

مهارت در جستجوی اطلاعات فارسی از اینترنت [۱]

محمد صابر راثی ساربانقلی [۲]

چکیده:

خط فارسی دارای مشکلات مختلفی می‌باشد که در جستجو و بازیابی اطلاعات مسائل و مشکلات فراوانی را فراوانی را قرار می‌دهد. به خصوص با رشد سریع انتشارات الکترونیکی بر روی وب در شکل‌های مختلف پایگاه‌های اطلاعاتی، وبلاگ و ... و اینکه هیچ قاعده مشخص و ثابتی برای رسم‌الخط فارسی وجود ندارد باعث شده است که جستجوگران مطالب فارسی با مشکلات فراوانی روبرو بشوند. این مقاله سعی دارد تا با اشاره به موارد مختلفی که می‌تواند در جستجو و بازیابی اطلاعات سرعت و دقت و جامعیت و مانعیت جستجو را بالا ببرد موجب افزایش مهارت کاربران اینترنت فارسی بشود.

کلید واژه ها: اینترنت، خط فارسی، جستجو و بازیابی اطلاعات.

مقدمه

اینترنت به عنوان یک محمل اطلاعاتی عظیم، منابع اطلاعاتی را در مقیاسی وسیع در دسترس مخاطبان بالقوه قرار داده است. اغلب سهولت دسترسی به منابع اطلاعاتی اعم از متن و سایر رسانه‌ها عمده‌ترین مزیت اینترنت محسوب می‌شود. اما این توانایی که هرکس ناشر آثار خود باشد عواقب ناخواسته‌ای را نیز در پی خواهد داشت و آشکارترین معضل، آن است که انبوهی از منابع بسیار متنوع و غیر قابل مدیریت را فراهم می‌آورد. افزایش سریع منابع اینترنتی نیازمند یک سازمان‌دهی مفید و موثر است. هرچند در حال حاضر راهنماهایی برای منابع اینترنتی تهیه شده است که براساس فایل‌های مقلوب ساخته شده توسط موتورهای جستجو و با استفاده از قابلیت‌های مختلف این موتورها از جمله: استفاده از عملگرهای بولی، جستجوی دقیق عبارت، محدود کردن یک جستجو به بخش خاصی از رکورد (مانند عنوان، آدرس)، کوتاه‌سازی کلمات، جستجوی نزدیک‌یابی واژه‌ها، ایجاد محدودیت زمانی و منطقه‌ای و زبانی، و ... به جستجوی اطلاعات کمک می‌کند، اما باید تاکید کرد که در امر بازیابی اطلاعات از اینترنت بدون نمایه‌سازی نظام یافته نمی‌توان انتظار بازیابی مفید و موثر را داشت. هرچند پیش‌تر اطلاعات موجود بر روی اینترنت به زبان انگلیسی است، ولی حجم اطلاعات به زبان فارسی نیز با سرعت در حال افزایش است و کاربران به دلایل مختلفی علاقه زیادی به اطلاعات فارسی نشان می‌دهند و از آنجائی‌که زبان غالب در اینترنت انگلیسی است جستجو به زبان‌های غیر انگلیسی از جمله فارسی، مسایل و مشکلات مختلفی را جدای از مشکلات عمومی اینترنت دارد.

خط فارسی

اشکال و نقصی که در همه خطوط جهان است دو علت دارد که یکی در اصل خط است و دیگری بر اثر تغییر و تحول زبان ایجاد می‌شود. دقت فراوان در ثبت همه دقایق تلفظ اغلب موجب دشواری شیوه خط است و این دقت زمانی ضرورت می‌یابد که زبانی توسعه بسیار بیابد و در کشورهای دیگری که به آن زبان سخن نمی‌گویند رایج شود. به عنوان مثال در خط عربی نقطه و علامت‌های حرکات وقتی به وجود آمد که زبان عربی نزد ملت‌های غیر عرب معمول شد، در خط یونانی نیز نشانه‌های آهنگ و تکیه [۳] پس از رواج آن زبان در مصر ایجاد شد تا کسانی که زبان مادری‌شان یونانی نبود و با تلفظ آن مانوس نبودند بتوانند کلمات و عبارات یونانی را هر چه درست‌تر ادا کنند. با این حال هیچ خطی هر قدر دقیق و شماره علامت آن فراوان باشد، ممکن نیست که کاملاً نشانه شیوه تلفظ باشد. و با کمک علامت متعدد علم حروف نیز تا کسی چگونگی تلفظ زبانی را نشنود نمی‌تواند عبارت و کلمات آنرا مانند اهل آن زبان ادا کند.

