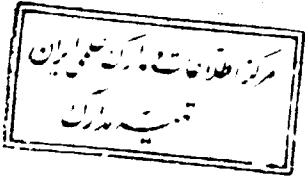


۲۹۹۶۷



دانشگاه علم و صنعت ایران

دانشکده مهندسی کامپیوتر

بازشناسی حروف فارسی با استفاده از شبکه عصبی
چند جمله ایهای جداکننده

سعید اسدی

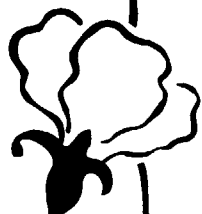
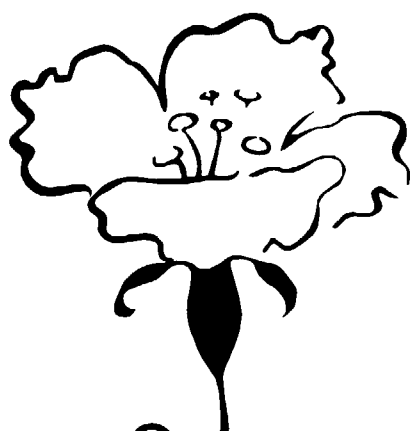
پایان نامه برای دریافت درجه کارشناسی ارشد

در رشته

۱۴۲۱۴

مهندسی کامپیوتر

استاد راهنما: دکتر محمود فتحی



تقدیم به مادرم

عزیزی که در نبود پدر وظیفه
پدری و مادری را بنحو احسن
جامه عمل پوشانید و با کوله بار
سنگین مرارت و زحمت رنج به
ثمر رسانیدن مرا به تنهایی بدوش
کشید.

و تقدیم به همسر عزیزم

چکیده

در این پایان نامه یک روش بازشناسی حروف تاییی فارسی بصورت ساختاری و باپارامتری فرض کردن حروف ارائه شده است. در این روش ویژگیهای مناسبی که بتواند شکل و حالات هریک از حروف را نشان دهد، استخراج میگردد. این ویژگیها بگونه ای مورد استفاده قرار میگیرند که بتوانند نمایانگر حروف مختلف بصورت یکتا باشند. با استفاده از جداول و درخت های تصمیم گیری، روشی جهت بازشناسی حروف تاییی فارسی ارائه میگردد. مسائل مهم در این پایان نامه عبارتند از: جداسازی حروف که براساس موقعیت خط زمینه صورت می پذیرد، استخراج ویژگیهای مناسب وبازشناسی حروف فارسی که با استفاده از بانک های اطلاعاتی ودرخت های تصمیم گیری انجام میگردد. روش دیگر بازشناسی که در این پایان نامه مورد استفاده قرار گرفته است، استفاده از (*Padaline's*) شبکه عصبی چند جمله ایهای جداکننده است. در این روش تعدادی از ویژگیهای استخراج شده تبدیل به ثابت های برای چند جمله ایهای مختلف میگرددو برای حالات مختلف حروف یک چند جمله ای ارائه میگردد. الگوریتم ها مستقل از اندازه حروف میباشند. روشهای ارائه شده در این پایان نامه بدون هیچگونه تغییری میتواند برای بازشناسی حروف تاییی لاتین نیز بکار رود. برای طبقه بندی نهایی در هر دو روش بازشناسی نیز از بانک های اطلاعاتی استفاده شده است.

تقدیر و تشکر

از زحمات استاد ارجمند جناب آقای دکتر محمود فتحی که راهنمای اینجانب در انجام پروژه و همچنین تهیه و ارائه این پایان نامه بوده اند.

و همچنین از هیئت محترم ممتحن آقایان دکتر کبیر و دکتر رحمانی و کلیه اساتید محترم دانشکده کامپیوتر.

