



## پیاده سازی یک غلط یاب املایی فارسی تحت وب

احمد یوسفان<sup>۱</sup> و بی بی صدیقه طباطبایی<sup>۲</sup>

<sup>۱</sup> کاشان، بلوار قطب راوندی، دانشگاه کاشان، دانشکده مهندسی، گروه مهندسی کامپیوتر، گروه مهندسی کامپیوتر،

<sup>۲</sup> دانش آموخته‌ی گروه مهندسی کامپیوتر دانشگاه کاشان، tabatabaeised@gmail.com

### چکیده

غلط یاب املایی فارسی یکی از ابزارهای مهمی است که در راستای کمک به نویسنده‌ی یک متن فارسی می‌تواند به او کمک شایانی در یافتن و درست کردن غلط فارسی نوشته شده در یک متن نماید. تاکنون غلط یاب‌های گوناگونی آماده شده و در برخی از ابزارهای نگارش به کار گرفته شده است. با این همه در این زمینه پژوهش هنوز ادامه دارد زیرا کارآیی بسیاری از این ابزارها در حد بالایی نیست و باید الگوریتم‌های تازه‌ای برای بهبود آنها پیشنهاد شود. به کارگیری ایده‌های نو در کنار دیگر کارهای پیشین می‌تواند به بهبود کار غلط یاب‌های کنونی کمک شایانی نموده و برخی از مشکل‌های آنها را برطرف نماید. در این مقاله در آغاز برخی از روش‌های متداول غلط‌یابی بررسی شده است سپس روشی ترکیبی برای غلط‌یاب املایی فارسی پیشنهاد شده و این پیشنهاد به کمک زبان javascript و php و html پیاده سازی شده است. تحت وب بودن این ابزار پیاده سازی شده کمک می‌کند تا بتوان به سادگی آن را آزمایش نمود و اشکال‌های آن را برطرف نمود.

### کلمات کلیدی

غلط‌یاب املایی، فارسی، پیاده سازی، وب

متائفه‌بهی‌دقتی و رواج بسیاری از خطاهای نگارشی به ویژه در صفحه‌های شخصی فارسی روی شبکه‌ی جهانی به اندازه‌ای زیاد شده است که بیم آن می‌رود زبان دیرپایی فارسی دچار مشکلات گوناگونی گردد و نتواند جایگاه خود را در میان زبان‌های زنده‌ی دنیا حفظ نماید.

امروزه رایانه‌ها برای آماده کردن نوشتۀ‌های گوناگون بسیار به کار برده می‌شوند و بیشتر کسانی که نوشتۀ‌ای را آماده می‌کنند از ابزاری در رایانه برای آماده کردن نوشتۀ‌ی خود کمک می‌گیرند. همانند دیگر زبان‌های دنیا باید ابزارهای گوناگونی به کمک نویسنده‌ی متن فارسی در دنیای کنونی بیاید تا بتواند به او در بهتر نوشتۀ‌من فارسی و همچنین پر هیز از خط امکن نمایند. همچنین افرون بر گزارش نادرست بودن یک کلمه پیشنهاد یا پیشنهادهایی برای کلمه‌ی جایگزین آن به کاربر بدهد تا او را در درست کردن آن کلمه کمک نماید.

ناهاینگی‌های گوناگونی در نگارش خط فارسی دیده می‌شود همچنین در نگارش رایانه‌ای متن فارسی قالب‌ها، ابزارها و سیستم عامل‌های گوناگون و روش‌های گوناگون کد کردن نوشتۀ‌ی فارسی دیده می‌شود که در [۱ ، ۳-۷] به برخی از ناهاینگی‌ها فرنگستان زبان و ادب فارسی کوشیده است برخی ناهاینگی‌ها را در خط فارسی برطرف نماید و استاندارد یکسانی برای نگارش خط فارسی پیشنهاد دهد در آدرس persianacademy.ir به روزترین نسخه از این استاندارد با نام دستور خط فارسی گذاشته شده است گرچه بسیاری از نویسنده‌گان و حتی برخی کتاب‌های درسی به خوبی این استاندارد را رعایت نمی‌کنند ولی به هر حال می‌تواند به عنوان پایه‌ای برای درست نویسی متن فارسی به کار گرفته شود.

