

۱۰۵ع

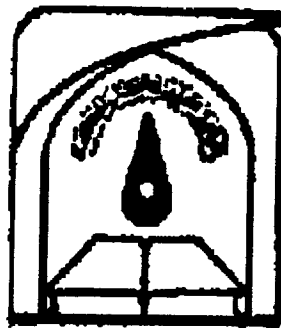
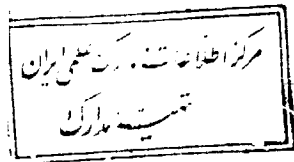
۶- دارد- ۳۵

بِسْمِ اللّٰهِ

الرّحمن

الرّحيم

۱۳۷۹ / ۹ / ۲۰



دانشگاه تربیت مدرس
دانشکده فنی و مهندسی

پایان نامه کارشناسی ارشد مهندسی برق
(الکترونیک)

پیش پردازش متون چاپی فارسی برای جداسازی حروف

حسین نظام آبادی پور

استاد راهنما

دکتر احسان‌اله کبیر

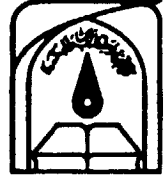
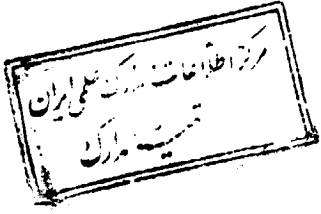
استاد مشاور

دکتر رضا عزمی

تابستان ۱۳۷۹

۸۶۲۸

۳۱۰۹۴



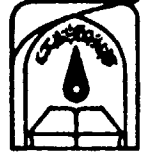
دانشگاه تربیت مدرس

تاییدیه هیات داوران

آقای حسین نظام آبادی پور پایان نامه ۶ واحدی خود را با عنوان پیش پردازش متون چاپی فارسی برای جداسازی حروف در تاریخ ۱۹/۶/۷۹ ارائه کردند. اعضای هیات داوران نسخه نهائی این پایان نامه را از نظر فرم و محتوی تایید و پذیرش آنرا برای تکمیل درجه کارشناسی ارشد رشته مهندسی برق باگرایش الکترونیک پیشنهاد می کنند. ۶۹ ب ۶

امضاء	نام و نام خانوادگی	اعضای هیات داوران
	آقای دکتر کبیری	۱- استاد راهنما:
	آقای دکتر عزمی	۲- استاد مشاور:
	آقای دکتر لطفی زاد	۳- استادان ممتحن:
	آقای دکتر صفا بخش	
	آقای دکتر نبوی	۴- مدیر گروه: (یا نماینده گروه تخصصی)

این نسخه به عنوان نسخه نهایی پایان نامه / رساله مورد تایید است.
امضای استاد راهنما:



بسمه تعالی

آیین نامه چاپ پایان نامه (رساله) های دانشجویان دانشگاه تربیت مدرس

نظر به اینکه چاپ و انتشار پایان نامه (رساله) های تحصیلی دانشجویان دانشگاه تربیت مدرس، مبین بخشی از فعالیت های علمی - پژوهشی دانشگاه است بنابراین به منظور آگاهی و رعایت حقوق دانشگاه، دانش آموختگان این دانشگاه نسبت به رعایت موارد ذیل متعهد می شوند:

ماده ۱ در صورت اقدام به چاپ پایان نامه (رساله) ی خود، مراتب را قبلاً به طور کتبی به «دفتر نشر آثار علمی» دانشگاه اطلاع دهد.

ماده ۲ در صفحه سوم کتاب (پس از برگ شناسنامه)، عبارت ذیل را چاپ کند:
و کتاب حاضر، حاصل پایان نامه کارشناسی ارشد / رساله دکتری نگارنده در رشته ررر - الرررررر است که در سال ۱۳۷۹ در دانشکده فنون و مهندسی دانشگاه تربیت مدرس به راهنمایی سرکار خانم / جناب آقای دکتر احسان گلگیر ، مشاوره سرکار خانم / جناب آقای دکتر رضا غریبی و مشاوره سرکار خانم / جناب آقای دکتر ... از آن دفاع شده است.

ماده ۳ به منظور جبران بخشی از هزینه های انتشارات دانشگاه، تعداد یک درصد شمارگان کتاب (در هر نوبت چاپ) را به «دفتر نشر آثار علمی» دانشگاه اهدا کند. دانشگاه می تواند مازاد نیاز خود را به نفع مرکز نشر در معرض فروش قرار دهد.