اما نقصی که بر اثر تحول زبان و به تدریج در خط حاصل می‌شود، مشکلی است که همه ملت‌ها با آن رو به رو هستند. بعضی از حروف و اصوات زبان در طی زمان تغییر می‌پذیرند و این تغییر در گفتار حاصل می‌شود، اما خط همیشه صورت کهن تلفظ را حفظ می‌کند، و از اینجا میان "گفتار" و "نوشتار" اختلاف روی می‌دهد. دیگر آن که هر زبانی ناگزیر لغاتی از زبان‌های دیگر به عاریت می‌گیرد و اگر علائم خط در این دو زبان یکی باشد کلمه خارجی به همان املاي اصلي در نوشتن به کار می‌رود که اغلب با املاي کلمه مشابه در زبان ثانوي تفاوت دارد و از اینجا برای اصوات واحد علائم خطي متعدد پدید می‌آید. در خط فارسي نمونه همه این موارد را می‌توان یافت. چون خط عربي برای نوشتن فارسي به کار رفت کلماتی که از آن زبان اخذ شده بود به همان صورت اصلي نوشته شد. حال آنکه به یقین در هیچ دوره‌اي حروف خاص عربي را فارسي زبان‌ها درست مثل اصل تلفظ نکرده‌اند. در زبان‌های دیگر نیز این گونه موارد نمونه‌های متعدد دارد. شاید دو زبان انگلیسی و فرانسه بیش از همه زبان‌های جهان دچار اختلاف تلفظ و خط باشند. به طور کلی نقائص و معایبی که در خطوط معمول جهان است را می‌توان به طریق زیر طبقه‌بندی کرد:

۱. شکل واحدی اصوات مختلف را بیان می‌کند. چنانکه در فارسي حرف "ي" را گاهی برای حرف لین بکار می‌بریم (یک) و گاهی برای حرف مد (بی) و گاهی به جای الف (عیسی) و گاهی برای نشان دادن مصوت مرکب (ری). و یا حرف «و» در کلمات (سوار، سود، تو)

۲. اصوات واحد به صورت‌های مختلف نوشته می‌شود. در فارسي حرف "س" سه صورت (س - ص - ث) و حرف "ز" چهار صورت (ز - ذ - ض - ظ) دارد؛ در زبان فرانسه حروفی که "سن" خوانده می‌شود پنج رسم الخط دارد که اگر صورت‌های جمع را نیز به حساب بیاوریم ده شکل می‌شود از این قرار (saint, ceint, sein, seing, sain)

۳. بسیاری از حروف نوشته می‌شود ولی خوانده نمی‌شود. یعنی علاماتی بی‌فایده در نوشتن به کار می‌رود در فارسي نوشتن "واو معدوله" و "هـاء غیر ملفوظ" از این قبیل است. در انگلیسی نمونه این مورد بسیار است مانند high که دو حرف آخر آن به کلی از تلفظ ساقط است. و یا "K" در کلمه "Know".

۴. اصواتی هستند که تلفظ می‌شود اما در خط نشانه‌ای برای آن‌ها نیست. در فارسي سه مصوت کوتاه (اَ ، اِ ، اُ) از این قبیل است هم چنین الف در کلمات اسحق و الله که در کتابت نمی‌آید. [۴]

زبان و خط فارسي نیز مشکلات خاصی را دارا می‌باشد و نظام نوشتاری فارسي برای ثبت دقیق گفتار، نارسائی دارد و قواعد نگارش آن مدون نیست، از این رو فاصله میان گفتار و نوشتار در فارسي قابل توجه است. بیشترین مشکلات نیز به جهت نبود يك رسم الخط واحد که عموم اساتید و اهل فن روی آن اجماع کرده باشند به وجود آمده است. به طوری که در حال حاضر جدای از چندین شیوه‌نامه رسمي همچون "شیوه‌نامه سمت، نشر دانشگاهی، فرهنگستان، آموزش و پرورش" به تعداد افراد جامعه، رسم الخط و شیوه نگارش زبان وجود دارد، هر ناشری برای خود به قاعده‌ای دلخواه عمل می‌کند که این تعددها موجب پریشانی و پراکندگی شده و با يك دیگر تفاوت‌هایی دارند. از دیگر دلایل می‌توان به عاریتی بودن خط فارسي و چاره‌اندیشی برای حرکات و عدم تطابق واج‌ها با حروف اشاره کرد. متصل و منفصل‌نویسی نیز یکی دیگر از حوزه‌های مورد اختلاف است از دیگر مشکلات: گوناگونی معادل‌های علمي، انواع مختلف ضبط اسامي خارجی، سرهم‌نویسی، جدانویسی، بی‌فصله‌نویسی، انواع جمع‌ها، صورت‌های مختلف نوشتاری، آوانویسی اسامي عناصر و ترکیبات شیمیایی، سرواژه‌ها و کوته‌نوشت‌ها می‌باشد.

