

تصحیح خودکار غلط‌های تایپی فارسی به کمک شبکه عصبی مصنوعی ترکیبی

امیرشهاب شاهمیری^۱ رضا صفابخش^۲ رسول دژکام^۳

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 املائی 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 املائی بر پایه الگوریتم‌های جست‌وجو 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 املائی 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 دقتی بین ۵۵٪ و ۱۰۰٪ درستی و شبکه پرسپترون چندلایه - که در این تحقیق عمل یادگیری را در قالب دسته‌بندی انجام می‌دهد - به دقتی بین ۸۰٪ و ۱۰۰٪ درستی دست یافته، که این مقدار با به‌کارگیری شبکه‌های ترکیبی به نزدیک به ۸۰٪ دقت درستی برای بیش از ۳۰۰۰ واژه افزایش یافته است.

واژه‌های کلیدی: صحیح خودکار غلط تایپی فارسی، غلط املائی، 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 : تهران - خیابان حافظ - پلاک ۴۲۴ - دانشگاه صنعتی امیر کبیر - دانشکده‌ی مهندسی کامپیوتر و فناوری

اطلاعات

۱- مقدمه

گفت‌وگویی معنی‌دار 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 روشنی را پیش روی پژوهشگران این عرصه قرار داده است.

مسأله مشابه دیگری که محققان هوش مصنوعی و متن‌کاوی^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 املائی یک واژه، جایگزین مناسبی به کمک یک پایگاه داده جامع از واژگان برای آن پیدا شود. برای نمونه به جای واژه "ویزایش" 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 و او یکدیگر را دید." کشف گردد و جمله‌ای مانند "دیروز من و 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 نتایج رضایت‌بخش رسیده است.

ویرایش مفهومی^۵: که در آن کوشش می‌شود تا به‌رغم درستی ساختار جمله از نظر دستورزبان، ناهماهنگی‌های مفهومی آن دست‌کم شناخته و تا جای ممکن تصحیح گردد. برای نمونه جمله "میز هواپیما را خورد." از دیدگاه املائی 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 آن‌که واژه مورد نظر در واژه‌نامه موجود است، اما به دلیل ضعف دانش زبانی نگارنده یا اشتباه وی، با توجه به مفهوم جمله ناهم‌نوشته شده است (مانند واژه "عریب" در جمله "من در این شهر قریب هستم"). این‌گونه خطاها 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 "a piece of cake") [۲]. این دسته از اشتباهات از جمله خطاهای ابهام‌گرا می‌باشند و تشابه آوایی^۶ هستند [۳] و در محدوده این پژوهش نمی‌گنجد، 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.

به‌طور کلی خطاها و غلط‌های متون و مستندات را می‌توان به سه

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 جابه‌جایی دو حرف همسایه^{۱۲} را در بر دارد [۳ و ۵]. این چهار گونه، نزدیک به ۸۰٪ خطاهای تایپی را در بر می‌گیرند [۶].

غلط املائی^{۱۳}: که ناشی از ضعف دانش زبانی نویسنده در تشخیص اصوات واژه و شیوه نگارش آن است.
خطای انتقال^{۱۴}: که در هنگام انتقال اطلاعات بر شبکه یا دیسک و بازنشاسی نوری حروف^{۱۵} رخ می‌دهد.

برای پیاده‌سازی یک سیستم مصحح املائی، دو کار اصلی باید انجام پذیرد: نخست باید واژه مورد نظر در فرهنگ واژگان جست‌وجو شود تا در صورت موجود نبودن غلط قلمداد گردد؛ سپس با استفاده از روش‌های گوناگون در حالت خودکار^{۱۶} یک، یا در حالت محاوره‌ای^{۱۷} چند واژه به‌عنوان 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 صد درصد 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 این‌که بیشتر روش‌های تصحیح غلط بر روی تنها یک خطا در واژه خوب کار می‌کنند، تا زمانی که روش‌های تصحیح خطاهای ترکیبی مانند جابه‌جایی رشته^{۱۹} و وارونه نویسی^{۲۰} به‌میان نیایند، نمی‌توان به آنها اطمینان چندانی داشت [۴]. از این‌رو به‌نظر 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 به‌ندرت در گروه واژگان آزمایشی جای می‌گیرند [۹ و ۵].

اغلب روش‌هایی که تاکنون برای تصحیح لغوی ارائه 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} به‌کار می‌گیرند و یک یا چند واژه جایگزین را پیشنهاد می‌دهند [۳-۹]. این فاصله ویرایشی بنا بر اختلافات تک‌تک حروف واژه غلط با کلمات فرهنگ واژگان به‌دست می‌آید و از این‌رو اغلب - به‌ویژه در هنگام ویرایش واژگان کم‌حرف - نتیجه دلخواهی به‌همراه ندارد. برای نمونه، اگر به‌جای واژه "word" رشته "wprd" تایپ شود، این روش ۷۱۰ واژه جایگزین هم‌سطح را پیشنهاد می‌دهد [۵].

