

الله اعلم



دانشگاه الزهرا(س)  
دانشکده فنی و مهندسی

پایان نامه  
جهت اخذ درجه کارشناسی ارشد  
رشته مهندسی کامپیوuter گرایش هوش مصنوعی

عنوان  
بازیابی بر اساس محتوای اسناد چاپی فارسی

استاد راهنما  
دکتر رضا عزمی

دانشجو  
زهرا بهمنی

۱۳۹۰ مهر

کلیه دستآوردهای این تحقیق متعلق به  
دانشگاه الزهراء(س) است

تقدیم به پدر و مادرم که همیشه در کنارم هستند

سپاس خداوندگارم که مرا در به پایان رساندن این پایان نامه یاری کرد تا برگی دیگر از دفتر زندگی ورق بخورد.  
بر خود واجب می دانم تا،  
از زحمات استاد راهنمای بزرگوارم جناب آقای دکتر عزمی که با دانش خود و همراهی بی-  
دریغشان در تمامی مراحل این پایان نامه مرا یاری نمودند،  
از جناب آقای دکتر شیری و آقای دکتر قلی زاده که داوری این پایان نامه را بعهده گرفتند،  
از دوست و همراهم جناب آقای وحید قلندری که آرام بخش لحظه هایم بودند، همچنین آقای  
صمصام احسانی به خاطر کمکشان،  
از خانواده ام که همیشه و همه جا پشتیبانم بودند، پدر و مادرم، خواهرانم مهین ، عصمت و  
فاطمه و برادرانم رضا و علی و همچنین محمد احمدی و شیروان علی زاده،  
از دوستان خوبیم فاطمه علمدار، سونیا اربابیان، محبوبه غلامپور و حمیده رفیعی،  
و از همه کسانی که در تمامی مراحل تحصیلم تا کنون یاری دهنده و امید راهم بودند، تشکر  
کنم. برای همه این عزیزان آرزوی توفیق روز افزون دارم.  
و به یاد عزیزی که سال هاست در کنار ما نیست برادر عزیزم هاشم.(روحش شاد)  
همچنین جا دارد از مرکز مخابرات ایران که از این پایان نامه حمایت کردند قدردانی کنم.

## چکیده

با افزایش کتابخانه‌های دیجیتال و برای دستیابی به هدف ادارات بدون کاغذ تعداد زیادی از کپی‌ها به دیجیتال تبدیل شده و در سیستم مدیریت اسناد ذخیره شده است. همچنین در حال حاضر میلیون‌ها سند دیجیتال دائماً بر روی اینترنت از یک نقطه به نقطه‌ی دیگر منتقل می‌شوند.

اگر چه تکنولوژی پردازش تصویر اسناد می‌تواند برای تبدیل اتوماتیک تصاویر دیجیتال این اسناد به فرمت متن قابل خواندن به وسیله کامپیوتر با استفاده از بازشناسی کارکتر نوری استفاده شود ولی این روش برای حجم عظیمی از اسناد بهینه و کارا نیست. با توجه به این شرایط یافتن راه حل بازیابی اسناد پرینت شده به رویی که نیاز به برگرداندن این اسناد به نسخه‌ی متنی نداشته باشد ضروری به نظر می‌آید.

روش‌های بازیابی و بازشناسی به دو دسته اصلی تقسیم می‌شوند. دسته اول بازیابی را بر اساس توصیف شکل کلی کلمات یا زیرکلمات انجام می‌دهند. دسته دوم کلمه را به حروف می‌شکنند و از هر تصویر حرف ویژگی استخراج می‌کنند. در روش‌های مبتنی بر جداسازی علاوه بر مشکلاتی مانند وجود نقاط و علائم و تنوع قلمها، مشکل جداسازی حروف نیز وجود دارد.

در این پژوهش ما از رویی برای استخراج ویژگی‌های کلمات استفاده کرده‌ایم که نیاز به تخمین نقاط جداسازی ندارد. برای این منظور از شناسایی اتصال دهنده‌های عناصر اصلی حروف استفاده شده است. با شناسایی این اتصال دهنده‌ها و حذف آنها عناصر اصلی که در این پژوهش با نام زیرحروف نامگذاری شده‌اند استخراج شده است.

در این پژوهش از سه روش زیر برای تشخیص زیرحروف استفاده شده است.

در روش اول زیرحروف با استفاده از درخت تصمیم و به کمک ویژگی‌های شکلی زیرحروف مانند ارتفاع، عرض، حفره‌ها، گودال‌ها، دره‌ها، فروفتگی‌ها و موقعیت عنصر نسبت به خط زیمینه تشخیص داده شده‌اند. دقت تشخیص در این روش برای زیرحروف بین ۸۰ تا ۱۰۰ درصد بوده است.

