پارس مورف: تحلیلگر صرفی زبان فارسی

وحید مواجی ۱، محرم اسلامی ۲ و بهرام وزیرنژاد ۱

ادانشگاه صنعتی شریف ۲دانشگاه زنجان

چکیدہ

در مقالهٔ حاضر می کوشیم مبنای نظری، نحوهٔ طراحی و عملکرد سامانهٔ تحلیل گر صرفی زبان فارسی با عنوان اختصاری "پارس مورف" را معرفی کنیم. پارس مورف سامانه ای مبتنی بر قواعد صرفی زبان فارسی است که ساخت درونی کلمات فارسی را با توجه به نظام تصریف و نظام واژهسازی زبان تجزیه و تحلیل می کند و مقولهٔ دستوری و نقش هر کدام از اجزای سازندهٔ کلمه را مشخص می کند. پارس مورف با استفاده از یک واژگان حدوداً ۴۵۰۰ واژه ای و نیز در چارچوب قواعد صرفی زبان فارسی که بر یک تحقیق جامع زبان شناختی استوار است می تواند واژه های پیچیده و نیز صورتهای ممکن تصریفی و حتی واژههای خارج از واژگان را تحلیل کند. دقت نسخهٔ اول پارس مورف حدود %۹۵ است که افزودن اطلاعات نحوی و مسائل مربوط به همنویسه ها و نیز لحاظ کردن ویژگی های خط فارسی می تواند این دقت را به %۱۰۰ نزدیک کند. از پارس مورف می توان در مطالعات محض زبانی و نیز پردازش ماشینی زبان فارسی استفاده

كلُّيدواژهها: يارس مورف، تحليل گر صرفي، زبان فارسي، واژهسازي، تصريف، اشتقاق، تركيب.

۱ مقدمه

صرف یا ساختواژه (morphology) به عنوان بخشی از زبانشناسی ناظر بر مطالعه ساخت درونی کلمات و روابط نظام مند صورت و معنا در کلمات است (2007) Booij . وقتی صحبت از ساخت (کلمه) می شود، تلویحاً می پذیریم که کلمه دارای اجزایی است و در نتیجه هرکدام از اجزا نقش و جایگاه مشخصی در ساخت کلمه دارد. در این تحقیق پس از صورت بندی ساخت درونی کلمه، سامانهای را طراحی کرده ایم که قادر است به طور خودکار ساخت درونی کلمات فارسی را تحلیل و اجزای سازنده کلمه را مشخص کند. پیش از ورود به بحث اصلی، در ادامه به طور مختصر به مبانی صرف اشاره می کنیم تا در جریان بحث پیرامون عملکرد پارس مورف ابهامی وجود نداشته باشد.

ساختواژه از دو نظام مستقل تصریف (inflection) و واژهسازی (word-formation) تشکیل می شود و در تصریف صورت-کلمهها (word-forms) مورد بحث قرار می گیرد. هر کدام از واژهها (lexemes) با توجه به جایگاهشان در جمله می توانند صورت-کلمههای متفاوت داشته باشند که در اصطلاح غیردقیق به آنها کلمه (word) نیز می گویند. صورت-کلمهها واژههای جدید نیستند، بلکه صورت تصریفی واژه به حساب می آیند. در نظام واژه سازی، واژه ها یا بسیط اند و یا پیچیده. واژه های بسیط می توانند با هم ترکیب شوند و واژه مرکب بسازند و از طرف دیگر واژه های بسیط و یا مرکب می توانند با عناصر واژه ساز (تکواژهای اشتقاقی) زبان ترکیب شوند و به ترتیب واژه های پیچیده مشتق و یا مرکب-مشتق بسازند. به نمونه های زیر توجه کنید:

```
الف) واژه بسیط: اکار ا
ب) صورت کلمه: کار، کارها، کارش، کارشه
ج) واژه مشتق: اکارگر، کارمند، کارایی، کاری(کارکن) ا
د) واژه مرکب: اکارخانهساز، کاربلد، کارآگاه، کارنامه ا
ه) واژه مرکب-مشتق: اکاردرمانی، کارخانه چی، کارمندمحور، کاردرمانگری ا
```

