

# تنوع نگارشی در زبان فارسی و تهیه خودکار دادگان املائی از پیکره زبانی مبتنی بر وب

مرضیه صنعتی<sup>‡</sup>

ساغر شریفی<sup>†</sup>

مسعود قیومی<sup>\*</sup>

\* دانشگاه آزاد برلین، برلین، آلمان

masood.ghayoomi@fu-berlin.de

† دانشکده زبان‌های خارجی، دانشگاه آزاد اسلامی واحد کرج، البرز، ایران

saghar.sharifi@kiaiu.ac.ir

‡ پژوهشکده زبان‌شناسی، سازمان میراث فرهنگی، صنایع دستی، و گردشگری، تهران، ایران

msanaati20@gmail.com

## چکیده

در عصر اطلاعات و ارتباطات، وب جایگاه ویژه‌ای پیدا کرده‌است، چراکه با کاربران بسیار متنوعی در تعامل بوده و می‌توان از آن به‌عنوان یک منبع اطلاعاتی غنی زبانی استفاده کرد. تهیه پیکره مبتنی بر وب می‌تواند برای پردازش‌های زبانی منبع مناسبی باشد. ولی استفاده از این منبع ساده نیست. از آنجا که کاربران مختلفی موجب خلق متن در وب می‌شوند، برخورد با پدیده تنوع نگارشی اجتناب‌ناپذیر خواهد بود. در این مقاله، به بررسی این پدیده در پیکره زبانی حاصل از وب برای زبان فارسی می‌پردازیم و با معرفی یک الگوریتم تلاش می‌کنیم تنوع نگارشی واژه‌ها را به‌طور خودکار استخراج کرده و براساس آن، دادگانی برای تنوع نگارشی واژه‌های فارسی تهیه کنیم. سپس به طبقه‌بندی تنوع نگارشی می‌پردازیم. این دادگان می‌تواند برای اتخاذ شیوه‌های آموزش زبان فارسی به غیرفارسی‌زبانان، و یا در زبان‌شناسی پیکره‌ای و پردازش زبان طبیعی مورد استفاده قرار گیرد.

## کلمات کلیدی

پردازش خودکار زبان فارسی، زبان‌شناسی پیکره‌ای، پیکره مبتنی بر وب، فاصله نوشتن، تنوع نگارشی، دادگان، طبقه‌بندی

## ۱ مقدمه

و بتوان از ابزار تهیه‌شده برای پردازش آن استفاده کرد. زبان فارسی از این قاعده مستثنی نیست، و حتی با مشکلاتی به‌مراتب بیش از زبان‌های دیگر روبه‌روست. یکی از مشکلات اصلی آن، مربوط به ویژگی‌های خط زبان فارسی است. از جمله چالش‌های تهیه پیکره و پردازش زبان فارسی، نحوه نگارش چندواژی واژه‌های مرکب، تفاوت سبک نوشتاری رسمی و محاوره‌ای، همزه و پایه آن، و همچنین واژه‌های خارجی است. آنچه در این مقاله به آن می‌پردازیم ارائه یک الگوریتم است که بتواند به‌طور خودکار تنوع نگارشی واژه‌ها را از پیکره زبانی استخراج کند و پس از طبقه‌بندی تنوع نگارشی به‌صورت یک دادگان سازمان‌دهی شود. ساختار مقاله حاضر چنین است: در بخش ۲ مقاله، به پیشینه پژوهش‌های انجام‌شده در حوزه معیارسازی خط زبان فارسی می‌پردازیم. در بخش ۳، الگوریتم استفاده‌شده را معرفی می‌کنیم. بخش ۴ ویژگی‌های پیکره‌های زبانی مورد استفاده را بیان می‌کند. بخش ۵، به یافتن واژه‌های مشابه و طبقه‌بندی داده می‌پردازد. در بخش ۶ طبقه‌بندی تنوع نگارشی بررسی و تحلیل می‌شود. در بخش ۷، به نتیجه‌گیری پرداخته خواهد شد.

