

بِسْمِ اللَّهِ الرَّحْمَنِ الرَّحِيمِ



دانشگاه سمنان

دانشکده مهندسی برق و کامپیوتر

پایان نامه کارشناسی ارشد مهندسی برق - الکترونیک

عنوان:

جداسازی خطوط در متون دست نویس فارسی

نگارش:

عاطفه سهرابی

استاد راهنما:

دکتر سعید مظفری

اسفند ۱۳۹۱

بسمه تعالی

اینجانب عاطفه سهرابی متعهد می‌شوم که محتوای علمی این نوشتار با عنوان *جد/سازی خطوط* در متون دست‌نویس فارسی که به عنوان پایان‌نامه کارشناسی ارشد رشته مهندسی برق گرایش *الکترونیک* به گروه مهندسی برق دانشکده مهندسی برق و کامپیوتر دانشگاه سمنان ارائه شده، دارای اصالت پژوهشی بوده و حاصل فعالیت‌های علمی این جانب می‌باشد. در صورتی که خلاف ادعای فوق هر زمانی محرز شود، کلیه حقوق معنوی متعلق به این پایان‌نامه از اینجانب سلب شده و موارد قانونی مترتب به آن نیز از طرف مراجع ذیربط قابل پیگیری است.

نام و نام خانوادگی

امضا

شماره دانشجویی

تقدیم به مهربان فرشتگانی که:

محطات ناب باور بودن، لذت و غرور دانستن، جسارت خواستن، عظمت رسیدن، و

تمام تجربه های زیبای زندگیم، دیون حضور سبز آنهاست.

تقدیم به پدر و مادر عزیز

و

همسر مهربانم

چکیده

هدف این رساله، تقسیم‌بندی یک سند دست‌نویس به خطوط مجزا از یکدیگر است. جداسازی و استخراج خط، اولین و مهمترین مرحله پیش‌پردازش برای آنالیز و بازیابی تصویر یک سند است. هنگامی که یک خط جداسازی شد، (خط ایزوله شده) برای مراحل جداسازی کلمه، شناسایی آن، بازیابی زیرکلمات و حروف و سایر مراحل که برای آنالیز یک سند مورد نیاز است، تحت بررسی قرار می‌گیرد. بنابراین استخراج صحیح خطوط، بمنظور درست انجام شدن سایر مراحل، دارای اهمیت فراوانی است. از طرفی جداسازی خطوط در اسناد دست‌نویس بدون محدودیت، به دلیل ویژگی‌های خاص این اسناد، یک مشکل اساسی است. تغییر در اندازه حروف، عدم وجود یک امتداد برای خط کرسی در سراسر یک خط، فاصله کم و متغیر بین خطوط و تماس و در هم فرو رفتگی خطوط مجاور و ... از جمله عواملی است که سبب می‌شود جداسازی خطوط با مشکلات زیادی مواجه گردد. در این رساله، ابتدا خط کرسی از یک خط با روش ورژن هموار شده نمودار منحنی طرح، تشخیص داده شده و سپس برای جداسازی کامل خط، به بررسی و نسبت دادن المان‌های موجود بین دو خط مجاور، به یکی از خطوط فوقانی یا پایینی پرداخته می‌شود، تا بدین ترتیب خط، جداسازی و استخراج گردد. یکی از مراحل اساسی که در این رساله برای استخراج خط کرسی و تخصیص مؤلفه‌های بین خطوط، مورد استفاده است، بکارگیری مدل مارکوف مخفی است. نتایج بدست آمده برای جداسازی خطوط به روش گفته شده، در زبان فارسی ۹۷.۷۵ درصد و برای سایر زبان‌ها ۹۸.۲۵ بدست آمده است.

کلید واژگان: اسناد دست‌نویس^۱، جداسازی خط^۲، خط کرسی^۳، بین خطوط^۴، مؤلفه‌های پیوسته^۵، مدل مارکوف مخفی^۶

¹ Handwritten Documents

² Line Segmentation

³ Base Line

⁴ Between lines

⁵ Connected Component

⁶ Hidden markov model

فهرست مطالب

صفحه	عنوان
	چکیده
۱	(۱) مقدمه
۱	(۱-۱) پردازش تصویر
۱	(۲-۱) تصویر سند
۲	(۳-۱) بازشناسی نوری حروف (OCR)
۳	(۴-۱) معرفی کلی یک سیستم بازشناسی متن
۵	(۵-۱) جداسازی خطوط
۸	(۶-۱) چالش‌های جداسازی
۸	(۱-۶-۱) نوسان و انحنای خطوط
۹	(۲-۶-۱) نزدیکی و مجاورت خطوط
۹	(۳-۶-۱) تکه‌تکه شدن دست‌نوشته
۱۰	(۷-۱) ویژگی‌های خط فارسی
۱۳	(۲) مروری بر تکنیک‌های موجود
۱۳	(۱-۲) مقدمه
۱۴	(۲-۲) تکنیک‌های موجود
۱۵	(۱-۲-۲) روش‌های بالا به پایین
۱۵	(۱-۱-۲-۲) روش‌های برمبنای طرح
۱۸	(۲-۱-۲-۲) روش‌های برمبنای مدل سند
۱۹	(۲-۲-۲) روش‌های پایین به بالا
۲۰	(۱-۲-۲-۲) کلاس‌بندی K_NN
۲۳	(۲-۲-۲-۲) تبدیل هاف
۲۳	(۳-۲-۲-۲) هموارسازی
۲۴	(۴-۲-۲-۲) شبکه جاذبه-دافعه
۲۵	(۵-۲-۲-۲) درخت پوشای کمینه
۲۵	(۶-۲-۲-۲) ریخت‌شناسی
۳۰	(۳-۲) خلاصه و نتیجه‌گیری