با توجه به نمونههای ارائهشده در مثال ۱ دستگاه صرف دو نقش اساسی برعهده دارد: تصریف و واژهسازی (طباطبائی، ۱۳۸۲). واژه بسیط (۱.الف) واژه ای است که از یک ریشه تشکیل شده است. صورت-کلمه (۱.ب) صورتهای تصریف شده یک واژه را می گویند. از ۱.ج تا ۱.ه واژههای پیچیده را می بینیم که در چارچوب امکانات واژهسازی زبان درست شده اند. واژه مشتق (۱.ج) واژه ای است که در ساختمان آن یک واژه بسیط یا واحد واژگانی بسیط و یک یا چند تکواژ اشتقاقی به کار رفته است. واژه مرکب واژه ای است که در ساختمان آن دو یا چند واژه بسیط یا واحد واژگانی واژگانی بسیط درست شده باشد. واژه مرکب-مشتق (۱.ه) واژه ای است که در ساختمان آن دو یا چند واژه بسیط یا واحد واژگانی بسیط و یک یا چند تکواژ اشتقاقی به کار رفته باشد. با توجه به مثالهای بالا می بینیم که بخشی از تکواژهای مقید واژه جدید نمی سازند (مانند ۱.ب) که به آنها تکواژهای غیراشتقاقی و تسامحاً کل آنها را تکواژهای تصریفی می گوییم. بخشی دیگر (در ۱.ج و ۱.ه) واژه جدید می سازند که به آنها تکواژهای اشتقاقی می گوییم. به عبارت دیگر در تصریف با واژه پردازی سروکار داریم، اما در اشتقاق با واژه سازی (طباطبائی، ۱۳۷۶). ستاک به عبارت دیگر در تصریف با طلاق می شود که وندهای تصریفی از آن حذف شده باشند. ستاک می تواند بسیط یا پیچیده باشد و ستاک بسیط همان ریشه است (2007). Booij (2007).

در طراحی پارس مورف، ساخت تصریفی کلمه در زبان فارسی (اسلامی و علیزاده، ۱۳۸۸) به صورت ساختمند در نظر گرفته شده است و جایگاه طبقات مختلف وندهای تصریفی در ساختمان انواع کلمات مشخص شده است. پارس مورف در تحلیل تصریفی کلمه ابتدا ستاک را شناسایی می کند و متناسب با ساخت تصریفی پیش بینی شده برای آن نوع ستاک به دنبال انواع وندهای تصریفی در جایگاه های خاص، با در نظر گرفتن صورتهای مختلف نوشتاری هرکدام از طبقات، می گردد. یافتن وندی مثلاً در جایگاه دوم پس از ستاک اسمی نشان می دهد که جایگاه اول خالی مانده است. در بخش ۲ در مقاله حاضر به طور کامل به موضوع ساخت تصریفی کلمه در زبان فارسی یرداخته ایم.

اگر ستاک پیچیده باشد، پارس مورف با اعمال قواعد واژهسازی زبان فارسی، ستاک مورد نظر را از حیث مشتق یا مرکب بودن تجزیه و تحلیل و اجزای سازنده آن را با ذکر مقوله دستوری و نقش آنها مشخص می کند. به دلیل معتبر نبودن فاصله به عنوان مرز کلمه در متون فارسی، پارس مورف در تجزیه ستاکهای پیچیده ترکیبهای بالقوه را نیز به عنوان گزینههای بعدی در اختیار ما می گذارد. مثلاً ستاک پیچیده "کارگر" در بخش اشتقاق به عنوان کلمه مشتق شناسایی می شود و یا در ترکیب گفته می شود که در فارسی ممکن است آن ترکیب "کار (N۱) + گر" (Adv) یعنی اسم + قید باشد که این دو کلمه بی فاصله در کنار هم آمدهاند. در بخش " به عملکرد پارس مورف در ریشه یابی و تعیین اجزای واژههای پیچیده خواهیم پرداخت.

در تمامی تحقیقات مربوط به پردازشهای خودکار زبانی در زبان فارسی، به خصوص در پردازش متن فارسی برای مقاصد مختلف اعم از ترجمه ماشینی، تبدیل متن به گفتار و غیره از نوعی تحلیل گر صرفی استفاده می شود؛ اگر چه در اغلب مواقع تحلیل گرهای صرفی در تحقیقات پیشین محدود، هدف محور و فاقد پشتوانه جامع زبان شناختی است. به عنوان مرور پیشینه پژوهش در ادامه تنها به مواردی اشاره می کنیم که در تحلیل ساخت درونی کلمه فارسی نگاه ساخت مند داشته اند و با یک رویکرد زبان شناختی - مهندسی به تحلیل صرفی کلمه فارسی پرداخته اند. در این خصوص ابتدا می توان به مطالعات دقیقی اشاره کرد که در پروژه شیراز (2000) Megerdoomian در طراحی تحلیل این خصوص ابتدا می توان به مطالعات دقیقی اشاره کرد که در میانه راه متوقف شد. دومین مورد از طراحی تحلیل خودکار فارسی مربوط به تحلیل گر تصریفی زبان فارسی است که در سامانه تبدیل متن به گفتار فارسی "گویا" به کار گرفته شد (اسلامی و دیگران، ۱۳۸۳) که تنها به تحلیل تصریفی کلمه محدود می شد. آماده سازی متن معیار برای زبان فارسی را از نظر صرفی دری و تحلیل کند (شمس فرد، ۱۳۸۸) که ۱۳۹۲ سامانه دیگری است که قادر است کلمات فارسی (اسلامی و دیگران، ۱۳۸۸) و دیگران، ۱۳۸۲) و قواعد صرفی که طراحان آن در نظر گرفته اند کار می کند.