ماده ۴ در صورت عدم رعایت ماده ۳، ۵۰٪ بهای شمارگان چاپ شده را به عنوان خسارت به دانشگاه تربیت مدرس، تأدیه کند.

ماده ۵ دانشجو تعهد و قبول می کند در صورت خودداری از پرداخت بهای خسارت، دانشگاه می تواند خسارت مذکور را از طریق مراجع قضایی مطالبه و وصول کند؛ به علاوه به دانشگاه حق می دهد به منظور استیفای حقوق خود، از طریق دادگاه، معادل وجه مذکور در ماده ۴ را از محل توقیف کتابهای عرضه شده نگارنده برای فروش، تأمین نماید.

ماده ۶ این جانب حسین نجات آبادی دانشجوی رشته برق - الکترنیس مقطع کارشناسی ارشد تعهد فوق و ضمانت اجرایی آن را قبول کرده، به آن ملتزم می شوم.

نام و نام خانوادگی: حسین نجات آبادی
تاریخ و امضاء:

تقدیم به:

پدر زحمتکش و مادر دلسوزم که موفقیت خود در تمامی مراحل زندگی به خصوص در امر تحصیل را مدیون آن ها میدانم؛ همچنین به اساتید دوران تحصیل به ویژه جناب آقای سید حسن جهان بین که رنج فراوانی را در راه تعلیم و تربیت اینجانب متحمل شده اند.

تشکر و قدردانی

لازم است از جناب آقای دکتر احسان اله کبیر که در طول انجام این تحقیق از راهنماییهای مفید و ارزنده ایشان استفاده فراوان برده ام، صمیمانه تشکر و قدردانی کنم. همچنین از جناب آقای دکتر رضا عزمی به خاطر رهنمودهای خردمندانه شان در طول انجام تحقیق تشکر و سپاسگزاری می کنم.

چکیده

پردازش مستندات یکی از جذاب ترین زمینه های بازشناسی الگو است و بازشناسی متون، محوری ترین بخش در پردازش مستندات است. یکی از مهمترین مراحل بازشناسی متون چاپی فارسی جداسازی حروف است.

ما در این پایان نامه با اصلاح الگوریتم عزمی که مبتنی بر کانتور بالایی است، الگوریتم جداسازی مناسبی برای متون چاپی قدیمی ارائه کرده ایم. برای حل مشکل نایکنواختی کرسی خط، روش مناسبی برای تعیین نوار زمینه پیشنهاد کرده ایم و با اصلاح روش بر چسب زنی کانتور بالایی و تکمیل قواعد جداسازی، دقت الگوریتم را برای متون قدیمی افزایش داده ایم. نرخ جداسازی درست حروف با استفاده از این الگوریتم برای یک مجموعه آزمایش که از متون قدیمی برگزیده شده است، بدون هیچ پس پردازشی ۹۷٪ است. برای جداسازی نقاط چسبیده به بدنه حروف نیز روشی ارائه شده است که نتایج رضایت بخشی داده است.

کلمات کلیدی: جداسازی حروف، متون چاپی فارسی، نوار زمینه محلی، کانتور بالایی،

هیستوگرام، پروفایل، جداسازی نقاط

فهرست

صفحه	عنوان
۱	فصل اول: مقدمه
۷	فصل دوم: جداسازی حروف در متون چاپی
۷	۱-۲) جداسازی حروف در متون چاپی فارسی و عربی
۸	۱-۱-۲) جداسازی بر اساس هیستوگرام عمودی
۱۱	۲-۱-۲) جداسازی بر اساس منحنی پیرامونی
۱۲	۳-۱-۲) جداسازی با استفاده از پروفایل بالایی
۱۵	۴-۱-۲) جداسازی بر اساس قواعد ساختاری نگارش فارسی
۱۹	۵-۱-۲) جداسازی با استفاده از کانتور بالایی کلمات
۲۳	۲-۲) جداسازی چسبیدگی های ناخواسته حروف در متون تایپی لاتین
۲۵	۱-۲-۲) جداسازی حروف با استفاده از شبکه عصبی
۲۸	۲-۲-۲) جداسازی حروف با استفاده از تکنیک پوسته محدب
۳۱	۳-۲-۲) جداسازی با استفاده از هیستوگرام عمودی نقاط سیاه و نمای پروفایل
۳۸	فصل سوم: الگوریتم جداسازی حروف و نقاط چسبیده به بدنه
۳۹	۱-۳) کلیات الگوریتم جداسازی
۴۲	۲-۳) تصویر برداری از متن
۴۲	۳-۳) جداسازی خطوط تصویر متن ورودی
۴۳	۴-۳) محاسبه پهنای قلم
۴۳	۵-۳) تعیین نوار زمینه
۴۵	۶-۳) جدا کردن بخش های همپوشان
۴۵	۷-۳) محاسبه پهنای قلم بصورت محلی برای هر بخش همپوشان

