است. در سایر موارد از برچسب‌دهی دستی استفاده شده است. روش جست‌وجوی ماشینی هم‌نگاره‌ها در پیکره وابستگی نحوی زبان فارسی و پیکره متنی زبان فارسی، که فایل آن‌ها موجود است، به صورت زیر توضیح داده می‌شود: جست‌وجو برای هم‌نگاره‌ها با یک دستور در زبان برنامه‌نویسی «پایتون»^۱ انجام شده است. به این صورت که کلمات موجود در پیکره، مانند پیکره متنی زبان فارسی، و برچسب‌هایی را که به آن کلمات تعلق گرفته است، بررسی می‌کند و کلماتی را که دارای بیش از یک برچسب باشند - به همراه فراوانی هر یک از برچسب‌ها - گزارش می‌کند. فرض کنید یک مجموعه داشته باشیم به نام M. برنامه با خواندن اولین کلمه در پیکره، بررسی می‌کند آیا این کلمه قبلاً در مجموعه M موجود است یا نه. اگر موجود نباشد، کلمه را به همراه برچسب آن به مجموعه M اضافه می‌کند، اگر موجود باشد، بررسی می‌کند که آیا برچسب کلمه‌ای که الان مشاهده کرده است، مشابه همان برچسب قبلی است؟ اگر بود، یک واحد به فراوانی برچسب قبلی اضافه می‌کند و اگر برچسب جدید

1. Python

بود، آن را در کنار برجسب قبلی قرار می‌دهد. عملکرد برنامه به زبان ساده - به عنوان مثال، برای پیکره «بی جن خان»^۱، که پیکره مبسوط تر آن، پیکره متنی زبان فارسی است - به صورت زیر است:

- ◇ هر سطر در پیکره متنی زبان فارسی شامل یک واژه و برجسب اجزای کلام^۲ آن است.
- ◇ یک مجموعه به اسم دیکشنری^۳ داریم که در واقع، یک ساختار داده در زبان برنامه نویسی «پایتون» است. ساختار آن بدین شکل است که هر واژه به علاوه برجسب‌های (POS) آن در مجموعه ذخیره می‌شود. نوع داده دیکشنری در زبان برنامه نویسی «پایتون»، به صورت فهرستی از کلیدها (key) و ارزش‌ها (value) است. در دیکشنری‌های «پایتون»، ترتیبی وجود ندارد و با دو مفهوم کلید (key) و ارزش (value) سروکار داریم (oruji.org). در مورد پژوهش حاضر، واژه‌ها همان key و برجسب اجزای کلام مربوط به هر واژه value آن در نظر گرفته می‌شوند.
- ◇ برنامه هر سطر از پیکره را می‌خواند. اگر واژه قید شده در آن سطر قبلاً در دیکشنری وجود نداشت، آن را به همراه POS آن به دیکشنری اضافه می‌کند. اگر قبلاً در دیکشنری وجود داشت، بررسی می‌کند که آیا POS قید شده در این سطر قبلاً در دیکشنری برای این واژه درج شده است یا نه. اگر درج نشده بود، آن را به لیست POS های این واژه در دیکشنری اضافه می‌کند.
- ◇ در نهایت، مجموعه دیکشنری بررسی می‌شود و واژه‌هایی که بیش از یک POS برای آن‌ها ثبت شده است، به عنوان خروجی برنامه ارائه می‌شوند.

۱. پیکره متنی استاندارد زبان فارسی / پیکره «بی جن خان» که در پژوهشکده پردازش هوشمند علائم تهیه شده است، مجموعه‌ای از متون نوشتاری و گفتاری زبان فارسی به صورت رسمی است که از منابع واقعی همچون روزنامه‌ها، سایت‌ها و مستندات از قبل تایپ شده، جمع‌آوری شده، تصحیح کرده و برجسب خورده است. حجم این دادگان حدوداً ۱۰۰ میلیون کلمه است و از منابع مختلف تهیه گردیده است و دارای تنوع بسیار زیادی است (Bijankhan et al. 2011).