در روش دوم از ترکیبی از درخت تصمیم نخکشی شده و شبکه عصبی RBF برای تشخیص زیرحروف استفاده شده است. در این روش علاوه بر ویژگی‌های شکلی کلمه، ویژگی نمایه‌ها در چهار جهت بکار برده شده است. نتایج روش دوم برای زیرحروف بین ۹۰ تا ۱۰۰ ارزیابی شده است.

سومین روش مورد استفاده شبکه عصبی چند سطحی می‌باشد. در این روش تشخیص تنها به وسیله ویژگی نمایه‌ها و در سه تا چهار سطح و با استفاده از شبکه عصبی انجام شده است. تشخیص در این روش برای اغلب زیرحروف‌ها بالای ۹۵ درصد بوده است.

در نهایت با کد کردن زیرحروف، زیرکلمات و کلمات موجود در تصویر سند کد شده و برای بازیابی کلمات کلیدی مورد نظر کاربر استفاده شده است. دقت سیستم برای کلمات با طول متغیر به طور متوسط بالای ۹۰ درصد ارزیابی شده است.

**کلمات کلیدی:** بازیابی تصویر سند، زیرحروف، درخت تصمیم، شبکه عصبی RBF، نمایه.

## فهرست مطالب

۱	۱- مقدمه
۲	۱-۱- تکنولوژی پردازش تصویر
۳	۱-۲- بازیابی به کمک بهبود روش‌های بازشناسی
۴	۱-۳- بازیابی مستقل از بازشناسی
۵	۱-۴- رویکردهای بازیابی متون
۵	۱-۵- اهداف و نوآوری
۷	۱-۶- ساختار پایان نامه
۸	۲- پیشینه تحقیق
۹	۲-۱- کتابخانه های دیجیتال
۱۰	۲-۲- بازیابی اسناد فارسی
۱۱	۲-۳- ۲- الگوریتم‌های بازیابی
۱۱	۲-۳-۱- روش HMM
۱۴	۲-۳-۲- خوشه‌بندی
۲۰	۲-۳-۳- روش ماشین بردار پشتیبان
۲۳	۲-۳-۳-۱- مکانیسم خوشه‌بندی بردار پشتیبان
۲۴	۲-۳-۴- شبکه عصبی
۲۷	۲-۴- روش‌های استخراج ویژگی
۲۷	۱-۴-۲- کدکردن شکل کلمه
۳۰	۲-۴-۲- هیستوگرام
۳۱	۳-۴-۲- کانتور
۳۳	۴-۴-۲- نمایه
۳۳	۴-۵- نازک سازی
۳۴	۱-۵-۴-۲- قالب تغییرپذیر

۳۵	-۲-۴-۵-۲- ویژگی‌های مجزا
۳۵	-۲-۴-۶- زنجیر کدها
۳۶	-۲-۵- بخش‌های گرافیکی
۳۷	-۲- دست خط
۳۷	-۲- بازیابی طرح و ظاهر
۳۷	-۲- جمع‌بندی
۳۸	-۳- سیستم بازیابی سند
۳۹	-۳- تفاوت شاخص گذاری سیستم بازیابی تصویر سند با سیستم بازیابی تصویر
۴۰	-۳- ۲- بخش‌های مختلف سیستم بازیابی سند
۴۰	-۳- ۲- ۱- پرس و جو
۴۳	-۳- ۲- ۲- واسط کاربر
۴۴	-۳- ۲- ۳- استخراج ویژگی
۴۴	-۳- ۲- ۴- مدل سند
۴۵	-۳- ۲- ۵- بازیابی ردبهتی شده
۴۶	-۳- ۳- ارزیابی سیستم بازیابی تصویر سند
۴۸	-۳- ۱- روش‌های موجود برای ارزیابی میزان تشخیص
۴۹	-۳- ۴- پایگاه داده تصویر اسناد
۴۹	-۳- ۵- جمع‌بندی
۵۰	-۴- چهارچوب پیشنهادی
۵۱	-۴- ۱- ویژگی‌های نوشتار فارسی
۵۱	-۴- ۲- چهارچوب کلی سیستم
۵۳	-۴- ۳- پیش پردازش
۵۴	-۴- ۴- پردازش
۵۵	-۴- ۴- ۱- بدنه زیرکلمات و علائم
۵۶	-۴- ۴- ۲- زیرحروف
۵۷	-۴- ۴- ۲- ۱- تشخیص عناصر اتصال دهنده زیرحروف

