

OCR

فناوری تبدیل تصویر متن به خود متن

(به زبان ساده)

● فرزاد دادرس*

کارشناس ارشد کتابداری و اطلاع‌رسانی و

مدرس گروه ICT دانشگاه صنعتی مالک اشتر

چکیده

این مقاله دربردارنده تجربه‌ای نوین در عرصه فناوری‌های رقوم‌سازی مدارک^۱ است. نویسنده به همراه گروهی از همکاران خود در عرصه استفاده کاربردی و مؤثر از این فناوری، تجربه‌ای چندساله را پشت‌سر گذاشته است امید می‌رود که بیان این تجربه مورد استفاده پژوهشگران و دانشجویان نیز باشد. در این مقاله داستانی، نخست مشکلاتی که سازمان‌ها و مراکز اطلاع‌رسانی با آن روبه‌رو هستند، مورد بحث قرار گرفته تا نیاز به استفاده بهینه از این فناوری برای همه خوانندگان به بهترین نحو آشکار شود. وقتی مشخص شد که چرا به فناوری OCR نیازمندیم، به توضیح در مورد آن پرداخته‌ام. در این مسیر، علاوه بر ارائه اطلاعات تخصصی به زبان ساده، فرایندهای مدیریتی این‌گونه مسائل را در قالب مدیریت پروژه‌های درون‌سازمانی، توضیح داده‌ام و سعی کرده‌ام تا بخشی از تجربیات کاربردی گذشته نیز در متن داستان به خوانندگان تقدیم شود. کیدواژه‌ها: او. سی. آر، تشخیص نوری، کارکترها، شناسایی نوری حروف

مقدمه

این مقاله را در دوره‌ای می‌نویسم که دیگر داد سخن دادن از برخی فناوری‌های نوین مانند دیجیتال‌سازی، می‌رود که به تدریج از مد و تب و تاب مقاله‌نویسان بیفتد. شاید یکی از دلایل از مد افتادن این موضوعات - در عرصه ارائه اطلاعات و دانش جدید - آن باشد که بسیاری از سخن‌پراکنان این نوشته‌ها، پس از برق‌انداختن واژه‌های قرضی، پر زرق و برق فناوری، قادر به حرکت با نیازهای عملیاتی سازمان‌ها و مراکزی که خواهان استفاده از این فناوری‌ها هستند نبوده و در عرصه ارائه تجربه‌های کاربردی، حرفی برای گفتن و تجربه‌ای برای به‌کارانداختن ندارند.

به هر صورت، زمانی که یورش دانش جدید که نه، یورش واژگان جدید توسط مترجمان غیرحرفه‌ای و گاه ناشی اما شیک‌پسند، انبوه این واژگان را به حوزه‌های خاص سرازیر می‌کند، موضوعات زخمی شده و به حال خود رها می‌شوند؛ در چنین شرایطی، تولد حلقه‌های معیوب تولید اطلاعات و دانش جدید، چندان دور از انتظار نیست. هرچند ممکن است برخی بگویند که مسئولیتی برعهده مترجمان به لحاظ اجرایی کردن و تسلط عملی و کاربردی به مطالبی که ترجمه می‌کنند نیست، اما داشتن تسلط موضوعی وظیفه اصلی مترجم است که در بسیاری از موارد دیده نمی‌شود. می‌خواهم بگویم که مقاله علمی و پژوهشی، الزاماً نشانه‌ای از دانش و تسلط موضوعی پدیدآورنده نیست.

ناامیدانه هم سخن نمی‌گویم، چرا که در این میان، پژوهشگران و تولیدکنندگان حقیقی دانش هم بسیارند و باری را که مصرف‌کنندگان و نمایش‌دهندگان واژه‌ها بر زمین می‌گذارند، بر دوش گذاشته و به دست اهلس می‌سپارند؛ مسئولیتی که بر عهده سازمان‌ها و انجمن‌های تخصصی است تا دست کم با تشخیص این دو گروه از یکدیگر - و نه با دادن امتیازی بیشتر - توان‌شناختی خود را به رخ هر دو دسته بکشند!

حاصل آنکه این نوشته درباره کارها و واژه‌های شیک، که فرنگی‌ها تولید کرده‌اند داد سخن نداده است، بلکه تنها یک معرفی ساده، کوتاه و داستان‌گونه از یک فناوری پیشرفته است به پشتوانه تجربیات کاربردی نویسنده و تنی چند از همکاران او بویژه در حوزه OCR فارسی.

اما داستان...