2. Part Of Speech (POS)

3. dictionary

از پایگاه دادگان زبان فارسی^۱ در تحقیق حاضر به این صورت استفاده شده است که فهرستی از کلمات پربسامد در پیکره در اختیار نگارنده قرار گرفت و از آنجا که فهرست مذکور فاقد هرگونه برچسبی است، نگارنده، خود به صورت دستی کلماتی را که می‌توانند بیش از یک برچسب داشته باشند، برچسب‌گذاری کرده است. برای این منظور، نگارنده هر کلمه موجود در فهرست مذکور را در جملات فرضی، به عنوان بافت، در نظر گرفته است تا امکان داشتن بیش از یک برچسب برای هر کلمه را با توجه به بافت بررسی کند. سپس، تنها کلماتی را که در بافت بیش از یک برچسب می‌توانسته‌اند داشته باشند، برچسب‌گذاری کرده است و هر تعداد برچسب ممکن برای هر یک از کلمات مذکور را کنار آن‌ها مشخص کرده است. وی برای کنترل و تأیید صحت روایی و درستی تشخیص برچسب‌ها، از یک متخصص زبان و ادبیات فارسی کمک گرفته است. جمع‌آوری داده‌ها در این مرحله دو کاربرد مهم داشت: ۱. نشان‌دهنده تعداد بسیار زیاد هم‌نگاره‌ها در پیکره‌های متنی فارسی که خود نشان می‌دهد در صورت رفع ابهام نکردن از هم‌نگاره‌ها در متون گوناگون، پردازش متن با مسائل و مشکلاتی روبه‌رو خواهد شد. ۲. فهرست تهیه‌شده از هم‌نگاره‌ها در پیکره‌های گوناگون می‌تواند در پژوهش‌های دیگر مورد استفاده قرار گیرد. به این صورت که می‌توان هم‌نگاره‌هایی را که فراوانی بالاتری دارند، در بافت نحوی مورد مطالعه قرار داد تا با توجه به ساخت نحوی که در آن به کار رفته‌اند، برچسب‌های درست به آن‌ها اختصاص داده شود.

قسمتی از فهرست تهیه‌شده در جدول ۱، آورده شده است:

۱. پایگاه دادگان زبان فارسی، مجموعه‌ای نرم‌افزاری برای ذخیره، پردازش و ارائه داده‌های زبانی فارسی است. این پایگاه دربرگیرنده پیکره‌های گوناگونی از زبان فارسی است که با وجود حجم عظیم و با گستردگی و گوناگونی‌های بسیار، دارای ساختاری بسامان و منطقی است و امکان هرگونه جست‌وجو و دستیابی سریع به آگاهی‌های مورد نیاز را در هر زمان فراهم آورده است. پیکره‌های این پایگاه می‌توانند همواره روزآمد شوند. پایگاه دادگان زبان فارسی فراگیر و متنوع است. در واقع، فراتر از یک یا چند پیکره خاص است و کاربران بر پایه نیاز و هدف پژوهشی خود می‌توانند پیکره مناسب را از آن برگزینند.

جدول ۱. قسمتی از فهرست هم‌نگاره‌های تهیه‌شده از پیکره‌های متنی فارسی

فرآوانی	کلمات	برچسب‌های هم‌نگاره‌ها
۱۰۳۲۳۶۱	به	P, ADV, N
۱۳۲۱۵۷	کرد	N, V, ADJ
۱۱۱۲۸۶	تا	P, N(NO-V)
۱۰۹۶۰۹	بر	P, N(NO-V)
۱۰۹۳۱۶	کند	V, V, ADV
۷۵۲۵۲	کرده	V, ADJ+V-PRE, ADJ-INO, N
۲۹۵۱۴	اجتماعی	N-IN, ADJ
۲۹۳۶۶	گرفته	V, ADJ, ADJ-INO
۱۷۶۷۰	شهرداری	N-IN, N, N+V-PRE, ADJ
۱۸۱۸	داوری	N-IN, N, N+V-PRE, ADJ

راهنمای برچسب‌ها:

=ADV قید	=P حرف اضافه
=N (NO-V) جزء غیر فعلی فعل مرکب	=ADJ صفت
=V-PRE فعل اسنادی	=V فعل
=N-IN اسم نکره	=ADJ-INO صفت مفعولی
	=N اسم