## ۲ پیشینه پژوهش‌های انجام‌شده در حوزه معیارسازی خط زبان فارسی

فرهنگستان زبان و ادب فارسی [۲] به‌عنوان متولی معیارسازی خط زبان فارسی دستور خط زبان فارسی را تهیه کرده است. اگرچه تلاش شده که با ارائه قواعد، معیارسازی خط زبان فارسی صورت پذیرد، برای هر قاعده فهرستی از استثناها ارائه شده‌است. این امر مبین این نکته است که معیارسازی خط زبان فارسی کاری ساده‌ای نیست. قیومی و ممتازی [۳]، قیومی و همکاران [۴] و شمس‌فرد [۹] به بررسی چالش‌های تهیه پیکره و پردازش زبان فارسی پرداخته‌اند.

شمس‌فرد و همکاران [۱۰] مجموعه ابزاری را تهیه کرده‌اند که با استفاده از آن می‌توان متن را با توجه به دستور خط فرهنگستان زبان و ادب

وب در ابتدای پیدایش خود به‌عنوان یکی از منابع مهم اطلاعات مورد استفاده قرار می‌گرفت. وبگاه‌های مختلف خبری یا بایگانی مقالات می‌توانست اطلاعات مفیدی را به کاربران ارائه دهد. در سال‌های اخیر این منبع اطلاعات نقش دیگری را نیز به‌عهده گرفته‌است و آن برقراری ارتباط میان کاربران می‌باشد. وبلاگ‌ها و گپ‌سراها نمونه‌هایی از نسل دوم وب می‌باشد که آنها این امکان را فراهم کرده‌اند تا انسان‌ها بتوانند به‌عنوان کاربر به تعامل با یکدیگر بپردازند و به این وسیله تفکرات، نظرات، و احساسات خود را درمورد مسایل گوناگون با استفاده از ابزار زبانی نوشتاری ابراز کنند. مسلماً داده حاصل از نسل جدید وب که روزانه به‌طور چشمگیری درحال گسترش است می‌تواند به‌عنوان یک منبع غنی‌تر اطلاعات مورد استفاده قرار گیرد. اما باتوجه به متنوع‌بودن کاربران به‌لحاظ تحصیلی، فرهنگی، اجتماعی، گویشی، جنسیتی، و سنی، تنوع سبکی در نوشتار آنها اجتناب‌ناپذیر است. این تنوع ممکن است با سبک نوشتار رسمی و معیار در کتاب، روزنامه، یا مجله متفاوت باشد.

یکی از کاربردهای عمده وب در پردازش زبان طبیعی می‌باشد، چراکه در پردازش زبان طبیعی یکی از منابع مهم مورد نیاز پیکره زبانی است. در پژوهش‌های پیشین انجام‌شده بر روی زبان‌های گوناگون، بیشتر از پیکره‌های استاندارد مانند وال‌استریت ژورنال<sup>۱</sup> برای زبان انگلیسی یا پیکره همشهری برای زبان فارسی استفاده شده‌است. براساس این داده‌های استاندارد، ابزار، روش، و الگوریتم‌های مورد نیاز پردازش زبان تهیه و معرفی شده‌است. پردازش داده وب با استفاده از ابزارهای تهیه‌شده یک چالش بسیار بزرگ و مهم محسوب می‌شود، زیرا داده وب بسیار پرنوفه است و پیش‌پردازش زیادی نیاز دارد تا داده به‌شکل استاندارد درآید