هدف این پژوهش ارائه 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

شناخت و ترکیب روش‌های مبتنی بر فاصله‌ی واژگان برای انواع خطاهای احتمالی در متن، تعاریفی سیستماتیک را همراه با الگوریتم‌های تصحیح ارائه دهند [۷].

۱-۲- ویژگی‌های مساله

در مساله‌ی غلط‌های املائی، احتمال پیش آمدن غلط تایپی ناشی از اشتباه کاربر، بیشتر از غلط‌های املائی متأثر از ضعف دانش زبانی نویسنده اصلی است؛ زیرا معقول‌تر آن است که فرض کنیم، کسی که خود متنی را به کمک رایانه تایپ می‌کند یا آن را به حروفچین می‌سپارد، به اندازه‌ی تسلط بر واژگان دارد که غلط‌های املائی متن وی ناچیز باشد، اما از سوی دیگر رخداد غلط تایپی امری رایج در هنگام ماشینی کردن متون است. این دو گونه اشتباه، ماهیتی متفاوت دارند و شناخت و استفاده از ویژگی‌های آنها می‌تواند به ما در تصحیح املائی متون کمک کند. این نکته مساله‌ای است که در روش‌های کلاسیک و نرم‌افزارهای تجاری تصحیح غلط املائی - که اغلب بر پایه‌ی الگوریتم‌های جست‌وجو استوارند - بدان توجه چندانی نشده و از این رو یکی از امتیازات این تحقیق به‌شمار می‌رود. البته باید در نظر داشت که نرم‌افزارهای تجاری از آن رو بر پایه‌ی روش‌های جست‌وجو در فرهنگ واژگان^{۳۴} کار می‌کنند که هدف آنها دستیابی سریع به پاسخ با کمترین حافظه‌ی مصرفی است؛ در حالی که روش‌های آکادمیک بیشتر به دستیابی به بهترین پاسخ می‌اندیشند [۴].

از آنجا که جست‌وجو در پایگاه داده‌ای با چند صد هزار واژه کاری وقت‌گیر است، روش‌های تصحیح غلط املائی چهار راهبرد اساسی را برای کاهش تعداد جست‌وجوها به کار می‌گیرند که آنها را روش‌های تطبیق کامل^{۳۵} می‌نامند [۴]:

بر پایه‌ی تعداد دفعات تکرار: بدین ترتیب که واژگانی که بیشترین مورد استفاده را در متون دارند، در یک فرهنگ واژگان کوچکتر گردآوری شده، مصحح در هنگام یافتن غلط، نخست این پایگاه را می‌جوید و در صورتی که واژه یافته نشد، پایگاه‌های واژگان کم‌استفاده‌تر در اولویت‌های بعدی را نیز ارزیابی می‌کند.

بر پایه‌ی طول واژه: فرهنگ واژگان 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 جست‌وجو می‌گردد.

بر پایه‌ی 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 فرزندی تا همین تعداد دارد. بنابراین برای یافتن واژه k حرفی به k جست‌وجو نیاز است. اما روشن است که این روش زمانی خوب کار می‌کند که حروف آغازین واژه درست نگاشته شده باشند.