۵. نتیجه‌گیری

یکی از بخش‌های مهم سامانه‌های برچسب‌دهی به کلمات، برچسب‌گذاری اسم‌ها و صفت‌ها در فارسی است. این سامانه‌ها جهت برچسب‌گذاری اجزای واژگانی کلام مانند: اسم، فعل، قید، صفت و ... تهیه می‌شوند. برچسب‌گذاری اجزای واژگانی کلام، در بسیاری از حوزه‌های پیشرفته‌تر پردازش زبان طبیعی از جمله ترجمه ماشینی، خطایاب، تبدیل متن به گفتار، بازیابی اطلاعات، موتورهای جست‌وجو و کمک به مدل‌های آماری اقدامی کاربردی است. با توجه به تعداد بسیار زیاد اسم‌ها در زبان فارسی و هم‌نگاره‌بودن

بسیاری از آن‌ها با صفت‌ها، به راحتی نمی‌توان اسم‌ها را در یک پیکره متنی فارسی برجسب‌گذاری کرد. در تهیه سامانه‌های برجسب‌گذاری برای اسم‌ها و صفت‌ها در زبان فارسی، بررسی ساخت‌واژی اسم‌ها و صفت‌ها از دو نظر حائز اهمیت است: ۱. اگر کلمه‌ای قبلاً در پیکره آموزشی ظاهر نشده باشد، نمی‌توان از پیکره آموزشی اطلاعات دقیقی راجع به آن کلمه به دست آورد و حتی از توزیع کلمات در پیکره نیز نمی‌توان استفاده کرد، زیرا توزیع کلمات ناشناخته کاملاً متفاوت از کلمات شناخته شده در متن است. بنابراین، اگر در یک پیکره متنی فارسی، کلمه‌ای در واژگان حضور نداشته باشد (کلمه خارج از واژگان)، نمی‌توان برجسب‌هایی را که کلمه به آن‌ها منتسب می‌شود، بازبایی کرد. در این صورت، برجسب کلمه را تنها می‌توان با توجه به شکل کلمه یا بافتی که کلمه در آن ظاهر می‌شود، یا هر دو، حدس زد. منظور از شکل کلمه، ویژگی‌هایی مانند انواع پیشوندها و پسوندهایی است که به کلمات متصل می‌شوند. بنابراین، در وهله اول، بررسی ساخت‌واژی اسم‌ها و صفت‌ها، فارغ از بافت، ضروری به نظر می‌رسد تا از این رهگذر بتوان اسم‌ها و صفت‌ها را در صورت دارابودن شاخص‌های صرفی و اشتقاقی شناسایی و سپس برجسب‌دهی کرد. ۲. زبان فارسی ظرفیت بالایی برای ساخت هم‌نگاره‌های جدید را که از ساخت‌واژه فارسی نشأت می‌گیرند نیز دارد. بسیاری از وندها، اشتقاقی و تصریفی، نمود نوشتاری مشابهی دارند و در اتصال به ستاک، هم‌نگاره‌های گوناگون می‌سازند. به‌عنوان مثال، نمود نوشتاری وند تصریفی نکره‌ساز /i/ و وند اشتقاقی صفت‌ساز /i/ یکسان است («-ی») و کلماتی مانند: «آسمانی» و «اسلامی»، فارغ از بافت، می‌توانند هم‌برجسب «اسم» داشته باشند و هم «صفت». بنابراین، بررسی ساخت‌واژی اسم‌ها و صفت‌ها، به منظور تفکیک آن‌ها از هم ضروری به نظر می‌رسد، زیرا اکثر صفت‌ها در بافت‌های گوناگون، با صورت نوشتاری یکسان، می‌توانند برجسب «اسم» بگیرند.

در پژوهش حاضر، پس از مقدمه، ابتدا به پیشینه تحقیق پرداخته شد. سپس، ساخت‌واژه اسم‌ها و صفت‌ها در فارسی بررسی شد تا از این رهگذر بتوان اسم‌ها و صفت‌ها را در صورت دارا بودن شاخص‌های صرفی و اشتقاقی، شناسایی و سپس برجسب‌دهی کرد. بخش چهارم، مطالعه هم‌نگاره‌ها در فارسی بود که به بررسی نظام نوشتاری زبان فارسی پرداخته شد تا از این رهگذر بتوان به شناسایی انواع هم‌نگاره‌ها در زبان فارسی پرداخت. سپس، انواع هم‌نگاره‌ها در زبان فارسی مورد مطالعه قرار گرفت؛ از جمله هم‌نگاره‌های اسمی و صفتی. در بخش پنجم به روش جمع‌آوری داده‌ها از پیکره‌ها