تشکر و امتنان خود را ابراز می نمایم

فهرست مطالب

<u>صفحه</u>	<u>عنوان مطالب</u>	<u>فصل</u>
۱	مقدمه	فصل اول:
۶	جداسازی	فصل دوم:
۶	مقدمه	
۸	۲-۱- روش های جداسازی حروف	
۱۴	۲-۲- روش پیشنهادی	
۱۹	خلاصه	
۲۰	شناسایی و استخراج ویژگی ها	فصل سوم:
۲۰	مقدمه	
۲۱	۳-۱- روش عددی - نماها	
۲۴	۳-۲- روش ساختاری	
۲۶	۳-۳- روشهای دیگر	
۲۸	۳-۴- روش پیشنهادی	
۳۷	خلاصه	
۳۸	بازشناسی حروف	فصل چهارم:
۳۸	مقدمه	
۳۹	۴-۱- بازشناسی ساختاری	
۴۶	۴-۲- بازشناسی با پارامتری فرض کردن حروف	
۴۷	۴-۲-۱- شبکه عصبی چندجمله ایهای جداکننده <i>padaline</i>	
۴۹	۴-۲-۲- بازشناسی با استفاده از <i>padaline</i>	
۵۷	خلاصه	
۵۸	نتایج	فصل پنجم:
۶۱	جداول	ضمیمه الف:
۶۶	راهنمای نرم افزار بازشناسی	ضمیمه ب:
۷۱	الگوریتم ها و لیست توابع مهم	
۹۶		منابع و مراجع

فهرست شکلها

<u>صفحه</u>	<u>عنوان شکلها</u>	<u>فصل</u>
		فصل دوم:
۱۰	۲-۱- حاصل روش مجموع فواصل در جداسازی	
۱۱	۲-۲- حاصل نیمرخهای یک کلمه	
۱۱	۲-۳- نمای ضخامت و حاصل جداسازی در این روش	
۱۲	۲-۴- نمای قائم $P1$	
۱۲	۲-۵- نمای قائم $P2$	
۱۳	۲-۶- جداسازی با استفاده از نمای ترکیبی	
۱۴	۲-۷- حاصل نمای عمودی	
۱۵	۲-۸- حاصل بدست آوردن خط زمینه با استفاده از نمای عمودی	
۱۵	۲-۹- حاصل نمای افقی	
۱۶	۲-۱۰- حاصل جداسازی حروف و کلمه های جدا از هم	
۱۶	۲-۱۱- حاصل علامت گذاری خط زمینه	
۱۶	۲-۱۲- حاصل علامت گذاری فقط تکرار خط زمینه	
۱۸	۲-۱۳- حاصل پردازش جداسازی	
۱۸	۲-۱۴- حاصل جداسازی	
		فصل سوم:
۲۱	۳-۱- نمای افقی و عمودی عدد ۴	
۲۲	۳-۲- نماهای افقی ، عمودی و مایل سه تصویر	
۲۲	۳-۳- نمای مایل (۴۵- درجه)	
۳۵	۳-۴- نمادهای حرف "ض"	
		فصل چهارم:
۳۹	۴-۱- حروف بزرگ فارسی	
۴۱	۴-۲- نمودار مرحله یادگیری	
۴۲	۴-۳- حرف "ز" در سه اندازه مختلف	
۴۴	۴-۴- نمودار بازشناسی در روش ساختاری	
۴۹	۴-۵- نمودار کلی استفاده از <i>padaline</i>	