۵۷	- تشخیص زیرحروف	۴-۴-۲-۲-۴
۸۹	- انتساب نقاط و علائم به زیرحروف	۴-۴-۳-۳
۸۹	- تولید کد	۴-۴-۵
۹۰	- بازیابی	۴-۶-۶
۹۱	- جمع بندی	۴-۷-۷
۹۲	- پیاده سازی و ارزیابی	۵-۵-پیاده
۹۳	- پیش پردازش	۵-۱-پیش
۹۳	- پردازش	۵-۲-پردازش
۹۳	- استخراج بدنی زیر کلمات و علائم	۵-۴-۱-استخراج
۹۴	- تشخیص عناصر اتصال دهنده زیرحروف	۵-۵-۲-تشخیص
۹۵	- تشخیص زیرحروف	۵-۵-۳-تشخیص
۹۵	- تشخیص زیرحروف به کمک کد شکل زیرحروف و درخت تصمیم	۵-۵-۲-۳-۱-تشخیص
۹۷	- طبقه بندی کننده ترکیبی درخت تصمیم و شبکه عصبی RBF	۵-۵-۲-۳-۲-طبقه
۹۸	- طبقه بندی کننده شبکه عصبی RBF چند سطحی	۵-۵-۳-۲-۳-طبقه
۱۰۰	- انتساب نقاط و علائم به زیرحروف	۵-۵-۴-۲-انتساب
۱۰۰	- تولید کد	۵-۵-۳-۵-تولید
۱۰۱	- بازیابی	۵-۴-۴-بازیابی
۱۰۴	- مقایسه روش بخشندهای پیشنهادی	۵-۵-۶-۱-مقایسه
۱۰۴	- نتیجه گیری	۵-۷-۷-نتیجه
۱۰۵	- نتیجه گیری و توسعه های آتی	۶-۱-نتیجه
۱۰۶	- نتیجه گیری	۶-۱-نتیجه
۱۰۸	- توسعه های آتی	۶-۲-توسعه
۱۰۹	مراجع	۹-۱-مراجع

## فهرست شکل

شکل ۲-۱: خوشه بندی نمونه های ورودی ..... ۱۴
شکل ۲-۲: تصاویر ماشین بردار پشتیبان ..... ۲۱
شکل ۲-۳: نمونه های آموزشی به همرا نمایش حاشیه، بردار پشتیبان و فاصله ..... ۲۲
شکل ۲-۴: طرح کلی خوشبندی بردار پشتیبان ..... ۲۴
شکل ۲-۵: شبکه عصبی ..... ۲۵
شکل ۲-۶: ویژگیهای شکلی مکانی ..... ۲۸
شکل ۲-۷: نقاط آغازین و پایانی اجراهای سفید عمودی و افقی ..... ۲۹
شکل ۲-۸: هیستوگرام عمودی و افقی ..... ۳۱
شکل ۲-۹: تعیین کانتورها در یک زیرکلمه ..... ۳۲
شکل ۲-۱۰: نمایه کلمه علی ..... ۳۳
شکل ۲-۱۱: اعمال الگوریتم نازکسازی ..... ۳۴
شکل ۲-۱۲: استاندارهای حرکت در زنجیر کد ..... ۳۵
شکل ۲-۱۳: بردار ویژگی هیستوگرام زنجیر برای یک قاب ..... ۳۶
شکل ۳-۱: اطلاعات ساختار فیزیکی و منطقی ..... ۴۱
شکل ۳-۲: استخراج بلوهای تصویر سند ..... ۴۲
شکل ۳-۳: نمونهای از رابط کاربر گرافیکی در یک سیستم بازیابی ..... ۴۴
شکل ۳-۴: مدل سند در بازیابی تصویر اسناد و انتزاع بوجود آمده بوسیله آن ..... ۴۵
شکل ۳-۵: شبه کد الگوریتم پیشنهادی ..... ۴۷
شکل ۳-۶: انطباق دو کلمه که در ابتدا و انتهای کلمه با هم اختلاف دارند ..... ۴۸
شکل ۴-۱: چهارچوب کلی سیستم و فرآیند پیشپردازش ..... ۵۴
شکل ۴-۲: شبه کد الگوریتم مربوط به برچسب زنی عناصر متصل ..... ۵۵
شکل ۴-۳: انواع نقاط ..... ۵۶
شکل ۴-۴: تشخیص عناصر اتصال دهنده ..... ۵۷
شکل ۴-۵: نمونه هایی از دسته بندی زیرحروف ..... ۵۸
شکل ۴-۶: زیرحروف و ویژگی های آنها ..... ۶۰
شکل ۴-۷: درخت تصمیم زیر حرف ابتدایی ..... ۶۲
شکل ۴-۸: درخت تصمیم زیر حروف میانی ..... ۶۳