اگرچه STeP۱ در تجزیه کلمه به اجزای سازنده آن تا حد زیادی موفق عمل می کند ولی مبنای علمی زبان شناختی

آن نیاز به بازنگری دارد تا از اشکالات فعلی پرهیز شود. به عنوان مثال برای کلمه "دانشگاهها" سه ریشه در نظر می گیرد، به ترتیب "دانشگاه، دانش، دان" که نشان می دهد حداقل منظور پدیدآورندگان آن از "ریشه" در مفهوم علمی آن اصطلاح نسیت. یا در کلمه "دانشگاههایمان" علاوه بر اختصاص سه ریشه فوق همزمان دو تحلیل زیر برای آن ارائه می شود:

. الف) اسم + علامت جمع + ى + ضمير ملكى اول شخص جمع كه در آن "دانشگاه" ريشه است، به صورت "دانشگاه + ها + ى + مان".

ب) اسم + گاه + علامت جمع + ی + ضمیر ملکی اول شخص جمع که در آن "دانش" ریشه است، به صورت "دانش + گاه + ها + ی + مان".

چنانچه میبینیم هیچ مبنای علمی برای ریشهبودن "دانشگاه" یا "دانش" در مثال فوق وجود ندارد، از طرف دیگر اجزای تجزیه شده لزوماً واحدهای نظام واژهسازی و نظام تصریف نیستند. به عنوان مثال "ی" در مثال بالا چه عنوان و نقش زبانی دارد؟ جز اینکه به خاطر شرایط واژ-واجی در زنجیره واجی کلمه ظاهر شده است. همچنین کلماتی مانند "کفشهایمه" یا "کفشهایمند" را STeP۱ نمی تواند تجزیه کند و هیچ ریشهای برای این قبیل تصریفها و صورت-کلمهها و همچنین برخی ستاکهای پیچیده و بسیط خارج از واژگان پیدا نمی کند. این در حالیست که یارس مورف با دقت کامل موارد فوق را تجزیه و تحلیل می کند.

پارس مورف، تحلیل گر صرفی زبان فارسی، بر یک مبنای کاملاً علمی زبانی استوار است و در چارچوب ساختمان صرفی کلمه فارسی به تجزیه و تحلیل و تعیین نقش هر کدام از اجزا در درون کلمه می پردازد. پس از صورت بندی دقیق اطلاعات نظام تصریف و واژهسازی در زبان فارسی سعی کردیم در مرحله اجرا و طراحی سامانه پارس مورف تمامی آن اطلاعات را به شکل دقیق به کار بگیریم. در حال حاضر پارس مورف با استفاده از آخرین ویرایش واژگان زایای زبان فارسی (اسلامی و دیگران، ۱۳۸۳) که حدود ۴۵۰۰۰ واژه در آن قرار دارد و در چارچوب قواعد صرفی زبان فارسی که در اختیار دارد، با دقت بالای %۹۵ می تواند ساخت درونی کلمات فارسی را تحلیل کند. نیز می تواند با استفاده از امکانات و اطلاعات صرفی که در اختیار دارد کلمات خارج از واژگان را نیز از حیث تصریف و واژهسازی تجزیه و تحلیل کند.

مراجع

Booij, G. (2007), The Grammar of Words, Oxford University Press, 2nd ed. . 1, 2

Megerdoomian, K. (2000), "Persian Computational Morphology, A Unification–Based Approach (Memoranda in Computer and Cognitive Science)," Tech. Rep. MCCS–00–321, Computing Research Laboratory, New Mexico State University. 2

اسلامی، محرم، شریفی آتشگاه، مسعود، علیزاده لمجیری، صدیقه، و زندی، طاهره (۱۳۸۳)، "واژگان زایای زبان فارسی، " در اولین کارگاه پژوهشی زبان فارسی و رایانه، دانشگاه تهران. ۲، ۲

اسلامی، محرم و علیزاده، صدیقه (۱۳۸۸)، "ساخت تصریفی کلمه در زبان فارسی،" زبان و ادب فارسی، ۵۲.

شمس فرد، مهرنوش (۱۳۸۸)، "STeP۱" تهیه متن معیار برای زبان فارسی، "گزارش طرح تحقیقی، آزمایشگاه پردازش زبان طبیعی، دانشگاه شهید بهشتی. ۲

> طباطبائی، علاءالدین (۱۳۷۶)، فعل بسیط فارسی و واژهسازی. مرکز نشر دانشگاهی. ۲ طباطبائی، علاءالدین (۱۳۸۲)، اسم و صفت مرکب در زبان فارسی. مرکز نشر دانشگاهی. ۲