



۱۷۲۹

سید

سید



دانشگاه تربیت مدرس

دانشکده فنی و مهندسی

پایان نامه کارشناسی ارشد مهندسی الکترونیک

شناسایی ساختاری حروف دستنویس فارسی

غلامرضا اردشیر بهرستاقی

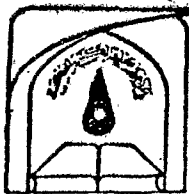
استاد راهنما:

دکتر احسان‌الله کبیر



۵۲۷

پاییز ۱۳۷۲



دانشگاه تربیت مدرس

بسمه تعالی

چکیده پایان نامه

عنوان پایان نامه: شناسایی ساختاری حروف دستنویس فارسی

نام نویسنده: غلامرضا اردشیر بهرستانی

استاد راهنما: دکتر احسان الله کبیر

اسانید مشاور:

دانشکده: فنی و مهندسی

رشته: برق - الکترونیک

تاریخ دفاع: ۷۲، ۷، ۲۱

تحقیق انجام شده در زمینه بازشناسی متون دستنویس است که یکی از شاخه‌های بازشناسی الگومی باشد. هدف، بازشناسی حروف دستنویس فارسی است که بطور مجزا توسط افراد مختلف نوشته باشند. فرض بر این است که حروف با دقت ۲۰۰ نقطه در اینج تصویربرداری و با کد MSP ذخیره شده اند.

الگوریتم طرح شده را می‌توان به دو بخش اساسی تقسیم نمود: در بخش اول، اقدام به استخراج ویژگیها می‌شود و در بخش دوم، حروف بر اساس یک درخت تصمیم به عنوان طبقه‌بندی‌کننده مورد بازشناسی قرار می‌گیرند. در مرحله استخراج ویژگیها، در کام اول یک مرحله پیش پردازش برای رفع نریذکی و حفره و نویز صورت می‌گیرد. سپس با اجرای الگوریتم برچسب‌زنی مولفه‌ها، نقطه و سرکش موجود از بدنه اصلی حرف جدا و نوع آن تشخیص داده می‌شود. در نهایت با گروه‌بندی بخشهای سطرهای مختلف ماتریس تصویر حرف، زیرحروف استخراج و مورد شناسایی قرار می‌گیرند.

برای یک مجموعه متشکل از ۶۱ نمونه از هر حرف، میزان بازشناسی ۹۱/۵ درصد بدست آمده است و برای مجموعه ای دیگر متشکل از ۴۰ نمونه از هر حرف ۹۲/۱۲ درصد حروف بطور صحیح بازشناسی شدند.

تحقیق انجام شده می‌تواند کاربردهای زیادی داشته باشد. برای مثال الگوریتم حاصل از این تحقیق را می‌توان به منظور بازشناسی خودکار حروف مجزایی استفاده نمود که در فرمهای گوناگون جمع آوری و ثبت اطلاعات و کدگذاریها به کار برده می‌شوند.

خواهشمند است این برگ را نشود

تقدیم به :

تمامی شهیدان راه حق

تقدیم به :

پدر و مادر و همسر مهربانم که همواره

مشوقم بوده اند.

تشکر و قدردانی:

با تشکر از استاد گرامی و ارجمند آقای دکتر
احسان‌الله کبیر به خاطر راهنمایی‌های ارزشمند
ایشان و دوست عزیزم آقای مهندس رضا عزمی و
سایر کسانی که در انجام این تحقیق مرا یاری
نمودند.

این تحقیق با پشتیبانی مالی و تجهیزاتی بخش برق
و الکترونیک سازمان پژوهش‌های علمی صنعتی ایران
انجام شده است.

چکیده:

تحقیق انجام شده در زمینه بازشناسی متون دستنویس است که یکی از شاخه‌های بازشناسی الگومی باشد. هدف، بازشناسی حروف دستنویس فارسی است که بطور مجزا توسط افراد مختلف نوشته^{شده} نباشند. فرض بر این است که حروف با دقت ۲۰۰ نقطه در اینج تصویربرداری و با کد MSP ذخیره شده‌اند.