من یک کتابدار هستم. هفته‌ای چند روز بعدازظهرها وقتم را در کتابخانه‌ای که اتفاقاً دوستم کتابدار آنجاست به مطالعه سپری می‌کنم. داستان مربوط می‌شود به ماجرای که دوست کتابدارم برایم تعریف کرد. او یکی از تکنیکی‌ترین کتابدارانی است که تاکنون دیده‌ام و البته یکی از دلسوزترین آنها. او سرش برای رفع نیاز اطلاعاتی مراجعه‌کنندگان درد می‌کند، و اما آغاز داستان از زبان او:

همان‌طور که در کتابخانه نشسته بودم و برای رفع خستگی، گاهی

ما به شکلی از فرمت الکترونیکی منابع، مخصوصاً کتاب‌ها نیاز داریم که به سادگی قابل جست‌وجو باشد تا کاربر به سادگی و خیلی سریع بتواند موضوع یا موضوع‌های مورد نظر خود را در کتاب پیدا کند و حتی بتواند آن را برای خود چاپ کند یا به صورت فایل مورد استفاده قرار دهد

آن را برای مشاوره با دیگران مطرح کنیم. در نتیجه ضروری دیدم در مورد راه‌حلی‌هایی که به ذهنم رسیده بود فکر کنم، سپس در صورت نیاز مشورت کنم.

راه‌حلی‌هایی که به ذهن کتابدار رسید

بدیهی بود که اولین راه‌حل باید می‌توانست سرعت چشمگیری در حل مشکل کتابخانه داشته باشد، چرا که به هیچ دلیلی نمی‌شود نیاز مراجعه‌کننده را معطل گذاشت. همچنین این راه‌حل باید همزمان مشکل فرسودگی کتاب‌ها به دلیل استفاده بیش از حد و مشکل ناکافی بودن تعداد کتاب‌ها را نسبت به تعداد زیاد مراجعه‌کننده‌ها حل می‌کرد.

در نتیجه جامع‌ترین راه‌حلی که در نگاه اول به ذهن او رسید، تهیه شکل یا فرمت الکترونیکی کتاب‌ها بود. به این ترتیب، هم کتاب‌های نفیس از فرسودگی در امان می‌ماندند و هم می‌شد به تعداد کافی و با همان کیفیت و به سرعت تعداد مناسبی از آنها را تکثیر کرد. اما مشکل اینجا بود که این کار یک فرایند طولانی و پرهزینه بود. طولانی بودن آن استفاده‌کننده‌ها را مدتی از استفاده محروم می‌کرد و پرهزینه بودن آن (تهیه سخت‌افزار مناسب و ...) برای کتابخانه امکان‌پذیر نبود. پس باید از روش دیگری برای رسیدن به هدف اقدام می‌کرد و همان‌طوری که خودش بعدها تعریف کرد، گفت می‌خواهد تصور کند که کسی پیش از او این مشکل را داشته و این کار را انجام داده است و من نباید آن را دوباره انجام دهم و برای آن هزینه کنم. در نتیجه دوست کتابدارم فکر کرد که شاید - بخوانید شاید -

کتابخانه ملی بتواند در این زمینه به من کمک کند یا فهرستی از مراکزی که اقدام به ساخت الکترونیکی منابع مکتوب می‌کنند و یا منابعی را الکترونیکی کرده‌اند داشته باشد و شاید کتاب‌های مورد نظرشان در این فهرست‌ها موجود باشد و ما بتوانیم آنها را بخریم. بنابراین به سرعت با کتابخانه ملی تماس گرفت و گفت که با مسئول بخش اشتراک منابع کاردارد.

مسئول بخش اشتراک منابع گفت که نسخه‌های الکترونیکی کتاب‌های مورد نظر ما در کتابخانه ملی موجود نیست، اما کتابخانه ملی می‌تواند با استفاده از تجهیزات سخت‌افزاری حرفه‌ای (اسکنر یا پوشگر)، کتاب‌های مورد نظر را به شکل الکترونیکی درآورده و در مقابل یکی از نسخه‌های الکترونیکی کتاب را در قبال این خدمت برای استفاده در آرشیو منابع الکترونیکی کتابخانه ملی نگهداری کند.