1 Corpus

2 <https://catalog.ldc.upenn.edu/LDC2000T43>

3 <http://ece.ut.ac.ir/dbrg/hamshahri/>

فارسی معیارسازی کرد. این ابزار متن‌باز<sup>4</sup> نبوده و به‌صورت رایگان در دسترس نیست. این ابزار سه کار تصحیح و تقطیع واژگانی، تحلیل واژگانی، و برچسب‌گذاری مقوله دستوری را انجام می‌دهد. همچنین سربابی و همکاران [۷] و سراجی و همکاران [۸] مجموعه ابزاری را تهیه کرده‌اند که پس از پیش‌پردازش داده ورودی، کار برچسب‌گذاری مقولات دستوری و همچنین تجزیه جملات در چارچوب دستور وابستگی انجام می‌شود. هزم<sup>5</sup> نیز مجموعه ابزار متن‌باز به زبان پایتون است که مانند دو مطالعه بالا عمل می‌کند. هدف قسمت پیش‌پردازش معیارسازی نسبی متن ورودی است. کارهای انجام‌شده مبین اهمیت جایگاه معیارسازی داده زبان فارسی و پردازش خودکار آن است. در مقاله حاضر به‌نوعی به‌دنبال معیارسازی متن فارسی هستیم، با این هدف که علاوه بر برآوردن نیاز افرادی که در حوزه پردازش زبان فارسی فعالیت می‌کنند، بتوان از آن برای زبان‌شناسی پیکره‌ای و همچنین اتخاذ روش‌هایی در حوزه آموزش زبان فارسی به غیرفارسی‌زبانان و رفع مشکلات در نوشتار این زبان استفاده کرد.

### ۳ معرفی الگوریتم یافتن واژه‌های مشابه

برای یافتن تنوع نگارشی واژه‌ها در این پژوهش، از الگوریتم فاصله لونشتاین [۵] استفاده کرده و مدل جدیدی را برای گسترش دادن این الگوریتم معرفی می‌نماییم. در الگوریتم اصلی فاصله لونشتاین، مدلی از برنامه‌نویسی پویا استفاده شده است که سه حالت جایگزینی، انطباق، یا درج حروف بین واژه‌های یک پیکره را در نظر می‌گیرد و با احتساب هر یک از سه حالت عددی که بیانگر فاصله دو واژه است ارائه می‌گردد. شبه‌کد این الگوریتم در زیر آمده است:

**ورودی:** دو واژه  $(W_1, W_2)$

**محاسبه فاصله لونشتاین:**

محاسبه طول واژه‌های  $W_1$  و  $W_2$ :  $M$  و  $N$

ساخت ماتریس دوبعدی  $(N+1) \times (M+1)$

پرکردن ردیف صفر ماتریس با نمایه ستون مرتبط با آن

پرکردن ستون صفر ماتریس با نمایه ردیف مرتبط با آن

**شروع حلقه**

شروع از ردیف یک و ستون یک تا ردیف  $N$  و ستون  $M$  به صورت ردیف به ردیف

محاسبه امتیاز سلول  $[i, j]$  براساس حالات ممکن

درج در واژه اول:

$$[i, j] = [i-1, j] + 1$$

درج در واژه دوم:

$$[i, j] = [i, j-1] + 1$$

انطباق در دو واژه:

$$[i, j] = [i-1, j-1]$$

جایگزینی یک حرف در دو واژه:

$$[i, j] = [i-1, j-1] + 1$$

انتخاب یکی از چهار حالت با کمترین امتیاز

**پایان حلقه**

**خروجی:** یک جفت واژه به‌همراه امتیاز محاسبه شده در سلول

$[N, M]$  فاصله لونشتاین

در نسخه گسترش‌یافته پیشنهادی ما، این الگوریتم به‌صورت وزن‌دهی شده مورد استفاده قرار می‌گیرد. به این معنا که در حالت درج یا جایگزینی همیشه عدد ۱ به مقدار سلول‌های مورد نظر اضافه نمی‌شود. بلکه بسته به اینکه چه حرفی درج شده و یا چه حرفی جایگزین یکدیگر شده‌است، وزن عددی متفاوتی به سلول مورد نظر اضافه می‌شود. در این