شکل ۴-۹: درخت تصمیم زیر حروف پایانی با حفره .....	۶۴
شکل ۴-۱۰: درخت تصمیم زیر حروف پایانی بدون حفره .....	۶۵
شکل ۴-۱۱: درخت تصمیم زیر حروف جدا با حفره .....	۶۶
شکل ۴-۱۲: درخت تصمیم زیر حروف جدا بدون حفره .....	۶۷
شکل ۴-۱۳: شبیه کد الگوریتم Kmens .....	۷۱
شکل ۴-۱۴: تصویر مربوط به این نمایه زیر حرف "ی" و اعمال تابع موجک گستته .....	۷۲
شکل ۴-۱۵: درخت تصمیم هرس شده و RBF مربوط به زیر حروف ابتدایی .....	۷۴
شکل ۴-۱۶: درخت تصمیم هرس شده و RBF مربوط به زیر حروف میانی .....	۷۵
شکل ۴-۱۷: درخت تصمیم هرس شده و RBF مربوط به زیر حروف پایانی با حفره .....	۷۶
شکل ۴-۱۸: درخت تصمیم هرس شده و RBF مربوط به زیر حروف پایانی بدون حفره .....	۷۷
شکل ۴-۱۹: درخت تصمیم هرس شده و RBF مربوط به زیر حروف جدا با حفره .....	۷۸
شکل ۴-۲۰: درخت تصمیم هرس شده و RBF مربوط به زیر حروف پایانی بدون حفره .....	۷۹
شکل ۴-۲۱: شبکه عصبی چند لایه برای زیر حروف ابتدایی .....	۸۶
شکل ۴-۲۲: شبکه عصبی چند لایه برای زیر حروف میانی .....	۸۶
شکل ۴-۲۳: شبکه عصبی چند لایه برای زیر حروف انتهایی .....	۸۷
شکل ۴-۲۴: شبکه عصبی چند لایه برای زیر حروف جدا .....	۸۸
شکل ۴-۲۵: شبیه کد الگوریتم نیو .....	۹۱

شکل ۵-۱: تشخیص عنصر اتصال دهنده .....	۹۴
شکل ۵-۶: نمودار تشخیص زیر حروف ابتدایی به وسیله طبقه بندی کننده ترکیبی .....	۹۹
شکل ۵-۷: نمودار تشخیص زیر حروف میانی به وسیله طبقه بندی کننده ترکیبی .....	۹۹
شکل ۵-۸: نمودار تشخیص زیر حروف انتهایی به وسیله طبقه بندی کننده ترکیبی .....	۹۹
شکل ۵-۹: نمودار تشخیص زیر حروف ج به وسیله طبقه بندی کننده ترکیبی .....	۹۹
شکل ۵-۱۰: خروجی برنامه برای جستجوی حرف "الف" .....	۱۱۸
شکل ۵-۱۱: خروجی برنامه برای جستجوی حرف "ج" .....	۱۱۹
شکل ۵-۱۲: خروجی برنامه برای جستجوی کلمه "در" .....	۱۲۰
شکل ۵-۱۳: خروجی برنامه برای جستجوی کلمه "لا" .....	۱۲۱
شکل ۵-۱۴: خروجی برنامه برای جستجوی کلمه "کتاب" .....	۱۲۲
شکل ۵-۱۵: خروجی برنامه برای جستجوی کلمه "دیجیتا" .....	۱۲۳

