

بِسْمِ اللَّهِ الرَّحْمَنِ الرَّحِيمِ



دانشکده مهندسی کامپیوتر و فناوری اطلاعات

رساله دکتری هوش مصنوعی

جستجوی کلمه در اسناد تصویری فارسی مبتنی بر خصیصه‌ها

نگارنده

مجید ایرانپور مبارکه

استاد راهنما

دکتر علیرضا احمدی فرد

استاد مشاور

دکتر حمید حسن پور

بهمن ۱۳۹۵

دانشکده : مهندسی کامپیوتر و فناوری اطلاعات

گروه : هوش مصنوعی

پایان نامه رساله دکتری آقای مجید ایرانپور مبارکه

تحت عنوان: جستجوی کلمه در اسناد تصویری فارسی مبتنی بر خصیصه‌ها

در تاریخ ۱۳۹۵/۱۱/۱۸ توسط کمیته تخصصی زیر جهت اخذ مدرک دکتری ارزیابی گردید و با درجه بسیار خوب مورد پذیرش قرار گرفت.

امضاء	اساتید مشاور	امضاء	اساتید راهنما
	نام و نام خانوادگی: دکتر حمید حسن پور		نام و نام خانوادگی: دکتر علیرضا احمدی فرد

امضاء	نماینده تحصیلات تکمیلی	امضاء	اساتید داور
	نام و نام خانوادگی:		نام و نام خانوادگی:
			نام و نام خانوادگی:
			نام و نام خانوادگی:

سپاس و ستایش مرخدای راجل و جلالة که آثار قدرت او بر چهره روز روشن، تابان است و انوار حکمت او در دل شب تار، در نشان.
آفریدگاری که خویش را به ما شناساند و درهای علم را بر ما گشود و عمری و فرصتی عطا فرمود تا بدان، بنده ضعیف خویش را در طریق علم و معرفت

بیازماید.

تقدیم به پدرم به استواری کوه،

مادرم به زلالی چشمه

و

همسرم به صمیمیت باران.

سپاس خدای را که سخنوران، دستوران او بمانند و شمارندگان، شمردن نعمت‌های او ندانند و کوشندگان، حق او را گزارش ندهند. و دور بر محمد و
 خاندان پاک او، طاهران معصوم، هم آنان که وجودمان و امدار وجودشان است؛ و نفرین پیوسته بردشمنان ایشان تا روز رستاخیز بدون شک جایگاه و
 منزلت معلم، اجل از آن است که در مقام قدردانی از زحمات بی‌شائبی او، با زبان قاصر و دست ناتوان، چیزی بنگاریم.
 اما از آنجایی که تجلیل از معلم، سپاس از انسانی است که هدف و غایت آفرینش را تا این می‌کند و سلامت امانت‌هایی را که به دستش سپرده‌اند، تضمین؛
 بر حسب وظیفه و از باب "من لم یشکر المنعم من المخلوقین لم یشکر الله عز و جل:"
 از پدر و مادر عزیزم، این دو معلم بزرگوارم، که همواره بر کوتاهی و درستی من، قلم عفو کشیده و گریانه از کنار غفلت‌هایم گذشته‌اند و در تمام عرصه‌های زندگی یاد
 و یاور بی‌چشم داشت برای من بوده‌اند؛
 از استاد با کمال و شایسته؛ جناب آقای دکتر علیرضا احمدی فرد که در کمال سع صدر، با حسن خلق و فروتنی، از هیچ کجی در این عرصه بر من دریغ ننمودند و
 زحمت راهنمایی این رساله را بر عهده گرفتند؛
 از استاد صبور و باتقوا، جناب آقای دکتر حمید حسن پور، مدیرت محترم دانشکده، که زحمت مشاوره این رساله را در حالی متمثل شدند که بدون مساعدت
 ایشان، این پروژه به نتیجه مطلوب نمی‌رسید؛

تعهد نامه

این جانب مجید ایرانیپور مبارکه دانشجوی دوره دکتری رشته مهندسی کامپیوتر/هوش مصنوعی دانشکده مهندسی کامپیوتر و فناوری اطلاعات دانشگاه صنعتی شاهرود نویسنده‌ی پایان‌نامه جستجوی کلمه در اسناد تصویری فارسی مبتنی بر خصیصه‌ها تحت راهنمایی دکتر علیرضا احمدی فرد متعهد می‌شوم :

- تحقیقات در این پایان‌نامه توسط این جانب 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 پایان‌نامه 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 پایان‌نامه تأثیرگذار بوده‌اند در مقالات مستخرج از پایان‌نامه رعایت می‌گردد.
- در کلیه مراحل انجام این پایان‌نامه، در مواردی که از موجود زنده (یا بافتهای آن‌ها) استفاده شده است ضوابط و اصول اخلاقی رعایت شده است.
- در کلیه مراحل انجام این پایان‌نامه، 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، برنامه‌های رایانه‌ای، نرم افزارها و تجهیزات ساخته شده است) متعلق به دانشگاه صنعتی شاهرود می‌باشد. این مطلب باید به نحو مقتضی در تولیدات علمی مربوطه ذکر شود.
- استفاده از اطلاعات و نتایج موجود در پایان‌نامه بدون ذکر مرجع مجاز نمی‌باشد.

چکیده

جستجوی کلمه، روشی برای اندیس‌گذاری، جستجو و بازیابی کلمات در اسناد تصویری است. در این مسئله هدف مشخص کردن موقعیت رخ داده‌های کلمه پرسشی مورد نظر کاربر در اسناد تصویری است. این زمینه تحقیقاتی در زمره چالش‌های مهم در حوزه آنالیز اسناد تصویری است که در سال‌های اخیر مورد توجه محققان قرار گرفته است. کارهای متنوعی در این حوزه بر روی اسناد لاتین انجام شده، اما در مورد زبان‌های فارسی و عربی کارهای محدودی انجام شده است که اکثراً در زمینه اسناد چاپی می‌باشد. هدف از این رساله ارائه یک سیستم کامل جهت جستجوی کلمات پرسشی دلخواه از اسناد تصویری فارسی است که وابسته به شیوه نگارش نویسنده نبوده و در اسناد دارای چند نویسنده نیز به خوبی کار می‌کند. علاوه بر این سیستم مذکور توانایی دریافت کلمه پرسشی به صورت متنی و تصویری را نیز دارا می‌باشد.

از آنجایی که تعداد کلاس‌ها (کلمات زبان) در این مسئله بسیار زیاد است و از طرف دیگر تنوع درون کلاسی ناشی از شیوه‌های نگارش متفاوت نیز وجود دارد، نیاز به استفاده از رده‌بندی مبتنی بر خصیصه احساس می‌شود. رده‌بندی مبتنی بر خصیصه با تبادل بین کلاسی، توانایی شناسایی کلاس‌های دیده نشده را دارد. در این روش، متن (برچسب کلاس کلمات) و تصویر در یک زیر فضای برداری مشترک تعبیه می‌شوند. این زیر فضا بر اساس چارچوب کاری که از رده‌بندی مبتنی بر خصیصه، تعبیه برچسب و رگرسیون زیر فضای مشترک ایجاد شده است، کار می‌کند. در این زیر فضای ایجاد شده، متن و تصویر یک کلمه یکسان، در زیر فضای مشترک آموزش داده شده، نزدیک به هم قرار می‌گیرند و مسئله جستجوی کلمات تبدیل به یک مسئله ساده نزدیک‌ترین همسایه می‌گردد.

در گام اول، این روش جهت جستجو و بازیابی کلمات دست‌نویس مجزا و اسناد قطعه‌بندی شده به کاررفته است. سپس بر اساس ساختار نوشتار زبان فارسی خصیصه‌های متناسب پیشنهاد شده است. دو روش بازیابی جهت خصیصه‌ها زبان فارسی معرفی شده که یکی بر اساس شکل متفاوت حروف در موقعیت‌های مختلف و دیگری بر اساس دسته‌بندی حروف مبتنی بر اشتراک بدنه‌ها، نقاط و سرکش‌ها در حروف فارسی پیشنهاد شده است. به منظور ارزیابی روش‌های پیشنهادی از مجموعه داده‌های فارسی که شامل کلمات مجزای دست‌نویس هستند، فارسا و ایرانشهر، استفاده شده است. میانگین دقت متوسط خصیصه‌های پیشنهادی در جستجوی کلمه برای این مجموعه داده‌ها به ترتیب ۹۵/۶۷ و ۹۶/۲۵ درصد و دقت بازیابی کلمات دست‌نویس به ترتیب ۹۶/۱۵ و ۹۷/۵۵ درصد به دست آمده است.

در گام دوم، جستجوی سند در سطح خط انجام شد. با استفاده از روشی پیشنهادی که مبتنی بر خوشه‌بندی سلسله مراتبی می‌باشد، خطوط هر سند استخراج می‌شود و سپس پنجره‌گذاری روی خطوط استخراج شده صورت

می‌پذیرد. برای فائق آمدن بر مشکلات پنجره گذاری، یک روش پنجره گذاری مبتنی بر اجزاء متصل پیشنهاد شده است. در این روش اندازه پنجره اهمیتی ندارد بلکه تعداد جزء متصلی که پنجره آن را در برمی‌گیرد مهم است. این پنجره‌ها به‌عنوان کلمات مستقل در نظر گرفته شده و بر اساس روش ارائه‌شده مبتنی بر زیرفضای مشترک برای کلمات قطعه‌بندی شده، بازیابی می‌شوند. به‌منظور ارزیابی روش پیشنهادی از مجموعه داده متون دست‌نویس فارسی (FHT)، که مجموعه‌ای از اسناد دست‌نویس فارسی است، استفاده شده است و در مورد کلمه پرسشی متنی و تصویری، میانگین دقت متوسط به‌دست‌آمده به ترتیب برابر $81/02$ و $78/14$ درصد می‌باشد.

از آنجایی که همیشه قطعه‌بندی کار پرخطا و زمان‌بری است، روش‌های مستقل از قطعه‌بندی از جایگاه ویژه‌ای برخوردارند. از این‌رو در ادامه روشی جهت جستجوی کلمات مستقل از قطعه‌بندی پیشنهاد شده است که به‌صورت دو مرحله‌ای کار می‌کند. در مرحله اول اجزاء متصل در کلمه پرسشی در سند جستجو می‌شود. سپس بر اساس اولویت اجزاء متصل بازیابی شده از سند و فاصله بین آن‌ها، نواحی کاندید برای کلمه مورد جستجو تعیین می‌شوند. جهت اطمینان نواحی کاندید گسترش داده می‌شوند. در مرحله بعد، رتبه‌بندی این نواحی کاندید بر اساس کلمه پرسشی ورودی انجام می‌پذیرد. این رتبه‌بندی با استفاده از سیستم پنجره‌گذاری استفاده شده در سطح خط صورت می‌پذیرد. جهت ارزیابی روش پیشنهادی در این مرحله نیز از مجموعه داده FHT استفاده شده است و نتایج تجربی به‌دست‌آمده برای کلمات پرسشی تصویری و متنی به ترتیب $75/07$ و $80/28$ درصد می‌باشد.

نکته حائز اهمیت در ارزیابی روش‌های پیشنهادی مبتنی بر قطعه‌بندی خط و مستقل از قطعه‌بندی این است که آموزش مدل با استفاده از مجموعه داده فارسی و ایرانشهر (کلمات مجزا) انجام شده است و مرحله آزمون با مجموعه داده FHT انجام پذیرفته است. این کار توانایی روش پیشنهادی در مقاوم بودن نسبت به تغییرات شیوه نگارش را نشان می‌دهد. علاوه بر این هیچ‌کدام از کلاس‌های کلمات مورد جستجو در مرحله آموزش مورد استفاده قرار نگرفته است، بنابراین روش پیشنهادی توانایی جستجوی هر کلمه دلخواه را به کاربر می‌دهد.

کلمات کلیدی: جستجوی کلمه در اسناد تصویری، اسناد دست‌نویس فارسی، تعبیه برچسب، رده‌بندی مبتنی بر خصیصه، اجزاء متصل

مقالات مستخرج از رساله

مقالات ژورنالی

- ✓ Isolated Persian/Arabic word spotting by label embedding, **Journal of Engineering Research** (JER), December 2016, Accepted.
- ✓ Segmentation-free Word Spotting on Persian Multi-writer Handwritten image documents. Ready to submit.
- ✓ Line-based Persian word spotting by label embedding. Ready to submit.

- ✓ ارائه خصیصه‌های خاص زبان فارسی جهت بازیابی و بازشناسی کلمات تصویری فارسی با استفاده از تعبیه برچسب. **مجله ماشین بینایی و پردازش تصویر**، چاپ شده.
- ✓ مروری بر روش‌های جستجوی کلمه در اسناد تصویری، **مجله محاسبات نرم**، تحت داوری.

مقالات کنفرانسی

- ✓ ارائه و مقایسه چندین روش تعبیه برچسب جهت بازیابی و بازشناسی کلمات در اسناد تصویری دست‌نویس فارسی (**مقاله برتر کنفرانس**) نهمین کنفرانس ماشین بینایی و پردازش تصویر ایران، ۱۳۹۴
- ✓ جداسازی خطوط متن و استخراج برون خط دنباله‌ای از مؤلفه‌های متصل در دست‌نوشته فارسی، نهمین کنفرانس ماشین بینایی و پردازش تصویر ایران، ۱۳۹۴
- ✓ مکان‌یابی متون فارسی در تصاویر صحنه طبیعی با استفاده امضای رنگ لبه‌های تصویر، هفتمین کنفرانس بین‌المللی فناوری اطلاعات و دانش، ۱۳۹۴
- ✓ جستجوی کلمات در یک تصویر دست‌نوشته فارسی با استفاده از روش تعبیه مشخصه‌ها، سومین کنفرانس ملی پژوهش‌های کاربردی در مهندسی کامپیوتر، ۱۳۹۴
- ✓ استخراج خط در اسناد دست‌نویس فارسی مبتنی بر خوشه‌بندی سلسله مراتبی، هشتمین کنفرانس بین‌المللی فناوری اطلاعات و دانش، ۱۳۹۵

فهرست مطالب

۱- مقدمه	۱
۱-۱- تعریف مسئله	۲
۲-۱- اهداف و چالش‌ها	۳
۳-۱- مجموعه داده‌ها	۱۰
۴-۱- ساختار پایان‌نامه	۱۳
۲- مروری بر کارهای پیشین	۱۵
۱-۲- مقدمه	۱۶
۲-۲- دسته‌بندی بر اساس روش جستجو	۱۷
۱-۲-۲- روش‌های مبتنی بر قطعه‌بندی	۱۷
۲-۲-۲- روش‌های مستقل از قطعه‌بندی	۲۵
۳-۲- دسته‌بندی بر اساس روش تشخیص	۲۸
۱-۳-۲- روش‌های مبتنی بر تطبیق الگو	۲۹
۲-۳-۲- روش‌های مبتنی بر یادگیری	۳۲
۴-۲- تقسیم‌بندی بر اساس نوع کلمه‌ی پرسشی	۳۴
۳- روش پیشنهادی اول (مبتنی بر قطعه‌بندی کلمه)	۴۱
۱-۳- مقدمه	۴۲
۲-۳- بازنمایی کلمه مبتنی بر خصیصه	۴۵
۱-۲-۳- تعبیه برچسب متنی	۴۶
۲-۲-۳- آموزش خصیصه‌ها بر اساس برچسب‌های تعبیه‌شده	۴۹
۳-۲-۳- زیرفضای مشترک بین برچسب‌ها و خصیصه‌ها	۵۴
۳-۳- نتایج تجربی و ارزیابی	۵۸

- ۳-۳-۱- معیارهای ارزیابی ۶۱
- ۳-۳-۲- نتایج تجربی ۶۲
- ۴- روش پیشنهادی دوم (مبتنی بر خط) ۶۷
- ۴-۱- مقدمه ۶۸
- ۴-۲- قطعه‌بندی خطوط ۶۹
- ۴-۲-۱- پیش‌پردازش ۷۲
- ۴-۲-۲- خوشه‌بندی ۷۵
- ۴-۲-۳- پس‌پردازش ۷۷
- ۴-۲-۴- نتایج تجربی ۷۸
- ۴-۳- پنجره‌گذاری مبتنی بر اجزاء متصل ۷۹
- ۴-۴- نتایج تجربی ۸۱
- ۵- روش پیشنهادی سوم (مستقل از قطعه‌بندی) ۸۵
- ۵-۱- مقدمه ۸۶
- ۵-۲- تعیین نواحی کاندید ۸۸
- ۵-۳- مرتب‌سازی مجدد ۹۱
- ۵-۴- نتایج تجربی ۹۲
- ۶- نتیجه‌گیری و کارهای آینده ۹۹
- ۶-۱- نتیجه‌گیری ۱۰۰
- ۶-۲- کارهای آینده ۱۰۱