تاکنون غلط‌یاب‌های گوناگونی برای زبان فارسی آماده شده است و پروژه‌های گوناگونی در این زمینه انجام شده است که در این میان Microsoft می‌توان به ابزار آماده شده برای خط‌یابی فارسی از

زبان فارسی دربردارنده گنجینه‌ی بزرگی از زیباترین سرودها و داستانها است. زبان فارسی یکی از پربارترین زبان‌های دنیا است. کتاب‌هایی چون مثنوی معنوی، دیوان حافظ، رباعیات خیام و ... به زبان‌های گوناگون گیتی برگردانده شده و بارها چاپ شده‌اند. برترین ویژگی این نوشتۀ‌ها، انسانی بودن آنها است بگونه‌ای که همه‌ی انسان‌ها گرایشی درونی به این نوشتۀ‌ها دارند. متائفه‌ای این درخت تنومند امروزه نیاز به توجه بیشتری دارد زیرا برای دنیا نوین آماده نشده است [۱].

این زبان که از نیمه‌ی سده‌ی سوم، آثاری از آن در دست است و از آن تاریخ به بعد، روز به روز گسترش یافته و آثار بی‌شماری در آن آفریده شده، در حیطه‌ی فرهنگ آسیایی و بلکه جهانی، تبدیل به زبانی شده است که شاهکارهای جهانی آفریده است؛ چنان که کمتر زبانی از زبان‌های زنده‌ی موجود می‌توان سراغ داشت که در این پنهانه با زبان فارسی کوس برابری زند. از جمله ویژگی‌های این زبان، دیرندگی آن است در طول بیش از هزار سال؛ بدین معنی که دیگرگونی‌ها در آن، در مقام سنجش با بسیاری از زبان‌های جهان، به نسبت کمتر بوده است؛ به طوری که فارسی زبانان امروز به آسانی می‌توانند شعر فردوسی را بخوانند و بهمند و از آن بالاتر این که در آثار بسیاری از گویندگان و نویسنده‌گان سده‌های اخیر و حتی معاصر می‌توان کاربردهای کهن هزار سال پیش را دید، حتی در کتاب‌های دوره‌ی دیستانی می‌توان شعر رودکی و فردوسی را گنجانید [۲].

رواچ صنعت چاپ و ماشین‌های تحریر و رایانه‌ها و شتاب زدگی در نوشتۀ‌ی بی‌دقتی را در درست و زیبا نوشتۀ دامن زده است. همچنین داده‌ها و صفحه‌های روی شبکه‌ی جهانی که برخی باید روزانه تغییر کنند، این دردرس را چندین برابر کرده‌اند [۱].



## جدول (۱): نتیجه‌های آماری از چهار نوع خطای اصلی

سندهای دولتی	لغت‌نامه Webster	نوع خطای اصلی
(۲/۶٪)۴	(۱۳/۱٪)۴۷	جابجایی
(۱۸/۷٪)۲۹	(۲۰/۳٪)۷۳	اضافه
(۳۱/۶٪)۴۹	(۳۴/۴٪)۱۲۴	حذف
(۴۰/۰٪)۶۲	(۲۹/۶٪)۹۷	جانشینی
(۹۲/۹٪)۱۴۴	(۹۴/۷٪)۳۴۱	کل

عامل‌هایی که در نمونه خطای مؤثر هستند عبارتند از:

- همسایگی حروف در صفحه کلید که منجر به خطاهای اساسی می‌شود.
- کاراکترهایی که برای تایپ آنها نیاز به فشردن کلید شیفت است.
- هم صدایی مانند اینکه به جای کلمه receive، کلمه recieve تایپ شود که هر دو آهنگی شبیه به هم دارند.
- شکل‌های شبیه به هم مانند اینکه به جای کلمه انتخاب، کلمه انتخاب تایپ شود.

## ۳- شناسایی خطای اصلی

شناسایی خطای شامل مراحل زیر است.