## فهرست جداول

جدول ۴-۱: شکل های مختلف حروف الفبای فارسی ..... ۵۲
جدول ۴-۲: حروف و زیرحروف مربوط به آن ..... ۵۹
جدول ۴-۳: نتایج شبکه عصبی آموزش دیده با همه زیرحروف ابتدایی ..... ۸۰
جدول ۴-۴: نتایج شبکه عصبی آموزش دیده با حذف زیرحروف شاخص و ادغام زیرحروف با تداخل زیاد(سطح اول) برای زیرحروف ابتدایی ..... ۸۱
جدول ۴-۵: نتایج شبکه عصبی آموزش دیده با حذف زیرحروف شاخص و ادغام زیرحروف با تداخل زیاد(سطح دوم) ..... ۸۱
جدول ۴-۶: نتایج شبکه عصبی آموزش دیده با همه زیرحروف میانی ..... ۸۲
جدول ۴-۷: نتایج شبکه عصبی آموزش دیده با حذف زیرحروف شاخص و ادغام زیرحروف با تداخل زیاد(سطح اول) برای زیرحروف میانی ..... ۸۲
جدول ۴-۸: نتایج شبکه عصبی آموزش دیده با حذف زیرحروف شاخص و ادغام زیرحروف با تداخل زیاد(سطح دوم) برای زیرحروف میانی ..... ۸۲
جدول ۴-۹: نتایج شبکه عصبی آموزش دیده با همه زیرحروف پایانی ..... ۸۳
جدول ۴-۱۰: نتایج شبکه عصبی آموزش دیده با حذف زیرحروف شاخص و ادغام زیرحروف با تداخل زیاد(سطح اول) برای زیرحروف پایانی ..... ۸۳
جدول ۴-۱۱: نتایج شبکه عصبی آموزش دیده با حذف زیرحروف شاخص و ادغام زیرحروف با تداخل زیاد(سطح دوم) برای زیرحروف پایانی ..... ۸۴
جدول ۴-۱۲: نتایج شبکه عصبی آموزش دیده با همه زیرحروف جدا ..... ۸۴
جدول ۴-۱۳: نتایج شبکه عصبی آموزش دیده با حذف زیرحروف شاخص و ادغام زیرحروف با تداخل زیاد(سطح اول) برای زیرحروف جدا ..... ۸۵
جدول ۴-۱۴: نتایج شبکه عصبی آموزش دیده با حذف زیرحروف شاخص و ادغام زیرحروف با تداخل زیاد(سطح دوم) برای زیرحروف جدا ..... ۸۵
جدول ۴-۱۵: دیکشنری تعریف شده برای چهار گروه زیرحروف، ابتدایی، میانی، انتهایی و جدا ..... ۹۰
جدول ۵-۱: دیکشنری نهایی ..... ۱۰۱
جدول ۵-۲: نتایج بازیابی حروف الفبا ..... ۱۰۲
جدول ۵-۳: درصد بازیابی صحیح کلمات (جدول ۱) ..... ۱۰۳
جدول ۵-۴: درصد بازیابی صحیح کلمات (جدول ۲) ..... ۱۰۳
جدول ۵-۵: مقایسه روش پیشنهادی با روش ارائه شده در [۱۰۵] ..... ۱۰۳

# فصل اول

مقدمه

تکنیک‌های مدرن امکان تولید، پردازش انتقال و ذخیره تصاویر را به طور مناسبی فراهم آورده است. در نتیجه مقدار اطلاعات تصویری<sup>۱</sup> با سرعت شتابنده‌ای در حال افزایش است. برای مثال برای دستیابی به هدف ادارات بدون کاغذ<sup>۲</sup> تعداد زیادی از کپی‌ها به دیجیتال تبدیل شده و در سیستم مدیریت اسناد ذخیره شده است. اما این پایگاه داده‌ها دارای اطلاعات شاخص مناسب و کافی نیستند. این مسائل بازیابی اطلاعات مورد علاقه کاربر از تصاویر را مشکلتر از بازیابی آنها از داده‌های متنی کرده است [۱]. همچنین هر روزه میلیون‌ها سند دیجیتالی به وسیله اینترنت از نقطه‌ای به نقطه دیگر جابجا می‌شود. قالب اکثر این اسناد، متن می‌باشد که در آن کاراکترها با کدهای قابل خواندن به وسیله کامپیوتر نمایش داده شده‌اند. همچنین اغلب اسنادی که به وسیله کامپیوتر ایجاد می‌گردد به همین قالب است. از طرف دیگر برای تبدیل میلیون‌ها سند سنتی و قدیمی موجود و در دسترس گذاشتن آنها از طریق اینترنت، این اسناد به وسیله ابزارهای الکترونیکی و دیجیتالی اسکن شده و به نسخه دیجیتال تبدیل می‌گردند. از سویی کتابخانه‌های دیجیتال بزرگی در حال حاضر به وجود آمده‌اند که مجموعه بزرگی از اسناد پرینت شده را بایگانی کرده‌اند. برای این که کاربران بتوانند به جستجو در این گونه پایگاه‌داده‌ها بپردازنند. باید با استفاده از روش‌های بازیابی براساس محتوای تصاویر اسناد امکان جستجو در این مجموعه‌ها را فراهم آورد.

## ۱-۱- تکنولوژی پردازش تصویر

پردازش تصاویر امروزه بیشتر به موضوع پردازش تصویر دیجیتال گفته می‌شود. پردازش تصاویر دارای دو شاخه عمده بهبود تصاویر<sup>۳</sup> و بینایی ماشین<sup>۴</sup> است. بهبود تصاویر دربرگیرنده روش‌هایی چون استفاده از فیلتر محوکننده و افزایش کنتراست برای بهتر کردن کیفیت دیداری تصاویر و اطمینان از نمایش درست آنها در محیط مقصد(مانند چاپگر یا نمایشگر رایانه) است، در حالی که بینایی ماشین به روش‌هایی می‌پردازد که به کمک آنها می‌توان محتوای تصاویر را درک کرد. هر چند از تکنولوژی پردازش تصاویر اسناد می‌توان کمک گرفت و تصاویر دیجیتالی اسناد موجود را به کمک تکنولوژی بازناسی نوری حروف<sup>۵</sup> (OCR) به قالب متنی قابل خواندن به وسیله کامپیوتر تبدیل کرد. پردازش حجم بزرگی از صفحات اسناد به وسیله این روش به دلایل