فهرست شکل‌ها

- شکل ۱-۱. برخی ویژگی‌های زبان فارسی [۶] ۶
- شکل ۲-۱. نمونه‌ای از اسناد دست‌نویس فارسی ۶
- شکل ۳-۱. فلوچارت کلی آموزش روش پیشنهادی مبتنی بر قطعه‌بندی ۸
- شکل ۴-۱. فلوچارت کلی روش پیشنهادی مبتنی بر قطعه‌بندی در سطح خطوط ۹
- شکل ۵-۱. فلوچارت کلی روش پیشنهادی مستقل از قطعه‌بندی ۱۰
- شکل ۶-۱. چند نمونه از تصویر کلمات مجموعه داده ایرانشهر ۱۱
- شکل ۷-۱. چند نمونه از تصویر کلمات مجموعه داده فارسا ۱۱
- شکل ۸-۱. چند نمونه از تصویر کلمات مجموعه داده FHT ۱۲
- شکل ۱-۲: دسته‌بندی روش‌های موجود جستجوی کلمه ۱۶
- شکل ۲-۲: نمایش حفره و مخازن آب در برخی از حروف نمونه انگلیسی [۱۰] ۱۸
- شکل ۳-۲: ساختار کلی روش ارائه‌شده در [۱۲] ۲۱
- شکل ۴-۲: نمونه‌ای از ویژگی‌های ساختاری چند کلمه فارسی [۱۲] ۲۱
- شکل ۵-۲: نمونه‌هایی از قطعه‌بندی وفقی کلمات در [۱۵] ۲۲
- شکل ۶-۲: شمای کلی شبکه عصبی با حافظه کوتاه‌مدت طولانی دووجهته [۳۰] ۲۴
- شکل ۷-۲: نمونه‌ای از قطعه‌بندی به بخشی کلمه عربی برای کلمه "تاریخ" ۲۴
- شکل ۸-۲: فرآیند آزمون رده‌بند سلسله مراتبی [۳۱]. الف) قطعه‌بندی. ب) نتایج رده‌بندی. ج) ساخت گراف. د) مدل‌های زبانی ۲۵
- شکل ۹-۲: مقایسه دقت و فراخوانی در روش مبتنی بر قطعه‌بندی و مستقل از آن در عربی [۳۲] ۲۶
- شکل ۱۰-۲: قطعه‌بندی کلمه عربی "نقاط" [۳۲] ۲۶
- شکل ۱۲-۲: شمای کلی روش معرفی‌شده در [۳۴] ۲۸
- شکل ۱۳-۲: تشخیص زیرکلمات بر اساس ویژگی‌های ساختاری برای کلمه پرسشی "آبادی" [۴۲] ۳۰
- شکل ۱۴-۲: نمونه‌ای از ترکیب نواحی یافت شده برای زیرکلمات کلمه پرسشی "آبادی" [۴۲] ۳۰
- شکل ۱۵-۲: تشریح روش ناحیه‌بندی جهت استخراج ویژگی [۱۴] ۳۱
- شکل ۱۶-۲: روند استخراج ویژگی از یک فریم از پنجره لغزان [۴۱] ۳۲
- شکل ۱۷-۲: روند کلی روش ارائه‌شده در [۵۴] ۳۴
- شکل ۱۸-۲: حرکت پنجره لغزان روی یک خط‌نمونه و ایجاد نواحی کاندید [۵۵] ۳۴
- شکل ۱۹-۲: فرآیند کلی روش پیشنهادشده در [۶۰] با نمایش تغییرات نسبت به کار قبلی در [۵۹] ۳۷

- شکل ۲-۲۰: تفاوت انتخاب نواحی مورد توجه در [۵۹] و [۶۰] ۳۷
- شکل ۲-۲۱: شمای کلی سیستم [۶۲] ۳۸
- شکل ۳-۱: شمای کلی روش رده‌بندی مبتنی بر خصیصه ۴۳
- شکل ۳-۲: شمای کلی روش معرفی‌شده در [۷, ۴۱, ۶۸] ۴۴
- شکل ۳-۳: روش تعبیه SPOC و مدل تعبیه برچسب در [۶۹] ۴۵
- شکل ۳-۴: فرآیند آموزش مدل خصیصه lam [۷] ۴۵
- شکل ۳-۵: بازنمای PHOC در سه سطح ۱ و ۲ و ۳ برای رشته متنی "beyond" [۷] ۴۷
- شکل ۳-۶: ساختار کلی آموزش مدل خصیصه‌ها ۵۱
- شکل ۳-۷: روند کلی استخراج ویژگی در روش پیشنهادی ۵۲
- شکل ۳-۸: یک نمونه از قطعه‌بندی تصویر کلمه (کلمه "شاهرود") ۵۳
- شکل ۳-۹: نگاهت امتیازات پیش‌بینی‌شده خصیصه‌ها و برچسب‌های تعبیه‌شده در یک زیرفضای دارای همبستگی بیشتر [۷] ۵۵
- شکل ۴-۱: ساختار کلی روش پیشنهادی مبتنی بر قطعه‌بندی خط ۶۸
- شکل ۴-۲: ساختار کلی روش پیشنهادی جهت قطعه‌بندی خطوط ۷۳
- شکل ۴-۳: الف) سند قطعه‌بندی شده به نوارهای عمودی ب) پروفایل افکنش افقی نوار مشخص‌شده با علامت ستاره سبز در تصویر الف به همراه نمایش قله‌ها. ۷۵
- شکل ۴-۴: یک سند نمونه خروجی مرحله خوشه‌بندی. ۷۷
- شکل ۴-۵: یک سند نمونه خروجی نهایی ۷۹
- شکل ۴-۶: یک سطر از یک سند نمونه با چهارضلعی‌های محاطی اجزاء متصل ۷۹
- شکل ۴-۷: نمونه‌های از خروجی پنجره به طول ۲ جزء متصل ۸۰
- شکل ۴-۸: نمونه‌های از خروجی پنجره به طول ۳ جزء متصل ۸۰
- شکل ۴-۹: نمونه‌های از خروجی پنجره به طول ۴ جزء متصل ۸۰
- شکل ۴-۱۰: نمونه‌های از خروجی پنجره به طول ۵ جزء متصل ۸۱
- شکل ۵-۱: مرحله برون‌خط روش پیشنهادی ۸۶
- شکل ۵-۲: مرحله برخط روش پیشنهادی ۸۷
- شکل ۵-۳: الف) نمونه‌ای از بازیابی اجزاء متصل در سند بر اساس کلمه پرسشی "تولید". هر یک از اجزاء متصل کلمه پرسشی در سند نمونه با رنگی متفاوت نشان داده شده است. ب) نواحی کاندید تعیین شده بر اساس اجزاء متصل بازیابی شده. ۸۹

- شکل ۴-۵: نمونه‌هایی از اجزاء متصل استفاده‌شده در آموزش سیستم جستجو و بازیابی اجزاء متصل ۹۰
- شکل ۵-۵: نمونه‌ای از به هم چسبیدگی اجزاء متصل. کلمات "سرخ"، "دارم" و "می‌گیرم" ۹۶
- شکل ۶-۵: نمونه‌ای از شکستگی اجزاء متصل. کلمات "بخش" و "دوستانی" ۹۶
- شکل ۷-۵: نمونه‌ای از خط خوردگی. کلمات "شبدر"، "حیوان" و "قبله‌ام" ۹۶
- شکل ۸-۵: نمونه‌ای از ساختار متفاوت با استانداردهای نگارشی. کلمات "تکه"، "کاشانم"، "دشت"، "باعث"، "من" و "فروتن" ۹۷

فهرست جدول‌ها

- جدول ۱-۱: شکل‌های مختلف حروف الفبای فارسی با توجه به محل قرار گرفتن در کلمه ۴
- جدول ۱-۲: کدگذاری شکل حروف انگلیسی ارائه شده در [۱۰] ۱۹
- جدول ۲-۲: کدگذاری به کار رفته در حروف عربی ارائه شده در [۱۱] ۲۰
- جدول ۳-۲: مقایسه روش‌های جستجوی کلمه با کلمه پرسشی متنی و تصویری ۳۶
- جدول ۱-۳۳: حروف دارای نقاط و سرکش مشترک ۴۹
- جدول ۲-۳: حروف دارای بدنه مشترک ۴۹
- جدول ۳-۳: نتایج دقت روش پیشنهادی جهت بازشناسی کلمه ۶۳
- جدول ۴-۳: نتایج میانگین دقت متوسط (MaP) روش پیشنهادی جهت جستجوی کلمه ۶۴
- جدول ۵-۳: نتایج دقت در اولین بازیابی ($P@1$) روش پیشنهادی جهت جستجوی کلمه ۶۵
- جدول ۱-۴: تعداد کلمات با فراوانی بین ۱ تا ۷ در هر سند مجزا ۸۲
- جدول ۲-۴: میانگین دقت متوسط با استفاده از روش مبتنی بر قطعه‌بندی خط ۸۳
- جدول ۳-۴: محاسبه دقت متوسط (AveP) برای چند کلمه پرسشی نمونه ۸۳
- جدول ۴-۴: لیست نمونه‌هایی از کلمات بازیابی شده توسط روش پیشنهادی با توجه به تنوع شیوه نگارش ۸۴
- جدول ۱-۵: پنجره‌های نمونه برای یک ناحیه کاندید برای کلمه پرسشی "تولید" ۹۲
- جدول ۲-۵: تعداد کلمات با فراوانی بین ۱ تا ۷ در هر سند مجزا ۹۳
- جدول ۳-۵: نرخ تشخیص بر اساس میزان اشتراک ناحیه کاندید با ناحیه واقعی ۹۵
- جدول ۴-۵: نرخ کاهش بر اساس معیارهای معرفی شده ۹۵
- جدول ۵-۵: میانگین دقت متوسط (MaP) برای کلمات پرسشی تصویری و متنی ۹۷
- جدول ۶-۵: میانگین دقت اولین بازیابی ($P@1$) برای کلمات پرسشی تصویری و متنی ۹۷

واژه‌نامه

Latent semantic indexing (LSI)	اندیس گذاری معنایی پنهان
Regularized discriminant analysis (RDA)	آنالیز تفکیکی منظم
Joint representation	بازنمایی مشترک
Content-based image retrieval (CBIR)	بازیابی محتوایی تصویر
Part of arabic word (PAW)	بخشی کلمه عربی
Fisher vector (FV)	بردار فیشر
Dymanic time wrapping (DTW)	پیچش زمان پویا
Compatibility function	تابع سازگاری
Canonical correlation analysis (CCA)	تحلیل همبستگی کانونی
Kernelized canonical correlation analysis (KCCA)	تحلیل همبستگی کانونی دارای هسته
Maximum likelihood estimation (MLE)	تخمین حداکثر درستنمایی
Optical character recognition (OCR)	تشخیص کاراکتر نوری
Document image analysis and recognition (DIAR)	تشخیص و تحلیل اسناد تصویری
Flexible sequence matching (FSM)	تطبیق انعطاف پذیر توالی
Label embedding	تعبیه برچسب
Input embedding	تعبیه ورودی
Word spotting	جستجوی کلمه
Out of vocabulary (OOV)	خارج از فرهنگ لغت
Attribute	خصیصه
Precision	دقت
Attribute-based classification	رده‌بندی مبتنی بر خصیصه
Common subspace regression (CSR)	رگرسیون زیرفضای مشترک
Bag-of-visual-word (BoVW)	سبد بصری کلمات
Scale invariant feature transform (SIFT)	سیفت
Recurrent neural network (RNN)	شبکه‌های عصبی بازگشتی
Bidirectional long short-term memory (BLSTM)	شبکه عصبی با حافظه کوتاه مدت بلند دو جهته
recall	فراخوانی

Codebook	کتابچه کد
Character shape coding(CSC)	کد گذاری شکل کاراکتر
Word shape coding (WSC)	کد گذاری شکل کلمه
Query by example (QBE)	کلمه پرسشی تصویری
Query by string (QBS)	کلمه پرسشی متنی
Min average horizontal distance (MAHD)	کمینه متوسط فاصله افقی
Min average vertical distance (MAVD)	کمینه متوسط فاصله عمودی
Exemplar SVM	ماشین بردار پشتیبان مبتنی بر یک نمونه
Ground Truth	مجموعه درستی
Discriminative model	مدل تمایزی
Generative model	مدل زایشی
Hidden markov model (HMM)	مدل مخفی مارکوف
Semi-continuous hidden markov model (SC-HMM)	مدل مخفی مارکوف شبه گسسته
Gussian mixture model (GMM)	مدل مخلوط گاوسی
Platts Scaling	مقیاس بندی پلاتس
Mean average precision (Map)	میانگین دقت متوسط
Connected components reduction rate (CCRR)	میزان کاهش تعداد اجزاء متصل
Detection rate (DR)	نرخ تشخیص
Pixel level hit rate (PLHR)	نرخ تطابق پیکسلی
Spatial pyramid of characters (SPOC)	هرم مکانی کاراکترها
Histogram of gradient (HOG)	هیستوگرام گرادیان
Persian pyramid histogram of characters (PPHOC)	هیستوگرام هرمی کاراکترهای فارسی
Persian pyramid histogram of positional characters (PPHOPC)	هیستوگرام هرمی کاراکترهای موقعیتی فارسی
Persian pyramid histogram of positional signs (PPHOPS)	هیستوگرام هرمی نشانه‌های موقعیتی فارسی
Patch	وصله
Zero-shot learning	یادگیری شات صفر

فصل

١- مقدمه

۱-۱- تعریف مسئله

از گذشته‌های بسیار دور متن به‌عنوان یکی از منابع مهم اطلاعاتی مورد استفاده قرار گرفته است و امروزه بسیاری از این نوع داده در کتابخانه‌های دیجیتال، اسناد، تصاویر و ویدئوها وجود دارد. پیشرفت فناوری دیجیتال باعث شده است که اسناد و کتاب‌ها اسکن شده و به‌صورت تصویر ذخیره شوند. اسناد وقتی که به‌صورت تصویری درمی‌آیند، ذخیره‌سازی و انتقال ساده‌تری خواهند داشت. بنابراین نگهداری اسناد و کتاب‌ها به‌صورت اسناد تصویری مورد توجه بسیاری قرار گرفته است. جهت پردازش و استفاده از محتوای متنی اسناد تصویری، روش‌های تشخیص و تحلیل اسناد تصویری^۱ مورد توجه قرار گرفته و توسعه‌یافته‌اند. جستجو و بازیابی اطلاعات همواره به‌عنوان یکی از مهم‌ترین شاخه‌های تحلیل داده، مورد توجه محققین بوده است. از این رو جستجو در اسناد تصویری نیز اهمیت خاصی را در حوزه تحلیل اسناد دارد. این مسئله در حوزه بازیابی محتوایی^۲ تصاویر قرار می‌گیرد.

رهیافت اولیه جهت جستجو و بازیابی اطلاعات متنی از اسناد تصویری، استفاده از روش‌های تشخیص کاراکتر نوری^۳ جهت تبدیل اسناد تصویری به متن و سپس بکار بردن روش‌های بازیابی متن است [۱]. روش‌های تشخیص کاراکتر نوری فعلی در مورد فونت‌های ساده با پس‌زمینه تمیز عملکرد قابل قبولی دارد، اما کارایی این روش با کاهش کیفیت اسناد یا در مورد اسناد دست‌نویس به شدت کاهش می‌یابد [۲]. علاوه بر این، روش مذکور در مورد پایگاه داده‌ها با تعداد زیادی سند تصویری، بسیار زمان‌بر است [۳].

بنابراین، محققین به دنبال روش‌های بازیابی بوده‌اند که توانایی جستجوی اطلاعات متنی بدون نیاز آنالیز طرح‌بندی^۴ و تبدیل کل سند تصویری به متن را داشته باشد. از این رو روش‌هایی که قادر به بازیابی تصاویر حاوی یک واژه مشخص باشند یا بتوانند لغات را در یک تصویر تشخیص دهند مورد توجه قرار گرفته‌اند. ایده جستجوی کلمات^۵، که در برخی از مقالات جستجوی کلمات کلیدی^۶ نیز خوانده می‌شود، در زمینه پردازش گفتار ارائه شده است [۴]. سپس این روش برای اندیس‌گذاری اسناد چاپی و دست‌نویس به کار گرفته شده است. این رهیافت موقعیت کلمه انتخاب شده توسط کاربر را در یک سند تصویری مشخص می‌کند. به‌عبارت‌دیگر، سیستم جستجو کلمه، کلمه پرسشی موردنظر کاربر را در سند تصویری جستجو می‌کند. جستجوی کلمه، روشی برای اندیس‌گذاری، جستجو و بازیابی کلمات در اسناد تصویری است. در این مسئله هدف مشخص کردن موقعیت رخدادهای کلمه پرسشی موردنظر کاربر در اسناد

^۱ Document image Analysis and Recognition (DIAR)

^۲ Content-Based Image Retrieval (CBIR)

^۳ Optical Character Recognition (OCR)

^۴ Layout analysis

^۵ Word spotting

^۶ Keyword spotting

تصویری است. در این سیستم کاربر کلمه پرسشی^۱ را از سند انتخاب کرده و به سیستم ارائه می‌کند. وظیفه سیستم یافتن کلمات مشابه در سند مورد جستجو می‌باشد. یکی از مسائل مهمی که برای سیستم‌های تشخیص کاراکتر نوری چالش مهمی محسوب می‌شود، مشکلات تبدیل طرح‌بندی صفحات سند است. تصاویر، گرافیک و تشخیص فونت و اندازه آن نیز به‌عنوان چالشی برای تبدیل تصویر به متن مشابه وجود دارد. استفاده از این روش باعث می‌شود که طرح‌بندی صفحات حفظ شود و بنابراین چالش‌های ذکر شده وجود نداشته باشند [۵].