1. شناسایی محدوده کلمه در متن باید در ابتدای صورت گیرد. عموماً کاراکتر فاصله و نقطه محدوده کلمه را در متن مشخص می‌کنند. البته شناسایی محدوده کلمه در زبان فارسی با دشواری‌های فراوانی روبرو است که در [۱، ۴-۳ و ۱۰] توضیح داده شده است. در این زمینه پژوهش‌ها و پیشنهادهای گوناگونی انجام شده است که یکی از آنها در [۱۱] آمده است و در [۱۰] همین روش با تغییرها و ساده‌سازی هایی برای ریشه‌یابی فارسی به کار گرفته شده است.
2. خطاهای real-word به راحتی نظری کردن کلمه با کلمات لغتنامه قابل شناسایی نیست بلکه نیاز به اطلاعاتی در رابطه با متن می‌باشد. خطاهای non-word به وسیله چک کردن کلمه در لغتنامه قابل شناسایی است. این روش دارای مشکلاتی است که عبارتند از:
  - نگهداری به هنگام یک چنین لغتنامه‌ای دشوار است.
  - گسترش لغتنامه به گونه‌ای که همه کلمات موجود در متن را پوشش دهد.
  - زمان پاسخ‌گویی سیستم پایین نگه داشته شود.
- یک راه حل سریعتر ولی غیر قابل اعتماد برای بررسی کردن متن از خطاهای این است که کنترل شود که همه n-gram های موجود در یک کلمه به ترکیبات مانندی n-gram تعلق داشته باشند. همچنین می‌توان به کمک یک ریشه‌یاب خوب لغتنامه کوچکتری را به کار برد تا کارآیی افزایش یابد و فضای کمتری نیز نیاز باشد.

## ۴- تصحیح خطای اصلی

تصحیح خطای اصلی به دو صورت خودکار یا فعل و انفعالی انجام شود. تصحیح خطای خودکار به این صورت است که به صورت

office ۲۰۱۰ تا ۲۰۰۳ office رسانی اشاره نمود. پروژه‌های گوناگون دیگری نیز در این زمینه انجام شده است که شرکت‌های گوناگونی در ایران آنها را آماده نموده‌اند. همچنین خطای‌های متربازی چون spell، ispell و hunspell تا کنون آماده شده و به روز شده‌اند که در آنها زبان فارسی نیز افزوده شده است. با این همه این ابزارها هنوز کاستی‌های گوناگونی دارند و باید این کاستی‌ها به کمک پیشنهادهای و پژوهش‌های تازه‌ای برطرف شود یا از مشکلات آنها کاسته شود. بنابراین کار پژوهشی هر چه بیشتر در این زمینه، به بهبود خطای‌های گونی کمک می‌کند و همچنین علاقمندان به این زمینه را با الگوریتم‌های غلط‌یابی آشنا می‌نماید. در این مقاله کوشش شده است که یک روش ترکیبی برای بهبود غلط‌یابی فارسی پیشنهاد شود و سپس به صورت تحت وب آمده شود.

## ۲- غلط املایی و انواع آن

خطاهای املایی به دو دلیل ممکن است به وجود آیند [۸]؛ دلیل اول این است که در هنگام تایپ متن ممکن است خطای اتفاق بیفتد مانند کلمه taht که حاصل از خطای تایپی بر روی کلمه that است و این گونه خطای حاصل از خطای نوشتن کلمه است. دلیل دوم این است که خطای رخ داده ناشی از ندانستن شکل صحیح املای کلمه است در این موارد می‌توان کلمه seprate را مثال زد که حاصل از ندانستن املای صحیح کلمه separate است.

دو نوع خطای املایی در متن‌ها ظاهر می‌شود [۹]. خطاهای word و خطاهای non-word است. خطاهای non-word خطاهایی هستند که شکل غلط آنها در زبان وجود ندارد مانند غلط املایی real-word که خطای املایی کلمه the است. خطاهای real-word که خطاهایی هستند که تنها با توجه به جمله می‌توان آنها را شناسایی کرد و شکل غلط آنها نیز در زبان وجود دارد برای نمونه form که به صورت from نوشته شود به طوری که هر دو شکل کلمه در زبان وجود دارد هم from و هم form در این نمونه باید در جمله مورد نظر کلمه form نوشته می‌شد ولی from نوشته شده است. این گونه خطاهای کاملاً وابسته به جمله است.