<sup>1</sup> . Visual

<sup>2</sup> . Paperless Office

<sup>3</sup> . Image Enhancement

<sup>4</sup> . Machine Vision

<sup>5</sup> . Optical Character Recognition

زیادی به صرفه و عملی نمی‌باشد. از جمله این دلایل اینکه تکنیک‌های تحلیل طرح‌بندی سند در مواجهه با اسناد با طرح‌بندی‌های پیچیده هنوز بسیار ساده و ابتدایی بوده و نتایج خروجی آنها مشکلات زیادی دارند. دلیل دیگر قدرت ذاتی پایین تکنولوژی بازشناسی حروف بخصوص در مواجهه با تصاویر اسناد با کیفیت پایین می‌باشد. تصحیح خروجی‌های بازشناسی نوری حروف به صورت دستی نیز در اغلب سیستم‌های پردازش تصویر سند ممکن نیست. بنابراین، می‌توان نتیجه گرفت که ذخیره اسناد سنتی و قدیمی به شکل تصویر راه حل عملی در بیشتر این موارد است. امروزه در اینترنت اسناد دیجیتالی زیادی را در قالب تصویری می‌توان یافت، مانند مقالات کنفرانس‌ها، مجلات اسکن شده، رساله دانشجویان و غیره. به علاوه بیشتر کتابخانه‌های دیجیتال و پورتال‌های وب مانند IEE, ACM, MEDLINE و غیره تصاویر اسکن شده اسناد را بدون نسخه معادل متنی آنها نگهداری می‌کنند.

کاربری که با تعداد زیادی تصویر سند در اینترنت برخورد می‌کند باید تمام آنها را دانلود کند. زیرا مکانیزمی مانند جستجوی کلمات کلیدی در تصاویر وجود ندارد که کاربر بتواند از محتواهای آنها اطلاع یابد. با توجه به این شرایط یافتن راه حلی مناسب برای بازیابی اسناد پرینت فارسی ضروری به نظر می‌آید.

## ۱-۲- بازیابی به کمک بهبود روش‌های بازشناسی

بازیابی تصاویر اسناد و پردازش تصویر سند رابطه نزدیکی با هم دارند، اما تفاوت‌های خاص بین آنها نیز وجود دارد. یک سیستم پردازش تصویر سند، بخش‌های مختلف متن را در یک صفحه سند تحلیل می‌کند، روابط بین این بخش‌های متنی مختلف را کشف می‌کند و آنها را به کمک یک سیستم بازشناسی نوری حروف به نسخه قابل خواندن توسط ماشین تبدیل می‌کند. برای زبان انگلیسی سیستم‌های بازشناسی نوری حروف زیادی با دقتهای بالای ۹۹ درصد در شناسایی حروف تولید شده‌اند و در دهه اخیر بسیاری از آنها به تولید تجاری و فروش رسیده‌اند.

می‌دانیم که دقتهای بازشناسی مورد نیاز در روش‌های بازیابی تصویر سند از دقتهای مورد نیاز در بسیاری از کاربردهای پردازش سند پایین‌تر است. با توجه به این نکته در سال‌های اخیر متدهای برای بازیابی اسناد پیشنهاد شده است که توانایی مواجهه با خطاهای بازشناسی که در بازشناسی نوری حروف اتفاق می‌افتد را دارند همچنین روش‌هایی برای بهبود روش‌های که از بازشناسی نوری استفاده می‌کنند معرفی شده است. پرسش اصلی که یک سیستم بازیابی

تصویر سند باید به آن پاسخ دهد این است که یک تصویر سند شامل کلمات مورد علاقه کاربر می‌باشد یا خیر، و سعی کند این پاسخ را بدون توجه و پردازش کلمات نامرتبط انجام دهد. به عبارت دیگر یک سیستم بازیابی تصویر سند با توجه به پرس‌وجو کاربر یک پاسخ بله یا خیر خواهد داد، اما مانند پردازش تصویر سند بازشناسی کامل را برای حروف و کلمات انجام نخواهد داد. در بازیابی اطاعات واحد اصلی با معنی کلمه می‌باشد نه حرف بنابراین، تطبیق مستقیم کلمات موجود در تصاویر اسناد راه حل دیگری برای بازیابی اطاعات از تصویر سند خواهد بود. به طور خلاصه بازیابی تصویر اسناد و پردازش تصویر سند نیازهای متفاوتی را پاسخگو هستند و هر کدام ارزش و اهمیت خاص خود را دارد. روش‌هایی که قصد دارند که اطاعات را به طور مستقیم از تصویر بازیابی کنند می‌توانند کارآیی بالاتری بر مبنای دقت، یادآوری و سرعت پردازش داشته باشند.

