

بِسْمِ اللَّهِ الرَّحْمَنِ الرَّحِيمِ



دانشگاه صنعتی شاهرود

دانشکده مهندسی برق و رباتیک

گروه الکترونیک

بازشناسی متن تایپی نوشته شده با قلم Iranian sans

دانشجو: زینب باقری

استاد راهنما:

جناب آقای دکتر خسروی

پایان نامه ارشد جهت اخذ درجه کارشناسی ارشد

زمستان ۱۳۹۲

تقدیم به:

همسر دلسوز و فداکارم

و

پدر و مادر مهربانم

تشر و قدردانی

با نثار تشر و قدردانی بی‌انتها از زحمات و راهنمایی‌های صبورانه استاد ارجمند، جناب آقای دکتر حسین خسروی، دوره کارشناسی ارشد برای من فرصتی بود که در محضر این بزرگوار و سایر اساتید گران‌قدر، درس انسانیت، درستی و تحقیق را فراگیرم.

تعهد نامه

- اینجانب زینب باقری دانشجوی دوره کارشناسی ارشد رشته برق دانشکده برق و رباتیک دانشگاه صنعتی شاهرود نویسنده پایان نامه بازنشاسی متن تاییپی نوشته شده با قلم Iranian sans تحت راهنمایی جناب آقای دکتر حسین خسروی متعهد می‌شوم.
- تحقیقات در این پایان نامه توسط اینجانب انجام شده است و از صحت و اصالت برخوردار است.
 - در استفاده از نتایج پژوهش‌های محققان دیگر به مرجع مورد استفاده استناد شده است.
 - مطالب مندرج در پایان نامه تاکنون توسط خود یا فرد دیگری برای دریافت هیچ نوع مدرک یا امتیازی در هیچ جا ارائه نشده است.
 - کلیه حقوق معنوی این اثر متعلق به دانشگاه صنعتی شاهرود می‌باشد و مقالات مستخرج با نام « دانشگاه صنعتی شاهرود » و یا « Shahrood University of Technology » به چاپ خواهد رسید.
 - حقوق معنوی تمام افرادی که در به دست آمدن نتایج اصلی پایان نامه تأثیرگذار بوده‌اند در مقالات مستخرج از پایان نامه رعایت می‌گردد.
 - در کلیه مراحل انجام این پایان نامه ، در مواردی که از موجود زنده (یا بافت‌های آن‌ها) استفاده شده است ضوابط و اصول اخلاقی رعایت شده است.
 - در کلیه مراحل انجام این پایان نامه، در مواردی که به حوزه اطلاعات شخصی افراد دسترسی یافته یا استفاده شده است اصل رازداری ، ضوابط و اصول اخلاق انسانی رعایت شده است.

تاریخ

امضای دانشجو

مالکیت نتایج و حق نشر

- کلیه حقوق معنوی این اثر و محصولات آن (مقالات مستخرج، کتاب، برنامه های رایانه ای، نرم افزار ها و تجهیزات ساخته شده است) متعلق به دانشگاه صنعتی شاهرود می باشد.
- این مطلب باید به نحو مقتضی در تولیدات علمی مربوطه ذکر شود.
- استفاده از اطلاعات و نتایج موجود در پایان نامه بدون ذکر مرجع مجاز نمی باشد.

چکیده

در دهه‌های اخیر تحقیقات گسترده‌ای در زمینه‌ی بازشناسی الگوهای نوشتاری شامل حروف، ارقام و سایر نمادهای متداول در اسناد مکتوب، به زبان‌های مختلف انجام شده است. با توجه به پیشرفت‌های حاصل شده در این زمینه، فناوری بازشناسی خودکار متون تحت عنوان بازشناسی نوری حروف یا ا.سی.آر شکل گرفته است. بازشناسی متن یکی از مهم‌ترین بخش‌های دولت الکترونیک به شمار می‌رود و در سال‌های اخیر در کشور ما نیز تقاضا برای یک سیستم بازشناسی متن فارسی، به شدت افزایش یافته است. با توجه به آنکه حجم زیادی از اسناد کاغذی موجود، توسط اسکنرها یا دوربین‌ها به اسناد تصویری دیجیتالی تبدیل می‌شوند؛ ذخیره‌سازی، بازیابی و مدیریت کارآمد این اسناد تصویری، در بسیاری از برنامه‌ها نظیر اتوماسیون اداری و کتابخانه‌های دیجیتالی اهمیت دارند.

به طور کلی سامانه بازشناسی متن شامل بخش‌های مختلفی از قبیل دریافت تصویر، پیش‌پردازش، آنالیز پیکربندی، تشخیص زبان، تشخیص قلم و در نهایت تشخیص متن می‌باشد. تحقیقات انجام شده در بعضی از این زمینه‌ها از قبیل پیش‌پردازش مستقل از زبان متن بوده و برای هر زبانی قابل استفاده است. لیکن برخی از قسمت‌های دیگر مانند تشخیص قلم و تشخیص متن به زبان متن وابسته بوده و نمی‌توان به طور مستقیم نتایج تحقیقات انجام شده برای سایر زبان‌ها را برای فارسی اعمال کرد. بیشتر تحقیقات انجام شده در زمینه‌ی بازشناسی متون فارسی روی تصاویری با درجه تفکیک زیاد، تصاویر متنی تمیز و غیرواقعی و شناسایی متن با چند قلم معروف بوده است. در تحقیقات انجام شده برای بازشناسی متون فارسی سه رویکرد عمده مبتنی بر جداسازی حروف، مبتنی بر بازشناسی شکل کلی زیرکلمات و روش ترکیبی وجود دارد.

در این پایان‌نامه هدف، بازشناسی متن تاپی نوشته شده با قلم Iranian sans با حداقل اندازه ۹ و درجه تفکیک ۳۰۰ نقطه بر اینچ است. این قلم با توجه به دو خصوصیت زیبایی و خوانایی، بسیار مورد توجه واقع شده و روز به روز بر حجم استفاده از آن در محیط رایانه و اینترنت افزوده می‌شود. این قلم

قابلیت جایگزینی با قلم تاهما که قلم پیش فرض سیستم عامل ویندوز است، را دارد. علیرغم خوانایی، فاصله استاندارد بین سطرها، زیبایی و سازگاری با لاتین، این قلم دارای پیچیدگی ساختاری خاصی بوده که این خود عمل بازشناسی آن را پیچیده می‌نماید.

در این پایان‌نامه ابتدا با تولید پایگاه داده مناسب، به آموزش دو طبقه‌بند برای حروف گسسته و پیوسته پرداخته شده و سپس با رفع مشکل همپوشانی زیرکلمات، از رویکرد مبتنی بر جداسازی برای جداسازی حروف استفاده می‌شود. طبقه‌بندهای مورد استفاده از نوع شبکه عصبی چند لایه می‌باشند. در نهایت، نتایج کارایی سامانه‌ی مذکور برای پردازش چند تصویر با متن چاپی، ارائه می‌شود که در آن در بخش جداسازی دقت ۹۶٪ و در بخش شناسایی دقت ۸۵٪ حاصل شد.

کلمات کلیدی: بازشناسی متن، قلم Iranian sans، رویکرد مبتنی بر جداسازی، طبقه‌بند شبکه عصبی

فهرست مطالب

صفحه	عنوان
د.....	فهرست شکل ها.....
ح.....	فهرست جدول ها.....
۱.....	فصل ۱: مقدمه.....
۲.....	۱-۱- مقدمه.....
۲.....	۱-۱-۱- معرفی بازشناسی نوری حروف.....
۳.....	۲-۱- کاربردهای اُسی.آر.....
۷.....	۳-۱- انواع اُسی.آر.....
۹.....	۴-۱- موانع موجود بر سر تکامل اُسی.آر فارسی.....
۱۰.....	۵-۱- معرفی ویژگی های خط فارسی.....
۱۳.....	۶-۱- نتیجه گیری.....
۱۵.....	فصل ۲: کارهای پیشین.....
۱۶.....	۱-۲- مقدمه.....
۱۶.....	۲-۲- روند تحقیقاتی اُسی.آر.....
۱۹.....	۱-۲-۲- نسل های مختلف نویسه خوان های نوری.....
۲۰.....	۲-۲-۲- سیر تحول اُسی.آر فارسی.....
۲۵.....	۳-۲- نتیجه گیری.....
۲۷.....	فصل ۳: مراحل پیاده سازی اُسی.آر.....
۲۸.....	۱-۳- مقدمه.....

- ۲۸..... ۲-۳ سیستم بازشناسی متن.....
- ۲۹..... ۱-۲-۳ پیش پردازش.....
- ۳۳..... ۲-۲-۳ قطعه بندی.....
- ۳۳..... ۳-۲-۳ تحلیل ساختار فیزیکی و منطقی.....
- ۳۴..... ۴-۲-۳ استخراج ویژگی.....
- ۳۵..... ۵-۲-۳ بازشناسی متن (با یک یا چند طبقه بند).....
- ۳۶..... ۶-۲-۳ پس پردازش.....
- ۳۶..... ۳-۳ روش های بازشناسی متن.....
- ۳۶..... ۱-۳-۳ روش های مبتنی بر جداسازی حروف.....
- ۳۸..... ۲-۳-۳ روش های مبتنی بر شکل کلی زیر کلمات.....
- ۳۹..... ۳-۳-۳ روش های ترکیبی (شکل کلی جداسازی).....
- ۴۰..... ۴-۳ قابلیت های سیستم های بازشناسی متن.....
- ۴۴..... ۵-۳ نتیجه گیری.....
- فصل ۴ : تولید پایگاه داده..... ۴۵**
- ۴۶..... ۱-۴ مقدمه.....
- ۴۷..... ۲-۴ تولید پایگاه داده حروف جدا.....
- ۵۵..... ۳-۴ تولید پایگاه داده حروف پیوسته.....
- ۵۷..... ۴-۴ استخراج ویژگی ها.....
- ۵۸..... ۵-۴ آموزش طبقه بندها.....
- ۵۹..... ۱-۵-۴ طبقه بند حروف گسسته و زیر کلمات.....
- ۶۰..... ۲-۵-۴ طبقه بند حروف پیوسته.....
- ۶۰..... ۶-۴ نتیجه گیری.....

فصل ۵: سامانه بازشناسی متن با قلم **Iranian sans** ۶۱

۱-۵- مقدمه ۶۲

۲-۵- بازشناسی متن ۶۳

۱-۲-۵- پیش پردازش ۶۳

۲-۲-۵- قطعه بندی ۶۵

۳-۲-۵- جداسازی ۶۹

۴-۲-۵- جداسازی حروف به روش پروفایل بالایی تعمیم یافته ۶۹

۱-۴-۲-۵- شناسایی زیرحروف نهایی ۷۴

۵-۲-۵- طبقه بندی نهایی ۸۱

۱-۵-۲-۵- طبقه بندی نهایی با استفاده از برجسب زنی مؤلفه ها ۸۴

۳-۵- ارزیابی الگوریتم برای یک تصویر کامل ۹۱

۴-۵- ارزیابی الگوریتم برای تصویرهای دیگر ۹۴

۱-۴-۵- نتیجه ارزیابی الگوریتم ۱۰۲

۵-۵- نتیجه گیری ۱۰۳

فصل ۶: نتیجه گیری و پیشنهادات ۱۰۵

۱-۶- نتیجه گیری ۱۰۶

۲-۶- نوآوری ها ۱۰۷

۳-۶- پیشنهادات ۱۰۷

فهرست مراجع ۱۰۹

فهرست شکل‌ها

صفحه	عنوان
۳	شکل ۱-۱ : مراحل مختلف ا.سی.آر
۶	شکل ۲-۱ : نمونه‌ای از یک فرم ثبت‌نام کنکور کارشناسی ارشد [۲]
۹	شکل ۳-۱ : سیستم‌های شناسایی برون خط و برخط
۱۰	شکل ۴-۱ : آشنایی با مفهوم حرف و زیر کلمه
۱۱	شکل ۵-۱ : شکل‌های مختلف یک حرف در کلمه
۱۲	شکل ۶-۱ : نمونه‌هایی از ادغام حروف مجاور در متون دست‌نویس
۱۲	شکل ۷-۱ : همپوشانی حروف
۱۲	شکل ۸-۱ : اتصال حروف از دو محل
۱۳	شکل ۹-۱ : اتصال حروف از سمت چپ
۱۶	شکل ۱-۲ : سامانه ا.سی.آر هندل [۵]
۱۹	شکل ۲-۲ : قلم NOF یا BICODES [۷]
۲۱	شکل ۳-۲ : جداسازی بر اساس کانتور بالایی [۱۰]
۲۲	شکل ۴-۲ : جداسازی بر مبنای برجسبزی به کانتور زیر کلمه [۱۲]
۲۳	شکل ۵-۲ : یافتن لبه های افقی پایین کلمات [۱۴]
۲۹	شکل ۱-۳ : مراحل سیستم بازشناسی متن
۳۱	شکل ۲-۳ : یافتن زاویه‌ی چرخش با اعمال چند عملگر مورفولوژی در زوایای مختلف [۲۸]
۳۳	شکل ۳-۳ : تصویر نهایی پس از رفع چرخش
۳۴	شکل ۴-۳ : الف) تصویر صفحه اول یک مقاله ب) ساختار فیزیکی ج) ساختار منطقی
۴۳	شکل ۵-۳ : ویژگی‌های تایپوگرافی [۴۰]
۴۳	شکل ۶-۳ : تشکیل یک بافت یکپارچه از روی بلوک متنی

- شکل ۱-۴ : روند کلی فرآیند جداسازی ۴۶
- شکل ۲-۴ : منابع اولیه تولید شده توسط نرم‌افزار ورود برای ساخت پایگاه داده ۴۷
- شکل ۳-۴ : تصویر اسکن شده حروف ۴۸
- شکل ۴-۴ : تولید پایگاه داده متنوع‌تر ۴۹
- شکل ۵-۴ : تصویر باینری شده با نویز فراوان ۵۰
- شکل ۶-۴ : حذف نویز از تصویر ۵۱
- شکل ۷-۴ : برچسب‌زنی مؤلفه‌ها برای تولید پایگاه داده ۵۲
- شکل ۸-۴ : مربعی کردن و ذخیره‌سازی حروف ۵۳
- شکل ۹-۴ : حذف مؤلفه‌های ۴و۳،۲ برای حذف نقاط ۵۳
- شکل ۱۰-۴ : چرخش ۱- و ۱ درجه برای تولید داده‌های متنوع‌تر ۵۵
- شکل ۱۱-۴ : تولید حروف ناقص و نامعتبر ۵۶
- شکل ۱۲-۴ : حروف ناقص تولید شده معتبر ۵۷
- شکل ۱۳-۴ : نمونه‌ای از زیرکلمات معتبر و موجود در زبان فارسی ۵۹
- شکل ۱-۵ : شمای کلی سیستم بازشناسی متن ۶۳
- شکل ۲-۵ : باینری کردن تصویر ۶۴
- شکل ۳-۵ : حذف نویز در دو مرحله ۶۵
- شکل ۴-۵ : اصلاح چرخش به روش افکنش افقی ۶۵
- شکل ۵-۵ : قطعه‌بندی جمله‌ها ۶۶
- شکل ۶-۵ : مشکل وجود همپوشانی در جداسازی زیرکلمه‌ها به روش هیستوگرام عمودی ۶۷
- شکل ۷-۵ : برطرف نمودن مشکل همپوشانی با استفاده از برچسب‌زنی مناسب مؤلفه‌ها ۶۸
- شکل ۸-۵ : برطرف نمودن مشکل همپوشانی ۶۸
- شکل ۹-۵ : شماره کلاس مطابق جدول ۱-۴ ۶۹

- شکل ۵-۱۰ : حذف نقاط، محاسبه لبه بالایی، تعیین مینیمم‌های محلی کاندید برای جداسازی ۷۰
- شکل ۵-۱۱ : عدم وجود مینیمم محلی بین حرف "ح" و "ک" به دلیل همپوشانی ۷۰
- شکل ۵-۱۲ : جداسازی نامناسب حرف ی به دلیل عدم وجود محدوده مناسب ۷۱
- شکل ۵-۱۳ : مشخص نمودن محدوده مجاز برای محاسبه پروفایل بالایی و محل جداسازی ۷۱
- شکل ۵-۱۴ : محاسبه پروفایل بالایی تصویر ۷۲
- شکل ۵-۱۵ : بهینه نمودن مینیمم‌های محلی ۷۲
- شکل ۵-۱۶ : تعیین نقاط جداسازی اولیه ۷۳
- شکل ۵-۱۷ : برطرف کردن مشکل همپوشانی ۷۴
- شکل ۵-۱۸ : جداسازی نادرست حرف سین در کلمه بستنی ۷۵
- شکل ۵-۱۹ : برچسب‌زنی مؤلفه و مشخص نمودن تعداد و محل دقیق نقاط ۸۴
- شکل ۵-۲۰ : کلاس نهایی حروف ۸۴
- شکل ۵-۲۱ : شماره کلاس نهایی حروف با وجود نقاط ۹۰
- شکل ۵-۲۲ : انجام عملیات پیش پردازش ۹۰
- شکل ۵-۲۳ : جداسازی اولیه زیرکلمات با وجود مشکل همپوشانی ۹۰
- شکل ۵-۲۴ : برطرف نمودن مشکل همپوشانی ۹۰
- شکل ۵-۲۵ : شناسایی اولیه و ارسال زیرکلمات کلاس ۴۰ برای جداسازی ۹۰
- شکل ۵-۲۶ : جداسازی زیرکلمات و کلاسه بندی بدنه حروف ۹۰
- شکل ۵-۲۷ : کلاسه‌بندی نهایی حروف ۹۱
- شکل ۵-۲۸ : سند تولید شده نهایی با استفاده از کلاسه‌بندی نهایی ۹۱
- شکل ۵-۲۹ : یک تصویر نمونه برای ارزیابی الگوریتم ۹۱
- شکل ۵-۳۰ : شناسایی جملات و جداسازی آن‌ها ۹۲
- شکل ۵-۳۱ : نتیجه شناسایی برای حالتی که جداسازی از یک سمت انجام می‌شود ۹۲

- شکل ۳۲-۵: نتیجه شناسایی برای حالتی که جداسازی از هر دو سمت انجام می‌شود ۹۳
- شکل ۳۳-۵: نتیجه نهایی شناسایی پس از حذف نقاط اضافی تولید شده ۹۴
- شکل ۳۴-۵: تصویر دوم، دارای چرخش اولیه و نویز ۹۴
- شکل ۳۵-۵: حذف نویز و چرخش و استخراج جملات تصویر دوم ۹۴
- شکل ۳۶-۵: نتیجه شناسایی تصویر دوم، جداسازی از یک سمت ۹۵
- شکل ۳۷-۵: نتیجه شناسایی تصویر دوم، جداسازی از هر دو سمت ۹۵
- شکل ۳۸-۵: نتیجه شناسایی تصویر دوم، جداسازی از هر دو سمت با تصحیح نقاط اضافی ۹۵
- شکل ۳۹-۵: نتیجه شناسایی تصویر سوم، جداسازی از هر دو سمت ۹۶
- شکل ۴۰-۵: حذف نویز و چرخش و استخراج جملات تصویر سوم ۹۶
- شکل ۴۱-۵: نتیجه شناسایی تصویر سوم، جداسازی از یک سمت ۹۷
- شکل ۴۲-۵: نتیجه شناسایی تصویر سوم، جداسازی از هر دو سمت ۹۷
- شکل ۴۳-۵: نتیجه شناسایی تصویر سوم، جداسازی از هر دو سمت با تصحیح نقاط اضافی ۹۷
- شکل ۴۴-۵: نتیجه شناسایی تصویر چهارم، جداسازی از هر دو سمت ۹۸
- شکل ۴۵-۵: حذف نویز و چرخش و استخراج جملات تصویر چهارم ۹۸
- شکل ۴۶-۵: نتیجه شناسایی تصویر چهارم، جداسازی از یک سمت ۹۹
- شکل ۴۷-۵: نتیجه شناسایی تصویر چهارم، جداسازی از هر دو سمت ۹۹
- شکل ۴۸-۵: نتیجه شناسایی تصویر چهارم، جداسازی از هر دو سمت با تصحیح نقاط اضافی ۹۹
- شکل ۴۹-۵: نتیجه شناسایی تصویر پنجم، جداسازی از هر دو سمت ۱۰۰
- شکل ۵۰-۵: حذف نویز و چرخش و استخراج جملات تصویر پنجم ۱۰۰
- شکل ۵۱-۵: نتیجه شناسایی تصویر پنجم، جداسازی از یک سمت ۱۰۱
- شکل ۵۲-۵: نتیجه شناسایی تصویر پنجم، جداسازی از هر دو سمت ۱۰۱
- شکل ۵۳-۵: نتیجه شناسایی تصویر پنجم، جداسازی از هر دو سمت با تصحیح نقاط اضافی ۱۰۱

فهرست جدول‌ها

صفحه	عنوان
۵۴	جدول ۱-۴ : کلاس‌های مورد استفاده در طبقه‌بند حروف گسسته‌ی فارسی.....
۵۵	جدول ۲-۴ : کلاس‌های مورد استفاده در طبقه‌بند حروف پیوسته فارسی.....
۷۵	جدول ۱-۵ : شناسایی و اعتبارسنجی نقاط جداسازی (مرحله اول).....
۷۶	جدول ۲-۵ : شناسایی و اعتبارسنجی نقاط جداسازی (مرحله دوم).....
۷۷	جدول ۳-۵ : شناسایی و اعتبارسنجی نقاط جداسازی (مرحله سوم).....
۷۸	جدول ۴-۵ : نتیجه کلی جداسازی کلمه بستنی.....
۷۹	جدول ۵-۵ : مرحله اول جداسازی از دو سمت کلمه بستنی.....
۸۰	جدول ۶-۵ : مرحله دوم جداسازی از دو سمت کلمه بستنی.....
۸۰	جدول ۷-۵ : مرحله سوم جداسازی از دو سمت کلمه بستنی.....
۸۱	جدول ۸-۵ : مرحله چهارم جداسازی از دو سمت کلمه بستنی.....
۸۲	جدول ۹-۵ : کلاس‌های بدنه حروف گسسته.....
۸۳	جدول ۱۰-۵ : کلاس‌های بدنه حروف پیوسته.....
۸۵	جدول ۱۱-۵ : کلاس‌های نهایی حروف گسسته.....
۸۷	جدول ۱۲-۵ : کلاس‌های نهایی حروف پیوسته.....
۹۳	جدول ۱۳-۵ : حذف نقاط اضافی در زمان جداسازی.....

فصل اول

مقدمه

۱-۱- مقدمه

در گذشته تنها وسیله ارتباطی بین مردم برای ثبت وقایع یا اطلاع‌رسانی، کاغذ بوده است. در نتیجه نسخ کاغذی بسیاری وجود داشت که مدیریت، بایگانی و جستجو در میان آن‌ها کار دشواری بود. علاوه بر این بسیاری از نسخ مهم و ارزشمند به دلایل مختلفی از جمله حوادث طبیعی از بین می‌رفتند و امکان بازیابی مجدد آن‌ها وجود نداشت. کم‌کم با پدید آمدن رایانه و آشنایی با قابلیت‌های این محیط الکترونیکی لزوم استفاده از سیستمی که بتواند اطلاعات کاغذی را به رایانه منتقل کرده و در آن‌ها امکان ویرایش به وجود بیاورد روز به روز بیشتر شد.

امروزه نیز مستندات کاغذی به دو صورت چاپی (شامل روزنامه‌ها، کتاب‌ها، مقالات و...) و دست‌نویس وجود دارد که تنوع فراوان، موارد استفاده، لزوم بایگانی و تهیهی آرشیوها باعث انجام تحقیقات گسترده‌ای در زمینه بازشناسی متون شده است.

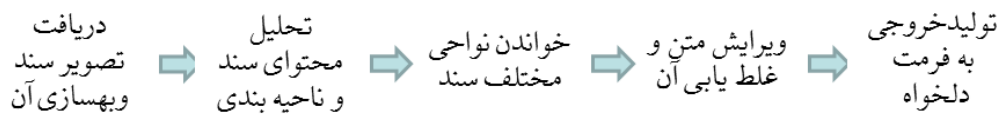
در این فصل ابتدا به معرفی سیستم بازشناسی نوری حروف پرداخته می‌شود و سپس ضمن بیان کاربردهای آن، تقسیم‌بندی این سیستم‌ها از نظر نوع ورودی و نحوه‌ی ورود اطلاعات مطرح شده و در آخر ویژگی‌های زبان فارسی و محدودیت‌های آن بیان می‌گردد.

۱-۱-۱- معرفی بازشناسی نوری حروف

ا.سی.آر مخفف بازشناسی نوری نویسه‌ها و یا نویسه‌خوان نوری است. یک روبشگر، تصویر را برای سامانه‌ی ا.سی.آر فراهم می‌کند، سپس این سامانه نواحی مختلف تصویر را شناسایی کرده و به فرمت مناسب ذخیره می‌نماید. در نتیجه یک فایل تصویری، به متن قابل ویرایش و جستجو تبدیل شده و حجم آن نیز کاهش می‌یابد. در فرآیند ا.سی.آر ورودی، تصویر سند مورد نظر و خروجی آن فایل دیجیتالی قابل ویرایش از اطلاعات سند مثلاً به فرمت وورد^۱ می‌باشد. شکل ۱-۱ فرآیند ا.سی.آر را با

¹ Word

جزئیات بیشتری نشان می‌دهد. در این شکل مراحل مختلف اُ.سی.آر تشریح شده است.



شکل ۱-۱: مراحل مختلف اُ.سی.آر

در این فرایند پس از دریافت تصویر سند، قسمت‌های مختلف تصویر از نظر معنایی و فیزیکی ناحیه‌بندی شده، بعد از شناسایی نواحی مختلف، ویرایش و غلط‌یابی صورت گرفته و سپس خروجی به فرمت دلخواه تولید می‌شود. در زمینه‌ی متون چاپی در زبان لاتین تحقیقات زیادی انجام شده ولی به علت تفاوت‌های عمده‌ای که در ساختار نوشتاری فارسی با لاتین وجود دارد (مثلاً در مراحل ماند تشخیص قلم^۱ و جداسازی) نمی‌توان آن نتایج را مستقیماً برای زبان فارسی به کار برد.

۱-۲- کاربردهای اُ.سی.آر

در ادامه برخی از کاربردهای سیستم‌های بازشناسی متن مطرح می‌شود.

۱-۲-۱ - دستگاه‌های چک‌خوان

از جمله کاربردهای سیستم‌های بازشناسی متن، دستگاه‌های چک‌خوان هستند که در بانک‌ها برای ورود داده استفاده می‌شوند. این سیستم‌ها برای خواندن بخش‌های کوچکی از تصویر چک شامل نویسه‌های چاپی و دست‌نویس مانند شماره حساب، امضای مشتری، شماره برگ چک و مبلغ چک طراحی می‌شوند. ساختار کاغذ چک معمولاً ثابت است و نواحی مشخصی از آن توسط ماشین خوانده می‌شود. دستگاه‌های چک‌خوان در حدود ۱۵۰۰۰ چک را در یک ساعت می‌خوانند و نرخ خطای آن‌ها کم است. این سیستم‌ها در کیفیت‌های پایین چاپ نیز به خوبی عمل می‌کنند.

^۱ Font

۱-۲-۲- ورود متن در سیستم‌های اتوماسیون اداری

از جمله نیازهای اساسی در سیستم‌های اتوماسیون اداری دسترسی به امکان ویرایش، اصلاح و جستجو کردن در سندهای اسکن شده می‌باشد که سیستم بازشناسی متن نقش مهمی در این امر دارد. این سیستم‌ها اسناد با قلم‌های محدود و کیفیت چاپ مشخصی را بازشناسی می‌کنند. این سیستم‌ها برای ورود حجم زیادی از متن به شکل قابل پردازش با نرم‌افزارهای پردازشگر متن استفاده می‌شوند. قابلیت این سیستم‌ها در ورود متن، رقیب قدرتمندی برای وارد کردن دستی اطلاعات است. کارایی این سیستم‌ها به کیفیت چاپ اسناد وابستگی زیادی دارد.

۱-۲-۳- خودکارسازی فرایند

در این کاربرد، کنترل یک فرایند مشخص، به بازشناسی دقیق متن نوشته شده اولویت دارد. به عنوان مثال در خودکارسازی فرایند جداسازی نامه در مراکز پستی، بازشناسی نام شهر مقصد، از بازشناسی دقیق کل متن نوشته شده بر روی پاکت نامه اولویت بالاتری دارد. نرخ بازشناسی این سیستم‌ها کاملاً به کیفیت تصویر پاکت نامه وابسته است. سرعت جداسازی نامه‌ها بر حسب آدرس، معمولاً در حدود چند ده هزار نامه در ساعت است.

۱-۲-۴- بازیابی اسناد^۱

هزینه‌ی بالای کار با پایگاه‌های داده بزرگ از تصاویر اسناد، نیاز به روش‌های خودکار و کارآمد برای دستیابی به اطلاعات داخل این تصاویر را به وجود آورده است. در تلاش برای حرکت به سوی ادارات بدون کاغذ، حجم بزرگی از اسناد کاغذی روبش شده و به صورت تصویری ذخیره می‌شوند، بدون اینکه برچسب اطلاعاتی مناسبی به آن‌ها زده شود یک راه برای برچسب‌زنی و بازیابی پایگاه داده، تبدیل کامل اسناد به شکل الکترونیکی است که می‌توانند به طور خودکار برچسب زده شوند. موانعی مثل

1 Document Retrieval

کیفیت پایین تصویر و هزینه بالای تبدیل در مقابل این روش وجود دارد. از طرفی چون بسیاری از اجزاء غیر متن را نمی‌توان به طور کامل در شکل تبدیل شده نمایش داد، داشتن یک کپی از سند در شکل تصویری لازم به نظر می‌رسد. نگهداری و دسترسی به یک پایگاه داده متنی مسئله‌ی مهمی است. در چنین پایگاهی، مستند جدیدی که وارد می‌شود لزوماً ساختاری مشابه مستندات قبلی ندارد، بنابراین تعریف ساختارهای مختلف از قبل، غیر ممکن است. توقعی که کاربران از یک سیستم پایگاه داده مستندات دارند، بازیابی دقیق مستند مورد نظر نیست، اما انتظار می‌رود که مکانیزمی برای بازیابی مستندات، به صورت مرتب شده بر اساس شباهت به مستند مورد پرس و جو، وجود داشته باشد. برای بازیابی اسناد مبتنی بر محتوا، می‌توان بلوک‌های متن تصویر سند را بازشناسی کرده و از آن‌ها برای بازیابی سند استفاده کرد [۱].

۱-۲-۵- کتابخانه دیجیتال

در یک کتابخانه‌ی دیجیتال برای دسترسی به کتاب مورد نظر از کلید واژه‌ها استفاده می‌شود. کلید واژه‌ها به صورت متن وارد می‌شوند و نرم افزار کتابخانه، کتاب‌های مرتبط با آن‌ها را بازیابی می‌کند. پیدا کردن کلمات کلیدی از تصاویر رویش شده‌ی کتاب‌ها، هسته‌ی اصلی نرم افزار جستجوی کتابخانه است.

۱-۲-۶- کتاب‌های الکترونیکی

با فناوری ا.سی.آر تصاویر صفحات کتاب به متن رایانه‌ای تبدیل شده و امکان ویرایش متن، جست‌وجوی متن، چاپ مجدد متن با کیفیت بالا و انتقال الکترونیکی سریع آن فراهم می‌شود. ا.سی.آر به ویژه برای دیجیتالی نمودن منابع و نسخ خطی بسیار مفید است، چرا که برای نسخ خطی، باید اصالت منبع با حفظ شکل و قالب اصلی متون و نه به صورت تایپ شده حفظ شود.

۱-۲-۷- کنترل ترافیک

برای تشخیص پلاک خودروها می‌توان از سیستم بازشناسی متن استفاده کرد. صحنه‌ی ترافیک

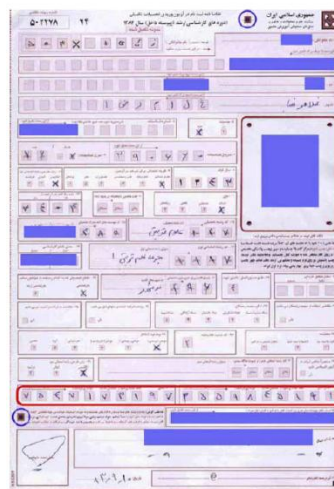
معمولاً با دوربین‌های سریع تصویربرداری می‌شود. از محدود بودن حروف استفاده شده در پلاک‌ها، می‌توان برای بهبود نتایج استفاده کرد. پس از جدا شدن تک تک عناصر موجود در پلاک، عناصر به صورت مجزا به سیستم بازشناسی متن جهت خواندن ارجاع داده شده و خوانده می‌شود.