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

نکاتی در آنها درج شده باشد که از دیدگاه دستوری یا مفهومی اشتباه به‌نظر 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 املائی 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 املائی انگلیسی غیرانگلیسی‌زبانان و با بررسی تفاوت‌های آوایی زبان ژاپنی و انگلیسی، سیستمی کمکی برای تصحیح غلط‌های املائی ژاپنی‌ها در هنگام نگارش به‌زبان انگلیسی پیشنهاد کرد که می‌توانست دقت نرم‌افزار تصحیح املائی کرکستار^{۳۶} را از ۶۰٪ به ۷۵٪ برساند [۱۰]. در سال ۱۹۹۶ "شانگ" و "مرتال" با بررسی چند الگوریتم دیگر، روش 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 به‌همراه یک تکنیک برنامه‌نویسی پویا برای تصحیح مقایسه‌ای رشته حروف بسازد که می‌توانست ۲۰۰ هزار واژه را در هر ثانیه تصحیح کند [۶]. "چاودوری" در سال ۲۰۰۲ با توجه به ویژگی‌های آوایی و غیرآوایی زبان هندی و با به‌کارگیری برخی از تکنیک‌ها - مانند روش m -گرام^{۳۸} - توانست سیستم مصحح املائی برای زبان هندی و بنگلایی^{۳۹} آرایه کند که به دقت درستی ۹۵٪ دست یافت [۱۲]. در سال ۲۰۰۰ "هاج" و "اوستین" توانستند با ترکیب روش m -گرام و رویکرد فاصله‌ی همینگ سیستمی را فراهم آورند که برای یک پایگاه با تعداد محدودی از واژگان و با سرعتی قابل قبول، هر چهار گونه از غلط املائی را با دقتی تا ۹۷٪ تصحیح کند که به‌طور متوسط ۸ واژه جایگزین را پیشنهاد می‌کرد [۵]. در سال ۲۰۰۰ "لی" و همکارانش توانستند یک روش تطابق تخمینی واژه فازی^{۴۰} را در قالب مصحح املائی اختصاصی برای زبان چینی آرایه کنند که علاوه بر غلط‌های رایج، جابه‌جایی 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 واژگان متون پزشکی به دقت درستی ۹۸٪ دست یابند [۳ و ۲]. "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^{۴۱} مانند شبکه‌ی هاپفیلد^{۴۲} و شبکه‌های پس‌انتشار^{۴۳} با دیگر روش‌ها، سیستم غلط‌یاب املائی برای واژگان کوتاه (۵ تا ۷ حرفی) و بلند (۱۰ تا ۱۲ حرفی) با تعداد گره‌های ورودی برابر با توانی از ۲۶ (به‌تعداد حروف الفبای انگلیسی) به مقدار n در الگوریتم m -گرام و گره‌های خروجی به‌تعداد واژه‌های ذخیره شده، بسازند و به دقت ۱۵ تا ۱۰۰ درصد برای انواع خطا و مقادیر n دست یابند [۴]. "گارفینکل" و همکارانش در سال ۲۰۰۲ کوشیدند تا با

ممکن بدن 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 جست‌وجو و مبتنی بر شبکه عصبی را چنین بیان کرد:

- با وجود این که شبکه عصبی در فاز یادگیری به زمان بالایی نیاز دارد، اما در فاز آزمایش به سرعت نتیجه می‌گیرد و بنابراین از هر روش دیگری سریع‌تر است.
- شبکه‌های عصبی مصنوعی ظرفیتی محدود دارند؛ در نتیجه برای یادگیری تعداد زیادی واژه که بتواند حداقل نیاز برای مصحح املائی 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 این‌رو مساله باید با شبکه‌های گوناگون آزموده شود، تا بهترین ترکیب به دست آید.
- با بهتر شدن و پیدایش انواع تازه و پرتوان‌تر شبکه‌های عصبی در آینده، نتایج کار نیز بهبود خواهند 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 شباهت 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 خطا 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 املائی 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

همان‌گونه که در مقدمه گفته شد، در هنگام مواجهه با غلط‌های املائی، 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 املائی متن وی ناچیز باشد، اما از سوی دیگر رخداد غلط تایپی در هنگام تایپ متن بسیار رایج است.

۲-۱- غلط املائی در زبان فارسی

غلط املائی 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. این‌گونه اشتباه معمولاً در مورد واژه‌های عربی دخیل در فارسی رخ می‌دهد. برای نمونه نویسنده ممکن است در واژه "اضطراب" حرف "ض" را با "ز"، "ذ" یا "ظ" و حرف "ط" را با "ت" اشتباه بگیرد و واژه "تقریظ" را "تقریظ"، "تقریض" یا به شیوه‌های 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) : معاذحه رخ می‌دهد.

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

۲-۲- غلط تایپی در زبان فارسی

اما غلط تایپی به سواد نویسنده بستگی ندارد، بلکه وابسته 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 پایپی در استانداردها و تعداد دکمه‌های صفحه کلید رخ داده و بخشی دیگر، به علت نبود اتفاق نظر میان استانداردارگذاران درون و حتی بیرون از کشور بر سر تعیین جایگاه حروف فارسی بر دکمه‌ها پدید آمده و این ناهماهنگی‌ها 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 گفت‌وگو 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 شباهت آوایی با هم دارند و نه همسایه نزدیک یکدیگر هستند، با هم اشتباه می‌شوند. مانند "پ" و "ز".

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

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

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

شکل (1): سه نمونه از جانمایی حروف فارسی بر صفحه کلید.

جابه‌جایی 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 : کیومرث / گیومرث".

انتقال کلید: یکی از حروف اصلی، با حرف همسایه یا چند حرف قبل یا بعدی جابه‌جا تایپ می‌شود. مانند: "کیومرث : کویمرث / کومرث". غلط املائی (جابه‌جایی ویژه): به جای یکی از حروف اصلی، یکی از دیگر کلیدهای صفحه کلید فشرده می‌شود؛ به گونه‌ای که در زبان فارسی آن حرف با حرف اصلی هم‌صدا است مانند: "کیومرث : کیومرس". این گونه غلط ناشی از ضعف دانش زبانی حروف چین است و بسار کم رخ می‌دهد، زیرا متن اصلی پیش روی حروف چین است.