۳۱	۳) مدل مارکوف مخفی
۳۱	۱-۳) مقدمه
۳۲	۲-۳) مدل مارکوف مخفی
۳۴	۳-۳) نمایش و ساختار مدل
۳۶	۴-۳) ارزیابی مدل
۳۸	۵-۳) تخمین پارامترهای مدل
۳۹	۶-۳) انتخاب دنباله حالات بهینه
۴۱	۷-۳) توزیع‌های مشاهدات
۴۱	۱-۷-۳) توزیع‌های گسسته
۴۲	۲-۷-۳) توزیع‌های پیوسته
۴۴	۳-۷-۳) توزیع‌های نیمه‌پیوسته
۴۵	۸-۳) اندکی داده‌ای آموزشی
۴۶	۹-۳) خلاصه و نتیجه‌گیری
۴۷	۴) جداسازی خط
۴۷	۱-۴) مقدمه
۴۸	۲-۴) استخراج خط کرسی
۴۹	۱-۲-۴) مجموعه اولیه مناطق خط و فاصله در نواحی عمودی
۵۲	۲-۲-۴) بهبود دادن جداکننده‌های اولیه خطوط متنی
۵۵	۳-۲-۴) الگوریتم رسم جداکننده‌های خطوط متنی
۵۷	۴-۲-۴) تخصیص دادن مؤلفه‌های پیوسته به خطوط متنی
۶۱	۵-۲-۴) روش پیشنهادی در بررسی موقعیت مکانی نقاط و اعراب
۶۳	۳-۴) زیرکلمات فارسی
۶۴	۱-۳-۴) لغت‌نامه بیژن‌خان
۶۵	۲-۳-۴) استخراج اطلاعات مورد نیاز از لغت‌نامه بیژن‌خان
۶۵	۱-۲-۳-۴) زیرکلمات
۶۶	۲-۲-۳-۴) اطلاعات آماری
۶۸	۴-۴) خلاصه و نتیجه‌گیری
۶۹	۵) بررسی نتایج
۶۹	۱-۵) مقدمه

۶۹	۲-۵) پایگاه داده
۷۳	۳-۵) اطلاعات مشخصه پایگاه داده فارسی
۷۴	۴-۵) نحوه ارزیابی
۷۸	۵-۵) بررسی عملکرد
۷۸	۶-۵) مقایسه الگوریتم ارائه شده و تکنیک‌های موجود
۷۹	۷-۵) کارهای پیشنهادی آینده
۸۰	۸-۵) خلاصه و نتیجه‌گیری
۸۱	منابع و مراجع
	واژگان
۸۶	واژگان انگلیسی به فارسی
۹۲	واژگان فارسی به انگلیسی



مقدمه

۱-۱) پردازش تصویر

هر نوع پردازش سیگنال که دارای ورودی و خروجی است، پردازش تصویر^۱ نامیده می‌شود. ورودی، یک تصویر مانند عکس است. خروجی نیز می‌تواند یک تصویر یا یک مجموعه از نشانه‌های ویژه و یا متغیرهای مربوط به تصویر باشد. اغلب تکنیک‌های پردازش تصویر، شامل برخورد با تصویر به عنوان یک سیگنال دو بعدی، و به کار بستن تکنیک‌های استاندارد پردازش سیگنال روی آنها است. پردازش تصاویر امروزه بیشتر به موضوع پردازش تصویر دیجیتال گفته می‌شود که شاخه‌ای از دانش رایانه است که با پردازش سیگنال دیجیتال که نماینده تصاویر برداشته شده با دوربین دیجیتال یا پویش شده توسط پویشگر^۲ است سروکار دارد. یکی از انواع تصاویر ورودی، تصویر سند است.