۱-۲-۸- تبدیل متن به صحبت

از یک سیستم بازشناسی متن در کنار یک سیستم پردازش صوت می‌توان در کاربردهای مختلف مانند کمک به نابینایان به منظور درک متون چاپی و خودکارسازی پاسخ به ارباب رجوع در ادارات استفاده کرد.

۱-۲-۹- فرم خوان‌ها

این سیستم‌ها همان سیستم شناسایی متون دست‌نویس گسسته است که برای خواندن فرم‌های خاص طراحی می‌شوند. این سیستم‌ها معمولاً از روبشگرهایی که رنگ خاصی را حذف می‌کنند استفاده کرده و به هنگام اسکن رنگ زمینه را حذف می‌کنند. حروف و ارقام به صورت مجزا در خانه‌هایی با رنگ مورد نظر نوشته می‌شوند و روبشگر هنگام روبش فقط این نواحی را اسکن کرده و به صورت خودکار رنگ زمینه‌ی این خانه‌ها را حذف می‌کند. بنابراین این سیستم‌ها نرخ بازشناسی بالایی دارند. در شکل ۲-۱ نمونه‌ای از یک فرم ثبت نام کنکور کارشناسی ارشد ارائه شده است.



شکل ۲-۱: نمونه‌ای از یک فرم ثبت نام کنکور کارشناسی ارشد [۲]

۱-۲-۱۰- نقشه خوانی خودکار

بازشناسی نویسه‌های روی نقشه‌ها با مسائلی مانند ادغام حروف با خطوط نقشه، چاپ شدن متن با زوایای مختلف، متون با قلم‌های متعدد، متون چاپ شده بر زمینه رنگی و وجود متن دست‌نویس روبرو است.

۱-۲-۱۱- نت خوانها

بازشناسی نت‌های موسیقی یکی دیگر از کاربردهای سیستم‌های بازشناسی متن است. نت‌های موسیقی به دو صورت چاپی و دست‌نویس وجود دارند.

۱-۳-۱- انواع ا.سی. آر

سیستم‌های بازشناسی عمدتاً از دو منظر نوع ورودی و نحوه‌ی ورود اطلاعات مورد بررسی و طبقه‌بندی قرار می‌گیرند.

۱-۳-۱- از نظر نوع ورودی

در یک تقسیم‌بندی کلی می‌توان سیستم‌های ا.سی. آر را از لحاظ نوع الگوی ورودی به دو گروه اصلی تقسیم کرد:

الف - سیستم‌های بازشناسی متون چاپی

ب - سیستم‌های بازشناسی متون دست‌نویس

در سیستم‌های بازشناسی متون چاپی، ورودی به صورت تصویری است که از دوربین دیجیتال یا اسکنر به دست آمده است. ولی سیستم‌های شناسایی متون دست‌نویس شامل تصویر متون نوشته شده است که این متون به دو بخش گسسته و پیوسته تقسیم می‌شود. متون دست‌نویس گسسته شامل متونی است که در آن حروف جدا از هم نوشته می‌شوند مانند حروفی که در فرم‌های آزمون‌ها در کادرهایی به صورت مجزا نوشته می‌شود، ولی متون دست‌نویس پیوسته در قالب نامه یا متن یا گزارش

توسط شخص روی کاغذ نوشته می‌شود. متون دست‌نویس پیوسته به علت مشکلاتی از قبیل کشیده یا شکسته نوشته شدن بعضی از حروف یا به طور کلی نوشتاری که مطابق قواعد نگارش رسم‌الخط فارسی نیست، هنوز به صورت یک سیستم نرم افزاری جامع درنیامده است و همچنان تلاش در ایجاد سیستم شناسایی متون دست‌نویس پیوسته با دقت قابل قبول وجود دارد. ولی برای شناسایی متون دست‌نویس گسسته به علت گسسته بودن حروف و همچنین متون تایی به علت اینکه عموماً از قواعد خاصی پیروی می‌کنند سیستم‌های نرم افزاری زیادی وجود دارند، که بسیاری از آنها با دقت بالایی عملیات شناسایی را انجام می‌دهند.

۱-۳-۲- تقسیم‌بندی از لحاظ نحوه ورود اطلاعات

از جنبه‌ی نحوه‌ی ورود اطلاعات، سیستم‌های ا.سی.آر به دو دسته‌ی سیستم‌های برخط و سیستم‌های برون خط تقسیم‌بندی می‌شوند.

۱-۳-۲-۱ سیستم‌های برخط^۱

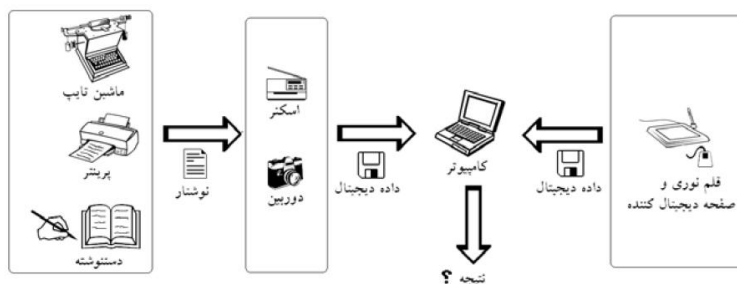
ورودی سیستم بازشناسی برخط، به صورت پیوسته همزمان با نوشتن نویسنده بر روی یک صفحه رقمی‌کننده به سیستم وارد می‌شود. بازشناسی برخط فقط در بازشناسی دست‌نوشته کاربرد دارد. بازشناسی برخط نوشتار به دلیل راحت‌تر بودن نوشتن از تایپ کردن، عدم امکان تایپ در بعضی مکان‌ها، عدم وجود یک صفحه کلید کامل روی رایانه‌های کوچک و سخت بودن تایپ نویسه‌ها در بعضی زبان‌ها به دلیل تعداد زیاد آنها، مورد توجه خاصی قرار گرفته است. در روش‌های برخط از داده‌های یک بعدی استفاده می‌شود. در این سیستم‌ها می‌توان قسمت‌های مختلف متن را ذخیره کرده و سپس عملیات شناسایی را توسط روش‌های برون خطی انجام داد. برای بازشناسی دست‌نوشته‌ی برخط

1 Online

از ویژگی‌هایی مانند مختصات مکانی حرکت قلم روی صفحه، سرعت نوشتن و میزان فشار قلم بر صفحه استفاده می‌شود.

۱-۳-۲ سیستم‌های برون خط^۱

در سیستم بازشناسی برون خط، تصویر اسکن شده به عنوان ورودی در نظر گرفته می‌شود. بازشناسی برون خط هم شامل بازشناسی نوشتار چاپی و هم نوشتار دست‌نویس می‌شود. روش‌های شناسایی برون خط شامل کاربرد شبکه عصبی، منطق فازی و غیره می‌باشد که بر روی تصاویر دوبعدی اعمال می‌شوند. سیستم برون خط دقت کمتری از سیستم‌های برخط دارند. نرخ بازشناسی برای سیستم‌های برخط بالاتر از این نرخ در سیستم برون خط گزارش شده است.



شکل ۱-۳: سیستم‌های شناسایی برون خط و برخط

۱-۴- موانع موجود بر سر تکامل ا.سی.آر فارسی

در دهه‌های اخیر تحقیقات زیادی در زمینه‌ی بازشناسی ارقام و حروف به زبان‌های مختلف، و از جمله زبان فارسی انجام شده است. بیشتر تحقیقات انجام شده در زمینه‌ی بازشناسی متون چاپی فارسی با محدودیت‌هایی از جمله درجه تفکیک زیاد (۴۰۰ و بیشتر)، استفاده از تصاویر متنی تمیز و غیر واقعی و شناسایی متن با چند قلم خاص مشابه همراه بوده است، هر چند در سال‌های اخیر تا حد امکان بسیاری از این محدودیت‌ها برطرف شده است. پیچیدگی‌های زبان فارسی، عدم سرمایه‌گذاری

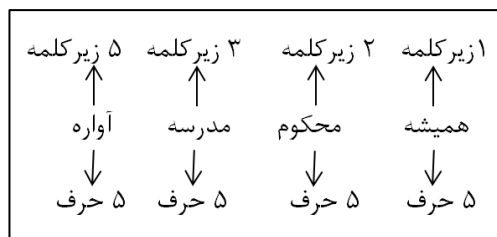
1 Offline

کافی، عدم آشنایی شرکت‌ها با ا.سی.آر، تنوع قلم‌ها، فقدان لغت نامه و پایگاه داده‌ی جامع و استاندارد و وجود جمعیت کم کاربر زبان فارسی از جمله عواملی هستند که مانع ایجاد سیستم جامع بازشناسی متون فارسی می‌شود.

۱-۵- معرفی ویژگی‌های خط فارسی

آشنایی با قواعد نگارش فارسی در انتخاب روش‌های مناسب برای بازشناسی متون نقش اساسی دارد. در زیر قواعد مهم نگارش فارسی بیان می‌شود.

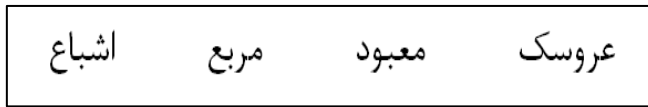
- خط فارسی بر خلاف لاتین از راست به چپ نوشته می‌شود. البته اعداد فارسی از چپ به راست نوشته می‌شوند.
- هر کلمه از یک یا چند زیرکلمه و هر زیرکلمه از یک یا چند حرف متصل به هم تشکیل می‌شود. مطابق قواعد نگارش زبان فارسی، باید بین کلمات فاصله وجود داشته باشد و بین زیرکلمات فاصله اضافه وجود نداشته باشد، که در بسیاری از موارد این قواعد رعایت نمی‌شود.



شکل ۱-۴: آشنایی با مفهوم حرف و زیرکلمه

- در کلمات فارسی، بیشتر حروف از دو طرف به حروف مجاور خود می‌چسبند. با این حال حروف {ا، د، ذ، ر، ز، ژ، و} به حرف بعدی خود نمی‌چسبند و فقط می‌توانند از سمت راست به حرف قبلی خود وصل شوند.
- حروف فارسی بسته به موقعیت آن‌ها در زیرکلمه می‌توانند تا چهار شکل متفاوت داشته

باشند: حروف ابتدایی، میانی، انتهایی و مجزا. برای مثال شکل‌های مختلف حرف "ع" در شکل ۵-۱ نشان داده شده است.

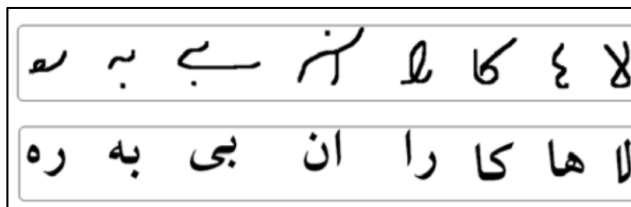


شکل ۵-۱: شکل‌های مختلف یک حرف در کلمه

- بعضی از حروف فارسی نقطه دارند که تعداد نقاط این حروف بین یک تا سه متغیر است. نقاط ممکن است در بالا "ز"، پایین "ب" یا داخل بدنه‌ی "ج" حرف قرار گیرند. حروف "ک" و "گ" سرکش دارند. روی حرف الف علامت مد "آ" نیز قرار می‌گیرد. حروف "ط" و "ظ" نیز دسته دارند.
- قسمتی از حرف یا زیرکلمه که نقاط آن حذف شده باشد را بدنه‌ی حرف یا زیرکلمه می‌گویند. برخی از حروف مانند "ب پ ت ث" دارای بدنه‌ی یکسان هستند و تفاوت آن‌ها در تعداد نقاط و محل قرار گرفتن آن می‌باشد.
- اندازه‌ی تمام حروف فارسی یکسان نیستند. مثلاً حروف (ب) و (س) در حالت جدا، اندازه‌ی بزرگ‌تری نسبت به حروف (د) و (ه) دارند. این تنوع در اندازه‌ی حروف، کار قطعه‌بندی حروف را مشکل می‌کند.
- در برخی از شیوه‌های نوشتاری زبان فارسی، ترکیب شدن دو یا چند حرف کنار هم، شکلی ایجاد می‌کند که شباهتی به حروف تشکیل دهنده‌ی آن ندارد. به این ترکیب، حروف ادغام شده^۱ می‌گویند. چنین مواردی نه‌تنها در نوشتار دست‌نویس، بلکه در متون تایپی در بعضی از قلم‌های خاص نیز وجود دارد. متداول‌ترین ترکیب در متون تایپی، ادغام دو حرف "ل" و "ا" به صورت "لا" است. در نوشته‌های دست‌نویس فارسی نیز با

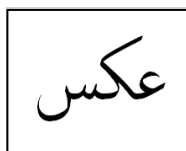
1 Ligature

توجه به سلیقه‌ی نویسنده شکل بعضی از حروف کنار هم، به کلی تغییر می‌کند (شکل ۶-۱).



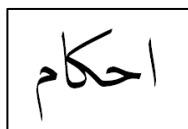
شکل ۶-۱: نمونه‌هایی از ادغام حروف مجاور در متون دست‌نویس

- حروف واقع در یک کلمه ممکن است همپوشانی داشته باشند، در این صورت نمی‌توان با رسم خطوط عمودی، حروف را به صورت کامل از یکدیگر مجزا نمود.



شکل ۷-۱: همپوشانی حروف

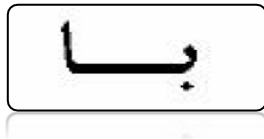
- در برخی از قلم‌ها بعضی از حروف، در دو محل (از یک سمت) به یکدیگر اتصال دارند.



شکل ۸-۱: اتصال حروف از دو محل

- حروف فارسی می‌توانند در بالا یا پایین بدنه‌ی خود دارای اعراب باشند. سه اعراب -َ -ِ -ُ در زبان فارسی اعراب‌های اصلی بوده و اعراب - َ در برخی کلمات عربی رایج در زبان فارسی دیده می‌شود. کلمات عربی دارای اعراب - و - ِ در زبان فارسی عمومیت نیافته‌اند. هر چند کاربرد اعراب در زبان فارسی نسبت به زبان عربی بسیار محدودتر است، اما در مواردی که کلمه‌ای نامتداول باشد و یا به دلیل تشابه نگارشی آن با کلمه‌ی دیگر، تأکید بر تلفظ صحیح آن باشد، مورد استفاده قرار می‌گیرند.

- در بالای بدنه‌ی یک حرف ممکن است علامت تشدید وجود داشته باشد.
- برخی از حروف شامل همزه هستند ("ئا"، "ا"، "ؤ").
- حروفی که از طرف چپ قابلیت اتصال به حرف مجاور خود را دارند، می‌توانند به صورت کشیده نوشته شوند.



شکل ۹-۱: اتصال حروف از سمت چپ

موارد مذکور نمایانگر پاره‌ای از مشکلات است که عمل بازشناسی را با محدودیت‌های زیادی مواجه می‌کند.

۱-۶- نتیجه‌گیری

در این فصل لزوم استفاده از سیستم‌های بازشناسی متن مطرح گردید و پس از تعریف مختصری درباره‌ی سیستم‌های بازشناسی متن، کاربردهای آن بیان شد. همچنین تقسیم‌بندی سیستم‌های بازشناسی متن از نظر نوع ورودی و نحوه‌ی ورود اطلاعات معرفی گردید. از طرفی با توجه به این که این پایان‌نامه در رابطه با بازشناسی متون فارسی است، به بررسی ویژگی‌های متون فارسی پرداخته شد و عوامل پیچیدگی بازشناسی زبان فارسی نیز ذکر گردید.

در فصل آتی به بررسی روند پیدایش سیستم‌های بازشناسی متن و تحقیقات انجام شده در زمینه‌ی ا.سی.آر پرداخته خواهد شد.

فصل دوم

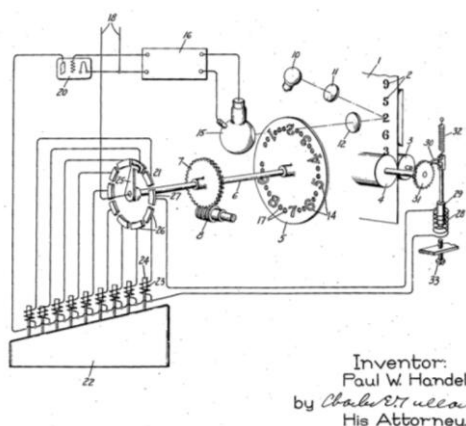
کارهای پیشین

۲-۱- مقدمه

در این فصل با آگاهی از مقدمات اولیه راجع به اُسی.آر و بیان ضرورت استفاده از آن در بسیاری از مواقع، به مروری بر سیر تحولی اُسی.آر پرداخته می‌شود که در آن نسل‌های مختلف اُسی.آر معرفی شده و روند تحقیقاتی آن معرفی می‌گردد. از آن جا که موضوع مورد بحث در مورد سیستم‌های بازشناسی متن به زبان فارسی است، به بررسی تحقیقاتی در زمینه‌ی اُسی.آر فارسی پرداخته می‌شود.

۲-۲- روند تحقیقاتی اُسی.آر

در ۱۹۰۰ میلادی، تیورین^۱، دانشمند روسی، وسیله‌ای برای خواندن متن برای نابینایان به ثبت رساند [۳]. در سال ۱۹۲۹ میلادی تاشک^۲ در آلمان و در سال ۱۹۳۲ هندل^۳ و همکارانش در آمریکا سامانه‌ی اُسی.آر خود را در زمینه‌ی سیستم بازشناسی حروف ثبت کردند (شکل ۱-۲) [۴]. روش آن‌ها، تطبیق الگو نام دارد و به این صورت کار می‌کند که به هر حرف، نور تابیده می‌شود و نور باز تابیده شده از حروف، از قالب‌های مکانیکی عبور داده می‌شود. هرگاه نوری از قالب عبور نکرد، حرف تشخیص داده می‌شود. این اختراع به دلیل تکنولوژی اپتومکانیکی مورد استفاده در آن کاربردی نبود.



شکل ۱-۲: سامانه اُسی.آر هندل [۵]

- 1 Tyurin
- 2 Tausheck
- 3 Handel

در دهه‌ی ۱۹۴۰ میلادی، اولین سامانه‌های کامپیوتری بازنشاسی نویسه معرفی شدند. این سامانه‌ها به بازنشاسی متون چاپی یا ارقام و حروف دست‌نویس مقید^۱ می‌پرداختند.

در دهه‌ی ۱۹۵۰، سیستم‌های تجاری ا.سی.آر برای خواندن ارقامی که با قلم خاص چاپ می‌شدند، به بازار آمدند. همچنین در همین دهه، اولین سامانه‌های تجاری برای بازنشاسی نویسه‌های برخط، که با قلم و صفحه الکترونیکی خاص نوشته می‌شدند، مطرح شدند [۳].

در سال ۱۹۵۹ میلادی اولین ا.سی.آر برای خواندن یک صفحه متن چاپی به بازار آمد که می‌توانست متون نوشته شده با یک قلم خاص در یک اندازه‌ی خاص را بخواند [۶].

در دهه‌ی ۱۹۷۰ [۴] تمرکز بیشتر مقالات روی روش‌های استخراج ویژگی از نظر یافتن ویژگی‌های مناسب و استخراج آن‌ها به منظور تمایز بین کلاس‌ها بود که به کلاس‌های ارقام، حروف و علائم محدود می‌شدند. در آن زمان به علت محدود بودن امکانات پیاده‌سازی، استفاده از تعداد زیادی ویژگی میسر نبود و امکانات محدود پردازشی، تأثیری تعیین‌کننده بر روش‌های استخراج ویژگی داشت.

در دهه‌ی ۱۹۸۰ تحقیقات ا.سی.آر به بازنشاسی حروف و ارقام دست‌نویس غیرمقید و حتی متون دست‌نویس گسترش یافت و نتایج خوبی بدست آمد. تا چند دهه، تحقیقات تنها بر زبان‌های لاتین متمرکز بود و از اوایل دهه‌ی ۱۹۸۰ پژوهش‌ها در زمینه‌ی عربی و فارسی نیز آغاز شد. همچنین روش‌های طبقه‌بندی جدیدی مثل شبکه‌های عصبی و روش‌های فازی بکار گرفته شدند.

در دهه‌ی ۱۹۹۰ موضوعات تازه‌ای در پژوهش‌های مربوط به ا.سی.آر مطرح شد. از جمله‌ی آن‌ها تحلیل سند است، به گونه‌ای که سیستم بازنشاسی متن می‌توانست ساختار فیزیکی صفحه را بیابد و صفحاتی با ساختارهای پیچیده که شامل عکس، گرافیک و جدول هستند، به عنوان ورودی بپذیرد. از دیگر موضوعات، ایده‌ی استفاده از شکل کلی کلمات و بازنشاسی با استفاده از آن‌ها در تصاویری بود که در مرحله‌ی جداسازی و شکستن کلمات به حروف مشکل جدی داشتند. موضوع دیگر، بازنشاسی

1 Constrained

متون چاپی بدون توجه به نوع قلم آن‌ها بود که تا حد زیادی برای قلم‌های گوناگون میسر شد. در این دهه کارایی سامانه‌های ا.سی.آر به حد قابل استفاده در کاربردهای واقعی رسید. روش‌های طبقه‌بندی توانمندتری به کار گرفته شدند و همچنین موضوع ترکیب نتایج حاصل از چند طبقه‌بند ساده به جای استفاده از یک طبقه‌بند قوی مطرح شد.

در دهه‌ی ۲۰۰۰ میلادی پیشرفت‌های زیادی حاصل شد. از جمله ترکیب طبقه‌بندها که در دهه‌ی گذشته مطرح شده بود، با توجه به نتایج خوبی که از آن به دست آمد، مورد توجه بیشتری قرار گرفت. از دیگر موضوعات مطرح شده در این دهه می‌توان بهسازی و بازشناسی مستندات تاریخی و شناسایی افراد از روی دست‌نوشته‌های آنان را نام برد که پژوهش‌های زیادی در این زمینه انجام شده است. در زمینه‌ی بازشناسی متون، توجه بیشتری به بازشناسی گرافیک و همچنین نمادهای ریاضی شده است. در این دوره نتایج قابل توجهی درباره تشخیص ساختار یک صفحه از متن چاپی حاصل شد.

با گسترش کاربرد دوربین‌های دیجیتال و افزایش کیفیت تصاویر آن‌ها، استفاده از این دوربین‌ها به جای روبشگرها در بازشناسی مستندات مورد توجه قرار گرفت. تفاوت بین تصاویر حاصل از این دو دستگاه، موضوعات پژوهشی تازه‌ای را مطرح کرد. بازشناسی متن در تصاویر ویدیویی نیز مورد توجه قرار گرفت. با رواج استفاده از تصاویر رنگی، موضوع به‌کارگیری اطلاعات رنگ در بازشناسی متون نیز مطرح شد. بازیابی تصویری مستندات^۱ را می‌توان یکی از مهم‌ترین موضوعات مطرح در سال‌های اخیر دانست. در یک سامانه‌ی بازیابی تصویر مستندات، می‌توان با روش‌های مختلف پرس و جو، مستند مورد نظر کاربر را در یک بانک تصویری بزرگ پیدا کرد. سابقه این موضوع به بازیابی اطلاعات و بازیابی تصویر برمی‌گردد. بازیابی تصویر از مهم‌ترین موضوعات پردازش تصویر در سال‌های اخیر است. امروزه سیستم‌های ا.سی.آر قادر به تشخیص دقیق کاراکترهای تائپتی با انواع قلم‌ها و در اندازه‌های متفاوت هستند و شناسایی متون دست‌نویس لاتین و یا قلم‌هایی که در آن‌ها از خطوط خمیده استفاده

1 Document Image Retrieval

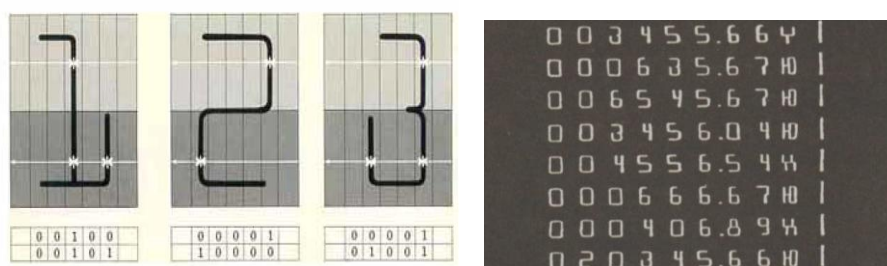
می‌شود (مثل فارسی و عربی)، رو به پیشرفت است.

۲-۲-۱- نسل‌های مختلف نویسه خوان‌های نوری

سیستم‌های تجاری عرضه شده در زمینه‌ی ا.سی.آر را می‌توان به چهار نسل تقسیم‌بندی نمود.

۲-۲-۱-۱ سیستم‌های نسل اول

این نویسه خوان‌ها در اوایل دهه‌ی ۱۹۶۰ میلادی به بازار آمدند. این سیستم‌ها فقط قابلیت تشخیص کاراکترهای خاص با اندازه و قلم مشخصی را داشتند. در آن‌ها از روش‌هایی استفاده شده بود که نسبت به تغییر جای کاراکترها، اندازه و دوران آن‌ها فوق‌العاده حساس بودند. بارزترین این نویسه خوان‌ها NCR420 بود که می‌توانست اعداد و پنج نماد دیگر را که با قلم خاص این نویسه خوان نوشته می‌شدند، بخواند. این قلم NOF یا Bicides نامیده می‌شد.



شکل ۲-۲: قلم NOF یا Bicides [۷]

۲-۲-۱-۲ سیستم‌های نسل دوم

این سیستم‌ها در اواسط دهه‌ی ۱۹۶۰ به بازار آمدند و فقط قادر به شناسایی دست‌نویس‌های مجزا مانند ارقام بودند. یک نویسه خوان شاخص از نسل دوم رتینا^۱ است. این سامانه می‌توانست مجموعه‌ی اعداد، پنج حرف بزرگ و دو علامت دست‌نویس را بشناسد. در مورد نویسه های چاپی، این سامانه

^۱ RETINA

می‌توانست ۴۰ نویسه را با قلم‌های مختلف بشناسد. البته لازم بود که کاربر نوع قلم را مشخص کند. سرعت خواندن این ماشین ۲۴۰۰ نویسه در ثانیه بود و علیرغم قیمت بالایش، کاربرد زیادی پیدا کرد.

۲-۲-۱-۳ سیستم‌های نسل سوم

مهم‌ترین مشخصه‌های ماشین‌های نسل سوم که در اواخر دهه‌ی ۱۹۶۰ میلادی مطرح شدند، پرداختن آن‌ها به نویسه‌های تایپی با کیفیت پایین و همچنین مجموعه‌های کامل نویسه‌های دست‌نویس لاتین بود. محصولات موفق تجاری این نسل، از سال‌های ۱۹۷۵ تا ۱۹۸۵ به بازار آمدند. سامانه‌های این نسل شامل سه بخش خواندن اسناد چاپی با کیفیت تصویری نامطلوب، خواندن اسناد حاوی نویسه‌های مجزای دست‌نویس و بسته‌های نرم افزاری بود.

۲-۲-۱-۴ سیستم‌های نسل چهارم

نویسه خوان‌های این نسل می‌توانند اسناد پیچیده شامل متن، جدول، تصویر، گرافیک و روابط ریاضی را پردازش کرده و ساختار اصلی آن‌ها را در خروجی حفظ کنند و همچنین اسناد رنگی و نویزی با کیفیت تصویری پایین را بازشناسی می‌کنند. این نسل تا به امروز هم ادامه دارد و فعالیت‌های گسترده‌ای در این زمینه در حال انجام است.

۲-۲-۲-۲ سیر تحول ا.سی.آر فارسی

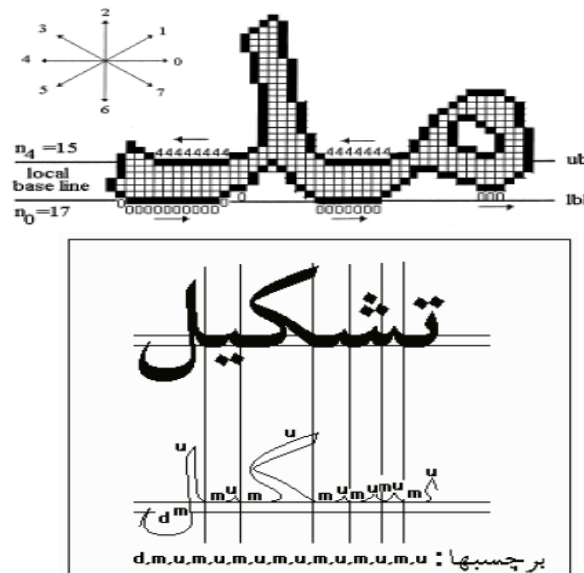
۲-۲-۲-۱ بررسی اجمالی برخی از تحقیقات

در سال‌های اخیر تلاش‌های قابل توجهی از سوی برخی شرکت‌های فعال در زمینه‌ی پردازش تصویر انجام شده که برخی از آن‌ها منجر به محصولات قابل قبولی شده است. برخی از تحقیقاتی که در زمینه‌ی ا.سی.آر فارسی انجام شده به شرح زیر می‌باشد.

در مورد بازشناسی متون چاپی فارسی مقاله‌ی [۸] روشی برای بازشناسی متون فارسی با قلم‌های درشت مانند تیرهای روزنامه‌ها ارائه کرده است.

در [۹] از ویژگی‌های شبه زرنیکی برای بازشناسی متون تاییبی و دست‌نویس با شبکه عصبی احتمالاتی استفاده شده است. این روش مبتنی بر شکل کلی زیرکلمات بوده و تعداد زیرکلمات ۱۰۰۰ فرض شده است. برای بازشناسی از شبکه عصبی استفاده شده و از ۲۸ ویژگی گشتاورهای شبه زرنیکی مرتبه ۴ برای توصیف شکل زیرکلمات استفاده شده است. با آزمایش این روش بر روی مجموعه‌ای شامل ۵۰۰ زیرکلمه دست‌نویس نوشته شده با خط نسخ، و یک صفحه متن چاپی در هر دو مورد، نرخ بازشناسی صحیح ۹۶٪ گزارش شده است.

عزمی یک الگوریتم جدید برای جداسازی زیرکلمات فارسی به حروف بر اساس برچسب‌زنی مقید به هر زیرکلمه ارائه کرده است. وی همچنین روشی برای میزان کردن خط پایه برای کانتور بالایی زیرکلمات ارائه کرده است [۱۰]. در این مقاله دقت جداسازی ۹۸/۹٪ گزارش شده است. برای بازشناسی از کدهای فریمن به عنوان ویژگی و از یک طبقه‌بند آماری برای طبقه‌بندی استفاده شده است. روش دومی که به کار برده شده استفاده از تبدیل هاف جهت استخراج ویژگی و شبکه عصبی برای طبقه‌بندی است.



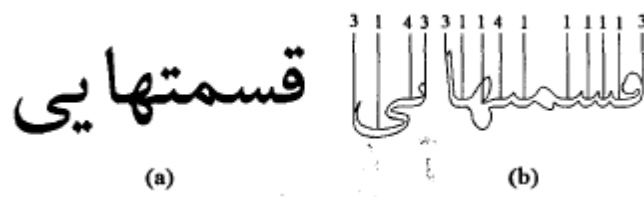
شکل ۲-۳: جداسازی بر اساس کانتور بالایی [۱۰]

عزمی همچنین از روش مبتنی بر شکل کلی کلمه برای بازشناسی استفاده کرده است. در این کار

تعداد ۵۷۳۷ زیر کلمه تهیه شده است که با حذف نقاط به ۲۷۹۰ زیر کلمه کاهش می‌یابند. در این

مقاله سه ویژگی مکان مشخصه، توصیفگر فوریه و کانتور بالایی آزمایش شده است [۱۱].

منه‌اج یک روش جداسازی و بازشناسی توأم برای متون چاپی فارسی ارائه کرده است. وی توصیفگرهای فوریه را به عنوان ویژگی و شبکه‌ی عصبی چند لایه را به عنوان طبقه‌بند استفاده کرده است [۱۲]. برای جداسازی از برجسب‌زنی به کانتور زیر کلمه استفاده شده است (شکل ۲-۴). این روش روی سه قلم تیترا، زر و تایمز نیو رومن آزمایش شده است. برای تصاویر تولیدی در کامپیوتر (کاملاً خالی از نویز) دقت ۱۰۰٪ گزارش شده است. برای متونی که هر سه قلم را داشته‌اند این دقت به ۹۳٪ کاهش یافته است. برای تصاویر اسکن شده با درجه تفکیک ۴۰۰، دقت ۹۰٪ گزارش شده است.