در نتیجه، کتابدار زیرک ما توانست با هماهنگی کتابخانه ملی و

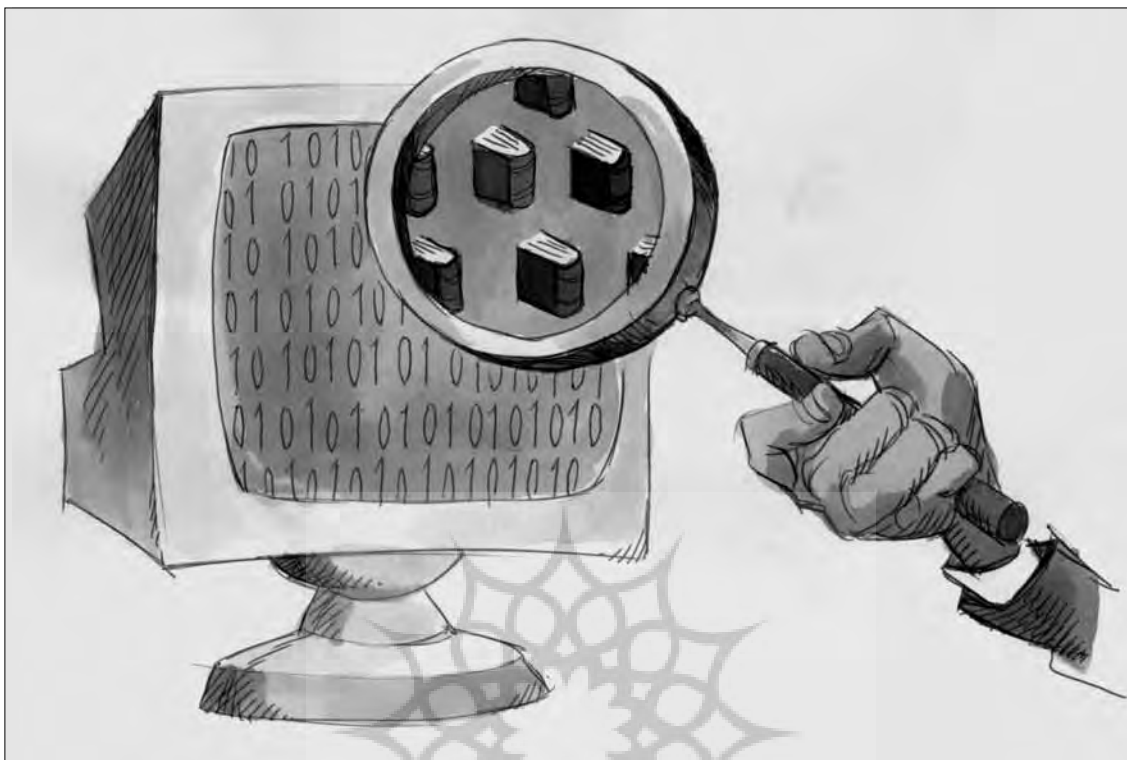
به سقف قشنگ کتابخانه و گاهی به رفتار مراجعه‌کننده‌هایی که برای پیدا کردن کتاب‌ها، مقالات یا عکس‌های مورد نظرشان به کتابخانه آمده بودند نگاه می‌کردم، چشمم به قفسه‌ای افتاد که عده‌ای را دور خود جمع کرده بود. گویا موضوع مربوط می‌شد به چند جلد کتاب ارزشمند و قدیمی که تعداد استفاده‌کننده‌های آنها به طرز چشمگیری بیشتر از تعداد جلد‌های موجود در کتابخانه بود. عجیب بود که روزهای بعد این صحنه را بارها و بارها در آن کتابخانه و چند جای دیگر هم دیدم. دوست کتابدار من هم به این جمع ملحق شد تا مسئله را حل و فصل کند. مسئله دقیقاً همان بود: کمبود منابع، تعداد انبوه استفاده‌کننده‌ها و مشکل اساسی‌تر اینکه: تهیه این منابع به دلیل نایاب بودن امکان‌پذیر نبود.

کتابدار چه باید می‌کرد؟ در آن لحظه به کدام یک از استفاده‌کننده‌ها باید کتاب را امانت می‌داد؟ مسئله حساس و کلیدی دیگری نیز وجود داشت: بر سر کتاب بعد از آن همه استفاده چه می‌آمد؟ طبیعی است که بعد از مدتی چیزی به شکل و هویت کتاب از آنها باقی نمی‌ماند. شاید به این دلیل که ماهیت کتاب هم مثل شمع است؛ نور می‌افشاند و می‌میرد.

بعدها دوست کتابدار جریان را این‌طور برایم تعریف کرد: پیش از هر تصمیمی، سعی کردم تا صورت مسئله را تا حد امکان به صورت روشن و بدون ابهام در ذهنم مرتب‌کنم و بعد هم آن را روی کاغذ بنویسم. این کار را به خاطر این کردم که کلک همیشگی ذهن را بلد بودم و چند بار این کلک را از ذهن خودم خورده بودم:

وقتی فکر می‌کنی چیزی را می‌دانی یا فهمیده‌ای، بهترین راه، پیش از منتقل کردن آن به دیگران، نوشتن آن مطلب است. مخصوصاً اگر پیچیدگی یا اهمیت خاصی داشته باشند. تازه در مرحله نوشتن است که متوجه می‌شوی کجای‌های مسئله را نفهمیده‌ای یا برایت مهم است و تلاش دوباره‌ها برای تجزیه و تحلیل مسئله باعث می‌شود که آن را بهتر درک کنی و در نتیجه بهتر بتوانی به دیگران منتقل کنی. بعد هم خدا را شکر می‌کنی که زودتر در مورد آن مطلب در جایی صحبت نکرده‌ای و حالا خوشحالی که می‌توانی حق مطلب را به خوبی ادا کنی.