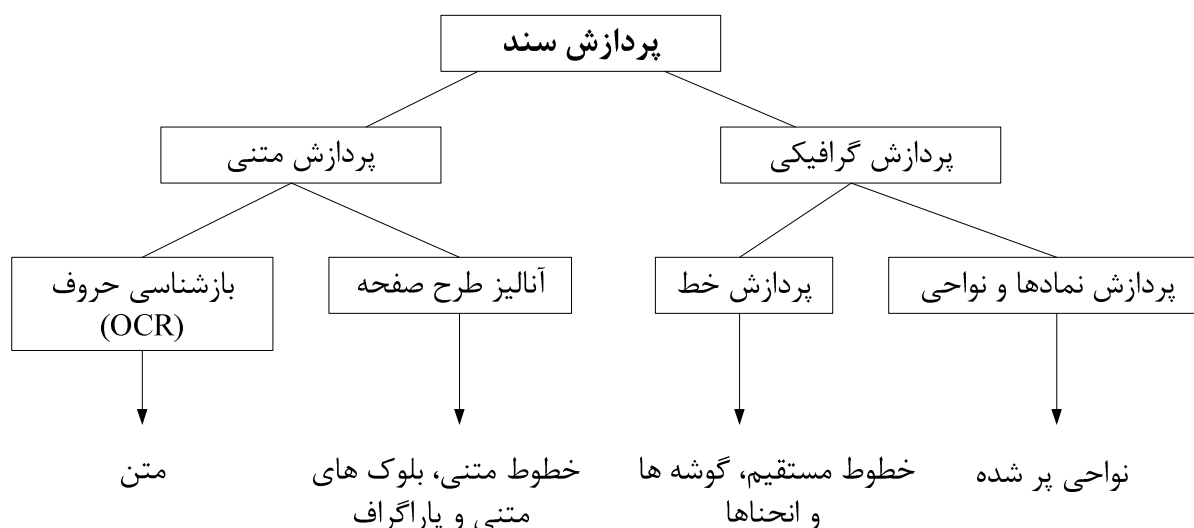
۲-۱) تصویر سند

روش متداول ذخیره‌سازی و ارسال اطلاعات، به وسیله اسناد کاغذی است. اما با توجه به مشکلات و سختی‌های سروکار داشتن با کاغذ، در چند دهه اخیر استفاده از کامپیوتر به منظور ذخیره‌سازی اسناد، با رشد چشم‌گیری همراه بوده است. شروع استفاده از این تکنولوژی از اوایل سال ۱۹۸۰ در آمریکا بود. یکی از مشکلات ذخیره‌سازی سند بصورت دیجیتال، سختی کار در بازیابی^۱ این اسناد است. بنابراین پردازش سند، بعنوان زیرمجموعه‌ای قابل توجه از پردازش تصویر مورد بررسی قرار گرفت. امروزه سیستم‌های پردازش اسناد دارای قابلیت‌های بسیار گسترده

^۱ Image processing

^۲ Scanner

هستند. به طور مثال این سیستم‌ها هم اکنون نوع سند را تشخیص می‌دهند، می‌توانند قسمت‌های اساسی آن را استخراج کنند، و در صورت نیاز آن را از یک فرمت به فرمت دیگر تبدیل کنند. انواع پردازش تصویر سند در شکل ۱-۱ نشان داده شده است [1].



شکل ۱-۱ سلسله بندی پردازش تصویر سند در نواحی مختلف

۳-۱) بازشناسی نوری حروف (OCR)

OCR^۲ بعنوان یک سیستم بازشناسی متن، معرفی می‌شود (بازشناسی متن، تبدیل تصویر به متن قابل ویرایش است). OCR کوتاه‌نوشت بازشناسی نوری نویسه‌ها^۳ یا نویسه‌خوان نوری^۴ است. OCR در ابتدا تنها در مورد بازشناسی ارقام و حروف چاپی بکار گرفته می‌شد. با گذشت زمان و پیشرفت قابل توجه در این زمینه، روش‌های بازشناسی دست‌نوشته و متون چاپی مطرح شدند که دامنه کار را به کلمات و عبارات رساندند. با وجود عدم تطبیق OCR با این موارد، این نام برای این روش‌ها و تا حدی برای بازشناسی دست‌نوشته‌های برخط هم استفاده شد و رواج پیدا کرد. بدین ترتیب، امروزه سیستم OCR برای بازشناسی کلیه مستندات چاپی و مکتوبات

¹ Recognition

² Optical Character Recognition

³ Optical Character Recognition

⁴ Optical Character Reader

تایپی^۱ مثل صفحات کتاب‌ها، مجله‌ها و نامه‌های چاپی و همچنین اسناد دست‌نویس^۲، هم در نوشته‌های برخط^۳ و هم در نویسه‌های برون‌خط^۴ بکار می‌رود.

یک سامانه OCR، مثل یک نفر ماشین‌نویس، یک متن را می‌خواند و آن را به قالب مناسب برای ذخیره در کامپیوتر تبدیل می‌کند. معمولاً یک پویشگر، تصویر متن را برای OCR فراهم می‌کند. این تصویر معمولاً از نقاط سیاه و سفید تشکیل می‌شود. سامانه OCR، اشیاء موجود در این تصویر را که ارقام، حروف، علائم، و کلمات هستند، بازشناسی کرده و نام آن‌ها را در قالب مناسب ذخیره می‌کند. یک فایل تصویری، حجم زیادی دارد و جستجوی متنی در آن ممکن نیست. این در حالی است که فایل خروجی یک سامانه OCR بسیار کم‌حجم و قابل جستجو است [۲]. بنابراین استفاده از سیستم‌های OCR دو مزیت عمده دارد:

۱. افزایش چشمگیر سرعت دسترسی به اطلاعات. زیرا در متن بر خلاف تصویر، امکان جستجو و ویرایش وجود دارد.

۲. کاهش فضای ذخیره سازی. زیرا حجم فایل متنی استخراج شده از یک تصویر، معمولاً بسیار کمتر از حجم خود فایل تصویری است.

۱-۴) معرفی کلی یک سیستم بازشناسی متن

طبق آنچه تاکنون بیان شد برای تبدیل اطلاعات داخل تصاویر متن، به شکل قابل فهم برای ماشین، نیاز به یک سیستم بازشناسی متن (OCR) است. در این بخش، مفاهیم، تعاریف و اصطلاحات بازشناسی متن معرفی می‌شود. اجزای اصلی یک سیستم بازشناسی متن عبارتند از:

جداسازی: جداسازی بلوک‌های تصویر مانند متن، گراف و شکل در این مرحله انجام می‌شود.

استخراج ویژگی: برای معرفی هر یک از اجزای تصویر به شکل قابل فهم برای رایانه، مشخصات هر یک از اجزای تصویر در این مرحله استخراج می‌شود.