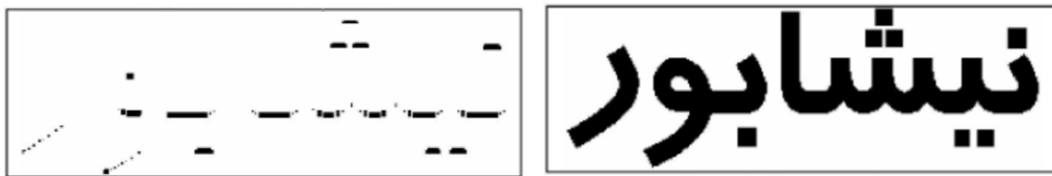


شکل ۲-۴: جداسازی بر مبنای برجسب‌زنی به کانتور زیر کلمه [۱۲]

در [۱۳] یک سیستم ا.سی.آر مبتنی بر جداسازی ارائه شده است. این سیستم شامل ماژول ناحیه بندی، جداسازی و بازشناسی است. برای جداسازی از اطلاعات کانتور بالای زیر کلمات استفاده شده است. برای بازشناسی از شبکه‌ی عصبی و SVM استفاده شده و ویژگی‌های مورد استفاده نیز یک سری ویژگی‌های آماری مثل افکنش عمودی، مشتق افکنش عمودی و فاصله‌ی حرف تا خط پایه است. برای بهبود نتایج بازشناسی از NLP استفاده شده است که در آن دقت ۹۸/۳٪ در سطح حرف و ۹۰/۱۷٪ در سطح کلمه گزارش شده است که این دقت با استفاده از NLP به ۹۴/۱۹٪ افزایش یافته است.

در [۱۴] از تبدیل موجک برای جداسازی زیر کلمات به حروف استفاده شده است. این تبدیل برای

آشکارسازی لبه‌های افقی پایین کلمات و یافتن خط پایه به کار رفته است (شکل ۲-۵).



شکل ۲-۵: یافتن لبه های افقی پایین کلمات [۱۴]

پس از یافتن لبه‌های پایین، از افکنش افقی روی آن‌ها استفاده شده تا موقعیت نقاط جداسازی مشخص شود. برای بازشناسی حروف پس از جداسازی، از شبکه‌ی عصبی استفاده شده است. این سیستم روی ۱۰۰۰ کلمه‌ی چاپی با قلم ترافیک آزمایش شده و دقت $97/83\%$ گزارش شده است. سرعت این روش ۳۹ حرف در ثانیه است.

در [۱۵] یک روش سریع و دقیق برای بازشناسی متن چاپی فارسی با درجه تفکیک ۳۰۰ نقطه بر اینچ معرفی شده است. این روش مبتنی بر جداسازی زیرکلمات به حروف و زیر حروف سازنده آن‌ها بوده و فرایند بازشناسی در چندین مرحله، با استفاده از طبقه‌بندهای شبکه‌ی عصبی تقویت شده انجام می‌گیرد. با توجه به این که جداسازی زیرکلمات، همواره یکی از مشکل‌ترین بخش‌های بازشناسی متون فارسی و عربی بوده است، کمترین اشتباه در فرایند جداسازی، موجب گسترش خطا در فرایند کلی بازشناسی می‌شود. در مقاله‌ی مذکور علاوه بر ارائه‌ی روش ساده و سریع برای جداسازی، با استفاده از نتایج مرحله‌ی بازشناسی، خطاهای مرحله‌ی جداسازی تصحیح می‌شود. به عبارتی، سیستم دارای یک حلقه‌ی بازخورد است که باعث افزایش قابلیت اعتماد آن می‌شود. داده‌های هدف در این تحقیق، متون فارسی با قلم‌های لوتوس، نازنین و میترا بوده است. البته الگوریتم به گونه‌ایست که برای سایر قلم‌ها قابل توسعه است. این روش روی ۸ صفحه متن فارسی با درجه‌ی تفکیک ۳۰۰ نقطه بر اینچ آزمایش شده و دقت بازشناسی 99% حاصل شده است.

در [۱۶] یک سیستم پیشنهادی جهت بازشناسی کاراکترهای چاپی فارسی با رویکردی بر مراحل بازنمایی و بازشناسی ارائه شده است. سیستم به صورت برون خط بوده و در ابتدا سیستم تصویر متنی را پیش‌پردازش می‌کند که شامل عملیات باینری‌سازی، کاهش نویز، اصلاح انحراف و نازک‌سازی است.

در مرحله‌ی قطعه‌بندی به ترتیب خطوط متن، نوار زمینه، زیرکلمات، کاراکترها و نقاط جداسازی می‌شوند. در مرحله‌ی بازنمایی برای اولین بار در یک سیستم بازشناسی، کاراکترهای چاپی چند فونتی فارسی از هفت گشتاور ثابت هندسی سازگار شده، با موفقیت استفاده شده که به طور ابتکاری شماره فرم کاراکتر (جداگانه، ابتدایی، میانی، انتهایی) نیز به بردار ویژگی‌ها اضافه می‌شود.

در [۱۷] یک روش جدید برای جداسازی خطوط متن، در تصاویر متون دست‌نوشته و چاپی ارائه شده است. در این روش ابتدا تکه‌های بهم پیوسته‌ی تصویر مشخص شده و برچسب‌گذاری می‌شود. سپس با در نظر گرفتن یک شش ضلعی پیرامون هر کدام از تکه‌های پیوسته، یک محدوده برای آن‌ها به دست می‌آید. با همپوشانی این محدوده‌ها با همدیگر، یک محدوده‌ی بزرگ‌تر برای هر کدام از خطوط متن حاصل می‌شود که به این ترتیب، هر کدام از خطوط متن در یک محدوده‌ی جدا از دیگر خطوط قرار می‌گیرد. در نهایت با جدا کردن این محدوده‌ها از یکدیگر، خطوط متن جداسازی می‌شوند. در آزمایش‌های انجام شده، درصد موفقیت این روش، برای متون دست‌نوشته‌ی فارسی و انگلیسی به ترتیب برابر با $87/2\%$ و $91/3\%$ و برای متون چاپی 100% به دست آمده است.

در [۱۸] برای غلبه بر مشکلات جداسازی حروف فارسی به علت اتصال حروف، روشی ترکیبی برای تشخیص حروف فارسی ارائه شده است که در آن جداسازی حروف (که منبع اصلی خطا در سیستم‌های تشخیص حروف فارسی کنونی است) به صورت کامل انجام نشده و فقط حرف اول هر زیرکلمه جدا شده و تشخیص داده می‌شود. سپس این حرف اول، معیاری برای دسته‌بندی زیرکلمات قرار می‌گیرد و هر زیرکلمه فقط در دسته‌ی مربوط تشخیص داده نمی‌شود. این روش روشی انعطاف‌پذیر و قدرتمند بوده که می‌تواند به آسانی توسعه داده شده و برای قلم‌های مختلف به کار برده شود. همچنین با استفاده از روش‌های پیش‌پردازش و پس‌پردازش مناسب، چارچوب ساده و کاملی پیشنهاد می‌دهد.

مقاله‌هایی درباره‌ی آ.سی.آر فارسی در مجله‌های معتبر خارجی منتشر شده است که از آن جمله

می‌توان به موارد زیر اشاره کرد.

موضوع مقاله‌های [۱۹] و [۲۰] بازنشاسی مجموعه‌ی از کلمات دست‌نویس است. در این مقالات سامانه‌ی جامع بازنشاسی متن فارسی ارائه شده است که در آن از مدل مخفی مارکوف و طبقه‌بند فازی استفاده شده است.

در [۲۱] با استفاده از روش‌های نوین، عملیات پیش‌پردازش که شامل دوسطحی کردن تصویر، حذف حاشیه‌های نامطلوب می‌باشد به همراه تشخیص تصاویر به صورت خودکار انجام شده است.

در [۲۲] با توجه به پیچیدگی‌های خط فارسی از جمله اشکال مختلف یک حرف در قسمت‌های مختلف کلمه، همپوشانی حروف و ...، از یک طبقه‌بند فازی برای شناسایی استفاده شده است.

روش‌های قبلی پیشنهاد شده برای تشخیص قلم اغلب با فیلترهای گابور و بر اساس تشخیص نوع قلم در یک بلوک از متن به جای یک خط و یا یک عبارت استفاده می‌شد، اما در [۲۳] خسروی با استفاده از فیلتر سوبل و رابرتز در هر خط و یا عبارت به تشخیص قلم پرداخته است.

در [۲۴] با توجه به شبیه بودن نقاط حروف فارسی به نویزهای موجود در تصویر، با تخمین اندازه‌ی نقاط در هر ناحیه، به حذف نویز در آن‌ها می‌پردازد.

در [۲۵] شیوه‌ای نوین برای تشخیص و پردازش معادلات ریاضی فارسی ارائه شده است که در بخش تشخیص کاراکتر، از یک سیستم فازی برای کلاسه‌بندی استفاده شده است.

در [۲۶] جهت افزایش سرعت آموزش و تست شبکه عصبی، از کارت‌های پردازش گرافیکی استفاده شده است.

۳-۲- نتیجه‌گیری

در این فصل به مروری بر سیر تحولی ا.سی.آر پرداخته شد که در آن نسل‌های مختلف ا.سی.آر معرفی شده و روند تحقیقاتی آن معرفی گردید. همچنین در پایان به معرفی تحقیقاتی در زمینه‌ی ا.سی.آر فارسی پرداخته شد و موانع موجود بر سر راه تکامل آن بیان گردید. در فصل آتی نیز به

معرفی مراحل مختلف پیاده‌سازی ا.سی.آر پرداخته خواهد شد.

مراحل پیاده‌سازی ا.سی.آر

۳-۱- مقدمه

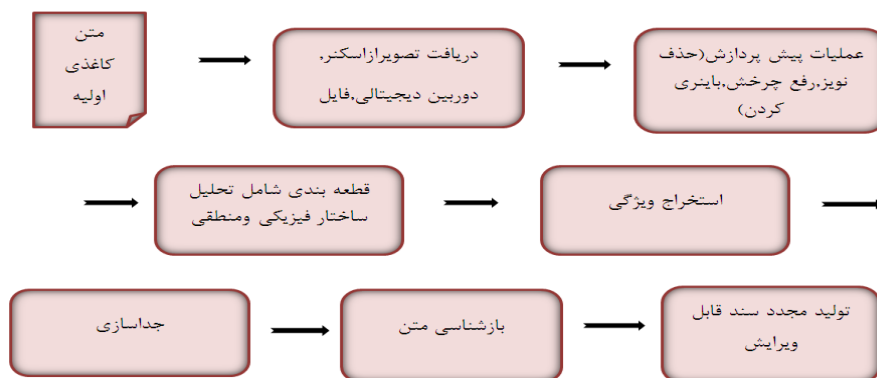
در فصل گذشته به معرفی سیر تکاملی ا.سی.آر و روش‌های گوناگون استفاده شده در طراحی سیستم‌های بازشناسی متن پرداخته شد. در این فصل تعاریف، مفاهیم، اصطلاحات بازشناسی متن و تمام مراحل مختلفی که برای یک سیستم بازشناسی متن لازم است، به صورت کلی شرح داده می‌شود. این مراحل شامل عملیات پیش‌پردازش، قطعه‌بندی، استخراج ویژگی و بازشناسی متن و روش‌های کلی برای پیاده‌سازی هر کدام از این مراحل است. پیاده‌سازی بعضی از این بخش‌ها ساده و بخش‌هایی از قبیل جداسازی و بازشناسی، کار پیچیده‌ای است.

۳-۲- سیستم بازشناسی متن

به طور کلی برای تبدیل هر گونه تصویر به فایل قابل ویرایش و قابل فهم برای رایانه، نیاز به یک سیستم بازشناسی متن می‌باشد.

اجزای اصلی یک سیستم بازشناسی متن عبارتند از:

- تصویربرداری از متن ورودی با سطوح خاکستری و درجه تفکیک مناسب
- پردازش‌های اولیه
- قطعه‌بندی
- استخراج ویژگی‌ها
- بازشناسی با یک یا چند طبقه‌بند
- به‌کارگیری اطلاعات جانبی مانند مجموعه لغات معتبر، اطلاعات آماری مربوط به رخداد حروف، اطلاعات دستوری و معنایی



شکل ۱-۳: مراحل سیستم بازشناسی متن

۳-۲-۱- پیش پردازش

بعد از اینکه تصویر متن از روبشگر به محیط الکترونیکی انتقال یافت به علت‌هایی از قبیل کثیف بودن صفحه کاغذ یا سطح روبشگر، استفاده از تصاویر کاربندی، کج گذاشتن صفحه، رنگی بودن آن و غیره باید عملیات اولیه‌ای روی آن انجام پذیرد تا برای مراحل بعدی بازشناسی آماده شود. در غیر این صورت به علت نواقص موجود در بازشناسی نتیجه خوبی گرفته نشده و باعث پایین آمدن دقت سیستم بازشناسی می‌شود. عملیات پیش‌پردازش می‌تواند شامل باینری کردن تصویر، حذف نویز و رفع چرخش^۱ باشد.

۳-۲-۱-۱ باینری کردن تصویر

بسیاری از الگوریتم‌های بازشناسی حروف از شکل دوسطحی^۲ تصاویر ورودی استفاده می‌نمایند. در نتیجه اگر تصویر ورودی خاکستری یا رنگی باشد باید ابتدا به تصویر دوسطحی تبدیل شود. برای دو سطحی کردن تصویر ورودی می‌توان از یک سطح آستانه‌ی ثابت یا یک سطح آستانه‌ی وفقی استفاده نمود. الگوریتم مورد نظر نباید باعث ایجاد نویز در مرز و داخل حروف، چسبیدگی، خوردگی یا

^۱ Skew correction

^۲ Binary

جداسازی آن‌ها شود.

۳-۲-۱-۲-۳ نویز

وقتی تصویر از یک منبع غیر دیجیتالی مانند کاغذ به دست می‌آید به علت‌های مختلفی از جمله جنس کاغذ، رنگ کاغذ، استفاده از کاغذ به همراه کاربن، استفاده از نسخه‌ی کاغذی پس از چند بار کپی گرفتن، وجود پارگی‌ها یا سوراخ دستگاه پانچ بر روی کاغذ و همچنین کثیف بودن سطح روبشگر و غیره نویز در تصویر ایجاد می‌شود. برای حذف نویز می‌توان از فیلتر نمک و فلفل، حذف بافت پس زمینه و ... استفاده کرد. حذف نویز باعث بهبود بازشناسی می‌گردد و از این رو حائز اهمیت است.

۳-۲-۱-۳ رفع چرخش

کاربر هنگام گذاشتن صفحات کاغذی بر روی سطح اسکنر ممکن است دقت کافی نداشته باشد و با کج گذاشتن صفحات در تصویر ایجاد شده چرخشی پدید می‌آید. همچنین ممکن است هنگام عکس‌برداری با کج نگه داشتن دوربین یا درست قرار نگرفتن تصویر در جلوی دوربین، تصویر حاصل نسبت به مبنا اندکی چرخش داشته باشد. در این صورت حاشیه‌های سیاه ناخواسته‌ای نیز اطراف تصویر ایجاد می‌شود. در نتیجه برای اینکه در مرحله بازشناسی با دقت بیشتری شناسایی صورت بگیرد، باید زاویه‌ی چرخش مربوطه را یافته و با روش‌های مناسب حذف کرد. همچنین حاشیه‌های سیاه اطراف تصویر نیز باید حذف شوند. برای رفع چرخش روش‌های زیادی متداول است که از آن جمله می‌توان افکنش افقی، استفاده از تبدیل هاف و استفاده از عملگرهای مورفولوژی را نام برد.

۳-۲-۱-۳-۱ افکنش افقی

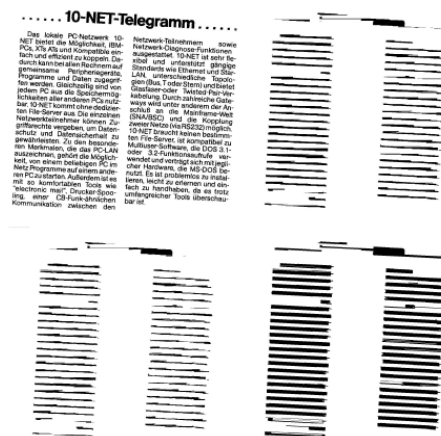
در [۲۷] ساده‌ترین روش‌ها با استفاده از افکنش افقی تصویر در جهات مختلف و یافتن بیشینه‌ی افکنش، زاویه‌ی چرخش یافت می‌شود. مزیت اصلی این روش، سادگی آن است ولی مشکل عمده‌ی این روش این است که برای محاسبه‌ی زاویه‌ی کلی تصویر اندکی کند است.

۳-۲-۱-۲-۳ تبدیل هاف

این تبدیل که اساساً برای یافتن خطوط در جهات مختلف استفاده می‌شود، برای یافتن زاویه‌ی چرخش سند، دو عیب عمده دارد: اول اینکه حجم محاسباتی بسیار بالا و در نتیجه سرعت بسیار کمی دارد. دوم اینکه در اسناد متنی معمولاً خط ممتد گرافیکی نداریم و این تبدیل تلاش می‌کند، خطوط پایه‌ی بلوک‌های متنی را پیدا کند که اگر فاصله‌ی بین خطوط کم باشد یا قلم متن خیلی بزرگ باشد، با مشکل روبرو می‌شود.

۳-۳-۱-۲-۳ عملگرهای مورفولوژی

در [۲۸] با استفاده از عملگر بستن^۱ که ترکیبی از عملگر گسترش^۲ و فرسایش^۳ است، هر خط به یک مؤلفه‌ی متصل و تقریباً یکنواخت تبدیل می‌شود. این کار به ازای عناصر ساختاری خاصی که برای زوایای مختلف حساب شده‌اند، انجام می‌شود. زاویه‌ای که در آن ضخامت خطوط مذکور بیشترین مقدار را داشته باشد، جایی است که عملگر مورد نظر در راستای خطوط متن قرار گرفته و از این رو به عنوان زاویه‌ی چرخش سند در نظر گرفته می‌شود (شکل ۲-۳)



شکل ۲-۳: یافتن زاویه‌ی چرخش با اعمال چند عملگر مورفولوژی در زوایای مختلف [۲۸]

- 1 Closing
- 2 Dilation
- 3 Erosion

این روش نسبت به سایر روش‌ها دقیق‌تر به نظر می‌رسد؛ لیکن از آنجا که برای هر زاویه، چندین عملگر با عناصر ساختاری تقریباً بزرگ بر روی تصویر اعمال می‌شود، مدت زمان اجرای الگوریتم زیاد است. البته روش‌هایی برای پیاده‌سازی بهینه‌ی عملگرهای مورفولوژی، خاصه روی تصاویر باینری، ابداع شده است که سرعت کار را بهبود می‌بخشد [۲۹].

۳-۲-۱-۳ یافتن چرخش بر اساس یافتن زاویه‌ی خطوط متن

در اسناد ورودی به سیستم بازشناسی متن همیشه تعدادی بلوک متنی وجود دارد که هر کدام از تعدادی خط متنی تشکیل شده‌اند. خطوط یک بلوک متنی غالباً در یک راستا هستند، در نتیجه اگر زاویه‌ی چرخش این خطوط را بیابیم، زاویه‌ی چرخش کل سند یافته می‌شود. در این روش [۳۰] ابتدا تصویر ورودی به ابعاد مناسبی تغییر مقیاس داده می‌شود. پس از دوسطحی‌سازی، تصویر کلمات یک خط با استفاده از عملگرهای مورفولوژی به صورت یک مؤلفه پیوسته درمی‌آید. با استفاده از الگوریتم برجسب‌زنی مؤلفه‌ها سه خطی که خواص خطوط متن را دارند انتخاب شده و زاویه‌ی چرخش هر کدام از این خطوط از روشی مثلاً افکنش افقی محاسبه می‌شود. زاویه‌ی نهایی از متوسط‌گیری سه زاویه‌ی به دست آمده، حاصل می‌شود. این روش به علت انتخاب چرخش از روی خطوط متن به بلوک‌های بزرگ در تصویر حساسیت ندارد و سرعت خوبی دارد. پس از یافتن زاویه‌ی چرخش سند، تصویر آن با همان زاویه و در خلاف جهت چرخانده می‌شود تا چرخش آن حذف شود. نتیجه‌ی این کار در شکل ۳-۳ نشان داده شده است. در این شکل ضمن چرخاندن تصویر، سعی شده است که ابعاد آن تغییر نکند، به عبارتی حاشیه‌ی اضافی که در حین چرخش ایجاد می‌شود، حذف شده است. تصویر نهایی سند در شکل ۳-۳ مشاهده می‌شود.



شکل ۳-۳: تصویر نهایی پس از رفع چرخش

۲-۲-۲- قطعه‌بندی

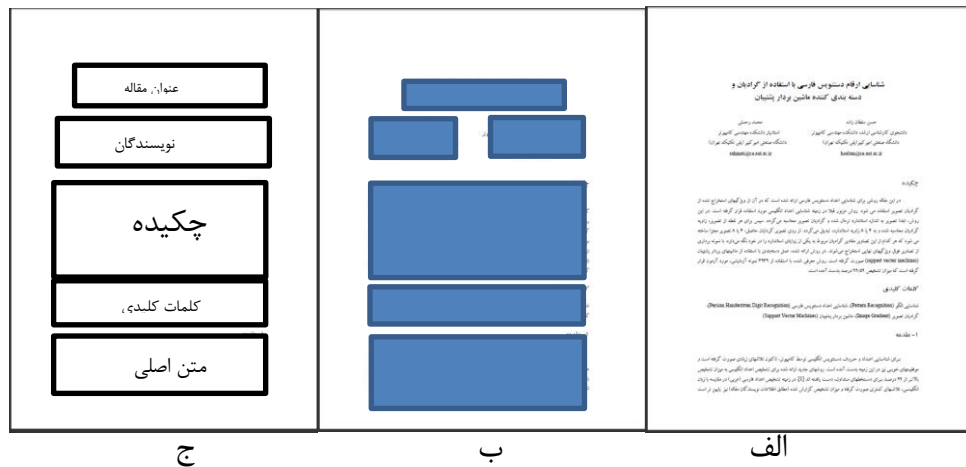
قطعه‌بندی عبارت است از روش‌هایی که بخش‌های مختلفی همچون جدول‌ها، تصاویر، پاراگراف‌ها، جملات یا کلمات و حروف را از تصویر سند استخراج می‌نمایند.

۳-۲-۳- تحلیل ساختار فیزیکی و منطقی

هر تصویر متن به طور کلی می‌تواند شامل متن، عکس، جدول، گراف و غیره باشد. بعد از مرحله‌ی پیش‌پردازش برای شناسایی بهتر می‌توان هر یک از بلوک‌های تصویر را از هم جدا کرده و سپس هر یک به صورت جدا به سیستم‌های بازشناسی سپرده شود. در نتیجه هر یک از بلوک‌ها با توجه به موقعیت آن‌ها در صفحه و همچنین موقعیت هر یک نسبت به همدیگر از هم جدا می‌شوند. همچنین از نظر معنایی می‌توان به هر یک از بلوک‌های متن متناسب با عملکرد معنایی آن‌ها برچسب زد.

پس از تحلیل ساختار اسناد، بلوک‌های مختلف تصویر از هم جدا می‌شوند. جداسازی با در نظر گرفتن ویژگی‌های هر بلوک انجام می‌شود. به عنوان مثال بلوک‌های متن را می‌توان با توجه به خاصیت قرار گرفتن منظم حروف، زیرکلمات و کلمات در کنار هم و تشکیل خطوط موازی تشخیص داد یا بلوک‌های گراف را می‌توان از فضاها سفید بین اجزای آن و داشتن خطوط راست در ساختار آن از

دیگر بلوک‌ها جدا کرد. پس از تعیین بلوک‌های متن، اجزای متن باید با توجه به روش بازشناسی از یکدیگر جدا و به عنوان ورودی به مرحله بازشناسی فرستاده شوند. در این مرحله تصویر بلوک متن به تصاویر خطوط متن تفکیک می‌شود. از تصویر خط، تصاویر کلمات، زیرکلمات و حروف استخراج می‌شوند.



شکل ۳-۴: تصویر صفحه اول یک مقاله (ب) ساختار فیزیکی (ج) ساختار منطقی

۳-۲-۴- استخراج ویژگی^۱

در این مرحله حروف و زیرکلمات به شکل مناسب‌تری که بازشناسی آن‌ها ساده‌تر باشد توصیف می‌شوند. از تصاویر حروف یا زیرکلمات، مشخصاتی که بیان‌کننده‌ی ویژگی هر یک از آن‌ها باشد و اختلاف بین آن‌ها را بهتر بیان کند، استخراج می‌شود. این مشخصات به صورت ویژگی‌های توصیف‌کننده‌ی حرف یا زیرکلمه یا هر جزء متن می‌باشد. بنابراین از این مرحله به بعد هر جزء متن با یک بردار ویژگی به رایانه معرفی می‌شود. در واقع استخراج ویژگی‌ها مجموعه‌ی کلیه‌ی محاسباتی است که روی الگوهای بدست آمده از مرحله‌ی پیش‌پردازش انجام می‌شود تا بردار ویژگی‌های متناظر با هر الگو تعیین گردد. روش استخراج ویژگی‌ها بر اساس دو نوع از ویژگی‌ها است: (۱) آماری (۲) ساختاری

1 Feature Extraction

از ویژگی‌های ساختاری می‌توان به تعداد و نوع پاره خط‌های مختلف (عمودی، افقی، بالا رونده و پایین رونده) اشاره کرد. گشتاورها و مکان‌های مشخصه، ناحیه‌بندی، افکنش‌ها و پروفایل‌ها، نقاط تقاطع و فاصله‌ها نمونه‌هایی از ویژگی‌های آماری هستند.

۳-۲-۵ - بازشناسی متن (با یک یا چند طبقه‌بند)

مرحله‌ی بازشناسی شامل روش‌هایی است که در آن بردار ویژگی مناسب از تصویر استخراج شده و در مرحله‌ی اول با بردارهای ویژگی موجود در پایگاه داده مقایسه شده و نهایتاً بر مبنای کم‌ترین فاصله با یکی از بردارهای مرجع متناظر می‌شود. بازشناسی هر یک از اجزای تصویر متن با طبقه‌بندی آن‌ها در یک مجموعه حروف، زیرکلمات یا کلمات معتبر انجام می‌شود. روش‌های طبقه‌بندی به سه شاخه‌ی روش‌های نحوی یا ساختاری، روش‌های آماری و روش‌های مبتنی بر شبکه‌های عصبی تقسیم می‌شوند [۳۱].

۳-۲-۵-۱ روش‌های نحوی یا ساختاری

در این روش‌ها یک الگو^۱ بر اساس اجزای تشکیل‌دهنده آن و روابط بین این اجزاء طبقه‌بندی می‌شود. طبقه‌بندی کننده در ابتدا اجزای تشکیل‌دهنده‌ی نویسه را بازشناسی کرده و با کنار هم قرار دادن آن‌ها بر اساس یک سری روابط نحوی، آن نویسه را بازشناسی می‌کند.

۳-۲-۵-۲ روش‌های آماری

در این روش‌ها، بردار ویژگی متناظر الگو با کلاس‌های آموزش داده شده به سیستم مقایسه می‌شود و نزدیک‌ترین کلاس مشخص می‌شود. روش نزدیک‌ترین همسایه متعلق به این دسته از روش‌های طبقه‌بندی است. در این روش، الگوی ورودی با الگوهای آموزش داده شده به سیستم در مرحله‌ی

1 Pattern

آموزش، مقایسه می‌شود. الگوی با نزدیک‌ترین فاصله به الگوی ورودی به عنوان کلاس الگوی ورودی تعیین می‌شود.

۳-۲-۵-۳ شبکه‌های عصبی

شبکه‌های عصبی از شبکه‌های بهم متصل و ساده با واحدهای پردازش غیرخطی تشکیل می‌شوند. شبکه‌های عصبی مصنوعی بر اساس توپولوژی شبکه، مشخصات نرون‌ها و الگوریتم آموزش آن‌ها از هم متمایز می‌شوند. از شبکه‌ی عصبی برای آموزش الگوهای مختلف به سیستم استفاده می‌شود. آموزش شبکه‌ی عصبی به کمک ویژگی‌های استخراج شده از الگوها انجام می‌شود. پس از آموزش شبکه، با ارائه‌ی ویژگی استخراج شده از یک الگوی ناشناخته به آن، کلاس متناظر آن الگو در خروجی شبکه عصبی مشخص می‌شود.

۳-۲-۶-۳ پس پردازش

با بازشناسی تمام تصاویر اجزای متن یک بلوک، متن متناظر با هر یک از آن‌ها مشخص می‌شود. برای بدست آوردن متن نهایی باید پس پردازش‌هایی روی این متن انجام شود مثلاً از به هم چسباندن حروف بازشناسی شده در مرحله‌ی قبل، زیرکلمات و کلمات تشکیل می‌شوند. تشکیل این زیرکلمات و کلمات با استفاده از موقعیت آن‌ها در تصویر، یک مرحله پس پردازش است. مرحله‌ی بعدی می‌تواند شامل مقایسه و تایید کلمات بازشناسی شده در مرحله‌ی بازشناسی با یک واژه‌نامه باشد.

۳-۳-۳ روش‌های بازشناسی متن

سیستم‌های بازشناسی متن را با توجه به روش بازشناسی می‌توان به سه دسته تقسیم کرد.

۳-۳-۱- روش‌های مبتنی بر جداسازی حروف

در این روش‌ها زیرکلمات و کلمات به حروف شکسته می‌شوند و حروف بازشناسی می‌شوند به

گونه‌ای که از بهم چسباندن این حروف، زیرکلمات و کلمات تشکیل می‌شوند.

روش‌های متداول موجود برای جداسازی حروف فارسی را می‌توان به پنج گروه عمده تقسیم کرد:

۳-۱-۳ روش‌های مبتنی بر افکنش عمودی

روش‌های مبتنی بر افکنش از ساده‌ترین روش‌های جداسازی هستند، در این روش‌ها از هیستوگرام عمودی تصویر استفاده می‌شود [۳۲]، [۳۳]. جاهایی که مقادیر هیستوگرام صفر است نقاط جداسازی محسوب می‌شوند. این روش در مواردی که بین حروف همپوشانی وجود دارد (به عنوان مثال حروفی که زیر سرکش کاف هستند) مناسب نیست.

۳-۱-۳ روش‌های مبتنی بر کانتور

روش‌های مبتنی بر آنالیز کانتور، با حرکت روی کانتور کلمات و یافتن مینیمم و ماکزیمم‌های محلی و تحلیل آن‌ها نقاط جداسازی را پیدا می‌کنند [۳۴]، [۱۰]، [۱۲] و [۳۵]. در زبان انگلیسی از این روش‌ها غالباً برای جداسازی حروف کلمات دست‌نویس استفاده می‌شود. در مورد زبان فارسی هم از این روش‌ها استفاده شده که نتایج قابل قبولی داشته است. این روش‌ها نیز غالباً در برخورد با همپوشانی‌های عمودی دچار اشکال می‌شوند.

۳-۱-۳ روش‌های مبتنی بر آنالیز پروفایل

روش‌های مبتنی بر آنالیز پروفایل شبیه به روش‌های مبتنی بر کانتورند با این تفاوت که به جای استفاده از کانتور کلمه، تنها از پروفایل بالایی کلمه استفاده می‌کنند [۳۶]. پروفایل بالایی مجموعه‌ی اولین نقاط بالایی تصویر در هر ستون آن است. در این روش‌ها، مشابه روش‌های مبتنی بر کانتور احتمال پنهان شدن برخی حروف در زیر سایر حروف زیاد است مثلاً در فارسی سرکش کاف یا سرکش الف چنین حالتی را ایجاد می‌کند و یا در انگلیسی ممکن است حروفی مثل i در زیر حروفی مثل T ناپدید شوند.

۴-۱-۳-۳ روش‌های مبتنی بر اسکلت

روش‌های مبتنی بر اسکلت که از دقت مناسب‌تری برخوردارند، ابتدا با نازک‌سازی تصویر زیرکلمه، اسکلت آن را تولید کرده و با یافتن خط پایه‌ی اسکلت حاصله و تحلیل آن، نقاط جداسازی را می‌یابند. اشکال این روش زمان گیر بودن و توزیع خطا در صورت یافتن اسکلت نادرست است.

۵-۱-۳-۳ روش‌های مبتنی بر اسکلت هوشمند یا جداسازی همراه با بازشناسی

در چند سال اخیر، الگوریتم‌های هوشمند به دلیل اینکه در آن‌ها امکان تصحیح نتایج جداسازی وجود دارد از محبوبیت بیشتری برخوردار شده‌اند [۳۷]، [۲۳]. در این الگوریتم‌ها، پس از جداسازی اولیه، یک موتور بازشناسی حروف، نتایج جداسازی را تأیید یا رد می‌کند.