فاصله اضافه: یک کاراکتر فاصله نابه‌جا درج می‌گردد. مانند: "کیومرث : کی ومرث". این غلط به‌ویژه در برخی حالت‌های ویژه مانند "ها"ی جمع و "می" استمرار و اغلب به‌جای فاصله مجازی^{۳۷} (کوتاه) رخ می‌دهد. مانند: "واژه‌ها : واژه‌ها" 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 : کیومرثسیامک". این غلط گونه‌ای ویژه از غلط‌های حذف است.

بر این پایه با نمونه‌برداری و تخمین آماری غلط‌های تایپی زبان فارسی در هنگام کار حروف چین 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^{۳۸} 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 جابه‌جا می‌شود. مانند "ا" به جای "آ" (که البته این مورد غلط املائی به‌شمار نمی‌رود)، یا "ی" و "ط"، "ز" به جای "ژ" (در برخی از صفحه کلیدها) و "ی"، "س" و "ش" به جای هر یک از اعراب.

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

گاه غلط املائی 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". از این‌رو انواع غلط‌های تایپی زبان فارسی به شرح زیر دسته‌بندی شد:

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

تکرار حرف (درج حرف همسایه): یعنی در هنگام نگارش، یکی از حروف اصلی به اشتباه دو بار تایپ می‌شود. مانند: "کیومرث : کیومرث". 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".

رخداد خطا، 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	جابجایی کلید	جابجایی	۳۹/۵
۵	انتقال کلید	انتقال	۵
۶	غلط املائی	جابجایی	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

۲-۳- تصحیح املائی بر پایه فاصله

امروزه بیشتر نرم‌افزارهای مصحح غلط املائی که به صورت تجاری در بازار ارایه می‌شوند، بر پایه فاصله (ویرایشی، همینگ، لونشتین و ...) کار می‌کنند، که البته در بسیاری از موارد - به ویژه در هنگام رخداد غلط‌های جابجایی - نتیجه دلچسبی به همراه ندارند. در این میان، مصحح املائی نرم‌افزار ورد از شرکت مایکروسافت^۳ یکی از کارآمدترین و پرطرفدارترین تصحیح‌کننده‌های املائی 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 به جای حرف "g"، حرف کناری آن "h" و بنابراین "holf" را تایپ کند، نرم‌افزار ورد واژگان زیر را به ترتیب اولویت پیشنهاد می‌دهد:

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

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

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

"hole": در این مورد هم "e" به جای "f" آمده، در صورتی که این دو حرف در دو سطر جدا قرار دارند و احتمال اشتباه تایپی این دو ناچیز است.

"holy" در اینجا "y" به جای "f" آمده در صورتی که این دو حرف در دو سطر 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 به جای حرف "s"، حرف کناری آن یعنی "a" و بنابراین "aear" را تایپ کند، نرم‌افزار ورد واژگان زیر را به ترتیب اولویت پیشنهاد می‌دهد:

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

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

"afar": در اینجا "f" به جای "e" آمده در صورتی که این دو حرف در دو سطر جدا قرار دارند و احتمال اشتباه تایپی این دو، به دلیل دوری از هم، ناچیز است.

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

"agar": باز هم در اینجا حرف "g" به جای "e" آمده در صورتی که این دو حرف در دو سطر جدا قرار دارند 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 بگیرید: اگر کاربر بخواهد واژه "that" را تایپ کند، اما به اشتباه به جای حرف "t" دوم، حرف کناری آن یعنی "y" و بنابراین "thay" را تایپ کند، نرم‌افزار ورد واژگان زیر را به ترتیب اولویت پیشنهاد می‌دهد:

"they": در اینجا "e" به جای "a" آمده در صورتی که این دو حرف در دو سطر جدا قرار دارند و احتمال اشتباه تایپی این دو ناچیز است.

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

"tray": در اینجا فرض شده که به جای حرف "r"، حرف "h" تایپ شده در صورتی که این دو حرف در دو سطر جدا قرار دارند و احتمال اشتباه تایپی این دو اندک است.

"that": در اینجا "t" به جای "y" آمده و این دو حرف بر کلیدهای مجاور قرار دارند و این انتخاب معقول‌ترین انتخاب به نظر می‌رسد.

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

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

۳-۳- فرهنگ واژگان به‌کار رفته در شبیه‌سازی

در شبیه‌سازی این پژوهش، ۶۰۰ واژه ۴، ۵ و ۶ حرفی به تعداد مساوی، برای آموزش شبکه‌ها و ۳۰۰ واژه نیز برای آزمایش به‌کار 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

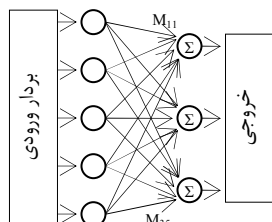
شبیه‌سازی و آزمون مقادیر و واژگان توسط نرم‌افزار متلب^{۴۰}، به‌همراه نرم‌افزار کمکی تحت زبان ویژوال بیسیک^{۴۱} برای تهیه مقادیر ورودی دو نوع شبکه عصبی هاپفیلد و پرسپترون چند لایه^{۴۲} انجام شد. الگوهای ورودی در دسته‌های ۴، ۵ و ۶ حرفی به شبکه عصبی وارد شد و پس از آموزش شبکه، هر بار الگوهای ورودی به‌دفعات و با تغییر تصادفی یکی از حروف واژه آزموده شد.