یکی از کاربردهای مطرح این روش، جستجوی کلمه در اسناد تاریخی تصویری می‌باشد؛ به طور مثال جستجوی شهر یا ناحیه خاصی در نقشه‌های اسکن شده قدیمی، یا جستجوی تصاویر با کلمه پرسشی "Figure" می‌تواند از کاربردهای مهم این روش ذکر شوند. علاوه بر این، جستجوی کلمه ابزار ارزشمندی برای کاربران جهت جستجو روی محتوای کتب قدیمی و دست‌نویس و بازیابی نامه‌های دست‌نویس خواهد بود.

۱-۲- اهداف و چالش‌ها

زبان فارسی یا پارسی، از دسته زبان‌های هندواروپایی محسوب می‌شود که در کشورهای ایران، افغانستان، تاجیکستان و ازبکستان به آن سخن گفته می‌شود. بیش از ۱۵۰ میلیون نفر از این زبان برای نوشتن استفاده می‌کنند. یکی از مهم‌ترین مشخصه‌های زبان فارسی ماهیت پیوستگی آن هم در متون تایپی و هم در متون دست‌نویس است، که باعث شده دقت تشخیص در متون تایپی فارسی کمتر از متون تایپی لاتین باشد. الفبای فارسی دارای ۳۲ حرف است که هرکدام از آن‌ها بسته به موقعیتش در کلمه ۲ تا ۴ شکل دارد. تمام متون فارسی از راست به چپ نوشته می‌شوند.

جدول ۱-۱ شکل‌های مختلف حروف الفبای زبان فارسی را نشان می‌داد که تفاوت ظاهری حروف بر اساس موقعیت مکانی در کلمه (ابتدا، وسط، انتها و منفرد) را نشان می‌دهد. مثلاً حرف -غ- که غین نامیده می‌شود دارای چهار شکل است. شکل آن در ابتدای کلمه -غ- و در وسط کلمه به‌صورت -گ- و در انتها به‌صورت -غ- می‌باشد. بعضی از حروف تنها دارای دو شکل متفاوت می‌باشند. مانند حرف -د- که دال نامیده می‌شود. در ابتدای کلمه به‌صورت -د- و در وسط و انتهای کلمه به‌صورت -د- می‌باشند.

¹ Query

جدول ۱-۱. شکل‌های مختلف حروف الفبای فارسی با توجه به محل قرار گرفتن در کلمه

مجزا	ابتدا	وسط	انتهای	
	آ	آ	آ	۱
	ا	ا	ا	۲
ب	ب	ب	ب	۳
پ	پ	پ	پ	۴
ت	ت	ت	ت	۵
ث	ث	ث	ث	۶
ج	ج	ج	ج	۷
چ	چ	چ	چ	۸
ح	ح	ح	ح	۹
خ	خ	خ	خ	۱۰
	د	د	د	۱۱
	ذ	ذ	ذ	۱۲
	ر	ر	ر	۱۳
	ز	ز	ز	۱۴
	ژ	ژ	ژ	۱۵
س	س	س	س	۱۶
ش	ش	ش	ش	۱۷
ص	ص	ص	ص	۱۸
ض	ض	ض	ض	۱۹
ط	ط	ط	ط	۲۰
ظ	ظ	ظ	ظ	۲۱
ع	ع	ع	ع	۲۲
غ	غ	غ	غ	۲۳
ف	ف	ف	ف	۲۴
ق	ق	ق	ق	۲۵
ک	ک	ک	ک	۲۶
گ	گ	گ	گ	۲۷
ل	ل	ل	ل	۲۸
م	م	م	م	۲۹
ن	ن	ن	ن	۳۰
ه	ه	ه	ه	۳۱
	و	و	و	۳۲
ی	ی	ی	ی	۳۳

تعدادی حروف در الفبای فارسی وجود دارند که از سمت چپ قابل اتصال نیست و اگر در ابتدا و وسط کلمه قرار بگیرند باعث جدا شدن کلمه و تشکیل زیر کلمات می‌شود. این مشخصه زبان فارسی باعث شده که فاز جداسازی در سیستم‌های مبتنی بر قطعه‌بندی با مشکل روبه‌رو شود، به این دلیل که فاصله‌گذاری مناسب بین کلمات و زیر کلمات در بیشتر مواقع رعایت نمی‌شود. این حروف شامل (ر، ز، ژ، د، ذ، و، ا) می‌باشد. مثلاً در کلمه "مدرسه" وجود دو حرف "د" و "ر" در کلمه باعث شده که کلمه "مدرسه" شامل سه زیر کلمه شود.

در زبان عربی و فارسی کلماتی وجود دارد که شکل یکسانی دارند اما تلفظ و معنی متفاوتی دارند، برای نشان دادن تلفظ و معنی موردنظر از یک سری علائم استفاده می‌شود. مثلاً کلمه فارسی "مرد" دارای دو تلفظ است که معنی‌های متفاوتی دارند: "مرد" به معنی از دنیا رفتن و "مرد" که به معنی آقا می‌باشد. البته در فارسی، در بیشتر مواقعی که در متن آورده می‌شود معنی آن با توجه به متن مشخص می‌شود. از دیدگاه تشخیص کلمه، وجود این علائم دقت تشخیص را کاهش می‌دهد.

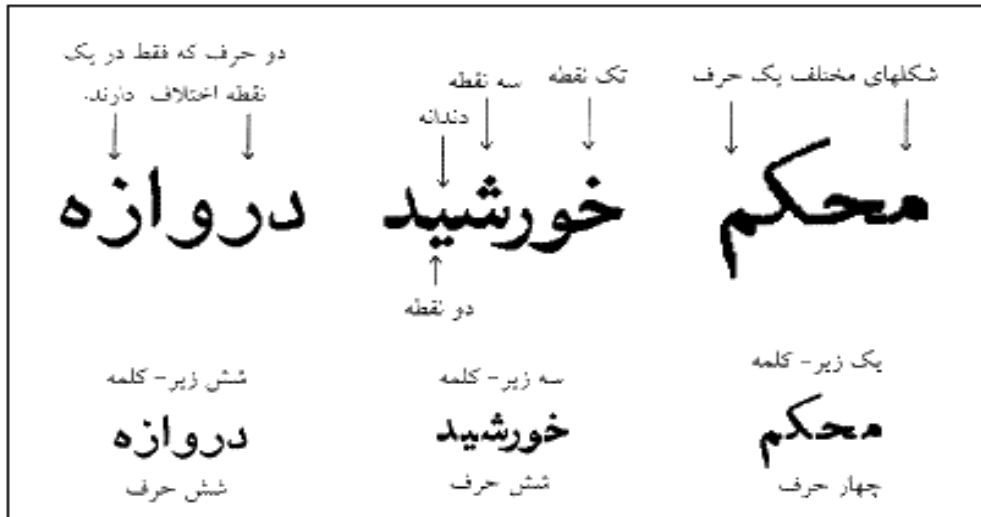
یکی دیگر از مشخصات زبان فارسی وجود نقطه در بعضی حروف است. ۱۸ حرف از ۳۲ حرف زبان فارسی شامل یک یا دو یا سه نقطه در بالا یا پایین خط اصلی می‌باشند:

- ده حرف فارسی دارای یک نقطه می‌باشند (ب، ج، خ، ذ، ز، ض، ظ، غ، ف، ن)
- سه حرف شامل دو نقطه می‌باشند (پ، ق، ت)
- پنج حرف شامل سه نقطه می‌باشند (پ، چ، ژ، ش، ث)

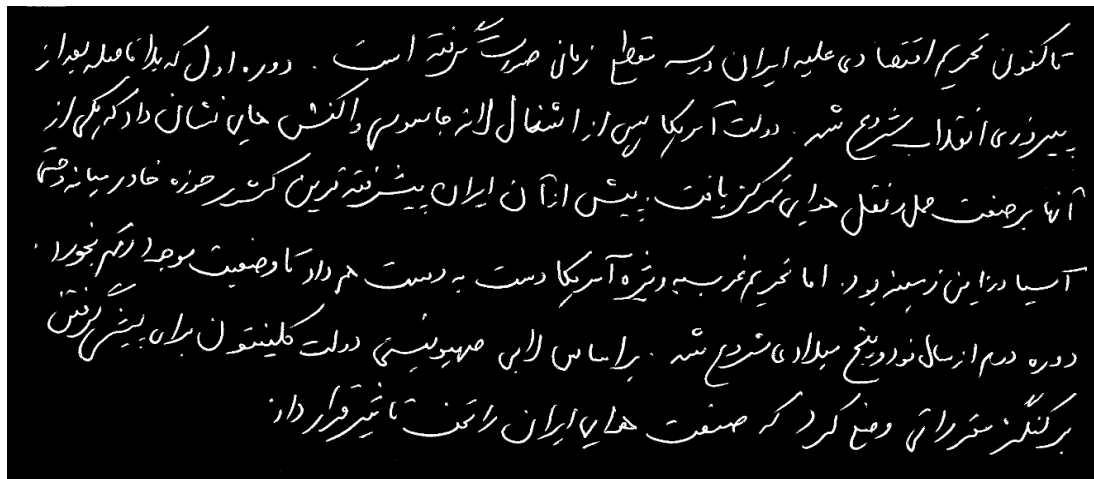
علاوه بر این در فارسی برخی از حروف دارای مکمل‌هایی به غیر از نقطه مانند سرکش حروف (ک و گ) یا آکلاد حرف (آ) می‌باشند.

وجود نقاط و علائم در کلمات، خصوصاً در اسناد دست‌نوشته، باعث می‌شود عمل تشخیص مشکل‌تر شود. دلیل آن را می‌توان این‌گونه بیان کرد که در متون دست‌نوشته اندازه و نحوه نوشتن نقاط، مکمل‌ها و علائم بسیار متغیر است و ممکن است بعضی از آن‌ها آن‌قدر کوچک باشد که در فاز پیش‌پردازش حذف شوند یا آن‌قدر بزرگ باشند که با بعضی حروف اشتباه گرفته شوند. مثال‌هایی که در شکل ۱-۱ و ۱-۲ آورده شده‌اند، برخی ویژگی‌های زبان فارسی را به ترتیب در سطح کلمه و سند نشان می‌دهند.

بیشتر تحقیقات انجام شده در این زمینه بر روی متون زبان لاتین ارائه شده است و تعدادی از آن‌ها در مورد زبان‌های دیگری مثل چینی، هندی و کره‌ای می‌باشند. تحقیقات انجام شده در زبان فارسی و عربی بسیار کم هستند که بیشتر آن‌ها نیز بر روی اسناد چاپی انجام شده است. به همین دلیل با توجه به اهمیت و کاربرد مسئله و اینکه طبق بررسی‌های صورت گرفته کار قابل توجهی در زبان فارسی انجام نشده است، نیاز به انجام تحقیق در این زمینه احساس می‌شود.



شکل ۱-۱. برخی ویژگی‌های زبان فارسی [۶]



شکل ۲-۱. نمونه‌ای از اسناد دست‌نویس فارسی

از آنجایی که نوشتارهای زبان‌های مختلف خواص متفاوتی دارند، در نظر گرفتن ساختار نوشتاری زبان در روش‌های جستجوی کلمه به‌عنوان یک نکته کلیدی می‌تواند در نظر گرفته شود. ساختار نوشتاری زبان فارسی پیچیدگی‌های بیشتری نسبت به انگلیسی در بسیاری از کارهای آنالیز اسناد مانند جستجوی کلمه ایجاد می‌کند. برخی از مهم‌ترین این ویژگی‌ها را می‌توان به‌صورت زیر ذکر کرد:

۱. ذات نوشتار فارسی پیوسته است.
۲. کلمه‌ها از مجموعه‌ای از یک یا چند زیر کلمه^۱ تشکیل می‌شوند.
۳. هر زیر کلمه می‌تواند یک حرف یا مجموعه‌ای از حروف متصل باشد.
۴. هر زیر کلمه می‌تواند از یک یا چندین جزء متصل ساخته شود.

^۱ Subword در برخی از مقالات به جای "زیر کلمه" از "بخشی از کلمه فارسی" (Part of Persian Word (PPW) استفاده شده است.

۵. تعداد زیادی از حروف از چندین جزء متصل، شامل بدنه حرف و نقاط یا سرکش‌ها، ساخته شده‌اند.

۶. شکل اکثر حروف بر اساس موقعیت آن در کلمه تغییر می‌کند.

۷. حدود کلمات در نوشتار فارسی، خصوصاً در مورد اسناد دست‌نویس، به راحتی قابل تشخیص نیست.

بر این اساس چالش‌های زیاد و متفاوتی در حوزه آنالیز اسناد در زبان فارسی وجود دارد که آن را متمایز با زبان‌های دیگری چون انگلیسی می‌کند. بنابراین با توجه به اینکه بر اساس تحقیقات صورت گرفته، پژوهش چشمگیری در زمینه جستجوی کلمه در اسناد دست‌نویس فارسی انجام نشده است، نیاز به انجام تحقیق در این زمینه احساس می‌شود.

به‌طور کلی سیستم‌هایی که روی دست‌نوشته کار می‌کنند حساسیت زیادی روی تنوع شیوه‌های نگارش (دستخط) دارند. از این رو بسیاری از کارهای اولیه انجام شده در زمینه جستجوی کلمه تک نویسنده است. روش‌هایی که اخیراً معرفی شده‌اند، سعی در فائق آمدن بر مسئله پیچیدگی و تنوع شیوه‌های نگارش داشته‌اند.

روش‌هایی مانند مدل مخفی مارکوف^۱ و پیچش زمان پویا^۲ در جستجوی کلمه در اسناد تصویری بسیار مورد استفاده قرار گرفته‌اند. اما این روش‌ها در سیستم‌هایی استفاده شده است که کلمه پرسشی به صورت تصویری به سیستم داده می‌شود. یکی از چالش‌های دیگر مربوط به کلمه پرسشی تصویری است. در کاربرد واقعی یافتن و ارائه کلمه پرسشی تصویری مورد نظر کاربر کاری مشکل و زمان‌بر خواهد بود. اما اگر کلمه پرسشی به صورت متنی به سیستم جستجوی کلمه ارائه شود، مشکلات مربوط به آماده‌سازی کلمه پرسشی وجود نخواهد داشت و ارائه کلمه پرسشی به سیستم ساده خواهد شد، اما برای این منظور سیستم پیچیده‌تری باید طراحی گردد.

با در نظر گرفتن تمامی موارد ذکر شده در بالا، در این رساله در گام اول یک سیستم مبتنی بر یادگیری استفاده شده است که بر اساس رده‌بندی مبتنی بر خصیصه^۳، تعبیه برچسب^۴ و رگرسیون زیرفضای مشترک^۵ کار می‌کند [۷]. در این رساله خصیصه‌هایی^۶ مبتنی بر شیوه نگارش در زبان فارسی پیشنهاد شده است که بهتر کار می‌کند و امکان به اشتراک‌گذاری اطلاعات بین کلاسی را افزایش می‌دهد. علاوه

^۱ Hidden Markov Model (HMM)

^۲ Dynamic Time Wrapping (DTW)

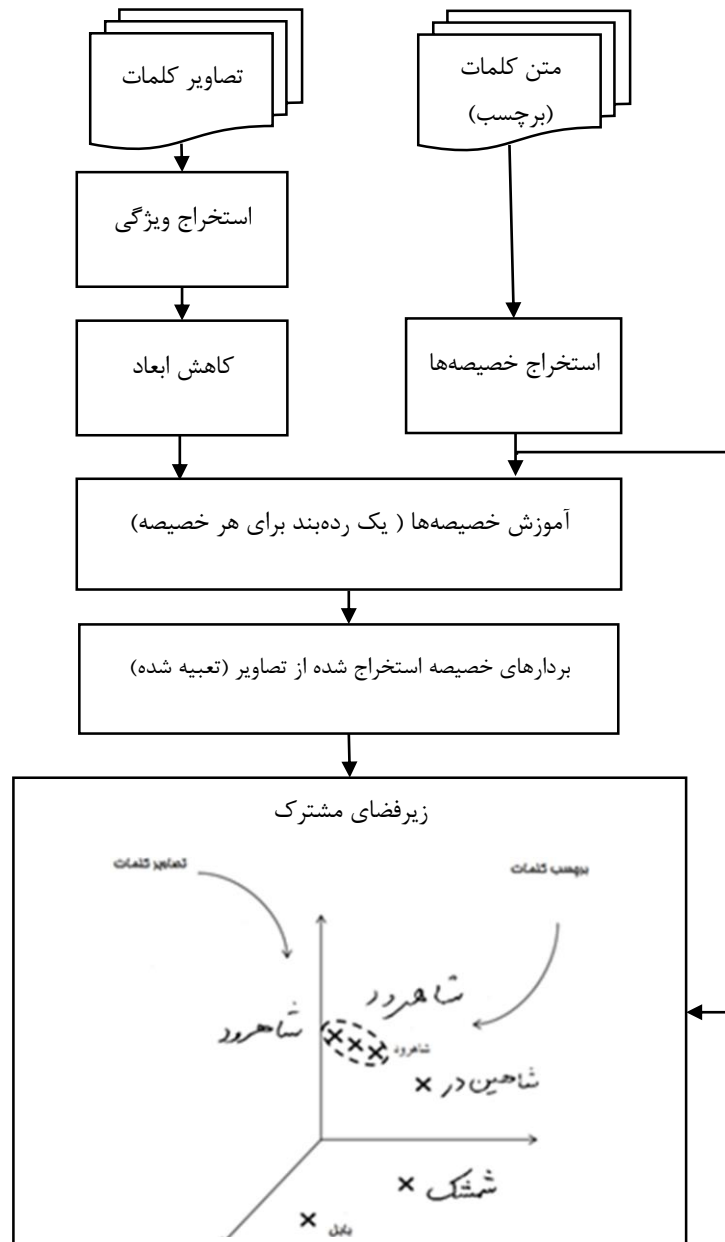
^۳ Attribute-based classification

^۴ Label embedding

^۵ Common Subspace Regression (CSR)