۲-۳-۳ روش‌های مبتنی بر شکل کلی زیرکلمات

در این روش کل زیرکلمه به صورت یک الگو در نظر گرفته می‌شود و ویژگی‌های لازم برای طبقه‌بندی از تصویر کل کلمه استخراج می‌شوند. ابتدا باید یک مجموعه پایگاه داده‌ی کامل از لغات پرکاربرد در نظر گرفته شود که تهیه‌ی این پایگاه به دلیل نبود یک مرجع، کار دشواری است. با توجه به اینکه بسیاری از زیرکلمات شباهت زیادی باهم دارند و تنها تفاوت آن‌ها در تعداد نقاط و موقعیت قرار گرفتن آن‌ها می‌باشد، برای بررسی تأثیر نقاط و علائم، دو رویکرد برای بازشناسی زیرکلمات در نظر گرفته می‌شود:

- بازشناسی زیرکلمات بدون حذف نقاط و علائم
- بازشناسی زیرکلمات با حذف نقاط و علائم

برای بازشناسی زیرکلمه‌ی ورودی باید تعیین شود کدام یک از زیرکلمات موجود بیشترین شباهت را به زیرکلمه‌ی ورودی دارد. به دلیل زیاد بودن تعداد کلاس‌های خروجی، نمی‌توان از روش‌های مرسوم بازشناسی الگو استفاده کرد. بنابراین باید از روش‌های خوشه‌یابی کمک گرفت. در نتیجه ابتدا

در فاز آموزش، ویژگی‌های مناسبی از زیرکلمات استخراج شده و با استفاده از الگوریتم‌های خوشه‌یابی، زیرکلمات تهیه شده برای آموزش به چندین خوشه تقسیم می‌شوند. سپس در فرایند بازشناسی زیرکلمه‌ی ورودی، ابتدا بر اساس ویژگی‌های مورد استفاده در خوشه‌یابی، نزدیک‌ترین خوشه‌ها به زیرکلمه‌ی ورودی تعیین می‌شود و پس از تعیین خوشه، خود زیرکلمه شناسایی می‌شود.

در این مرحله، ویژگی‌های دیگری که توصیف بهتری شامل جزئیات زیرکلمه را داشته باشند، استخراج شده و با ویژگی‌های متناظر از زیرکلمات کاندید مقایسه می‌شود. نزدیک‌ترین زیرکلمه به زیرکلمه‌ی ورودی به عنوان برنده انتخاب می‌شود.

در رویکرد بازشناسی با حذف نقاط و علائم، ابتدا با استفاده از روش برچسب‌زنی مؤلفه‌ها نقاط را حذف کرده و بدنه‌ی زیرکلمه شناسایی می‌شود که در این صورت پیکسل‌های نويز هم از بین می‌رود. با حذف نقاط و علائم حجم پایگاه داده به حدود نصف کاهش یافته و با این کار فرایندهای خوشه‌یابی و بازشناسی ساده‌تر و در عین حال دقیق‌تر می‌شود و همچنین می‌توان زیرکلماتی را که در پایگاه داده نبوده‌اند ولی بدنه‌ی آن‌ها مشابه یکی از بدنه‌های موجود باشد را شناسایی کرد. پس از بازشناسی اطلاعات می‌توان نقاط را به بدنه اضافه کرد. برای تولید پایگاه تصویری زیرکلمات بدون نقطه، مشابه حالت با نقطه عمل می‌شود با این تفاوت که یک مرحله‌ی حذف نقاط قبل از افزودن به پایگاه داده قرار می‌گیرد.

مزیت استفاده از شکل کلی در جاهایی است که جداسازی امکان پذیر نیست. مثلاً برای قلم‌های خاص که جداسازی آن‌ها خیلی مشکل است یا در مواردی که نقاط و علائم به بدنه چسبیده‌اند و یا اندازه‌ی قلم خیلی کوچک است.

۳-۳-۳ - روش‌های ترکیبی (شکل کلی-جداسازی)

روش‌های ترکیبی بدین ترتیب است که با توجه به نوع قلم و زبان متن در مواردی که روش جداسازی پاسخ مناسب‌تری برای شناسایی بدهد از این روش استفاده کرده و در سایر موارد که شکل

کلی مناسب‌تر باشد و یا جداسازی به راحتی امکان پذیر نباشد، برای بازشناسی از شکل کلی استفاده می‌کند. در این پایان‌نامه، رویکرد پیش‌بینی شده بازشناسی بر اساس جداسازی خواهد بود.

۳-۴- قابلیت‌های سیستم‌های بازشناسی متن

در بررسی یک سیستم بازشناسی متن، قابلیت آن در بازشناسی متون مختلف سنجیده می‌شود که از آن جمله می‌توان به متون چاپ شده با چند قلم محدود و نامحدود، بازشناسی متون داخل جداول، بازشناسی متون چند زبانی، بازشناسی متون چند ستونی، بازشناسی متون دست‌نویس داخل متن چاپی، بازشناسی فرمول، بازشناسی متون رنگی و خاکستری، بازشناسی متون دست‌نویس با محدودیت و بدون محدودیت و بازشناسی دست‌نوشته برخط با صفحات رقمی کننده مختلف اشاره کرد. در ادامه برخی از موارد بالا به اختصار بررسی می‌شوند.

۳-۴-۱- متن‌های تک قلمی

سیستم بازشناسی متن ممکن است برای بازشناسی یک قلم مشخص طراحی شود. در چنین مواردی کارایی این سیستم در بازشناسی متون با قلم‌های دیگر پایین‌تر از قلم آموزش داده شده به آن خواهد بود.

۳-۴-۲- متن‌های چند قلمی^۱

سیستم بازشناسی متن ممکن است برای بازشناسی چند قلم محدود طراحی شده و آموزش داده شود. کارایی این سیستم در بازشناسی متون با قلم‌هایی که به سیستم آموزش داده نشده‌اند پایین‌تر از متون چاپ شده با قلم‌های آموزش دیده خواهد بود. در مورد قلم‌های مشابه با قلم‌های آموزش، ممکن است کارایی سیستم بهتر باشد.

1 Multi-Font

۳-۴-۳- متن‌های دست‌نویس مقید^۱

در یک سیستم بازشناسی متون دست‌نویس (برخط و برون‌خط) ممکن است محدودیت‌هایی برای نوشتن بعضی علائم و نقاط، حروف، زیرکلمات یا کلمات فرض شود. به عنوان مثال، سه نقطه باید به صورت یک نیم دایره نوشته شود و سه نقطه‌ی جدا از هم به عنوان سه تک نقطه بازشناسی می‌شوند. نرخ بازشناسی در این سیستم‌ها بسته به محدودیت‌های اعمال شده می‌تواند بسیار بالا باشد.

۳-۴-۴- متن‌های دست‌نویس بدون قید^۲

تنوع دستخط افراد مختلف و کیفیت معمولاً پایین تصاویر، بازشناسی متون دست‌نوشته را به کاری پیچیده تبدیل کرده است. در حالی که این متون بدون محدودیت نیز نوشته شوند، دشواری بازشناسی آن‌ها دو چندان می‌شود. به این دلیل نرخ بازشناسی سیستم‌های بازشناسی متون بدون محدودیت نسبت به سیستم‌های با محدودیت پایین‌تر است.

۳-۴-۵- متن‌های داخل جداول و فرم‌ها

بسیاری از اسناد اداری و متون علمی شامل جداول و فرم‌هایی با متون چاپی هستند. در برخی موارد متن داخل جدول به خطوط آن می‌چسبد. توانایی یک سیستم در جداسازی جدول از تصویر و بازشناسی این متون، قابلیت مهمی برای یک سیستم بازشناسی متن است.

۳-۴-۶- متن‌های چند زبانی

در تصویر یک سند ممکن است متون با چند زبان وجود داشته باشد یا هر سند به یک زبان نوشته شده باشد. یک سیستم بازشناسی متن که قابلیت بازشناسی زبان‌های مختلف در تصاویر مجزا را دارد،

1 Constrained

2 Unconstrained

ممکن است قادر به بازشناسی متون زبان‌های مختلف در یک تصویر کنار هم نباشد. بنابراین در انتخاب یک سیستم بازشناسی متن چند زبانی باید به این قابلیت دقت شود.

۳-۴-۷- متن‌های شامل فرمول

اگر سیستم بخواهد در سطح وسیع عملیات بازشناسی را انجام بدهد توجه به نکاتی از قبیل وجود فرمول‌های زیاد در متون علمی نباید نادیده گرفته شود. در نتیجه برای بازشناسی متون علمی، باید سیستمی با قابلیت بازشناسی فرمول انتخاب شود. باید توجه داشت که به دلیل حضور علائم و برخی نویسه‌ها در فرمول‌های ریاضی، قابلیت بازشناسی فرمول با قابلیت بازشناسی متون لاتین متفاوت است.

۳-۴-۸- متن‌های چند ستونی

برخی متون کتاب‌ها، مجلات و روزنامه‌ها به صورت چند ستونی چاپ می‌شوند. بنابراین برای بازشناسی متون این اسناد، سیستم بازشناسی باید قابلیت جداسازی متون چند ستونی را داشته باشد.

۳-۴-۹- تصویرهای خاکستری و رنگی

همه‌ی تصاویر روبش شده‌ی متون، دودویی نبوده و سهم عمده‌ای از آن‌ها به صورت خاکستری و بخشی از آن‌ها نیز به صورت رنگی هستند. بنابراین یکی از ویژگی‌های سیستم بازشناسی متن، قابلیت بازشناسی متون خاکستری و رنگی است.

۳-۴-۱۰- بازشناسی قلم

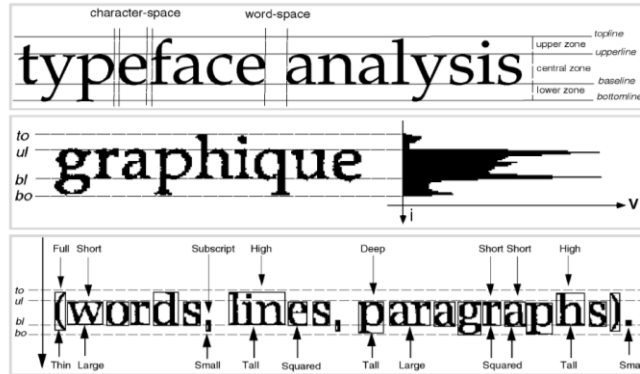
دو رویکرد مرسوم در زمینه‌ی بازشناسی قلم وجود دارد:

- بازشناسی قلم بر مبنای ویژگی‌های تایپوگرافی

روش‌های تایپوگرافی نوع قلم را در سطح یک خط یا عبارت تشخیص می‌دهند [۳۸] و [۳۹].

ویژگی‌های تایپوگرافی، مانند شیب حروف، ضخامت حروف، عرض فاصله‌ی بین حروف و افکنش در نواحی بالا، پایین و وسط خط معمولاً مشخصات یک قلم را به خوبی نمایندگی می‌کنند، ولی استخراج

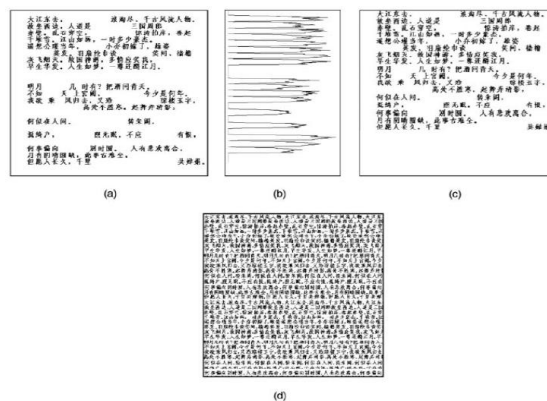
این ویژگی‌ها نیازمند این است که سند ورودی خالی از نویز باشد و با درجه تفکیک بالا، مثلاً ۳۰۰ نقطه بر اینچ، اسکن شده باشد (شکل ۳-۵).



شکل ۳-۵: ویژگی‌های تایپوگرافی [۴۰]

• با استفاده از ویژگی‌های بافتی

ویژگی‌های بافتی بیش از سایر ویژگی‌ها برای بازشناسی قلم استفاده شده‌اند. این ویژگی‌ها غالباً با استفاده از فیلترهای گابور و تبدیل موجک حاصل می‌شوند. غالب کارهای صورت گرفته با استفاده از این ویژگی‌ها، در سطح یک بلوک متنی انجام شده است نه در سطح یک خط یا کلمه. در این روش‌ها ابتدا یک بلوک متنی پردازش شده و بعد از حذف فضاهای خالی بین خطوط و بین حروف، یک بافت یکپارچه حاصل می‌شود. سپس ویژگی‌های بافتی استخراج شده و قلم متن بر اساس این ویژگی‌ها بازشناسی می‌شود.



شکل ۳-۶: تشکیل یک بافت یکپارچه از روی بلوک متنی

مسئله‌ی مهم دیگر در بازشناسی قلم، پیچیدگی محاسباتی روش است که مستقیماً سرعت سامانه‌ی نویسه خوان را تحت تأثیر قرار می‌دهد. برخی ویژگی‌های بافتی مانند فیلترهای گابور که به کرات استفاده شده است بسیار زمان‌بر هستند و استفاده از آن‌ها در سیستم ا.سی.آر سرعت کلی سیستم را به شدت کاهش می‌دهد.

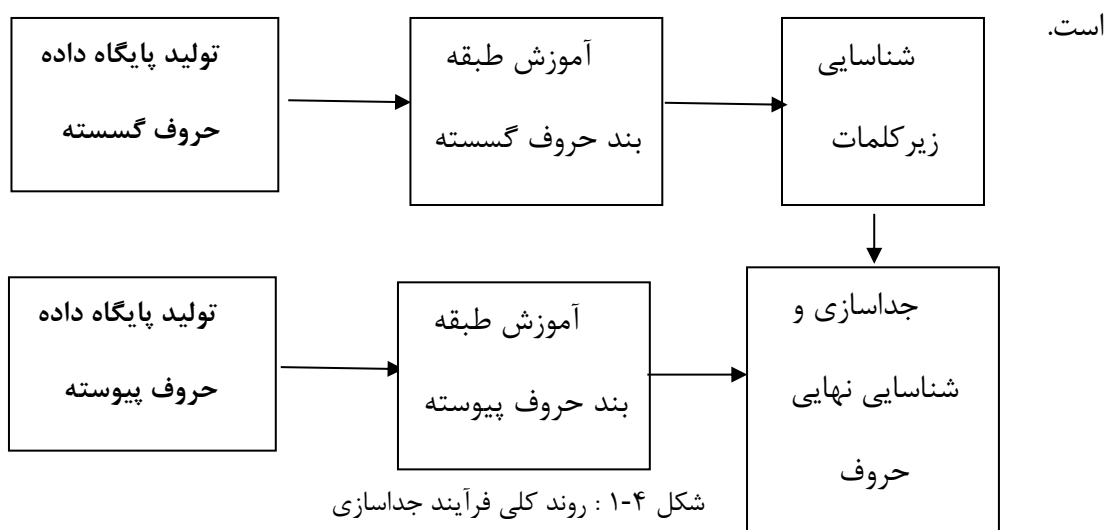
۳-۵- نتیجه‌گیری

در این فصل مفاهیم و اصطلاحات و مراحل مختلف یک سیستم بازشناسی متن به صورت کلی شرح داده شده و هرکدام به اختصار توضیح داده شد. انواع روش‌های بازشناسی متن و قابلیت‌های یک سیستم بازشناسی متن بیان شد. قابل ذکر است در این پایان‌نامه بعضی از این موارد پیاده‌سازی شده و برخی فقط جهت کامل بودن عملیات یک سیستم بازشناسی متن بیان گردیده است.

فصل چهارم

تولید پایگاه داده

در فصل گذشته راجع به مراحل کلی سیستم بازشناسی متن توضیحاتی داده شد. در این پایان نامه هدف، تولید سیستم بازشناسی است که متون نوشته شده با قلم Iranian sans را شناسایی کند. به منظور شناسایی از دو طبقه‌بند حروف گسسته و پیوسته استفاده می‌شود. برای آموزش این طبقه‌بندها نیاز به پایگاه داده می‌باشد. روند کلی کار بدین صورت است که ابتدا با استفاده از طبقه‌بند حروف گسسته زیرکلمات شناسایی شده و برای جداسازی به سمت طبقه‌بند حروف پیوسته فرستاده می‌شود. برای آموزش طبقه‌بندهای فوق نیاز به پایگاه داده می‌باشد. این روند در شکل ۴-۱ نشان داده شده

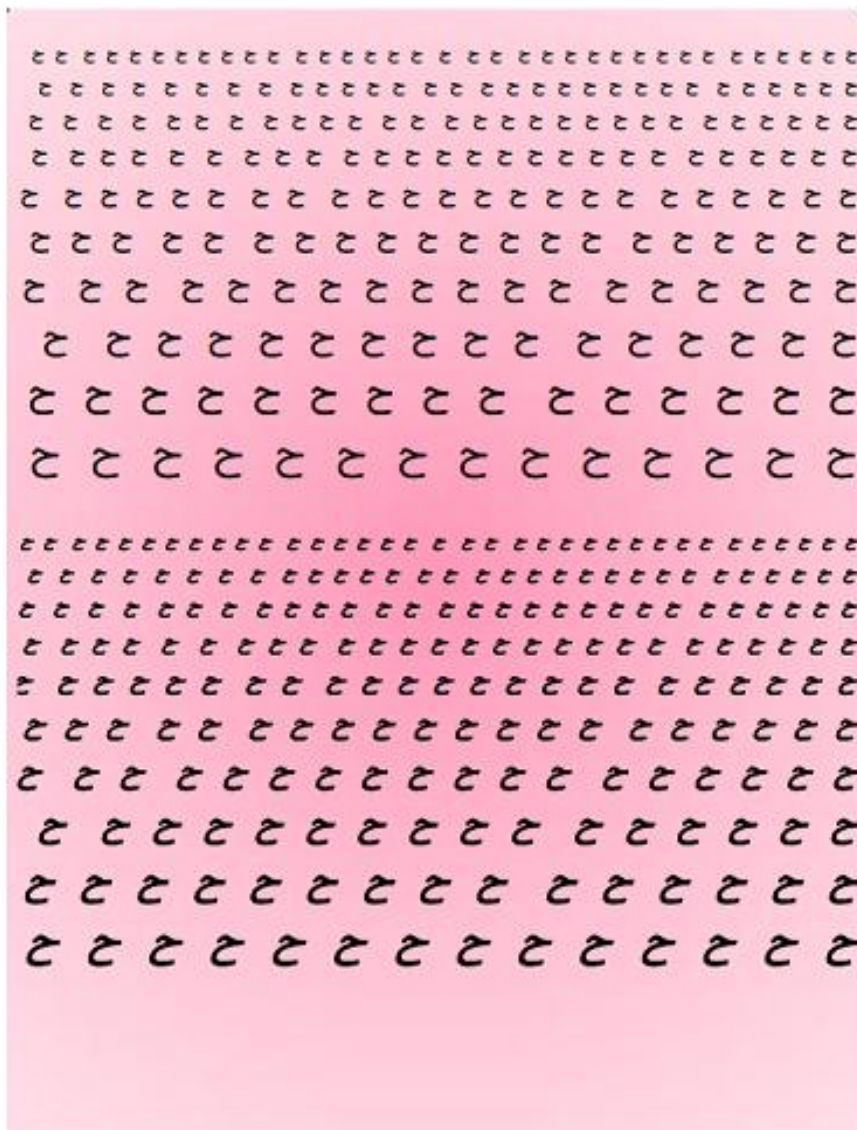


با توجه به اینکه هیچ پایگاه داده اولیه‌ای آماده‌ای در مورد قلم مذکور وجود ندارد، ابتدا به تولید پایگاه داده پرداخته و شبکه آموزش داده می‌شود. پایگاه داده برای حروف پیوسته و گسسته تولید می‌شود. اندازه‌های قلم بین ۹ تا ۲۴ بوده و حروف در حالت‌های ساده، پررنگ و زاویه‌دار نوشته شده‌اند. پس از تولید داده، مجموعه‌ی حروف بدون توجه به نقاط آن‌ها ابتدا با یک چاپگر لیزری مدل HP چاپ و سپس با یک اسکنر مدل HP با درجه تفکیک ۳۰۰ نقطه بر اینچ اسکن شده و به صورت فایل دیجیتالی درآمده‌اند. در حین چاپ و اسکن سعی شده با تغییر سطح روشنایی و استفاده از کاغذهایی با پس‌زمینه‌های مختلف تصاویری با کیفیت متفاوت ایجاد شود. در نهایت پس از اعمال پیش‌پردازش بر روی تصاویر فوق، مؤلفه‌های مختلف پایگاه داده با برچسب‌زنی مؤلفه‌ها تشخیص داده

شده و پس از استخراج ویژگی، نمونه‌ها برای آموزش به شبکه داده می‌شوند.

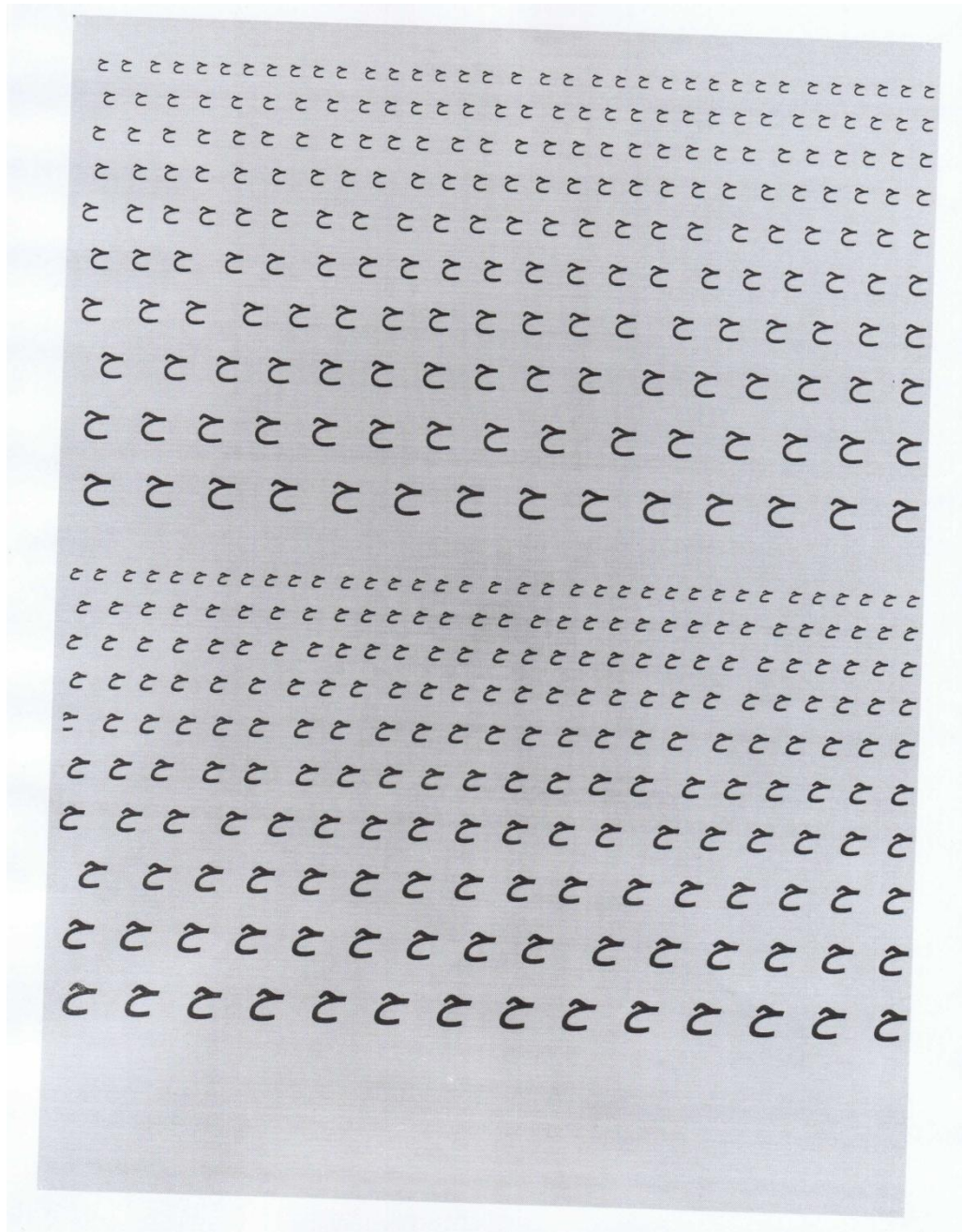
۲-۴- تولید پایگاه داده حروف جدا

در مرحله‌ی نخست با نگارش مناسب هر حرف الفبای فارسی با قلم مذکور مطابق قالب زیر منابع اولیه برای ساخت پایگاه داده فراهم آورده شد که در سطرهای آن از اندازه ۹، ۱۰، ۱۱، ... ۲۴ و سپس ضخیم شده و اندکی مورب شده‌ی آن استفاده شده است که نمونه‌ی آن در شکل ۲-۴ آمده است.



شکل ۲-۴: منابع اولیه تولید شده توسط نرم‌افزار ورود برای ساخت پایگاه داده

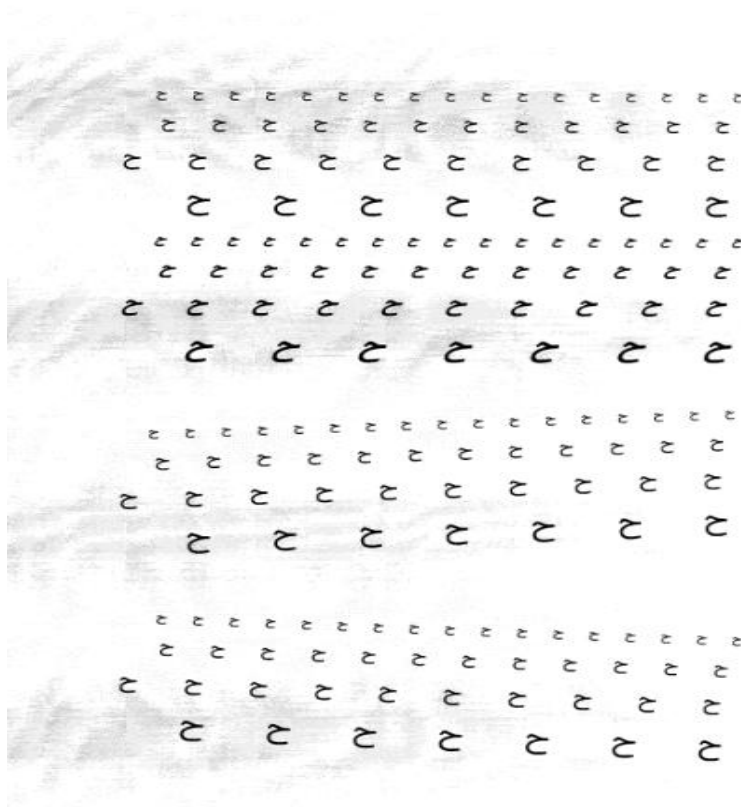
این کار برای حروف دیگر نیز تکرار گردید. سپس در مرحله‌ی بعد نوشته‌های قسمت قبل پرینت گرفته شده و با اندکی چرخش و با عمق ۳۰۰ نقطه بر اینچ اسکن گردید که نتایج آن در شکل ۳-۴ آمده است.



شکل ۳-۴ : تصویر اسکن شده حروف

همان‌طور که مشاهده می‌شود در ابتدا در زمان اسکن چرخش اولیه‌ای معادل ۱۰ درجه در تصویر ایجاد شد و همچنین از رنگ زمینه به منظور ایجاد نویز در تصویر استفاده شد.

جهت شناسایی هرچه بهتر و دقیق تر و با توجه به اینکه پایگاه داده قبلی در مرحله شناسایی برای جملاتی که به صورت زاویه دار نوشته شده بودند با مشکل مواجه می شد، علاوه بر حروف قبل چند سطر به صورت زاویه دار مطابق شکل ۴-۴ ایجاد می شود.



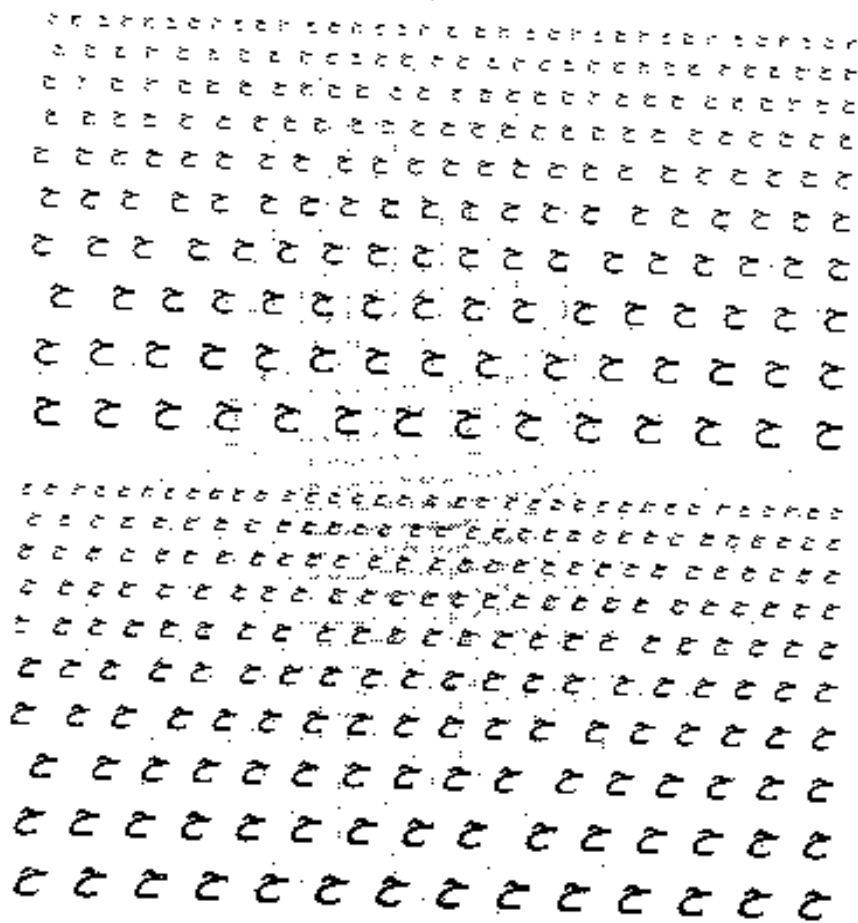
شکل ۴-۴: تولید پایگاه داده متنوع تر

۴-۲-۱- پیش پردازش

پیش پردازش شامل کلیه‌ی اعمالی است که روی تصویر اولیه انجام می شود تا موجب بهتر شدن روند مراحل بعدی شود. از جمله اعمال پیش پردازش می توان به باینری کردن تصویر، حذف نویز، رفع چرخش و غیره اشاره کرد.

۴-۲-۱-۱ باینری کردن تصویر

جهت پردازش مناسب بر روی تصویر اسکن شده، باید تصویر حاصل به صورت یک ماتریس با مقادیر ۰ و ۱ تبدیل شود.



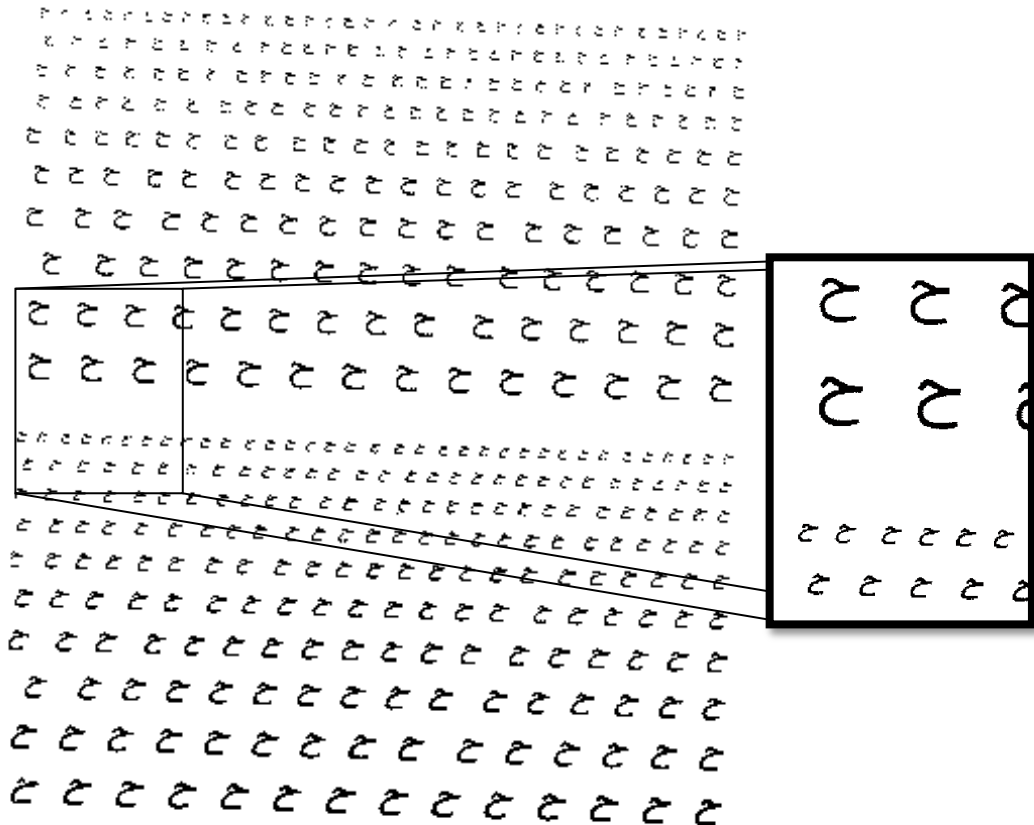
شکل ۴-۵: تصویر باینری شده با نویز فراوان

۴-۲-۱ حذف نویز

با توجه به اینکه زمینه‌ی تصویر رنگی می‌باشد پس از اسکن کردن و باینری کردن آن، نقاط تیره‌ای در تصویر به وجود آمده است؛ لذا برای شناسایی بهتر در مرحله بازشناسی باید به طریق مناسب این نقاط تیره از تصویر حذف شوند. جهت رسیدن به این هدف در این مرحله ابتدا سعی شد با استفاده از روش برچسب‌زنی مؤلفه‌ها و حذف مؤلفه‌های کوچک‌تر از مؤلفه‌های اصلی، نویزها حذف شوند ولی این روش به علت وجود نویز زیاد در تصویر و مدت زمان زیاد پردازش برای برچسب‌زنی همه‌ی مؤلفه‌های نویز استفاده نشد. در نتیجه از دو روش زیر برای حذف اولیه‌ی نویز استفاده شد:

روش اول: ابتدا تقسیم بلوکی تصویر به بلوک‌های k در k انجام شده و سپس با محاسبه‌ی تعداد

پیکسل‌های سیاه در هر بلوک و تعیین سطح آستانه^۱ برای آن به صورت مقدار L ، پی به وجود نقاط نویز در بلوک‌ها برده شده که اگر تعداد آن از مقدار L کمتر باشد، آن نقاط نویز تلقی شده و از تصویر حذف خواهند شد. حاصل این کار در شکل ۴-۶ آمده است.