فهرست جداول

<u>صفحه</u>	<u>عنوان جدول</u>	<u>فصل سوم:</u>
۲۴	۳-۱- عملکرد نماها در شناسایی حروف تایی فارسی	
۲۵	۳-۲- ویژگی های استفاده شده در روش ساختاری	
۲۵	۳-۳- چند نمونه از قوانین توصیف حروف	
۲۸	۳-۴- ویژگی های استفاده شده در این پروژه	
۲۹	۳-۵- وضعیت حروف نسبت به خط زمینه	
۳۰	۳-۶- تقسیم بندی حروف براساس حفره	
۳۲	۳-۷- تقسیم بندی حروف براساس تعداد قطعات در حالت ساده	
۳۳	۳-۸- تقسیم بندی حروف براساس تعداد قطعات	
۳۵	۳-۹- ویژگی های حرف "ض"	
فصل چهارم:		
۴۰	۴-۱- تقسیم بندی حروف بزرگ فارسی	
۴۵	۴-۲- ویژگی های استخراج شده در متال	
۴۵	۴-۳- حاصل استخراج موارد مشابه در تعداد قطعات	
۴۶	۴-۴- حاصل استخراج موارد مشابه از بافر با قوس های مساری	
۵۰	۴-۵- ویژگی انتخابی جهت استفاده از روش <i>padaline</i>	
۵۰	۴-۶- ویژگی انتخابی در دسته اول	
۵۱	۴-۷- ویژگی انتخابی در دسته دوم	

فهرست فرمولها

<u>صفحه</u>	<u>عنوان فرمول ها</u>	<u>فصل</u>
۴۷	۱-۴- نمایش <i>Padaline</i> در حالت ساده	فصل چهارم:
۴۷	۲-۴- روش محاسبه ثابت های چندجمله ایها	
۴۸	۳-۴- روش محاسبه $exp(Ai)$	
۴۸	۴-۴- روش محاسبه ثابت های چندجمله ایها بصورت کلی	
۵۵	۵-۴- فرمولهای بدست آمده از آموزش شبکه عصبی	

مقدمه

از زمانیکه کامپیوتر توسط بشر ساخته شده، بشر همواره در آرزوی این بوده است که بگونه ای خواندن و نوشتن را با کمک کامپیوتر سهل تر بنماید، در زمینه نوشتن موفقیت های چشمگیری برای بشر حاصل شده است، ولی برای خواندن هنوز راه طولانی در پیش است.

هدف کلی بازشناسی حروف (*optical character recognition*) تشخیص

حروف مختلف توسط کامپیوتر میباشد. نمونه های از کاربردهای بازشناسی، تشخیص حروف زبانهای مختلف بصورت تایپی و دستنویس، تشخیص پلاک خودروها و غیره میباشد. بدیهی است جواب مناسب در زمینه هر کدام از موارد اشاره شده کاربردهای متعدد و موثری در زندگی روزمره خواهد داشت. بعنوان مثال تشخیص پلاک خودروها، باتوجه به تصاویر ارائه شده از دوربین های مختلف که در محل های مشخص نصب شده است، کمک موثری در ردیابی رانندگان متخلف میباشد. هم چنین بازشناسی حروف تایپی و چاپی برای ورود اطلاعات چاپ شده به کامپیوتر با سرعت و دقت و درعین حال با هزینه کمتر کاربرد دارد. بطور متوسط میتوان گفت که دقت یک تایپیست حرفه ای که حروف را تشخیص داده و سپس تایپ میکند در حدود ۹۵٪ میتواند باشد، در حالیکه دقت نرم افزارهای حرفه ای بازشناسی حروف تایپی که اخیراً به بازار معرفی شده اند در حدود ۹۹٪ است [1]. سرعت نرم افزارهای مذکور بر روی یک کامپیوتر ۴۸۶ بیشتر از ده برابر یک تایپیست حرفه ای است [2]. بدیهی است با بهینه شدن روشهای بازشناسی و سرعت کامپیوترها، نسبت فوق نیز افزایش خواهد یافت.

بازشناسی حروف توسط کامپیوتر را بطور کلی به دو دسته *on line* و *off line* میتوان تقسیم بندی کرد. در مدل *on line* کاربرد با استفاده از وسایلی نظیر موشواره و قلم نوری حروف و یا کلمه مورد نظر خود را مینویسد و در لحظه نوشتن کامپیوتر حروف نوشته شده را تشخیص خواهد داد. در مدل *off line* تصاویری که شامل حروف مورد نظر جهت بازشناسی میباشد در حافظه و یا دیسک ذخیره میشوند و سپس پردازش بازشناسی حروف بر روی این تصاویر صورت خواهد پذیرفت. در این پایان نامه مدل *off line* جهت تشخیص حروف تایی مورد نظر است.