^۶ Attributes

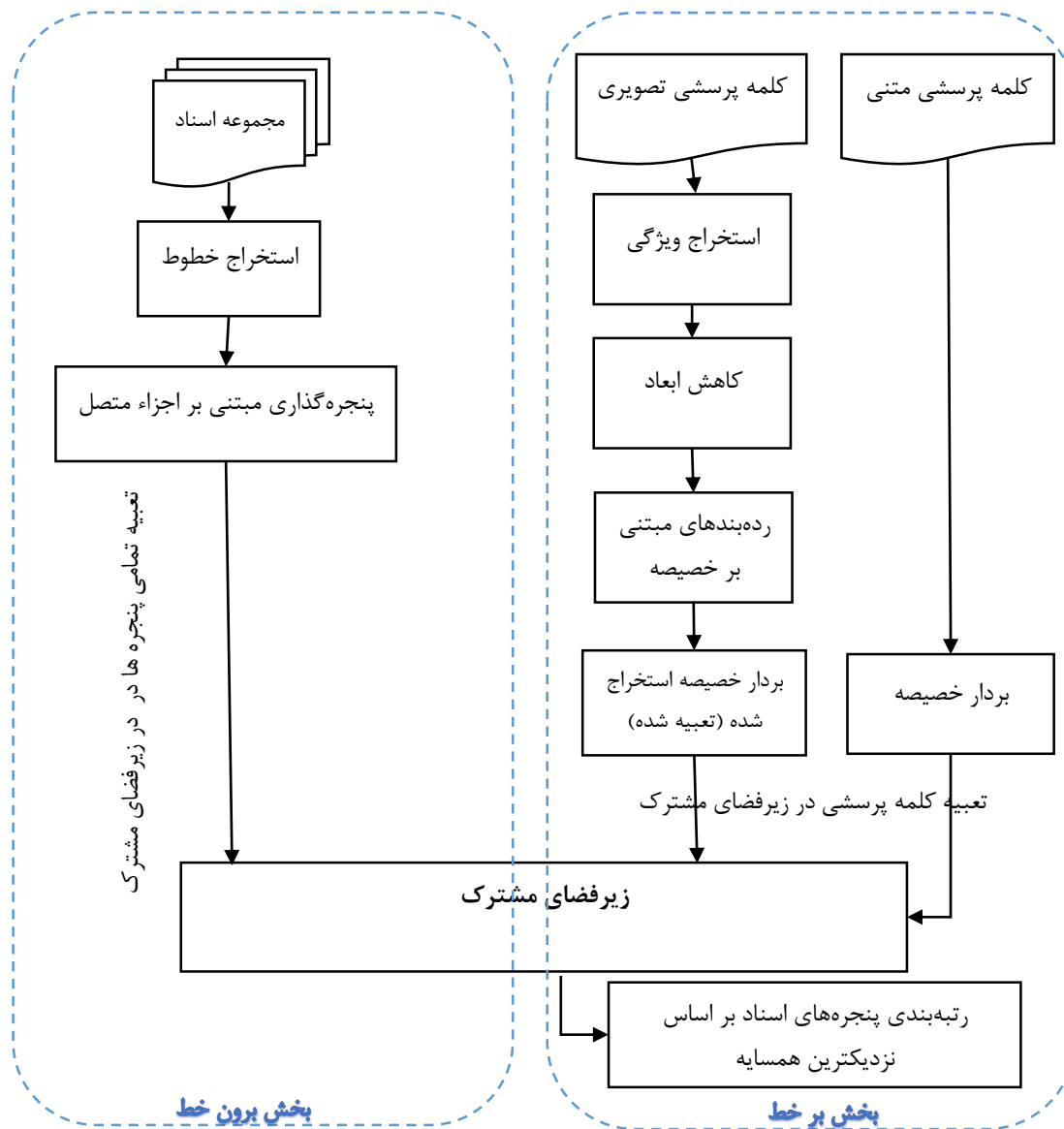
بر این خصیصه‌های پیشنهادی امکان جستجوی زیرکلمات و اجزاء متصل مربوط به کلمات را نیز فراهم می‌کند. خصیصه‌های پیشنهادی و زیرفضای مشترک آموزش داده شده و برای جستجوی کلمات مبتنی بر قطعه‌بندی مورد استفاده قرار گرفته است. شکل ۱-۲ ساختار کلی روش مذکور را در فاز آموزش نشان می‌دهد.



شکل ۱-۳. فلوجارت کلی آموزش روش پیشنهادی مبتنی بر قطعه‌بندی

سپس برای فائق آمدن بر خطای زیاد مربوط به قطعه‌بندی کلمات، که در فارسی به دلیل مشخص نبودن محدوده کلمات در کارهای آنالیز اسناد متنی مشکلات زیادی را ایجاد می‌کند، یک روش مبتنی بر

قطعه‌بندی در سطح خط پیشنهاد شده است که از اجزاء متصل موجود در هر خط به‌عنوان معیاری جهت پنجره‌گذاری استفاده شده است (شکل ۴-۱).

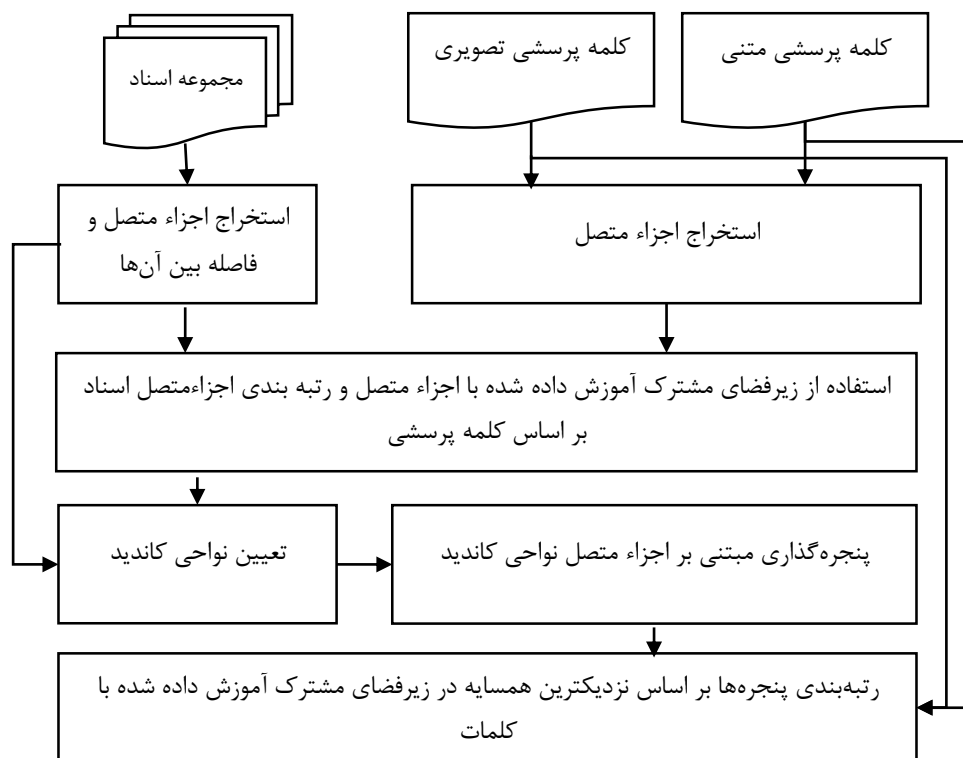


شکل ۴-۱. فلوجارت کلی روش پیشنهادی مبتنی بر قطعه‌بندی در سطح خطوط

در گام بعدی نیز یک چهارچوب کاری سلسله‌مراتبی با استفاده از قابلیت‌های خصیصه‌های پیشنهادی ارائه شده است که بر اساس جستجوی اجزاء متصل سازنده کلمه پرسشی ورودی نواحی کاندید را تعیین می‌کند. در این روش نواحی کاندید با استفاده از سیستم جستجوی کلمات مبتنی بر قطعه‌بندی رتبه‌بندی می‌شوند. از آنجایی که در روش مبتنی بر خط تعداد مثبت اشتباه^۱ زیاد است و کاهش تعداد آن به‌طور معناداری باعث افزایش کارایی سیستم جستجوی کلمه می‌گردد، این روش جدید که تعداد

¹ False positive

کاندیدها را کاهش می‌دهد پیشنهاد شده است. علاوه بر این، روش مذکور به قطعه‌بندی سند نیز نیاز ندارد (شکل ۱-۵).



شکل ۱-۵. فلوجارت کلی روش پیشنهادی مستقل از قطعه‌بندی

۱-۳- مجموعه داده‌ها

بررسی روش پیشنهادی بر اساس نتایج حاصل شده بر روی مجموعه داده‌های استاندارد همواره به‌عنوان یکی از روش‌های ارزیابی مهم به حساب می‌آید. در روش‌هایی که در این رساله پیشنهاد شده است از ۳ مجموعه داده استاندارد زبان فارسی استفاده شده است که در زیر معرفی می‌شود.

مجموعه داده ایران‌شهر: اسامی ۵۰۳ شهر ایران در یک مجموعه‌ای بیش از ۱۷۰۰۰ نمونه‌ای را فراهم آورده است. در هر کلاس از کلمات تصویری شامل نام شهرهای ایران حداقل ۲۵ نمونه وجود دارد. شکل ۱-۶ نمونه‌هایی از کلمات موجود در این مجموعه داده را نشان می‌دهد.

ارومیا	بسطام	بهبهان	آباداخ	خاش
اروسه	سالم	بهبهان	آبادان	خاش
ترت حیدریه	شهرضا	زاهد	قصر قند	نشارود
ترت حیدریه	شهرضا	زاهد	قصر قند	نشارود
گکتاب	گنرود	نزد	فرشهر	هفت تپه
گکتاب	گنرود	نزد	فرشهر	هفت تپه

شکل ۱-۶. چند نمونه از تصویر کلمات مجموعه داده ایران شهر

مجموعه داده فارسا [۸]: مجموعه داده فارسا که در دانشگاه صنعتی شاهرود جمع آوری شده است مجموعه‌ای از ۳۰۰ کلمه پر کاربرد فارسی است که توسط ۱۰۰ نویسنده مختلف به صورت دست‌نویس نوشته شده است که مجموعه‌ای مشتمل بر ۳۰۰۰۰ تصویر کلمه دست‌نوشته می‌باشد. نمونه‌هایی از برخی از کلمات موجود در این مجموعه داده در شکل ۱-۷ نمایش داده شده است.

مردم	تولید	کم	منطقه	خرد
مردم	تولید	کسر	منطقه	خرد
اسان	تحقیق	مادر	صد	با
اسان	تحقیق	مادر	صد	با

شکل ۱-۷. چند نمونه از تصویر کلمات مجموعه داده فارسا

مجموعه داده اسناد دست‌نویس فارسی (FHT) [۹]: این پایگاه داده شامل ۱۰۰۰ سند می‌باشد که دارای ۴۰ دسته متن متفاوت است که توسط ۲۵ نویسنده با سطح سواد متفاوت نوشته شده‌اند. در واقع از هر فرد تقاضا شده است ۴ متن متفاوت را بنویسد که در مجموع این پایگاه داده، ۲۵۰ نویسنده دارد. هر متن به‌طور میانگین دارای ۶/۴۵ خط، ۱۰۶/۶ کلمه، ۲۳۰/۱۷۵ زیر کلمه، ۱۴۰۶ حرف و ۱۳۲/۱

نقطه است. یکی از مزیت‌های این پایگاه داده وجود خطوط صحیح^۱ برای هر خط می‌باشد. این اسناد هم شامل دستخط خوانا و خوش خط و هم غیر خوانا و بد خط است. شکل ۱-۸ دو متن نمونه از این مجموعه داده را نمایش می‌دهد.

کینه از این معانیها پیروزها کار روزی بود. کارش نوعی سرای بران زبان بود. او قسم
 می‌خواند. ضلع طرف و فضول بود اتفاقاً شیرین زبان و حال هم بود. ما بچم ها اورا
 خیل دوست داشتیم. و همگی می‌آمدند گفت ما به راه بود بچم ها قصه می‌گفت.
 ما هم هم تصنیف می‌خواند و همه در خانه گفت می‌زدند. فدای ما می‌ان عزیزان بود
 یک عینک هم داشت. از آن عینک ما سه بار این شخص قدیم. یک روز که
 پیروزان نبود رفتیم سرانجام. اول کتاب ما را به هم رسیدیم. به از
 روی سینه و شلوار عینک مرصوف را از هیبه در آوردیم. آن را به چشم
 گذاشتیم که بروم دید این رحمت صفتک سر به سر خواصم به آرام
 سرز فرانسوا بنامکم. بران ما حفظه الله عصب و عقیسی بود.

کار رویت نظیر یک دعاوت است. حکایتی در بیار تدریز شک من رود در ره بیماری است
 من پرسید نظیر یک صید اصطلاح روشن را که من بر دو سینه من وید داروها تو اینجا است جمله
 وقتها ما چیزی از آن اصطلاحات سر در نمی‌آوریم اما علاوه بر آن از تصرف در واحد
 چند نفع هم بر سبب سیدانی داروهای که تصرف می‌کنیم چه بود و چه می‌باشیم. اما در این
 با روش تفکری این گونه نیست. چون داروی که روش تفکر ما هم در همان اصطلاحات است و
 در اینجا هم اصطلاحات است که یک کارایی دارو می‌انجامد. نکته‌ای که روش تفکر ایرانی
 صندان به آن توجه ندارند.

شکل ۱-۸. چند نمونه از تصویر کلمات مجموعه داده FHT

¹ Ground Truth

۱-۴- ساختار پایان نامه

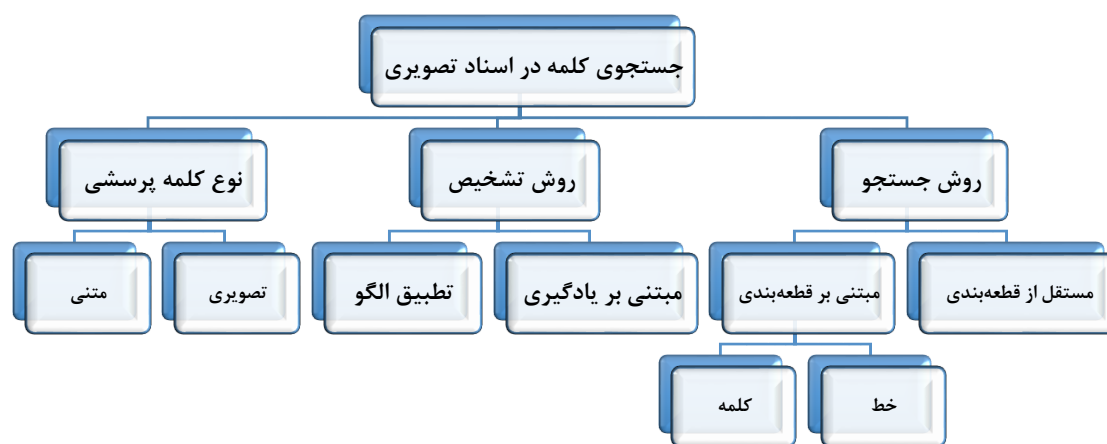
در فصل بعدی کارهای انجام شده در زمینه جستجوی کلمات مورد بررسی قرار گرفته است. در این فصل انواع دسته‌بندی‌های ارائه شده از منظرهای مختلف آورده شده است. فصل سوم روش پیشنهادی جهت جستجوی کلمات فارسی به صورت قطعه‌بندی شده، در سطح کلمه ارائه شده است. این فصل شامل پیشنهاد مدل اصلی استفاده شده در این رساله و خصیصه‌های پیشنهادی می‌باشد. ضمناً در انتهای فصل ارزیابی و نتایج تجربی آورده شده است. از آنجایی که از این روش جهت بازشناسی دست‌نوشته‌های فارسی نیز می‌توان بهره برد، مقایسه‌ای نیز با بهترین کارهای انجام شده در این حوزه آورده شده است. جستجوی کلمه در روشی مبتنی بر قطعه‌بندی در سطح خطوط در فصل چهارم ارائه شده است. این بخش شامل روش پیشنهادی جهت استخراج خطوط و پنجره‌گذاری مبتنی بر اجزاء متصل می‌باشد و مدل یادگیری آن مدل پایه‌ی پیشنهادی است. در فصل پنجم یک روش مستقل از قطعه‌بندی پیشنهاد می‌شود. این روش چهارچوب کاری متفاوتی را بر اساس سیستم پایه معرفی می‌کند که با تعیین نواحی کاندید فضای جستجو را به طور چشم‌گیری کاهش می‌دهد. فصول چهارم و پنجم نیز شامل ارزیابی و نتایج تجربی سیستم معرفی شده در آن فصل هستند. نهایتاً فصل ششم به نتیجه‌گیری و ارائه پیشنهادات اختصاص داده شده است.