به طور کلی نقص‌هایی که برای زبان فارسي شمرده‌اند به شرح زیر می‌توان عنوان کرد:

۱. سه مصوت کوتاه یعنی حرکات زیر و زیر و پیش (اَ ، اِ ، اُ) را از نوشتن ساقط می‌کنیم. و این باعث می‌شود به جای این که از خط و نوشتار پی به معنی ببریم بایستی از معنی کلمه و جایگاه آن در جمله آنرا درست بخوانیم مانند کلمات (کَرَم، کَرَم، کَرَم، کِرْم، کُرْم، کِرْم) و (مَلِك، مَلِك، مَلِك، مَلِك) و یا سه کلمه (حَكَم، حُكَم، حِکَم) و نیز نوشتن مصوت‌های کوتاه در داخل متن باعث می‌شود که برای تلفظ صحیح اجباراً لاتین کلمات به صورت پانویس متن آورده شود که همین امر باعث اتلاف وقت و انرژی می‌شود. که البته همین لاتین‌نویسی هم قاعده خاصی ندارد و هر ناشر و نویسنده‌ای سلیقه خاص خودش را برای آوانویسی حروف فارسي به لاتین دارد، که به عنوان نمونه برای نشان دادن حرکت فتحه و الف و آ هیچ‌گونه هماهنگی در کتاب‌ها و خصوصاً فرهنگ‌های مختلف دیده نمی‌شود. "هر چند برخی معتقدند همین نوشتن حرکات مزیتی است و موجب تندنویسی می‌شود" [۵].

۲. برای يك حرف چند علامت مختلف داریم مانند علامت‌های (س، ص، ث) که هر سه در فارسي یکسان خوانده می‌شوند و هم چنین (ذ، ز، ض، ظ) و نیز (ت، ط). البته این امر در زبان انگلیسی هم وجود دارد چنان که «ف» ممکن است به شکل‌های «F. GH. PH. V» باشد.

۳. يك علامت را برای دلالت بر چند حرف مختلف استعمال می‌کنیم مانند "و" که پنج مورد نوشتن دارد یکی برای بیان ضمه در کلمات "خوش" و "تو". دیگر بیان مصوت ممدود یا "واو ماقبل مضموم" مانند "شور" و "او". سوم بیان حرف صامت "واو" در

کلماتی چون "آواز" و "والی" و "عفو". چهارم بیان حرف مصوت مرکبی که در کلمات "نو" و "جوشن" و مانند آنهاست. پنجم حرفی که در زبان کنونی خوانده نمی‌شود مانند "واو معدوله" در کلمات "خواهر" و "خواستن" و "واو" در کلمه "عمرو" [۶].
 ۴. حرف‌هایی هم هست که در کلمات خاصی از نوشتن حذف می‌شود مانند "الف" در کلمات "اسحق" و "اسمعیل" و "الله".
 ۵. نقطه‌هایی متعدد در بالا و پائین حرف که هم سبب دشواری و هم موجب اشتباه در خواندن می‌شود. اهمیت بیش از حد نقطه در خط فارسی هنگام تشخیص نوری کاراکترها [۷] تولید اشکال اساسی می‌کند. به عنوان مثال در نظر بگیرید که تفاوت <ر> و <ز> و یا تفاوت <د> و <ذ> و یا تفاوت <ب> <ت> <پ> <ث> فقط در نقطه است و چون نقطه جزء بسیار کوچکی است در این امر مشکلات زیادی را فرا روی متخصصین قرار می‌دهد. و یا کلمات زیر را در نظر بگیرید که با یک یا چند نقطه عوض می‌شوند (بر، پُر، پَر، تَر، پَز، پَر، بَر، تَز).

۶. یک عیب دیگر هم که برای خط فارسی ذکر کرده‌اند این است که از راست به چپ نوشته می‌شود. و برای این مورد دلایل مختلفی ذکر شده است از جمله عدم هماهنگی و ایجاد مشکل در نوشتن متون ریاضی و شیمی و نت‌های موسیقی و دستورات شطرنج و این که خط تصویری یعنی علائم گرافیکی که در کل جهان استفاده می‌شود مانند علائم راهنمایی و رانندگی تماماً از چپ خوانده می‌شوند.