از سایر جنبه ها ، بازشناسی حروف تایی را به دسته بندیهای مختلفی می توان تقسیم کرد. بطور مثال تقسیم بندی بر اساس اتصال یا انفصال حروف در کلمه هارا میتوان ذکر کرد. در زبانهای نظیر فارسی و عربی حروف با اتصال بهم تشکیل کلمه ها را میدهند ، در عین حال شکل بعضی از حروف در زمان اتصال تغییر خواهد کرد و درحالیکه در زبان های نظیر لاتین و چینی حروف تشکیل دهنده کلمه ها جدا از هم میباشد. میتوان ساختار حروف را نیز از نظر ثابت یا متغییر بودن مدنظر قرار داد . ساختار ثابت مانند یک متن تایی با اندازه ثابت حروف می باشد و ساختار متغییر مانند متون دستنویس و یا متون تایی با اندازه و قلم های متفاوت میباشد.

در مقایسه با تشخیص حروف لاتین ، در زبان فارسی برای بازشناسی حروف در کلمه ها باید ابتدا عملیات جداسازی صورت پذیرد ، که این امر عملیات بازشناسی را پیچیده تر و طولانی تر خواهد کرد. بدیهی است صحت عملیات بازشناسی ارتباط مستقیم با جداسازی

جداسازی دارد. بطور کلی بازشناسی حروف فارسی و زبانهای مشابه از قسمت های زیر

تشکیل شده است: ۱- ورود تصویر به کامپیوتر (*Digitizing*)

۲. پردازش تصویر (*Image Processing*)

۳. جداسازی حروف (*Segmentation*)

۴. استخراج ویژگیها (*Feature extraction*)

۵. طبقه بندی و بازشناسی حروف (*Classification & Recognition*)

در زبانهایی که حروف تشکیل دهنده کلمه ها جدا از هم میباشند ، عملیات جداسازی

صورت نخواهد پذیرفت و مراحل کلی بازشناسی به چهار مرحله خلاصه میگردد. در عین حال

هر کدام از مراحل اشاره شده ممکن است به چندین قسمت دیگر تقسیم شوند.

۱- در مرحله اول ابتدا باید تصویر با ابزاری نظیر اسکنر، برنامه های نقاشی و غیره

بصورت دیجیتالی وارد کامپیوتر گردد. البته در این پروژه تصاویر حاوی اطلاعات از برنامه

گرافیکی (نقاشی) تحت ویندوز دریافت میشود بنابراین حداقل کیفیت مورد آزمایش در این

پروژه 75dpi خواهد بود. این تصاویر به صورت سیاه و سفید ذخیره میشوند . بنابراین در این

پروژه مرحله پردازش تصویر مورد نظر نمی باشد.

۲- یکی از مطالب مهم در بازشناسی حروف فارسی عملیات جداسازی حروف است.

اولین قدم در جداسازی حروف تایی فارسی در این پروژه پیدا کردن خط زمینه است . بطوریکه

مشخص است بیشترین تراکم نقاط سیاه در یک جمله در خط زمینه میباشد. یکی از روشهای

مناسب در بدست آوردن خط زمینه ، تراکم نقاط سیاه در یک جمله مورد پردازش میباشد، این

تراکم با شمردن تعداد نقاط سیاه در هر سطر (در صورتیکه قرارگیری نقاط تشکیل دهنده یک جمله را در یک ماتریس فرض کنیم) بدست می آید. بیشترین تراکم نقاط سیاه در واقع در خط زمینه یا در نزدیکی خط زمینه است. با استفاده از موقعیت و قطر خط زمینه روشی جهت جداسازی ارائه میگردد. این بخش از پروژه در فصل "جداسازی حروف" شرح داده شده است.