پس صورت مسئله را نوشتم و برای درک کامل ابعاد مختلف آن به قدر کافی وقت گذاشتم و آن را چندین بار اصلاح و تکمیل کردم تا مطمئن شدم تمام زوایای آن برایم روشن شده و هیچ نقطه ابهامی وجود ندارد. بعد به خودم گفتم که حالا می‌توانم مسئله را حل کنم یا



داریم، اما نمی‌دانستیم که چگونه می‌توانیم این خدمات را فراهم کنیم و ارائه بدهیم. نیاز دقیق ما این بود:

«ما به شکلی از فرمت الکترونیکی منابع، مخصوصاً کتاب‌ها نیاز داریم که به سادگی قابل جست‌وجو باشد تا کاربر به سادگی و خیلی سریع بتواند موضوع یا موضوع‌های مورد نظر خود را در کتاب پیدا کند و حتی بتواند آن را برای خود و چاپ کند یا به صورت فایل مورد استفاده قرار دهد، اما چگونه؟»

گوشی را برداشتم و شماره تلفن مشاور کتابخانه در فناوری‌های نوین را شماره‌گیری کردم.

مشاوران کتابخانه‌ها

در بسیاری از مواقع، مشاوران کتابخانه‌ها و مراکز اطلاع‌رسانی و بقیه سازمان‌هایی که به نوعی درگیر استفاده از فناوری‌های جدید برای بهبود مسائل و مشکلات خود می‌باشند، افراد جوان، تندذهن و کارآزموده‌ای هستند که تجربه‌های کارآمد و اثربخشی را در حوزه دانش مربوط به فناوری و پیاده‌سازی آن کسب کرده‌اند. پس طبیعی بود که صدای جوانی از آن طرف خط به کتابدار کتابخانه پاسخ داده باشد. دوست کتابدارم گفت که مسئله را به دقت و کاملاً روشن برای مشاور تعریف کردم و قرار شد تا در اولین فرصت در کتابخانه جلسه‌ای را برای تبادل نظر و بررسی راهکارها برگزار کنیم.

راهکارهای دیگر فناوری از زبان مشاور کتابخانه

مشاور کتابخانه هم یک کتابدار مسلط به مسائل فناوری‌های نوین و مرتبط به کتابخانه و مراکز اطلاع‌رسانی و نیز یک تحلیلگر کارآزموده سیستم بود. مشاور جوان، مطالب خود را این‌گونه آغاز کرد: با توجه به اینکه کتابدار دقیق و موشکاف کتابخانه، صورت مسئله

در طول چند روز تعطیلی که پیش آمده بود کتاب‌های مورد نظر را به کتابخانه ملی منتقل کرده و آنها را اسکن کند و شکل الکترونیکی آنها را به دست آورد؛ لازم به توضیح است که کتابخانه ملی هیچ وقت تعطیل نیست.

مشکلات بعدی کتابخانه

مشکل اول حل شد. حالا هم مسئله فرسودگی منابع حل شده بود و هم مراجعه‌کنندگان به راحتی و بدون هیچ محدودیتی می‌توانستند به هر تعداد از نسخه‌های کتاب‌های مورد نظر خود دسترسی داشته باشند. حتی اگر ده‌ها و صدها مراجعه به این منابع صورت بگیرد، هیچ محدودیتی در تکثیر ثانیه‌ای این منابع وجود ندارد^(۲).

اما مشکل جدید، زمانی خود را نشان داد که مراجعه‌کنندگان می‌خواستند به صفحات خاصی از کتاب خود مراجعه کنند. به همین دلیل، ناچار بودند که با استفاده از صفحه فهرست مندرجات کتاب و برداشت شماره صفحه مورد نظر خود، صفحات کتاب را - که حالا تبدیل به عکس شده بودند - تک تک ورق بزنند تا به مطلب خود برسند؛ چه کار طولانی و پر دردسری!