الگوریتم طرح شده را می‌توان به دو بخش اساسی تقسیم نمود: در بخش اول، اقدام به استخراج ویژگیها می‌شود و در بخش دوم، حروف بر اساس یک درخت تصمیم به عنوان طبقه‌بندی‌کننده مورد بازشناسی قرار می‌گیرند. در مرحله استخراج ویژگیها، در گام اول یک مرحله پیش‌پردازش برای رفع بریدگی و حفره و نویز صورت می‌گیرد. سپس با اجرای الگوریتم برجسبازنی مولفه‌ها، نقطه و سرکش موجود از بدنه اصلی حرف جدا و نوع آن تشخیص داده می‌شود. در نهایت با گروه‌بندی بخشهای سطوهای مختلف ماتریس تصویر حرف، زیرحروف استخراج و مورد شناسایی قرار می‌گیرند.

برای مجموعه تمرین متشکل از ۱۶ نمونه از هر حرف، میزان بازشناسی ۹۱/۵ درصد بدست آمده است و برای مجموعه آزمایش متشکل از ۴۰ نمونه از هر حرف ۹۲/۱۲ درصد حروف بطور صحیح بازشناسی شده‌اند.

تحقیق انجام شده می‌تواند کاربردهای زیادی داشته باشد. برای مثال الگوریتم حاصل از این تحقیق را می‌توان به منظور بازشناسی خودکار حروف مجزایی استفاده نمود که در فرمهای گوناگون جمع‌آوری و ثبت اطلاعات و کدگذاریها به کار برده می‌شوند.

صفحه	عنوان
	فصل اول : مقدمه
۱	۱-۱- کلیات
۲	۲-۱- عمل با زشناسی
۳	۳-۱- وضعیت با زشناسی کا راکترها درزبا نهایی مختلف
۴	۴-۱- موضوع این تحقیق
۵	۵-۱- سازمان کلی رساله
	فصل دوم : با زشناسی الگو
۶	۱-۲- الگو
۷	۲-۲- با زشناسی الگو
۸	۳-۲- روشهای مختلف در با زشناسی الگو
۸	۱-۳-۲- روشهای تحلیلی
۸	الف (روشهای آماری
۱۰	ب (روشهای معین
۱۰	۲-۳-۲- روشهای نحوی یا ساختاری
۱۲	۳-۳-۲- مقایسه روشهای تحلیلی و نحوی
۱۴	۴-۲- با زشناسی کا راکترهای دستنویس
۱۶	۵-۲- نمونه‌هایی از کارهای انجام شده
۱۶	۱-۵-۲- استخراج ویژگیها
۱۸	۲-۵-۲- طبقه‌بندی
۱۹	۳-۵-۲- نتایج

(ب)

فهرست مطالب

صفحه

عنوان

فصل سوم : مراحل بازشناسی حروف دستنویس فارسی

۲۱	۱-۳-۱- مقدمه
۲۲	۲-۳-۲- تعریف بخش
۲۳	۳-۳-۳- پیش پردازش
۲۴	۱-۳-۳-۱- رفع بریدگی
۲۴	۲-۳-۳-۲- رفع حفره ونویز
۲۷	۴-۳-۴- تشخیص نقطه وسرکش
۳۲	۵-۳-۵- تعیین نوع نقطه وسرکش
۳۵	۶-۳-۶- پیدا کردن ضخامت قلم
۳۵	۷-۳-۷- زیرحرف
۳۶	۸-۳-۸- تعریف گروه وگروه بندی
۳۹	۹-۳-۹- اصلاح گروه بندی
۴۰	۱۰-۳-۱۰- شناسایی نوع گروه یا زیرحرف
۴۰	۱-۱۰-۳-۱- شناسایی کمان
۴۱	۲-۱۰-۳-۲- شناسایی گروه " S "
۴۲	۳-۱۰-۳-۳- شناسایی خط مورب
۴۲	۴-۱۰-۳-۴- شناسایی گروه " T "
۴۴	۵-۱۰-۳-۵- شناسایی خطوط افقی وعمودی
۴۵	۶-۱۰-۳-۶- الگوریتم شناسایی نوع گروه
۴۵	۱-۳-۱۱-۱- شناسایی منحنی بسته
۴۷	۲-۳-۱۲-۲- شناسایی دندان