مدل، دو عدد  $0.1$  و  $1$  را در نظر می‌گیریم؛ به این ترتیب که برای حروف خاصی که در فارسی به‌دلیل تنوع نگارشی درج می‌شود عدد  $0.1$  اضافه شده و برای سایر حروف مانند نسخه اصلی الگوریتم عدد  $1$  اضافه می‌شود. همین‌طور برای جایگزینی حرفی که معمولاً در فارسی به دلیل تنوع نگارشی جایگزین هم می‌شود عدد  $0.1$  و برای سایر حروف عدد  $1$  اضافه می‌شود. این وزن‌دهی شامل واژه‌های ساده و مرکب می‌باشد. این وزن‌دهی سبب می‌شود راحت‌تر بتوان واژه‌های دارای تنوع نگارشی را جستجو کرد، چراکه فاصله چنین واژه‌هایی کمتر از  $1$  می‌باشد. نمونه‌هایی که برای آنها وزن کمتری لحاظ می‌گردد از قرار زیر است:

#### • جایگزینی

۱: جایگزینی حرف «آ» و «ا»: در پیکره مورد استفاده، واژه‌هایی یافته می‌شود که در آنها به‌جای «آ» از «ا» استفاده شده است، مانند: آرام و ارام: الآن و الان.

۲: جایگزینی همزه با معادل کرسی آن سبب تنوع نگارشی می‌شود:

۱-۲: جایگزینی «ا» و «آ»، مانند: راس و رأس.

۲-۲: جایگزینی «ا» و «ع»، مانند: مسأله و مسئله.

۳-۲: جایگزینی «ا» و «إ»، مانند: انشاءالله و إنشاءالله.

۴-۲: جایگزینی «و» و «ؤ»، مانند: مسوول و مسؤول.

۵-۲: جایگزینی «و» و «ع»، مانند: مسوول و مسئول.

۶-۲: جایگزینی «ی» و «ئ»، مانند: رییس و رئیس.

۳: جایگزینی حروف در گونه رسمی و گفتاری

۱-۳: جایگزینی شناسه «د» و «ه» در آخر فعل‌هایی مانند: می‌خورد و می‌خوره.

۲-۳: جایگزینی «ا» و «و» در واژه قبل از حروف «م» یا «ن»، مانند: نان و نون.

#### • درج

۱: درج «ء» در واژه، مانند: انشاءالله و انشاءالله، سوءاستفاده و سواستفاده

۲: درج «ی» در واژه، مانند: آئینه و آینه

۳: درج فاصله مجازی، مانند می‌شود و میشود و می‌شود

### ۴ پیکره‌های زبانی استفاده‌شده

در پژوهش حاضر، از پیکره تاک‌بانک<sup>۸</sup> استفاده شده‌است. این پیکره از وب تهیه شده و شامل وبلاگ و متون نوشتاری و روزنامه‌ای در وب با حجم بیش از ۵۰۰ میلیون واژه است. قاعدتاً پردازش این حجم داده که از منابع گوناگون جمع‌آوری شده کار بسیار دشوار و پیچیده‌ای است، چراکه با متن معیار فاصله زیادی دارد. با استخراج بسامد اولیه از پیکره، ویژگی چندنگارشی واژه‌ها خود را به‌راحتی نشان می‌دهد، زیرا سبب می‌شود که در فهرست واژگان استخراج‌شده چندین مدخل با بسامد متفاوت در فهرست یافت. برای مطالعه جامع‌تر تنوع نگارشی واژه‌ها، پیکره تاک‌بانک را با پیکره بی‌جن‌خان [۱] و همچنین واژه‌های پیکره وابستگی [۶] ادغام کردیم. لازم به ذکر است این دو پیکره از منابع نوشتاری و روزنامه تهیه شده و استاندارد شده‌اند.

### ۵ یافتن واژه‌های مشابه و طبقه‌بندی داده

در بخش ۳، مدل پیشنهادی ما برای محاسبه فاصله لونشتاین توصیف شد که از آن برای یافتن واژه‌های مشابه استفاده می‌شود. واژگان پیکره ادغام‌شده که در بخش ۴ ذکر گردید استخراج شده و به‌عنوان داده ورودی به الگوریتم داده می‌شود تا فاصله لونشتاین تمامی واژه‌ها محاسبه شود. در مرحله بعد، جفت واژه‌هایی که فاصله لونشتاین آنها کمتر از ۱

## ۶-۲ تنوع ناشی از تغییرات زبان‌شناسی

گاهی شاهد تنوعات نگارشی هستیم که از تغییرات آوایی حاصل می‌شود. این تغییرات آوایی اغلب به‌صورت قاعده‌مند رخ می‌دهد. در زیر به برخی از مواردی اشاره می‌شود که نتیجه این تنوعات و تغییرات است.