فصل

۲- مروری بر کارهای پیشین

۲-۱- مقدمه

روش‌های جستجوی کلمه یک الگوی ساختاری یکسان را دنبال می‌کنند؛ در این الگو ابتدا اسناد در یک فاز برون خط اندیس‌گذاری و ذخیره می‌شوند، سپس کاربر یک پرسش را مطرح می‌کند و سیستم شباهت آن را با کلمات اسناد ارزیابی می‌کند. به‌طور کلی این سیستم‌ها، کلمه موردنظر کاربر را در سند تصویری جستجو می‌کنند. با توجه به تنوع کارهای انجام شده روی اسناد، روش‌های موجود در جستجو کلمه را می‌توان از چند دیدگاه تقسیم‌بندی نمود (شکل ۲-۱).



شکل ۲-۱: دسته‌بندی روش‌های موجود جستجوی کلمه

در این دسته‌بندی روش‌های موجود از سه منظر متفاوت دسته‌بندی شده‌اند: اولین مورد روش‌ها را بر اساس نوع کلمه پرسشی دسته‌بندی می‌کند. روش‌های موجود کلمه پرسشی را به‌صورت تصویری^۱، متنی^۲ یا به هر دو صورت می‌پذیرند. در اکثر روش‌هایی که کلمه پرسشی صرفاً تصویری وارد می‌شود، کلمه موردنظر از سند بریده می‌شود و به‌عنوان کلمه پرسشی مورد استفاده قرار می‌گیرد، اما روش‌های جدید از رشته متنی ورودی به‌عنوان کلمه پرسشی بهره می‌برند. دسته‌بندی دیگر بر اساس روش جستجو می‌باشد. روش‌هایی وجود دارند که ابتدا سند را به خطوط و سپس به کلمات سازنده قطعه‌بندی کرده و سپس عمل جستجو را انجام می‌دهند، این دسته روش‌ها مبتنی بر قطعه‌بندی در سطح کلمه نامیده می‌شوند. برخی دیگر قطعه‌بندی را فقط در سطح خطوط انجام می‌دهند. برای پیشگیری از وقوع خطای قطعه‌بندی، روش‌های مستقل از قطعه‌بندی نیز ارائه شده است. آخرین دسته‌بندی بر اساس روش تشخیص صورت گرفته است. عمل تشخیص می‌تواند مبتنی بر یادگیری یا بر اساس روش‌های تطبیق الگو باشد. روش‌های تطبیق الگو معمولاً یا عمل تطبیق را روی داده‌های خام و در سطح پیکسل و یا بر اساس ویژگی‌های استخراج‌شده از تصاویر انجام می‌دهند.

¹ Query By Example (QBE)

² Query By String (QBS)

۲-۲- دستهبندی بر اساس روش جستجو

در زمینه جستجوی کلمات، روش‌های موجود ممکن است سند مورد جستجو را در سطح خط یا کلمه قطعه‌بندی کنند یا بدون قطعه‌بندی کل سند را مورد جستجو قرار دهند. در صورتی که نیاز به قطعه‌بندی کلمه وجود داشته باشد، ابتدا باید خطوط و سپس کلمات هر خط از سند جدا شود و در نوع دیگر فقط خطوط استخراج می‌شوند. به‌طور کلی قطعه‌بندی از دسته مراحل پیش‌پردازشی است که میزان خطای بالایی دارد و این خطا تأثیر زیادی بر روی عملکرد مراحل بعدی خواهد داشت. در آخرین مورد هیچ عمل قطعه‌بندی انجام نمی‌شود و از بروز این خطا پیشگیری می‌کند، اما چون نیاز به بکار بردن روش‌هایی چون پنجره لغزان وجود دارد پردازش پرهزینه خواهد بود.

۲-۱- روش‌های مبتنی بر قطعه‌بندی

یکی از پیش‌پردازش‌هایی که در بسیاری از روش‌های پردازش اسناد تصویری انجام می‌شود عمل قطعه‌بندی است. این عمل می‌تواند در سطوح مختلف از قبیل خط، کلمه یا حرف انجام شود. نکته قابل‌ذکر در مورد قطعه‌بندی این است که، این مسئله تا حدود زیادی وابسته به زبان است. در برخی از زبان‌ها با توجه به حروف و ساختار نگارش آن‌ها قطعه‌بندی ممکن است ساده‌تر باشد، از این‌رو در بسیاری از کارهایی که روی زبان‌هایی مانند چینی و برخی کارهای انگلیسی صورت گرفته است، قطعه‌بندی در سطح حروف دیده می‌شود. در این روش‌های پایه اصلی روش مبتنی بر کدگذاری شکل کاراکتر^۱ است. به این صورت که پس از استخراج کاراکترها، با مجموعه کدهای از پیش تعریف‌شده آن‌ها را تفسیر می‌کنند. محدودیت اصلی این روش‌ها حساسیت به خطای قطعه‌بندی است. به علت مشکلات دستگاه‌های تصویربرداری و گسسته‌سازی و یا در اثر کیفیت پایین اسناد و یا در اسناد دست‌نویس، امکان به هم چسبیدن یا شکسته شدن حروف وجود دارد. این به هم چسبیدگی یا شکستگی باعث می‌شود که کدگذاری اشتباه انجام شود. به دلیل خطای بالای قطعه‌بندی در سطح حروف در زبان‌هایی که ساختار نگارشی پیوسته‌ای دارند، مانند زبان فارسی و عربی، از این دسته روش‌ها استفاده نشده است.

کدگذاری شکل کلمه^۲، به جای قطعه‌بندی کلمه به حروف و سپس کد کردن حروف، به‌صورت کل‌نگر کدگذاری را روی کلمه اجرا می‌کند. در این دسته روش‌ها تصویر کلمه به‌عنوان کوچک‌ترین جزء در نظر گرفته می‌شود. کدگذاری‌هایی که معمولاً در سطح کلمه در نظر گرفته می‌شود از ساختار و شکل حروف

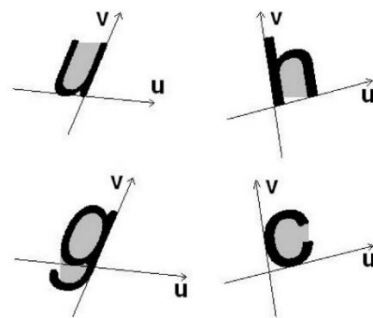
^۱ Character Shape Coding(CSC)

^۲ Word Shape Coding (WSC)

استفاده می‌کنند که معمولاً ویژگی‌هایی مانند حفره‌ها، پایین‌رونده و بالارونده‌ها، نقاط و سایر ویژگی‌های ساختاری وابسته به شکل نگارش زبان بهره می‌برند. معمولاً این روش‌ها از تطبیق الگو یا معیارهای شباهت جهت سنجش میزان شباهت کلمه پرسشی تصویری و کلمات کاندید در سند استفاده می‌کنند.

روش کدگذاری شکل کاراکتر در تحقیق [۱۰] استفاده شده تا موقعیت کلمات در تصاویری که با استفاده از دوربین دیجیتال اخذ شده است را مشخص نماید. در مرحله پیش‌پردازش استخراج خطوط و کلمات با استفاده از آنالیز اجزاء متصل انجام شده است. حروف l و 1 به‌عنوان حروف مبنا در نظر گرفته شده‌اند که جهت و میزان کجی خطوط با استفاده از آن‌ها سنجیده می‌شود و با استفاده از آن خط مبنا شناسایی می‌گردد. ویژگی‌های ساختاری حفره، مخازن آب‌ها^۱ و بالارونده و پایین‌رونده برای کد کردن کلمات استفاده می‌شوند. شکل ۲-۲ نمونه‌هایی از وقوع این ویژگی‌ها را در برخی از حروف انگلیسی نشان می‌دهد. تعداد حفره‌های درون کاراکتر و محل قرار گرفتن مخازن آب در چهار جهت چپ، راست، پایین و بالا ۵ عدد اول دنباله ۶ رقمی کد کلمه را می‌سازند. بخش آخر از کد، بالارونده و پایین‌رونده بودن را با اختصاص یکی از ارقام ۱، ۰ و ۱- مشخص می‌کند. جدول ۱-۲ کدهای حروف مختلف انگلیسی را با استفاده از این روش کدگذاری ۶ رقمی نشان می‌دهد.

برای تعیین میزان شباهت، فاصله همینگ^۲ کد کلمه پرسشی با کد کلمات با طول یکسان در سند محاسبه شده است. متوسط فراخوانی و دقت در ارزیابی تجربی به ترتیب ۷۵/۷٪ و ۶۹/۰۸٪ گزارش شده است.



شکل ۲-۲: نمایش حفره و مخازن آب در برخی از حروف نمونه انگلیسی [۱۰]

^۱ Water reservoirs

^۲ Hamming distance

جدول ۲-۱: کدگذاری شکل حروف انگلیسی ارائه شده در [۱۰]

Characters	Codes	Characters	Codes
a	110010	bd	100001
c	001000	e	101100
ft	001001	g	11010-1
h	000011	ir	000000
j	00000-1	k	001111
l	000001	m	000020
n	000010	o	100000
pq	10000-1	s	011110
uv	000100	w	000210
x	011110	y	01010-1
z	011000		

روش کدگذاری بسیار وابسته به روش نگارش در زبان است. در زبان فارسی و عربی که ذات ساختار نگارشی پیوسته است استفاده از کدگذاری شکل حروف جایگاهی ندارد اما به جای حروف می‌توان از کدگذاری شکل کلمات بهره برد. این عمل در تحقیق [۱۱]، با استفاده از ویژگی‌های ساختاری روی زبان عربی انجام شده است. جدول ۲-۲ کدهای استفاده شده برای حروف عربی سازنده کلمات در این تحقیق را نشان می‌دهد.

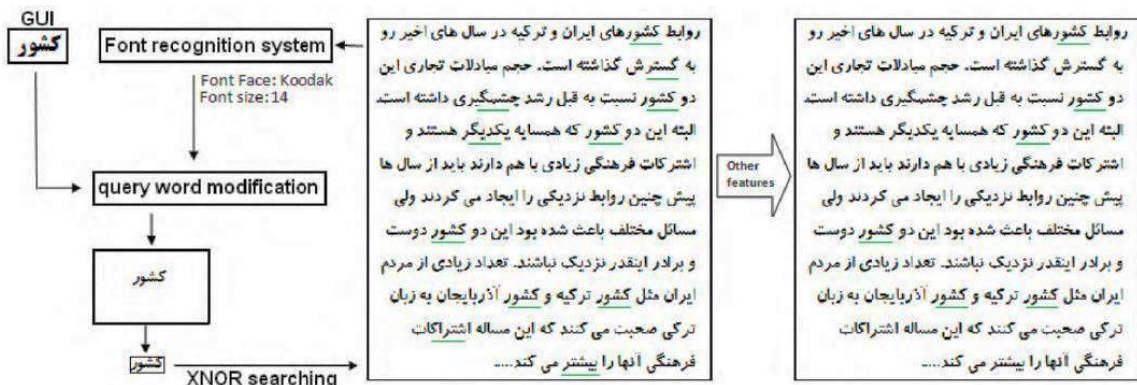
کدگذاری شکل کلمات در مورد اسناد چاپی فارسی نیز مورد استفاده قرار گرفته شده است [۱۲-۱۴]. عملکرد این روش بر اساس تشخیص فونت سند مورد جستجو و سنتز کلمه پرسشی مطابق با فونت و ساینز تشخیص داده شده استوار است. بعد از ایجاد کلمه پرسشی مطابق با سند مورد جستجو، خطوط سند با استفاده از افکنش استخراج می‌شوند و جستجو روی خطوط با استفاده از پنجره‌گذاری انجام می‌گیرد. جهت اندازه‌گیری میزان مشابهت نیز از XNOR استفاده شده است. این مرحله در واقع کاندیدها را استخراج می‌کند و جهت بازیابی دقیق به مرحله بعدی می‌فرستد. در مرحله نهایی از ویژگی‌های ساختاری مانند تعداد حفره‌ها، نقاط و بالارونده‌ها و پایین‌رونده‌ها استفاده شده است. شکل ۲-۳ و ۲-۴ ساختار کلی روش و نمونه‌هایی از ویژگی‌های ساختاری استفاده شده در مورد چند کلمه نمونه را نشان می‌دهد.

جدول ۲-۲: کدگذاری به کاررفته در حروف عربی ارائه شده در [۱۱]

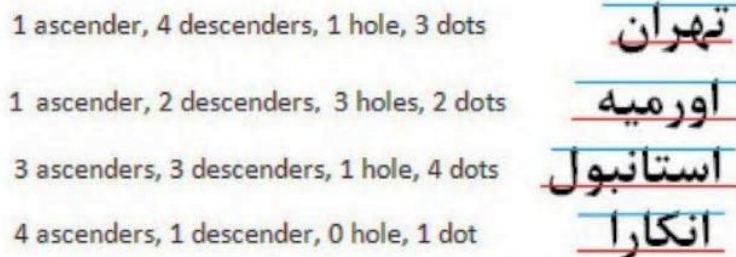
<i>Character</i>	<i>Code</i>	<i>Designation</i>
ا - ل - ک - ج	h	Ascender
!	hq	Ascender +Down Diacritic
آ - ل	ph	Up Diacritic +Ascender
ل	hj	Ascender +Descender
ط	bh	Loop +Ascender
ظ	bph	Loop +Up Diacritic +Ascender
لا	hbh	Ascender +Loop+ Ascender
ك	hp	Ascender +Up Diacritic
ي	jq	Descender +Down Diacritic
غ - خ - ذ - ت - ث - د	p	Up Diacritic
غ	jp	Descender +Up Diacritic
ث - ت	pp	Up Daicritic +Up Diacritic
ن - ز - خ - ی	jp	Descender +Up Diacritic
س	ppj	Up Diacritic + Up Diacritic + Descender
ض	bpj	Loop + Up Diacritic + Descender
ض - ف - ه - غ - ق	bp	Loop +Up Diacritic
ق	pbj	Up Diacritic + Loop + Descender
ی - ج - ب	q	Down Diacritic
ح - ع - س - ر - ی	j	Descender
ج	jq	Descender + Down Diacritic
ع - م - ص - ه - ه	b	Loop
ه	bb	Loop +Loop
ح - و - ص - م	bj	Loop +Descender
لا	hh	Ascender +Ascender
ة	pb	Up Diacritic +Loop
غ	pbj	Up Diacritic + Loop + Descender
و	bjp	Loop + Descender + Up Diacritic
لا	hbhp	Ascender + Loop + Ascender + Up Diacritic
لا	hbqh	Ascender+Loop+Down Diacritic+ Ascender

جهت ارزیابی روش مذکور، ۶۶۴ سند که بخشی از آن به صورت اسکن شده و بخش بزرگتری از آن به صورت کامپیوتری تهیه شده بود استفاده گردید. دقت در مورد اسنادی که به صورت کامپیوتری تهیه شده بودند به دلیل عدم وجود مشکلات سخت‌افزاری و نویز در زمان تبدیل سند به تصویر بالاتر بوده و

برابر ۹۷/۵٪ است در حالی که برای اسناد اسکن شده دقتی برابر ۸۷/۶٪ گزارش شده است. همین مسئله در مورد معیار فراخوانی نیز صادق است و به ترتیب برابر ۹۲/۱٪ و ۷۹/۳٪ می باشد.



شکل ۲-۳: ساختار کلی روش ارائه شده در [۱۲]



شکل ۲-۴: نمونه‌ای از ویژگی‌های ساختاری چند کلمه فارسی [۱۲]

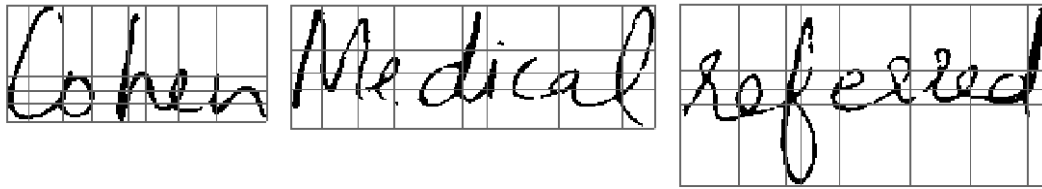
در [۱۵] و [۱۶] از ویژگی‌های مبتنی بر پروفایل دودویی با نام "گرادیانی، ساختاری و خمیدگی" استفاده شده است. برای استخراج این بردار ویژگی، هر تصویر به پنجره‌های ۴ در ۸ افقی تقسیم می‌شود (شکل ۲-۵). بردار ویژگی نهایی شامل ۱۹۲ ویژگی گرادیانی، ۱۹۲ ویژگی ساختاری و ۱۲۸ ویژگی تقعر برای هر کلمه است. در این روش از معیار همبستگی^۲ برای محاسبه میزان شباهت بین دو بردار باینری استفاده شده است.

$$D^b(X, Y) = \frac{1}{2} - \frac{S_{11}S_{00} - S_{10}S_{01}}{2((S_{10} + S_{11})(S_{01} + S_{00})(S_{11} + S_{01})(S_{00} + S_{10}))^{1/2}} \quad (1-2)$$

$D^b(X, Y)$ معیار عدم شباهت متناظر با معیار همبستگی است. در این رابطه، S_{ij} تعداد تطبیق i در الگوی اول با j الگوی دوم در موقعیت‌های متناظر است.