شکل ۴-۶: حذف نویز از تصویر

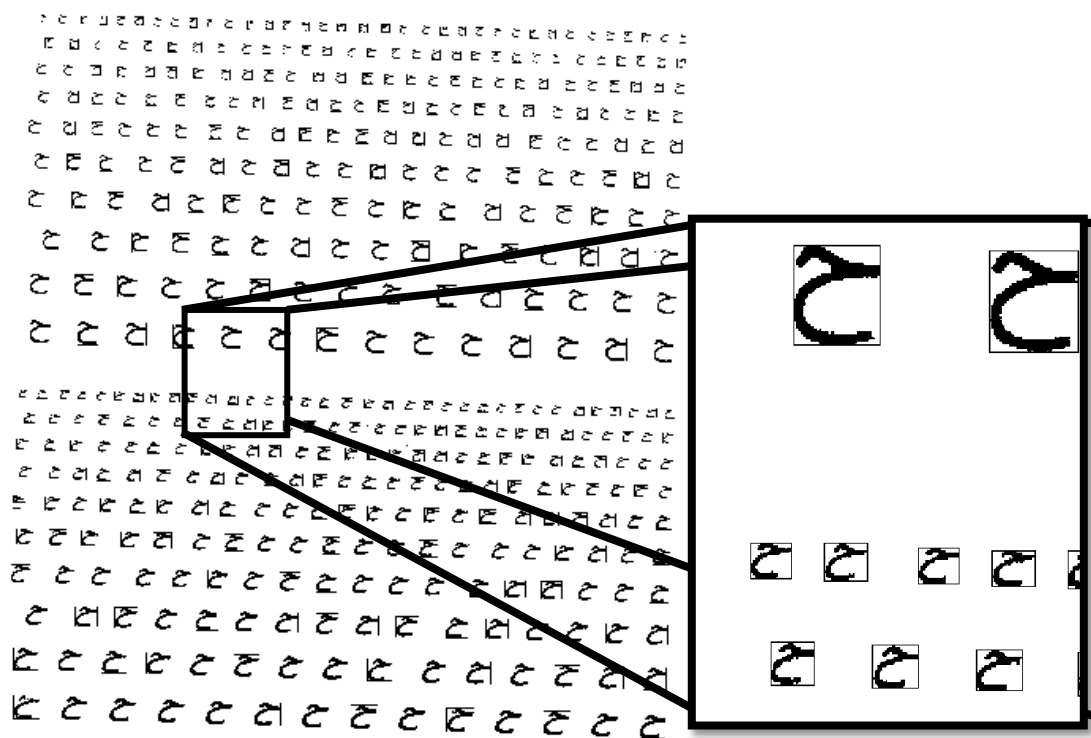
روش دوم: در این روش نیز ابتدا تقسیم بلوکی تصویر به بلوک‌های ۳ در ۳ انجام شده و سپس با محاسبه تفاضل پیکسل‌های موجود در هر بلوک از مقدار پیکسل مرکزی و تعیین سطح آستانه برای آن به صورت مقدار L ، پی به وجود نقاط در بلوک‌ها برده شده که اگر مقدار تفاضل از مقدار L کمتر باشد، آن نقاط نویز تلقی شده و از تصویر حذف خواهند شد. که نتایج حاصل از این روش نیز مشابه روش قبل بود با این تفاوت که در آن بجای معیاری برای تعداد نقاط سیاه در هر بلوک از نوعی معیار

¹ Threshold

گرادیان استفاده شده و سریع‌تر است. پس از انجام این دو روش در نهایت از برچسب‌زنی برای حذف نقاط باقی مانده به عنوان نویز استفاده شد. از این مرحله نمی‌توان قبل از این دو روش استفاده نمود چون قبل از این مراحل با توجه تعداد زیاد مؤلفه‌های نویز مدت زمان زیادی برای حذف آن‌ها در مرحله اول لازم است.

۴-۲-۲- استخراج و جدا کردن حروف

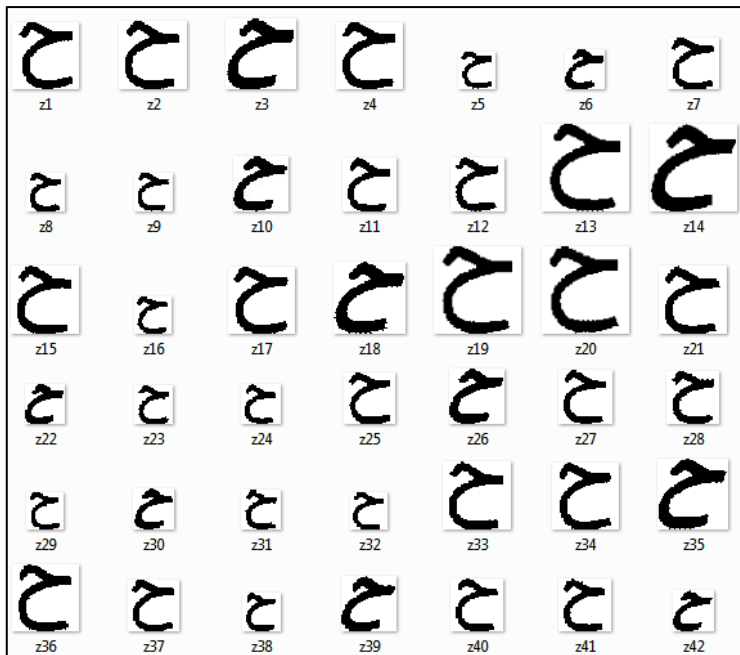
در این مرحله از برچسب‌زنی مؤلفه‌ها برای استخراج حروف استفاده می‌شود، که با استفاده از آن مکان و اندازه‌ی تصویر هر حرف مشخص شده و در متغیری مناسب جهت تولید پایگاه داده ذخیره می‌شود.



شکل ۴-۷: برچسب‌زنی مؤلفه‌ها برای تولید پایگاه داده

۴-۲-۳- شروع عملیات جداسازی

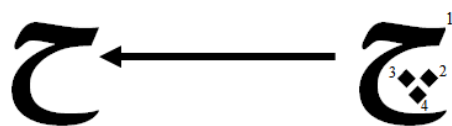
مؤلفه‌های برجسب زده شده در قسمت قبل جهت استانداردسازی، به صورت مربعی ذخیره می‌شوند. به این صورت که با توجه به اندازه هر تصویر به طول یا عرض آن فضای سفید اضافه شده تا مربعی شده و در تصویری مجزا ذخیره می‌شود.



شکل ۴-۸: مربعی کردن و ذخیره‌سازی حروف

۴-۲-۴- حذف نقطه

با توجه به اینکه شناسایی اولیه بر اساس بدنه‌ی حروف (حروف بی‌نقطه) است، برای حذف نقطه‌ی حروف از برجسب‌زنی مؤلفه‌ها به گونه‌ای استفاده می‌شود که در هر برجسب بزرگ‌ترین اندازه، تصویر اصلی شناخته شده و سایر برجسب‌هایی که اندازه‌ی آن‌ها از ۲ برابر عرض قلم کوچک‌تر باشند حذف می‌گردد.



شکل ۴-۹: حذف مؤلفه‌های ۳، ۲ و ۴ برای حذف نقاط

۴-۲-۵- حروف گسسته

در این پایان‌نامه تمام حروف گسسته بسته به شکل بدنه‌ی آن‌ها و برخی از مهم‌ترین علائم و اعداد به ۴۰ کلاس تقسیم شده‌اند. تعدادی از حروف الفبای فارسی مانند «ب»، «پ»، «ت»، «ث» بدنه‌ی مشابهی دارند و تفاوت آن‌ها تنها در تعداد و جای قرار گرفتن نقاط آن‌ها است. در این صورت حروف با بدنه یکسان با یکی از این حروف به عنوان نماینده جایگزین می‌شوند و عمل استخراج ویژگی تنها برای بدنه‌ی حروف (بدون نقطه) انجام می‌شود. برخی کلاس‌ها مانند کلاس ۱۲ تنها شامل ۱ حرف و برخی دیگر مانند کلاس ۱۳ شامل ۴ حرف هستند. البته زیرکلماتی چون لا، لله، سی، صی به دلیل اینکه شکل خاصی دارند و در مرحله‌ی بازشناسی به راحتی قابل جداسازی نیستند، در حروف گسسته به عنوان کلاس‌هایی مجزا در نظر گرفته می‌شوند.

جدول ۴-۱: کلاس‌های مورد استفاده در طبقه‌بند حروف گسسته‌ی فارسی

شماره کلاس	حرف	شماره کلاس	حرف	شماره کلاس	حرف	شماره کلاس	حرف
۱	آ	۲۱	ف	۳۱	لله		
۲	ا	۲۲	ق	۳۲	لا		
۳	ب پ ت ث	۲۳	گ	۳۳	ء		
۴	ج ح خ چ	۲۴	ک	۳۴	.		
۵	د ز	۲۵	ل	۳۵	صی		
۶	ر ز ژ	۲۶	م	۳۶	شی سی		
۷	س ش	۲۷	ن	۳۷	؟		
۸	ص ض	۲۸	و	۳۸	!		
۹	ط ظ	۲۹	ه	۳۹	,		
۱۰	ع غ	۳۰	ی	۴۰	زیرکلمه		

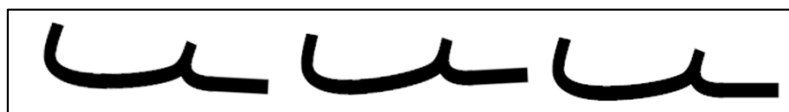
۳-۴- تولید پایگاه داده حروف پیوسته

برخلاف حروف گسسته، حروف پیوسته از لحاظ نحوه‌ی اتصال به حروف دیگر شامل تنوع بیشتری می‌باشد که در جدول زیر این تنوع نمایان شده است. همچنین یک کلاس به عنوان کلاس حروف نامعتبر در نظر گرفته می‌شود که شامل ترکیباتی است که تشکیل یک حرف نمی‌دهند یا شامل بیش از یک حرف هستند. اصولاً حروف پیوسته به سه نوع تقسیم می‌شوند: آغازین، میانه و پایانی که در جدول ۲-۴ نمایش داده شده است.

جدول ۲-۴: کلاس‌های مورد استفاده در طبقه‌بند حروف پیوسته فارسی

شماره کلاس	حرف	شماره کلاس	حرف	شماره کلاس	حرف	شماره کلاس	حرف	شماره کلاس	حرف
۴۱	L	۵۱	گ	۶۱	ک	۷۱	د	۸۱	ط
۴۲	ب	۵۲	گ	۶۲	ک	۷۲	ر	۸۲	ط
۴۳	ح	۵۳	گ	۶۳	ل	۷۳	ص	۸۳	و
۴۴	ع	۵۴	ق	۶۴	ل	۷۴	ص	۸۴	ی
۴۵	ع	۵۵	ح	۶۵	ل	۷۵	ص	۸۵	ی
۴۶	ع	۵۶	ح	۶۶	م	۷۶	ص	۸۶	س
۴۷	د	۵۷	ه	۶۷	م	۷۷	سی	۸۷	لا
۴۸	ف	۵۸	ه	۶۸	م	۷۸	س	۸۸	ناقص
۴۹	ف	۵۹	ه	۶۹	س	۷۹	س		
۵۰	ف	۶۰	ک	۷۰	ر	۸۰	ط		

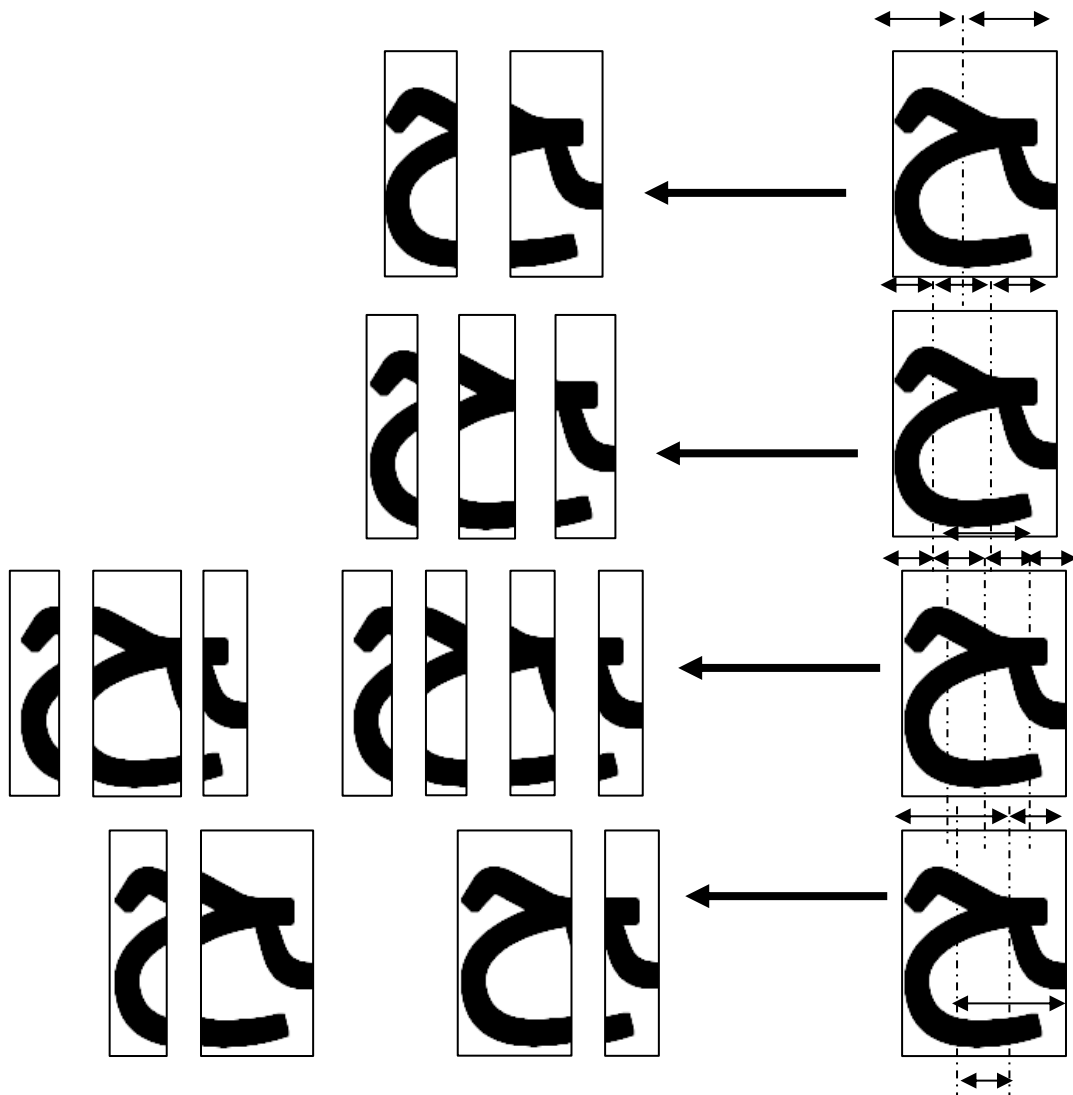
با توجه به اینکه نقاط جداسازی کلمه ممکن است در موقعیت‌های مختلف اتصال دو حرف قرار بگیرد، سعی می‌شود نمونه داده‌های تولید شده با طول‌های مختلفی باشند و همچنین برای تولید نمونه‌های متنوع‌تر از حروف، اندکی چرخش نیز در بعضی نمونه‌ها ایجاد می‌شود.



شکل ۴-۱۰: چرخش ۱- و ۱ درجه برای تولید داده‌های متنوع‌تر

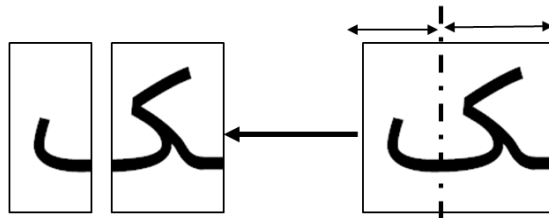
۴-۳-۱- تولید حروف ناقص پیوسته

تولید حروف ناقص بدین علت انجام می‌شود که در زمان جداسازی بتوان صحت جداسازی را رد و یا تایید کرد. برای تشکیل کلاس حروف نامعتبر از ناقص کردن حروف پیوسته‌ی معتبر استفاده می‌شود. کلاس نامعتبر شامل تمام حروف و زیرحروفی است که تشکیل یک حرف نمی‌دهند که نحوه‌ی تولید آن در شکل ۴-۱۱ نمایش داده شده است. مثلاً حرفی مانند "ح" را ابتدا به قسمت‌های مختلف تقسیم کرده و ترکیبی از این قسمت‌ها به عنوان یک زیرکلمه‌ی نامعتبر شناخته شده و در کلاس حروف ناقص قرار داده می‌شود.



شکل ۴-۱۱: تولید حروف ناقص و نامعتبر

در هنگام تولید حروف ناقص با الگوریتم بالا مشاهده شد که برخی از اشکال تولید شده خود عضوی از حروف پیوسته معتبر است، در نتیجه باید از حروف ناقص حذف شوند مانند حرف کاف. همان طور که مشاهده می شود حرف کاف تولید شده دیگر ناقص نبوده و باید از این گروه حذف شود.



شکل ۴-۱۲: حروف ناقص تولید شده معتبر

۴-۴- استخراج ویژگی ها

پس از تولید داده برای کلاس های ۱ تا ۴۰ و ۴۱ تا ۸۸، عملیات استخراج ویژگی برای داده های هر کدام از کلاس ها به صورت مستقل انجام می شود. استخراج ویژگی شامل کلیه محاسباتی است که روی نمونه های به دست آمده از مرحله ی پیش پردازش انجام می شود تا بردار ویژگی های متناظر با هر نمونه حاصل شود.

در اینجا برای استخراج ویژگی از روش هیستوگرام گرادیان استفاده شده است. این ویژگی که برای بازشناسی ارقام دست نویس فارسی معرفی شده است [۴۱]، با انجام تغییراتی برای زیرکلمات فارسی به کار گرفته شد. برای استخراج ویژگی ابتدا تصویر به اندازه ی $N \times N$ نرمال می شود، سپس گرادیان تصویر با استفاده از ماسک های سوبل محاسبه می شود. در این روش g_x گرادیان در جهت افقی و g_y گرادیان در جهت عمودی با استفاده از دو ماسک مجزا به دست می آید. برای هر نقطه از تصویر با استفاده از ماسک های فوق بردار گرادیان به صورت $[g_x \ g_y]$ محاسبه می گردد. شدت و زاویه ی گرادیان برای هر نقطه از تصویر، با به دست آوردن اندازه و جهت بردار فوق به دست می آید. پس از به دست آوردن زاویه ی گرادیان لازم است مقادیر زاویه به چند زاویه ی استاندارد تبدیل گردند. نحوه ی استخراج ویژگی برای زیرکلمات فارسی به صورت زیر است:

(۱) نرمال کردن تصویر ورودی به اندازهی 60×60 به گونه ای که نسبت ابعاد تصویر حفظ شود.

(۲) اعمال فیلتر سوبل که به صورت زیر تعریف می شود:

$$g_x(x, y) = f(x+1, y-1) + 2f(x+1, y) + f(x+1, y+1) - f(x-1, y-1) - 2f(x-1, y) - f(x-1, y+1)$$

$$g_y(x, y) = f(x-1, y+1) + 2f(x, y+1) + f(x+1, y+1) - f(x-1, y-1) - 2f(x, y-1) - f(x+1, y-1)$$

رابطه (۴-۱)

از روی این دو مؤلفه‌ی گرادیان، اندازه و جهت گرادیان محاسبه می شود. برای ترکیب g_x و g_y ، 16

حالت اتفاق می افتد که شامل زوایای $0, \pm\frac{\pi}{8}, \pm\frac{\pi}{4}, \pm\frac{3\pi}{8}, \pm\frac{\pi}{2}, \pm\frac{5\pi}{8}, \pm\frac{3\pi}{4}, \pm\frac{7\pi}{8}, \pm\pi$ شده و یک

حالت زمانی رخ می دهد که g_x و g_y هر دو صفر بوده و زاویه تعریف نشده وجود داشته باشد.

(۳) تصویر ورودی به 16 بلوک مساوی 15×15 پیکسلی تقسیم می شود.

(۴) در هر بلوک مقدار شدت گرادیان به ازای هر 16 زاویه‌ی ذکر شده با هم جمع می شوند. به ازای

زاویه‌ی تعریف نشده چون شدت گرادیان همواره صفر می شود و ویژگی خاصی به دست نمی دهد در

نظر گرفته نمی شود. به این ترتیب در نهایت بردار ویژگی شامل $16 \times 4 \times 4$ یعنی 256 مؤلفه می شود.

عمل استخراج ویژگی روی 300 تصویر مجزای موجود در هر کلاس انجام شده و برای آموزش

طبقه‌بندها استفاده می شود.

۴-۵- آموزش طبقه‌بندها

از ویژگی‌های استخراج شده در مرحله‌ی قبل برای آموزش طبقه‌بندها استفاده می شود.

طبقه‌بندهای مورد استفاده در این پایان‌نامه از نوع شبکه‌های عصبی پرسپترون چند لایه^۱ می باشد. دو

نوع طبقه‌بند، یکی برای طبقه‌بندی حروف گسسته و زیرکلمات و دیگری برای حروف پیوسته آموزش

^۱ MLP

داده خواهد شد، که وظیفه‌ی هر یک از آن‌ها مقایسه‌ی کلاس ورودی و نمونه‌ی مورد آزمایش و تشخیص کلاس است. طبقه‌بند گسسته شامل ۴۰ کلاس و طبقه‌بند پیوسته حاوی ۴۸ کلاس است.

۴-۵-۱- طبقه‌بند حروف گسسته و زیرکلمات

همان‌طور که در جدول ۴-۱ نشان داده شد، طبقه‌بند حروف گسسته در ۴۰ کلاس آموزش داده می‌شود، که وظیفه‌ی آن تشخیص کلاس ورودی حروف گسسته و زیرکلمات است. در این طبقه‌بند کلاس‌های ۱ تا ۳۹ برای حروف گسسته بوده و کلاس ۴۰ برای زیر کلمات می‌باشد. برای استخراج ویژگی کلاس ۴۰ از ۶۶۵۰ زیرکلمه معتبر و موجود در زبان فارسی استفاده شده است. نمونه‌ای از این زیرکلمات در شکل ۴-۱۳ آمده است.

سل	مگا	کما	سها	لم	هم	بد	ها
IM20104	IM20105	IM20106	IM20107	IM20108	IM20109	IM20110	IM20111
لی	سس	حد	به	سو	حبا	حو	ب
IM20112	IM20113	IM20114	IM20115	IM20116	IM20117	IM20118	IM20119
همه	حر	حد	نو	وا	هر	وفا	کر
IM20120	IM20121	IM20122	IM20123	IM20124	IM20125	IM20126	IM20127
بکر	لب	کو	سا	که	یکی	فب	س
IM20128	IM20129	IM20130	IM20131	IM20132	IM20133	IM20134	IM20135
سب	عر	حس	حب	سا	بما	بی	م
IM20136	IM20137	IM20138	IM20139	IM20140	IM20141	IM20142	IM20143

شکل ۴-۱۳: نمونه‌ای از زیرکلمات معتبر و موجود در زبان فارسی

روند کار بدین ترتیب است که اگر ورودی به صورت زیرکلمه (کلاس ۴۰) شناسایی شود، برای جداسازی به مرحله‌ی بعد فرستاده می‌شود و اگر ورودی حرف یا علامت باشد، یکی از کدهای ۱ تا ۳۹ فعال می‌شوند.

۴-۵-۲- طبقه‌بند حروف پیوسته

طبقه‌بند حروف پیوسته مطابق جدول ۴-۲ شامل ۴۸ کلاس (کلاس‌های ۴۱ تا ۸۸) است. از این طبقه‌بند برای جداسازی زیرکلمات کلاس ۴۰ استفاده می‌شود. روش کار بدین صورت است که ابتدا زیرکلمات کلاس ۴۰ توسط طبقه‌بند حروف گسسته شناسایی شده و برای جداسازی به طبقه‌بند حروف پیوسته فرستاده می‌شود. در این مرحله با اعمال روش‌های مناسب، نقاط کاندید اولیه‌ی جداسازی مشخص شده و صحت این نقاط جداسازی با استفاده از طبقه‌بند حروف پیوسته مشخص می‌شود. اگر نقطه‌ی جداسازی معتبر باشد یکی از کلاس‌های ۴۱ تا ۸۷ انتخاب می‌شود و اگر نقطه‌ی جداسازی نامعتبر باشد کلاس ۸۸ انتخاب شده و نقطه‌ی جداسازی جدیدی انتخاب می‌شود.

۴-۶- نتیجه‌گیری

در این فصل با توجه به موجود نبودن پایگاه داده‌ی مناسب برای قلم Iranian sans اقدام به تولید پایگاه داده‌ای مناسب جهت آموزش طبقه‌بندهای حروف پیوسته و گسسته شد. در این فرایند پس از تولید سندهای مناسب و اسکن نمودن آن‌ها عملیات پیش‌پردازش بر روی آن‌ها انجام گرفت و با استفاده از روش هیستوگرام‌گرادیان ویژگی‌های لازم جهت آموزش طبقه‌بندها استخراج شد. طبقه‌بندهای مورد استفاده در این پایان‌نامه از نوع شبکه‌ی عصبی Mlp می‌باشد.

در فصل آتی از طبقه‌بند حروف گسسته برای طبقه‌بندی حروف جدا و شناسایی زیرکلمات استفاده کرده و از طبقه‌بند حروف پیوسته برای جداسازی زیرکلمات شناسایی شده در مرحله‌ی قبل استفاده خواهد شد.

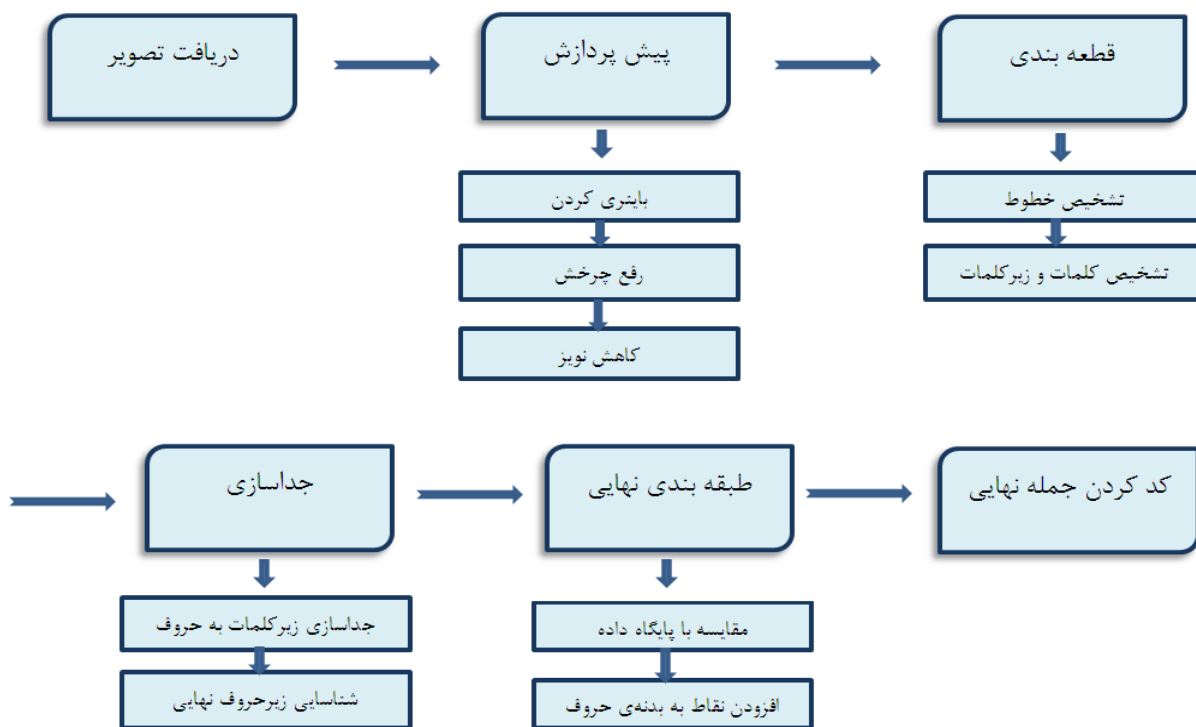
فصل پنجم

سامانه بازشناسی متون

Iranian sans با قلم

در فصل گذشته به تولید پایگاه داده جهت آموزش طبقه‌بندها پرداخته شد. بدین صورت که یک طبقه‌بند برای حروف گسسته و زیرکلمات و یک طبقه‌بند برای حروف پیوسته آموزش داده شد.

در این فصل به توضیح روند اجرای سیستم بازشناسی متن پرداخته می‌شود. همان‌طور که در شکل زیر نشان داده شده است، روند کار بدین صورت است که در مرحله اول پس از دریافت تصویر پیش‌پردازش‌هایی شامل باینری کردن، اصلاح چرخش و حذف نویز بر روی تصویر اعمال می‌شود و در مرحله جداسازی با شناسایی حدود مکانی جمله و با استفاده از روش‌های مناسب و رفع مشکلات همپوشانی به جداسازی حروف گسسته و زیرکلمات پرداخته می‌شود. با استفاده از طبقه‌بند حروف گسسته، حروف مجزا از زیر کلمات جدا شده و با توجه به شماره‌ی کلاس آن ذخیره می‌شود. زیرکلمات شناسایی شده در این مرحله به قسمت جداسازی حروف رفته و با روش‌های مختلف این جداسازی انجام می‌شود. سپس با استفاده از طبقه‌بند حروف پیوسته به شناسایی حروف پرداخته می‌شود. در نهایت این حروف در کنار یکدیگر قرار گرفته و کد می‌شوند و کلمه و در نهایت جمله بازیابی می‌شود.



شکل ۵-۱: شمای کلی سیستم بازشناسی متن

۵-۲- بازشناسی متن

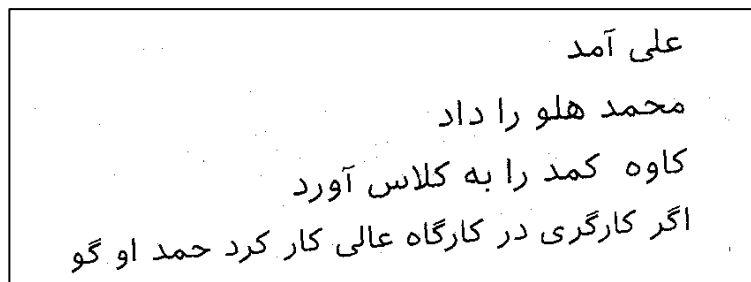
هر یک از اجزای متن با یک بردار ویژگی مشخص می‌شود. هر سیستم بازشناسی متن، این بردارها را در مرحله‌ی آموزش فرا می‌گیرد و به عنوان داده‌های معتبر در حافظه خود نگهداری می‌کند (آموزش طبقه‌بند). برای بازشناسی تصویر یک جزء متن، بردار ویژگی متناظر با آن با همان روش استخراج در مرحله‌ی آموزش، از این تصویر ناشناخته استخراج می‌شود. این بردار ویژگی با داده‌های سیستم مقایسه و بازشناسی می‌شود.

