

تصحیح خودکار غلط‌های تایپی فارسی به کمک شبکه عصبی مصنوعی ترکیبی

رسول دژکام^۳

رضا صفابخش^۲

امیرشہاب شاھمیری^۱

۱- دانشآموخته کارشناسی ارشد دانشکده مهندسی کامپیوتر- دانشگاه صنعتی امیرکبیر - تهران- ایران

amir@shahmiri.ir

۲- استاد دانشکده مهندسی کامپیوتر- دانشگاه صنعتی امیرکبیر - تهران- ایران

safa@ce.aut.ac.ir

۳- دانشآموخته کارشناسی ارشد دانشکده مهندسی کامپیوتر- دانشگاه صنعتی امیرکبیر- تهران- ایران

dezhkam@ce.aut.ac.ir

چکیده: ارایه راهی برای تصحیح غلط‌های املایی نگاشته شده توسط انسان یکی از اهداف مورد توجه در دانش هوش مصنوعی، متن کاوی و پردازش زبان طبیعی است. بیشتر روش‌های موجود برای تصحیح غلط‌های املایی بر پایه الگوریتم‌های جستجو در فرهنگ واژگان و تعیین نسبت شباهت واژگان درست موجود در فرهنگ واژگان با واژه نادرست مورد نظر کار می‌کنند.

در این پژوهش طراحی، پیاده‌سازی و ارزیابی یک مصحح املایی به کمک شبکه‌های عصبی مصنوعی هاپفیلد و پرسپترون چند لایه با رویکرد ویژه به غلط‌های تایپی کاربر ارایه می‌شود. نتایج به دست آمده نشان می‌دهند که برای یادگیری واژه‌نامه‌ای مشتمل بر ۴ تا ۲۵۶ واژه ۴ تا ۶ حرفی و تصحیح غلط‌های مربوط به آنها، شبکه هاپفیلد به دقیقی بین٪۰.۵۵ و٪۱۰۰ درستی و شبکه پرسپترون چندلایه - که در این تحقیق عمل یادگیری را در قالب دسته‌بندی انجام می‌دهد - به دقیقی بین٪۸۰ و٪۱۰۰ درستی دست یافته، که این مقدار با به کارگیری شبکه‌های ترکیبی به نزدیک به٪۸۰ دقت درستی برای بیش از ۳۰۰۰ واژه افزایش یافته است.

واژه‌های کلیدی: صحیح خودکار غلط تایپی فارسی، غلط املایی، شبکه عصبی مصنوعی هاپفیلد، پرسپترون چند لایه، فاصله کلید.

تاریخ ارسال: مقاله : ۱۳۸۶/۸/۴

تاریخ پذیرش مقاله : ۱۳۸۷/۱۰/۲

نام نویسنده‌ی مسئول : رضا صفابخش

نشانی نویسنده‌ی مسئول : ایران - تهران - خیابان حافظ - پلاک ۴۲۴ - دانشگاه صنعتی امیر کبیر - دانشکده‌ی مهندسی کامپیوتر و فناوری

اطلاعات



۱- مقدمه

دسته تقسیم کرد[۴]:

غلط تایپی^۸، که چهار گروه اصلی درج حرف اضافه^۹، حذف حرف^{۱۰}، نگارش یک حرف اشتباه به جای حرف اصلی^{۱۱} و جایه جایی دو حرف همسایه^{۱۲} را در بر دارد[۳و۵]. این چهار گونه، نزدیک به ۸۰٪ خطاهای تایپی را در بر می‌گیرند[۶].

غلط املایی^{۱۳}، که ناشی از ضعف دانش زبانی نویسنده در تشخیص صوات واژه و شیوه نگارش آن است.

خطای انتقال^{۱۴}، که در هنگام انتقال اطلاعات بر شبکه یا دیسک و بازناسانی نوی حروف^{۱۵} رخ می‌دهد.

برای پیاده‌سازی یک سیستم مصحح املایی، دو کار اصلی باید انجام پذیرد: نخست باید واژه موردنظر در فرهنگ واژگان جستجو شود تا در صورت موجود بودن غلط قلمداد گردد؛ سپس با استفاده از روش‌های گوناگون در حالت خودکار^{۱۶} یک، یا در حالت محاوره‌ای^{۱۷} چند واژه به عنوان جایگزین پیشنهاد گردد[۷]. از آنجا که بسیاری از واژگان با تغییر یک حرف به واژه‌ای دیگر تبدیل می‌شوند، هیچگاه نمی‌توان انتظار داشت که دقت تصحیح لغوی^{۱۸} به صدرصد برسد و دستیابی بهترین روش‌های مصحح به این دقت نیز در عمل - حتی با توجه به درک مفهوم جملات بر پایه روش‌های پردازش زبان طبیعی - برای تصحیح تمامی واژگان یک زبان امکان‌ناپذیر خواهد بود. گذشته از این، با توجه به این که بیشتر روش‌های تصحیح غلط بر روی تنها یک خطای واژه خوب کار می‌کنند، تا زمانی که روش‌های تصحیح خطاهای ترکیبی مانند جایه جایی رشته^{۱۹} و وارونه نویسی^{۲۰} به میان نیایند، نمی‌توان به آنها اطمینان چندانی داشت[۴]. از این‌رو به نظر می‌رسد که دقت درستی ۹۰٪ برای یک مصحح خودکار مطلوب باشد؛ در حالی که این مقدار برای انسان ۷۵٪ است[۸]. از این‌رو این‌گونه غلطها به ندرت در گروه واژگان آزمایشی جای می‌گیرند[۹].

اغلب روش‌هایی که تاکنون برای تصحیح لغوی ارایه شده‌اند، تطبیق رشته حروف^{۲۱} واژه غلط^{۲۲} را - که در فرهنگ واژگان موجود نیست - با نزدیکترین واژه در فرهنگ واژگان، بر پایه فاصله ویرایشی^{۲۳}، فاصله همینگ^{۲۴} یا فاصله لونشتین^{۲۵} به کار می‌گیرند و یک یا چند واژه جایگزین را پیشنهاد می‌دهند[۹-۳]. این فاصله ویرایشی بنا بر اختلافات تک‌تک حروف واژه غلط با کلمات فرهنگ واژگان به دست می‌آید و از این‌رو اغلب - به‌ویژه در هنگام ویرایش واژگان کم حرف - نتیجه دلخواهی به‌همراه ندارد. برای نمونه، اگر به جای واژه "رشته" "word" را تایپ شود، این روش ۷۱۰ واژه جایگزین هم‌سطح را پیشنهاد می‌دهد[۵].

هدف این پژوهش ارایه روشی جدید برای ویرایش لغوی متون با رویکرد ویژه به غلط‌های تایپی است. دلیل این انتخاب آن بوده که متونی که توسط کاربر تایپ می‌شوند، خود بیشتر توسط نویسنده (که ممکن است خود کاربر باشد) و با آگاهی وی از ساختار جملات و به‌ویژه مفهوم آنها شکل گرفته‌اند و بنابراین نیاز چندانی به ویرایش دستوری و مفهومی ندارند و حتی ممکن است به‌دلایل ویژه، به‌عمد

گفت‌وگوی معنی دار میان انسان و ماشین یکی از آرزوهای دانشمندان علوم رایانه و یکی از مباحث مورد توجه در زمینه هوش مصنوعی است. در این زمینه، درک جمله بیان شده توسط انسان به زبان طبیعی از سوی ماشین و نیز تولید جمله درست و با معنی توسط ماشین، از مهم‌ترین اهداف علوم پردازش و فهم زبان طبیعی^۱ و هوش مصنوعی است که هرچند تاکنون به موفقیت‌های چشمگیری دست نیافته، اما افق‌های روشی را پیش روی پژوهشگران این عرصه قرار داده است.

مسئله مشابه دیگری که محققان هوش مصنوعی و متن کاوی^۲ بدان توجه دارند، ویرایش یا درک متن نوشته شده به دست انسان، توسط ماشین است که معمولاً سه گام زیر را در بر می‌گیرد[۱]:

ویرایش لغوی^۳: که در آن کوشش می‌شود تا در صورت اشتباه بودن نگارش املایی یک واژه، جایگزین مناسبی به کمک یک پایگاه داده جامع از واژگان برای آن پیدا شود. برای نمونه به جای واژه "ویرایش" واژه "ویرایش" پیشنهاد گردد.

ویرایش دستوری^۴: که در آن کوشش می‌شود تا اشتباهات دستوری جملات شناسایی و تصحیح گردد. برای نمونه، رخداد زمانی میان فعل و قید جمله "دیروز من او یکدیگر را دید." کشف گردد و جمله‌ای مانند "دیروز من او یکدیگر را دیدیم." پیشنهاد گردد. این امر به کمک مجموعه قواعد دستورزبانی و روش‌های گوناگون تجزیه واژگان جمله در مبحث پردازش و فهم زبان طبیعی به نتایجی رضایت‌بخش رسیده است.

ویرایش مفهومی^۵: که در آن کوشش می‌شود تا به رغم درستی ساختار جمله از نظر دستورزبان، ناهمانگی‌های مفهومی آن دست کم شناخته و تا جای ممکن تصحیح گردد. برای نمونه جمله "میز هوایما را خورد." از دیدگاه املایی و دستوری اشتباهی ندارد، اما مفهوم درستی از آن بر نمی‌آید.

روشن است که در هر گونه پردازش متن نخستین گام، تصحیح غلط‌های واژگان متن است. غلط‌های واژه‌ای یک متن به دو دسته تقسیم می‌شوند: یکی آن که واژه مورد نظر در واژه‌نامه موجود است، اما به‌دلیل ضعف دانش زبانی نگارنده یا اشتباه وی، با توجه به مفهوم جمله نابه جا نوشته شده است (مانند واژه "عربی" در جمله "من در [ین شهر قریب هستم"). این‌گونه خطاهای در زبان فارسی کمیاب است و بیشتر در واژه‌های دخیل از زبان‌های دیگر، به‌ویژه زبان عربی، رخ می‌دهد و در زبان‌های غیرآوایی مانند انگلیسی نیز بسیار دیده می‌شود (مانند "piece" در "peace of cake"). این دسته از اشتباهات از جمله خطاهای ابهام گرامری^۶ و تشابه آوایی^۷ هستند[۳] و در محدوده این پژوهش نمی‌گنجند، زیرا تنها با تحلیل معنا و مفهوم جمله می‌توان آنها را کشف کرد. دسته دوم غلط‌ها نیز واژگانی هستند که در واژه‌نامه موجود نیستند و از این‌رو غلط به‌شمار می‌آیند. بدطور کلی خطاهای و غلط‌های متون و مستندات را می‌توان به سه



شناخت و ترکیب روش‌های مبتنی بر فاصله و ازگان برای انواع خطاهای احتمالی در متن، تعاریفی سیستماتیک را همراه با الگوریتم‌های تصحیح ارایه دهنده‌اند [۷].