این بار که استفاده‌کننده متوجه حضور یک کتابدار کارشناس و باهوش در کتابخانه شده بودند، خودشان سراغ او رفتند و مشکل را توضیح دادند. حل این مسئله دیگر در حوزه تجربه و تخصص کتابدار نبود. او حالا لازم می‌دید که باب مشورت را با دیگر کارشناسان باز کند. او این‌طور ادامه داد:

مسئله بسیار روشن بود و ما تجربه‌ای برای حل این مسئله نداشتیم و تنها چیزی که می‌دانستیم - و البته مهم‌ترین نکته بود - این بود که می‌دانستیم به چه خدماتی در مورد این مسئله خاص نیاز

بخش رقومی سازی هر مرکز یا کتابخانه قبل از اقدام به شروع پویش و رقومی سازی مدارک خود، باید بر حسب شرایط، امکانات و توانایی های خود، فرایند هدفمند و قابل پیاده سازی را برای سازمان خود طراحی و بومی سازی کند. در این فرایند، مراحل و چگونگی انجام کارها، فرایندهای کنترلی و نیروی انسانی مورد نیاز به دقت تدوین شده اند و تازه بعد از این مراحل، آغاز کار منطقی است

را بسیار حرفه ای، جامع و روشن تعریف کردند، ما می توانیم بدون اتلاف وقت، برای ایجاد یک برداشت یکسان از مسئله، به سرعت وارد بحث اصلی یعنی راهکارهایی بشویم که فناوری پیش روی ما گذاشته است. اما پیش از شروع بحث مایل هستم تا دیگر مزایا و معایب استفاده از منابع پویش شده را توضیح دهم تا ضرورت نیاز به راه حل جدید، بیشتر برای شما روشن و قطعی شود و بعد برویم سراغ راه حل اصلی.

همان طور که خودتان هم به این نتیجه درست رسیده بودید، فناوری های رایجی که برای تولید خروجی های الکترونیکی وجود دارند قادرند توسط پوششگرها صفحات چاپی را به تصاویری از متن تبدیل کنند. باید گفت که این حرکت از اولین قدم ها برای ایجاد فرایند رقومی سازی منابع در کتابخانه است. گفتیم از اولین قدم ها و نه اولین قدم، زیرا اولین قدم، طراحی و تدوین خود فرایند در یک محیط انتزاعی است؛ حال می خواهد این فرایند مربوط به رقومی سازی مدارک باشد و یا هر فرایند دیگری، به هر حال طراحی آن قبل از اجرا ضروری است. بخش رقومی سازی هر مرکز یا کتابخانه قبل از اقدام به شروع پویش و رقومی سازی مدارک خود، باید بر حسب شرایط، امکانات و توانایی های خود، فرایند هدفمند و قابل پیاده سازی را برای سازمان خود طراحی و بومی سازی کند. در این فرایند، مراحل و چگونگی انجام کارها، فرایندهای کنترلی و نیروی انسانی مورد نیاز به دقت تدوین شده اند و تازه بعد از این مراحل، آغاز کار منطقی است.

طبیعی است که مجریان و ناظران طرح، باید دقیقاً مطابق فرایندهای طراحی شده، اجرای کار را دنبال کنند. در این صورت، ضریب اطمینان رسیدن به اهداف از پیش تعیین شده بسیار بالا رفته و احتمال انحراف و یا خطا در اجرا بسیار کاهش می یابد. این ضریب افزاینده در کیفیت مدیریت و اجرا، موجب می شود تا سازمان در هزینه های خود صرفه جویی کرده و حتی یک شیب کاهنده ایجاد کند و مهم تر اینکه به دلیل بد انجام شدن کارها ناچار به صرف هزینه های اضافی یا پذیرش نقص در خدمات رسانی نباشد. این، در حالی است که در بیشتر سازمان ها و مراکز و کتابخانه ها، درست عکس این داستان اتفاق می افتد: اول کار انجام می شود، بعد روی آن فکر می کنند! در این گونه سازمان ها، پروژه هایی که با فناوری ارتباط دارند، با اتکا به دانش ذهنی و تجربیات افرادی محدود انجام می شوند و هیچ رویکرد روشمندی که نگاهی نظام مند به مسائل داشته باشد و مدیران را ملزم به مکتوب کردن استراتژی های انجام کارها کند وجود ندارد. با این توضیحات، مروری خواهیم داشت بر دیگر مزایا و معایب روش

انتخاب شده:

مزایا:

- امکان ایجاد دسترسی کامل کاربران به منابع مورد نیاز (بدون محدود شدن به متغیرهایی مانند تعداد درخواست کننده ها)
- امکان حفظ و نگهداری بهتر مدارک اصلی

• تکثیر آسان

- امکان راه اندازی خدمات اشتراک منابع الکترونیکی
- امکان راه اندازی خدمات سفارش و دریافت الکترونیکی

معایب:

- محدودیت چشمگیر در تهیه رسانه های ذخیره سازی مدارک پویش شده
- عدم امکان جستجوی متن، بویژه در منابع فارسی در مورد مزایا به طور مشخص، می توان گفت که می توانید خدمات جدیدی را برای کاربران به فهرست خدمات پیشین خود بیفزایید. اما در مورد معایب و مشکلات این روش:

مشکل اول:

حجم فایل های تصویری از حجم فایل های متنی بیشتر است. این مسئله موجب می شود تا سازمان شما برای ذخیره سازی مدارک در حجم انبوه، هزینه های سنگینی بپردازد و این هزینه ها برای بسیاری از کتابخانه ها، مراکز اسناد و سازمان هایی که دارای حجم بالایی از مدارک هستند، مانعی جدی بر سر راه رقومی سازی و الکترونیکی کردن مدارک تلقی می شود.