پرداخته شد. پیکره‌های استفاده‌شده برای این منظور، شامل پیکره متنی زبان فارسی، پایگاه دادگان زبان فارسی و پیکره وابستگی نحوی زبان فارسی است و جست‌وجو به دو روش جست‌وجوی ماشینی و دستی انجام گرفت و فهرست مبسوطی از هم‌نگاره‌های به‌دست آمده از پیکره‌های مذکور تهیه شد. نمونه‌ای از این فهرست در جدول ۱، آورده شده است. بررسی کلی هم‌نگاره‌ها در پیکره‌های مورد مطالعه نشان می‌دهد که تعداد هم‌نگاره‌ها در پیکره‌ها قابل توجه است و حتی می‌توان گفت بیشتر هم‌نگاره‌ها کلماتی هستند که فراوانی آن‌ها در پیکره‌ها زیاد است. اکثر این هم‌نگاره‌ها به نظر می‌رسد در اثر یکسان‌بودن نمود نوشتاری تکواژ یای نکره، یای اسم‌ساز (اسم مکان، اسمی که دال بر شغل یا محافظت و دارندگی است، اسم معنا یا اشیا، تصغیر و تحییب، اسم مصدر یا حاصل مصدر)، شناسهٔ دوم‌شخص مفرد و یای صفت‌ساز (صفت فاعلی و مفعولی، صفتی که دال بر نسبت است) و یای متصل به گروه اسمی باشند. مانند:

الف. یای نکره: کشاورزی که محصول خود را به‌موقع درو می‌کند.

ب. یای اسم‌ساز: وی چند سال کشاورزی را دنبال می‌کرد.

پ. یای شناسه دوم‌شخص مفرد: تو خود کشاورزی.

و یا در اثر یکسانی تکواژهای ضمیر متصل سوم‌شخص مفرد با تکواژ اسم‌ساز به‌وجود آمده‌اند؛ مانند:

الف. ضمیر متصل سوم‌شخص مفرد: به رویش خندیدم.

ب. تکواژ اسم‌ساز: رویش گل‌ها را بین.

جمع‌آوری داده‌ها در این مرحله دو کاربرد مهم داشت: ۱. نشان‌دهندهٔ تعداد بسیار زیاد هم‌نگاره‌ها در پیکره‌های متنی فارسی که خود نشان می‌دهد در صورت رفع ابهام‌نکردن از هم‌نگاره‌ها در متون گوناگون، پردازش متن با مسائل و مشکلاتی روبه‌رو خواهد شد. ۲. فهرست تهیه‌شده از هم‌نگاره‌ها می‌تواند در پژوهش‌های گوناگون مورد استفاده قرار گیرد. به این صورت که هم‌نگاره‌هایی که فراوانی بالاتری دارند در بافت نحوی مورد مطالعه قرار گیرند و با توجه به ساخت نحوی که در آن به کار رفته‌اند، برچسب‌های درست به آن‌ها تخصیص داده شود.

۶. پیشنهاد برای پژوهش‌های بیشتر

در پژوهش حاضر به بررسی ساخت‌واژی اسم‌ها و صفت‌ها در فارسی پرداخته شده است. همچنین، پس از بررسی تعداد و انواع هم‌نگاره‌ها در پیکره‌های متنی فارسی، این نتیجه به‌دست آمد که تعداد زیادی از هم‌نگاره‌ها در اثر یکسان‌بودن نمود نوشتاری/املائی و ندهای اشتقاقی و تصریفی، به‌ویژه وندهایی با نمود نوشتاری «ی-ی» ایجاد شده‌اند. پیشنهاد می‌شود این هم‌نگاره‌ها که فراوانی بالایی در پیکره‌های متنی فارسی دارند، در پیکره‌ها در بافت نحوی مورد مطالعه قرار گیرند. ضمناً از آنجا که در فارسی هم‌نگاره‌های دیگری نیز وجود دارد که در اثر عدم نمایش واژه‌های کوتاه در خط فارسی یا نبود تناظر یک-به-یک میان واج‌ها و نگاره‌ها ایجاد می‌شوند، پیشنهاد می‌شود چنین هم‌نگاره‌هایی نیز در پیکره‌های متنی فارسی مورد مطالعه قرار گیرند.

