



# پیاده‌سازی یک سیستم کنترل خطای املایی در زبان فارسی بر اساس کدگذاری Soundex

انور بهرامپور<sup>۱</sup>، فردین اخلاقیان طاب<sup>۲</sup>، جلال سجادی<sup>۳</sup>، وفا بارخدا<sup>۴</sup>

<sup>۱</sup>گروه فناوری اطلاعات، دانشگاه آزاد اسلامی- واحد سنترج، سنترج، ایران [Bahrampour@iausdj.ac.ir](mailto:Bahrampour@iausdj.ac.ir)

<sup>۲</sup>گروه کامپیوتر و فناوری اطلاعات، دانشگاه کردستان، سنترج، ایران [Akhlaghian,J.sajadi,Barkhoda@uok.ac.ir](mailto{Akhlaghian,J.sajadi,Barkhoda@uok.ac.ir})

## چکیده

در بسیاری از کاربردها مانند پردازشگرهای متن، موتورهای جستجو، لغتتامه‌های الکترونیک، کاربردهای تلفن همراه و ... عمل کنترل املا بخشی از عملکرد کلی سیستم بشمار می‌رود. از آنجا که کارایی الگوریتم‌های موجود در زبانهای مختلف متفاوت بوده و این راهکارها نیازمند تعریف ساختارها و الگوریتم‌های مناسب با زبان است، بررسی ساختارها و الگوریتم‌های کنترل املا در زبان فارسی نیز بسیار ضروری است. در این مقاله انواع خطاهای املایی، ساختار سیستمهای کنترل املا و تعدادی از الگوریتم‌های مورد استفاده در تشخیص و تصحیح خطاهای املایی بررسی شده است. سپس با تعریف یک سیستم کدگذاری بر اساس روش Soundex برای زبان فارسی، یک سیستم کنترل املا در زبان فارسی با کارایی قابل قبول پیاده سازی شده است. سیستم پیشنهادی برای ذخیره کردن لغات در فرهنگ‌لغت از درخت B استفاده می‌نماید.

## کلمات کلیدی

سیستم کنترل خط، تشخیص خط، تصحیح خط، روش کدگذاری Soundex

روی ویرایشگرهای مختلف و نیز قابلیت تصحیح جملات از نظر گرامری پا به عرصه وجود گذاشتند.

یک سیستم کنترل املا ساده هر واژه را در محتواهای فرهنگ لغت جستجو می‌کند، اگر آن را در فرهنگ لغت پیدا نمود آن واژه و یا ریشه آن یک کلمه معتبر بوده در غیر این صورت آنرا عنوان یک خطای املایی شناسایی می‌کند؛ بدیهی است که تمامی واژه‌های زبان در فرهنگ لغت وجود ندارد، برای مثال ریشه افعال در فرهنگ لغت قرار گرفته و شکل‌های مختلف آنرا باید با اعمال قواعد زبان تولید نمود، بنابراین عملکرد یک سیستم کنترل املا فقط یک جستجوی ساده در فرهنگ لغت نیست [۱,2,3] البته عمل تشخیص خطای املایی محدود به استفاده از فرهنگ لغت نیست و روش‌های دیگری نیز برای تشخیص خطای املایی در نظر گرفته و در بخش‌های بعدی به تعدادی از آنها اشاره شده است. با شناسایی واژه نادرست در یک متن تلاش برای یافتن واژه صحیح در سیستم کنترل املا شروع می‌گردد. اگر کلمه‌ای در لغتتامه یافت نشد، آن را به عنوان کاندیدای خطای در نظر گرفته و در این صورت امکان دارد عملکرد سیستم اینگونه باشد که لیستی از کلمات مشابه و نزدیک به آن کلمه را به ترتیب میزان مشابهت برای کاربر تولید کند. برای تعیین میزان مشابهت کلمات از الگوریتم‌های مختلفی استفاده می‌شود که در این میان

سیستم‌های کنترل خط<sup>۱</sup> به برنامه‌هایی گفته می‌شوند که کلماتی را که امکان دارد از لحاظ املایی غلط باشند، علامت گذاری می‌کنند و به ازای هر یک از کلمات نادرست تعدادی واژه مشابه به عنوان جایگزین پیشنهاد می‌دهند. عملکرد این برنامه‌ها شامل دو بخش تشخیص خط و تصحیح خط است [۱]. این برنامه‌ها ممکن است به عنوان نرم‌افزاری مستقل کار تشخیص و تصحیح خط را انجام دهند و یا اینکه بخشی از برنامه‌های بزرگتر همانند پردازشگرهای متن، لغتتامه‌های الکترونیک و یا موتورهای جستجو باشند. سیستم‌های کنترل املا برای یافتن و تصحیح خودکار خطاهای املایی در متن زبان‌های طبیعی، بسیار مورد ملاحظه هستند و به طور پیوسته بر محبوبیت آنها افزوده می‌شود. به این برنامه‌ها با عنوان تصحیح کننده خودکار خط<sup>۲</sup> نیز ارجاع می‌شود [۲]. تلاش برای یافتن الگوریتم‌ها و تکنیکهایی برای تشخیص و تصحیح خودکار خطاهای املایی با گسترش ویرایشگرهای متنی، در اوایل دهه ۱۹۶۰ شروع شد [۱] و به تدریج پس از ظهور سیستمهای با قابلیت محدود تشخیص خط، سیستمهای با امکان تشخیص و تصحیح خط و قابل نصب بر

<sup>1</sup> Spell Checker

<sup>2</sup> Automatic Error Corrector