¹ Machine-Printed

² Handwritten

³ Online

⁴ Offline

بازشناسی یا طبقه‌بندی متن : بلوک‌های متن استخراج شده در مرحله قبل، در این مرحله بازشناسی می‌شوند. بازشناسی هر یک از اجزای تصویر متن با طبقه‌بندی آن‌ها در یک مجموعه حروف، زیرکلمات یا کلمات معتبر انجام می‌شود.

اجزای مکمل برای یک سیستم بازشناسی متن عبارتند از:

پیش‌پردازش : پردازش‌های لازم برای بهبود اولیه کیفیت تصویر، مانند رفع کجی و حذف نویز در این مرحله انجام می‌شود.

پس‌پردازش : کلمات بازشناسی شده در مرحله بازشناسی، با یک واژه نامه مقایسه و تأیید می‌شود.

بازشناسی فرمول : روش بازشناسی که قادر به تبدیل تصویر فرمول به متن قابل ویرایش باشد. تشخیص متن داخل جدول : قابلیت سیستم بازشناسی در استخراج متون داخل جدول و بازشناسی آن‌ها، و تبدیل این متون با شکل جدول متنی.

تحلیل ساختار فیزیکی : استخراج بلوک‌های داخل سند به نحوی که موقعیت فیزیکی هر یک در سند ثبت شده و پس از بازشناسی محتوای آن‌ها، برای بازتولید سند متنی، قابل استفاده باشد. تحلیل ساختار منطقی : تعیین جنس هر یک از بلوک‌های تصویر مانند عنوان، چکیده، مقدمه و متن اصلی برای بلوک‌های متنی.

غلط‌یاب املایی : قابلیت سیستم بازشناسی در تشخیص غلط‌های املایی.

پس‌پردازش نحوی : قابلیت سیستم بازشناسی در تشخیص غلط‌های نحوی.

بازسازی سند¹ : پس از بازشناسی تمام اجزای تصویر، شکل ظاهری سند با استفاده از اجزای بازشناسی شده، مانند تصویر اصلی بازتولید می‌شود.

تکامل و توسعه سیستم OCR، با سرعت چشمگیری انجام شد، اما علی‌رغم پیشرفت‌های زیادی که در زمینه پردازش متون تایپی و دست‌نویس بدست آمده است، قدرت ماشین‌ها در خواندن متون، هنوز از توانایی بشر فاصله زیادی دارد. با این وجود، خواندن خودکار متون تایپی و دست‌نویس در بسیاری از زبان‌ها به مراحل بلوغ رسیده است. تاکنون کارهای زیادی بر روی متون به زبان‌های لاتین، چینی و ژاپنی انجام شده است، اما برای متون فارسی و عربی و امثال آن‌ها، کارهای انجام

¹ Document Regeneration

شده نسبتاً کم و پراکنده بوده است. ذکر این نکته ضروری است که بیش از یک سوم جمعیت جهان به زبان‌هایی صحبت می‌کنند که برای نوشتن آن‌ها از نویسه‌های عربی یا شبیه به آن استفاده می‌کنند. از زبان‌های غیرعربی می‌توان به فارسی، دری، پشتو، اردو، کردی و جاوی اشاره کرد. با استناد به این مسئله، مشخص است که رشد OCR در زبان‌های عربی و مشابه آن، بسیار کند بوده و افق‌های دستیابی به تکامل آن چندان نزدیک به نظر نمی‌رسد. از جمله علل این امر را می‌توان عدم سرمایه‌گذاری کافی، خصلت پیوسته بودن این خطوط، و فقدان پایگاه داده‌ها و لغت-نامه‌های استاندارد و جامع برای متون فارسی و عربی و امثال آن‌ها ذکر کرد [۲].

۵-۱) جداسازی خطوط

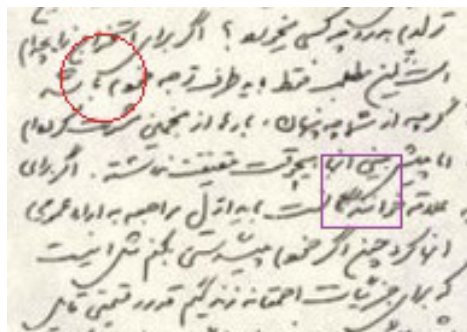
در این بخش با توجه به عنوان و هدف رساله، تمرکز بر روی زمینه‌های محدودتری از سیستم کلی OCR است. جداسازی خطوط متنی، یکی از زیرمجموعه‌های پیچیده و پرکاربرد آنالیز تصاویر اسناد است. این شاخه از آنالیز تصاویر، اطلاعات بسیار قاطع و تعیین‌کننده‌ای، بمنظور تصحیح کجی‌ها و انحنایها، تقسیم‌بندی نواحی و بازشناسی حروف را فراهم می‌کند. بنابراین جداسازی خطوط، مرحله‌ای مهم بمنظور آنالیز برون خط اسناد است، چرا که جداسازی نادرست خطوط، منجر به بروز خطاهایی در مرحله شناسایی خواهد گشت. اگرچه جداسازی خطوط متنی در اسناد تایپی یا دست‌نویس، معمولاً بعنوان یک مسئله حل شده به نظر می‌رسد [3]، اما خطوط متنی اسناد دست‌نویس آزاد^۱، هنوز بعنوان یک مشکل قابل توجه وجود دارد. در جدول ۱-۱، لیستی از چند مقایسه بین جزئیات اسناد تایپی و دست‌نویس فهرست شده است [4].