۵-۲-۱- پیش پردازش

عملیات پیش پردازش فرآیندی برای ارتقای تصویر ورودی است که آن را برای تحلیل و بازشناسی در مراحل بعد مهیا می‌سازد. از مهم‌ترین مراحل پیش پردازش می‌توان به دوسطحی سازی، بهبود تصویر با حذف نویز و اصلاح چرخش نام برد.

۵-۲-۱-۱ باینری کردن تصویر

عموماً تصاویر پس از اسکن شدن به صورت خاکستری یا رنگی هستند با توجه به اینکه بسیاری از الگوریتم‌های پردازش تصویر از شکل دوسطحی تصاویر استفاده می‌کنند، نیاز است تصویر خاکستری به دوسطحی تبدیل شود. در اینجا با استفاده از یک سطح آستانه ثابت تصویر ورودی به دوسطحی تبدیل می‌شود.



شکل ۵-۲: باینری کردن تصویر

۵-۲-۱-۲ کاهش نویز

با توجه به اینکه تصاویر استفاده شده از کاغذ و توسط اسکنر پدید آمده است، به دلایل مختلف (که برخی از آن‌ها در قبل توضیح داده شد) دارای نویز می‌باشد. با توجه به اینکه حذف نویز گام موثری در بهبود بازشناسی دارد از روشی مشابه فصل گذشته برای حذف نویز استفاده می‌شود که شامل دو مرحله است. در مرحله اول با روش میانگین‌گیری بلوکی و تفاضل پیکسل مرکزی از پیکسل‌ها مقداری از نویز کاسته می‌شود و سپس در مرحله دوم با استفاده از روش برچسب‌زنی مؤلفه‌ها نویز به صورت کامل حذف می‌شود. استفاده از مرحله دوم به صورت مستقیم به دلیل وجود نویز فراوان نیاز به پردازش زیادی داشته و زمان‌بر است؛ لذا ابتدا در مرحله اول مقداری از نویز حذف شده و در مرحله دوم نویز به صورت کامل حذف می‌شود.

علی آمد
محمد هلو را داد
کاوه کمد را به کلاس آورد
اگر کارگری در کارگاه عالی کار کرد حمد او گو

شکل ۳-۵: حذف نویز در دو مرحله

۳-۱-۲-۵ اصلاح چرخش

تصویر اغلب اوقات توسط اسکنرها دریافت می‌شود که ممکن است به دلیل خطای انسانی در گذاشتن صفحات، تصویر سند با محورهای تصویر تراز نباشد یا ممکن است تصویر روبروی دوربین درست قرار نگیرد که باعث ایجاد زاویه‌ای در تصویر می‌شود و می‌تواند در نرخ بازشناسی تأثیر قابل ملاحظه‌ای بگذارد. برای رفع چرخش در تصویر دریافتی از افکنش افقی تصویر استفاده شده است. پس از یافتن زاویه‌ی چرخش، تصویر با همان زاویه و در خلاف جهت چرخانده می‌شود تا چرخش تصویر حذف شود.

علی آمد
محمد هلو را داد
کاوه کمد را به کلاس آورد
اگر کارگری در کارگاه عالی کار کرد حمد او گو

شکل ۴-۵: اصلاح چرخش به روش افکنش افقی

۲-۲-۵-۲-۲-۵ قطعه‌بندی

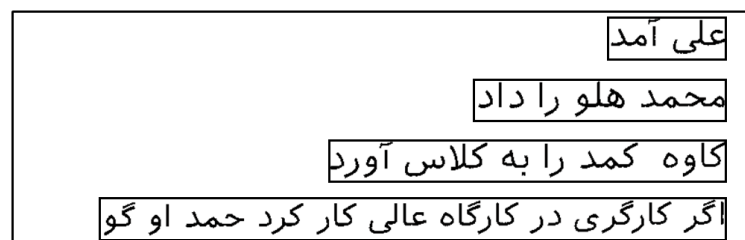
قطعه‌بندی مرحله‌ای مهم است که در آن اجزای تصویر متنی که باید به مرحله شناسایی تحویل داده شوند، از یکدیگر جدا می‌شوند. قطعه‌بندی در سه سطح انجام می‌شود:

- تشخیص خطوط
- تشخیص کلمات یا زیرکلمات
- تشخیص کاراکترها

۵-۲-۱-۲-۱ تشخیص خطوط

متون فارسی اغلب به صورت مجموعه‌ای از خطوط موازی هستند که با فاصله نوشته می‌شوند. الگوریتم جداسازی خطوط، معمولاً از این فاصله‌ی طبیعی استفاده می‌کند. هیستوگرام افقی با شمارش عناصر سیاه تصویر در هر ردیف محاسبه شده و ردیف‌هایی که مقدار هیستوگرام در آن‌ها صفر یا از حد آستانه‌ای کمتر باشد به عنوان ردیف‌های بین خطوط در نظر گرفته می‌شوند. در واقع هر ردیفی از متن که بیشترین فراوانی را داشته باشد نمایانگر خطوط اصلی متن است.

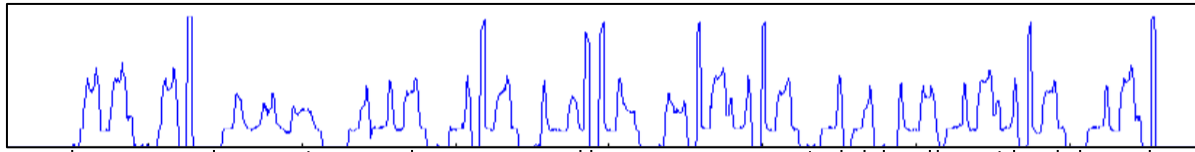
با استفاده از این روش، نقاط و علائمی که بین آن‌ها و خطوط اصلی متن ردیف‌های خالی وجود دارد، به عنوان یک خط جدید شناسایی می‌شوند. این خط با در نظر گرفتن حد آستانه‌ای برای ارتفاع هیستوگرام افقی خطوط متن برطرف می‌گردد و بخش‌هایی که هیستوگرام افقی آن‌ها از حد آستانه‌ای کمتر باشد به نزدیک‌ترین خط مجاور خود ملحق می‌شوند.



شکل ۵-۵: قطعه بندی جمله‌ها

۵-۲-۲-۲ تشخیص کلمات و زیرکلمات

بعد از جداسازی خطوط، سیستم باید هر خط را به زیرکلمات تشکیل دهنده‌ی آن تفکیک نماید. بدین صورت که کلمات و کلمات فرعی درون خط با استفاده از هیستوگرام عمودی تعیین می‌گردد. در این حالت ستون‌هایی که مجموع پیکسل‌های آن صفر هستند محل جداسازی کلمات و زیرکلمات را تعیین می‌کند. اشکال این روش این است که اگر چند زیرکلمه با هم همپوشانی داشته باشند (شکل ۵-۶)، کل آن‌ها یک زیرکلمه در نظر گرفته می‌شوند، چون هیستوگرام صفر ایجاد نمی‌شود.



اگر کارگری در کارگاه عالی کار کرد حمد او گو

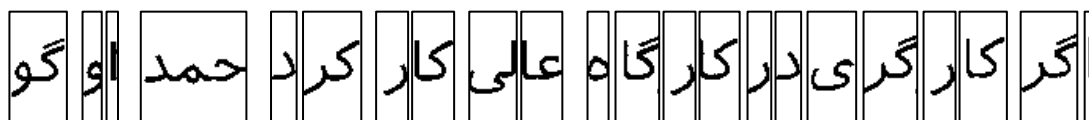
اگر کارگری در کارگاه عالی کار کرد حمد او گو

شکل ۵-۶: مشکل وجود همپوشانی در جداسازی زیر کلمه‌ها به روش هیستوگرام عمودی

همان‌طور که مشاهده می‌شود به دلیل وجود همپوشانی کلمات، زیر کلمات به درستی تفکیک نمی‌شوند. در این حالت برای برطرف نمودن این مشکل از برچسب‌زنی مؤلفه‌ها استفاده می‌شود. در بدترین حالت این همپوشانی بین سه زیر کلمه متوالی روی می‌دهد. در این حالت ابتدا برچسب‌زنی مؤلفه‌ها انجام شده و سپس درصد همپوشانی مؤلفه‌ها با مؤلفه‌ی اصلی محاسبه می‌شود. این کار برای آن است که ابتدا تشخیص داده شود این مؤلفه جزئی از خود زیر کلمه است (مانند نقطه و یا سرکش گاف و علامت مد در آ) و یا زیر کلمه‌ی مجاور است. در حقیقت با قواعد مناسبی نقاط و علائم به بدنه‌ی زیر کلمه نسبت داده می‌شود. در صورتی که این همپوشانی نزدیک به ۱۰۰٪ باشد، معلوم می‌شود که این مؤلفه جزئی از زیر کلمه است و در غیر این صورت این مؤلفه زیر کلمه‌ی مجاور بوده که دارای همپوشانی شده است. در این حالت ابتدا تصویر مؤلفه‌ی مرکزی و یا به عبارتی بدنه‌ی زیر کلمه را از تصویر جدا کرده تا عامل همپوشانی حذف گردد و با استفاده از هیستوگرام عمودی جداسازی مؤلفه اول و آخر نیز انجام می‌شود.



شکل ۷-۵: برطرف نمودن مشکل همپوشانی با استفاده از برچسب زنی مناسب مؤلفه‌ها



شکل ۸-۵: برطرف نمودن مشکل همپوشانی

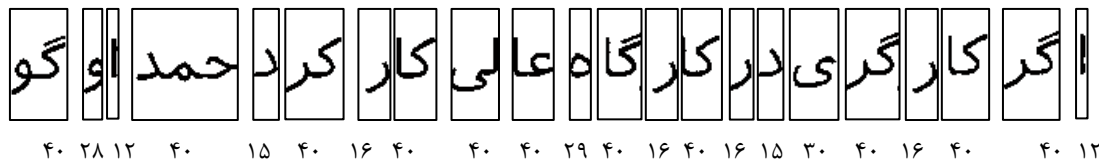
بدین صورت عمل جداسازی به زیرکلمات و حروف گسسته انجام می‌شود. حال قطعات تولید شده‌ی مذکور حذف نقطه می‌شوند. در این فرایند از عرض قلم به عنوان معیار مناسبی برای حذف نقطه استفاده می‌شود.

۵-۲-۲-۱ عرض قلم

برای یافتن عرض قلم کافیست در هر ستون تصویر از برچسب‌زنی مؤلفه‌ها استفاده کرد که در این صورت طول عمودی هر مؤلفه (تعداد پیکسل‌های سیاه) حاصل می‌شود. برداری تشکیل داده می‌شود که هر درایه‌ی آن بیانگر تعداد تکرارهای طول می‌باشد. طولی که تعداد تکرار آن از همه بیشتر باشد به عنوان عرض قلم در نظر گرفته می‌شود. بدین ترتیب با محاسبه‌ی عرض قلم و برچسب‌زنی مؤلفه‌ها، مؤلفه‌هایی که دارای طول و عرضی کمتر از عرض قلم باشد از تصویر حذف می‌شود تا برای شناسایی به قسمت طبقه‌بند حروف گسسته فرستاده شود.

۵-۲-۲-۲ شناسایی حروف گسسته و زیرکلمات

پس از برطرف نمودن مشکل همپوشانی و حذف نقاط، حروف و زیرکلمات تولید شده با استفاده از طبقه‌بند حروف گسسته شناسایی می‌شود. بدین ترتیب که حروف گسسته با توجه به جدول ۴-۱، از شماره ۱ تا ۳۹ کلاسه‌بندی شده و همچنین تصاویر شناسایی شده که در کلاس ۴۰ (زیرکلمه) قرار می‌گیرند جهت جداسازی به قسمت جداسازی فرستاده می‌شود.



شکل ۵-۹: شماره کلاس مطابق جدول ۴-۱

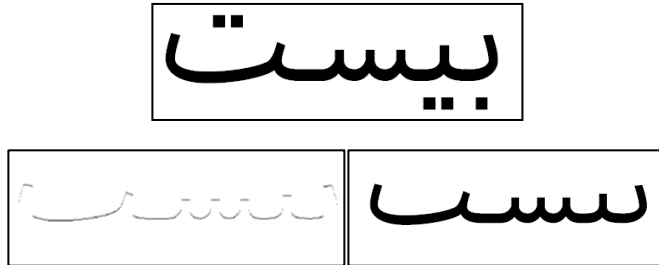
۵-۲-۳- جداسازی

به طور کلی برای بازشناسی متون سه رویکرد عمده وجود دارد: روش‌های مبتنی بر جداسازی، روش‌های مبتنی بر شکل کلی و روش‌های ترکیبی. در این پایان نامه از رویکرد مبتنی بر جداسازی استفاده شده است. جداسازی به این معنی است که یک زیرکلمه به حروف سازنده‌اش تجزیه شود. برای یافتن نقاط جداسازی روش‌های مختلفی وجود دارد که به دلیل پاره‌ای از محدودیت‌ها از جمله همپوشانی حروف، شباهت بسیاری از زیرحروف به حروف و عوامل دیگر، تعیین نقاط دقیق جداسازی کار ساده‌ای نیست. در این روش برای جداسازی حروف از روش پروفایل بالایی تعمیم یافته استفاده شده است. برای اینکه نقاط جداسازی مطمئن باشد، باید ابتدا نقاط اولیه‌ی جداسازی تعیین شود، سپس به وسیله‌ی موتور بازشناسی (طبقه‌بند حروف پیوسته) که در مرحله‌ی تولید داده آموزش داده شد، نقاط نامعتبر حذف شده و نتایج جداسازی بهبود می‌یابد.

۵-۲-۴- جداسازی حروف به روش پروفایل بالایی تعمیم یافته

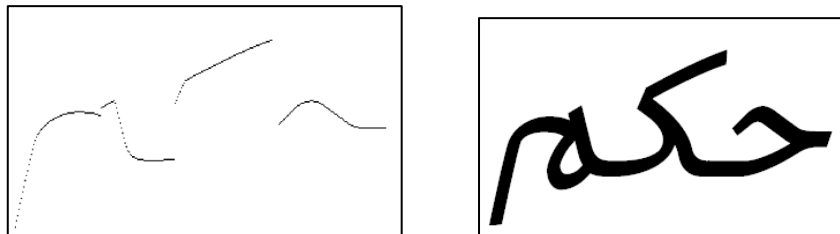
در راستای جداسازی حروف به روش پروفایل بالایی پس از حذف نقاط و علائم ابتدا با استفاده از روش گذر از سفید به سیاه به تشکیل پروفایل بالایی و لبه‌های بالایی تصویر پرداخته می‌شود. روش کار بدین صورت است که پس از تشکیل پروفایل بالایی و محاسبه‌ی مینیمم‌های محلی و استفاده از روش برجسبزی مؤلفه‌ها به بررسی نقاط کاندید برای جداسازی پرداخته می‌شود. به عنوان مثال

برای زیرکلمه‌ی "بیست" در مرحله‌ی اول پس از حذف نقاط به محاسبه‌ی پروفایل بالایی پرداخته می‌شود.



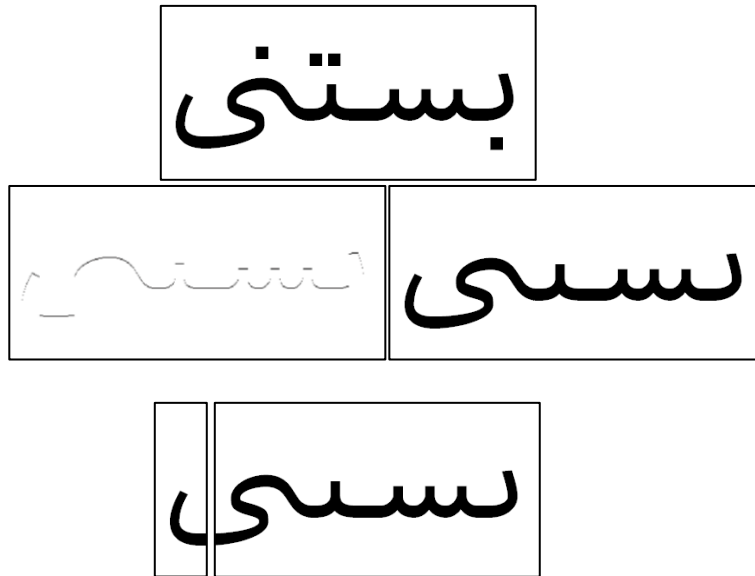
شکل ۵-۱۰: حذف نقاط، محاسبه لبه بالایی، تعیین مینیمم‌های محلی کاندید برای جداسازی

در استفاده از این روش به تنهایی چند مشکل عمده وجود دارد. اول آنکه با توجه به قرار گرفتن بعضی حروف در زیر حروف دیگر (همپوشانی حروف) پس از اعمال روش پروفایل بالایی نمی‌توان عمل جداسازی را به درستی انجام داد. مثلاً در کلمه‌ای مانند حکم به علت وجود سرکش حرف "ک" که روی حرف "ح" را پوشانده، محاسبه‌ی مینیمم محلی برای جداسازی، عملی غیر ممکن است و منجر به جداسازی نادرست می‌شود.



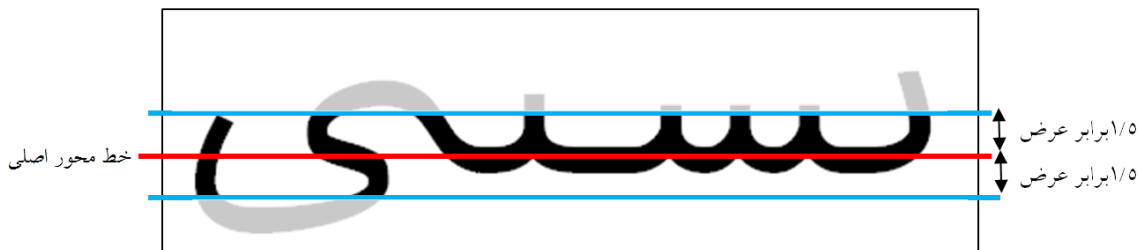
شکل ۵-۱۱: عدم وجود مینیمم محلی بین حرف "ح" و "ک" به دلیل همپوشانی

از طرفی با توجه به آنکه نقطه‌ی جداسازی باید در محدوده‌ی خط زمینه باشد، برای تعیین نقاط جداسازی مناسب، باید محدوده‌ای مشخص شود. در غیر این صورت تعیین نقاط جداسازی به درستی انجام نمی‌شود. مثلاً مطابق شکل ۵-۱۲ حرف "ی" به صورت نامناسب جدا خواهد شد.



شکل ۵-۱۲: جداسازی نامناسب حرف ی به دلیل وجود محدوده مناسب

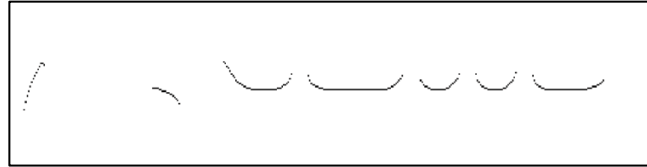
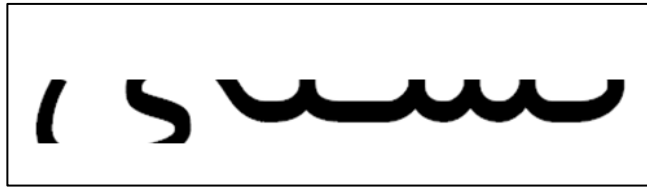
با توجه به مشکلات مذکور باید ابتدا محدوده ای مناسب برای جداسازی مشخص شود. به این ترتیب که ابتدا خط اصلی متن با استفاده از هیستوگرام افقی محاسبه شده و سپس با محاسبه‌ی عرض قلم، محدوده‌ای به اندازه یک و نیم برابر عرض قلم پایین و بالای خط کرسی، به عنوان محدوده‌ی مجاز برای محاسبه‌ی نقاط جداسازی در نظر گرفته شده و پروفایل بالایی فقط برای این محدوده از تصویر در نظر گرفته می‌شود. بدین ترتیب علاوه بر برطرف نمودن مشکل همپوشانی، از جداسازی نامناسب در خارج محدوده‌ی خط مبنا جلوگیری به عمل خواهد آمد.



شکل ۵-۱۳: مشخص نمودن محدوده مجاز برای محاسبه پروفایل بالایی و محل جداسازی

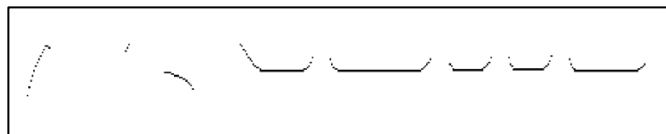
پس از این مرحله محاسبه‌ی پروفایل بالایی برای تصویر حاصل انجام می‌شود که حاصل این کار

در شکل ۵-۱۴ آمده است.



شکل ۱۴-۵: محاسبه پروفایل بالایی تصویر

پس از این مرحله جهت محاسبه‌ی کاندیدهای مناسب جداسازی در محدوده‌ی مینیمم محلی، مکان‌هایی که به اندازه یک پیکسل با مرکز مینیمم محلی فاصله دارند را در یک راستا قرار داده تا در مرحله‌ی بعد نقاط جداسازی بهینه‌تری حاصل شود.



شکل ۱۵-۵: بهینه نمودن مینیمم‌های محلی

سپس در این مرحله با اعمال هیستوگرام افقی، سه سطری که بیشترین فراوانی‌ها را داشته باشند در نظر گرفته می‌شود. هر سطر را به صورت برداری ذخیره نموده که حاوی تکه خط‌های افقی می‌باشد. سرانجام با استفاده از برجسب‌زنی مؤلفه‌ها وسط تکه خط‌های افقی ایجادشده به عنوان کاندید اولیه‌ی جداسازی در نظر گرفته می‌شوند و با این شرط که فاصله‌ی این نقاط از هم نباید از عرض قلم کمتر باشد، این نقاط به عنوان نقاط جداسازی خواهند بود. بدین ترتیب اگر فاصله‌ی دو نقطه‌ی جداسازی از عرض قلم کمتر بود، آن که امتداد بزرگ‌تری دارد، نقطه‌ی جداسازی در نظر گرفته می‌شود.

سیسی

سیسی

سیسی

سیسی

سیسی

سیسی

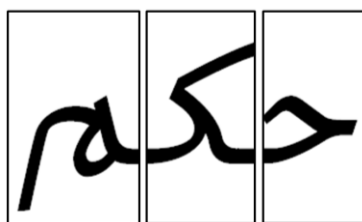
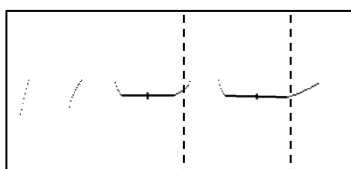
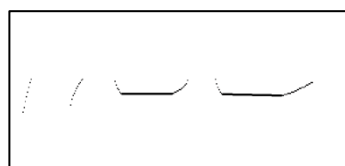
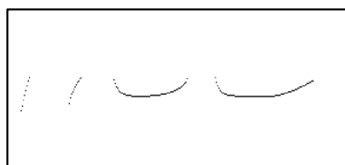
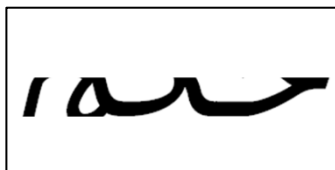
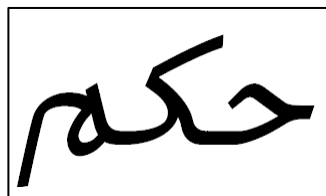
سیسی

شکل ۵-۱۶: تعیین نقاط جداسازی اولیه

بدین ترتیب نقاط جداسازی اولیه تعیین می‌شود و حروف و زیرحروف ایجادشده به شبکه‌ی

بازشناسی حروف پیوسته سپرده می‌شود تا صحت این جداسازی تعیین و در نهایت کلاسه‌بندی شده

و به مرحله‌ی بعد فرستاده شود. در حقیقت با مشخص کردن ناحیه‌ی مجاز برای جداسازی مطابق شکل ۱۷-۵ هم مشکل همپوشانی (شکل ۵-۱۱) و هم مشکل جداسازی حروفی مانند "ی" که نقاط جداسازی غیرمجاز (شکل ۵-۱۲) تولید می‌کند، حل می‌شود.

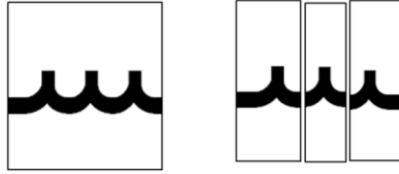


شکل ۵-۱۷: برطرف کردن مشکل همپوشانی

۱-۴-۲-۵ شناسایی زیرحروف نهایی

در مرحله‌ی قبل نقاط جداسازی اولیه مشخص شده و تصویر به تعدادی زیرکلمه تقسیم می‌شود.

اما هنوز معتبر بودن نقطه‌ی جداشده مشخص نشده است. مثلاً مطابق شکل ۵-۱۸ حرف "س" به تعدادی دندان تقسیم شده است که این جداسازی معتبر نمی‌باشد.



شکل ۵-۱۸: جداسازی نادرست حرف سین در کلمه بستنی

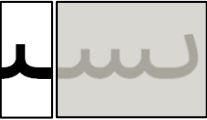


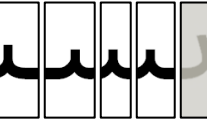

بنابراین لازم است صحت جداسازی مرحله‌ی قبل بررسی شود. در این مرحله ابتدا جداسازی حروف کلمه‌ی "بستنی" مطابق شکل ۵-۱۵ و با استفاده از روش پروفایل بالایی انجام می‌شود، سپس از سمت چپ قطعات تولید شده در کنار یکدیگر قرار گرفته و تشکیل یک تصویر جدید می‌دهد. سپس تصویر حاصله برای شناسایی به شبکه بازشناسی حروف پیوسته فرستاده شده و اعتبارسنجی می‌شود. این کار برای چهار تصویر اول انجام می‌شود و اعتبار هر مرحله ثبت می‌شود.

جدول ۵-۱: شناسایی و اعتبارسنجی نقاط جداسازی (مرحله اول)

شماره مرحله	کلاس شناسایی شده مطابق جدول ۴-۲	دقت شناسایی	تصویر حاصل
اول	کلاس ۸۵ (زیر بدنه "یی")	%۹۹	
دوم	کلاس ۸۸ (حرف ناقص)	%۹۰	
سوم	کلاس ۸۸ (حرف ناقص)	%۹۵	
چهارم	کلاس ۸۸ (حرف ناقص)	%۹۸	
نتیجه	کلاس ۸۵ (زیر بدنه "یی")	%۹۹	

در مرحله‌ی اول جداسازی مشاهده می‌شود که با دقت ۹۹٪ محل جداسازی درست بوده و مطابق جدول ۲-۴ کلاس ۸۵ شناخته می‌شود. در مرحله‌ی دوم با اضافه نمودن قطعه‌ی بعدی به این تصویر، تصویر جدید با دقت ۹۰٪ نامعتبر و جزو حروف ناقص شناخته می‌شود. این کار تا ۴ مرحله انجام می‌شود و با دقت ۹۵ و ۹۸ درصد، تصویر ناقص تشخیص داده می‌شود. بنابراین در مجموع این ۴ مرحله، تصویر مرحله اول به عنوان کلاس ۸۵ شناخته می‌شود و از تصویر اصلی جدا شده و این روند برای باقیمانده‌ی تصویر تکرار می‌شود.

جدول ۲-۵: شناسایی و اعتبارسنجی نقاط جداسازی (مرحله دوم)

شماره مرحله	کلاس شناسایی شده مطابق جدول ۲-۴	دقت شناسایی	تصویر حاصل
اول	کلاس ۷۱ (دندانه)	٪۹۷	
دوم	کلاس ۸۸ (حرف ناقص)	٪۹۰	
سوم	کلاس ۷۹ (سس)	٪۸۰	
چهارم	کلاس ۸۸ (حرف ناقص)	٪۹۸	
نتیجه	کلاس ۷۱ (دندانه)	٪۹۷	

در مرحله‌ی دوم، جداسازی با دقت ۹۷٪ درست بوده و مطابق جدول ۲-۴ تصویر حاصل کلاس ۷۱ شناخته می‌شود. در مرحله دوم با اضافه نمودن قطعه‌ی بعدی به این تصویر، تصویر جدید با دقت ۹۰٪ نامعتبر و جزو حروف ناقص شناخته می‌شود. در مرحله‌ی سوم تصویر حاصله با اعتبار ۸۰٪ حرف (سس) و کلاس ۷۹ تشخیص داده می‌شود (به دلیل تشابه زیاد با حرف س) و همچنین در مرحله

چهارم با اضافه شدن قطعه بعدی، با دقت ۹۸٪ حرف ناقص شناخته می‌شود. بدین ترتیب در مجموع با اعتبار ۹۷٪ کلاس ۷۱ شناخته و ذخیره شده و از تصویر اصلی جدا می‌شود و بقیه‌ی تصویر به مرحله‌ی بعد فرستاده می‌شود.

جدول ۳-۵: شناسایی و اعتبارسنجی نقاط جداسازی (مرحله سوم)

شماره مرحله	کلاس شناسایی شده مطابق جدول ۲-۴	دقت شناسایی	تصویر حاصل
اول	کلاس ۷۱ (دندانه)	۹۰٪	
دوم	کلاس ۸۸ (حرف ناقص)	۹۷٪	
سوم	کلاس ۷۹ (حرف -س-)	۱۰۰٪	
چهارم	کلاس ۸۸ (حرف ناقص)	۹۸٪	
نتیجه	کلاس ۷۹ (حرف -س-)	۱۰۰٪	

در مرحله‌ی سوم، جداسازی اولیه با دقت ۹۰٪ بوده و مطابق جدول ۲-۴ تصویر حاصل کلاس ۷۱ و یا به عبارتی به صورت یک دندانه شناخته می‌شود. در مرحله‌ی دوم با اضافه نمودن قطعه‌ی بعدی به این تصویر، تصویر جدید با دقت ۹۷٪ نامعتبر و جزو حروف ناقص شناخته می‌شود. در مرحله‌ی سوم تصویر حاصله با اعتبار ۱۰۰٪ حرف (-س-) و کلاس ۷۹ تشخیص داده می‌شود و همچنین در مرحله‌ی چهارم با اضافه شدن قطعه‌ی بعدی، با دقت ۹۸٪ حرف ناقص شناخته می‌شود. بدین ترتیب در مجموع با اعتبار ۱۰۰٪ کلاس ۷۹ و یا حرف (-س-) شناخته، ذخیره و از تصویر اصلی جدا شده و

بقیه‌ی تصویر به مرحله بعد فرستاده می‌شود. از بقیه‌ی تصویر تنها ۱ قطعه باقیمانده که با دقت ۹۵٪ کلاس ۷۰ شناخته می‌شود.

ر	دقت شناسایی	کلاس شناسایی شده	مرحله
	۹۵٪	کلاس ۷۰	اول

نتیجه نهایی جداسازی و طبقه‌بندی تصویر در جدول ۴-۵ آمده است.

جدول ۴-۵: نتیجه کلی جداسازی کلمه بستنی

شماره تصویر	کلاس شناسایی شده مطابق جدول ۲-۴	دقت شناسایی	تصویر حاصل
اول	کلاس ۸۵	۹۹٪	ی
دوم	کلاس ۷۱	۹۷٪	س
سوم	کلاس ۷۹ (س)	۱۰۰٪	س
چهارم	کلاس ۷۰	۹۵٪	ر

بستی


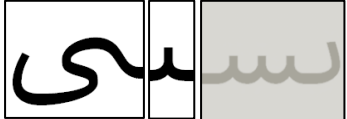

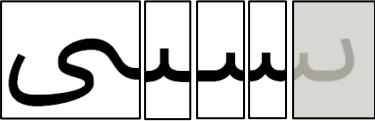

شماره کلاس زیر بدنه: ۸۵ ۷۱ ۷۹ ۷۰

بستی

شماره کلاس حروف با نقطه: ؟ ؟ ؟ ؟

برای بهبود عملکرد می‌توان عمل جداسازی را در هر مرحله از طرف دیگر ادامه داد. بدین صورت که ابتدا از سمت چپ شروع کرده و با شناسایی و جداسازی تصویر اول، شناسایی در مرحله‌ی بعد باید از سمت راست ادامه یابد. این کار تا آخرین مرحله‌ی جداسازی حروف انجام می‌شود. مزیت این روش آن است که در صورت بروز اشکال در یکی از مراحل جداسازی، عمل جداسازی متوقف نمی‌شود و از سمت دیگر ادامه می‌یابد.