در ادامه منظور از هر گونه خطای املایی، خطای non-word است، زیرا در این مقاله فقط به این نوع از خطای پرداخته شده است. خطاهای real-word مستلزم پردازش جمله و متن به صورت کامل و در صورت نیاز درک کامل آن از سوی برنامه است. هنوز راه درازی تا انجام این کار در زبان فارسی وجود دارد و غلط‌یابهای کنونی هیچکدام این نوع خطاهای را به خوبی شناسایی نمی‌کنند. انواع اصلی خطاهای را می‌توان به صورت زیر ذکر کرد.

- خطاهایی که حاصل از اضافه شدن یک حرف است مانند across که می‌تواند خطای در نوشتن کلمه cress باشد.
- خطاهایی که حاصل از حذف شدن یک حرف در کلمه است، مانند across که می‌تواند خطای در نوشتن کلمه actress باشد.
- خطاهایی که حاصل از عوض شدن یک حرف با حرف دیگر است، مانند across که می‌تواند خطای در نوشتن کلمه باشد. که در این نمونه حرف o به حرف e تبدیل شده است.
- خطاهای جایه جایی، یک حرف کلمه با یک حرف دیگر از همان کلمه جایه جای می‌شود؛ مانند across که می‌تواند خطای در نوشتن کلمه caress باشد، که در این نمونه دو حرف c و a جایه جای شده‌اند.

نتایج آماری از این نوع خطاهای برای زبان انگلیسی انجام شده است در جدول ۱ یکی از این نتایج نشان داده شده است [۸].



## دانشگاه سمنان - دانشکده مهندسی برق و کامپیوتر

خودکار کلمه غلط با کلمه درست جایگزین می‌شود بدون آنکه کاربر مداخله‌ای داشته باشد. بنابراین کلمه روی کلمه با املای خوب داشته باشد که بهترین کلمه درست پیشنهادی را به دست آورد. تصحیح خطای فعل و انفعالی می‌تواند یکی از پیشنهادهایی که برای تصحیح کلمه غلط داده شده است را خود انتخاب کند. بنابراین غلط یا ب می‌تواند چندین پیشنهاد را برای تصحیح کلمه غلط به کاربر بدهد.

تصحیح کلمات بدون در نظر گرفتن متى که درونش هستند شامل سه مسئله می‌شود [۱۲]:

الف- شناسایی خطأ

ب- ساخت پیشنهادهایی برای تصحیح

ج- رتبه‌بندی پیشنهادها

شناسایی خطأ پیش از این شرح داده شد. ماتریس‌های n-gram در یک متن می‌توانند پیشنهادهای تصحیح و رتبه‌بندی را انجام می‌دهند. در اینجا اشاره‌گری نیاز است که ارتباط میان یک عنصر ماتریس n-gram را با داده‌ای که آن غنصر درون آن است برقرار نماید. برای داده‌های زیاد (تعداد کلمه‌ها و اسناد زیاد) این راه حل عملی به نظر نمی‌رسد. اگر تعداد کلمه‌ها زیاد باشد یک روش برای کاهش تعداد کلمه‌های پیشنهاد شده برای تصحیح این است که زیر مجموعه‌ای از کلمه‌ها جستجو شود. با این کار احتمال از دست دادن پیشنهاد کلمه‌ی درست وجود دارد. بنابراین باید بر روی هر الگوریتم پیشنهادی دو عامل زمان پاسخ‌گویی و کیفیت پاسخ بررسی انجام شود.

معمولًا رتبه‌بندی میزان شباهت میان کلمه‌ی مورد بررسی و کلمه‌ی پیشنهادی را می‌سنجد و برای این کار از نتیجه‌ی الگوریتم شباهت کمک می‌گیرد. رتبه‌بندی می‌تواند ضعیف یا قوی باشد. در رتبه‌بندی ضعیف کلمه‌های پیشنهادی رتبه‌بندی یکسانی داشته باشند و بنابراین شایستگی یکسانی را برای جانشینی دارند. افراطی ترین گونه‌ی رتبه‌بندی ضعیف آن است که همه کلمه‌های پیشنهادی رتبه‌ی یکسانی داشته باشند یا به عبارت بهتر اصلاً رتبه‌گذاری انجام نشود و فقط از میان کلمه‌های لغتنامه یا کلمه‌ای پیشنهاد شود یا خیر (تصمیم‌گیری دودویی). در رتبه‌بندی قوی به هر کلمه‌ی پیشنهادی یک رتبه‌ی یکتا نسبت داده می‌شود. در برنامه‌های تصحیح خودکار رتبه‌بندی قوی مورد نیاز است.