۱-۲-۱- ویژگی‌های مساله

در مساله غلط‌های املایی، احتمال پیش آمدن غلط تایپی ناشی از اشتباہ کاربر، بیشتر از غلط‌های املایی متاثر از ضعف دانش زبانی نویسنده اصلی است؛ زیرا معمول تر آن است که فرض کنیم، کسی که خود متنی را به کمک رایانه تایپ می‌کند یا آن را به حروفچین می‌سپارد، به اندازه‌ای تسلط بر ازگان دارد که غلط‌های املایی متن وی ناچیز باشد، اما از سوی دیگر رخداد غلط تایپی امری رایج در هنگام ماشینی کردن متون است. این دو گونه اشتباہ، ماهیتی متفاوت دارند و شناخت و استفاده از ویژگی‌های آنها می‌تواند به ما در تصحیح املایی متون کمک کند. این نکته مساله‌ای است که در روش‌های کلاسیک و نرم‌افزارهای تجاری تصحیح غلط املایی - که اغلب بر پایه الگوریتم‌های جستجو استوارند - بدان توجه چندانی نشده و از این‌رو یکی از امتیازات این تحقیق به شمار می‌رود. البته باید در نظر داشت که نرم‌افزارهای تجاری از آن‌رو بر پایه روش‌های جستجو در فرهنگ و ازگان^۴ کار می‌کنند که هدف آنها دستیابی سریع به پاسخ با کمترین حافظه مصرفی است؛ در حالی که روش‌های آکادمیک بیشتر به دستیابی به بهترین پاسخ می‌اندیشند [۴].

از آنجا که جستجو در پایگاه داده‌ای با چند صد هزار واژه کاری وقت‌گیر است، روش‌های تصحیح غلط املایی چهار راهبرد اساسی را برای کاهش تعداد جستجوها به کار می‌گیرند که آنها را روش‌های تطبیق کامل^۵ می‌نامند [۴] :

بر پایه تعداد دفعات تکرار؛ بدین ترتیب که ازگانی که بیشترین مورد استفاده را در متون دارند، در یک فرهنگ و ازگان کوچکتر گردآوری شده، مصحح در هنگام یافتن غلط، نخست این پایگاه را می‌جوبد و در صورتی که واژه یافته نشد، پایگاه‌های و ازگان کم استفاده‌تر در اولویت‌های بعدی را نیز ارزیابی می‌کند.

بر پایه طول واژه؛ فرهنگ و ازگان به فرهنگ‌های ۲ تا n حرفی بخش می‌شود و در هنگام برخورد با واژه غلط k حرفی، نخست فرهنگ k حرفی و سپس در صورت نیافتن واژه صحیح، فرهنگ‌های $k+1$ و $k-1$ حرفی جستجو می‌گردد.

بر پایه نخستین حرف واژه‌نامه با ساختار درختی؛ در این ساختار، ریشه به تعداد حروف خط زبان مورد نظر (انگلیسی ۲۶ و فارسی ۳۲) گره دارد و هر گره نیز فرزندانی تا همین تعداد دارد. بنابراین برای یافتن واژه k حرفی به k جستجو نیاز است. اما روشی است که این روش زمانی خوب کار می‌کند که حروف آغازین واژه درست نگاشته شده باشند.

بر پایه فشرده‌سازی فرهنگ و ازگان؛ با منظور کردن این نکته که بسیاری از ازگان یک زبان ریشه‌ای مشترک دارند، می‌توان فرهنگ را به این و ازگان ریشه‌ای محدود ساخت و در عوض قوانینی برای تولید ازگان

نکاتی در آنها درج شده باشد که از دیدگاه دستوری یا مفهومی اشتباہ به نظر برسد و بنابراین اعمال تغییرات بر آنها خود موجب پیش آمدن اشتباهات بیشتر گردد. از سوی دیگر، با افزایش سواد و دانش عمومی و نیز با گسترش کاربرد رایانه در جهان، امروزه دیگر مشکل غلط املایی رو به کاهش و در عوض غلط تایپی رو به افزایش است.

۱-۱- پیشینه پژوهش

"فوروگوری" در سال ۱۹۹۰ با توجه به تفاوت ماهیت غلط‌های املایی انگلیسی، سیستمی کمکی برای تصحیح غلط‌های املایی زبان پاپی و انگلیسی، سیستمی کمکی برای تصحیح غلط‌های املایی زبان پاپی ها در هنگام نگارش به زبان انگلیسی پیشنهاد کرد که می‌توانست دقت نرم‌افزار تصحیح املایی کرک استار^۶ را از ۶۰٪ به ۷۵٪ برساند [۱۰]. در سال ۱۹۹۶ "شانگ" و "مرتل" با بررسی چند الگوریتم دیگر، روش ترای^۷ را برای تخمین نزدیکی دو رشته حروف بهبود و گسترش دادند [۱۱]. "لاونیه" در سال ۱۹۹۲ توانست تراشهای در قالب آرایه‌ای ۲-بعدی از ۶۹ پردازنده را به همراه یک تکنیک برنامه‌نویسی پویا برای تصحیح مقایسه‌ای رشته حروف بسازد که می‌توانست ۲۰۰ هزار واژه را در هر ثانیه تصحیح کند [۱۴]. "چاودوری" در سال ۲۰۰۲ با توجه به ویژگی‌های آوابی و غیرآوابی زبان هندی و با به کارگیری برخی از تکنیک‌ها - مانند روش n -گرام^۸ - توانست سیستم مصحح املایی برای زبان هندی و بنگالای^۹ ارایه کند که به دقت درستی ۹۵٪ دست یافت [۱۲]. در سال ۲۰۰۰ "هاج" و "اوستین" توانستند با ترکیب روش n -گرام و رویکرد فاصله همینگ سیستمی را فراهم آورند که برای یک پایگاه با تعداد محدودی از ازگان و با سرعتی قابل قبول، هر چهار گونه از غلط املایی را با دقیقی تا ۹۷/۵٪ تصحیح کند که به طور متوسط ۸ واژه جایگرین را پیشنهاد می‌کرد [۵]. در سال "ای" و همکارانش توانستند یک روش تطابق تخمینی واژه فازی^{۱۰} را در قالب مصحح املایی اختصاصی برای زبان چینی ارایه کنند که علاوه بر غلط‌های رایج، جایه‌جایی رشته را نیز تصحیح می‌کرد [۱۳]. "راش" و همکارانش در سال ۲۰۰۱ توانستند با ترکیبی از روش‌ها - مانند تطابق رشته به رشته و مدل مخفی مارکف - در تصحیح ازگان متون پژوهشی به دقت درستی ۹۸٪ دست یابند [۱۰-۱۳]. "هوانگ" و "پاورز" در سال ۲۰۰۱ با ترکیبی از روش‌ها و با در نظر گرفتن برخی غلط‌های تایپی توانستند در متون حجیم به دقت تصحیح ۷۴٪ دست یابند [۱۴].

"چرکاسکی" و همکارانش در سال ۱۹۹۰ توانستند با ترکیب برخی از شبکه‌های عصبی یادآور^{۱۱} مانند شبکه هاپفیلد^{۱۲} و شبکه‌های پس انتشار^{۱۳} با دیگر روش‌ها، سیستم غلط‌یاب املایی برای ازگان کوتاه (۵ تا ۷ حرفی) و بلند (۱۰ تا ۱۲ حرفی) با تعداد گره‌های ورودی برابر با توانی از ۲۶ (به تعداد حروف الفبای انگلیسی) به مقدار n در الگوریتم n -گرام و گره‌های خروجی به تعداد واژه‌های ذخیره شده، بسازند و به دقت ۱۵ تا ۱۰۰ درصد برای انواع خطأ و مقادیر n دست یابند [۴]. "گارفینکل" و همکارانش در سال ۲۰۰۲ کوشیدند تا با

۱-۲- غلط املایی در زبان فارسی

غلط املایی در زبان فارسی معمولاً به دو شکل پیش می‌آید: یکی این که نگارنده یکی از حروف هم صدا را به جای دیگری به کار برد. این گونه اشتباه معمولاً در مورد واژه‌های عربی دخیل در فارسی رخ می‌دهد. برای نمونه نویسنده ممکن است در واژه "اضطرار حرف "ض" را با "ز"، "ذ" یا "ظ" و حرف "ط" را با "ت" اشتباه بگیرد و واژه "تفصیل" را "تفیریط"، "تفیریض" یا بهشیوه‌های دیگر بنگارد. این اشتباه عموماً در مورد واژه‌های هم آوا در خط فارسی مانند "ت : ط"، "ح : ه"، "ز : ذ : ض : ظ"، "ق : غ" و گاهی نیز "ع" (مانند مواخذه : معاخذه) رخ می‌دهد.

دوم اشتباه در نگارش حروف صدادار (مصطفوت) یا آوای واژگان است که آن هم بیشتر در مورد واژگان دخیل از زبان‌های اروپایی - مانند انگلیسی و فرانسه - و نیز زبان عربی رخ می‌رهد. برای نمونه برای واژگان "Flat" و "Float" در "Flat" و "Float" دو نگارش "فlot" و "flet" را به کار می‌برند یا واژه "Robot" را "ربات"، "روبات" و یا "ربوت" و یا "مسئله" از عربی را "مساله"، "مسئله" یا بهشیوه‌های دیگر می‌نویسند. هر دوی این اشتباهات اغلب برای نویسنده‌گان کم سن و سال (تا مقطع دبستان و راهنمایی) پیش می‌آید و با افزایش سن و تجربه کم کم از بین می‌رود.