جدول ۵-۵: مرحله اول جداسازی از دو سمت کلمه بستنی

شماره مرحله	کلاس شناسایی شده مطابق جدول ۲-۴	دقت شناسایی	تصویر حاصل
اول	کلاس ۸۵ (زیر بدنه "یی")	%۹۹	
دوم	کلاس ۸۸ (حرف ناقص)	%۹۰	
سوم	کلاس ۸۸ (حرف ناقص)	%۹۵	
چهارم	کلاس ۸۸ (حرف ناقص)	%۹۸	
نتیجه	کلاس ۸۵ (زیر بدنه "یی")	%۹۹	

جدول ۵-۶: مرحله دوم جداسازی از دو سمت کلمه بستنی

شماره مرحله	کلاس شناسایی شده مطابق جدول ۲-۴	دقت شناسایی	تصویر حاصل
اول	کلاس ۷۰	%۹۵	
دوم	کلاس ۸۸ (حرف ناقص)	%۹۳	
سوم	کلاس ۸۸ (حرف ناقص)	%۹۵	
چهارم	کلاس ۸۸ (حرف ناقص)	%۱۰۰	
نتیجه	کلاس ۸۵ (زیر بدنه "یی")	%۹۹	

جدول ۵-۷: مرحله سوم جداسازی از دو سمت کلمه بستنی

شماره مرحله	کلاس شناسایی شده مطابق جدول ۲-۴	دقت شناسایی	تصویر حاصل
اول	کلاس ۷۱ (دندانه)	%۹۷	
دوم	کلاس ۸۸ (حرف ناقص)	%۹۰	
سوم	کلاس ۷۹ (حرف س-)	%۸۰	
چهارم	کلاس ۸۸ (حرف ناقص)	%۹۸	
نتیجه	کلاس ۷۱ (دندانه)	%۹۷	

جدول ۵-۸: مرحله چهارم جداسازی از دو سمت کلمه بستنی

شماره مرحله	کلاس شناسایی شده مطابق جدول ۴-۲	دقت شناسایی	تصویر حاصل
اول	کلاس ۷۱ (دندانه)	٪۹۷	
دوم	کلاس ۸۸ (حرف ناقص)	٪۱۰۰	
سوم	کلاس ۷۹ (حرف س)	٪۹۸	
نتیجه	کلاس ۸۵ (زیر بدنه "بی")	٪۹۹	

پس از جداسازی حروف پیوسته به روش‌های فوق و شناسایی کلاس بدنه‌ی حروف، باید به شناسایی حروف اصلی که در آن نقاط دیگر حذف نشده‌اند، پرداخت. به عنوان مثال کلاس بدنه‌ی ۷۹ هم می‌تواند مربوط به حرف "س" و هم مربوط به حرف "ش" باشد. که در مرحله‌ی طبقه‌بندی نهایی حروف این تمایز انجام می‌شود.

۵-۲-۵- طبقه‌بندی نهایی

در مراحل گذشته ابتدا با انجام عملیات پیش‌پردازش، به جداسازی زیرکلمه و بر طرف نمودن مشکل همپوشانی پرداخته شد و حروف گسسته در ۳۹ کلاس طبقه‌بندی شدند که در آن زیرکلمه‌ی متناظر با کلاس ۴۰ به بخش جداسازی حروف پیوسته فرستاده شد.

جدول ۵-۹: کلاس‌های بدنه حروف گسسته

شماره کلاس	حرف	شماره کلاس	حرف	شماره کلاس	حرف
۱	ا	۱۶	ر ز ژ	۳۱	لله
۲	ب	۱۷	س ش	۳۲	لا
۳	پ	۱۸	ص ض	۳۳	ء
۴	ت	۱۹	ط ظ	۳۴	.
۵	ث	۲۰	ع غ	۳۵	صی ضی
۶	ج	۲۱	ف	۳۶	شی سی
۷	چ	۲۲	ق	۳۷	؟
۸	ح	۲۳	گ	۳۸	!
۹	خ	۲۴	ک	۳۹	،
۱۰	د	۲۵	ل	۴۰	زیرکلمه
۱۱	ذ	۲۶	م		
۱۲	ر	۲۷	ن		
۱۳	ز	۲۸	و		
۱۴	س	۲۹	ه		
۱۵	ش	۳۰	ی		

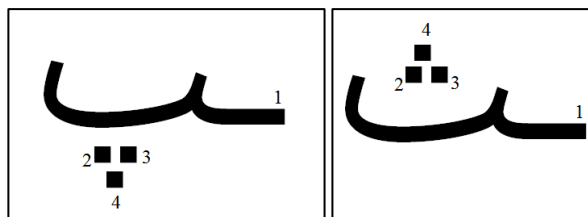
زیرکلمه‌های مرتبط با کلاس ۴۰ نیز با استفاده از روش پروفایل بالایی تعمیم یافته به حروف پیوسته تجزیه و از کلاس ۴۱ تا ۸۷ طبقه‌بندی می‌شوند. اما این زیرحروف نماینده باید در حالتی که در آن نقطه‌ها حذف نشده‌اند نیز شناسایی شوند. به عنوان مثال زیر کلاس ۴۳ می‌تواند مرتبط به حرف ج، چ و یا خ باشد که باید به روش مناسب با استفاده از برجسب‌زنی مؤلفه‌ها این طبقه‌بندی را انجام داد.

جدول ۵-۱۰ : کلاس‌های بدنه حروف پیوسته

شماره کلاس	زیرحرف نماینده	شماره کلاس	زیرحرف نماینده	شماره کلاس	زیرحرف نماینده
۴۱	ا	۵۷	ه	۷۳	ص
۴۲	ب	۵۸	هـ	۷۴	صد
۴۳	بج	۵۹	هـج	۷۵	صدج
۴۴	بم	۶۰	کم	۷۶	صم
۴۵	بم	۶۱	ک	۷۷	سم
۴۶	بک	۶۲	کک	۷۸	سک
۴۷	بل	۶۳	کل	۷۹	سک
۴۸	بف	۶۴	ل	۸۰	ط
۴۹	ف	۶۵	ل	۸۱	ط
۵۰	م	۶۶	م	۸۲	ط
۵۱	گ	۶۷	م	۸۳	و
۵۲	گ	۶۸	م	۸۴	و
۵۳	گ	۶۹	س	۸۵	و
۵۴	ف	۷۰	ر	۸۶	س
۵۵	ح	۷۱	ر	۸۷	لا
۵۶	ح	۷۲	ر	۸۸	حروف ناقص

۵-۲-۱ طبقه‌بندی نهایی با استفاده از برچسب‌زنی مؤلفه‌ها

در این قسمت با استفاده از روش برچسب‌زنی مؤلفه‌ها به تشخیص حروف اصلی از بدنه‌ی حروف پرداخته می‌شود. روند کار بدین صورت است که با استفاده از برچسب‌زنی مؤلفه‌ها، تعداد نقاط و محل آن‌ها مشخص می‌شود و متناسب با این معیار حرف اصلی مشخص شده و در کلاس مناسب طبقه‌بندی می‌شود. به عنوان مثال برای شناسایی زیرحرفی که در کلاس ۴۲ قرار دارد ابتدا مؤلفه‌های موجود در تصویر برچسب‌زنی می‌شوند.



شکل ۵-۱۹: برچسب‌زنی مؤلفه و مشخص نمودن تعداد و محل دقیق نقاط

اگر تعداد مؤلفه‌ها برابر ۲ باشد یعنی ۱ مؤلفه برای زیرحرف نماینده بوده و دیگری نقطه‌ی مربوط به حرف است و مشخص می‌شود که حرف مورد نظر **پ** است. اگر تعداد مؤلفه‌ها برابر ۳ باشد نشانه‌ی آن است که حرف مورد نظر **ت** است. اما زمانی که تعداد مؤلفه‌ها برابر ۴ است دو حالت ممکن است رخ داده باشد (**پ**، **ت**). حال اگر میانگین عرضی مکان نقطه‌ها پایین‌تر از خط زمینه باشد حرف مورد نظر **پ** بوده و اگر بالاتر از خط زمینه باشد حرف **ت** خواهد بود.

پ ت پ ت

۴۲۰۳
۴۲۰۲
۴۲۰۱
۴۲۰۰

شکل ۵-۲۰: کلاس نهایی حروف

بنابراین زیرحرف کلاس ۴۲ با توجه به تعداد و مکان نقاط، به صورت کلاس ۴۲۰۰، ۴۲۰۱، ۴۲۰۲ و ۴۲۰۳ خواهد بود. برای حروف دیگر نیز کلاسه‌بندی به همین صورت انجام شده که نتیجه‌ی آن در جدول ۵-۷ آمده است.

جدول ۵-۱۱: کلاس‌های نهایی حروف گسسته

شماره کلاس نهایی	زیرحرف	شماره کلاس نماینده حرف	شماره کلاس نهایی	زیرحرف	شماره کلاس نماینده حرف
۱۵۰۰	د	۱۵	۱۰۰	۱	۱
۱۵۰۱	ذ		۲۰۰	۲	۲
۱۶۰۰	ر	۱۶	۳۰۰	۳	۳
۱۶۰۱	ز		۴۰۰	۴	۴
۱۶۰۲	ژ		۵۰۰	۵	۵
۱۷۰۰	س	۱۷	۶۰۰	۶	۶
۱۷۰۱	ش		۷۰۰	۷	۷
۱۸۰۰	ص	۱۸	۸۰۰	۸	۸
۱۸۰۱	ض		۹۰۰	۹	۹
۱۹۰۰	ط	۱۹	۱۰۰۰	۰	۱۰
۱۹۰۱	ظ		۱۱۰۰	آ	۱۱
۲۰۰۰	ع	۲۰	۱۲۰۰	ا	۱۲
۲۰۰۱	غ		۱۳۰۰	ب	۱۳

۲۱۰۰	ف	۲۱	۱۳۰۱	ت	۱۴
۲۲۰۰	ق	۲۲	۱۳۰۲	پ	
۲۳۰۰	ک	۲۳	۱۳۰۳	ث	
۲۴۰۰	گ	۲۴	۱۴۰۰	ح	
۲۵۰۰	ل	۲۵	۱۴۰۱	خ	
۲۶۰۰	م	۲۶	۱۴۰۲	ج	
۲۷۰۰	ن	۲۷	۱۴۰۳	چ	

شماره کلاس نهایی	زیرحرف	شماره کلاس نماینده حرف	شماره کلاس نهایی	زیرحرف	شماره کلاس نماینده حرف
۳۵۰۰	ص	۳۵	۲۸۰۰	و	۲۸
۳۵۰۱	ضی		۲۹۰۰	ه	۲۹
۳۶۰۰	س	۳۶	۳۰۰۰	ی	۳۰
۳۶۰۱	شی		۳۱۰۰	له	۳۱
۳۷۰۰	؟	۳۷	۳۲۰۰	لا	۳۲
۳۸۰۰	!	۳۸	۳۳۰۰	ء	۳۳
۳۹۰۰	،	۳۹	۳۴۰۰	.	۳۴

جدول ۵-۱۲: کلاس‌های نهایی حروف پیوسته

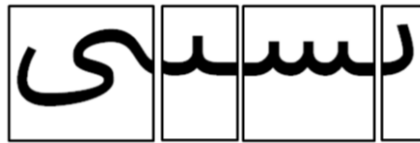
شماره کلاس نهایی	زیرحرف	شماره کلاس نماینده حرف	شماره کلاس نهایی	زیرحرف	شماره کلاس نماینده حرف
۴۸۰۰	ف	۴۸	۴۱۰۰	ا	۴۱
۴۹۰۰	فا	۴۹	۴۲۰۰	با	۴۲
۴۹۰۱	قا		۴۲۰۱	تا	
۵۰۰۰	فا	۵۰	۴۲۰۲	پا	
۵۰۰۱	قا		۴۲۰۳	تا	
۵۱۰۰	گا	۵۱	۴۳۰۰	جا	۴۳
۵۲۰۰	گا	۵۲	۴۳۰۱	خا	
۵۳۰۰	گا	۵۳	۴۳۰۲	چا	
۵۴۰۰	قا	۵۴	۴۳۰۳	چا	
۵۵۰۰	حا	۵۵	۴۴۰۰	عا	۴۴
۵۵۰۱	خا		۴۴۰۱	غا	
۵۵۰۲	جا		۴۵۰۰	ما	۴۵
۵۵۰۳	چا		۴۵۰۱	غا	

۵۶۰۰	ح	۵۶	۴۶۰۰	ع	۴۶
۵۶۰۱	خ		۴۶۰۱	غ	
۵۶۰۲	ج		۴۷۰۰	د	۴۷
۵۶۰۳	چ		۴۷۰۱	ذ	

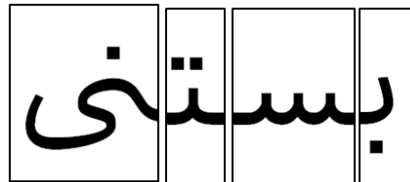
شماره کلاس نهایی	شماره کلاس	زیرحرف	شماره کلاس نهایی	زیرحرف	شماره کلاس اولیه
۷۱۰۰	ب	۷۱	۵۷۰۰	ه	۵۷
۷۱۰۱	ن		۵۸۰۰	هـ	۵۸
۷۱۰۲	پ		۵۹۰۲	هـ	۵۹
۷۱۰۳	ت		۶۰۰۰	ک	۶۰
۷۱۰۴	پ		۶۱۰۰	ک	۶۱
۷۱۰۵	ث		۶۲۰۰	ک	۶۲
۷۲۰۰	ر	۷۲	۶۳۰۰	ل	۶۳
۷۲۰۱	ز		۶۴۰۰	ل	۶۴
۷۲۰۲	ژ		۶۵۰۰	ل	۶۵
۷۳۰۰	ص	۷۳	۶۶۰۰	م	۶۶
۷۳۰۱	ض		۶۷۰۰	م	۶۷
۷۴۰۰	ص	۷۴	۶۸۰۰	م	۶۸

۷۴۰۱	ضد		۶۹۰۰	نن	۶۹
۷۵۰۰	صد	۷۵	۷۰۰۰	ب	۷۰
۷۵۰۱	ضد		۷۰۰۱	نا	
۷۶۰۰	صی	۷۶	۷۰۰۲	یا	
۷۶۰۱	ضی		۷۰۰۳	تا	
۷۷۰۰	سی	۷۷	۷۰۰۴	پا	
۷۷۰۱	شی		۷۰۰۵	ثا	
شماره کلاس نهایی	شماره کلاس	زیرحرف	شماره کلاس نهایی	زیرحرف	شماره کلاس اولیه
۸۴۰۰	سی	۸۴	۷۸۰۰	سا	۷۸
۸۵۰۰	جی	۸۵	۷۸۰۱	شا	
۸۵۰۱	نی		۷۹۰۰	سا	۷۹
۸۵۰۲	چی		۷۹۰۱	شا	
۸۵۰۳	تی		۸۰۰۰	طا	۸۰
۸۵۰۴	ثی		۸۰۰۱	ظا	
۸۵۰۵	پی	۸۱۰۰	طا	۸۱	
۸۶۰۰	س	۸۶	۸۱۰۱		ظا
۸۶۰۱	ش		۸۲۰۰	طا	۸۲
۸۷۰۰	لا	۸۷	۸۲۰۱	ظا	
۸۸۰۰	ناقص	۸۸	۸۳۰۰	و	۸۳

مطابق با جدول بالا شناسایی نهایی کلمه بستنی مطابق شکل ۵-۲۱ خواهد بود.



شماره کلاس زیر بدنه: ۷۰ ۷۹ ۷۱ ۸۵



شماره کلاس حروف با نقطه: ۷۰۰۰ ۷۹۰۰ ۷۱۰۳ ۸۵۰۱

شکل ۵-۲۱: شماره کلاس نهایی حروف با وجود نقاط

مراحل کامل کار برای جملهی "مرام او مرا رام کرد." به صورت زیر خواهد بود.

مرام او مرا رام کرد.

شکل ۵-۲۲: انجام عملیات پیش پردازش

م | ر | ا | م | ا | و | م | ر | ا | ر | ا | م | ک | ر | د | .

شکل ۵-۲۳: جداسازی اولیه زیرکلمات با وجود مشکل همپوشانی

م | ر | ا | م | ا | و | م | ر | ا | ر | ا | م | ک | ر | د | .

شکل ۵-۲۴: برطرف نمودن مشکل همپوشانی

م | ر | ا | م | ا | و | م | ر | ا | ر | ا | م | ک | ر | د | .

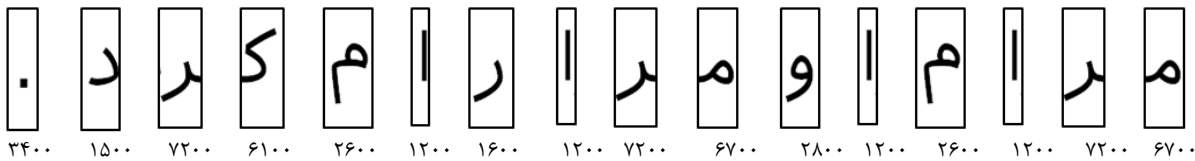
شماره کلاس: ۴۰ ۱۲ ۲۶ ۱۲ ۲۸ ۴۰ ۱۲ ۱۶ ۱۲ ۲۶ ۴۰ ۱۵ ۳۴

شکل ۵-۲۵: شناسایی اولیه و ارسال زیرکلمات کلاس ۴۰ برای جداسازی

م | ر | ا | م | ا | و | م | ر | ا | ر | ا | م | ک | ر | د | .

۳۴ ۱۵ ۷۲ ۶۱ ۲۶ ۱۲ ۱۶ ۱۲ ۷۲ ۶۷ ۲۸ ۱۲ ۲۶ ۱۲ ۷۲ ۶۷

شکل ۵-۲۶: جداسازی زیرکلمات و کلاسه بندی بدنه حروف



شکل ۵-۲۷: کلاسه بندی نهایی حروف

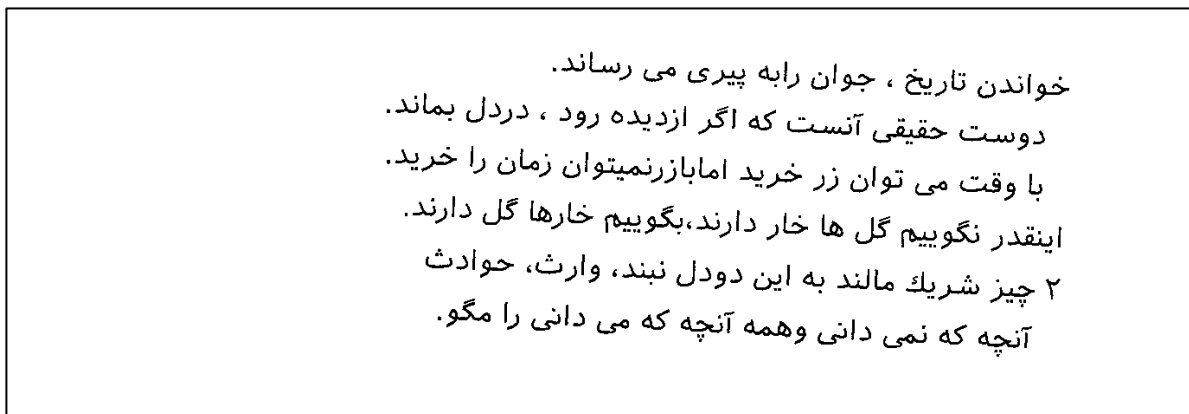
مرام او مرا رام کرد.

شکل ۵-۲۸: سند تولید شده نهایی با استفاده از کلاسه بندی نهایی

بنابراین با توجه به شماره کلاس نهایی تولید شده برای هر حرف جمله مذکور، مطابق جداول طبقه بندی نهایی، سند دیجیتالی به صورت متن تکست^۱ ذخیره خواهد شد.

۵-۳- ارزیابی الگوریتم برای یک تصویر کامل

این ارزیابی بر روی مجموعه ای از تصاویر اسکن شده انجام شده است که در آن تصاویر به گونه ای در نظر گرفته می شود که هم دارای چرخش اولیه بوده و هم دارای اندکی نویز است. نمونه ای از این تصاویر مطابق شکل ۵-۲۹ است.



شکل ۵-۲۹: یک تصویر نمونه برای ارزیابی الگوریتم

¹ Text

خواندن تاریخ ، جوان رابه پیری می رساند.

دوست حقیقی آنست که اگر ازدیده رود ، دردل بماند.

با وقت می توان زر خرید اما باز نمیتوان زمان را خرید.

اینقدر نگوییم گل ها خار دارند، بگوییم خارها گل دارند.

۲ چیز شريك مالند به این دودل نبند، وارث، حوادث

آنچه که نمی دانی وهمه آنچه که می دانی را مگو.

شکل ۵-۳۰: شناسایی جملات و جداسازی آنها

در این قسمت ابتدا با تکنیک‌های ذکر شده در فصول قبل، چرخش‌های لازم انجام شده و نویزهای موجود حذف می‌شود. سپس خطوط جملات موجود در متن تشخیص داده شده و تفکیک می‌شود. پس از جداسازی، جملات به بخش استخراج زیرکلمات فرستاده شده تا برای هر جمله، زیرکلمات آن استخراج شود. پس از شناسایی، زیرکلمات به بخش جداسازی زیرحروف فرستاده شده و با استفاده از روش پروفایل بالایی کاندیدهای اولیه‌ی جداسازی انتخاب شده و با استفاده از طبقه‌بند حروف پیوسته نقاط جداسازی بدنه اصلی زیرکلمه مشخص می‌شود.

در مرحله‌ی اول، جداسازی زیرکلمه فقط از یک سمت انجام می‌گرفت. مشکل این روش آن است که اگر روند جداسازی از یک سمت انجام شود در صورت عدم شناسایی حروف اولیه یک زیرکلمه، شناسایی حروف بعدی آنها نیز انجام نمی‌شود. نتیجه‌ی جداسازی برای یک جمله در شکل ۵-۳۱ آمده است.

خواندن تاریخ، جوان رابه پیری رساند.
دوست حقیقی آنست که اگر ازدیده رود، دردل بماند.
با وقت می توان زر خرید اما باز نمیتوان زمان را خرید.
اینقدر نگوییم گل ها خار دارند، بگوییم خارها گل دارند.
۲ چیز شريك مالند به این دودل نبند، وارث، حوادث

شکل ۵-۳۱: نتیجه شناسایی برای حالتی که جداسازی از یک سمت انجام می‌شود

در این روند جداسازی و شناسایی با خطای بالایی همراه است. دو مشکل اساسی در این روند وجود دارد. اول آنکه جداسازی و شناسایی فقط از یک سمت انجام شده و در صورت بروز مشکلی در

شناسایی حروف ابتدایی، عمل جداسازی زیرکلمه متوقف می‌شود. بنابراین برای غلبه بر این مشکل باید جداسازی از هر دو سمت انجام شود. نتیجه‌ی بهبود این جداسازی در شکل ۵-۳۲ آمده است.

خواندن تاریخ، جوان را به پیزی رساند.
 دست حقیقی آنست که اگر از دیدهِ ر_د، در دل بما ند.
 با وقت م_توان زر خزیذا ما با زر نمیشوان زمان را خزیذ.
 انتقدر نگوییم گل‌ها خار دارند، بگو پیم خارها گل دارند.
 ۲ چیز یز يك مالـد به این دو دل نیتد، وارث، حوادث

شکل ۵-۳۲: نتیجه شناسایی برای حالتی که جداسازی از هر دو سمت انجام می‌شود

مشکل دیگری نیز وجود دارد که در شناسایی نهایی موجب غلط‌های املائی می‌شود. این مسئله از آنجا ناشی می‌شود که ممکن است در زمان جداسازی زیرکلمات، نقطه‌های تصویر نیز تقسیم شده و در دو تصویر مجاور محاسبه می‌شوند. برای غلبه بر این مشکل با استفاده از برچسب‌زنی مؤلفه‌ها و در نظر گرفتن آستانه‌ی مناسب، نقطه‌های اضافی تولید شده حذف می‌شود. روند تصحیح و حذف نقاط در جدول ۵-۹ آمده است. با اعمال روش مذکور به جملات فوق، غلط‌های املائی بوجود آمده تا حد قابل قبولی برطرف می‌شود. نتیجه‌ی نهایی در شکل ۵-۳۳ آمده است.

جدول ۵-۱۳: حذف نقاط اضافی در زمان جداسازی

شناسایی درست	حذف نقاط اضافی	شناسایی نادرست	ایجاد نقطه اضافی در زمان جداسازی
پیری	پیری	پیری	پیری
دیده	دیده	دیده	دیده
نمیتوان	نمیتوان	نمیتوان	نمیتوان
خرید	خرید	خریذ	خرید

خواندن تاریخ ، جوان را به پیری رساند.
دست حقیقی آنست که اگر از دیدِ رَد، در دل بماند.
با وقت می توان زر خرید اما با زر نمیتوان زمان را خرید .
اینقدر نگو بییم گل ها خار دارند ، بگو بییم خارها گل دارند .
۲ چیز شریک مالند به این دو دل نبند ، وارث ، حوادث

شکل ۳۳-۵ : نتیجه نهایی شناسایی پس از حذف نقاط اضافی تولید شده

۵-۴- ارزیابی الگوریتم برای تصویرهای دیگر

در ادامه الگوریتم پیشنهادی برای چند تصویر دیگر نیز مورد آزمایش قرار گرفت و نتایج آن ثبت

گردید. این نتایج روند بهبود شناسایی تصویر را در سه مرحله نشان می دهد.

عجله و مزاح ، وقار را میبرد ، سکوت و آرامش هیبت می آورد.
دروغ درسازش بین مردمان بهتر از راستی درفتنه انگیزی میان آنهاست .
مردان حق گفتارشان ذکر ، سکوتشان فکر و نگاهشان عبرت است .
آزموده را آزمودن خطاست.
برسرسفره غذا بامهمانت صبورباش وپس از اوبرخیز .
دوست ، آینه دوست است با انعکاس شفاف خوبیها و بدیها.
بدرستی که جوانا پیر و پیران جوان دردنیا بسیارند.
وقتی آنچه داریم می بخشیم ، آنچه نیازمندیم دریافت خواهیم کرد

شکل ۳۴-۵ : تصویر دوم، دارای چرخش اولیه و نویز

عجله و مزاح ، وقار را میبرد ، سکوت و آرامش هیبت می آورد.

دروغ درسازش بین مردمان بهتر از راستی درفتنه انگیزی میان آنهاست .

مردان حق گفتارشان ذکر ، سکوتشان فکر و نگاهشان عبرت است .

آزموده را آزمودن خطاست.

برسرسفره غذا بامهمانت صبورباش وپس از اوبرخیز .

دوست ، آینه دوست است با انعکاس شفاف خوبیها و بدیها.

بدرستی که جوانا پیر و پیران جوان دردنیا بسیارند.

وقتی آنچه داریم می بخشیم ، آنچه نیازمندیم دریافت خواهیم کرد

شکل ۳۵-۵ : حذف نویز و چرخش و استخراج جملات تصویر دوم

عجله و مزاج، قار را میبرد، سکوت و آرامش هیبت می آورد. دروغ در سازش بین مردمان بهتراستی در فتنه انگیزی میان آنهاست. مردان حق گفتارشان ذکر، سکوتشان فکر و نگاهشان عبرت است. آزموده را آزمون خطاست. بر سر سفره غذا با می نث صبور باش پس از او بر خیز. دوست، آینه دوست است با انعکاس شفاف خود بدید. بدستی که جوانا پیر و پیران جوان در دنیا بسپارند. وقتی آنچه داریم می بخشیم، آنچه نیازمندیم دریافت خواهیم کرد.

شکل ۵-۳۶: نتیجه شناسایی تصویر دوم، جداسازی از یک سمت

عجله و مزاج، قار را میبرد، سکوت و آرامش هیبت می آورد. دروغ در سازش بین مردمان بهتراستی در فتنه انگیزی میان آنهاست. مردان حق گفتارشان ذکر، سکوتشان فکر و نگاهشان عبرت است. آزموده را آزمون خطاست. بر سر سفره غذا با مهمان صبور باش پس از او بر خیز. دوست، آینه دوست است با انعکاس شفاف خود بیها بدیها. بدستی که جوانا پیر و پیران جوان در دنیا بسپارند. وقتی آنچه داریم می بخشیم، آنچه نیازمندیم دریافت خواهیم کرد.

شکل ۵-۳۷: نتیجه شناسایی تصویر دوم، جداسازی از هر دو سمت

عجله و مزاج، قار را میبرد، سکوت و آرامش هیبت می آورد. دروغ در سازش بین مردمان بهتراستی در فتنه انگیزی میان آنهاست. مردان حق گفتارشان ذکر، سکوتشان فکر و نگاهشان عبرت است. آزموده را آزمون خطاست. بر سر سفره غذا با مهمان صبور باش پس از او بر خیز. دوست، آینه دوست است با انعکاس شفاف خود بیها بدیها. بدستی که جوانا پیر و پیران جوان در دنیا بسپارند. وقتی آنچه داریم می بخشیم، آنچه نیازمندیم دریافت خواهیم کرد.

شکل ۵-۳۸: نتیجه شناسایی تصویر دوم، جداسازی از هر دو سمت با تصحیح نقاط اضافی

انسان مانند زمین است، برای رسیدن به خوشبختی باید سختی بکشد. صاحب همت در پیچ و خم های زندگی هیچ گاه با یاس روبه رو نخواهد شد. بهترین افراد کسانی هستند که از خوشبختی دیگران خوشحالند. کسانی که بیش از اندازه فکرمی کنند فاقد اراده و تصمیم می باشند. آنان که نمی توانند خود را اداره کنند ناچار به پیروی از دیگرانند. قرض کردن از کسی شایسته نیست مخصوصا از تازه به دوران رسیده ها. بدبین، همچون نابینایی است که از خورشید فقط گرمایش را حس می کند. درموقع جنگ نمی توان زره پوشید که از وقتش گذشته است. دینی که به درد دنیایت نخورد، به درد آخرت هم نخواهد خورد. تغافل از دانسته هایش از آنکه از ما پرسیده شود شایسته است! دنیا و آخرت همچون باختر و خاورند دوری از یکی، نزدیکی به دیگری است. بریدن از نادانی که به اصلاحش امید نیست همانا پیوند با خردمند است.

شکل ۵-۳۹: نتیجه شناسایی تصویر سوم، جداسازی از هر دو سمت

انسان مانند زمین است، برای رسیدن به خوشبختی باید سختی بکشد.
صاحب همت در پیچ و خم های زندگی هیچ گاه با یاس روبه رو نخواهد شد.
بهترین افراد کسانی هستند که از خوشبختی دیگران خوشحالند.
کسانی که بیش از اندازه فکرمی کنند فاقد اراده و تصمیم می باشند.
آنان که نمی توانند خود را اداره کنند ناچار به پیروی از دیگرانند.
قرض کردن از کسی شایسته نیست مخصوصا از تازه به دوران رسیده ها.
بدبین، همچون نابینایی است که از خورشید فقط گرمایش را حس می کند.
درموقع جنگ نمی توان زره پوشید که از وقتش گذشته است.
دینی که به درد دنیایت نخورد، به درد آخرت هم نخواهد خورد.
تغافل از دانسته هایش از آنکه از ما پرسیده شود شایسته است!
دنیا و آخرت همچون باختر و خاورند دوری از یکی، نزدیکی به دیگری است.
بریدن از نادانی که به اصلاحش امید نیست همانا پیوند با خردمند است.

شکل ۵-۴۰: حذف نویز و چرخش و استخراج جملات تصویر سوم

انسان مانند زمی است، برای رسیدن به خوشی باید سختی بکشد. صاحب همت در پیچ و خم های زندگی هیچ گاه با یاس و بهر و نخو اهد شد. بهترین افراد کسانی هستند که از خوشی دیگری خوشحالند. کسانی که بیش از انداز فکر میکنند فاقد اراده و تصمیم باشند. آنان که نمی توانند خود را در آوار کنده ناچار به پیروی از دیگرانند. قرض کردن از کسی شایسته نیست مخصوصاً از تازه به دوران رسیده ها. بدبین، همچون نابینا است که از خورشید فقط گرمایی را حس میکند. در موقع جنگ نمیتوان زره پوشید که از وقت گذشته است. دینی که یه درددنیات نخورد، به دردت آخرت همدردی نخواهد خورد. تغافل از دانسته های پیش از آنکه از ما پرسید شود شایسته است! دنیا آخرت همچون باختر و خا رنددوری از یک، نزدیکی به دیگری است. بریدن از نادانی که به اصلاحش امید نیست همانا پیوند با خردمند است.

شکل ۴۱-۵: نتیجه شناسایی تصویر سوم، جداسازی از یک سمت

انسان مانند زمین است، برای رسیدن به خوشی باید سختی بکشد. صاحب همت در پیچ و خم های زندگی هیچ گاه با یاس و بهر و نخو اهد شد. بهترین افراد کسانی هستند که از خوشی دیگری خوشحالند. کسانی که بیش از انداز فکر میکنند فاقد اراده و تصمیم باشند. آنان که نمی توانند خود را در آوار کنده ناچار به پیروی از دیگرانند. قرض کردن از کسی شایسته نیست مخصوصاً از تازه به دوران رسیده ها. بدبین، همچون نابینا است که از خورشید فقط گرمایی را حس میکند. در موقع جنگ نمیتوان زره پوشید که از وقت گذشته است. دینی که یه درددنیات نخورد، به دردت آخرت همدردی نخواهد خورد. تغافل از دانسته های پیش از آنکه از ما پرسید شود شایسته است! دنیا آخرت همچون باختر و خا رنددوری از یک، نزدیکی به دیگری است. بریدن از نادانی که به اصلاحش امید نیست همانا پیوند با خردمند است.