جدول ۱-۱ مقایسه اسناد دست‌نویس و تایپی

خط متن	فاصله بین خطوط همسایه	طرح (حاشیه)
اسناد تایپی	مستقیم	قابل توجه
اسناد دست‌نویس	منحني	ناچیز
		نامنظم

با توجه به جدول فوق، مرحله جداسازی خطوط برای اسناد تایپی، تا حدود زیادی به تکامل خود دست یافته است، اما در خطوط دست‌نویس موجود در اسناد، که از قالب خاصی تبعیت نمی‌کنند،

همچنان یک مشکل مهم تلقی می‌شود و همچنان بعنوان یکی از مشکلات پیچیده در بهبود OCR موثق، باقی مانده است. علت این امر نیز همان‌طور که در جدول فوق دیده می‌شود؛ طبیعت متون دست‌نویس است که فرآیند جداسازی خطوط را بسیار پرچالش می‌کند، زیرا در اسناد دست‌نویس فاصله داخلی بین خطوط و میزان انحراف خط کرسی متغیر است. مشکلات زمانی حادث می‌شود که دو خط مجاور در متن سند، در محل‌هایی دچار تماس^۲ یا هم‌پوشانی^۳ شده باشند. نمونه آن در شکل ۱-۲ نشان داده شده است.



شکل ۱-۲ هم‌پوشانی (مستطیل) و تماس (دایره) دو کلمه

متون تاییبی دارای ویژگی‌های مناسبی برای جداسازی هرچه بهتر خطوط هستند. از جمله مهمترین این ویژگی‌ها، یکسان بودن فاصله خطوط از یکدیگر در تمام امتداد دو خط مجاور، و همچنین متصل بودن تمام اجزای یک حرف به هم است. خلأ وجود این ویژگی‌ها، در اکثریت قریب به اتفاق متون دست‌نویس، مشکلاتی را نظیر هم‌پوشانی و یا اتصال اجزا در دو خط مجاور ایجاد می‌کند.

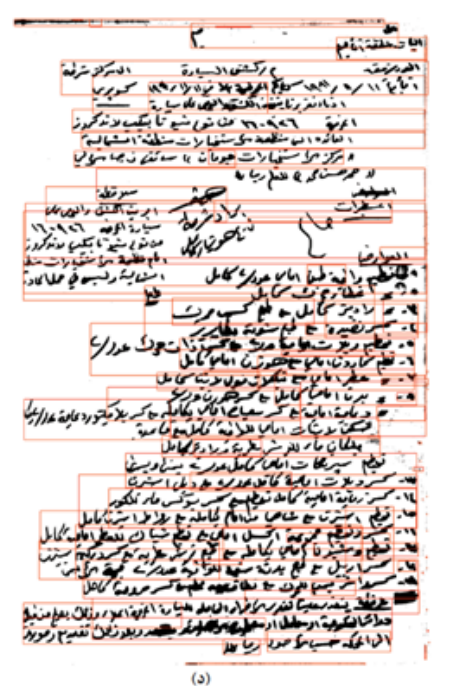
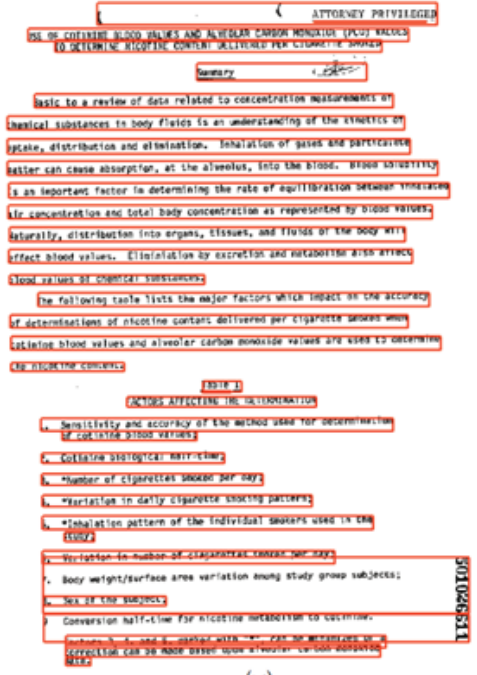
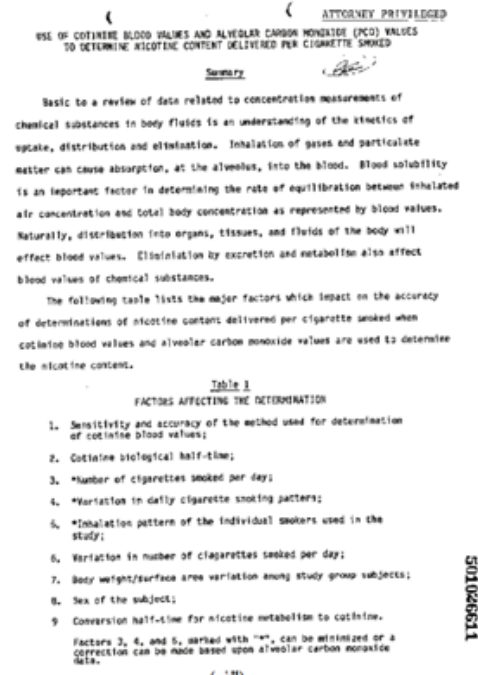
شکل ۱-۳، تأییدی بر توصیفات فوق است. یکی از روش‌هایی که بصورت گسترده در جداسازی مطلوب نواحی، خطوط، کلمات و حروف در متون تاییبی مورد استفاده قرار می‌گیرد، الگوریتم Docstrum (الگوریتم کادرها و مستطیل‌های محیطی^۱ است) [5] و [6]. با توجه به شکل، این روش در جداسازی خطوط موجود در اسناد دست‌نویس، به نتیجه غیرقابل قبول نشان داده شده منجر می‌شود. علت این امر را علاوه بر دست‌نویس بودن اسناد، باید در ویژگی‌های خاص نوشتار زبان فارسی نیز، جستجو کرد که در دست‌نویس شدن این زبان، این ویژگی‌ها، موضوع اصلی را با