۲- غلط تایپی در زبان فارسی

اما غلط تایپی به سواد نویسنده ندارد، بلکه وایسته به مهارت و دقت حروف‌چین و موقعیت کلیدهای صفحه کلید است. شکل ۱، سه نمونه از موقعیت‌های حروف بر انواع صفحه کلید را نشان می‌دهد. بخشی از این ناهمگونی‌ها به دلیل تغییرات پیاپی در استانداردها و تعداد کمکهای صفحه کلید رخ داده و بخشی دیگر، به علت نبود اتفاق نظر میان استانداردگذاران درون و حتی بیرون از کشور بر سر تعیین جایگاه حروف فارسی بر دکمه‌ها پدید آمده و این ناهمانگی‌ها موجب افزایش خطاهای تایپی میان کاربران فارسی‌زبان گشته است.

در این زمینه برای شناخت بهتر ماهیت غلط تایپی فارسی با چند فرد خبره در حروف‌چینی گفت و گو شد و اطلاعات آنان به شکل زیر طبقه‌بندی گردید:

اشتباه در تایپ یک واژه، معمولاً بین کلیدهای همسایه رخ می‌دهد. مانند: "ج" و "ح"، "ن" و "ت" یا "ی" و "س".

اشتباه تایپی اغلب بین کلیدهای یک سطر از صفحه کلید رخ می‌دهد و اشتباه با کلیدهای سطر بالا یا پایین کم پیش می‌آید. برای نمونه حرف "ه" با "خ" یا "ع" اشتباه گرفته می‌شود، اما اشتباه آن با "ت" یا "ن" که در همسایگی سطر پایین‌تر قرار دارد، بسیار کمیاب است.

به دلیل استاندارد نبودن جای برخی از حروف بر صفحه کلید، برخی از حروف که نه شباخت آویزی با هم دارند و نه همسایه نزدیک یکدیگر هستند، با هم اشتباه می‌شوند. مانند "پ" و "ز".

ممکن بدان افزود. روش فشرده‌سازی فرهنگ، به دلیل ویژگی‌های واژه‌سازی زبانی، در زبان عربی بهتر از فارسی و انگلیسی کار خواهد کرد.

در مجموع می‌توان ویژگی‌های روش‌های تصحیح کلاسیک مبتنی بر جستجو و مبتنی بر شبکه عصبی را چنین بیان کرد:

□ با وجود این که شبکه عصبی در فاز یادگیری به زمان بالای نیاز دارد، اما در فاز آزمایش به سرعت نتیجه می‌گیرد و بنابراین از هر روش دیگری سریع‌تر است.

□ شبکه‌های عصبی مصنوعی طرفیتی محدود دارند؛ در نتیجه برای یادگیری تعداد زیادی واژه که بتواند حداقل نیاز برای مصحح املایی را برآورده سازد، با مشکل کاهش دقت روبرو می‌شوند.

□ یافتن بهترین شبکه با تعداد لایه و نرون، خود مساله‌ای نسبی و تجربی است و تصمیم‌گیری در مورد آن دقیق نیست. از این رو مساله باید با شبکه‌های گوناگون آزموده شود، تا بهترین ترکیب بدست آید.

□ با بهتر شدن و پیدایش انواع تازه و پرتوان تر شبکه‌های عصبی در آینده، نتایج کار نیز بهبود خواهد یافت.

□ در شبکه‌های عصبی با ساختار و طراحی ساده، با ورود یک واژه (درست یا نادرست) خروجی در نهایت یک واژه است، اما در روش‌های کلاسیک می‌توان بنا بر نسبت شباخت واژه نادرست به واژگان درست در فرهنگ واژگان، تعدادی واژه جایگزین را پیشنهاد کرد.

□ در صورتی که واژه‌ای بدون بررسی درست یا نادرست بودنش (بنا بر بود یا نبود آن واژه در فرهنگ واژگان) به شبکه عصبی وارد گردد، به دلیل ماهیت شبکه‌های عصبی که همواره در صدی از خط را به همراه دارد، ممکن است که واژه درست هم به واژه‌ای نادرست نگاشت شود.

در ادامه این مقاله نخست در بخش ۲ به تعریف و شناخت ماهیت غلط تایپی می‌پردازیم و تفاوت‌های آن با غلط املایی در زبان فارسی را بررسی خواهیم کرد. سپس در بخش ۳ به روش‌های ارایه شده برای تصحیح این اشتباهات خواهیم پرداخت و پس از آن در بخش ۴ نتایج تصحیح را بررسی خواهیم نمود و سرانجام در بخش ۵ مزایا و معایب جمع‌بندی و اهداف آینده بیان خواهد شد.

۲- غلط تایپی

همان گونه که در مقدمه گفته شد، در هنگام مواجهه با غلط‌های املایی، بهتر است فرض کنیم که غلط موجود ناشی از اشتباه تایپی بوده است و نه ضعف دانش زبانی نویسنده؛ زیرا کسی که متنی را به کمک رایانه تایپ می‌کند یا آن را به حروف‌چین می‌سپارد، به اندازه کافی بر زبان و واژگانش آشنایی دارد که غلط‌های املایی متن وی ناچیز باشد، اما از سوی دیگر رخداد غلط تایپی در هنگام تایپ متن بسیار رایج است.



'	1	2	3	4	5	6	7	8	9	0	-	=	Back Space
Tab	ض Q	ص W	ث E	ق R	ف T	غ Y	ع U	ه I	خ O	ح P	ز [ة]	
Caps Lock	ش A	س S	ى D	پ F	ج G	ه H	ت J	ن K	م L	ک ;	گ ‘		Enter
Shift	ظ Z	ط X	ج C	ر V	ذ B	د N	ئ M	و	،	.	.	/	Shift
پ ‘	1	2	3	4	5	6	7	8	9	0	-	=	Back Space
Tab	ض Q	ص W	ث E	ق R	ف T	غ Y	ع U	ه I	خ O	ح P	ز [ة]	
Caps Lock	ش A	س S	ى D	پ F	ج G	ه H	ت J	ن K	م L	ک ;	گ ‘		Enter
Shift	ظ Z	ط X	ج C	ر V	ذ B	د N	ئ M	و	،	.	.	/	Shift
پ ‘	1	2	3	4	5	6	7	8	9	0	-	=	٪ Back
Tab	ض Q	ص W	ث E	ق R	ف T	غ Y	ع U	ه I	خ O	ح P	ز [ة]	
Caps Lock	ش A	س S	ى D	پ F	ج G	ه H	ت J	ن K	م L	ک ;	گ ‘		Enter
Shift	ظ Z	ط X	ج C	ر V	ذ B	د N	ئ M	و	،	.	.	/	Shift

شکل (۱): سه نمونه از جانشانی حروف فارسی بر صفحه کلید.

جایه‌جایی کلید: به جای یکی از حروف اصلی، یکی از کلیدهای دیگر بر صفحه کلید فشرده می‌شود. کلید اشتباه اغلب نزدیک به کلید اصلی است. مانند: "کیومرث: کبومرث / گیومرث".

انتقال کلیدی: یکی از حروف اصلی، با حرف همسایه یا چند حرف قبل یا بعدی جایه‌جا تایپ می‌شود. مانند: "کیومرث" : کویمرث / کومریث".
غلط املایی (جایه‌جایی ویژه): به‌جای یکی از حروف اصلی، یکی از دیگر کلیدهای صفحه کلید فشرده می‌شود؛ به‌گونه‌ای که در زبان فارسی آن حرف با حرف اصلی هم صدا است مانند: "کیومرث" : کیومرس". این گونه غلط ناشی از ضعف دانش زبانی حروف‌چین است و بسیار کم رخداده، زیرا متن اصلی پیش روی حروف‌چین است.
فاصله اضافه: یک کاراکتر فاصله نابه‌جا درج می‌گردد. مانند: "کیومرث" : کی و مرث". این غلط به‌ویژه در برخی حالت‌های ویژه مانند "ها" ^۳ ای جمع و "می" استمرار و اغلب به‌جای فاصله مجازی ^۴ (کوتاه) رخداده. مانند: "وازه‌ها" و "می‌رود" : می رود". این غلط گونه‌ای ویژه از غلط‌های درج است.

فاصله کم؛ کاراکتر فاصله میان دو واژه جداگانه درج نمی‌گردد. مانند: "کیومرث سیامک"؛ این غلط گونه‌ای ویژه از غلط‌های حذف است.

بر این پایه با نمونه برداری و تخمین آماری غلط‌های تایپی زبان فارسی در هنگام کار حروف‌چین خبره به شرح جدول ۱ به دست آمد. برای بررسی صحت و سقم هر یک از موارد این جدول از هر یک از افراد خواسته شد تا متنی را بدون استفاده از کلید بازگشت^{۲۸} در هنگام

ارتباطی ناچیز هم در اشتباه بین حروف هم‌صدا وجود دارد.^{۳۶} مانند: "ض" و "ظ" یا "غ" و "ق".

گاه در حروفی که نیاز به فشردن کلید شیفت دارند، آن حرف با حرف اصلی جایه‌جا می‌شود. مانند "ا'" به جای آ"' (که البته این مورد غلط املایی به شمار نمی‌رود)، یا "ی" و "ط"، "ز" به جای "ژ" (در برخی از صفحه کلیدها) و "س" و "ش" به جای هر یک از اعراب.

گاه دو حرف از یک واژه جایه‌جا تایپ می‌شود، مانند: "زیبا : زیبا". این اتفاق نیز بیشتر در مورد حروف همسایه رخ می‌دهد.

گاه غلط املایی بدین شکل رخ می‌دهد که یکی از حروف واژه زده نمی‌شود یا یکی از کلیدهای جانبی یا خود کلید یکی از حروف واژه، اضافه زده می‌شود. مانند: "هشدار : هشدار : هشار".

درج حرف: در هنگام نگارش، یک حرف اضافه به اشتیاه تایپ می‌شود.
مانند: "کیومرث : کیومنرث".