مشکل دوم:

مشکل دوم را شما خودتان به خوبی و با دقت کافی توصیف کردید. یک مشکل بسیار جدی که حتی در صورتی که کتابخانه شما یا هر سازمان دیگری تصمیم به پرداخت هزینه های ذخیره سازی بگیرد، در خدمات رسانی نهایی به کاربر با آن روبه روست و مشکلات بسیاری را برای کاربران ایجاد می کند: مسئله جستجو در متون الکترونیکی. باید بگویم که در حال حاضر، فناوری، یک راه حل مطلوب برای حل این مسئله ارائه کرده است؛ راه حلی که شما را قادر می سازد تا خروجی های الکترونیکی متنی و کم حجمی تولید کنید که به راحتی مانند یک فایل متنی ساده قابل جستجو و تکثیر باشد و حجم بسیار پایین تری نسبت به فایل های تصویری داشته باشد: فناوری OCR.

فناوری OCR^۳

OCR فناوری تبدیل تصاویر گرفته شده از متن به خود متن است. خروجی حاصل از این فناوری، دو اشکال بزرگ ذخیره و بازیابی مدارک را که در قسمت قبل به آن اشاره شد حل می کند

OCR فناوری تبدیل تصاویر گرفته شده از متن به خود متن است. خروجی حاصل از این فناوری، دو اشکال بزرگ ذخیره و بازیابی مدارک را که در قسمت قبل به آن اشاره شد حل می کند.

نرم افزارهای OCR قادرند تصاویر پویش شده از متن را به میلیون ها نقطه نوری تبدیل کنند. به این ترتیب، نرم افزار می تواند یک درک محاسباتی از صفحه پویش شده برای خود ایجاد کند و نقاط سفید صفحه را صفر و نقاط تاریک یا همان نوشته های متنی را یک فرض کند. حال، نرم افزار می تواند نقاط تاریک یا نوشته ها را با الگوهایی از حروف که از قبل در نرم افزار ذخیره شده یا به آن آموزش داده شده، مقایسه کرده و حدس بزند که هر کدام از این سلسله نقاط تاریک شبیه کدام حرف از حروف فارسی یا انگلیسی است و آنگاه آن حرف را در یک فایل متنی قرار داده و برای تمام نقاط تاریک عکس، همین فرایند را ادامه می دهد تا متن کامل از تصویر به دست آید. طبیعتاً این نرم افزار از نقاط سفید یا صفر صفحه برای ایجاد حاشیه هایی که در صفحات وجود دارند و همچنین برای فاصله گذاری میان کلمات استفاده می کند.

به این ترتیب، نرم افزار تصاویر رقمی شده را به یک تصویر قابل پردازش تبدیل می کند. حال، ما می توانیم هر چیزی را به سادگی در متن جدید تولید شده جست و جو کنیم و مثلاً ببینیم که کلمه مورد نظر ما در چه جاهایی از متن به کار رفته و چند بار؛ کاری که به هیچ وجه نمی توانستیم با تصاویر ایجاد شده از متن انجام دهیم. اگر بخواهیم مروری بر دیگر کاربردهای این فناوری داشته باشیم، باید به کاربردهای زیر اشاره کنیم:

- آرشیو الکترونیکی متون مجلات
- آرشیو الکترونیکی پایان نامه ها و مقالات
- آرشیو اسناد و تصاویر متنی

ویژگی همه این آرشیوها وجود قابلیت جست و جو و حجم پایین به واسطه استفاده از فایل های متنی است.

انواع OCR

OCR یا روی متون تایپی انجام می شود یا روی متون دست نویس. یعنی یا باید متنی را که از قبل تایپ و چاپ شده و فایل متنی آن موجود نیست و ما فقط نسخه چاپی آن را داریم OCR کنیم یا متن های دست نویس را. متن های دست نویس هم دو دسته اند؛ یا گسسته اند یا پیوسته. متن گسسته دست نویس، مانند فرم هایی که مثلاً برای نام نویسی کنکور استفاده می کنیم و شما باید حرف های نام و نام خانوادگی تان را در مربع های کوچک بنویسید. منظور از متن پیوسته دست نویس هم همان دست نوشته های خودمان است که هر روز در کلاس دانشگاه یا دفتر کارمان روی کاغذ می نویسیم.