صفحه	عنوان
۴۷	۸-۳) بدست آوردن منحنی پیرامونی زیر کلمه
۴۸	۹-۳) تعیین نوار زمینه محلی
۴۹	۱-۹-۳) تعیین نوار زمینه محلی با استفاده از کانتور زیر کلمات
۵۲	۲-۹-۳) تعیین نوار زمینه محلی با استفاده از هیستوگرام های $h_{T_i}(n)$ و $h_T(n)$
۵۳	۳-۹-۳) تعیین نوار زمینه محلی با استفاده از ماکزیمم های اول و دوم هیستوگرام های $h_{T_i}(n)$ و $h_T(n)$ و موقعیت نوار زمینه کلی
۵۶	۴-۹-۳) تصحیح نوار زمینه در زیر کلمات شامل دندان‌ها زیاد
۵۷	۱۰-۳) بر چسب زدن نقاط کانتور بالایی
۵۹	۱۱-۳) برچسب زدن پاره مسیرها و آشکارسازی نقاط جداسازی
۶۳	۱۲-۳) معیارهایی برای محک زدن نوار زمینه محلی
۶۳	۱-۱۲-۳) سطح مشترک میان نوارهای زمینه محلی و نوار زمینه کلی
۶۴	۲-۱۲-۳) میانگین فواصل نوارهای زمینه محلی از نوار زمینه کلی
۶۵	۳-۱۲-۳) انتخاب نوار زمینه محلی مرجع
۶۶	۱۳-۳) جداسازی نقاط چسبیده به بدنه
۶۶	۱-۱۳-۳) جداسازی نقاط پایینی
۶۹	۲-۱۳-۳) جداسازی نقاط بالایی
۷۳	فصل چهارم: نتایج آزمایشها و بررسی علل خطاها
۷۳	۱-۴) معرفی مجموعه های تمرین و آزمایش
۷۶	۲-۴) معرفی خطاها و بررسی علل آنها
۹۳	۳-۴) نتایج آزمایشها
۹۵	۴-۴) اثرات افزایش درجه تفکیک

صفحه	عنوان
۹۷	فصل پنجم: نتیجه گیری و پیشنهادها
۱۰۰	مراجع
۱۰۳	واژه نامه فارسی - انگلیسی
۱۰۵	واژه نامه انگلیسی - فارسی

فصل اول

مقدمه

اتوماتیک کردن کارها و سپردن فعالیت های انسانی به ماشین از جمله مباحثی است که در دنیای امروز توجه دانشمندان و محققان را به خود جلب کرده است. سیستم های قدیمی دیگر قادر به پاسخگویی به احتیاجات انسان ها و سرویس دهی به موقع و منظم نیستند. از دیگر عواملی که خودکار کردن سیستم ها را شتاب می دهد، مقرون به صرفه کردن آنها و جایگزینی ماشین خستگی ناپذیر به جای انسان است. در چنین شرایطی است که انسان خود را ملزم به بکار گیری ماشین در پردازش اطلاعات می داند. از جمله قسمتهای مهم و پیچیده یک سیستم پردازش اطلاعات، بخش بازشناسی الگو است.

پردازش مستندات از جمله پرکاربردترین امور در زمینه بازشناسی الگو است. در این زمینه آنچه که قبل از هر چیز باید به آن پرداخته شود، درک هندسی ماشین از مستند تصویر برداری شده است. در این مرحله پس از پیش پردازش هایی مثل رفع نویز و برطرف کردن کجی تصویر، محل گراف ها، جدول ها، تصاویر و متون مشخص می شود. بازشناسی متون محوری ترین بخش در پردازش مستندات است. این بخش، اطلاعات تصویری متن را به اطلاعات نمادین حروف و کلمات تبدیل می کند.