تکرار حرف (درج حرف همسایه): یعنی در هنگام نگارش، یکی از حروف اصلی به اشتیاه دو بار تایپ می‌شود. مانند: "کیومرث" : کیومرمث". این گونه غلط خود نوعی درج حرف است که بیش از گونه‌های دیگر رخ می‌دهد.

حذف حرف: در هنگام نگارش، یکی از حروف اصلی تایپ نمی‌شود.
مانند: "کیومرث : کیمرث".

"*holy*" در اینجا "y" بهجای "j" آمده در صورتی که این دو حرف در سطح جدا قرار دارند و احتمال اشتباه تایپی این دو به دلیل دوری از هم بسیار ناچیز است.

همان‌گونه که مشاهده می‌شود، در این نرمافزار واژه "golf" اصلاً پیشنهاد نشده و بهترین پیشنهاد در اولویت سوم قرار دارد.

به عنوان نمونه‌ای دیگر، در صورتی که کاربر بخواهد واژه "sear" را تایپ کند، اما به اشتباه بهجای حرف "s"، حرف کناری آن یعنی "a" و بنابراین "ear" را تایپ کند، نرمافزار ورد واژگان زیر را به ترتیب اولویت پیشنهاد می‌دهد:

: در اینجا نرمافزار، حرف "r" دو خانه به عقب برده و فرض کرده که کاربر جای آن را دو خانه جایه‌جا تایپ کرده که البته این اشتباه چندان محتمل نیست.

: در اینجا فرض شده که کاربر بهجای حرف "i" دو حرف "ea" را به اشتباه تایپ کرده و البته این اشتباه هم چندان منطقی نیست.

: در اینجا "afar" بهجای "e" آمده در صورتی که این دو حرف در سطح جدا قراردارند و احتمال اشتباه تایپی این دو، به دلیل دوری از هم، ناچیز است.

: در اینجا هم مانند مورد قبلی حرف "j" بهجای "e" آمده در صورتی که این دو حرف در دو سطح جدا قرار دارند و احتمال اشتباه تایپی این دو اندک است.

: باز هم در اینجا حرف "g" بهجای "e" آمده در صورتی که این دو حرف در دو سطح جدا قراردارند و احتمال اشتباه تایپی این دو کم است.

همان‌گونه که دیده می‌شود، در این نرمافزار واژه "sear" که منطقی‌ترین گزینه است، اصلاً پیشنهاد نشده و بهترین پیشنهاد در رده سوم قرار دارد.

اکنون نمونه‌ای دیگر را در نظر گیرید: اگر کاربر بخواهد واژه "that" را تایپ کند، اما به اشتباه بهجای حرف "t" دوم، حرف کناری آن یعنی "u" و بنابراین "thay" را تایپ کند، نرمافزار ورد واژگان زیر را به ترتیب اولویت می‌دهد:

: در اینجا "e" بهجای "a" آمده در صورتی که این دو حرف در دو سطح جدا قراردارند و احتمال اشتباه تایپی این دو ناچیز است.

: در این مورد فرض شده که کاربر بهجای حرف "e"، به اشتباه دو حرف "ay" را تایپ کرده که امری بسیار بعید است.

: در اینجا "tray" بهجای "tray" آمده در صورتی که این دو حرف شده در صورتی که این دو حرف در دو سطح جدا قراردارند و احتمال اشتباه تایپی این دو اندک است.

: در اینجا "that" بهجای "y" آمده و این دو حرف بر کلیدهای مجاور قرار دارند و این انتخاب معقول‌ترین انتخاب بهنظر می‌رسد.

: در این مورد هم "w" بهجای "y" آمده، در صورتی که این دو حرف در همسایگی دوری نسبت بهم جای دارند و احتمال اشتباه تایپی این دو کم است.

رخداد خطأ، تایپ کنند. نتایج به دست آمده در تایید موارد فوق بود و خلاف آن را نشان نمی‌داد.

جدول (۱): انواع غلط‌های تایپی و نسبت رخداد آنها در زبان فارسی.

ردیف	نوع غلط	گونه اصلی	درصد
۱	درج حرف	درج	۸
۲	تکرار حرف	درج	۱۴/۵
۳	حذف حرف	حذف	۱۹
۴	جایه‌جایی کلید	جایه‌جایی	۳۹/۵
۵	انتقال کلید	انتقال	۵
۶	غلط املایی	جایه‌جایی	۰/۵
۷	فاصله افزوده	درج	۲/۵
۸	فاصله کم	حذف	۵
۹	موارد دیگر	ترکیبی	۶/۵
۱۰	نسبت حروف غلط به حروف درست	-	۱/۹۵
۱۱	نسبت واژگان غلط به واژگان درست	-	۵/۳۱

۳-۲- تصحیح املایی بر پایه فاصله

امروزه بیشتر نرمافزارهای مصحح غلط املایی که به صورت تجاری در بازار ارایه می‌شوند، بر پایه فاصله (ویرایشی، همینگ، لونشتن و ...) کار می‌کنند، که البته در بسیاری از موارد - به ویژه در هنگام رخداد غلط‌های جایه‌جایی - نتیجهٔ دلچسبی به همراه ندارند. در این میان، مصحح املایی نرمافزار ورد از شرکت مایکروسافت^۹ یکی از کارآمدترین و پر طرفدارترین تصحیح‌کننده‌های املایی متن را ارایه کرده است. با وجود آن که این نرمافزار خود برای تایپ و صفحه‌بندی طراحی شده، به غلط تایپی توجهی ندارد، بلکه احتمالاً با توجه به پایگاه داده، تک‌تک واژگان تایپ شده را با آن پایگاه داده می‌سنجد و در صورتی که آن واژه در پایگاه موجود نباشد، آن را مشخص کرده و به کمک برخی قوانین کلاسیک برای یافتن نزدیکترین واژه، چند واژه را به ترتیب اولویت از روی سازگاری به کاربر پیشنهاد می‌دهد.

برای نمونه، در صورتی که کاربر بخواهد واژه "golf" را تایپ کند، اما به اشتباه بهجای حرف "g"، حرف کناری آن "h" و بنابراین "half" را تایپ کند، نرمافزار ورد واژگان زیر را به ترتیب اولویت پیشنهاد می‌دهد:

: در اینجا "a" بهجای "o" آمده در صورتی که احتمال اشتباه تایپی این دو به دلیل دوری از هم بسیار ناچیز است.

: در اینجا "o" بهجای "l" آمده در صورتی که این دو حرف در دو سطح جدا قرار دارند و احتمال اشتباه تایپی این دو به دلیل دوری از هم بسیار ناچیز است.

: در اینجا "d" بهجای "f" آمده و این دو حرف بر کلیدهای مجاور قرار دارند و این انتخاب معقول‌تر به نظر می‌رسد.

: در این مورد هم "e" بهجای "f" آمده، در صورتی که این دو حرف در دو سطح جدا قراردارند و احتمال اشتباه تایپی این دو ناچیز است.



ض 1F00	ص 0F00	ث 0700	ق 0300	ف 0100	غ 0000	ع 0040	ه 00C0	خ 00E0	ح 00F0	ج 00F8	ج 00FC
ش 9F01	س 8F01	ی 8701	ب 8301	ل 8101	ا 8001	ت 8041	ن 80C1	م 80E1	ک 80F1	گ 80F9	
ظ DF03	ط CF03	ز C703	ر C303	ذ C103	د C003	ئ C043	و C0C3		پ C706	ژ C702	

شکل (۲): کد اختصاص داده شده به هر یک از حروف فارسی بنا بر فاصله همسایگی.

کناری آن به گونه‌ای تنظیم شده است که قواعد شاعع همسایگی رعایت شود. برای نمونه، کلید کناری آن حرف "ع" کد ۰۰۴۰_(H) یا ۰۰۰۰۰۰۰۱۰۰۰۰۰۰₍₂₎ بیت اختلاف دارد و حرف "ه" نیز کد ۰۰۰۰۰۰۰۱۱۰۰۰۰۰۰₍₂₎ یا ۰۰۰۰۰۰۰۱۱۰۰۰۰۰۰_(H) را می‌گیرد که با حرف "ع" یک اختلاف و با حرف "غ" دو اختلاف دارد و این اختلافات نشانگر فاصله همسایگی میان آنها است.

شکل ۲ حروف فارسی و جایگاه آن بر صفحه کلید را به مراره کدهای شانزده‌شانزدهی مربوطه نشان می‌دهد. ورودی شبکه عصبی هر یک از حروف فارسی با شرح بیشتر در جدول ۲ درج گردیده است.

جدول (۲): حروف فارسی و کلید و کد آنها بنا بر فاصله همسایگی.

کد دودوبی (ورودی شبکه عصبی)	کد هنگر	کد هنگر	کلید	حروف فارسی	ردیف
.....1	8001	h	ا	۱	
.....11.....1	8301	f	ب	۲	
.....111.....110	C706	\ یا /	پ	۳	
.....11.....1.....1	8041	j	ت	۴	
.....111.....	0700	e	ث	۵	
.....1111.....	00F8	[ج	۶	
.....11111.....	00FC]	چ	۷	
.....11111.....	00F0	p	ح	۸	
.....111.....	00E0	o	خ	۹	
.....11.....11	C003	n	د	۱۰	
.....11.....11	C103	b	ذ	۱۱	
.....11.....11	C303	v	ر	۱۲	
.....111.....11	C703	c	ز	۱۳	
.....111.....10	C702	C یا \	ژ	۱۴	
.....1111.....1	8F01	s	س	۱۵	
.....11111.....1	9F01	a	ش	۱۶	
.....1111.....	0F00	w	ص	۱۷	
.....11111.....	1F00	q	ض	۱۸	
.....11111.....11	CF03	x	ط	۱۹	
.....11111.....11	DF03	z	ظ	۲۰	
.....1.....	0040	u	ع	۲۱	
.....1.....	0000	y	غ	۲۲	
.....1.....	0100	t	ف	۲۳	
.....11.....	0300	r	ق	۲۴	
.....1111.....1	80F1	:	ک	۲۵	
.....11111.....1	80F9	'	گ	۲۶	
.....11.....1	8101	g	ل	۲۷	
.....111.....1	80E1	l	م	۲۸	
.....11.....1	80C1	k	ن	۲۹	
.....11.....11	C0C3	,	و	۳۰	
.....11.....	00C0	i	ه	۳۱	
.....111.....1	8701	Z و d	ی	۳۲	
.....11.....11	C043	m	ئ	۳۳	

همان‌گونه که دیده می‌شود، در این مثال گزینه درست در اولویت چهارم پیشنهاد شده است و برخی از گزینه‌ها غیر منطقی به نظر می‌رسند.