فهرست مطالب

صفحه	عنوان
۵۱	۳-۱۳- شناسایی مجموعه‌گروه‌هایی که به صورت کمان هستند
۵۲	۳-۱۴- شناسایی مجموعه‌گروه‌هایی که به صورت خط مورب هستند
۵۳	۳-۱۵- شناسایی فرورفتگی در "ی" و حروف "ک، گ" با سرکش متصل به بدنه
فصل چهارم : الگوریتم بازشناسی حروف	
۵۵	۴-۱- مقدمه
۵۵	۴-۲- دسته‌بندی حروف
۵۵	۴-۲-۱- دسته‌بندی حروف براساس پایین‌ترین گروه
۵۸	۴-۲-۲- دسته‌بندی حروف براساس حرکت در دو طرف گروه \mathcal{G}
۶۱	۴-۲-۳- دسته‌بندی حروف براساس تعداد دگروه‌ها در دو طرف گروه \mathcal{G}
۶۳	۴-۲-۴- دسته‌بندی حروف براساس مرکز ثقل
۶۴	۴-۲-۵- دسته‌بندی حروف براساس تعداد دگروه‌های موجود در بدنه آنها
۶۶	۴-۳- الگوریتم بازشناسی حروف
۸۹	۴-۴- نتایج
فصل پنجم : خلاصه، نتیجه‌گیری، پیشنهادهای	
۹۴	۵-۱- مقدمه
۹۴	۵-۲- شرایط تحقیق انجام شده و مشخصات مجموعه داده
۹۵	۵-۳- خلاصه‌ای از مراحل اساسی الگوریتم و نتایج بدست آمده
۹۶	۵-۴- پیشنهادها
۹۸	منابع

فصل اول

۱-۱ کلیات

تحقیقات در زمینه بازشناسی الگو تقریباً از سال ۱۹۵۰ شروع شده است. اولین تلاشها در زمینه بازشناسی الگو، برای تهیه برنامه‌های کامپیوتری جهت تصمیم‌گیری خودکار و توسعه سخت‌افزارهای خاص برای خواندن الگوها — مانند کاراکترهای ارقام و حروف تایپ شده بوده است. در اواخر دهه ۱۹۵۰، — Resenblatt یک مدل اولیه در الگوریتم پرسپترون^۱ برای آرایش و ذخیره اطلاعات در مغز عرضه کرد. اکثر روشهای عرضه شده برای مسئله بازشناسی الگو در آن زمان بر اساس تئوری تصمیم‌آماری^۲ و اصول منطق آستانه‌های^۳ بنا شده بود. در طی دهه ۱۹۶۰ تحقیقات در زمینه طراحی سیستمهای بازشناسی الگو مقدار زیادی پیشرفت داشت. با استفاده از تئوری زبان کامپیوتر^۴ و استفاده از از تواناییهای پردازش کامپیوتر، روش نحوی^۵ به عنوان یک مکمل برای تکنیکهای تحلیلی^۶ در حل مسائل بازشناسی الگوهای تصویری معرفی شد (Toe ۸۱). — مفاهیم بازشناسی الگو بطور فزاینده‌ای در طراحی سیستمهای مدرن اطلاعاتی کامپیوتری به عنوان یک عامل مهم در نظر گرفته شده‌اند. جذابیت در این زمینه هنوز هم با سرعت زیاد در شدمی‌کند. در شاخه‌های مختلف مانند مهندسی، علوم کامپیوتر، علوم اطلاعات^۷، فیزیک، شیمی، زبان، روانشناسی، زیست‌شناسی، فیزیولوژی، پزشکی و زبان‌شناسی موضوع بازشناسی الگو مطرح است. هر کدام از این شاخه‌ها به جنبه خاصی از مسئله بازشناسی الگو اهمیت می‌دهند.

1. Pattern Recognition
2. Perceptron
3. Statistical
4. Threshold Logic Principles
5. Computer Language Theory
6. Syntactic
7. Analytical

بازشناسی الگوارویزگیهای انسان و دیگر موجودات زنده است که در تمام لحظاتی که در بیداری هستند، بطور خودکار آن را انجام می دهند. انسان از این نظر که از توانائی بالائی در بازشناسی برخوردار است یک سیستم کاملاً "خبره" و ماهر است.