از زمینه های کاربردی بازشناسی حروف می توان به خواندن اتوماتیک چک های بانکی، خواندن اتوماتیک پرسشنامه ها، تبدیل متون معمولی به متون مخصوص نابینایان، خواندن کدپستی و آدرس نامه ها و طبقه بندی اتوماتیک آنها اشاره کرد. شرکتها، ادارات بیمه، قضات، مهندسان

گرافیک ، کارمندان دولت، روزنامه نگاران و تمامی کسانی که میل دارند اطلاعات را بدون آنکه لازم باشد آنرا تایپ کنند در کامپیوتر خود داشته باشند، از بازشناسی اتوماتیک متون کمک می گیرند. اولین نرم افزار بازشناسی حروف توسط شرکت ماشین های هوشمند¹ در سال ۱۹۵۹ ساخته شد [WWW.OCR.com]

تحقیقات درباره بازشناسی حروف با استفاده از تصویر برداری نوری را میتوان به دو دسته مهم تقسیم کرد:

بازشناسی حروف تایپ شده بازشناسی حروف دستنویس

در مورد بازشناسی حروف انگلیسی ، ژاپنی و چینی تحقیقات قابل ملاحظه ای صورت گرفته است [Mor,92][Bok,92]. الگوریتم های بکار رفته در بازشناسی متون تایپ شده و دستنویس برای این زبان ها دارای دقت بالایی هستند. در هر زبان روش بازشناسی متناسب با ویژگیهای خاص الفبای آن زبان است. به دلیل قواعد خاص نگارش متون فارسی، هیچیک از روش های موجود در بازشناسی متون به زبان های دیگر ، به طور مستقیم برای الفبای فارسی و عربی قابل استفاده نیست. شروع تحقیقات درباره بازشناسی حروف فارسی و عربی به سال ۱۹۸۰ بر می گردد [par,81]. تحقیقات در این زمینه اگر چه دیر آغاز شده اما خوشبختانه در سالهای اخیر توجه شایانی به بازشناسی متون فارسی در دانشگاههای ایران شده است و نتایج الگوریتم های حاصل از این تحقیقات امیدوار کننده است.

بازشناسی متون چاپی فارسی را به سه روش کلی می توان انجام داد.

۱- شکستن کلمات به حروف و بازشناسی آنها (رویکرد بر اساس جداسازی)

۲- بازشناسی زیر کلمات به عنوان الگوهای پایه (رویکرد بدون جداسازی)

۳- رویکرد ترکیبی

¹ Intelligent machine

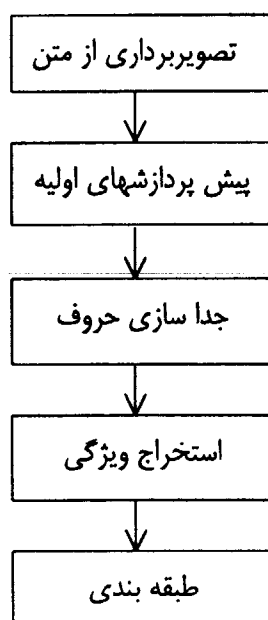
در روش اول، از حروف به عنوان الگوی بازشناسی استفاده می شود. بنابراین یک مرحله اساسی در این روش ، جداسازی حروف است. منظور از جداسازی حروف پیدا کردن چار چوب هر حرف به صورت مجزا است. در این روش ابتدا زیر کلمات به حروف شکسته می شوند و سپس حروف جدا شده بازشناسی می شوند.

در روش دوم سعی بر این است که هر زیر کلمه بصورت یک واحد در نظر گرفته شود و بدون نیاز به جداسازی حروف آن شناسایی شود. به علت تعداد بیش از حد زیر کلمات فارسی، این روش فقط در بازشناسی حجم محدود کلمات کاربرد دارد [رفیعی ، ۷۴].

در روش سوم اطلاعات مربوط به جداسازی حروف کلمه و اطلاعات مربوط به شکل کلمه در قالب یک سیستم ترکیبی جدا سازی _ بازشناسی بکار گرفته می شود.

روش دوم که مبتنی بر بازشناسی بدون جداسازی است به علت حجم زیاد زیر کلمات فارسی از نظر حافظه و سرعت مقرون به صرفه نیست. روش اول و سوم در بازشناسی متون فارسی بیشترین کاربرد و استفاده را دارند . در این روش ها مرحله جداسازی حروف یکی از اساسی ترین مراحل است. به گونه ای که نتایج نهایی به شدت وابسته به نتایج این مرحله است.

مراحل مختلف یک سیستم بازشناسی حروف در بلوک دیاگرام شکل ۱-۱ آمده است.