۳- تصحیح غلط تایپی به کمک شبکه عصبی

شبکه عصبی مصنوعی روش و ابزاری محاسباتی است که بر روی مقادیر عددی کار می‌کند. بنابراین نخستین کار در چنین مسائلای این است که مفاهیم کیفی مانند حروف الفبا و دوری یا نزدیکی آنها به یکدیگر، به مقادیر کمی و عددی قابل محاسبه تبدیل شود. از این‌رو به هر یک از حروف صفحه کلید یک کد عددی به گونه‌ای اختصاص می‌یابد که مفهوم دوری / نزدیکی کلیدها را نیز در برداشته باشد.

۳-۱- تعریف فاصله (شعاع همسایگی) کلیدها

برای داشتن معیاری برای سنجش دوری یا نزدیکی دو حرف بر صفحه کلید، فاصله یا شاعع همسایگی دو کلید را چنین تعیین می‌کنیم:

- فاصله هر کلید با کلید همسایه چپ یا راستش برابر با یک است.
- فاصله هر کلید با کلید kام کنارش برابر با k است.
- فاصله هر کلید با کلید همسایه‌اش در سطر بالا یا پایین برابر با ۲ است.

با توجه به این‌که در طولانی‌ترین سطر صفحه کلید ۱۲ حرف فارسی گنجانیده شده، یک کد ۱۲ بیتی می‌تواند فاصله کلیدهای هر ردیف را نشان دهد و از آنجا که ۳ ردیف حروف بر صفحه کلید موجود است و هر دو ردیف باید دست کم فاصله ۲ با هم داشته باشند، طول این کد به ۱۶ بیت افزایش می‌یابد. بنابراین با اختصاص این کدها (شکل ۲ و جدول ۲) می‌توان به هر یک از کلیدهای صفحه کلید یک کد ۱۶ بیتی منحصر به فرد اختصاص داد.

۲-۳- کدهای ورودی شبکه عصبی

برای تبدیل هر یک از حروف الفبای فارسی به یک مقدار عددی که به عنوان ورودی برای شبکه عصبی مصنوعی مناسب باشد، به جای این که هر یک از حروف با یکی از اعداد ۰ تا ۳۲ متناظر شود، به‌هر یک از آنها یک کد ۱۶ بیتی اختصاص یافته است. بدین منظور حرف "غ" مرکز کلیدهای فارسی با کد دودوبی ۰۰۰۰۰۰۰۰۰۰۰۰۰۰۰۰₍₂₎ و شانزده‌شانزدهی ۰۰۰۰_(H) در نظر گرفته شده است و کد کلیدهای

البته در عمل، مقادیر ورودی شبکه‌های عصبی نه به‌طور باینری، بلکه

به‌صورت دوقطی در نظر گرفته شده‌اند تا بازدهی شبکه بهتر باشد.

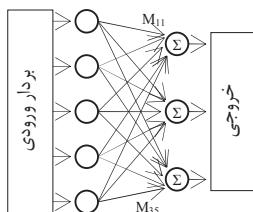
پس هر یک از حروف واژگان، یک کد با ۱۶ ورودی ۱ یا ۱- را به‌خود

اختصاص می‌دهد و از این‌رو در لایه ورودی شبکه عصبی، ۱۶ نرون به

هر یک از آنها اختصاص خواهد یافت.

$$\begin{aligned} [11010] &= "اب" \\ [01110] &= "بر" \\ [10101] &= "یاد" \end{aligned}$$

شکل ۳-۳-(الف): کدگذاری و ذخیره واژگان به‌عنوان ورودی شبکه



شکل ۳-۳-(ب): شبکه عصبی حافظه یادآور

$$R = M \cdot S = \begin{bmatrix} 11010 \\ 01110 \\ 10101 \end{bmatrix} \begin{bmatrix} 1 \\ 0 \\ 1 \\ 0 \\ 1 \end{bmatrix} = \begin{bmatrix} 1 \\ 1 \\ 3 \end{bmatrix}$$

شکل ۳-۳-(پ): بازیابی واژگان در شبکه حافظه یادآور به‌کمک ماتریس همبستگی

شبکه عصبی مصنوعی هایپفیلد شاخص ترین نوع شبکه‌های یادآور و از گونه حافظه خودی‌یادآور^{۴۴} است که کاربردهای بسیاری در دانش‌های گوناگون دارد. این شبکه نخستین بار در سال ۱۹۸۲ میلادی توسط جان هایپفیلد^{۴۴} ارایه شد^[۱۵]. وی در سال ۱۹۸۵ این شبکه را به‌کمک تنک^{۴۸} گسترش داد و با آن مساله فروشنده دوره‌گرد^{۴۹} را با در نظر گرفتن ده شهر و صد نرون با کارایی بهتری حل کرد. فروشنده دوره‌گرد یک مساله بهینه‌سازی معروف است که در زمرة مسائل بسیار مشکل قرار می‌گیرد و با روش‌های معمولی نمی‌تواند در زمانی معقول پاسخی بهینه را به دست آورد. هایپفیلد و تنک مساله خود را تا ۳۰ شهر با موفقیت گسترش دادند و پس از گذشت ۲۰ سال هنوز هم روش آنها جزو بهترین الگوریتم‌های شبکه عصبی برای حل مساله فروشنده دوره‌گرد است^[۱۶].

شبکه عصبی هایپفیلد در قالب یک سیستم پویا توسط یکتابع انرژی که باید تعادلی میان اهداف تابع مساله - که باید حداقل شود - ایجاد می‌کند^[۱۷]. پس از این موفقیت شبکه عصبی هایپفیلد، بسیاری از مسائل مهندسی در قالب تابع انرژی که باید کمینه شود، ارایه گردید. چنین راه حلی بسیار جذاب است، زیرا پردازش موازی را برای حل مسائل امکان پذیر می‌سازد^[۱۸].

بنابراین این شبکه می‌تواند با یادگیری و حفظ تعدادی واژه در حافظه خود، ورودی همراه با نویز یا همان غلط املایی را به نزدیک ترین الگو نگاشت کند و صورت درست واژه را در خروجی تداعی کند. ساختار شبکه هایپفیلد در شکل ۴ و تابع فعالیت^{۵۰} آن در شکل ۵

۳-۳- فرهنگ واژگان به کار رفته در شبیه‌سازی

در شبیه‌سازی این پژوهش، ۶۰۰ واژه، ۴، ۵ و حرفی به تعداد مساوی، برای آموزش شبکه‌ها و ۳۰۰ واژه نیز برای آزمایش به کار گرفته شده‌اند. بیشتر واژگان از اسامی اعلام (کسان و جای‌ها) یا کلمات پرکاربرد در زبان فارسی و از فرهنگ واژه‌های معتبر فارسی، مانند "لغت‌نامه دهخدا" و "فرهنگ معین"، برگزیده شده‌اند.

۴- شبیه‌سازی

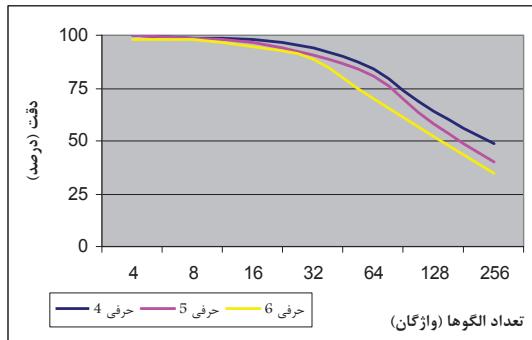
شبیه‌سازی و آزمون مقادیر و واژگان توسط نرم‌افزار متلب^{۴۰}، به‌همراه نرم‌افزار کمکی تحت زبان ویژوال بیسیک^{۴۱} برای تهیه مقادیر ورودی دو نوع شبکه عصبی هایپفیلد و پرسپترون چند لایه^{۴۲} انجام شد. الگوهای ورودی در دسته‌های^{۴۳} ۵ و ۶ حرفی به شبکه عصبی وارد شد و پس از آموزش شبکه، هر بار الگوهای ورودی به‌دفعات و با تغییر تصادفی یکی از حروف واژه آزموده شد.

۴-۱- شبیه‌سازی با شبکه عصبی هایپفیلد

هنگامی که درباره تصحیح املایی واژگان سخن می‌گوییم، شبکه‌های گونه حافظه یادآور^{۴۴} نخستین شبکه‌های مناسب به‌نظر می‌رسند. شبکه حافظه یادآور سیستمی است که می‌تواند داده‌های ذخیره شده (الگوهای^{۴۴}) را حتی با دیدن ورودهای همراه با غلط بازیابی (یادآوری) کند^[۴۵]. به‌عنوان نمونه‌ای کوچک برای آشنازی با روش کار این شبکه خط فارسی را به ۵ حرف ("ا", "ب", "ر", "ش" و "ی") محدود می‌کنیم و می‌خواهیم سه واژه "اب", "آرش" و "بیش" را بیاد گرفته، یادآوری نماییم. بر این پایه هر واژه را می‌توان بنا بر حروف مورد استفاده در آن به‌صورت یک ماتریس ۵ عنصری مانند شکل ۳-الف نمایش داد. این مساله به شبکه‌ای با ۵ گره ورودی و ۳ گره خروجی مانند شکل ۳-ب نیاز دارد. پاسخ آزمون واژه‌ای مانند "اید" - که همان واژه "یاد" با رخداد خطای انتقال است - از ضرب دو ماتریس فرهنگ واژگان در داده آزمایشی به‌دست آمده، سطر خروجی بزرگتر، پاسخ سیستم است. ماتریسی که فرهنگ واژگان را در بر دارد، ماتریس همبستگی^{۴۶} می‌نامند.



نشان داده شده است.