۴-۱- شبیه‌سازی با شبکه عصبی هاپفیلد

هنگامی که درباره تصحیح املائی واژگان سخن می‌گوییم، شبکه‌های گونه حافظه یادآور^{۴۳} نخستین شبکه‌های مناسب به‌نظر می‌رسند. شبکه حافظه یادآور سیستمی است که می‌تواند داده‌های ذخیره شده (الگوها^{۴۴}) را حتی با دیدن ورودهای همراه با غلط بازیابی (یادآوری) کند [۴]. به‌عنوان نمونه‌ای کوچک برای آشنایی با روش کار این شبکه خط فارسی را به ۵ حرف ("ا"، "ب"، "ر"، "ش" و "ی") محدود می‌کنیم و می‌خواهیم سه واژه "ابر"، "آرش" و "بیش" را یاد 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^{۴۵} می‌نامند.

$$\begin{aligned} \text{ی ر د ب ا} \\ \text{"ابر"} &= [1 \ 1 \ 0 \ 1 \ 0] \\ \text{"برد"} &= [0 \ 1 \ 1 \ 1 \ 0] \\ \text{"یاد"} &= [1 \ 0 \ 1 \ 0 \ 1] \end{aligned}$$

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



شکل ۳-ب): شبکه عصبی حافظه یادآور

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

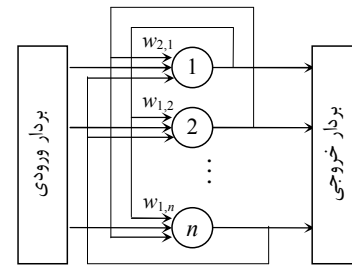
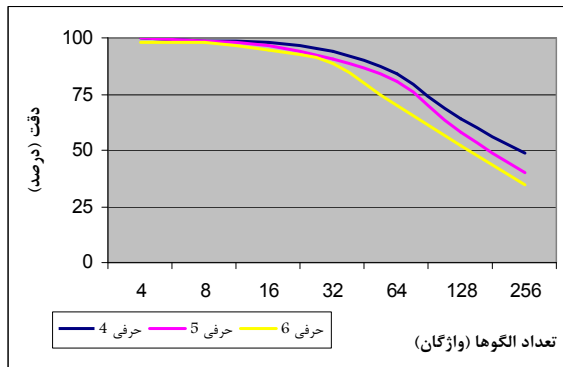
شکل ۳-پ): بازیابی واژگان در شبکه حافظه یادآور به‌کمک ماتریس همبستگی

شبکه عصبی مصنوعی هاپفیلد شاخص‌ترین نوع شبکه‌های یادآور و از گونه حافظه خودیادآور^{۴۶} است که کاربردهای بسیاری در دانش‌های گوناگون دارد. این شبکه نخستین بار در سال ۱۹۸۲ میلادی توسط جان هاپفیلد^{۴۷} ارائه 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 مسایل بسیار مشکل قرار می‌گیرد و با روش‌های معمولی نمی‌تواند در زمانی معقول پاسخی بهینه را به‌دست آورد. هاپفیلد 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 هاپفیلد در قالب یک سیستم پویا توسط یک تابع انرژی که باید تعادلی میان اهداف تابع مساله - که باید حداقل شود - ایجاد می‌کند [۱۷]. پس از این موفقیت شبکه عصبی هاپفیلد، بسیاری از مسایل مهندسی در قالب تابع انرژی که باید کمینه شود، ارائه گردید. چنین راه حلی بسیار جذاب است، زیرا پردازش موازی را برای حل مسایل امکان‌پذیر می‌سازد [۱۸].

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

نشان داده شده است.

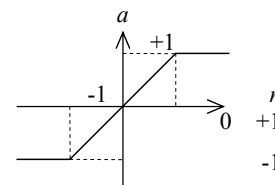


شکل (۴): نمای کلی شبکه هاپفیلد.

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

همان گونه که در این جدول مشاهده می شود، با افزایش تعداد الگوها، 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^۱) در هنگام آموزش نیاز است در صورتی که دقت تصحیح شبکه از ۵۵٪ فراتر نمی رود.