¹ Gradient, Structural and Convexity (GSC)

² Correlation measure



شکل ۲-۵: نمونه‌هایی از قطعه‌بندی افقی کلمات در [۱۵]

در [۱۷] از توصیفگرهای مبتنی بر کانتور برای بازنمایی کلمات و از تطبیق کانتور کشسان استفاده می‌شود، درحالی‌که استفاده از ویژگی‌های گابور و فاصله اقلیدسی در [۱۸] بررسی شده است. ویژگی‌های کل‌نگر کلمه با یک مدل آماری حاشیه‌نویسی^۱ در [۱۹] پیشنهاد داده شده است که می‌تواند کلمات دلخواه را جستجو کند.

برخی از رهیافت‌های محلی رایج، یک کلمه را به‌صورت توالی از ویژگی‌ها که با استفاده از یک پنجره لغزان روی کلمه استخراج شده است، نمایش می‌دهد. جهت مقایسه چنین توالی‌هایی، پیچش زمان پویا رایج‌ترین راهکار استفاده شده در جستجوی کلمات می‌باشد [۲۰، ۲۱] و به‌طور گسترده‌ای مورد استفاده قرار گرفته است [۲۲]. یک مقایسه بین پیچش زمان پویا و روش‌های پیکسلی در [۲۳] ارائه شده و مقایسه روش‌های مبتنی بر "گرادیانی، ساختاری و خمیدگی" با پیچش زمان پویا در [۱۵] و [۲۴] ارائه شده است.

تمامی رهیافت‌هایی که تا اینجا ذکر شد، نیازمند این بود که قبل از شروع عملیات جستجوی کلمه در مرحله پیش‌پردازش اسناد تصویری به کلمات تشکیل دهنده آن، قطعه‌بندی شود. با توجه به خطای بالای این نوع قطعه‌بندی، سناریوی متفاوت این است که سند فقط به خطوط قطعه‌بندی گردد و محل وقوع کلمه پرسشی در خطوط بررسی گردد یا اینکه خطوط بر اساس کلمه پرسشی رتبه‌بندی شوند. در [۲۵] روشی معرفی شده است که با استفاده از ویژگی‌های مبتنی بر پروفایل هر خط را به‌صورت یک سیگنال یک‌بعدی بازنمایی می‌کند و عمل جستجوی کلمه را با استفاده از پیچش زمان پویا انجام می‌دهد.

مدل مخفی مارکوف و تعمیم‌های مختلف آن، تکنیک‌های پرکاربردی در بازشناسی کلمات تصویری هستند [۲۶] که در جستجوی کلمه نیز مورد استفاده قرار گرفته‌اند. در [۲۷] یک مدل مخفی مارکوف شبه دویبعدی جهت جستجوی کلمه در اسناد چاپی دارای کیفیت پایین پیشنهاد شده است و ادعا شده که نسبت به کجی و اندازه فونت مقاوم است. تفاوت مدل مخفی مارکوف تعمیم‌یافته ارائه شده در [۲۸] نسبت به مدل پایه در انتخاب مدل انتقال است. در مدل مخفی مارکوف هر حالت مخفی تنها می‌تواند یک مشاهده داشته باشد درحالی‌که در مدل تعمیم‌یافته برای هر حالت مخفی مجموعه‌ای از مشاهده‌ها

¹ Annotation

می‌تواند اتفاق بیافتد. در تحقیق [۲۹] پیشنهادی ارائه شده است که از مزایای مدل مخفی مارکوف شبه گسسته^۱ استفاده می‌کند تا مدلی وفق پذیر با شیوه‌های نگارش مختلف بسازد. خاصیت مدل مخفی مارکوف شبه گسسته این است که پارامترهای مدل کلمه را جدا می‌کند. این پارامترها به یک کتابچه کد^۲ مشترک از اشکال و مجموعه‌ای از پارامترهای خاص هر کلمه، تجزیه می‌شود. در واقع ادعا شده است که پارامترهای وابسته به کلمه از پارامترهای وابسته به نویسنده با استفاده از این روش جدا می‌شود.

مدل‌های مخفی مارکوف چندین عیب شناخته‌شده دارند. یکی از این عیوب این است که در مدل مخفی مارکوف فرض می‌شود که احتمال هر مشاهده فقط به حالت جاری وابسته است؛ از این رو مدل کردن اثرات محتوایی مشکل است. علاوه بر این، مدل مخفی مارکوف یک مدل زایشی^۳ است، در صورتی که مدل‌های تمایزی^۴ کارایی بهتری را در رده‌بندی و برچسب‌گذاری دارند.

نوعی از شبکه‌های عصبی بازگشتی^۵ نیز برای جستجوی کلمه مورد استفاده قرار گرفته است. در این روش یک شبکه عصبی بازگشتی برای برچسب‌گذاری یک سری آموزش داده می‌شود. قابلیت این نوع شبکه عصبی این است که یک حافظه از ورودی‌های قبلی در حالات درونی شبکه باقی می‌ماند و استفاده از محتوای قبلی را ممکن می‌سازد، بنابراین قابلیت پذیرش سری را دارا می‌باشد. میزان اطلاعات محتوایی که شبکه عصبی بازگشتی پایه می‌تواند حفظ کند بسیار محدود است (مشکل شیب میرا^۶). برای حل این مشکل، شبکه عصبی با حافظه کوتاه‌مدت بلند^۷ معرفی شده است که لایه‌های مخفی آن شامل زیر شبکه‌هایی با اتصالات بازگشتی است که بلوک حافظه نامیده می‌شوند. شبکه عصبی با حافظه کوتاه‌مدت بلند دوجهته^۸ هم در جستجوی کلمه استفاده شده است [۳۰]. این کار در دو بخش انجام شده است که در بخش اول از شبکه عصبی مذکور استفاده شده است. در این مرحله شبکه عصبی هر موقعیت از سری ورودی را به یک بردار نگاشت می‌کند که احتمال وقوع هر کاراکتر را در آن موقعیت خاص نشان می‌دهد. در مرحله دوم الگوریتم گذر نشانه^۹ CTC این سیستم توانایی بازیابی و جستجوی کلمات خارج از فرهنگ لغت^{۱۰} را نیز داراست. علاوه بر این برای آموزش نیازی به کلمات و حروف مجزا نیز ندارد. شمای کلی عملکرد این سیستم در شکل ۲-۶ نمایش داده شده است. در این روش از شبکه عصبی بازگشتی استفاده

¹ Semi-continuous hidden Markov model (SC-HMM)

² Codebook

³ Generative

⁴ Discriminative

⁵ Recurrent Neural Network (RNN)

⁶ Vanishing gradient problem

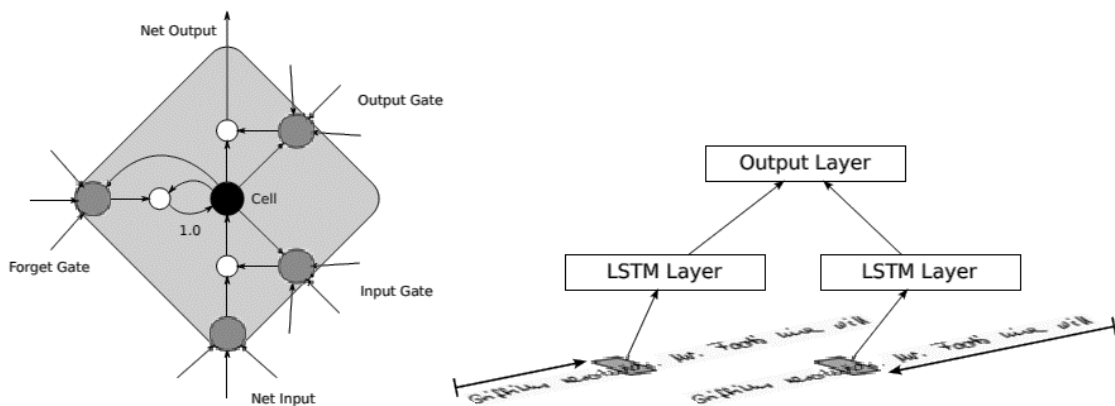
⁷ Long Short-Term Memory (LSTM)

⁸ Bidirectional long short-term Memory (BLSTM)

⁹ CTC Token Passing Algorithm

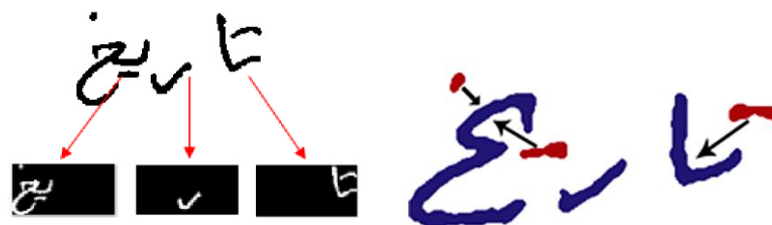
¹⁰ Out of Vocabulary (OOV)

شده است و مانند بسیاری دیگر از روش‌های موجود بزرگ‌ترین چالشی که دارد تفاوت در نوع دست‌خط نویسنده‌ها می‌باشد.



شکل ۲-۶: شمای کلی شبکه عصبی با حافظه کوتاه مدت طولانی دوجهته [۳۰]

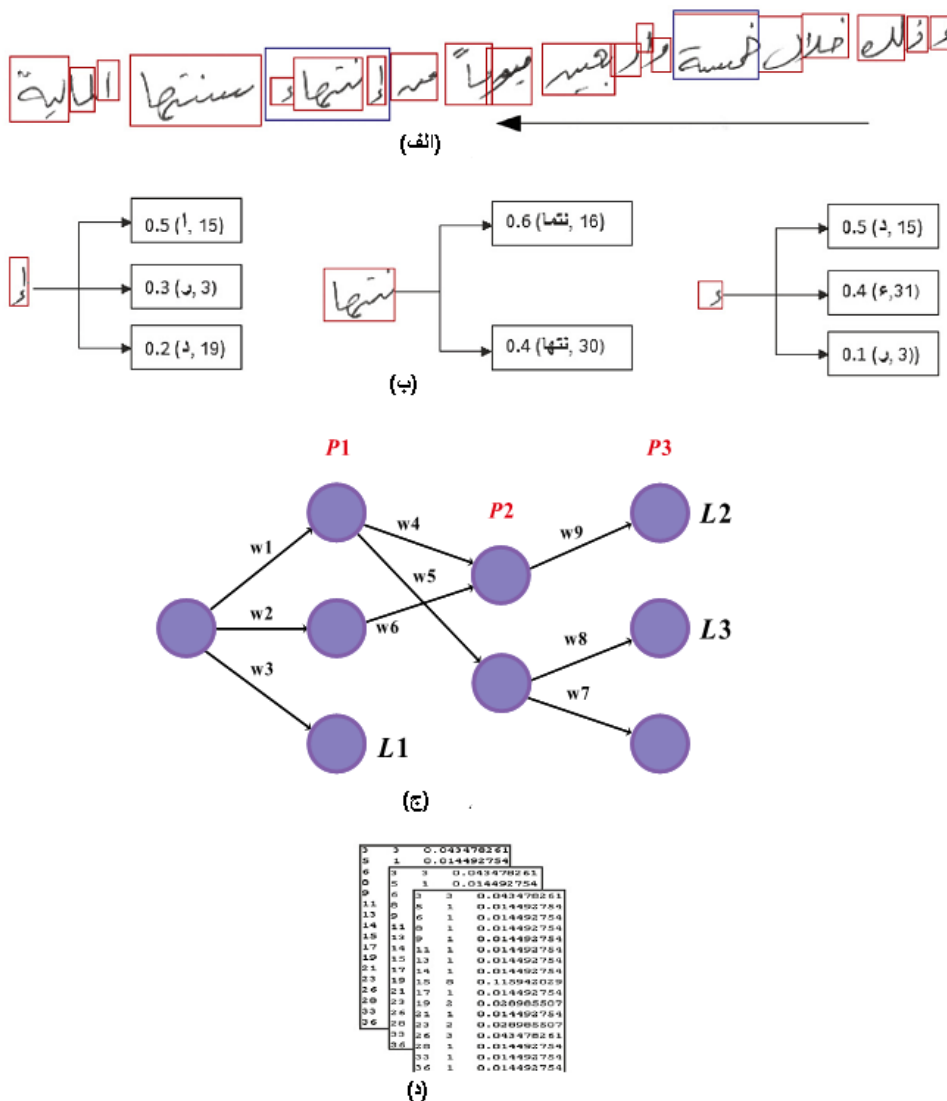
در رهیافتی جدید روشی در سطح خط و مبتنی بر یادگیری ارائه شده است که از مدل زبانی جهت جستجوی کلمه در اسناد عربی بهره می‌برد [۳۱]. این سیستم در ابتدا "بخشی از کلمه عربی" را جدا می‌کند. بخشی از کلمه عربی که در برخی اسناد زیرکلمه گفته می‌شود شامل یک جزء متصل از بدنه اصلی کلمه به همراه نقاط و نشانه‌های آن می‌باشد (شکل ۲-۷). سپس با استفاده از مدل زبانی عمل بازیابی کلمات دست‌نویس انجام می‌شود. در این روش خطوط جداشده از سند و کلمه پرسشی تصویری به بخشی از کلمات عربی قطعه‌بندی می‌شود. ویژگی‌های گرادایانی و رده‌بند آنالیز تفکیکی منظم^۲ جهت بازشناسی بخشی کلمات عربی استفاده شده است. برای هر کدام از بخش‌های استخراج‌شده از کلمه پرسشی یک رده‌بند آموزش داده می‌شود. در نتیجه یک توالی از رده‌بندها برای هر کلمه پرسشی وجود خواهد داشت. بخشی کلمات موجود در هر خط از سند به این رده‌بند سلسله مراتبی داده می‌شود. گراف‌هایی برای بخشی کلمات با ضریب اطمینان بیشتر از یک حد آستانه ساخته می‌شود. سپس مسیره‌های گراف جهت پذیرش یا رد کلمه موردنظر ارزیابی می‌شود (شکل ۲-۸).



شکل ۲-۷: نمونه‌ای از قطعه‌بندی به بخشی کلمه عربی برای کلمه "تاریخ"

¹ Part of Arabic Word (PAW)

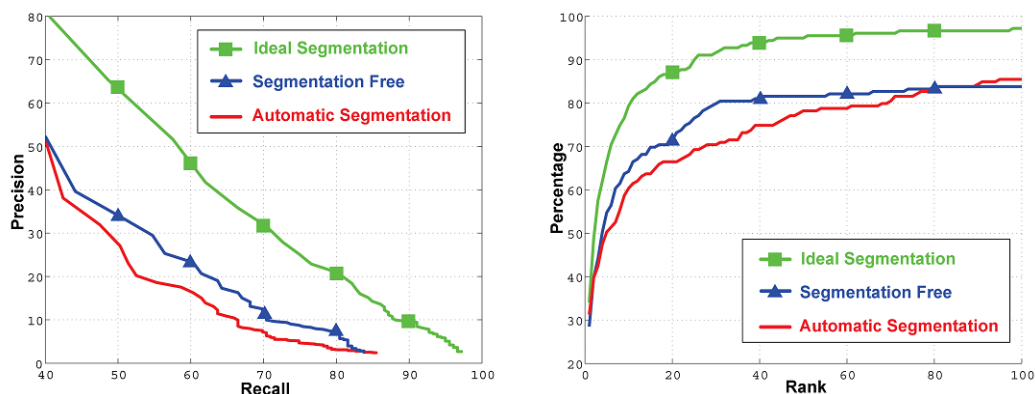
² Regularized Discriminant Analysis (RDA)



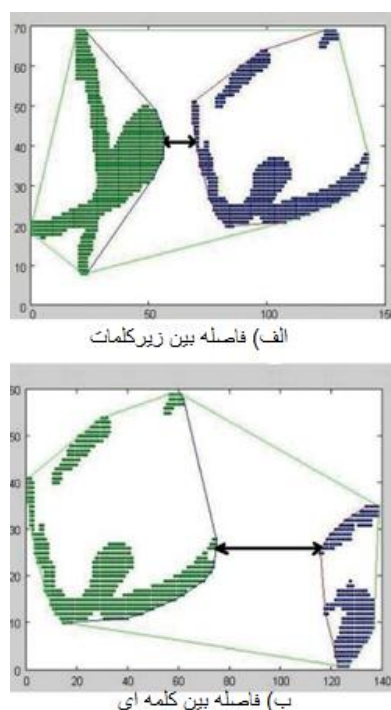
شکل ۲-۸: فرآیند آزمون رده‌بند سلسله مراتبی [۳۱]. الف) قطعه‌بندی. ب) نتایج رده‌بندی. ج) ساخت گراف. د) مدل‌های زبانی

۲-۲-۲- روش‌های مستقل از قطعه‌بندی

روش‌های جدید جستجوی کلمه سعی کرده‌اند تا به علت پرخا بودن عمل قطعه‌بندی از قطعه‌بندی کلمه یا به‌طور کلی قطعه‌بندی اجتناب کنند. جهت مقایسه در [۳۲، ۳۳] دو روش مبتنی بر قطعه‌بندی مستقل از قطعه‌بندی بر روی دست‌نوشته‌های عربی مورد بررسی قرار گرفته است و نشان داده شده است که روش‌های مستقل از قطعه‌بندی دقتی در حدود ۵ تا ۱۰ درصد بالاتر از روش‌های مبتنی بر قطعه‌بندی در دست‌نوشته‌های به زبان عربی دارند. شکل ۲-۹ مقایسه روش مبتنی بر قطعه‌بندی اتوماتیک، مستقل از قطعه‌بندی و مبتنی بر قطعه‌بندی ایده آل را نشان می‌دهد.