۷. پیوسته‌نویسی و جدانویسی کلمات مرکب که در اکثر موارد به صورت سلیقه‌ای عمل می‌شود مانند تنوع استفاده از <می> چسبان و غیر چسبان و یا تنوع نحوه به کار بردن «علامت‌های جمع <ها، ان، جات>، هم، هیچ، که، (ضمایر شخصی متصل مان، تان، شان)، شناسی، را، چه، چون، تر، ترین، بی (پیشوند نفی)، به، ای (نشانه ندا)، آن و این» در کلمات به صورت پیوسته و یا جدا گانه: (آنچه، آن چه)، (همچنانکه، همچنان که)؛ (جناب‌عالی، جناب‌عالی)؛ (هیچکس، هیچکس)؛ (میتواند، می‌تواند)؛ (آن‌ها، آنها) در این مورد کلماتی که پیشوند و یا پسوند دارند نیز در شکل‌های مختلف نوشته می‌شوند. برخی از کلمات در دو شکل متصل‌نویسی و منفصل‌نویسی به دو شکل مختلف ظاهر می‌شوند، مانند «علاقمند و علاقه‌مند؛ اندیشمند و اندیشه‌مند». مصدرها و فعل‌های مرکب و اسم‌های مشتق از آنها نیز به دو صورت متصل و منفصل نوشته می‌شوند مانند «نگه‌داشتن و نگهداشتن». در جستجوی مطالب از اینترنت این مورد تولید اشکال می‌کند چنانکه جستجوی «هیچکس» نتایج متفاوتی را با جستجوی «هیچکس» می‌آورد و یا جستجوی «کتاب‌شناسی» و «کتاب‌شناسی» در موتور جستجوی گوگل نتایج متفاوتی را ارائه می‌کند. این گونه کلمات با این که در خواندن متن اشکال کمی به وجود می‌آورند و هر آشنای به زبان فارسی به راحتی می‌تواند آن را بخواند اما در فن‌آوری امروزه و تجزیه و تحلیل کلمات به کمک رایانه اشکال اساسی تولید می‌کند و شاید اگر قاعده‌ای جامع و مانع برای آن وضع گردد، بتوان گفت بزرگ‌ترین مشکل خط فارسی حل شده است. منظور این که، برای مثال خواندن سه کلمه «بی‌حوصلگی، بی‌حوصلگی، بی‌حوصله‌گی» مشکلی ایجاد نمی‌کند. اما در محیط الکترونیکی و شبکه اینترنت برای بازیابی این کلمه بایستی برای تمام اشکال این کلمه، جستجو را انجام دهیم، البته اگر آگاهی از تمام اشکال نوشتاری آن داشته باشیم. آ

۸. سی و دو حرف الفبای فارسی همراه با چهار علامت مد، همزه، تنوین، تشدید به ۱۳۰ شکل مختلف ظاهر می‌شوند و تفاوت این اشکال در اتوماسیون خط فارسی تولید اشکال می‌کند. «تنوع و تعدد نویسگان، یادگیری زبان و خط فارسی را برای آموزگار و آموزنده دشوار و برای نوآموز توان‌فرسا می‌سازد. تعداد زیاد نویسگان در رابطه با اتوماسیون زبان توسط رایانه مشکلاتی در خصوص تعداد و ترتیب قرار گرفتن نویسگان در جداول کد ایجاد می‌نماید و طراحان کد در جای دادن این تعداد نویسه در جداول با مساله کمبود جا رو به رو هستند. هر چند که مشکل جا با کد ۱۶ بیتی حل شده است اما مسایل دیگری همچنان باقی می‌مانند که احتیاج به برطرف شدن دارند» [۸]

۹. نوشتن ك و گ (ك گ ك گ ك گ ك) در اشکال مختلف نیز باعث سردرگمی و عدم جستجوی صحیح می‌شود.

۱۰. در اغلب اوقات يك فاصله اضافی معنی متفاوتی و یا متضادی را می‌دهد (مثل مادر، ما در).

۱۱. سه کرسی مختلف برای حرف‌های مختلف الفبا باعث می‌شود که در مقایسه با اکثر زبان‌ها تعداد سطرهای هر صفحه به مراتب بیشتر گردد چون برخی حروف روی خط کرسی قرار می‌گیرند و برخی پائین خط کرسی و برخی بالای خط کرسی مثل (ا ب م)

۱۲. از آنجائیکه حروف در نوشتن غالباً به صورت چسبیده و پیوسته نوشته می‌شوند و این امر تشخیص حرف به حرف نوشته به وسیله رایانه را، دچار مشکل می‌کند.

۱۳. در او. سی. آر. فارسی هم چنین اعداد نیز مشکل ساز هستند چنانچه صفر در فارسی يك نقطه کوچک است که می‌تواند رایانه را به اشتباه بیاندازد و نیز اعداد ۱ و ۲ و ۳ بسیار شبیه هم هستند و تفاوت‌شان در يك دندان کوچک است.

۱۴. تنوع املائی یا تنوع در رسم الخط بعضی از کلمات که همه شکل‌های آن نیز درست است مانند (اتاق و اطاق) و یا (امپراتور و امپراطور). و کلماتی که فقط يك شکل آنها صحیح می‌باشد ولی شکل ناصحیح آن نیز زیاد استفاده می‌شود مانند «ذغال و زغال؛ خوشنود و خشنود». البته این جدای از تنوع در مفهوم کلمات است که در دیگر زبان‌ها نیز وجود دارد، یعنی