## ۵- روش‌های بررسی شباهت میان کلمات

در اینجا باید بتوان کلمه‌ها را با هم تطبیق داد؛ بنابراین به روش‌هایی برای بررسی و شباهت و نزدیکی کلمات و کلید شباهت پرداخته می‌شود. در این قسمت الگوریتم‌های متفاوتی تقسیم بندی شده‌اند. که همه یک حقیقت مشترک را دارند و آن این است که آنها سعی دارند که معیاری را برای فاصله میان کلمات بدهنند. این معیارهای فاصله می‌تواند انواع زیادی داشته باشند.

## ۱-۵- معیار فاصله Damereau

چهار نوع از خطاهای املایی تک‌کاراکتری رایج‌تر در متون تایپ شده با ماشین را در [۹] شناسایی شده و یک الگوریتم برای تصحیح آنها ارائه شده است. این خطاهای اضافه، حذف، جانشینی و جابه‌جایی هستند و تخمین زده شد که ۸۰ درصد از خطاهای

موجود در متن از این چهار نوع می‌باشند. بنابراین همه اضافه‌ها، حذف‌ها، جانشینی‌ها و جایه‌جایی‌های ممکن روی کلمه با املای غلط آزمایش می‌شود و سپس کلمات تغیر یافته با لغتنامه نظری می‌شوند. این روش یک مجموعه‌ای از  $1 - k*(2^{n+1} + n)$  کلمه صحیح را ارائه می‌دهد. به طوری که k تعداد حروف الفبای زبان و n طول کلمه با املای غلط است. برای یک کلمه ۷ کاراکتری در زبان سوئدی (۹ حرف) یک مجموعه ۳۴۱ کلمه‌ای باید با لغتنامه بررسی شوند. اگر ترکیباتی از حروف که نمی‌تواند وجود داشته باشند، حذف شوند؛ نتیجه می‌تواند به نصف کاهش باید. به طور مثال در [۱۳] فیلتر اعتباریابی letter-N-Gram به صورت موثری در این مورد کار کرد. هزینه‌های زمانی می‌تواند با الگوریتم‌های تصفیه شده حذف شود به عنوان نمونه Domeij Domeij این روش جستجوی لغتنامه، به صورت کارآمد اجرا کرد. همچنین این روش به وسیله Pollock & Zamora با لغتنامه کاهیده شده سازگار شده است. این لغتنامه کاهیده شده به وسیله روش‌های کلید شباهت خوشنان بسته آمده است.

اصل علمی که این روش دنبال می‌کند کار مستقیم با کلمه غلط، تولید تصحیحات ممکن و سپس مقایسه آنها با لغتنامه است. به این روش نام دیگر فاصله ویرایشی معکوس نیز داده اند. بیشتر روش‌های دیگری که وجود دارند این گونه کار می‌کنند که در ابتدا دو رشته را داده و سپس شباهت آنها را بررسی می‌کنند.

## Damereau-Levenshtein - ۲-۵

روشن است که راه حل Damereau خطاهای املایی را در کلمات با حداقل یک خطأ حل می‌کند. اصل Damereau از چهار خطای املایی متقابل اضافه کردن، حذف، جانشینی، جایه‌جایی بهره می‌گیرد که برای مقایسه فاصله ویرایشی بین دو کلمه طراحی شده است. بنابراین یک پایگاه داده از نام‌های درست و کامل می‌تواند بر طبق شباهت و نزدیکی به کلمه و رویی رتبه‌بندی شود. Damereau-Levenshtein با یک خطای تنها یا به نوعی دیگر فاصله ویرایشی ۱ متوقف نمی‌شود. در عوض مستقل از اینکه تغیرات چه چیز را می‌سازد، کوتاهترین فاصله بین یک کلمه (s1) و کلمه دیگر (s2) را محاسبه می‌شود. اندازه این فاصله با فرمول بازگشتی زیر محاسبه می‌شود که d یک ماتریس دو بعدی است.

$$d(0,0) = 0$$

$$d(i,j) = \min\{$$

$$\begin{aligned} & \text{اضافه کردن , } d(i-1,j)+1 \\ & \text{حذف کردن , } d(i,j-1)+1 \\ & \text{جانشینی , } d(i-1,j-1)+c(i,j) \\ & \text{جایه‌جایی } d(i-2,j-2)+c(i,j-1)+c(i-1,j)+1 \end{aligned}$$