شکل ۴۲-۵: نتیجه شناسایی تصویر سوم، جداسازی از هر دو سمت

انسان مانند زمین است، برای رسیدن به خوشی باید سختی بکشد. صاحب همت در پیچ و خم های زندگی هیچ گاه با یاس و بهر و نخو اهد شد. بهترین افراد کسانی هستند که از خوشی دیگری خوشحالند. کسانی که بیش از انداز فکر میکنند فاقد اراده و تصمیم باشند. آنان که نمی توانند خود را در آوار کنده ناچار به پیروی از دیگرانند. قرض کردن از کسی شایسته نیست مخصوصاً از تازه به دوران رسیده ها. بدبین، همچون نابینا است که از خورشید فقط گرمایی را حس میکند. در موقع جنگ نمیتوان زره پوشید که از وقت گذشته است. دینی که یه درددنیات نخورد، به دردت آخرت همدردی نخواهد خورد. تغافل از دانسته های پیش از آنکه از ما پرسید شود شایسته است! دنیا آخرت همچون باختر و خا رنددوری از یک، نزدیکی به دیگری است. بریدن از نادانی که به اصلاحش امید نیست همانا پیوند با خردمند است.

شکل ۴۳-۵: نتیجه شناسایی تصویر سوم، جداسازی از هر دو سمت با تصحیح نقاط اضافی

قدرزمان حال را بدانید که گذشته هرگز برنمی‌گردد و آینده شاید نیاید.
دانشگاه تمام استعداد‌های افراد از جمله بی‌استعدادی آن‌ها را آشکار می‌کند.
اگر صخره در مسیر رود نبود، رود هیچ آوازی از خود سر نمی‌داد.
حکمت و خرد را از همه کس حتی بدان فراگیر.
هیچ می‌دانی فرصتی که از آن بهره نمی‌گیری، آرزوی دیگران است.
هر انسانی مرتکب اشتباه می‌شود.
اما فقط انسان‌های احمق اشتباه خود را تکرار می‌کنند.
خداوند از ما نمی‌خواهد کارهای بزرگ را به ثمر برسانیم.
ایمنی از دیگران، آزادی است و ترس از آنها بردگی.
او فقط از ما می‌خواهد کارهای کوچک را با عشقی شگرف انجام دهیم.
اگر تو نمی‌توانی باور کنی، برای کسی که باور دارد همه چیز امکان‌پذیر
است.

شکل ۴۴-۵: نتیجه شناسایی تصویر چهارم، جداسازی از هر دو سمت

قدرزمان حال را بدانید که گذشته هرگز برنمی‌گردد و آینده شاید نیاید.

دانشگاه تمام استعداد‌های افراد از جمله بی‌استعدادی آن‌ها را آشکار می‌کند.

اگر صخره در مسیر رود نبود، رود هیچ آوازی از خود سر نمی‌داد.

حکمت و خرد را از همه کس حتی بدان فراگیر.

هیچ می‌دانی فرصتی که از آن بهره نمی‌گیری، آرزوی دیگران است.

هر انسانی مرتکب اشتباه می‌شود.

اما فقط انسان‌های احمق اشتباه خود را تکرار می‌کنند.

خداوند از ما نمی‌خواهد کارهای بزرگ را به ثمر برسانیم.

ایمنی از دیگران، آزادی است و ترس از آنها بردگی.

او فقط از ما می‌خواهد کارهای کوچک را با عشقی شگرف انجام دهیم.

اگر تو نمی‌توانی باور کنی، برای کسی که باور دارد همه چیز امکان‌پذیر

است.

شکل ۴۵-۵: حذف نویز و چرخش و استخراج جملات تصویر چهارم

قد ر زمان حال را بد اتید که گذ شد هر گز بر نمِ گردد _ آیتد ه شاید نیاید .
 دانش _ تمام استغداد های افراد از جمله ی استغدادی آن ها را آیش ر می کند .
 اگر صخره در مسیر رود نبود ، رود هیچ آوازی از خود سر نمی داد .
 حکمت و خرد را از همه کس حتی بدان فرا بگیر .
 هیچ می دان فرستی که از آن بهر _ نمی گئی ، آرزوی دیگران است .
 هر انسانی مر تکب اشته می شود ،
 اما فقط انسان های احمق اشته ه خود را تکرار می کنند .
 خدا _ نداز ما نمی خواهد کارهای بزرگ را به ثمر برساند .
 ایمنی از دیگران ، آزادی است _ تیزس از آنها بردگی .
 _ فقط از _ می خواهد کارهای کوچک را با عشق شگ انجام دهد .
 اگر تو نمی توانی با _ رکنی ، برای کسی که با _ ر دارد همه چیز امکان پذیر است .

شکل ۴۶-۵ : نتیجه شناسایی تصویر چهارم، جداسازی از یک سمت

قد ر زمان حال را بد اتید که گذشته هر گز بر نمِ گردد _ آیتد ه شاید نیاید .
 دانشگا _ تمام استغداد های افراد از جمله ی استغدادی آن ها را آشکار می کند .
 اگر صخره در مسیر رود نبود ، رود هیچ آوازی از خود سر نمی داد .
 حکمت و خرد را از همه کس حتی بدان فرا بگیر .
 هیچ می دان فرستی که از آن بهر _ نمی گئی ، آرزوی دیگران است .
 هر انسانی مر تکب اشتباه می شود ،
 اما فقط انسان های احمق اشتباه خود را تکرار می کنند .
 خدا _ نداز _ نمی خواهد کارهای بزرگ را به ثمر برساند .
 ایمنی از دیگران ، آزادی است _ تیزس از آنها بردگی .
 _ فقط از ما می خواهد کارهای کوچک را با عشقی شگرف انجام دهد .
 اگر تو نمی توانی با _ رکنی ، برای کسی که با _ ر دارد همه چیز امکان پذیر است .

شکل ۴۷-۵ : نتیجه شناسایی تصویر چهارم، جداسازی از هر دو سمت

قد ر زمان حال را بد انید که گذشته هر گز بر نمِ گردد _ آیتد ه شاید نیاید .
 دانشگا _ تمام استعداد های افراد از جمله ی استعدادی آن ها را آشکار می کند .
 اگر صخره در مسیر رود نبود ، رود هیچ آوازی از خود سر نمی داد .
 حکمت و خرد را از همه کس حتی بدان فرا بگیر .
 هیچ می دان فرستی که از آن بهر _ نمی گئی ، آرزوی دیگران است .
 هر انسانی مر تکب اشتباه می شود ،
 اما فقط انسان های احمق اشتباه خود را تکرار می کنند .
 خدا _ نداز ما نمی خواهد کارهای بزرگ را به ثمر برساند .
 ایمنی از دیگران ، آزادی است _ ترس از آنها بردگی .
 _ فقط از _ می خواهد کارهای کوچک را با عشقی شگرف انجام دهد .
 اگر تو نمی توانی با _ رکنی ، برای کسی که با _ ر دارد همه چیز امکان پذیر است .

شکل ۴۸-۵ : نتیجه شناسایی تصویر چهارم، جداسازی از هر دو سمت با تصحیح نقاط اضافی

دنیا را محبت ، نجات می دهد.
صبر ، ضامن پیروزی است.
افراد شجاع فرصت می آفرینند.
ترسوها و ضعفا منتظر فرصت می نشینند.
هیچ وقت به خدا نگوئید من یک مشکل بزرگ دارم.
به مشکلاتان بگوئید من یک خدای بزرگ دارم.
اگر زمین سختی زمستان را نکشد ، بهار نمی شود.
زندگی بدون عشق مثل دشت بی باران است.
قشنگ ترین اسارت زندگی است.
تنها مرگ است که دروغ نمیگوید.
هر ساختمان بزرگ ، زمانی فقط یک نقشه ساده بوده.
تافرمان نبری ، فرمانروایی نتوان کرد.

شکل ۵-۴۹: نتیجه شناسایی تصویر پنجم، جداسازی از هر دو سمت

دنیا را محبت ، نجات می دهد.

صبر ، ضامن پیروزی است.

افراد شجاع فرصت می آفرینند.

ترسوها و ضعفا منتظر فرصت می نشینند.

هیچ وقت به خدا نگوئید من یک مشکل بزرگ دارم.

به مشکلاتان بگوئید من یک خدای بزرگ دارم.

اگر زمین سختی زمستان را نکشد ، بهار نمی شود.

زندگی بدون عشق مثل دشت بی باران است.

قشنگ ترین اسارت زندگی است.

تنها مرگ است که دروغ نمیگوید.

هر ساختمان بزرگ ، زمانی فقط یک نقشه ساده بوده.

تافرمان نبری ، فرمانروایی نتوان کرد.

شکل ۵-۵۰: حذف نویز و چرخش و استخراج جملات تصویر پنجم

د نیا را محبے ، نجات مے دھد .
 صبر ، ضامن پیر و زی است .
 افراد شجاع فرصت مے آفرینند .
 ترسوها و ضعفا منتظر فرصت مے نشیے .
 هیچ وقت به خدا تگو یید من یک مشیے بزرگ دارم .
 به مشکللتان یگو یید من یک خدای بزرگ دارم .
 اگر زمین سختی زمستان را نکیے ، بهار نمی شود .
 زندگی بدین عشق مثل دشت بی باران است .
 قشنگ ترین اسارت زندگی است .
 تنها مرگ است که در غمیید .
 هر ساختمان بزرگ ، زمانی فقط یک نقشه ساده بود .
 تا فرمانبری ، فرمانری ای نتوان کرد .

شکل ۵-۵۱ : نتیجه شناسایی تصویر پنجم، جداسازی از یک سمت

د نیا را محبت ، نجات مے دھد .
 صبر ، ضامن پیر و زی است .
 افراد شجاع فرصت مے آفرینند .
 ترسوها و ضعفا منتظر فرصت مے نشینند .
 هیچ وقت به خدا تگو یید من یک مشکل بزرگ دارم .
 به مشکللتان یگو یید من یک خدای بزرگ دارم .
 اگر زمین سختی زمستان را نکشد ، بهار نمی شود .
 زندگی بدین عشق مثل دشت بی باران است .
 قشنگ ترین اسارت زندگی است .
 تنها مرگ است که در غمیگو یید .
 هر ساختمان بزرگ ، زمانی فقط یک نقشه ساده بود .
 تا فرمانبری ، فرمانری ای نتوان کرد .

شکل ۵-۵۲ : نتیجه شناسایی تصویر پنجم، جداسازی از هر دو سمت

د نیا را محبت ، نجات مے دھد .
 صبر ، ضامن پیر و زی است .
 افراد شجاع فرصت مے آفرینند .
 ترسوها و ضعفا منتظر فرصت مے نشینند .
 هیچ وقت به خدا تگو یید من یک مشکل بزرگ دارم .
 به مشکللتان یگو یید من یک خدای بزرگ دارم .
 اگر زمین سختی زمستان را نکشد ، بهار نمی شود .
 زندگی بدین عشق مثل دشت بی باران است .
 قشنگ ترین اسارت زندگی است .
 تنها مرگ است که در غمیگو یید .
 هر ساختمان بزرگ ، زمانی فقط یک نقشه ساده بود .
 تا فرمانبری ، فرمانری ای نتوان کرد .

شکل ۵-۵۳ : نتیجه شناسایی تصویر پنجم، جداسازی از هر دو سمت با تصحیح نقاط اضافی

۵-۴-۱- نتیجه ارزیابی الگوریتم

در قسمت قبل جهت اعتبارسنجی، روند بازشناسی بر روی ۵ صفحه‌ی متن فارسی با درجه تفکیک ۳۰۰ نقطه بر اینچ آزمایش شد و نتیجه بازشناسی در سه مرحله بهبود یافت.

در مرحله اول شناسایی از یک سمت انجام می‌شود. در این مرحله اگر در زیرکلمه‌ای جداسازی برای حروف اولیه زیرکلمه به درستی انجام نشود عمل جداسازی برای بقیه حروف آن زیر کلمه ادامه نمی‌یابد.

در مرحله‌ی دوم جهت بر طرف نمودن مشکل مرحله‌ی قبل و شناسایی و جداسازی هرچه بهتر، جداسازی از دو سمت انجام می‌شود. بدین صورت که در صورت بروز مشکل در روند شناسایی در یک زیرکلمه عملیات شناسایی و جداسازی از سمت دیگر انجام می‌شود که این خود باعث بهبود قابل ملاحظه‌ی درصد جداسازی و شناسایی می‌شود. در مرحله‌ی سوم با تصحیح نقاط اضافی تولید شده غلط‌های املائی بوجود آمده تصحیح می‌شوند.

همانطور که گفته شد عمل بازشناسی متن در چند فرایند انجام می‌شود که در هر مرحله ممکن است خطایی رخ داده و درصد شناسایی نهایی کاهش یابد. در مرحله‌ی اول پس از عملیات پیش پردازش، فرایند جداسازی زیرکلمات و حروف انجام شده که جهت رفع مشکل همپوشانی ممکن است خطایی رخ دهد. پس از برطرف نمودن مشکل همپوشانی، حروف و زیرکلمات به طبقه‌بند حروف گسسته فرستاده می‌شود. که در این مرحله نیز ممکن است خطایی رخ دهد. مثلاً برخی زیرکلمات بصورت حرف تشخیص داده شده و برای جداسازی به بخش جداسازی فرستاده نمی‌شود و یا برخی حروف به درستی شناسایی نمی‌شوند مانند شناسایی حرف "ه" بجای حرف "ه".

زیرکلمات شناسایی شده به بخش جداسازی فرستاده می‌شود. در این بخش ابتدا نقاط اولیه‌ی کاندید برای جداسازی انتخاب شده و سپس توسط طبقه‌بند حروف پیوسته نقاط جداسازی نهایی مشخص می‌شود که در هر دو بخش (بخش تعیین نقاط جداسازی اولیه و بخش جداسازی با طبقه‌بند حروف پیوسته) ممکن است خطایی رخ دهد.

با توجه به این که شناسایی هم شامل شناسایی زیرکلمات و حروف گسسته بوده و هم شناسایی برای حروف پیوسته در زمان جداسازی، لذا خطای رخ داده در مرحله شناسایی، مجموع خطای رخ داده در مرحله جداسازی بعلاوه خطای شناسایی حروف گسسته می‌باشد. که در نهایت و در مرحله سوم دقت در سطح جداسازی، ۹۶٪ و در سطح شناسایی نهایی ۸۵٪ حاصل شد.

۵-۵- نتیجه‌گیری

در این فصل به توضیح روند کلی سیستم بازشناسی متن پرداخته شد. در مرحله‌ی اول پیش‌پردازش‌هایی بر روی تصویر ورودی اعمال شد و در مرحله‌ی جداسازی با شناسایی حدود مکانی جمله و با استفاده از روش‌های مناسب و رفع مشکلات همپوشانی به جداسازی حروف گسسته و زیرکلمات پرداخته شد. با استفاده از طبقه‌بند حروف گسسته، حروف مجزا از زیرکلمات جدا شده و با توجه به شماره کلاس آن ذخیره گردید و زیرکلمات شناسایی شده در این مرحله به قسمت جداسازی حروف رفته و با روش‌های مختلف این جداسازی انجام شد. سپس با استفاده از طبقه‌بند حروف پیوسته به شناسایی بدنه حروف پرداخته شد و سپس بدنه‌ی این حروف به همراه نقطه‌های آن تحلیل و در کلاس‌های نهایی طبقه‌بندی گردید. در نهایت این حروف در کنار یکدیگر قرار گرفت و جمله بازیابی شد. جهت اعتبارسنجی، روند مذکور بر روی ۵ صفحه‌ی متن فارسی با درجه تفکیک ۳۰۰ نقطه بر اینچ آزمایش شده است که در آن دقت در سطح جداسازی، ۹۶٪ و در سطح شناسایی نهایی ۸۵٪ می‌باشد. این درصد در مقایسه با [۱۵] که برای قلم‌های میترا، نازنین و لوتوس انجام شده دارای خطای بیشتری بوده که این خطا ناشی از پیچیدگی‌های ساختاری و نوشتاری قلم Iranian sans می‌باشد.

فصل ششم

نتیجه‌گیری و پیشنهادات

۶-۱- نتیجه گیری

امروزه با افزایش تقاضا در همه‌ی زبان‌ها روز به روز نیاز به یک سامانه‌ی بازشناسی متن قوی افزایش می‌یابد. در زبان فارسی با توجه به تنوع بسیار زیاد قلم‌ها برای متون تایپی و همچنین وجود دست‌نوشته‌هایی که افراد با توجه به سلیقه‌ی شخصی و به اشکال متنوع می‌نویسند، برای دستیابی به یک سیستم بازشناسی قوی راه طولانی در پیش است. در این پایان‌نامه در فصل اول تعاریف اولیه‌ای از سیستم بازشناسی متن و کاربردهای آن مطرح شد، سپس بعضی از ویژگی‌های زبان فارسی که عملیات بازشناسی را با مشکل روبه‌رو می‌کند و بعضی از علت‌های عدم تکامل سیستم بازشناسی زبان فارسی مطرح شد. در فصل دوم به معرفی تحقیقاتی که در زمینه‌ی پیدایش ا.سی.آر انجام شده، پرداخته شد و همچنین روند تحقیقات ا.سی.آر بیان شد. سپس به مرور بعضی از مقالات در زمینه‌ی ا.سی.آر فارسی و همچنین روند تحقیقات فارسی پرداخته شد. در فصل سوم در مورد مراحل یک سیستم بازشناسی متن و روش‌های هر یک از آن‌ها به طور کلی توضیحاتی داده شد. در فصل چهارم به منظور فقدان پایگاه داده برای قلم Iranian sans، یک پایگاه داده‌ی نسبتاً کامل از حروف گسسته و پیوسته‌ی این قلم خاص به همراه اعداد و بعضی از علائم ایجاد شد تا به عنوان نمونه‌ی آموزش برای طبقه‌بندیها به کار گرفته شوند. در فصل پنجم مراحل مختلف سیستم بازشناسی متن به کار رفته در این پایان‌نامه به طور کامل توضیح داده شد. در این فصل ابتدا با اعمال عملیات پیش‌پردازش، سند ورودی شبکه آماده می‌شود، سپس با قطعه‌بندی جمله‌ها، کلمات و زیرکلمات عملیات جداسازی آغاز می‌شود. جداسازی با استفاده از روش پروفایل بالایی تعمیم یافته و سپس جاروب از طرفین زیرکلمه انجام می‌شود. سپس با افزودن نقاط به بدنه‌ی حروف و اختصاص شماره کلاس‌های مربوطه طبقه‌بندی نهایی انجام شده و در نهایت با توجه به شماره‌ی کلاس‌ها، زیرکلمه‌ها، کلمه‌ها و در نهایت جمله بازیابی می‌شود. در نتیجه یک سیستم بازشناسی متن برای قلم Iranian sans ارائه شد.

۶-۲- نوآوری‌ها

- ✓ ارائه‌ی ساختاری مناسب جهت حذف نویزهای موجود در تصویر
- ✓ تولید پایگاه داده برای حروف گسسته و پیوسته قلم Iranian sans
- ✓ پیاده‌سازی و آموزش طبقه‌بند حروف پیوسته و گسسته برای قلم Iranian sans
- ✓ غلبه بر مشکلات پیچیدگی‌های ساختاری قلم مذکور با افزودن زیرکلماتی مانند "سی"، "صی"، "لا"، "جی" "لله" به کلاس‌های حروف جدا و حروف پیوسته
- ✓ پیاده‌سازی روشی مناسب برای جداسازی مناسب زیرکلمات و غلبه بر مشکلات همپوشانی
- ✓ پیاده‌سازی الگوریتم مناسب جهت جداسازی بهینه‌ی زیرکلمات به حروف سازنده‌ی آن با استفاده از روش پروفایل بالایی تعمیم یافته و جاروب از طرفین زیرکلمه
- ✓ ایجاد کلاس نامعتبر برای طبقه‌بند حروف پیوسته جهت شناسایی بهتر
- ✓ ارائه‌ی الگوریتم مناسب برای حذف نقاط اضافی ایجاد شده در هنگام جداسازی و تصحیح نتایج شناسایی

۶-۳- پیشنهادات

- برای ادامه‌ی کار در زمینه‌ی سیستم‌های بازشناسی متن می‌توان موارد زیر را پیشنهاد داد.
- تلاش برای تقویت هر یک از مراحل بازشناسی متن، مثلاً در مرحله پیش‌پردازش برای باینری کردن، حذف نویز، اصلاح چرخش و تقویت سایر مراحل دیگر
 - استفاده از چند طبقه بند بصورت تقویت شده با الگوریتم آدابوست برای بالا بردن دقت شناسایی

- تلاش برای تولید سیستم‌های بازشناسی متن که بتواند متونی نوشته شده با بیشترین تعداد قلم و انواع زبان‌ها را شناسایی کند، که این امر نیاز به آشنایی با نوشتار زبان‌های مختلف و قلم‌های آن‌ها دارد.
- تلاش برای ساخت سامانه‌ای که نسبت به درجه تفکیک حساس نباشد و بتواند هر تصویر با هر درجه تفکیک را شناسایی کند. که طراحی چنین سامانه‌ای کار بسیار سختی است.
- طراحی سامانه‌ای که قابلیت شناسایی تصاویر خاکستری و رنگی را (بدون دوسطحی کردن) داشته باشد.
- برخی حروف یا ارقام در یک قلم چاپی شبیه به حروف یا ارقامی از قلم‌های دیگر هستند و نرم افزارهای اُ.سی.آر را به خصوص در متونی که با چند نوع قلم نوشته شدند به اشتباه می‌اندازند که یک راه حل می‌تواند بازشناسی نوع قلم پیش از بازشناسی متن باشد.
- طراحی روشی برای جداسازی زیرکلماتی از جمله "لم"، "می"، "سی"، "صی" برای قلم Iranian sans که به راحتی و با روش ذکر شده قابل جداسازی نیستند.

فهرست مراجع

- [1] D. Doermann, "The indexing and retrieval of document images: a survey," *Computer Vision and Image Understanding*, Vol. 70, No. 3, pp. 287-298, 1998.
- [2] H. Khosravi, E. Kabir, "Introducing a very large dataset of handwritten Farsi digits and a study on their varieties," *Pattern Recognition Letters*, Vol. 28, No. 10, pp. 1133-1141, Jul 2007.
- [3] N. Arica, F. T. Yarman-Vural, "An overview of character recognition focused on off-line handwriting," *IEEE Trans. Systems, Man, and Cybernetics*, part C, Vol. 31, No. 2, pp. 216-233, 2001.
- [4] S. Mori, C. Y. Suen, and K. Yamamoto, "Historical Review of OCR Research and Development," *Proc. of IEEE*, Vol. 80, No. 7, pp. 1029-1058, July 1992.
- [5] <http://www.uspto.gov>
- [6] M. Bokser, "Omnidocument Technologies," *Proceedings of the IEEE*, Vol. 80, No. 7, pp. 1066-1078, 1992.
- [7] <http://www.ncr.co.jp/library.html>
- [8] B. Parhami, and M. Taraghi, "Automatic recognition of printed Farsi texts," *Pattern Recognition Letters*, Vol. 14, No. 1-6, pp. 395-403, 1981.
- [۹] م. شیرعلی شهرضا و ک. فائز، "تشخیص کلمات و ارقام دست‌نویس فارسی به وسیله شبکه‌های عصبی (خط نسخ)"، رساله دکترای مهندسی برق - کامپیوتر، دانشگاه صنعتی امیرکبیر، ۱۳۷۴.
- [10] R. Azmi, E. Kabir, "A new segmentation technique for omnifont Farsi text," *Pattern Recognition Letters*, Vol. 22, No. 2, pp. 97-104, Feb. 2001.
- [۱۱] ر. عزمی، "بازشناسی متون چاپی فارسی"، رساله دکترای، بخش مهندسی برق، دانشگاه تربیت مدرس، ۱۳۷۸.

- [12] M. B. Menhaj, and M. Adab, "Simultaneous segmentation and recognition of Farsi/Latin printed texts with MLP," *International Joint Conference on Neural Networks*, Vol. 2, pp. 1534-1539, May 2002.
- [13] R. Mehran, H. Pirsiavash, F. Razzazi, "A Front-end OCR for Omni-font Persian/Arabic Cursive Printed Documents," *Digital Image Computing: Techniques and Applications*, pp. 385-392, Dec 2005.
- [14] A. Broumandnia, J. Shanbehzadeh, et al, "Segmentation of Printed Farsi/Arabic Words," *IEEE/ACS International Conference on Computer Systems and Applications*, pp.761-766,2007.
- [۱۵] ح. خسروی، و.ا. کبیر، "بازشناسی متن چاپی فارسی بر مبنای جداسازی هوشمند،" سومین کنفرانس بین‌المللی فناوری اطلاعات و دانش، مشهد، دانشگاه فردوسی مشهد، ۱۳۸۶.
- [۱۶] ح. بویری، م. عباسی دزفولی، و م. یکتایی، "بازشناسی نوری کاراکترهای چاپی فارسی با استفاده از شبکه عصبی فازی،" اولین کنفرانس دانشجویی فناوری اطلاعات ایران، سنندج، دانشگاه کردستان، ۱۳۸۹.
- [۱۷] ا. کیومرثی، ح. نظام آبادی پور، و ی. نوروززاده، "یک روش جدید برای جداسازی خطوط متن دست نوشته"، هفتمین کنفرانس ماشین بینایی و پردازش تصویر، تهران، دانشگاه علم و صنعت، ۱۳۹۰.
- [۱۸] م. زند، "تشخیص حروف چاپی فارسی با استفاده از روش ترکیبی،" همایش منطقه ای علوم کامپیوتر، مهندسی کامپیوتر و فناوری اطلاعات، دورود، دانشگاه آزاد اسلامی واحد دورود، ۱۳۹۱.
- [19] M. Dehghan, K. Faez, M. Ahmadi, M. Shridhar, "Unconstrained Farsi handwritten word recognition using fuzzy vector quantization and hidden markov models," *Pattern Recognition Letters*, Vol. 22 , No. 2, pp. 209-214, Feb 2001.
- [20] M. Dehghan , K. Faez, M. Ahmadi, M. Shridhar, "Handwritten Farsi (Arabic) word recognition: A holistic approach using discrete HMM," *Pattern Recognition* ,Vol. 34, No. 5 , pp. 1057-1065, May 2001.
- [21] S. Shirali-Shahreza , M.T. Manzuri-Shalmani, M.H. Shirali-Shahreza, "Preparing Persian/Arabic Scanned Images for OCR," *Information and Communication Technologies*, Vol.1, pp.1332-1336, 2006.

- [22] R. Halavati, S.B. Shouraki, "Recognition of Persian online handwriting using elastic fuzzy pattern recognition," *International Journal of Pattern Recognition and Artificial Intelligence*, Vol. 21 ,No. 3 , pp. 491-513,2007.
- [23] H. Khosravi, E. Kabir, "Farsi font recognition based on Sobel–Roberts features," *Pattern Recognition Letters*, Vol. 31 ,No.1 , pp. 75-82 ,January 2010.
- [32] M.H. Shirali-Shahreza, S. Shirali-Shahreza, "Removing Noises Similar to Dots from Persian Scanned Documents," *International Colloquium on IEEE*, Vol. 2, pp.313-317,2008.
- [25] S. Khalighi,P. Tirdad, H.R. Rabiee, M. Parviz, "A Novel OCR System for Calculating Handwritten Persian Arithmetic Expressions," *International Conference on IEEE*, pp.755-758,2009.
- [26] E. Arianyan, S.A. Motamedi, I. Arianyan, " Efficient Optical Character Recognition on Graphics Processing Unit, " *Sixth International Symposium on IEEE*, pp.789-793, 2012.
- [27] A. and J. Kanai, "Projection profile based skew estimation algorithm for JBIG," *compressed images, ICDAR*,pp. 401- 405,1997.
- [28] L. Najman, "Using mathematical morphology for document skew estimation," *SPIE Document Recognition and Retrieval XI*, pp. 182-191, 2004 .
- [29] R. v. d. Boomgaard, and R. v. Balen, "Methods for fast morphological image transform using bitmapped binary images", *Graphical Models and Image Processing*, pp.252-258 , 1992.
- [۳۰] ح. خسروی، و ا. کبیر، "یافتن زاویه چرخش سند مبتنی بر زوایای خطوط متن با استفاده از عملگر مورفولوژی"، سومین کنفرانس بین‌المللی فناوری اطلاعات و دانش، مشهد، دانشگاه فردوسی مشهد، ۱۳۸۶.
- [31] B. Al-Badr, S. A. Mahmoud, "Survey and bibliography of Arabic optical text recognition," *Signal Processing*, Vol. 41, pp. 49-77,1995.

[32] H. Nashida, and S. Mori, "An Algebraic Approach to Automatic Construction of Structured Models," *Pattern Analysis and Machine Intelligence*, Vol. 15, No.12, pp. 1298-1311, 1993.

[33] H. Y. Abdelazim, and M. A. Hashish, "Arabic Reading Machine", *Proc. of the 10th National Computer Conference*, Jeddah, pp. 733-744, 1988.

[34] Houle, G. and M. Shridhar, "Handwritten Word Recognition With OCR-based Segmenter," *Proc. of the Workshop on Document Image Analysis*, pp. 51-58, 1997.

[۳۵] ح. نظام آبادی پور، ا. کبیر و ر. عزمی، "الگوریتم اصلاح شده جداسازی حروف در متون چاپی با برچسب زدن به کانتور بالایی کلمات" استقلال ۱۳۸۰.

[36] F. Kimura, M. Shridhar, et al., "Improvements of a Lexicon Directed Algorithm for Recognition of Unconstrained Handwritten Words," *Proc. of 2nd ICDAR Conf*, pp.18-22, 1993.

[37] A. Cheung, M. Bennamoun, et al., "An Arabic optical character recognition system using recognition-based segmentation," *Pattern Recognition*, Vol. 34, pp. 215-233, 2001.

[38] B.B. Chaudhuri, and U. Garain, "Automatic Detection of Italic, Bold and All-Capital words in Document Images". *14th International Conference on Pattern Recognition*, pp. 610-612, 1998.

[39] C.B. Jeong, H. K. Kwag, et al., "Identification of Font Styles and Typefaces in Printed Korean," *6th International Conference on Asian Digital Libraries*, Vol. 2911, pp. 666-669, 2003.

[40] A. Zramdini, and R. Ingold, "Optical Font Recognition Using Typographical Features," *Pattern Analysis and Machine Intelligence*, Vol. 20, No. 8, PP. 877-882, 1998.

[۴۱] ح. خسروی و ا. کبیر، معرفی دو ویژگی سریع و کارآمد برای بازشناسی ارقام دستنویس فارسی، چهارمین کنفرانس ماشین بینایی و پردازش تصویر، مشهد، دانشگاه فردوسی مشهد، ۱۳۸۵.

Abstract

In recent decades, extensive research, writing recognition patterns include letters, numbers and other symbols commonly used in written documents in different languages is done. Given the progress made in the field of automatic text recognition technology called optical character recognition or OCR is formed. Text recognition is considered as an important part of e-government in our country and in recent years as the demand for a Persian text recognition system greatly enhanced. Due to the large amount of paper documents, digital image documents are converted by the scanner or camera, storage, efficient management and retrieval of these documents, files, in many applications, including office automation, and digital libraries are important.

In general, text recognition system includes several parts , such as receiving image preprocessing, configuration analysis, diagnostics language, font and finally text recognition. Research conducted in some topics , such as preprocessing is independent of the text language and can be used with any language . But some other topics , such as font recognition depends on the context and results of research conducted for other languages can not directly be applied to Persian . Most researches in the field of Persian literature on the recognition of images with high resolution images and text clean and false and document identification are with some known font. In research conducted for the recognition of Persian texts there are three approaches that based on separation words based on the overall shape recognition and mix of them.

This thesis aims at recognizing typed text written with the Iranian sans font, with a minimum size of 9 and the resolution is 300 dpi. according to The font style and readability, it is more addressed and every day the volume of computing and Internet environment is enhanced. This font style is replacement for the default Windows operating system font likeTahma, Despite readability, standard spacing between rows, beauty and consistency of Latin, this font has a structural complexity that complicates the process of recognition.

In this thesis after the production of database suitable, discrete and continuous characters classification were trained. then with solving the problem of overlapping of the words, separation approach is used to separate letters. Finally, the performance result of the system for processing some images of printed text, is provided where the separation accuracy of 96% and an accuracy of 85% was achieved in the identification.

Keywords: text recognition pen Iranian sans, separation-based approach, neural network classification