۳- مهمترین قسمت بازشناسی، شناسایی و استخراج ویژگیهای مناسب است و در بخش "استخراج ویژگیها" درباره آن بحث خواهیم کرد. این ویژگیها میتوانند بصورت کمیت های عددی و یا توصیفی بیان شوند و بصورت کلی و عمومی روش ثابت و تعریف شده ای برای شناسایی و استخراج ویژگیهای مناسب وجود ندارد. نوع ویژگی های که انتخاب میشوند تا حد زیادی دقت و سرعت بازشناسی را نشان میدهند. در بخش "استخراج ویژگیها"، خصوصیات مختلفی که در مقالات داخلی و خارجی برای شناسایی حروف تایی فارسی معرفی شده، ارائه میگردد. در این پروژه ویژگی های جدیدی انتخاب شده است که برخی از آنها شامل: درصد نقاط سیاه در بالای خط زمینه، تعداد حفره ها در هر حرف، تعداد قوسها یا دهانه ها در جهت های اصلی و در جهت های اصلی با ۴۵ درجه چرخش حرف، تعداد قطعه ها در هر حرف، میباشند. ویژگی های مذکور پس از پردازش اولیه، در مراحل آموزش و بازشناسی مورد استفاده قرار میگیرند.

۴- شیوه های مختلفی نظیر شبکه های عصبی، الگوریتم های آماری، سیستم های خبره، بانک های اطلاعاتی و غیره برای کلاسه بندی وجود دارد. در این پروژه از دو روش

شبکه های عصبی چند جمله ای های جداکننده (*Padaline's*) و بانک های اطلاعاتی (*Microsoft Access*) برای کلاسدی و بازشناسی استفاده شده است . البته روش شبکه های عصبی به همراه ترکیبی از بانک های اطلاعاتی (*Microsoft Access*) مورد استفاده قرار گرفته اند . این بخش از پایان نامه در فصل "بازشناسی حروف" مورد بحث قرار گرفته است .

در انتهای پایان نامه الگوریتم های بکار برده شده در بخشهای جدا سازی ، استخراج ویژگی ها و بازشناسی شرح داده میشوند . در بخش ضمازم، جداول تقسیم بندی حروف و نتایج حاصل از استخراج ویژگی ها از یک نمونه آزمایشی ارائه شده است . همچنین لیست توابع مهم برنامه که به زبان *Visual Basic* و تحت ویندوز میباشد، در "ضمیمه ب" آورده شده است .

فصل دوم

جداسازی حروف

Character Segmentation

مقدمه:

یکی از مهمترین قسمت های بازشناسی حروف فارسی و یا زبان های مشابه ، جداسازی حروف در کلمه ها میباشد . صحت بازشناسی ارتباط مستقیم با عملیات جداسازی دارد و چنانچه جداسازی صحیح انجام نپذیرد ، نمی توان انتظار نتایج درست از پردازش های بازشناسی داشت . بصورت کلی روش های مختلفی برای جداسازی وجود دارد که در این فصل به شرح آنها پرداخته میشود و همچنین روش پیشنهادی در این پروژه که مورد آزمایش قرار گرفته است ، بطور کامل شرح داده میشود. چند ویژگی مهم حروف فارسی که تاثیر بسیار زیادی در عملیات جداسازی و همچنین در بازشناسی دارند بصورت مختصر در زیر بیان میشود:

- حروف فارسی دارای دو ، سه و یا چهار نماد میباشند . شکلهای مختلف یک حرف را اصطلاحاً نماد می نامند ، این نمادها بستگی به محل قرارگیری در کلمه ها مشخص میشوند. برخی از نمادهای حروف کاملاً جدا هستند و بعضی از دو طرف و بعضی فقط از راست به حروف دیگر متصل میشوند و محل اتصال حروف به یکدیگر