}

اگر z مین حرف از کلمه اول برای z مین حرف از کلمه دوم باشد آنگاه مقدار  $d(z,i)$  برای صفر و گرنه برای یک خواهد بود. این یک فاصله یا متریک استاندارد است که اغلب به وسیله برنامه‌نویسی پویا قابل اجرا است. رویه‌های دیگر برای حل مشکلات منظور شده‌اند که به طور مثال می‌توان از نمونه ای که توسط Brodder یا نمونه دیگری که توسط Du & chan معرفی شده نام برد. برای بهبود سرعت اجرا می‌توان از یک معیار برای بررسی استفاده کرد، منظور، برش قطع روند کار است؛ که قادر است هرگاه خطاهای زیادی شناسایی شد فرایند مقایسه را ناتمام رها کند.



## ۶- روش پیشنهادی

### ۷- پیاده سازی تحت وب

- پروندهای نوشته شده برای این غلط یاب **pellChecker.html** است:
- پرونده **spellChecker.php**: این پرونده، یک صفحه است که در مرحله اول به کاربر نمایش داده می شود تا متن خود را در آن وارد کند.
- پرونده **spellChecker.php**: این پرونده متى را که کاربر وارد کرده است را از **spellChecker.html** دریافت کرده و عملیات غلط یابی و پیشنهاد تصحیح را انجام داده و نتیجه را به کاربر نمایش می دهد.
- پرونده **lib.php**: این پرونده تابع هایی را دربردارد که برای غلط یابی و تصحیح خطای پرونده **spellChecker.php** به آن نیاز دارد.
- پرونده **dictionary.txt**: این پرونده شامل کلمات صحیحی است که به عنوان لغتنامه از آن استفاده می شود. این لغتنامه در بردارنده همه کلمه هایی است که در زبان فارسی درست هستند و چون ریشه هایی در این غلط یاب به کار برده می شود همه کلمه ها و مشتقات آن در این پرونده گذاشته شده است. این پرونده بر پایه کارهای پیشین آماده شد و به صورت خودکار همه مشتقات کلمه ها ساخته شده بود و با کمک کلمه های گرفته شده از روزنامه ها با فراوانی بالاتر کلمه های درست برگزیده شده بودند و به این ترتیب می توان مشتقات کلمه ها ساخته شد. از سوی دیگر برخی از ترکیب ها که درستی آن ها روش ن است بدون بررسی در کلمه های درون روزنامه ها به این مجموعه از کلمه ها افزوده شدند. برای کاهش سرعت جستجو در این کلمه ها می توان آن ها را به صورت مرتب درون پرونده گذاشت و سپس به کمک جستجوی دودویی روی آن ها جستجو کرد تا سرعت به دلیل اندازه بزرگ این پرونده کاهش نیابد.
- غلط یاب به این ترتیب کار می کند که می توان کلمه یا کلمه هایی را وارد کرد و غلط یاب کلمه های نادرست را نشان می دهد و به ترتیب پیشنهادهای برای آن کلمه را نیز نشان می دهد سپس کاربر می تواند از میان پیشنهادها که به ترتیب به او نشان داده می شود یکی را برگزیند. در جدول زیر پیشنهادهای داده شده برای تعدادی از کلمه های داده شده نشان داده شده است.

جدول (۲): خروجی خطایاب

کلمه	پیشنهادهای داده شده
سلام	سلام، صلاح، صلاحیم، غلام، کلام، آلام، صنام، صلح، صلیم، صلاح، لام، صلا
هال	ستنا، سمع، جسمما، رسمما، سیما، سرما، سمهما، سماق، سماک، هما، شما، بما، اما، نما، دما، سقا، سزا، سرا، سوا، سمج، سمش، سمت، سمم، ما، سم، سام
حه	صحه، حقه، حصه، حبه، چه، سه، به، له، نه، مه، که، زه، ره، ده، آه، وه، حج، حق، حس، حی، حب، حل، حک، حر، حد
توره	توره، دوره، طفره، خوره، غوره، سوره، سوره، اوره، کوره، پوره، طوفه، طورش، طوری، طورت، طورم، طره، طور

طی بررسی هایی که، نه به صورت آماری بلکه از روی نتایج به دست آمده از دیگر زبان هایی که رسم الخط آنها شبیه فارسی است نظیر زبان اردو، انجام شد، می توان انواع خطاهایی که بیشتر در زبان فارسی اتفاق می افتاد را به شکل زیر معرفی کرد.