اما دو نقص بزرگ شبکه هاپفیلد ظرفیت پایین آن (در حدود ۱۵٪ اندازه شبکه یا تعداد گره ها) و همچنین هزینه محاسباتی بالای آن است [۱۹]. از این رو انتظار می رود که با افزایش تعداد الگوهای یاد داده شده، دقت شبکه کاهش یابد و بنابراین شبکه هاپفیلد نامزدی مناسب برای سیستم غلطیاب املائی نیست [۴].



شکل (۵): تابع فعالیت (دو قطبی) شبکه هاپفیلد.

نتایج این آزمون و شبیه سازی آن در جدول ۳ و شکل ۶ نمایش داده شده است.

۴-۲- شبیه سازی با شبکه عصبی پرسپترون چند لایه

شبکه عصبی پرسپترون چند لایه، شبکه ای است که در اصل برای مسایل دسته بندی و تخمین تابع طراحی شده و در این کارها از موفق ترین شبکه های عصبی بوده و توانسته است با استفاده از قانون انتشار خطا به عقب^۲ بسیاری از مسایل غیر قابل حل توسط شبکه های دیگر را حل کند [۱۹]. اما با توجه به این که کار اصلی این شبکه دسته بندی و تخمین تابع است، بنابراین به نظر نمی رسد که در مساله تصحیح غلط تایپی - که می توان گفت مساله یادآوری است - به کار آید.

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

ردیف	تعداد الگوها (واژگان)	تعداد حروف واژه	حدافل 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	۳۲	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

نکته ساده ای 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 با یک دسته در فضای n بعدی متناظر کنیم، برای یادگیری ۲^n واژه، مساله به مساله دسته بندی ۲^n دسته ای با n نرون در لایه خروجی شبکه متناظر می شود. پس برای نمونه، برای یادگیری مجموعه ای از واژگان با حدود ۶۵ هزار واژه، به شبکه ای با ۱۶ گره خروجی نیاز است. جدول ۴ نمونه ای از این روش را برای دسته بندی ۸ واژه چهار حرفی نشان می دهد. واژه های مورد آزمایش به عمد به گونه ای برگزیده شده اند که به هم نزدیک باشند و با دگرگونی یکی-دو حرف به واژه های دیگر در واژه نامه تبدیل شوند، تا نتایج تصحیح بر روی آنها بهتر نمایان شود.

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

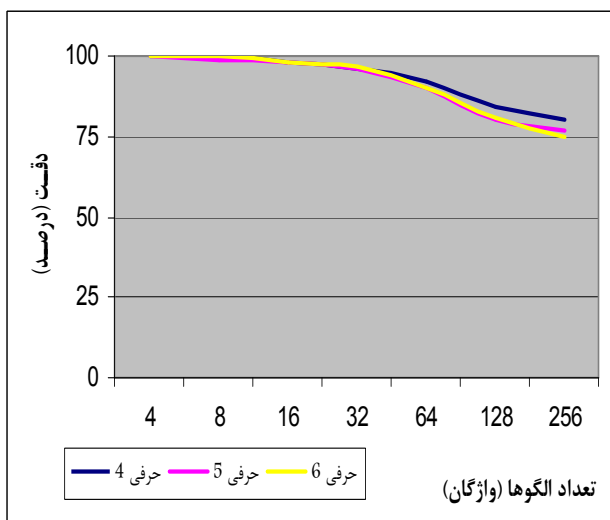
ردیف	واژه	کد باینری (دسته)
۱	نامی	۰۰۰
۲	مانی	۰۰۱
۳	مینا	۰۱۰
۴	امین	۰۱۱
۵	میان	۱۰۰
۶	ایمن	۱۰۱
۷	نیام	۱۱۰
۸	نیما	۱۱۱

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 ایک‌ها 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 شبکه با تعداد دسته، لایه و نرون بسیار سنگین می‌شود (آموزش مورد ردیف ۱۹ بیش از ۳ روز زمان برده است!) کار با ایک‌های کمتر خاتمه یافته است.

همان گونه که مشاهده می‌شود، با افزایش تعداد الگوها، 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 تاپیی واژگان.

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



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

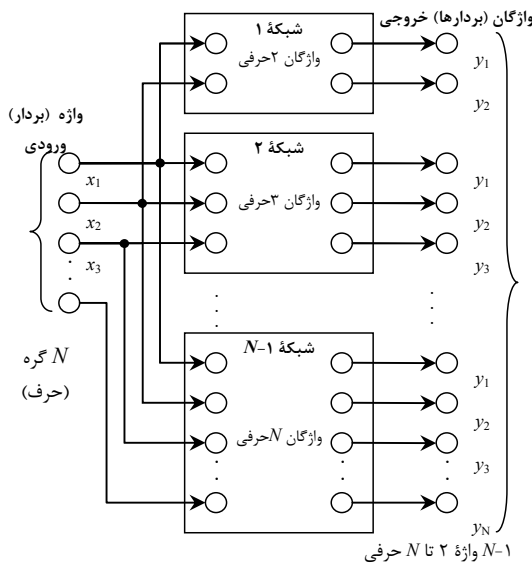
جدول (۶): مقایسه میانگین بازدهی دو شبکه هاپفیلد و پرسپترون چند لایه

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	۴۱/۳	۷۷/۳

۴-۳- 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.