شکل ۲-۹: مقایسه دقت و فراخوانی در روش مبتنی بر قطعه‌بندی و مستقل از آن در عربی [۳۲]



شکل ۲-۱۰: قطعه‌بندی کلمه عربی "نقاط" [۳۲]

روش مبتنی بر قطعه‌بندی ارائه‌شده ابتدا زیرکلمات عربی را استخراج می‌کند. این کار بر اساس اجزاء متصل صورت می‌پذیرد. به این صورت که اجزاء متصل استخراج شده و اجزاء کوچکی مانند نقاط به اجزاء بزرگ‌تر نسبت داده شده و یک خوشه را تشکیل می‌دهند. نمونه‌ای از این خوشه‌ها در شکل ۲-۱۰ نشان داده شده است. حرف ابتدای کلمه در بسیاری از کلمات در زبان عربی "الف" است. این مسئله به‌عنوان معیاری جهت تشخیص اول کلمه در این روش استفاده شده است. بعد از این باید تشخیص داده شود که فواصل بین خوشه‌ها مربوط به فاصله بین کلمات است یا فاصله بین زیرکلمات تشکیل‌دهنده یک کلمه،

از این رو جهت رده‌بندی فواصل بین خوشه‌ها از شبکه عصبی مصنوعی استفاده شده است. شکل ۲-۱۰ دو خوشه‌ی مجاور از کلمه "نقاط" را نشان می‌دهد. ویژگی‌هایی که جهت تشخیص این فاصله مورد استفاده قرار می‌گیرند عبارت‌اند از: طول خوشه‌ها، وجود یا عدم وجود حرف "الف" در خوشه‌ها، تعداد اجزاء متصل خوشه‌ها، تفاوت کادر محاط کننده^۱ محتوای دو خوشه، حداقل فاصله بین بدنه‌ی محدب^۲ خوشه‌ها و نسبت مجموع بدنه‌ی محدب در هر دو خوشه به بدنه‌ی محدب کلمه‌ای که هر دو خوشه را شامل می‌شود.

در تحقیقی که اخیراً گزارش شده است [۳۴] اسناد با استفاده از شبکه‌ای از توصیفگرهای هیستوگرام گرادیان^۳ بازنمایی شده است. اسناد به پنجره‌های مربعی شکل با طول ثابتی به اندازه ۵ پیکسل تقسیم‌بندی شده و ویژگی‌های هیستوگرام گرادیان از آن استخراج می‌شود. رهیافت پنجره لغزان راهکاری است که جهت جستجوی کلمه استفاده شده است. با استفاده از روش بدون ناظر ماشین بردار پشتیبان مبتنی بر یک نمونه^۴ و بردار فیشر^۵ [۳۵] نواحی بازیابی شده مجدداً مرتب‌سازی می‌شوند. روش فیشر با توجه به فشردگی زیاد اجازه می‌دهد تا تعداد سند بیشتری به‌طور هم‌زمان در حافظه اصلی قرار گیرد و سرعت جستجو را افزایش می‌دهد. مراحل این سیستم در شکل ۲-۱۲ نشان داده شده است. در این روش اسناد و کلمات پرسشی به پنجره‌های کوچک هم‌اندازه‌ای تقسیم شده و ویژگی هیستوگرام گرادیان از هر پنجره استخراج می‌شود. امتیاز نواحی اسناد از کانولوشن کلمه پرسشی نسبت به نواحی سند به دست می‌آید. این عمل با استفاده از ضرب داخلی توصیفگرهای هیستوگرام گرادیان، به‌عنوان یک معیار مشابهت، صورت می‌پذیرد. علاوه بر این از ماشین بردار پشتیبان مثالی، برای یادگیری بدون ناظر در زمان جستجو با استفاده از کلمه پرسشی استفاده شده است.

روش مستقل از قطعه‌بندی دیگری با بهره‌گیری از قابلیت توصیفگر سیفت^۶ در یک چارچوب کاری مبتنی بر وصله^۷ تصویر یک مدل سبب بصری کلمات^۸ ایجاد کرده است [۳۶]. توصیفگرهای هر وصله تصویر با اعمال اندیس‌گذاری معنایی پنهان^۹ پالایش شده‌اند.

تأثیر قطعه‌بندی نامناسب در روش‌های جستجوی کلمه نیز مورد بررسی قرار گرفته است [۳۷]. در این تحقیق زمانی که برچسب کلمات و یک روش قطعه‌بندی مناسب وجود داشته باشد، استفاده از بردار

¹ Bounding Box

² Convex Hull

³ Histogram of Gradient (HOG)

⁴ Exemplar SVM

⁵ Fisher Vector

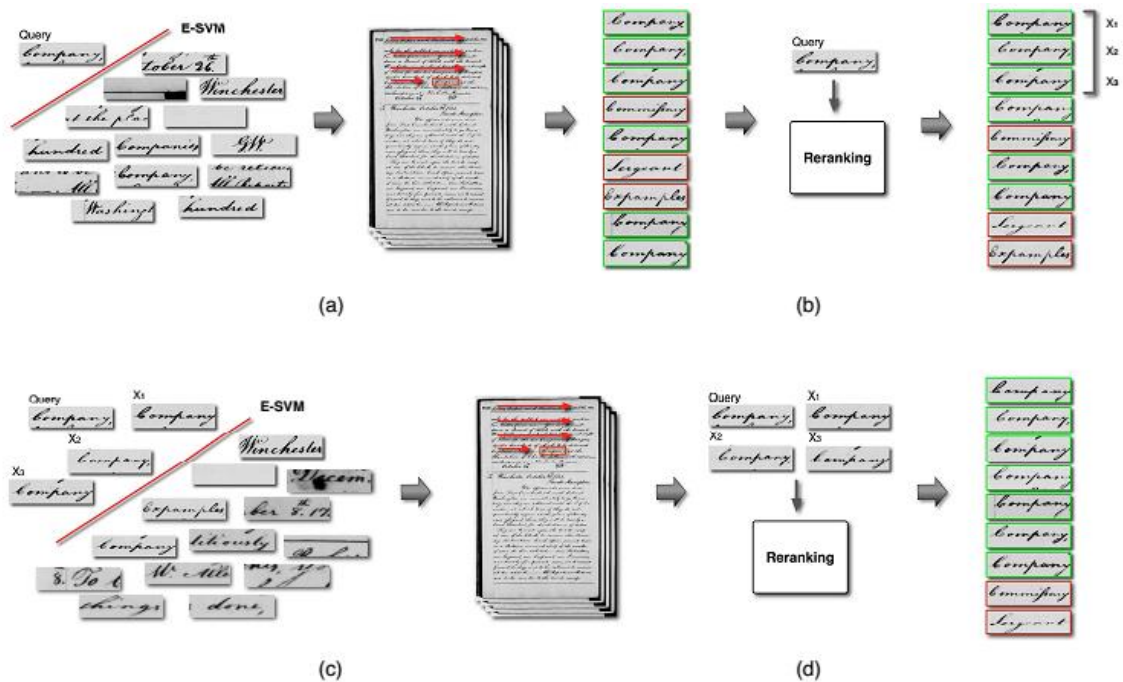
⁶ Scale invariant feature transform (SIFT)

⁷ Patch

⁸ Bag-of-Visual-Word (BoVW)

⁹ Latent Semantic Indexing (LSI)

فیشر بهترین انتخاب معرفی شده است. سید بصری کلمات نیز برای روش‌های مستقل از یادگیری بهترین نتایج را گزارش کرده است [۳۷].



شکل ۲-۱۲: شمای کلی روش معرفی شده در [۳۴]

۲-۳- دسته‌بندی بر اساس روش تشخیص

روش‌های جستجوی کلمه در برخی از مقالات بر اساس روش تشخیص دسته‌بندی شده‌اند. روش تشخیص می‌تواند بر اساس تطبیق قالب در سطح پیکسل یا بر اساس ویژگی‌های استخراج شده انجام شود. عمل تطبیق کلمه پرسشی تصویری با کلمات سند بر اساس یک معیار مشابهت صورت می‌پذیرد. این معیارهای مشابهت تنوع بسیار زیادی در کارهای جستجوی کلمه بر اساس تطبیق دارند. از معیار شباهت ساده‌ای مانند یای انحصاری^۱ [۱۲، ۳۸] تا معیارهای مشابهت معروفی مانند فاصله اقلیدسی و معیار کساین جهت ارزیابی مشابهت استفاده شده است. همچنین معیارهایی برای حل این مسئله خاص ساخته شده‌اند. ویژگی‌های مبتنی بر پروفایل و پیچش زمان پویا از جمله روش‌های قدیمی و پرکاربرد در این زمینه است [۲۳]. اگرچه روش‌های دیگری نیز معرفی شده‌اند، اما روش پیچش زمان پویا و نسخه‌های تغییر یافته آن به‌عنوان یکی از معیارهای مشابهت پرکاربرد در این دسته از روش‌ها

¹ XOR

مورد استفاده قرار گرفته است. بزرگ‌ترین مشکل این دسته روش‌ها، حساس بودن زیاد به تنوع شیوه نگارش و تبدیلات آفاین^۱ است. همچنین روش‌های بیان‌شده بسیار وابسته به نویسنده است و با تغییر نویسنده یا تنوع آن دقت نتایج به‌طور شگفت‌آوری کاهش می‌یابد.

استفاده از ویژگی‌هایی که با پنجره‌گذاری به دست می‌آیند و طول متغیری وابسته به طول و شیوه نگارش کلمه دارند، بسیار رایج است. ویژگی‌های با طول متغیر به دلیل توانایی بیشتر در وفق‌پذیری با تنوع شیوه نگارش و طول کلمه، انعطاف‌پذیری بسیار بیشتری نسبت به روش‌های با ویژگی‌های با طول ثابت از خود نشان می‌دهند [۳۹، ۴۰]. از این‌رو، زمانی که به دنبال استفاده از مدل‌های یادگیری در مسئله جستجوی کلمه بوده‌اند، مدل مخفی مارکوف همواره به‌عنوان یکی از پرکاربردترین روش‌ها مطرح بوده است. مشکلات ویژگی با طول متغیر و استفاده از مدل مخفی مارکوف به شرح زیر است: اولاً، محدودیت در بازیابی و جستجوی کلمه پرسشی وجود دارد. در این حالت فقط توانایی جستجوی مجموعه محدودی از کلمات که در فاز آموزش استفاده شده است امکان‌پذیر خواهد بود و توانایی شناسایی کلمات خارج از مجموعه آموزشی وجود نخواهد داشت. ثانیاً، روش‌هایی که شباهت بین توالی‌ها را ارزیابی می‌کنند معمولاً روش‌های بسیار کندی هستند که این مسئله زمانی که تعداد اسناد افزایش پیدا می‌کند به‌صورت چشم‌گیری خود را نشان می‌دهد. معمولاً در این دسته روش‌ها نیاز به قطعه‌بندی در سطح کلمه وجود دارد. مدل مخفی مارکوف که یکی از پرکاربردترین رده‌بندیهایی است که در چنین سیستم-هایی بکار گرفته شده است، زمانی برای اسناد چند نویسنده نتایج خوبی را به همراه خواهد داشت که مجموعه داده آموزشی بزرگی فراهم شده باشد [۴۱].

۲-۳-۱- روش‌های مبتنی بر تطبیق الگو

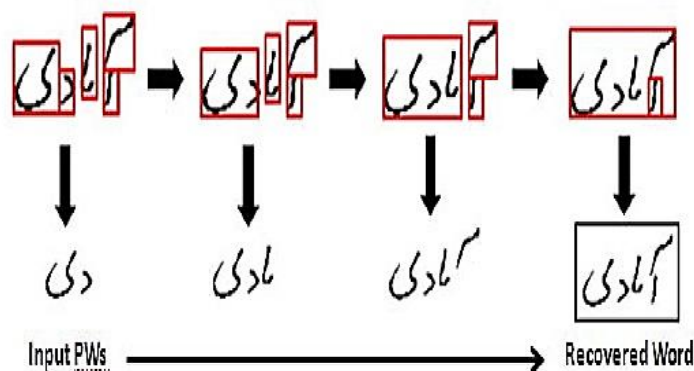
تحقیق [۴۲] نیز روشی را جهت جستجوی کلمه روی اسناد دست‌نویس اردو پیشنهاد داده است. ساختار اصلی و حروف در نوشتار زبان اردو، عربی و فارسی بسیار به هم نزدیک هستند. از آنجایی که ساختار نوشتاری در این زبان‌ها ذات پیوسته‌ای دارد و هر کلمه در این زبان می‌تواند از کنار هم قرار گرفتن زیرکلمات تشکیل شود، روش ارائه‌شده از دو بخش اصلی اندیس‌گذاری و بازیابی تشکیل شده است. در مرحله اندیس‌گذاری قطعه‌بندی در سطح زیرکلمات انجام می‌شود و ویژگی‌های زیرکلمات قطعه‌بندی شده هر سند ذخیره می‌شود. ویژگی‌های استخراج‌شده ساختاری هستند و شامل موارد اسکالر مانند نسبت طول به عرض و مساحت ناحیه محدب، ویژگی‌های برداری مانند پروفایل‌ها، افکنش‌ها و گذر جوهر به پس‌زمینه می‌باشد.

¹ Affine Transform

در زمان جستجو یک فرآیند سه مرحله‌ای انجام می‌شود. مرحله اول کاهش فضای جستجو است که برخی از زیرکلمات از فاز جستجو حذف می‌گردند. این مرحله بر اساس ویژگی‌های اسکالر که بسیار کلی هستند و مربوط به اندازه و طول و عرض ناحیه محدب زیرکلمه می‌باشند انجام می‌شود. پیش از زمان پویا و ویژگی‌های ساختاری برداری در مرحله بعد برای تطبیق زیرکلمات سند و کلمه پرسشی استفاده می‌شود. این مرحله روی زیرکلماتی از اسناد انجام می‌گیرد که از مرحله اول گذر کرده و حذف نشده‌اند. مرحله نهایی ترکیب زیرکلمات کنار هم و به دست آوردن کلمه پرسشی مورد جستجوی کاربر است. شکل ۲-۱۳ و ۲-۱۴ مراحل ذکر شده را نشان می‌دهد.

دل کی آبادی و رونق کا بکرو سا کیا ہے
 ن کی ن میں ویران بھی ہو سکتا ہے
 سوچے سوچے تجھ کو بھینڈا جائے
 اور پھر فوٹا بٹرا نینڈا دینا ہے

شکل ۲-۱۳: تشخیص زیرکلمات بر اساس ویژگی‌های ساختاری برای کلمه پرسشی "آبادی" [۴۲]

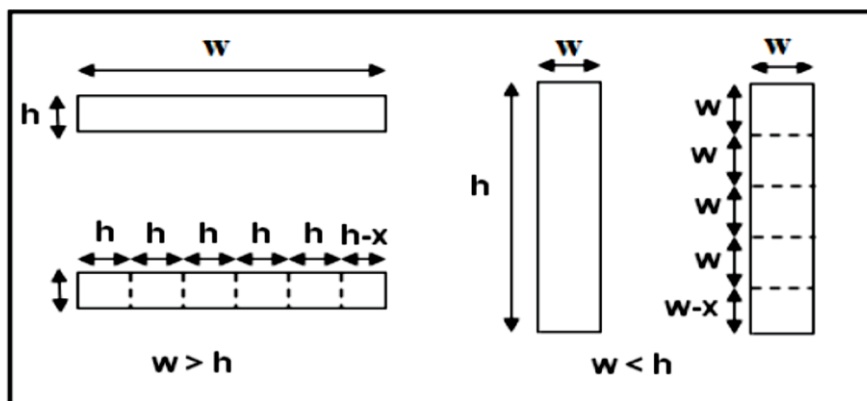


شکل ۲-۱۴: نمونه‌ای از ترکیب نواحی یافت شده برای زیرکلمات کلمه پرسشی "آبادی" [۴۲]

اولین تحقیق در زمینه‌ی جستجوی کلمات فارسی در اسناد چاپی، بر اساس ادعای نویسنده، در [۱۴] گزارش شده است که با بهره‌گیری از شناسایی اندازه و فونت سند مورد جستجو کار می‌کند. ویژگی‌های

استفاده شده جهت بازنمایی کلمات مستقل از فونت است و شامل مواردی مانند تعداد زیرکلمات و نسبت طول به عرض و ویژگی‌هایی که با استفاده از ناحیه‌بندی تصویر به دست می‌آیند، می‌باشد. ویژگی‌هایی که با استفاده از ناحیه‌بندی تصویر به دست می‌آیند تصویر را در جهتی که گستردگی بیشتری دارد بلوک‌بندی می‌کنند. اندازه این بلوک‌ها به اندازه جهتی است که اندازه کوچک‌تری دارد. در هر ناحیه تعداد پیکسل‌های قلم محاسبه می‌گردد. شکل ۲-۱۵ نحوه انجام این بلوک‌بندی را نمایش می‌دهد.