OCR انگلیسی و فارسی

OCR انگلیسی سال هاست که به بهره برداری رسیده و در حال حاضر استفاده های زیادی دارد. OCR فارسی در ایران دو سه سالی است که مورد توجه قرار گرفته و چند گروه تحقیقاتی و چند شرکت کوشیده اند تا نرم افزار مستقلی برای این کار تولید کنند، اما به دلایل متعدد، تاکنون نتوانسته اند تا در حد مطلوبی از پس مشکلات

OCR فارسی برآیند.

البته کارهایی که روی OCR گسسته فارسی انجام شده، به ظاهر در مراحل پایانی خود قرار دارد، ولی همان طور که اشاره شد، OCR پیوسته فارسی روی متون تایپی هنوز سال های بسیاری کار دارد. البته در مورد این مسئله استثنای بسیار مهمی وجود دارد، چرا که یکی از شرکت های کم سروصدایی که سال هاست روی OCR پیوسته متون تایپی فارسی کار کرده است به نتایج بسیار مهمی رسیده و در آینده نزدیک، خدمت بسیار جدیدی را در عرصه رقمی سازی مدارک به سازمان ها ارائه خواهد کرد و من در قسمت های بعد به آن اشاره خواهیم کرد.

مشکلات OCR فارسی

بسیاری از مشکلات OCR فارسی به ویژگی های زبان فارسی مربوط می شود. مثلاً ما در فارسی حروفی داریم که شباهت بسیاری به هم دارند مثل: «د» «ذ» «ر» «ز» «ژ» «ت» «ث». اعداد نیز دارای چنین شباهت هایی با یکدیگر هستند. مثلاً اعداد «۱» «۲» «۳» که تنها در یک دندان به هم تفاوت دارند و صفر در سیستم نوشتاری ما شبیه یک نقطه است. همچنین هر یک از حروف فارسی را می توان به اشکال مختلف نوشت. مثلاً حرف س در کلمات مختلف، یک جا س را با دندان می نویسند، یک جا بدون دندان و دیگری س را می کشند. حال، وقتی این حروف متفاوت به یکدیگر می چسبند، مشکلات بیشتر می شوند. تشخیص صحیح و دقیق حروف و اعداد فارسی و در نتیجه کلمات را دشوار کرده و ضریب خطای خروجی های نرم افزار OCR فارسی را افزایش می دهد.

مرور کلی

حال می خواهیم مروری کلی به فرایند OCR داشته باشیم و مراحل آن را دقیق تر به خاطر بسپاریم. مراحل انجام OCR به شرح زیر است:

- بازشناسی حروف
- بازشناسی الگو
- بازشناسی دنباله ها
- مدل سازی و پردازش زبانی

• بازشناسی حروف

در این مرحله، نرم افزار مکان قرار گرفتن حروف را از دیگر اجزای

کلمات غلط را با نمونه‌های صحیح آنها در بانک اطلاعات تطبیق داده و شکل صحیح کلمه را جایگزین نموده یا به اپراتور سیستم پیشنهاد می‌دهند. با قراردادن این کنترل‌ها اکنون نرم‌افزارهای OCR انگلیسی خروجی‌های مناسبی را ارائه می‌دهند، اما نرم‌افزارهای OCR فارسی که در داخل کشور طراحی و پیاده‌سازی شده‌اند، همچنان راه پیشرفت را طی می‌کنند و امید است که تا چند سال آینده به تولید انبوه و صنعتی برسند.

وضع کنونی OCR فارسی در ایران

با اینکه چند گروه تخصصی در صدد ایجاد نرم‌افزارهای ساخت داخل هستند و هنوز راه طولانی تا رسیدن به کیفیت مطلوب برای متون فارسی تایپی پیوسته وجود دارد، اما این مسئله موجب نشده است تا گروهی از کارشناسان یکی از شرکت‌های داخلی برای دادن خدمات OCR فارسی در انتظار تکمیل نرم‌افزارهای داخلی برای آغاز سرویس‌دهی بمانند. این گروه به نتایج چشمگیری در این زمینه رسیده‌اند و به زودی اطلاع‌رسانی عمومی خود را آغاز می‌کنند و شما هم می‌توانید از خدمات OCR فارسی در کتابخانه و مرکز اطلاع‌رسانی خود استفاده کنید.