### ۱-۳-بازیابی مستقل از بازشناسی

در سال‌های اخیر محققان تلاش‌های بسیاری کرده‌اند تا روش‌هایی طراحی کنند که در بازیابی تصاویر اسناد از تکنیک‌های بازشناسی حروف استفاده نکنند [۲و۳]. برای مثال «چن» و «بلومبورگ» [۴] روشی برای انتخاب جملات و کلمات کلیدی برای ایجاد خلاصه از یک متن پیشنهاد کرده‌اند. روش آنها برای ایجاد خلاصه از تصویر سند نیاز به بازشناسی حروف هر کلمه ندارد. آنها کلاس‌های کلمات همارز را با استفاده از تبدیل hit-miss برای مقایسه کلمات ایجاد کردن و از یک طبقه‌بندی آماری برای تعیین شباهت هر جمله با جمله خلاصه استفاده کرده‌اند. کارهای [۴و۵و۷و۸] بازیابی اسناد تصویری را از طریق تطبیق شکل کلمه انجام می‌دهند. کارهای [۹و۱۰و۱۱و۱۲و۱۳و۱۴و۱۵و۱۶و۱۷] بازیابی اسناد تصویری را از طریق کردن شکل کلمه انجام می‌دهند. «شیجیان لو» و «چیو لیم تان» هم تصویر سند را از طریق کد کردن شکل کلمه، بازیابی می‌کنند [۱۸]. در واقع هر کلمه تصویر را به کد شکل کلمه با استفاده از ویژگی‌های شکل کلمه تبدیل می‌کند. همین افراد شناسایی زبان را با استفاده از کد کردن شکل کلمه انجام می‌دهند [۲۰و۱۹]. آنها همچنین جستجوی کلمات را بر اساس کد کردن شکل کلمه انجام می‌دهند [۲۱و۲۲]. «اسپایتز» و «مقبولة» از کد شکل کلمه برای طبقه‌بندی سند استفاده می‌کند [۲۳]. در [۲۳] دو کاربرد بازیابی تصویر سند را بر اساس کد کردن شکل کلمه برای استخراج ویژگی از هر شیء تصویر کلمه نشان می‌دهد. اولین کاربرد آن یک سیستم بازیابی مبتنی بر وب می‌باشد که تصاویر سند را به صورت بلادرنگ بر پایه یک

مجموعه از کلمات پرسش ورودی، از کتابخانه دیجیتال بازیابی می‌کند. دومین کاربرد آن یک ابزار جستجوی plug-in تعبیه شده در Acrobat Reader است.<sup>[۲۴]</sup> روشی را بر پایه توزيع چگالی دو بعدی پیشنهاد می‌کند که "Pseudo relevance feedback" را برای بسط دادن یک پرسش اولیه جهت بهتر شدن کارایی بازیابی به کار می‌گیرد.

#### ۱-۴- رویکردهای بازیابی متون

روش‌های بازیابی و بازناسی متون از دو رویکرد مبتنی بر جداسازی<sup>۱</sup> کلمات به حروف و زیرحروف<sup>۲</sup> و رویکرد مبتنی بر شکل کلی کلمات، استفاده می‌کنند [۲۵].

روشهای مبتنی بر جداسازی خود به دو گروه تقسیم می‌شوند. گروه اول روشهایی هستند که کلمه را به حروف تشکیل دهنده آن جداسازی می‌کنند و هر حرف را بطور مجزا، با استفاده از ویژگیهای استخراج شده از آن بازناسی می‌کنند. در این روشهای بازناسی کلمه با استفاده از حروف شناخته شده و پسپردازش‌هایی مبتنی بر واژه نامه انجام می‌شود.

گروه دوم، روشهایی هستند که کلمه را به حروف می‌شکنند و از هر تصویر حرف ویژگی استخراج می‌کنند. سپس با کنار هم قرار دادن این ویژگیها یک توصیفگر کلمه استخراج می‌کنند و به کمک آن در واژه نامه تصویری بهترین گزینه را انتخاب می‌کنند. در روشهای مبتنی بر جداسازی، علاوه بر مشکلاتی مانند وجود نقاط و علائم و تنوع قلمها، مشکل جداسازی حروف نیز وجود دارد. مناسب نبودن کیفیت سند، پایین بودن درجه تفکیک یا کجی تصویر سند نیز مشکلاتی هستند که به سختی کار می‌افزایند.

در رویکرد مبتنی بر شکل کلی کلمات یا زیرکلمات، توصیفگرها مستقیماً از تصویر کلمه یا زیرکلمه استخراج می‌شوند و با استفاده از یک واژه نامه تصویری کلمه مناسب استخراج می‌شود. برای بازناسی کلمه یا زیرکلمه می‌توان از روشهای مختلف توصیف شکل استفاده کرد.