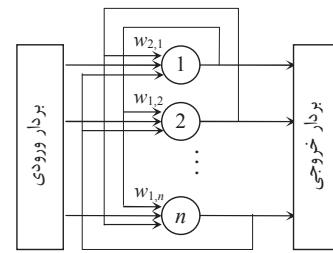
شکل (۶): نتایج شبکه هایپفیلد در تصحیح غلط تایپی واژگان ۴، ۵ و ۶ حرفی

همان گونه که در این جدول مشاهده می شود، با افزایش تعداد الگوهای دقت شبکه کاهش یافته است. همچنین افزایش تعداد حروف واژگان مورد آزمایش اندکی از دقت شبکه کاسته و برای یادگیری و تصحیح تنها ۲۵۶ واژه به بیش از ۲۵۰۰۰ دور گردش شبکه (اپک^۳) در هنگام آموزش نیاز است در صورتی که دقت تصحیح شبکه از ۵۵٪ فراتر نمی رود.

۲-۴- شبیه‌سازی با شبکه عصبی پرسپترون چندلایه

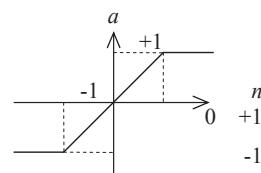
شبکه عصبی پرسپترون چند لایه، شبکه‌ای است که در اصل برای مسایل دسته‌بندی و تخمين تابع طراحی شده و در این کارها از موفق‌ترین شبکه‌های عصبی بوده و توانسته است با استفاده از قانون انتشار خطا به عقب^۳ بسیاری از مسایل غیرقابل حل توسط شبکه‌های دیگر را حل کند [۱۹]. اما با توجه به این‌که کار اصلی این شبکه دسته‌بندی و تخمين تابع است، بنابراین به‌نظر نمی‌رسد که در مساله تصحیح غلط تایپی - که می‌توان گفت مساله یادآوری است - به‌کار آید.

نکته ساده‌ای که باید برای کارآمدی شبکه پرسپترون چند لایه در مساله تصحیح واژه غلط به کار گرفت، این است که: اگر هر یک از الگوهایی که باید یاد گرفته شود را بایک دسته در فضای n بعدی متناظر کنیم، برای یادگیری 2^n واژه، مساله به مساله دسته‌بندی 2^n دسته‌ای با n نرون در لایه خروجی شبکه متناظر می‌شود. پس برای نمونه، برای یادگیری مجموعه‌ای از واژگان با حدود ۶۵ هزار واژه، به شبکه‌ای با ۱۶ گره خروجی نیاز است. جدول ۴ نمونه‌ای از این روش را برای دسته‌بندی ۸ واژه چهار حرفی نشان می‌دهد. واژه‌های مورد آزمایش به‌عمدّه‌ای برگزیده شده‌اند که به‌هم نزدیک باشند و با دگرگونی یکی-دو حرف به واژه‌ای دیگر در واژه‌نامه تبدیل شوند، تا نتایج تصحیح بر روی آنها بهتر نمایان شود.



شکل (۴): نمای کلی شبکه هایپفیلد.

اما دو نقص بزرگ شبکه هایپفیلد ظرفیت پایین آن (در حدود ۱۵٪ اندازه شبکه یا تعداد گره‌ها) و همچنین هزینه محاسباتی بالای آن است [۱۹]. از این‌رو انتظار می‌رود که با افزایش تعداد الگوهای یاد داده شده، دقت شبکه کاهش یابد و بنابراین شبکه هایپفیلد نامزدی مناسب برای سیستم غلط‌یاب املایی نیست [۴].

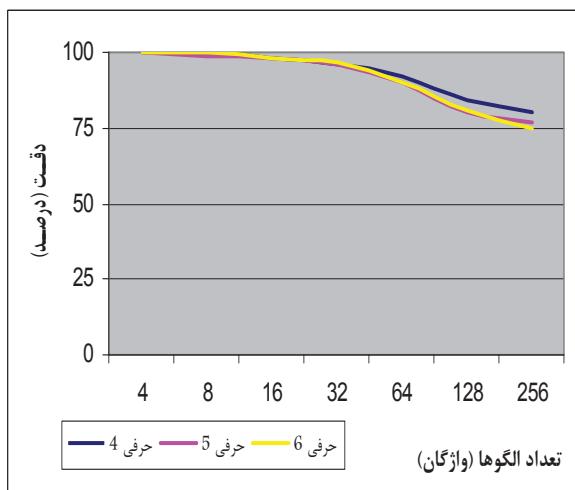


شکل (۵): تابع فعالیت (دو قطبی) شبکه هایپفیلد.
نتایج این آزمون و شبیه‌سازی آن در جدول ۳ و شکل ۶ نمایش داده شده است.

جدول (۳): ارزیابی نتایج شبکه هایپفیلد در تصحیح غلط تایپی واژگان.

ردیف	تعداد الگوها (واژگان)	تعداد حروف واژه	حداقل اپک لازم برای همگرا شدن شبکه	دقت (درصد)
۱	۴	۴	۲۰	۱۰۰
۲	۴	۵	۵۰	۱۰۰
۳	۴	۶	۶۰	۹۸
۴	۸	۴	۱۰۰	۹۹
۵	۸	۵	۱۵۰	۹۹
۶	۸	۶	۲۰۰	۹۸
۷	۱۶	۴	۳۰۰	۹۸
۸	۱۶	۵	۵۰۰	۹۷
۹	۱۶	۶	۸۰۰	۹۵
۱۰	۳۲	۴	۱۰۰۰	۹۴
۱۱	۲۲	۵	۲۵۰۰	۹۱
۱۲	۳۲	۶	۴۰۰۰	۸۹
۱۳	۶۴	۴	۹۰۰۰	۸۴
۱۴	۶۴	۵	۱۲۰۰۰	۸۱
۱۵	۶۴	۶	۱۲۰۰۰	۷۰
۱۶	۱۲۸	۴	۲۰۰۰۰	۶۴
۱۷	۱۲۸	۵	۲۰۰۰۰	۵۸
۱۸	۱۲۸	۶	۲۰۰۰۰	۵۲
۱۹	۲۵۶	۴	۲۰۰۰۰	۴۹
۲۰	۲۵۶	۵	۲۰۰۰۰	۴۰
۲۱	۲۵۶	۶	۲۵۰۰۰	۳۵

جدول (۴): نمونه‌ای از ۸ واژه آزموده با شبکه عصبی و کد ویژه آنها.



شکل (۷): نتایج شبکه پرسپترون چند لایه در تصحیح واژگان ۴، ۵ و ۶ حرفی

جدول (۶): مقایسه میانگین بازدهی دو شبکه هاپفیلد و پرسپترون چند لایه

ردیف	تعداد الگوها	هاپفیلد	پرسپترون چندلایه
۱	۴	۹۹/۳	۱۰۰
۲	۸	۹۸/۷	۹۹/۷
۳	۱۶	۹۶/۷	۹۸
۴	۳۲	۹۱/۳	۹۶/۳
۵	۶۴	۷۸/۳	۹۰/۷
۶	۱۲۸	۵۸	۸۱/۷
۷	۲۵۶	۲۵۶	۷۷/۳

۳-۴- ارزیابی نتایج

جدول‌های ۳ و ۵ به خوبی نشان می‌دهند که کارکرد شبکه پرسپترون چند لایه بسیار بهتر از شبکه هاپفیلد بوده است. جدول ۶ و شکل ۸ نیز این تفاوت را به صورت میانگین برای واژگان ۴ تا ۶ حرفی نشان می‌دهند.

ردیف	واژه	کد باینری (دسته)
۱	نامی	...
۲	مانی	۰۰۱
۳	مینا	۰۱۰
۴	امین	۰۱۱
۵	میان	۱۰۰
۶	ایمن	۱۰۱
۷	نیام	۱۱۰
۸	نیما	۱۱۱

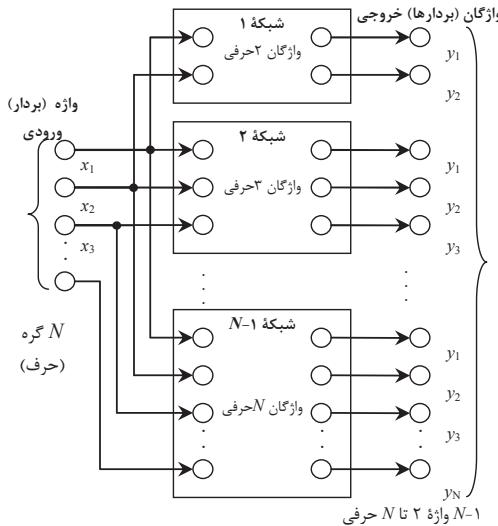
روشن است که این مجموعه نیاز به شبکه‌ای با ۳ نرون خروجی دارد. نتایج آزمایش با این روش با تعداد واژگان مختلف در جدول ۵ و شکل ۷ آمده است. در این آزمون‌ها، تعداد اپک‌ها به گونه‌ای تنظیم شده که خطای شبکه به صفر برسد، اما از آنجا که گاه این امر میسر نمی‌شود و با توجه به این که کار شکله با تعداد دسته، لایه و نرون بسیار سنگین می‌شود (آموزش مورد ردیف ۱۹ بیش از ۳ روز زمان برده است!) کار با اپک‌های کمتر خاتمه یافته است. همان‌گونه که مشاهده می‌شود، با افزایش تعداد الگوهای دقت شبکه کاهش می‌یابد. همچنین افزایش تعداد حروف واژگان مورد آزمایش اندکی از دقت شبکه کاسته است.

جدول (۵): ارزیابی نتایج شبکه پرسپترون چندلایه در تصحیح غلط تایپی واژگان.