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.



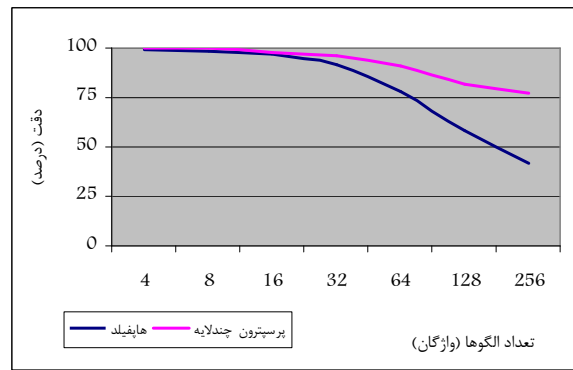
شکل (۹): به کارگیری شبکه‌های موازی در پردازش واژگان N تا ۲ حرفی

جدول ۷ نتایج آزمون روش تقسیم‌بندی بر پایه طول واژه را بر واژگان ۴ تا ۶ حرفی با ۳ زیرشبکه 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.

ردیف	تعداد الگوهای هر شبکه	تعداد کل واژگان	تقسیم‌بندی طول واژه	
			هافیلد	پرسپترون چندلایه
۱	۴	۱۲	۹۹/۳	۱۰۰
۲	۸	۲۴	۹۸/۲	۹۹/۵
۳	۱۶	۴۸	۹۵/۵	۹۷/۲
۴	۳۲	۹۶	۸۹/۱	۹۴/۱
۵	۶۴	۱۹۲	۷۳	۸۸/۷
۶	۱۲۸	۳۸۴	۵۰/۵	۷۹/۸
۷	۲۵۶	۷۶۸	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 به‌طور متوسط N برابر افزایش می‌دهد و گذشته از آن، امکان پیشنهاد چندین واژه جایگزین در خروجی را هم پدید می‌آورد. روشن است که گره‌های لایه ورودی و خروجی این شبکه نیز، که در شکل ۱۰ نمایش داده شده است، هر یک



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

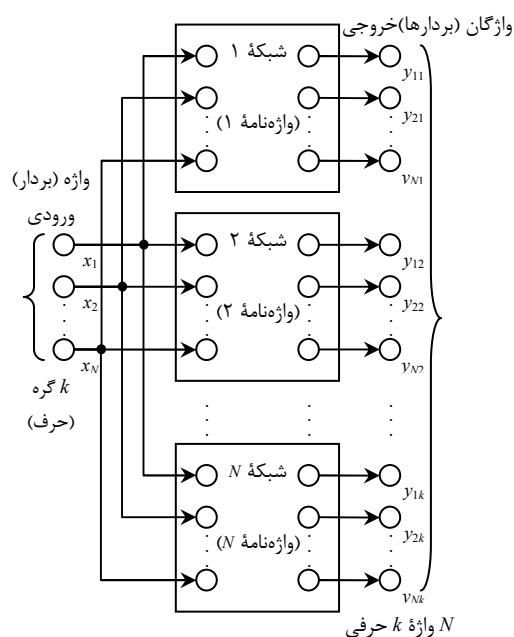
۴-۴- ترمیم مشکل ظرفیت شبکه

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 به‌نظر نمی‌رسد؛ زیرا شبکه هافیلد 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 مرز نیاز دور است. شبکه پرسپترون چند لایه نیز - به‌رغم عملکرد ممتازش نسبت به دیگر شبکه‌ها - با توجه به ماهیت مسأله این پژوهش، امکان افزایش چشمگیر ظرفیت با تنظیمات کنونی را ندارد.

البته هر چند که مشکل ظرفیت مانعی اصولی بر راه انجام و پیشبرد پژوهش در به‌کارگیری شبکه‌های عصبی مصنوعی برای اجرای طرح تصحیح غلط املائی و تایپی نخواهد بود و با پیدایش شبکه‌های پرتوان‌تر یا استفاده بهینه‌تر از شبکه‌های موجود، این مسأله حل خواهد شد، اما به‌ر حال رایج 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 گسسته متشکل از واژگان ۲ تا N حرفی (حدود ۱۰) با ۳۳ نوع حرف ("ا" تا "ی" 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 قطبی تشکیل می‌شود. این تدبیر ظرفیت کل سیستم را به‌طور 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 $i-1$ حرفی انجام خواهد شد.

نشانیگر 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 (۱۰): بخش واژگان در شبکه‌های موازی.