شکل ۱-۱ بلوک دیاگرام یک سیستم بازشناسی حروف

ابتدا از متن ورودی، تصویری با درجه تفکیک مناسب (حدود 300 dpi) با سطوح خاکستری گرفته می شود. در مرحله بعد، عمل دو سطحی کردن تصویر با استفاده از یک آستانه کلی یا آستانه های محلی انجام می شود. این مرحله از اهمیت ویژه ای برخوردار است زیرا خطاهای این مرحله موجب از بین رفتن اطلاعاتی می شود که در مراحل بعد، سیستم های جداساز حروف و بازشناسی را به اشتباه می اندازد. در مرحله جداسازی، ابتدا باید خطوط متن از یکدیگر جدا و نوار زمینه تعیین شود. بعد از آن در هر خط متن تصاویر زیر کلمات از یکدیگر جدا می شوند و الگوریتم جدا سازی هر زیر کلمه را به حروف آن می شکنند. سپس ویژگیهای مناسب استخراج می شوند و در مرحله طبقه بندی نوع حروف مشخص می شود. آنچه گفته شد، شرح کوتاهی بر یک سیستم بازشناسی حروف چاپی بود. اما در عمل سیستم های بازشناسی متون فارسی با مشکلات فراوانی روبرو می شوند.

از جمله مشکلات مشترک میان بازشناسی حروف فارسی و سایر زبان ها میتوان وجود نویز و بریدگی حروف را نام برد. چسبیدن نقاط به بدنه و یا چسبیدن حروف در غیر نقاط متعارف به یکدیگر از دیگر عواملی است که بازشناسی و جداسازی را با مشکل روبرو می کند. جداسازی حروف در این متون بدون پشتیبانی الگوریتم های بازشناسی مشکل می شود. باید در این موارد بازشناسی و جداسازی با یکدیگر توأم باشند [Lu,95,96].

حال مشکلاتی را مطرح می کنیم که به قواعد و ساختار زبان و نگارش فارسی در ارتباط است.

الف) اندازه هندسی حروف فارسی مختلف است.

ب) اندازه و شکل هر حرف در قلم های مختلف، متفاوت است.

ج) حروف با توجه به موقعیت آن ها در کلمه شکل های مختلفی به خود می گیرند. برخی حروف مثل "ع" چهار شکل مختلف به خود می گیرند.

د) در بعضی از قلم های حروف ممکن است با یکدیگر همپوشانی عمودی داشته باشند که جداسازی آنها مشکل می شود. مثل کلمات "محمد، اسلام"

ذ) بعضی از حروف در دو محل به یکدیگر می چسبند، مثل "کا".

ر) حروف مجزای فارسی دارای ارتفاع یکسانی نیستند، بنابراین تشخیص نوع حروف مشکل می شود.

و) بعضی از حروف بدنه مشابه دارند و تنها در تعداد نقاط متفاوت هستند. این موضوع شناسایی را مشکل می کند.

همانگونه که قبل از این نیز اشاره کردیم، جداسازی حروف در میزان بازشناسی نقش تعیین کننده ای دارد. اگر چه الگوریتم هایی که تا کنون ارائه شده اند از توانایی بالایی برخوردارند، اما هر کدام به نحوی خاص دارای نقاط ضعف هستند. یکی از معایب الگوریتم های جداسازی که قبل از این ارائه شده است عدم توجه به جابجایی کرسی خط می باشد. این موضوع باعث شده است تا در متون قدیمی تر که به دلیل ابتدایی بودن ماشین های تایپ کلمات بروی یک خط زمینه ثابت نوشته نمی شدند، عمل جداسازی و به تبع آن بازشناسی با نرخ مناسبی انجام نشود.

هدف پایان نامه

مهمترین موضوعی که ما در این پایان نامه دنبال می کنیم، مسئله جابجایی کرسی نوار زمینه است. سعی کرده ایم تا این موضوع را به نحو شایسته و ممکن قبل از مرحله بازشناسی یعنی در مراحل پیش پردازش و جداسازی حل کنیم. توجه زیادی به نوار زمینه محلی کرده ایم و چگونگی تعیین این نوار به نحو شایسته ای بررسی شده است. در ضمن تا حدی نیز به تصحیح و بهبود کیفیت الگوریتم جداسازی خاصی که مبتنی به نوار زمینه بوده است پرداخته ایم [عزمی، ۷۸].

خوشبختانه نتایج الگوریتم نهایی امیدوار کننده است و نشان از پیشرفت الگوریتم جداساز دارد. این موضوع امکان استفاده از روش اول در بازشناسی متون چاپی فارسی که مبتنی به جداسازی است را افزایش می دهد.