¹ Freestyle

² Touch

³ Overlap

چالش‌های زیادی مواجه می‌کنند. این شاخص‌ها، خاص زبان‌های شبه عربی مانند فارسی هستند و در ادامه بطور جامع بیان خواهند شد.



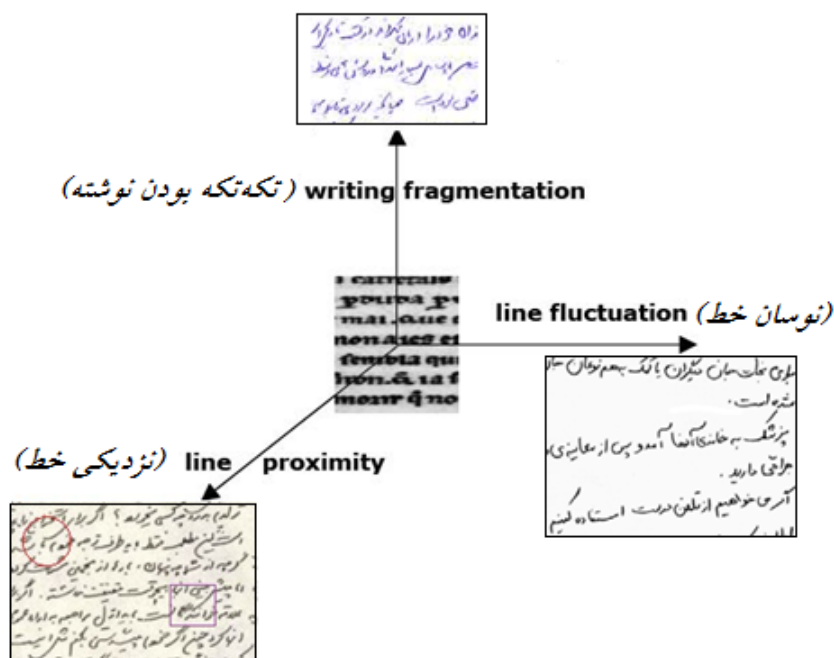
شکل ۱-۳ (الف) تصویر یک سند تایپی. (ب) نتیجه جداسازی خطوط. (ج) تصویر یک سند دست‌نویس. (د) نتیجه جداسازی خطوط.

¹ Bounding Box

۶-۱) چالش‌های جداسازی

مشکلات جداسازی خطوط در متون دست‌نویس را، می‌توان بصورت خلاصه در شکل زیر

مشاهده کرد [7].



شکل ۴-۱ پیچیدگی در اسناد دست‌نویس

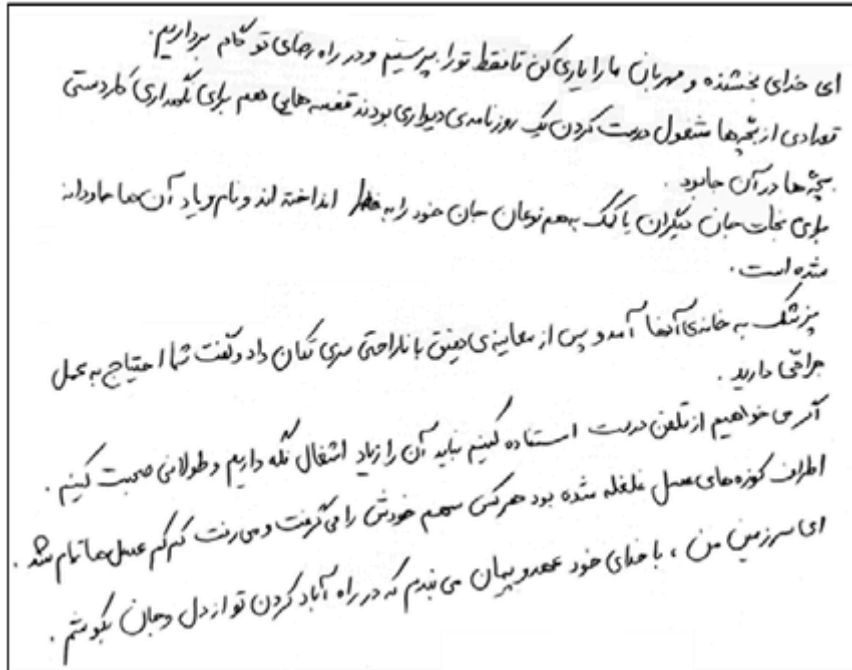
۱-۶-۱) نوسان و انحناهای خطوط

خطوط یک متن دست‌نویس عموماً مستقیم نیست. در حالت کلی سه نوع انحنا در متن وجود دارد:

- انحناهای سراسری: تمام قطعات صفحه، جهت کجی یکسانی دارند.
- انحناهای چندگانه: پاراگراف‌های نامنظم با میزان کجی متفاوت در قطعه‌های مختلف صفحه سند که تنها می‌توان آن‌ها را بصورت غیر خطی نظم بخشید.
- انحناهای خطوط (یا تغییر شیب خط): یک خط دارای چندین شیب متفاوت است، مانند خط منحنی‌وار.