### ۶-۲-۱ سبک متفاوت

گاهی هنگامی که تغییرات آوایی در انواع مقولات دستوری ایجاد می‌شود واژه‌هایی با سبک متفاوت اعم از رسمی، محاوره‌ای عامیانه یا ادبی پدید می‌آید، مانند: «گه» (اگر)، «داره» (دارد) و «همون» (همان)، «نون» (نان)، «بام» (باهام)، و «برا» (برای).

در فعل‌ها، تغییرات آوایی و در پی آن، تغییر سبکی ممکن است در چند بخش ایجاد شود، مثلاً در ستاک فعل مانند «می\_s\_گم» (می‌گویم)، یا در شناسه آن مانند: «نمیده / نمی\_s\_ده» (نمی‌دهد)، و یا در هر دو،

مانند: «نمی\_s\_ره» (نمی‌رود) و «می\_s\_شه» (می‌شود). همان‌گونه که از نمونه‌ها برمی‌آید این تغییر سبکی عمدتاً از رسمی و معیار به محاوره‌ای یا عامیانه میل می‌کند.

از سوی دیگر، برخی از صورت‌های محاوره‌ای، عامیانه و ادبی در نتیجه تغییر آوایی به‌دست نیامده‌است. واژه‌های «ازش» (از+ش) و «اوا» (در صورتی که نشانه تعجب باشد) و «اوی» (هنگامی که برای صدا کردن به کار می‌رود) از آن جمله است. به عبارتی، این واژه‌ها معادل تغییر یافته واژه دیگری نیست.

### ۶-۲-۲ چندمعنایی و هم‌نویسی

در مواردی، تغییرات آوایی باعث ایجاد چندمعنایی یا هم‌آوا هم‌نویسی در واژه‌ها می‌گردد، مثلاً «بوم» می‌تواند صورت محاوره‌ای «بام» یا به معنای «بوم» نقاشی و یا به معنای جغد باشد.

نمونه دیگر «نمیره» است که هم به معنای «نمیرد» و هم «نمی‌رود» است. در تمام این موارد ممکن است همزمان با تغییر معنا، شاهد تفاوت در سبک واژه‌های متفاوت باشیم. همچنین این پدیده معنایی در تمام مقولات دستوری رخ می‌دهد و محدود به مقوله خاصی نیست.

### ۶-۲-۳ ساخت صرفی نحوی

همچنین، در برخی موارد، تغییرات آوایی سبب ایجاد تغییر در ساخت صرفی یا نحوی واژه‌ها می‌گردد. «منو» که به‌صورت «من\_س\_و» هم آمده در واقع از دو تکواژ «من» و «را» تشکیل شده است که در آن، «را» دستخوش تغییر آوایی و به واژه‌بست «و» (O) تبدیل شده است. «بم» یا «به\_س\_ام» نیز که صورت تغییر یافته «به+م» (=به من) است نیز چنین وضعیتی دارد.

### ۶-۲-۴ سایر موارد

در پیکره نمونه‌هایی از قبیل «ا\_س\_ن\_س\_ها»، «ر\_س\_ها» و «میکه» (که احتمالاً صورت درست آن «میکه» است) یافت شده‌است که به ماهیت خط فارسی (مبنی بر چسبیده بودن یا جدانویسی) مربوط می‌شود. علت وجود چنین مواردی می‌تواند شتابزدگی، بی‌دقتی، و البته بی‌اطلاعی شخص پدیدآورنده متن باشد. این اشکالات، سبب به‌وجود آمدن موارد بی‌معنی بسیار زیاد و متنوعی در پیکره شده‌است.