فهرست منابع

- آخشیبک، سمیه، و رحمت‌الله فتاحی. ۱۳۹۱. تحلیل چالش‌های پیوسته‌نویسی و جدانویسی واژگان فارسی در ذخیره و بازیابی اطلاعات در پایگاه‌های اطلاعاتی. «فصلنامه کتابداری و اطلاع‌رسانی» ۱۵ (۳): ۹-۳۰.
- صادقی، علی اشرف ۱۳۷۰. شیوه‌ها و امکانات واژه‌سازی در زبان فارسی معاصر (۱). نشر دانش ۶۴: ۱۲-۱۸.
- _____. ۱۳۷۰. «شیوه‌ها و امکانات واژه‌سازی در زبان فارسی معاصر (۲)». نشر دانش ۶۵: ۶-۱۲.
- _____. ۱۳۷۰. «شیوه‌ها و امکانات واژه‌سازی در زبان فارسی معاصر (۳)». نشر دانش ۶۷: ۲۸-۳۳.
- _____. ۱۳۷۱. «شیوه‌ها و امکانات واژه‌سازی در زبان فارسی معاصر (۴)». نشر دانش ۶۹: ۲۱-۲۵.
- _____. ۱۳۷۱. «شیوه‌ها و امکانات واژه‌سازی در زبان فارسی معاصر (۵)». نشر دانش ۷۰: ۳۹-۴۵.
- _____. ۱۳۷۱. «شیوه‌ها و امکانات واژه‌سازی در زبان فارسی معاصر (۶)». نشر دانش ۷۱: ۱۵-۱۹.
- _____. ۱۳۷۱. «شیوه‌ها و امکانات واژه‌سازی در زبان فارسی معاصر (۷)». نشر دانش ۷۲: ۱۹-۲۳.
- _____. ۱۳۷۱. «شیوه‌ها و امکانات واژه‌سازی در زبان فارسی معاصر (۸)». تهران: نشر دانش. شماره ۷۴: ۹۸-۱۰۵.
- _____. ۱۳۷۲. «شیوه‌ها و امکانات واژه‌سازی در زبان فارسی معاصر (۹)». تهران: نشر دانش. شماره ۷۵: ۹-۱۵.
- _____. ۱۳۷۲. «شیوه‌ها و امکانات واژه‌سازی در زبان فارسی معاصر (۱۰)». تهران: نشر دانش. شماره ۷۶: ۲۳-۱۵.
- _____. ۱۳۷۲. «شیوه‌ها و امکانات واژه‌سازی در زبان فارسی معاصر (۱۱)». تهران: نشر دانش. شماره ۷۷: ۲۵-۲۱.
- _____. ۱۳۷۲. «شیوه‌ها و امکانات واژه‌سازی در زبان فارسی معاصر (۱۲)». تهران: نشر دانش. شماره ۷۹ و

۸۰: ۱۲-۱۵.

- عاصی، مصطفی. ۱۳۸۳. پردازش دستوری زبان فارسی با رایانه. *ویژنامه نامه فرهنگستان (دستور)* ۱: ۲۹-۵۱.
- علایی، الهام، و محمود بی‌جن‌خان. ۱۳۹۲. «عمق خط فارسی». پژوهش‌های زبانی. *مجله سابق دانشکده ادبیات و علوم انسانی دانشگاه تهران* ۴(۱): ۱-۱۹.
- قطره، فریبا. ۱۳۸۶. مشخصه‌های تصریفی در زبان فارسی امروز. *دستور* ۳: ۵۲-۸۱.
- کشانی، خسرو. ۱۳۷۱. *اشتقاق پسوندی در زبان فارسی امروز*. تهران: مرکز نشر دانشگاهی.
- لازار، ژ. ۱۹۵۷. *دستور زبان فارسی معاصر*. ترجمه مهستی بحرینی و توضیحات و حواشی هرمز میلانیان. ۱۳۸۹. تهران: انتشارات هرمس.
- محسنی، مهدی. ۱۳۸۷. سیستم برچسب‌گذاری و ابهام‌زدایی خودکار اجزای کلام برای پیکره متنی زبان فارسی. تهران: دانشگاه علم و صنعت. دانشکده مهندسی کامپیوتر.
- _____. و بهروز مینایی بیدگلی. ۱۳۸۸. سیستم برچسب‌گذاری اجزای واژگانی کلام در زبان فارسی. دو فصلنامه پردازش علائم و داده‌ها ۲. پیاپی ۱۲: ۱۳-۲۶.
- مرادی، مهدی و بهرام وزیرنژاد. ۱۳۹۱. ساخت پیکره‌های نشانه‌گذاری شده با رویکرد وب به‌عنوان پیکره. *مجموعه مقالات دومین هم‌اندیشی زبان‌شناسی رایانشی*. تهران: انجمن زبان‌شناسی ایران.
- مسعودی، بابک، سعید قوچانی، و اعظم استاجی. ۱۳۸۹. یک روش بیزی برای رفع ابهام معنایی کلمات در زبان فارسی با تأکید بر ویژگی‌های محلی کلمه. *اولین کنفرانس ملی محاسبات نرم و فناوری اطلاعات*. ماهشهر: دانشگاه آزاد اسلامی، واحد ماهشهر.