باید دقت کرد که در اکثر موارد، خطوط در اسناد دست‌نویس هر دو نوع کجی دسته دوم و سوم را دارند، بعبارتی علاوه بر اینکه خط منحنی‌وار است و عموماً تخمین خطی و یا تکه‌ای-خطی در

مورد آن صحیح نیست، کجی یک فرمت مشخص ندارد و در خطوط متفاوت، جهت متفاوتی نسبت به خط‌های دیگر پیدا می‌کند. شکل ۱-۵ این واقعیت را نشان داده است.



شکل ۱-۵ یک صفحه دست‌نویس فارسی با کجی‌های متفاوت در خطوط مختلف

۱-۶-۲) نزدیکی و مجاورت خطوط

فاصله کم بین خطوط مجاور، سبب هم‌پوشانی یا تماس قسمتهایی از کلمات می‌شود که حروفشان نسبت به خط، دارای قسمت‌های بالارونده یا پایین‌آینده هستند.

۱-۶-۳) تکه‌تکه شدن دست‌نویسته

برخی حروف از بیش از یک جزء ساخته شده‌اند. حروف نقطه‌دار در زبان‌های فارسی و عربی (نقطه‌های جداکننده)، مثالی از این حروف هستند.

مشکلات اشاره شده، مشکلاتی هستند که در تمام اسناد دست‌نویس و مستقل از نوع زبان نوشتاری‌شان قابل مشاهده‌اند.

۷-۱) ویژگی‌های خط فارسی

پیچیدگی‌های موجود در اسناد دست‌نویس فارسی، حادث‌تر و بارزتر خواهد بود. علت این امر، خصوصیات خاص نگارشی این زبان است. در این بخش برخی ویژگی‌ها و شاخص‌های نوشتاری فارسی بطور مختصر بیان می‌شود. شاید بتوان مهمترین ویژگی نوشتاری زبان فارسی را پیوستگی ذاتی و جهت نگارش آن از راست به چپ دانست که این زبان را از سایر زبان‌ها متمایز کرده است. در ادامه تعدادی از ویژگی‌های مهم نوشتاری خط فارسی بطور خلاصه ذکر می‌گردد [۲].

الف) اشکال متفاوت یک حرف

شکل حروف در زبان فارسی تابعی از محل قرار گرفتن آن‌ها در کلمه است و هر حرف بر حسب موقعیت‌های مختلف در کلمه (اول، آخر و جدا) می‌تواند اشکال مختلفی داشته باشد (جدول ۱-۲). بعنوان مثال، شکل حرف (م) در حالت جدا، باید به شکل (م) و در ابتدا یا وسط کلمه، به شکل (م)، نوشته شود. شکل این حرف در انتهای کلمه و در یک متن دست‌نویس ممکن است موجب هم‌پوشانی و یا تماس بین دو خط متوالی گردد، که عمل جداسازی دو خط را با مشکل مواجه می‌کند. بعنوان مثال تداخل حرف (م) در حالت دست‌نویس با سرکج حروف (ک) یا (گ) و یا مد (آ) در خط متوالی بعدی رخ می‌دهد.

ب) نقطه‌دار بودن حروف

بیش از نیمی از حروف فارسی نقطه‌دار هستند. بعبارت دقیق‌تر، ۱۰ حرف دارای یک نقطه، ۳ حرف دارای دو نقطه و ۵ حرف دارای سه نقطه هستند. بحث مهم نقطه‌دار بودن کلمات نیست بلکه محل قرار گرفتن نقاط است که مشکل‌ساز می‌شود. مثلاً با توجه به دو حرف {ب، ن} یا {چ، ژ} می‌توان به وجود این مشکل در متون دست‌نویس فارسی پی برد که قرار گرفتن حداقل یک نقطه بین دو خط متوالی، امری کاملاً طبیعی است ولی تشخیص اینکه نقطه مربوط به خط بالایی است یا پایینی، تصمیمی است که روند جداسازی خطوط را سخت می‌کند.

ج) وجود اعراب

در نوشتار فارسی، هنگامی که احتمال تلفظ اشتباه یک کلمه وجود داشته باشد، صدا بصورت اعراب به بالا یا پایین برخی از حروف اضافه می‌شود. قرار گرفتن اعراب بین خطوط نیز، احتمال بروز مشکلی نظیر وجود نقطه‌ها را که در فوق توضیح داده شد، بوجود می‌آورد. علاوه بر اعراب، در

نوشتار فارسی علائمی از قبیل تشدید، تنوین، همزه و مد (در صورت نوشتاری حرف آ) نیز وجود دارند که جدا کردن خطوط را در متون دست‌نویس مشکل می‌سازند. علت بروز این مشکل علاوه بر سخت بودن تشخیص تعلق علائم به کدام خط (خط بالا یا خط پایین)، می‌تواند ناشی از تماس حروف موجود روی خطوط بالا یا پایین با این علائم باشد.