ردیف	تعداد دسته‌ها	تعداد حروف	حداقل اپک (لاز) یا نرون‌های لایه	تعداد	دققت (درصد)
۱	۴	۴	۱۰	۴	۱۰۰
۲	۴	۵	۱۵	۲	۱۰۰
۳	۴	۶	۲۰	۱	۱۰۰
۴	۸	۴	۵۰	۱	۱۰۰
۵	۸	۵	۶۰	۱	۹۸
۶	۸	۸	۸۰	۱	۱۰۰
۷	۱۶	۴	۱۰۰	۲	۹۸
۸	۸	۵	۱۴۰	۱	۹۸
۹	۱۶	۶	۲۰۰	۲	۹۸
۱۰	۲۲	۴	۵۰۰	۲	۹۶
۱۱	۳۲	۵	۸۰۰	۲	۹۶
۱۲	۳۲	۶	۱۳۰۰	۲	۹۸
۱۳	۶۴	۴	۲۰۰۰	۲	۹۲
۱۴	۶۴	۵	۳۰۰۰	۶	۹۰
۱۵	۶۴	۶	۴۰۰۰	۶	۹۰
۱۶	۱۲۸	۴	۱۰۰۰۰	۲	۸۴
۱۷	۱۲۸	۵	۱۲۰۰۰	۲	۸۰
۱۸	۱۲۸	۶	۱۵۰۰۰	۲	۸۱
۱۹	۲۵۶	۴	۱۰۰۰۰	۳	۸۰
۲۰	۲۵۶	۵	۱۲۵۰۰	۳	۷۷
۲۱	۲۵۶	۶	۱۵۰۰۰	۳	۷۵



در این شبکه، بردار واژه k حرفی به شبکه k و در صورت نیاز به شبکه‌های $k+1$ یا $k-1$ وارد می‌شود و بنابراین پاسخ تنها در خروجی شبکه‌های k ، $k+1$ و $k-1$ ظاهر می‌گردد.



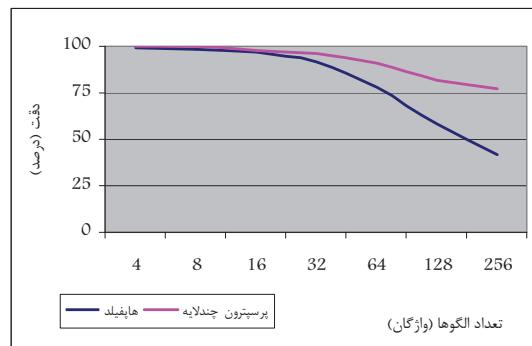
شکل (۹): به کار گیری شبکه‌های موازی در پردازش واژگان ۲ تا N حرفی

جدول ۷ نتایج آزمون روش تقسیم‌بندی بر پایه طول واژه را بر واژگان ۴ تا ۶ حرفی با ۳ زیر شبکه که هر یک از آنها ۴ تا ۲۵۶ واژه را یاد گرفته‌اند، نشان می‌دهد. در صورتی که واژه درست مرتبط با واژه غلط مورد آزمایش، دست کم در یکی از خروجی‌های سیستم ظاهر شده باشد، عملکرد شبکه درست در نظر گرفته شده است.

جدول (۷): افزایش ظرفیت با روش تقسیم‌بندی طول واژه.

تعداد شبکه	ردیف الگوهای هر واژگان	تعداد کل واژگان	پرسپترون چندلایه	هایپفیلد	تعداد کل پرسپترون
۱	۱	۱			۱
۲	۳	۲۴	۸	۹۸/۲	۹۹/۵
۳	۶	۴۸	۱۶	۹۵/۵	۹۷/۲
۴	۱۲	۹۶	۳۲	۸۹/۱	۹۴/۱
۵	۲۴	۱۹۲	۶۴	۷۳	۸۸/۷
۶	۴۸	۳۸۴	۱۲۸	۵۰/۵	۷۹/۸
۷	۹۶	۷۶۸	۲۵۶	۳۵	۷۲

تقسیم‌بندی بر پایه نوع واژه: در هنگام آموزش و دسته‌بندی واژگان k حرفی، می‌توان واژگان را به جای یک شبکه، در N شبکه پخش کرد. این ترفند ظرفیت کل سیستم را به طور متوسط N برابر افزایش می‌دهد و گذشته از آن، امکان پیشنهاد چندین واژه جایگزین در خروجی را هم پیدید می‌آورد. روشن است که گره‌های لایه ورودی و خروجی این شبکه نیز، که در شکل ۱۰ نمایش داده است، هر یک



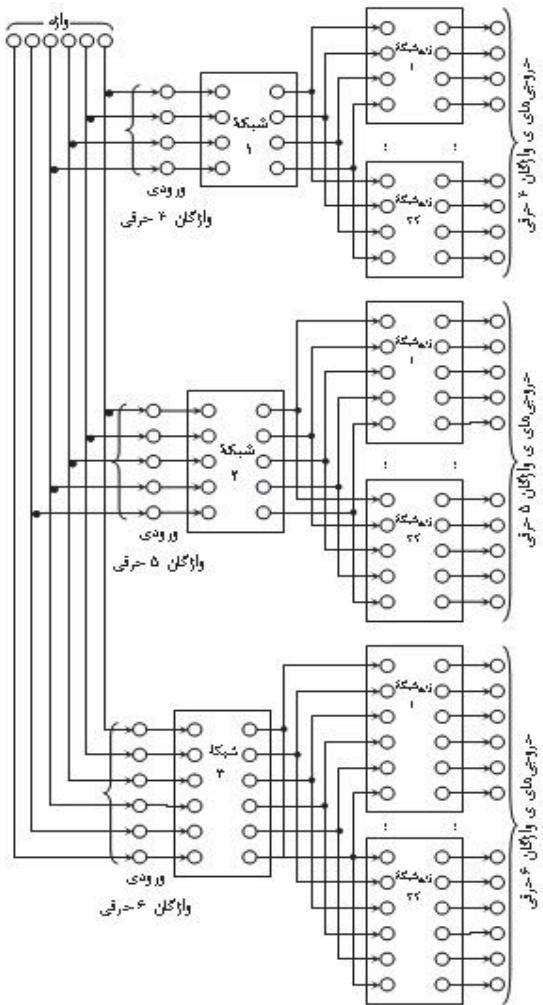
شکل (۸): بازدهی دو شبکه هایپفیلد و پرسپترون چند لایه.

۴-۴- ترمیم مشکل ظرفیت شبکه

همان‌گونه که از نتایج شبیه‌سازی بر می‌آید، پایین بودن ظرفیت این دو شبکه عصبی بزرگترین مشکل بر راه اجرایی شدن طرح است و برطرف ساختن آن ساده به نظر نمی‌رسد؛ زیرا شبکه هایپفیلد اصولاً ظرفیتی پایین دارد و تنها با اعمال برخی تغییرات ساختاری بر آن می‌توان ظرفیت شبکه را تا چند برابر افزایش داد، که آن نیز از نزدیکترین مزد نیاز دور است. شبکه پرسپترون چند لایه نیز - به رغم عملکرد ممتاز نسبت به دیگر شبکه‌ها - با توجه به ماهیت مسئله این پژوهش، امکان افزایش چشمگیر ظرفیت با تنظیمات کنونی را ندارد.

البته هر چند که مشکل ظرفیت مانع اصولی بر راه انجام و پیشبرد پژوهش در به کار گیری شبکه‌های عصبی مصنوعی برای اجرای طرح صحیح غلط املایی و تایپی نخواهد بود و با پیدایش شبکه‌های پرتوان‌تر یا استفاده بهینه‌تر از شبکه‌های موجود، این مسئله حل خواهد شد، اما به هر حال ارایه روش‌هایی مستقل از نوع شبکه‌ها برای افزایش ظرفیت سیستم سودمند خواهد بود که در ذیل به دو مورد اشاره می‌شود:

تقسیم‌بندی بر پایه طول واژه: با توجه به ماهیت شبکه که ورودی‌هایی گیسته متشکل از واژگان ۲ تا N حرفی (حدود ۱۰) با ۳۳ نوع حرف ("ا" تا "ي" و "ء") دارد، می‌توان آن را مانند شکل ۹ به $N-1$ شبکه مستقل برای پردازش واژگان ۲ تا N حرفی تقسیم کرد. روشن است که هر یک از گره‌های لایه ورودی یا خروجی این شبکه برای نمایش حروف، خود از شانزده نرون باینری یا دو فطی تشکیل می‌شود. این تدبیر ظرفیت کل سیستم را به طور متوسط $N-1$ برابر افزایش می‌دهد. دیگر مزیت این کار، رفع مشکل درج و حذف حرف در واژگان است که بسته به نیاز، با آزمون واژگان i حرفی در شبکه‌های $i+1$ و $i-1$ انجام خواهد شد.

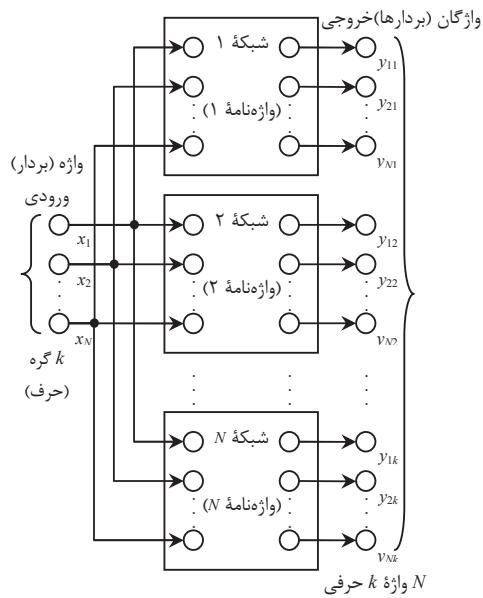


شکل (۷): افزایش ظرفیت با تکیب شبکه‌ها.