### ۷ نتیجه‌گیری

در این مقاله، به بررسی تنوع نگارشی در پیکره مبتنی بر وب پرداخته شد. برای استخراج تنوع نگارشی، تلاش کردیم با وزن‌دهی به حروف خاصی در الگوریتم لونشتاین، به‌طور خودکار واژه‌هایی که تنوع نگارشی داشت را مشخص کنیم. سپس با ارزیابی کیفی اولیه داده، طبقه‌بندی‌ای را با توجه به دلایل ایجاد تنوع نگارشی تهیه کردیم. این طبقه‌بندی، سه طبقه کلی

است از خروجی الگوریتم انتخاب می‌گردد. در مرحله آخر، واژه‌های مشابه کنار هم مرتبط می‌شود تا تحلیل داده و طبقه‌بندی آن ساده‌تر گردد. از آنجا که هدف این مقاله ارائه یک طبقه‌بندی از تنوع نگارشی واژه‌ها به‌همراه بسامد آن برای نمایش مقبولیت نگارش واژه توسط گویشوران فارسی است، داده به‌دست‌آمده به‌صورت دستی ارزیابی کیفی شده است. در گام نخست تحلیل داده، ۵۰۰ واژه که دارای تنوع نگارشی است توسط دو زبان‌شناس تحلیل گردید. در مواردی که اتفاق نظر وجود نداشت با تجزیه و تحلیل موردی، تصمیم‌گیری شده است. نتایج اولیه بررسی این واژه‌ها از سه جنبه بررسی شده است که در قسمت بعدی توضیح داده می‌شود.

## ۶ تحلیل طبقه‌بندی تنوع نگارشی

با بررسی ۵۰۰ واژه از فهرست واژه‌های مشابه، مواردی که سبب پدید آمدن تنوع نگارشی در واژه‌های مشابه شده‌است به شرح زیر دسته‌بندی می‌شود.

### ۶-۱-۱ تنوع ناشی از ویژگی‌های خط فارسی

#### ۶-۱-۱-۱ عدم رعایت مرز درون واژه‌ای

در واژه‌های مشتق و غیربسیط (اعم از تصریفی یا اشتقاقی) که دارای پیشوند و پسوندی مانند «می»، «ی»، «ای»، «ها» است، استفاده نادرست از فاصله مجازی و یا عدم استفاده از آن، مثال‌هایی مانند

«حرفها» و «حرف\_س\_ها»، «نمیدهد» و «می\_س\_گو\_ید» را به دست داده است. این مسئله در واژه‌های مرکب و مشتق مرکب مانند

«راه\_س\_اندازی» نیز صادق است. البته جدانویسی می‌تواند به این دلیل باشد که شخص، نمونه مورد نظر را دو واژه فرض کرده است، مانند

«ازنو» که به‌صورت «از\_س\_نو» هم آمده یا «برصد» که به‌صورت «بر\_س\_صد» نیز نوشته شده است.

#### ۶-۱-۱-۲ همزه

در مواردی، نوشتن همزه بر روی کرسی‌های گوناگون باعث پدید آمدن صورت‌های متنوع شده است، مانند «سوالهای» و «سوالهای». در برخی از واژه‌ها، نوشتن همزه به جای «ی» سبب تفاوت در تلفظ می‌شود، مانند: «امریکایی» و «امریکائی» و یا «سوئی» و «سوئی». همچنین، نوشتن یا نوشتن همزه عامل دیگری در گوناگونی واژه‌هاست، مانند: «رای»، «رای»، و «رای».

#### ۶-۱-۱-۳ نگارش «ا» به جای «آ»

حضور یا عدم حضور «آ» یا مَدّه، در بعضی واژه‌ها سبب تفاوت در تلفظ واژه می‌شود، مانند «امریکا» و «امریکا» یا «الان» و «الان». درحالی‌که در «انها» و «انها» یا «ابان» و «ابان» تفاوت در تلفظ نداریم. این امر ممکن است به چند دلیل باشد: سلیقه و سبک مؤلفان، شتابزدگی سازندگان پیکره یا بی‌اطلاعی آنها.