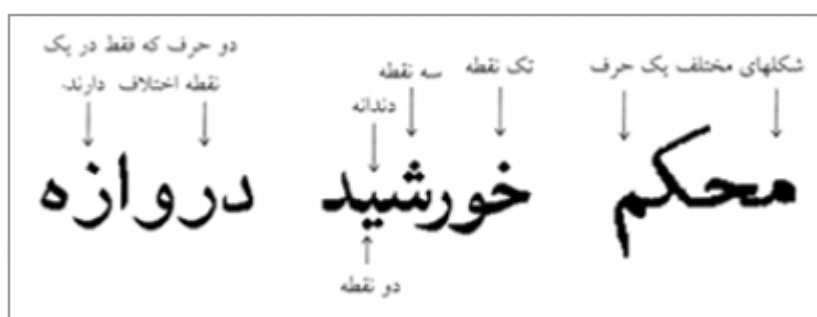
جدول ۱-۲ شکل‌های مختلف حروف فارسی با توجه به محل قرار گرفتشان در زیرکلمه

مجزا	ابتدا	میان	انتهای		مجزا	ابتدا	میان	انتهای	
ا یا آ	ا یا آ	---	ا	۱۷	ص	ص	ص	ص	۱
ب	ب	ب	ب	۱۸	ض	ض	ض	ض	۲
پ	پ	پ	پ	۱۹	ط	ط	ط	ط	۳
ت	ت	ت	ت	۲۰	ظ	ظ	ظ	ظ	۴
ث	ث	ث	ث	۲۱	ع	ع	ع	ع	۵
ج	ج	ج	ج	۲۲	غ	غ	غ	غ	۶
چ	چ	چ	چ	۲۳	ف	ف	ف	ف	۷
ح	ح	ح	ح	۲۴	ق	ق	ق	ق	۸
خ	خ	خ	خ	۲۵	ک	ک	ک	ک	۹
د	---	---	د	۲۶	گ	گ	گ	گ	۱۰
مجزا	ابتدا	میان	انتهای		مجزا	ابتدا	میان	انتهای	
ذ	---	---	ذ	۲۷	ل	ل	ل	ل	۱۱
ر	---	---	ر	۲۸	م	م	م	م	۱۲
ز	---	---	ز	۲۹	ن	ن	ن	ن	۱۳
ژ	---	---	ژ	۳۰	و	---	---	و	۱۴
س	س	س	س	۳۱	ه	ه	ه	ه	۱۵
ش	ش	ش	ش	۳۲	ی	ی	ی	ی	۱۶

د) تنوع فراوان در شیوه‌های نگارشی

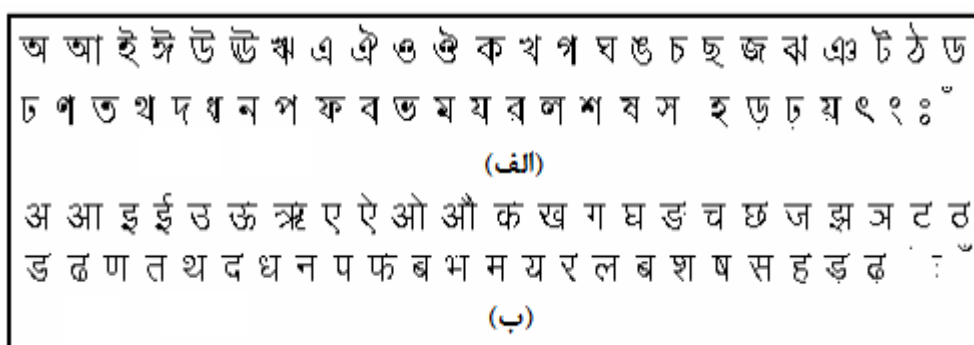
در مقایسه با نوشتار لاتین، فونت‌های فارسی در حالت چاپی دارای واحدهای نوشتاری بسیار زیاد (شامل ارقام، اشکال مختلف حروف، علائم خاص، اجزاء کوچک دارای اهمیت مانند دندانه) و در حالت دست‌نویس دارای سبک‌های متعدد می‌باشد. امری کاملاً واضح است که تعدد سبک نوشتاری اسناد دست‌نویس، فرآیند جداسازی را دشوارتر می‌سازد.

شکل زیر نیز نمونه‌هایی از ویژگی‌های خط نوشتاری فارسی را بخوبی بیان می‌کند.



شکل ۱-۶ ویژگی‌های خط نوشتاری فارسی

علاوه بر آنچه تاکنون ذکر شد، می‌توان به این نکته توجه کرد که در برخی زبان‌ها، وجود ویژگی‌های خاص برای حروف الفبای آن زبان، جداسازی خطوط را بسیار آسان می‌کند. مثلاً شکل ۱-۷ دو نمونه از این زبان‌ها را نشان می‌دهد که پیدا کردن خط افقی موجود در قسمت فوقانی حروف، منجر به اجرای الگوریتم موفق‌تری بمنظور جداسازی خطوط، خواهد بود اما با توجه به تنوع اشکال حروف در زبان فارسی، پیدا کردن چنین نشانه‌ای که تمام حروف را پوشش دهد امری سخت و بدون اغراق تقریباً ناممکن است.



شکل ۱-۷ وجود یک خط افقی در اکثریت حروف (الف) bangla و (ب) devnagari (از زبان‌های هندی)