- خطاهایی که حاصل از شباهت شکل دیداری کلمات است.
- خطاهایی که به دلیل تاثیر کلید شیفت اتفاق می افتد.
- خطاهایی که حاصل از حذف فاصله است.

درصد از خطاهای مربوط به خطاهای Real-word می باشد.

مانند اسلامی که به صورت غلط یاب استفاده شده است، تکنیک های که برای طراحی این غلط یاب استفاده شده است، تکنیک های Edit Distance و Soundex در این پژوهش به کار گرفته شده است. و برای بررسی کردن اعتبار رشته های تصحیح تولید شده از لغتنامه استفاده شده است، به گونه ای که رشته تولیدی اگر جزء کلمات لغتنامه باشد معتبر شناخته می شود. تغییراتی که در روش Soundex داده شده است به قرار زیر بوده است:

- حرف اول نیز کد می شود.
- تعداد گروه های حرفی افزایش یافته است.
- از کدهای هگزا سیمال به جای کدهای دسیمال استفاده شده است.

ابتدا تمام پیشنهادهای ممکن را با اجرای روش فاصله ویرایشی معکوس که در بالا شرح داده شد، به دست آورده و سپس برای رتبه بندی، ابتدا پیشنهادهایی که از روش Edit Distance برای حذف فاصله مجازی به دست آمداند، دارای بالاترین رتبه قرار داده و سپس مابقی این پیشنهادها را برای رتبه بندی به الگوریتم Soundex داده و آن پیشنهادهایی که از نظر این الگوریتم پذیرفته شده اند در رتبه دوم قرار داده و مابقی پیشنهادهایی که مانده است را در رتبه چهارم. ازین این پیشنهادهای رتبه بندی شده، ۱۰ مورد که دارای بیشترین رتبه می باشند، انتخاب کرده و به کاربر نمایش داده شده است.

دلیل انتخاب این روش ها این است که می خواستیم هم تاثیر هم صدایی ها را در ایجاد خطاهای املایی لحظه کنیم و هم معیاری را برای میزان تفاوت دو کلمه داشته باشیم. شاید در نظر شما الگوریتم های Metaphone یا Editex این کار را به خوبی انجام می دهند ولی باید توجه داشته باشید که این الگوریتم ها در ابتدا به قوانین موجود در زبان می پردازند در صورتی که در زبان فارسی قوانین زیادی برای صدای کلمات نداریم.



## دانشگاه سمنان - دانشکده مهندسی برق و کامپیوتر

- [6] Kokabi, Mortaza. *The Newer, the Worse: the Status of Farsi Word Processing Softwares in Iran*, Proceedings ELPUB2005 Conference on Electronic Publishing – Kath. Univ. Leuven – June 2005
- [7] یوسفان، احمد. ساخت یک نرم افزار خطایاب املایی و ریشه یاب خودکار فارسی. طرح پژوهشی داخلی در دانشگاه کاشان، ۱۳۸۸.
- [8] Kukich, Karen. *A Comparison of Some Novel and Traditional Lexical Distance Metrics for Spelling Correction*. In: Proceedings of INNC-90-Paris, Paris, France, July 1990, pp 309-313.
- [9] F. J. Damerau. 1964. A technique for computer detection and correction of spelling errors. In Communications of the ACM, volume 7(3), pages 171–176.
- [10] شیخ زادگان. جواد مرکز هوشمند عالم. ساخت یک نرم افزار خطایاب املایی و ریشه یاب خودکار فارسی. طرح پژوهشی داخلی در دانشگاه کاشان، ۱۳۸۸.
- [11]
- [12] Sannella, M. J., *Constraint Satisfaction and Debugging for Interactive User Interfaces*, Ph.D. Thesis, University of Washington, Seattle, WA, 1994.
- [13] Zachman, John A., "A Framework for Information Systems Architecture", IBM Systems Journal, Vol. 26, No. 3, 1987.
- [14] Plamondon, R., Lorette, G., "Automatic Signature Verification and Writer Identification - The State of the Art", Pattern Recognition, Vol. 22, pp. 107-131, 1989.
- [15] Object Management Group. *Unified Modeling Language: Superstructure*, Version 2.0, ptc/03-07-06, July 2003, <http://www.omg.org/cgi-bin/doc?ptc/2003-08-02>.