References

- Assi, M. and M. Haji Abdohosseini. 2000. Grammatical tagging of a Persian corpus". *International journal of corpus linguistics* 5 69-81 :1(.
- Bijankhan, M., J. Sheykhzadegan, M. Bahrani, & M. Ghayoomi. 2011. Lessons from Building a Persian Written Corpus: Peykare. *Language Resources and Evaluation* 45 (2): 143-164.
- Elhaj, M. 2009. Statistical part of speech tagger for traditional Arabic text". *Journal of computer science* 5 (11): 794-800.
- Frodl, M. 2013. *Part of speech tagging using Neural Networks*. Czech Republic: Masaryk University, faculty of informatics. <https://github.com/sobhe/hazm> <https://oruij.org/python/dictionary>
- Hearst, M. A. 1991. Noun homograph disambiguation using local context in large text corpora. *Proceedings of the 17th annual conference of the University of Waterloo Centre for the new OED and text research*. Oxford.
- Hu, Y., I. Matveeva, J. Goldsmith, and C. Sprague. 2005. Using morphology and syntax together in unsupervised learning". In *Proceeding of the workshop on Psychocomputational Models of Human Language Acquisition*. PMHLA 05. 20-27. Stroudsburg: PA, USA. Association for computational linguistics.
- Megerdooomian, K. 2000. Unification-based Persian morphology. In *Proceedings of CICLing, Centro de investigacion en computacion- IPN, Mexico* 311-318 :.
- _____. 2004. Developing a Persian part-of-speech tagger". *Proceedings of the 1st workshop on Persian*

language and computer. 99-105. University of California. San Diego.

Merialdo, B. 1994. Tagging English text with a probabilistic mode". *Computational linguistics* 20: 155-171.

Schuetze, Hinrich. 1995. Distributional Part-of-Speech Tagging, From Texts to Tags: Issues in Multilingual Language Analysis. Online Proceedings of the ACL SIDGAT Workshop. On the Internet at <http://xxx.lanl.gov/find/cmp-lg>

الهام علایی ابوذر

دانش‌آموخته دکتري تخصصی در رشته زبان‌شناسی همگانی در سال ۱۳۹۲ از دانشگاه تهران و کارشناسی ارشد زبان‌شناسی همگانی در سال ۱۳۸۶ از همان دانشگاه و کارشناسی زبان و ادبیات انگلیسی در سال ۱۳۸۳ از دانشگاه گیلان است. وی همکاری خود را با پژوهشگاه علوم و فناوری اطلاعات ایران (ایرانداک) از سال ۱۳۹۴، به‌عنوان عضو هیئت علمی آغاز نموده است و هم‌اکنون استادیار پژوهشی این پژوهشگاه است. بیشتر فعالیت‌های وی تاکنون در حوزه زبان‌شناسی رایانه‌ای (سیستم‌های تبدیل متن به گفتار، نظام‌های نوشتاری، یادگیری ماشینی و پردازش زبانی)، زبان‌شناسی نظری و پیکره‌ای بوده است و در حوزه‌های ذکر شده با پژوهشگاه علوم و فناوری اطلاعات ایران همکاری می‌نماید.