نتیجه‌گیری مشاور

پس از تمام شدن صحبت‌های مشاور که بیشتر شبیه یک کارگاه آموزشی بود، همه راضی به نظر می‌رسیدند، چرا که تمام مطالب گفته شده در جلسه را آموخته بودند و علاوه بر آن، بینش و نگرش آنها در مورد فرایندهای رقمی‌سازی و طراحی این فرایندها به کلی تغییر کرده بود. این مطالب را می‌شد هم در مدیریت پروژه‌های فنی کتابخانه و هم در کارهای اجرایی به کار برد. ضمن اینکه مطالب تازه اشتیاق خاصی را در مخاطبان ایجاد کرده بود.

سرانجام، قرار شد تا مشاور کتابخانه با حضور شرکتی که خدمات OCR فارسی ارائه می‌کرد، مدیر کتابخانه و کتابدار متخصص جلسه‌ای بگذارد و در آن راهکارهای اجرایی کار را بررسی کنند، چرا که نیاز شدیدی در این قسمت از خدمات کتابخانه به کاربران استفاده‌کنندگان آن احساس می‌شد.

شاید شما هم به زودی تصمیم بگیرید که از خدمات ارائه متن کامل با قابلیت جست‌وجوی متنی و با کمک فناوری OCR فارسی بهره‌مند شوید.

پی‌نوشت‌ها

* Dadrass_f@yahoo.com

1. Document Digitalization

۲. البته ممکن است مشکلات مالکیت معنوی (کپی رایت) محدودیت‌هایی را در تکثیر منابع ایجاد کند.

3. Optical Character Recognition

4. Optical Character Reader

5. Spell Checker

منبع

۱. نشانی اینترنتی OCR: آموزش الفبای فارسی به رایانه.

<http://www.iritn.com/index.php?action=show&type=news&id=6080>

نرم‌افزارهای OCR قادرند تصاویر پویشی شده از متن را به میلیون‌ها نقطه نوری تبدیل کنند. به این ترتیب، نرم‌افزار می‌تواند یک درک محاسباتی از صفحه پویشی شده برای خود ایجاد کند و نقاط سفید صفحه را صفر و نقاط تاریک یا همان نوشته‌های متنی را یک فرض کند

صفحه (مثلاً خطوط کادر) تشخیص می‌دهد و حروف را از یکدیگر جدا می‌کند.

• بازشناسی الگو

در این مرحله نرم‌افزار با تعدادی، شرط می‌فهمد که حرف جدا شده مثلاً «گ» هست یا نه؟ در حقیقت، در این مرحله «گ» جدا شده توسط نرم‌افزار با «گ» های نمونه خود نرم‌افزار مقایسه و بازشناسی می‌شود.

• بازشناسی دنباله‌ها

کار با تشخیص صحیح حروف تمام نمی‌شود. در این مرحله، نرم‌افزار باید بفهمد که این حرف فقط «گ» است یا به حروف دیگری نیز چسبیده است؟ به این ترتیب نرم‌افزار دنباله‌های حروف را نیز در کلمات تشخیص می‌دهد.

• مدل‌سازی یا پردازش زبانی

حروف به هم چسبیده که کلمات را درست می‌کنند، باید معنی دار شناخته شوند. در این مرحله، بررسی می‌شود که چه کلماتی در زبان وجود دارند و چه ترکیبی از کلمات مجاز هستند. برای تشخیص ترکیب‌های مجاز یک کلمه یا معنی دار بودن کلمه، به تهیه بانک‌های اطلاعاتی نشان استاندارد نیاز است. در مراحل پیشرفته‌تر ایده‌آل مدل‌سازی گرامری و مدل‌سازی معنایی که نشان می‌دهد جمله‌ها از نظر دستوری و معنایی درست هستند یا نه انجام می‌شود. البته فرایندهای مرحله آخر، بیشتر در مرحله آزمایشگاهی قرار دارند.

فرایندهای کنترل در نرم‌افزارهای OCR

برخی از نرم‌افزارهای OCR از امکانات خاصی برای کنترل خروجی‌های خود استفاده می‌کنند، زیرا در این سیستم‌ها همیشه احتمال خطا هست.

• کنترل Reject

در این فرایند کنترلی نرم‌افزار آن دسته از حروفی را که نتوانسته تشخیص دهد، علامت‌گذاری و اعلام می‌کند. در این صورت، اپراتور نرم‌افزار یا به صورت دستی و یا با استفاده از نرم‌افزارهای کنترل املايي اقدام به تصحیح خطاها می‌کند.

• کنترل Spell Checker

در برخی دیگر از سیستم‌های OCR نرم‌افزارهای کنترل و تصحیح املايي کلمات قرار داده شده‌اند. این نرم‌افزارها که درحقیقت باید یک بانک اطلاعات استاندارد از کلمات یک زبان خاص، باشند