#### ۶-۱-۱-۴ حذف

در مواردی، حذف یک حرف از واژه را داریم، مانند حذف الف از فعل «است» در «آبست» (=آن است) و «آن\_س\_ست» و «آند» (= آن اند/هستند). حذف در این حالت تغییر معنایی در واژه ایجاد نمی‌کند و در یک دوره‌زمانی، بیشتر مشاهده می‌شود که نویسنده‌ها تابع سرهم‌نویسی بوده‌اند.

9 برای نمایش فاصله مجازی از علامت \_s\_ استفاده شده‌است.

[10] Shamsfard, Mehrnoush and Hoda Sadat Jafari and Mahdi Ilbeygi, “**STeP-1: A Set of fundamental tools for Persian Text Processing**”, In Proceedings of the 7th International Conference on Language Resources and Evaluation. Valletta, Malta, May 19-21, pp: 859-865, 2010.

برای تنوع نگارشی را دارد که می‌تواند ناشی از ویژگی‌های رسم الخط زبان فارسی، ویژگی‌های زبان‌شناختی، و یا سایر موارد باشد. نتیجه این پژوهش می‌تواند در آموزش زبان به غیرفارسی‌زبانان، زبان‌شناسی پیکره‌ای، و یا پردازش خودکار زبان فارسی مورد استفاده قرار گیرد. در ادامه این کار پژوهشی، علاوه بر بررسی تعداد بیشتری از واژه‌های مشابه، و سازمان‌دهی آن به صورت یک دادگان، به انطباق دستور خط فرهنگستان با واقعیت زبانی و بررسی آماری تنوع نگارشی بین گویشوران خواهیم پرداخت.

#### منابع

- [1] بی‌جن‌خان، محمود، نقش پیکره زبانی در نوشتن دستور زبان: معرفی یک نرم‌افزار رایانه‌ای، مجله زبان‌شناسی، ۱۹ (۲): ۴۸-۶۷، ۱۳۸۳.
- [2] فرهنگستان زبان و ادب فارسی، دستور خط فارسی، فرهنگستان زبان و ادب فارسی، تهران، ۱۳۸۹.
- [3] Ghayoomi, Masood and Saeedeh Momtazi, “**Challenges in developing Persian corpora from on-line resources**”, In Proceedings of 2009 IEEE International Conference on Asian Language Processing, Singapore, pp: 108-113, 2009.
- [4] Ghayoomi, Masood and Saeedeh Momtazi and Mahmood Bijankhan, “**A study of corpus development for Persian**”, In International Journal on Asian Language Processing, 20 (1): 17-33, 2010.
- [5] Levenshtein, Vladimir I., “**Binary codes capable of correcting deletions, insertions, and reversals**”, Soviet Physics Doklady, 10 (8): 707–710, 1966.
- [6] Rasooli, Mohammad Sadegh and Manouchehr Kouhestani and Amirsaeid Moloodi, “**Development of a Persian syntactic dependency treebank**”, In Proceedings of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Atlanta, USA, June 9-14, pp: 306–314, 2013.
- [7] Sarabi, Zahra and Hooman Mahyar and Mojgan Farhoodi, “**ParsiPardaz: Persian language processing toolkit**”, In Proceedings of IEEE 3rd International eConference on Computer and Knowledge Engineering, Mashad Ferdowsi University, October 31-November 1, pp. 73-79, 2013.
- [8] Seraji, Mojgan and Beáta Megyesi and Joakim Nivre, “**A basic language resource kit for Persian**”, In Proceedings of the 8th International Conference on Language Resources and Evaluation. Istanbul, Turkey, May 19-21, pp: 2245-2252, 2012.
- [9] Shamsfard, Mehrnoush, “**Challenges and open problems in Persian text processing**”, In Proceedings of the 5th Language & Technology Conference: Human Language Technologies as a Challenge for Computer Science and Linguistics, November 25-27, Poznań, Poland, pp: 65-69, 2011.