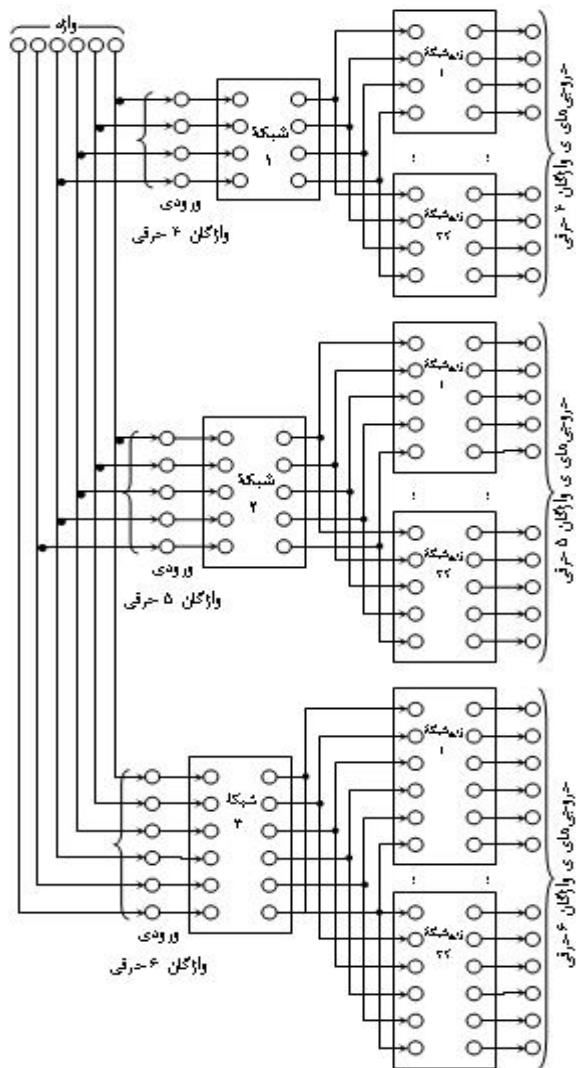
جدول ۸ نتایج آزمون روش تقسیم‌بندی بر پایه نوع واژه را بر واژگان ۵ حرفی با شبکه‌هایی متشکل از ۴ تا ۳۲ زیرشبکه که هر یک از آنها ۳۲ واژه را یاد گرفته‌اند، نشان می‌دهد.

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

ردیف	تعداد زیرشبکه	تعداد کل واژگان	هافیلد	پرسپترون 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 هافیلد 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 متناظرش وارد می‌شود و هر واژه ۳۲ خروجی خواهد داشت.



شکل (۷): افزایش ظرفیت با تکب شبکه‌ها.

۵- نتیجه‌گیری

در این پژوهش کوشش شد تا به کمک شبکه عصبی مصنوعی روشی برای کشف و تصحیح غلط‌های املائی برآمده از خطاهای تایپی کاربر ارائه 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 هافیلد و پرسپترون چند لایه آزموده شدند. نتایج این آزمون‌ها نشان داد که شبکه هافیلد 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 چندلایه چنین فرهنگ واژگانی را با بیش از ۷۷٪ دقت درستی غلطیابی می‌کند. از این‌رو با اطمینان می‌توان گفت که شبکه پرسپترون چند لایه، در این مسأله یادآوری بهتر از شبکه هافیلد کار می‌کند. افزون بر این، گذشته از آن که دقت عملکرد دو شبکه در هنگام افزایش واژگان کاهش

- correction of substitution, deletion, insertion and reversal errors in words”, The Computer Journal 20 (2) (1977) 141–147.
- [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.

می‌یافت، با افزایش حروف واژگان نیز با کاهش محسوس روبرو می‌گردید.

همچنین در گام بعدی با تقسیم واژگان واژه‌نامه در گروه‌های کوچکتر بر پایه طول و نوع واژگان و آموزش آنها در شبکه‌های جداگانه و سپس ترکیب این شبکه‌ها، ظرفیت سیستم غلط‌یاب - با حفظ نسبی دقت - به اندازه‌ای چشمگیر افزایش یافت؛ به‌گونه‌ای که بیش از ۳۰۰۰ واژه با دقت درستی ۸۰ درصد تصحیح گردید.

البته روش‌های ارایه شده در این پژوهش، به‌ترتیب در صورتی خوب کار می‌کند که حروف واژه مورد آزمون اشتباه، 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 دست‌کم 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 لایه‌ها و نرون‌های هر لایه در شبکه پرسپترون چندلایه
- یافتن نقطه بهینه برای حداکثر تعداد اپک‌های شبکه
- ادغام واژگان با تعداد حروف مختلف در یک شبکه
- بهینه‌سازی کدهای حروف صفحه‌کلید و کاهش طول آنها.

مراجع

زیر نویس‌ها

- 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 گونه‌دوم 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 املائی 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