۵- نتیجه‌گیری

در این پژوهش کوشش شد تا به کمک شبکه عصبی مصنوعی روشی برای کشف و تصحیح غلط‌های املایی برآمده از خطاهای تایپی کاربر ارایه شود. بدین منظور بنا بر یک قرارداد مبتنی بر فاصله دکمه‌های صفحه کلید، کدی ۱۶ بیتی به هر یک از حروف اختصاص داده شد تا شبکه عصبی از روی آن همسایگی را تشخیص دهد. سپس نخست واژگان در گروه‌های ۴، ۵ و ۶ حرفی و دسته‌های ۴ تا ۲۵۶ واژه‌ای با دو نوع شبکه عصبی هاپفیلد و پرسپترون چند لایه آزموده شدند. نتایج این آزمون‌ها نشان داد که شبکه هاپفیلد از دقت عملکرد ۱۰۰٪ درستی برای تصحیح غلط‌های تایپی فرهنگ لغتی با ۴ واژه، به دقت درستی کمتر از ۴۲٪ برای تصحیح غلط‌های تایپی فرهنگ لغتی با ۲۵۶ واژه می‌رسد؛ در حالی که شبکه عصبی پرسپترون چند لایه چنین این‌رو با اطمینان می‌توان گفت که شبکه پرسپترون چند لایه، در این مساله یادآوری بهتر از شبکه هاپفیلد کار می‌کند. افزون بر این، گذشته از آن که دقت عملکرد دو شبکه در هنگام افزایش واژگان کاهش

نشانگر شانزده نرون هستند. در این شبکه، بردار واژه k حرفی به تمام شبکه‌ها وارد می‌شود و بنابراین در همه خروجی‌های N شبکه، پاسخ ظاهر می‌گردد.



شکل (۸): پخش واژگان در شبکه‌های موازی.

جدول ۸ نتایج آزمون روش تقسیم‌بندی برای ایجاد نوع واژه را بر واژگان ۵ حرفی با شبکه‌ایی مشتمل از ۴ تا ۳۲ زیرشبکه که هر یک از آنها ۳۲ واژه را یاد گرفته‌اند، نشان می‌دهد.

جدول (۸): افزایش ظرفیت با روش تقسیم‌بندی نوع واژه.

ردیف	تعداد زیرشبکه	تعداد واژگان	تعداد هاپفیلد	برسپترون چندلایه
۱	۴	۱۲۸	۹۰/۱	۹۵
۲	۸	۲۵۶	۸۸/۱	۹۴/۷
۳	۱۶	۵۱۲	۸۲/۱	۹۰/۹
۴	۳۲	۱۰۲۴	۷۷/۷	۸۶/۱

شبکه ترکیبی: با ترکیب این دو شبکه بالا، در حالتی که سیستم غلط‌یاب از ۳ شبکه ۴ تا ۶ حرفی، هر یک با ۳۲ زیرشبکه تقسیم شده بر پایه ۳۲ واژه (در مجموع ۳۰۷۲ واژه) شکل گرفته است، کارایی سیستم به دقت ۷۰ درصد در شبکه هاپفیلد و ۷۹/۶ درصد در شبکه پرسپترون چند لایه رسید. شکل ۱۱ نمای کلی این شبکه را نشان می‌دهد.

در این سیستم، هر واژه k حرفی تنها به ورودی شبکه متناظر شود و هر واژه ۳۲ خروجی خواهد داشت.



می یافته، با افزایش حروف واژگان نیز با کاهشی محسوس روبرو می گردید.

- [10] Teiji Furugori, "Improving spelling checkers for Japanese users of English". IEEE Transaction on Professional Communication, 33 (3) (1990) 138–142.
- [11] H. Shang and T. H. Merrettal, "Tries for approximate string matching". IEEE Transaction on Knowledge and Data Engineering, 8 (4) (1996) 540–547.
- [12] Bidyut Baran Chaudhuri, "Towards Indian language spell-checker design". IEEE Proceedings of the Language Engineering Conference (LEC'02), (2002) 139-146.
- [13] Z. Lei, Z. Ming, H. Changning and S. Maosong, "Automatic Chinese text error correction approach based-on fast approximate Chinese word-matching algorithm", Proceedings of the 3rd World Congress on Intelligent Control and Automation. (2000) 2739-2743.
- [14] Jin Hu Huang and David Powers, "Large scale experiments on correction of confused words". IEEE Proceedings 24th Australasian Computer Science Conference ACSC2001 (2001) 77-82.
- [15] V. Parisi, E. Garcia, J. Cabestany, J. Font, , and J. Salas, "A Hopfield neural network to track drifting buoys in the ocean", IEEE OCEANS '98 Conference Proceedings 2 (1998) 1010-1016.
- [16] E.M. Cochrane, J.E. Beasley, "The co-adaptive neural network approach to the Euclidean Travelling Salesman Problem", Neural Networks 16 (2003) 1499–1525.
- [17] M.A.S. Monfared, M. Etemadi, "The impact of energy function structure on solving generalized assignment problem using Hopfield neural network", European Journal of Operational Research, (2004).
- [18] Wenjing Li, Tong Lee, "Projective invariant object recognition by a Hopfield network", Neurocomputing 60 (2004) 1–18.
- [19] Robert Hecht-Nielsen, Neurocomputing, Addison-Wesley Publishing Company, 1989.
- همچنین در گام بعدی با تقسیم واژگان و ازهانمه در گروههای کوچکتر بر پایه طول و نوع واژگان و آموزش آنها در شبکههای جدآگانه و سپس ترکیب این شبکههای ظرفیت سیستم غلطیاب - با حفظ نسبی دققت - به اندازهای چشمگیر افزایش یافت؛ به گونهای که بیش از ۳۰۰۰ واژه با دقت درستی ۸۰ درصد تصحیح گردید.
- البته روشهای ارایه شده در این پژوهش، به ترتیب در صورتی خوب کار می کند که حروف واژه مورد آزمون اشتباہ، کم یا اضافه تایپ شده باشد و در مورد خطاهای ناشی از انتقال حروف در واژه، دیگر روشهای موجود کارآمدتر هستند.
- همچنین اهداف زیر دستور کارهای آینده قرار دارد:
- آزمایش تعداد واژگان بیشتر تا دست کم ۱۶ هزار کلمه که مقداری مناسب برای فرهنگ واژگان است.
 - آزمایش شبکههای عصبی دیگر برای دستیابی به نتایج بهتر، بهویژه شبکههای عصبی فازی ^{۵۳}
 - تنظیم بهتر تعداد لایه ها و نرون های هر لایه در شبکه پرسپترون چندلایه
 - یافتن نقطه بهینه برای حداکثر تعداد اپکهای شبکه
 - ادغام واژگان با تعداد حروف مختلف در یک شبکه
 - بهینه سازی کدهای حروف صفحه کلید و کاهش طول آنها.

مراجع

زیرنویس‌ها

¹ Natural Language Processing / Understanding

² Text Mining

³ Spell Checking

⁴ Syntax Checking

⁵ Concept Checking

⁶ Grammatical Confusion / Grammos

⁷ Phonological Similarity / Phonos

⁸ Typing Errors / Keyboard Mistyping / Typus

⁹ Insertion

¹⁰ Deletion

¹¹ Substitution

¹² Transposition (Interchange)

¹³ Spelling Errors

¹⁴ Transmission and Storage Errors

¹⁵ Optical Character Recognition (OCR)

¹⁶ Automatic Mode

¹⁷ Interactive Mode

¹⁸ Spelling Correction

¹⁹ String Substitution

²⁰ Reversal Errors

- [1] Allen James, Natural Language Understanding, The Benjamin/Cummings Publishing Co., 2nd Edition, 1994.
- [2] Patrick Ruch, Robert Baud and Antoine Geissbuhler, "Toward filling the gap between interactive and fully-automatic spelling correction using the linguistic context", IEEE International Conference on Systems, Man, and Cybernetics 1 (2001) 199-204.
- [3] Patrick Ruch, Robert Baud and Antoine Geissbuhler, "Using lexical disambiguation and named-entity recognition to improve spelling correction in the electronic patient record", Artificial Intelligence in Medicine 29 (2003) 169-184.
- [4] V. Cherkassky, N. Vassilas and G. L. Brodt, "Conventional and associative memory-based spelling checkers", Proceedings of the 2nd International IEEE Conference on Tools for Artificial Intelligence (1990) 138-144.
- [5] V.J. Hodge, J. Austin, "A comparison of a novel neural spell checker and standard spell checking algorithms", Pattern Recognition 35 (11) (2002) 2571–2580.
- [6] Dominique Lavenier, "A high performance systolic chip for spelling correction", Euro ASIC '92, Proceedings 1 (1992) 381-384.
- [7] R. Garfinkel., E. Fernandez, R. Gopal, "Design of an interactive spell checker: optimizing the list of offered words", Decision Support Systems 35 (2003) 385–397.
- [8] K. Kukich, "Techniques for automatically correcting words in text", ACM Comput Surveys 24 (4) (1992) 377–439.
- [9] J.R. Ullman, "A binary n-gram technique for automatic

-
- ²¹ String Matching
 - ²² Erroneous Word
 - ²³ Edit Distance
 - ²⁴ Hamming Distance
 - ²⁵ Levenshtein Distance
 - ²⁶ CorrectStar
 - ²⁷ Trie Algorithm
 - ²⁸ n-gram
 - ²⁹ Bangla
 - ³⁰ Fuzzy Approximate Word-Matching
 - ³¹ Associative Memory
 - ³² Hopfield
 - ³³ Backpropagation Networks
 - ³⁴ Dictionary Look-up Methods
 - ³⁵ Exact Matching

36 حروفچین‌ها به دو دسته تقسیم می‌شوند: یکی آنها که هر چه را که می‌بینند، تایپ می‌کنند و دیگر، آنان که خواسته یا ناخواسته متن را می‌فهمند و سپس تایپ می‌کنند. هر چند که یافتن فردی از گونه دوم برای تایپ موهبتی است و مزایای بسیار دارد، اما احتمالاً از آنجا که فرد نخست واژه را در مغز خود پردازش می‌کند و سپس آنرا تایپ می‌کند و گاه سرعت دست از سرعت پردازش برخی از واژگان در مغز بیشتر می‌شود، احتمال خطای املایی نیز در این افراد پیش می‌آید.

- ³⁷ Virtual Space
- ³⁸ Back Space
- ³⁹ Microsoft Word
- ⁴⁰ Matlab
- ⁴¹ Visual Basic
- ⁴² Multi-Layer Perceptron
- ⁴³ Neural Associative Memories
- ⁴⁴ Pattern
- ⁴⁵ Correlation Matrix
- ⁴⁶ Auto-Associative Memory
- ⁴⁷ John J. Hopfield
- ⁴⁸ Tank D. W.
- ⁴⁹ Traveling salesman problem
- ⁵⁰ Activation Function
- ⁵¹ Epochs
- ⁵² Back-Propagation
- ⁵³ Neuro-Fuzzy Networks