این به این معنی است که تنها کلمات با همان شکلی که هستند بررسی می شوند و در لغتنامه نیز همه کلمات وجود دارد؛ به طور مثال هم باید کلمه "خوردند" و هم کلمه "می خورند"، "خوردم" و "خوردید" و تمام ترکیبات خوردن در این لغتنامه وجود داشته باشد؛ و هر چه حجم کلمات لغتنامه بیشتر، زمان اجرای عملیات غلطیابی نیز افزایش داده می شود. بنابراین یکی از مواردی که اگر در آینده به آن پرداخته شود در بهبود کارایی این غلط یاب بسیار موثر است این است که قبل و بعد از عملیات غلط یابی روش های دیگر نیز اجرا شوند که وظیفه این روش ها ریشه یابی کلمه و روایی و پس از غلط یابی، بیشندهایی که مطرح می شود، حالات مختلفی که می توانند با این ریشه پیشنهادی به کاربر داده شود، داده شود. یکی دیگر از بهبودهای ممکن برای این غلط یاب، بررسی روش های سرعی جستجو در پرونده است یا بررسی بهترین حالتی که بتوان لغتنامه را نگهداری کرد.

## ۸- نتیجه

زبان فارسی به عنوان زبان رسمی ایران به عنوان یکی از نمادهای فرهنگ ایران به حساب می آید. کتابهای نوشته شده گشتگانمان به این زبان به خوبی نشان دهنده تعلق این آثار گرانبهای ایران و فرهنگ ایران است. با وجود زبان ها و لهجه های گوناگونی که در ایران هستند که نباید این زبان ها و لهجه ها از میان بروند و باید از آنها نگهداری کرد و به کار برد ولی نگهداری و به کارگیری زبان فارسی به عنوان یکی از نمادهای یکپارچگی و وحدت ایران ارزش بالایی دارد و کمرنگ شدن زبان فارسی به هر دلیلی، خیانت به یکی از فرهنگ های کهن گیتی است. این مقاله به یکی از مشکلات رو در روی زبان فارسی، خط در نوشتار، پرداخته است. برخی از پژوهش های انجام شده در این زمینه بررسی شده است و به بهرگیری از کارهای پیشین و ترکیب و تغییر برخی از آنها خطایاب ساده ای تحت وب طراحی و پیاده سازی شده است. گرچه این خطایاب توانایی بالای ندارد ولی اینده های به کار رفته در آن می تواند در کارهای آینده به کار گرفته شود و یا کارهای دیگر را بهبود بخشد.

## مراجع

- [۱] یوسفان، احمد. یک سیستم بازیابی اطلاعات متنی برای زبان فارسی بر پایه نمایه گذاری معانی پنهان. پایان نامه کارشناسی ارشد، شیراز : دانشگاه شیراز، شهریور ۱۳۸۲
- [۲] احمدی گیوی ۱۳۸۰، حسن. دستور تاریخی فعل. تهران: نشر قطره، ۱۳۸۰، جلد اول، ص ۲۷.
- [۳] یوسفان، احمد، سمية صالحی، بهروز مینایی بیدگلی. دشواری های ریشه یابی فارسی و روشی برای ریشه یابی فعل های ساده فارسی، دومین کارگاه پژوهشی زبان فارسی و ریانه، تهران: دانشگاه تهران، تیرماه ۱۳۸۵
- [۴] یوسفان، ۱۳۸۴، احمد، و منصور ذوالقدری جهرمی، مهدی فارسی. کنفرانس بین المللی اطلاعات و دانش (IKT2005)، تهران: امیرکبیر، خرداد ۱۳۸۴
- [۵] احمدی بیرجندی ، احمد (۱۳۶۸). شیوه‌ی آموزش املای فارسی و نگارش. تهران: انتشارات وزارت فرهنگ و ارشاد اسلامی . چاپ اول بهار ۱۳۶۸