الگو توصیفی از یک موضوع است. از دیدگاههای مختلف، تقسیمات مختلفی می توان برای بازشناسی الگو انجام داد. بسته به طبیعت الگوهای که با یستی بازشناسی شوند، ممکن است عمل بازشناسی را به دو نوع عمده "بازشناسی نمونه های مجرد و انتزاعی" و نمونه های ملموس تقسیم نمود. هنگامی که یک بحث قدیمی را به خاطر می آوریم، یا راه حل مسئله ای را پیدا می کنیم یک بازشناسی الگوی مجرد و انتزاعی انجام داده ایم. در بازشناسی نمونه های ملموس، الگوهای زمانی^۳ و مکانی^۴ بازشناسی می شوند. بازشناسی کاراکترها، بازشناسی آشکارسنگت، خواندن نقشه های هواشناسی، موضوعات فیزیکی و تصاویر، نمونه های از بازشناسی الگوهای مکانی هستند. از نمونه های الگوهای زمانی میتوان به شکل موجهای صحبت، الکتروکاردیوگرامها و الکتروانسفالوگرامها اشاره کرد.

مطالعه در مورد بازشناسی الگو ممکن است به دو زمینه عمده تقسیم شود:

- ۱- مطالعه توانائی بازشناسی الگو توسط انسان و دیگر موجودات زنده.
- ۲- مطالعه تئوریهها و تکنیکهای مناسب برای طراحی دستگاہهایی که قادر به بازشناسی الگو باشند.

زمینه اول به روانشناسی، فیزیولوژی و زیست شناسی مربوط می شود

-
1. Abstract
 2. Concrete
 3. Temporal
 4. Spatial

وزمینہ دوم اصولاً "با مهندسی، کامپیوتر و علوم اطلاعات سروکار دارد. موضوع این تحقیق در گروه دوم قرار دارد. به زبان ساده با زشناسی الگو می توانند به عنوان دسته بندی داده های ورودی در کلاسهای قابل با زشناسی از طریق استخراج ویژگیها و خواص با معنی آنها تعریف شود. سیستم با زشناسی کا راکتریک سیستم با زشناسی الگومی باشد که تصویر کا راکترها، به عنوان داده های ورودی آن و نام کا راکترها، خروجی های آن محسوب می شود.

۳-۱- وضعیت با زشناسی کا راکترها در زبانهای مختلف

در زمینہ با زشناسی کا راکترهای انگلیسی تحقیقات قابل ملاحظه ای صورت گرفته و مقالات زیادی چاپ شده است. الگوریتمهایی که برای کا راکترهای تایپ شده یا دستنویس طراحی شده، قادر به با زشناسی کا راکترها با دقت خوبی می باشد. در مورد حروف چینی و ژاپنی نیز تحقیقاتی انجام گرفته است. در زمینہ با زشناسی کا راکترهای عربی، فارسی و هندی تحقیقات کمی صورت گرفته است. تفاوت بین کا راکترهای لاتین و کا راکترهای عربی و فارسی اجازه استفاده مستقیم از الگوریتمهای مربوط به با زشناسی کا راکترهای لاتین را برای کا راکترهای عربی و فارسی نمی دهد.

از تحقیقات انجام شده برای شناسایی کا راکترهای عربی می توان موارد زیر را نام برد. کامبیز بدیع در سال ۱۹۸۰ حدود ۹۰ درصد از ۱۲ حرف بزرگ (۱۰ تا ۱۲ نمونه از هر حرف) دستنویس عربی را با زشناسی نمود (Bad ۸۰ و Bad ۸۲). آدنان امین در سال ۱۹۸۰ تعداد ۷۳ کا راکتر را که توسط ۳ نویسنده نوشته شده بودند با دقت ۹۲ تا ۹۵ درصد با زشناسی کرد (Ami ۸۰ و Ami ۸۲). یوسف در سال ۱۹۸۸ حدود ۹۹ درصد از کا راکترهای دستنویس عربی را با زشناسی نمود که شامل ۵۰ نمونه در هر کلاس نوشته شده توسط ۲۵ نفر بوده است (You ۸۸).