#### ۱-۵- اهداف و نوآوری

در این پژوهش سعی شده است که راهکار مناسبی برای بازیابی اسناد چاپی فارسی ارائه شود که از دقت قابل قبولی برخوردار باشد. با توجه به این که اغلب کارهای انجام شده در زمینه بازیابی اسناد فارسی از طریق شکل کلی کلمات و یا زیرکلمات<sup>۳</sup> بوده است و اغلب از دقت

<sup>1</sup>. Segmentation

<sup>2</sup>. Sub-letter

<sup>3</sup>. Sub-word

پایینی برخوردار هستند، ارائه راهکاری که بازیابی را به روش جداسازی انجام می‌دهد امری ضروری به نظر می‌آید.

همچنین با توجه به مشکلاتی که در زمینه بخش‌بندی کلمات فارسی به حروف موجود در آن وجود دارد و کار بازیابی به این روش را با خطای تشخیص نقطه جدا کننده روبرو کرده است، یافتن راه حلی که بتواند بخش‌بندی را بدون نیاز به تشخیص نقطه جداسازی انجام دهد بهبود قابل توجهی را در فرآیند بازیابی در اسناد فارسی خواهد داشت.  
در زیر موارد بررسی شده و راهکارهای ارائه شده در این پژوهش ذکر شده است.

- حروف فارسی بررسی شده، عناصر اصلی سازنده حروف فارسی شناسایی و نحوه اتصال آنها مورد بررسی قرار گرفته‌اند.
- روشی برای استخراج این عناصر بدون نیاز به تشخیص نقطه انفصال دو حرف ارائه شده است. در این روش حروف شامل دو بخش هستند. بخش‌های اصلی حروف که در این پژوهش با نام زیرحروف نامگذاری شده‌اند و شامل بخش‌های سازنده حرف است که ویژگی‌های اصلی حرف مورد نظر را در بر می‌گیرد و اتصال‌دهنده که این عناصر را به هم متصل می‌کند.
- با توجه به عناصر استخراج شده دیکشنری تصویری تعریف شده است که شامل همه عناصر اصلی موجود در حروف فارسی است که با اتصال‌دهنده‌ها به هم متصل شده‌اند.
- زیرحروف استخراج شده با سه روش زیر تشخیص داده شده‌اند.
  - با استفاده از درخت تصمیم<sup>۱</sup> و به کمک ویژگی‌های شکلی.
  - با استفاده از ترکیبی از درخت تصمیم نخ‌کشی شده و شبکه عصبی<sup>۲</sup> RBF و RBF<sup>۳</sup> و به کمک ویژگی‌های شکلی<sup>۴</sup> و نمایه.
  - به کمک شبکه عصبی RBF<sup>۵</sup> چند سطحی<sup>۶</sup> و با استفاده از ویژگی نمایه.<sup>۷</sup>
- زیرحروف تشخیص داده شده با استفاده دیکشنری تصویری کد شده و کد زیرکلمات با کناره‌هم گذاشتن کد زیرحروف ساخته شده‌اند، با ترکیب آن‌ها کد کلمات و در نتیجه کد محتوای سند تولید شده است.

<sup>1</sup>. Decision Tree

<sup>2</sup>. Radial-Basis Function Neural Network

<sup>3</sup>. Topological Shape

<sup>4</sup>. Cascade RBF

<sup>5</sup>. Profile

- کلمات کلیدی مورد نظر با توجه به حروف موجود در کلمه و دیکشنری تعریف شده به کد تبدیل شده، حاصل آن در کدمحتوای اسناد جستجو شده است.

## ۱-۶-ساختار پایان نامه

در این پایان نامه سایر بخش‌ها به صورت زیر سازماندهی می‌شود:

در فصل دوم به بررسی کارهای انجام شده در زمینه بازیابی سند پرداخته می‌شود و الگوریتم‌ها و راهکارهای ارائه شده در این زمینه مورد ارزیابی قرار می‌گیرد. فصل سوم به بررسی ساختار کلی یک سیستم بازیابی سند پرداخته و یک ساختار کلی برای سیستم بازیابی تصویر اسناد تشریح می‌شود. در فصل چهارم روش پیشنهادی برای بازیابی تصویر اسناد فارسی مطرح شده و جزئیات هر بخش توضیح داده می‌شود. فصل پنجم به نحوه پیاده‌سازی سیستم شرح داده شده می‌پردازد در اینجا نتایج پیاده‌سازی بیان شده است. سرانجام در فصل شش نتیجه‌گیری و راه-کارهای پیشنهادی برای کارهای آتی ذکر شده است.

# فصل دوم

پیشینه تحقیق