از تحقیقات انجام شده در زمینه حروف فارسی می توان به کار عزمی اشاره نمود که حدود ۲۰۰ نمونه از ۳۳ کلاس در نظر گرفته شده برای حروف دستنویس فارسی را که توسط افراد مختلف نوشته شده بودند با دقت ۹۰٪ بازشناسی کرده است (عزمی ۷۲).

۴-۱- موضوع این تحقیق

موضوع این تحقیق، بازشناسی حروف دستنویس فارسی است. حروف به صورت مجزا فرض شده و توسط افراد مختلف نوشته شده اند. به منظور تصویربرداری از یک روبشگزنوری استفاده شده است که از حروف با دقت ۲۰۰ dpi تصویربرداری نموده و با کد MSP در حافظه ذخیره می نماید.

برای بازشناسی حروف از ویژگیهای ساختاری آنها استفاده شده است. ویژگیهای اصلی عبارتند از: کمان، خطوط مورب، خطوط عمودی و افقی و شکلهای "S" و "T". بعضی از ویژگیهای ساختاری حروف فارسی نیز مورد توجه قرار گرفته اند. دندان در حروف "س" و "ش"، منحنی بسته در حرف "ص" و غیره و فرورفتگی در حروف "ی" و "ک" و "گ" از این جمله اند.

برای بازشناسی از یک درخت تصمیم به عنوان طبقه بندی کننده استفاده شده است. این درخت تصمیم به صورت تجربی طرح شده است. در این تحقیق، از ۶۱ نمونه برای هر کلاس به عنوان مجموعه تمرین استفاده شده که میزان بازشناسی ۹۱/۵ درصد حاصل شده است. برای یک مجموعه دیگر از حروف (مجموعه آزمون) شامل ۴۰ نمونه از هر کلاس میزان بازشناسی ۹۲/۲۲ بدست آمده است.

۱-۵- سازمان کلی رساله

- در فصل دوم این رساله، ابتدا در مورد الگو و بازشناسی الگو توضیح داده میشود.
- سپس در مورد روشهای مختلف بازشناسی الگو و تقسیم بندی آنها مطالبی بیان میشود.
- در آخر فصل دوم درباره بازشناسی کاراکترهای دستنویس و نمونه کارهای انجام شده توضیحات کوتاهی آورده میشود.
- در فصل سوم، مراحل مختلف استخراج ویژگیهای ساختاری بیان می شود.
- در فصل چهارم الگوریتم بازشناسی حروف را شرح خواهیم داد.
- فصل پنجم به بررسی نتایج و پیشنهادهای اختصاصی دارد.

فصل دوم : بازشناسی الگو

در این فصل ابتدا مفهوم الگو و بازشناسی الگو مطرح می شود . سپس در مورد روشهای مختلف بازشناسی الگو و تقسیم بندی آنها به روشهای آماری و نحوی (ساختاری) مطالبی بیان می شود. در آخر درباره بازشناسی کاراکترهای دستنویس و نمونه کارهای انجام شده توضیحات کوتاهی آورده می شود .

۱-۲- الگو :

در متون مهندسی بازشناسی الگو^۲ ، توصیفهای مختلفی از الگو ارائه شده است (Kab ۹۰) : "یک الگو توصیف یک شی است" ، "درکلیترین عبارت الگوها وسائلی هستند که ما جهان را بوسیله آنها تفسیر می کنیم" ، "بسیاری از اطلاعاتی که در زندگی با آنها سروکار داریم به شکل الگوهای پیچیده هستند" ، "یک الگو یک توصیف کمی یا ساختاری از یک شی است" ، "یک الگو ترکیبی از n مورد است که n می تواند برابر با ۲ یا بزرگتر از آن باشد" .

در یک بررسی دقیق ، کلمه الگو به دو صورت تعریف می شود : در -

تعریف اول الگو به عنوان یک پدیده یا موضوع الگو شده تعریف می شود . به عبارت دیگر پدیده یا موضوع خودش به عنوان یک الگو در نظر گرفته می شود . بنا براین بازشناسی الگو همان طبقه بندی است . در تعریف

1. Pattern

2. Pattern Recognition