در فاز جستجو، اولین مرحله تشخیص فونت و اندازه سند مورد جستجو می‌باشد. پس از شناسایی فونت و اندازه، کلمه پرسشی بر اساس فونت و اندازه تشخیص داده شده بازنویسی می‌شود و ویژگی‌های مذکور استخراج می‌شوند.



شکل ۲-۱۵: تشریح روش ناحیه‌بندی جهت استخراج ویژگی [۱۴]

در مرحله جستجو عمل تطبیق به صورت چند سطحی انجام شده است. ویژگی نسبت طول به عرض به عنوان یک ویژگی ساده که قابلیت کاهش فضای جستجو را دارد، مورد استفاده قرار گرفته است. در مرحله نهایی فقط عمل تطبیق با کلمات سند موجود در فضای جستجوی کاهش یافته انجام می‌گیرد. در این مرحله ویژگی‌های پروفایل مورد استفاده قرار می‌گیرد. مشکل اصلی گزارش شده در این روش مسئله خطای قطعه‌بندی سند به خطوط و کلمات است.

در تحقیق [۴۳] روشی جهت جستجوی کلمات در اسناد تاریخی فارسی ارائه شده است که بر اساس اندیس‌گذاری زیرکلمات کار می‌کند و از ویژگی‌های ساختاری مبتنی بر پروفایل و پیچش زمان پویا جهت جستجوی کلمه بهره برده است. در این تحقیق ارزیابی روش پیشنهادی روی یکی از کتب تاریخی دانشگاه بیرجند گزارش شده است. برخی نیز سعی کرده‌اند با ارائه روش‌های تطبیق دقیق‌تر و سریع‌تر مشکلات موجود در این حوزه را حل نمایند. روش پیچش زمان پویای دوطرفه با استفاده از ویژگی‌های هیستوگرام گرادیان یکی از روش‌های پیشنهادی در این زمینه بوده است [۴۴]. اخیراً روش تطبیق

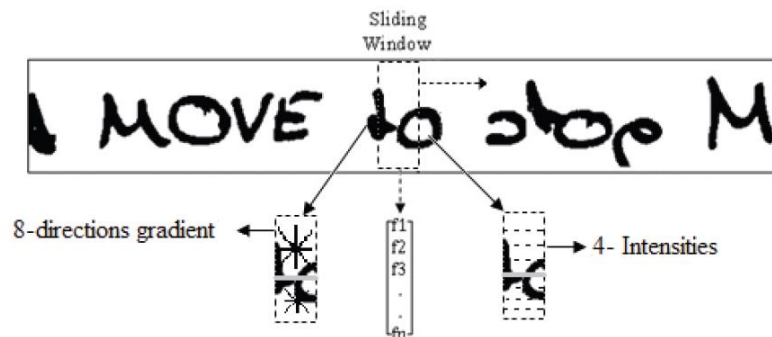
انعطاف‌پذیر توالی^۱ ارائه شده است که توانایی بازیابی کلماتی را دارد که بخشی از آن مشابه با کلمه پرسشی مورد نظر کاربر است [۴۵].

۲-۳-۲- روش‌های مبتنی بر یادگیری

روش‌های جستجوی کلمه مبتنی بر یادگیری جهت کار با اسناد چند نویسنده پیشنهاد شده و نتایج قابل قبولی را نشان داده‌اند. علاوه بر این، پایگاه داده‌های بزرگ نیاز به سیستم‌های مبتنی بر یادگیری دارند. مدل مخفی مارکوف رایج‌ترین رده‌بندی است که در سیستم‌های جستجوی کلمه مورد استفاده قرار گرفته است [۴۶، ۴۷]. روش دیگری که قبلاً شرح داده شده از شبکه عصبی با حافظه کوتاه‌مدت بلند دو جهت بهره می‌برد [۴۸].

در تحقیق [۴۹] مدل مخفی مارکوف به‌عنوان یک مدل یادگیری جهت جستجوی کلمه در اسناد دست‌نویس مورد استفاده قرار گرفته است. در این تحقیق، مدل مخفی مارکوف روی حروف آموزش داده شده است و ترکیب این مدل‌ها و استفاده از مدل فیلتر، روشی مستقل از قطعه‌بندی کلمه را ارائه کرده است که روی زبان‌های مختلفی از جمله انگلیسی، عربی و دوانگاری^۲ مورد ارزیابی قرار گرفته است.

استخراج ویژگی در سند با استفاده از حرکت پنجره لغزان روی خطوط انجام می‌پذیرد. این پنجره لغزان در جهت نگارش و با هم‌پوشانی ۸۵ درصدی حرکت می‌کند. ویژگی‌هایی که از هر فریم از پنجره استخراج می‌شود، ویژگی‌های سطح روشنایی و گرادیانی است. ویژگی‌های سطح روشنایی با تقسیم‌بندی فریم به چهار ناحیه و محاسبه نسبت پیکسل‌های سیاه و سفید محاسبه می‌گردد. به‌صورت شماتیک روش استخراج ویژگی در شکل ۲-۱۶ نشان داده شده است.



شکل ۲-۱۶: روند استخراج ویژگی از یک فریم از پنجره لغزان [۴۱]

¹ Flexible sequence matching (FSM)

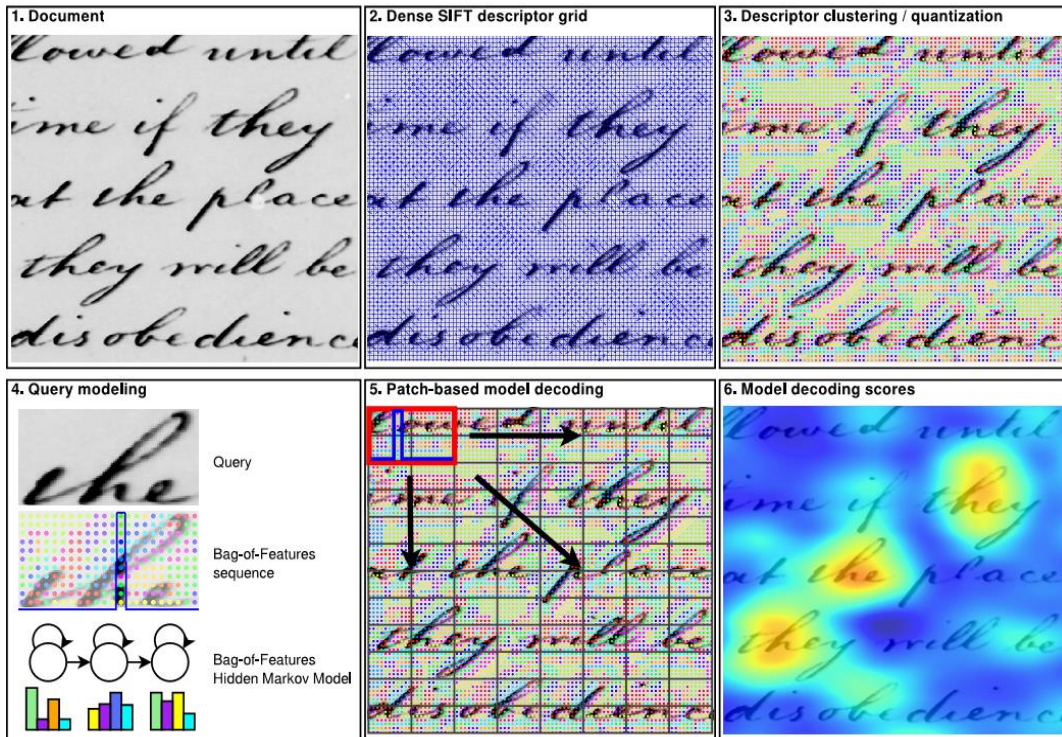
² Devanagar

رهیافت سبب بصری کلمات یکی از روش‌های رایج و پرکاربرد در بینایی ماشین است [۵۰، ۵۱]. رهیافت سبب بصری کلمات روشی است که از مفهوم سبب کلمات در پردازش زبان طبیعی الگوبرداری شده است. در تصاویر، مفهوم کلمه همانند پردازش زبان طبیعی وجود ندارد. بنابراین بخش‌هایی از تصویر که شامل اطلاعات مهم‌تری هستند و می‌توانند تصویر را به‌خوبی توصیف کنند به‌عنوان یک کلمه بصری در نظر گرفته می‌شود. معمولاً این بخش‌های حاوی اطلاعات مهم به‌صورت اتوماتیک، استخراج و خوشه‌بندی شده و نماینده هر خوشه کلمات بصری را می‌سازند. به این منظور، نقاط کلیدی استخراج‌شده از تصویر با استفاده از یکی از توصیفگرهای مناسب بازنمایی می‌شوند و بر اساس این ویژگی‌های استخراج‌شده خوشه‌بندی می‌شوند. در این صورت هر خوشه معرف یک کلمه‌ی بصری^۱ است و مجموعه خوشه می‌تواند به‌عنوان یک فرهنگ لغت بصری در نظر گرفته شود. هر تصویر با استفاده از بردار هیستوگرامی که رخدادهای هر کلمه را در سند تصویری نشان می‌دهد، توصیف می‌گردد [۵۲]. از آنجایی که کیف بصری کلمات در جوامع آنالیز اسناد تصویری محبوبیت بسیاری پیدا کرده است، اخیراً کارهای زیادی با استفاده از این چهارچوب ارائه شده است [۵۳].

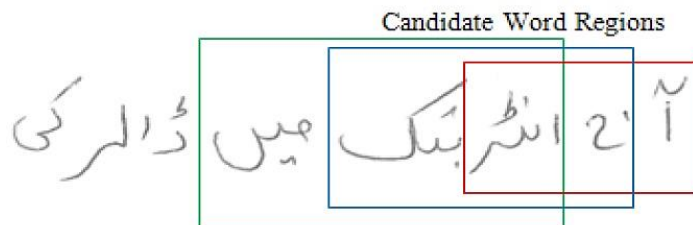
ترکیبی از کیف بصری کلمات و مدل مخفی مارکوف جهت جستجوی کلمه در اسناد دست‌نویس مستقل از قطعه‌بندی در [۵۴] پیشنهاد شده است. این روش برای کلمات با طول کوتاه بهبود خوبی را نشان می‌دهند. علاوه بر این تنها با داشتن یک نمونه می‌تواند مدل کلمه پرسشی را تخمین بزند. ساختار کلی این روش در شکل ۲-۱۷ نشان داده شده است.

تحقیق [۵۵]، روشی را جهت جستجوی کلمه روی اسناد دست‌نویس اردو پیشنهاد کرده است. این روش مبتنی بر قطعه‌بندی در سطح خط انجام شده است. پنجره لغزان بر اساس اجزاء متصل تشکیل‌دهنده زیرکلمات روی خطوط از راست به چپ حرکت می‌کند و نواحی کاندید را ایجاد می‌نماید. شکل ۲-۱۸ نمونه‌ای از حرکت پنجره لغزان جهت ایجاد نواحی کاندید روی خطوط را نمایش می‌دهد. ویژگی‌های گرادینانی و ساختاری از نواحی کاندید استخراج شده و با استفاده از رده‌بند ماشین بردار پشتیبان، جستجوی کلمه انجام می‌شود. یک مرحله اضافی با استفاده از رتبه خروجی ماشین بردار پشتیبان تعداد پیکسل‌های سیاه و تعداد نقاط، نتیجه خروجی مرحله قبل را واریسی و تائید می‌نماید. دقتی برابر ۷۵/۵۰٪ به همراه فراخوانی ۱/۷۰٪ جهت اسناد دست‌نویس اردو گزارش شده است.

¹ Visual word



شکل ۲-۱۷: روند کلی روش ارائه شده در [۵۴]



شکل ۲-۱۸: حرکت پنجره لغزان روی یک خط نمونه و ایجاد نواحی کاندید [۵۵]

۲-۴- تقسیم بندی بر اساس نوع کلمه ی پرسشی

در اکثر روش های جستجوی کلمه در اسناد تصویری، کلمه پرسشی به صورت تصویری در نظر گرفته شده که از تصویر سند اصلی بریده شده است؛ بر این اساس آماده کردن کلمه پرسشی کار زمان بر و طاقت فرسایی است. علاوه بر این، کارایی در سیستم هایی که کلمه پرسشی را فقط به صورت تصویری می پذیرند، نتایج مطلوبی را بر روی اسناد چند نویسنده ارائه نمی کنند. علاوه بر این، در روش هایی که جستجو در آن ها بر اساس یادگیری انجام شده است، مشکل وابستگی به فرهنگ لغت وجود دارد و توانایی جستجو و بازیابی کلمات خارج از فرهنگ لغت وجود ندارد. جستجو با این نوع کلمه پرسشی

می‌تواند به صورت یک مسئله کلاسیک بازیابی تصویر^۱ فرمول‌بندی شود. تصویر کلمه‌ی پرسشی با مجموعه‌ای از تصاویر کاندید در پایگاه داده مقایسه شده و سیستم شبیه‌ترین تصاویر را از پایگاه داده برمی‌گرداند و یا به صورت کلی‌تر تمام کاندیدها را با توجه به میزان شباهت به تصویر کلمه‌ی پرسشی به صورت نزولی رتبه‌بندی می‌کند. از جمله مهم‌ترین معایب این روش می‌توان به این دو مسئله اشاره کرد: اولاً کاربر برای یافتن یک نمونه از کلمه مورد جستجو و مشخص کردن محدوده آن، باید سند را به صورت دستی مورد جستجو قرار دهد. دوماً وابستگی شدید به شیوه نگارش وجود دارد. دو مؤلفه مهم در چنین سیستمی، روش بازنمایی تصاویر و معیار شباهت بکار رفته می‌باشد.

در روش‌های جدیدتر، کلمه پرسشی به صورت یک رشته متنی در نظر گرفته شده است که به این دسته از روش‌ها مبتنی بر کلمه پرسشی متنی گفته می‌شود. این دسته از سیستم‌ها به عنوان کاربردی از بازیابی محتوایی تصویر شناخته می‌شوند. هر یک از سیستم‌ها با کلمه پرسشی تصویری یا متنی مزایای خاص خودشان را دارند. از طرفی سیستم‌هایی که با کلمه پرسشی متنی کار می‌کنند تمام مزایای سیستم‌های شناسایی دستخط مثل انعطاف‌پذیری در جستجوی هر کلمه‌ی کلیدی یا امکان استفاده از حجم زیاد داده‌های برچسب‌دار موجود را دارند. سیستم‌هایی که با کلمه پرسشی تصویری کار می‌کنند به طور خاص برای مواردی که داده‌های برچسب‌دار در دسترس نیستند یا با هزینه بالا به دست می‌آیند مناسب هستند. از طرف دیگر تحقیق [۳۸] نشان می‌دهد که رهیافت‌های کلمه پرسشی تصویری منحصراً بر اساس تطبیق تصویر کار می‌کنند؛ یعنی فاز آموزش با نظارت نیاز ندارد و می‌تواند دقت کافی مورد نیاز سناریوهای عملی را به دست آورد. با توجه به مقایسه فوق انتظار می‌رود که روش‌های مبتنی بر کلمه پرسشی متنی با توجه به توانایی بازیابی فرم‌ها و نوشتارهای مختلف از کلمه ورودی، فراخوانی بهتری داشته باشد. مقایسه این دو دسته سیستم در جدول ۲-۳ آمده است.

نکته قابل ذکر دیگر این است که سیستم‌های مبتنی بر تطبیق الگو با کلمه پرسشی تصویری کار می‌کنند و برخی از سیستم‌های مبتنی بر یادگیری توانایی بازیابی مبتنی بر کلمه پرسشی متنی را دارا هستند. علاوه بر این، برای رهایی از مشکل آماده کردن کلمه پرسشی در این حالت در برخی از کارها از سنتز کردن کلمه پرسشی استفاده کرده‌اند. به این ترتیب که کلمه پرسشی را به صورت متنی دریافت می‌کنند و تصویر آن را ایجاد کرده و به عنوان ورودی به سیستم جستجوی کلمه ارائه می‌دهند [۵۶-۵۸]. سنتز کردن کلمه پرسشی به شیوه نگارش بستگی دارد و در صورت لزوم جستجوی اسناد با چندین نویسنده در اسناد دست‌نویس یا فونت‌های متنوع در اسناد چاپی باید چند مدل متفاوت برای سنتز کردن کلمات ایجاد شود. در زبان فارسی نیز روشی برای جستجوی کلمه روی اسناد چاپی ارائه شده که از سنتز کردن

¹ Image retrieval

کلمه پرسشی استفاده کرده است و تصویر کلمه را بر اساس فونت و اندازه‌ای که از سند مورد جستجو تشخیص داده است، می‌سازد [۱۲].

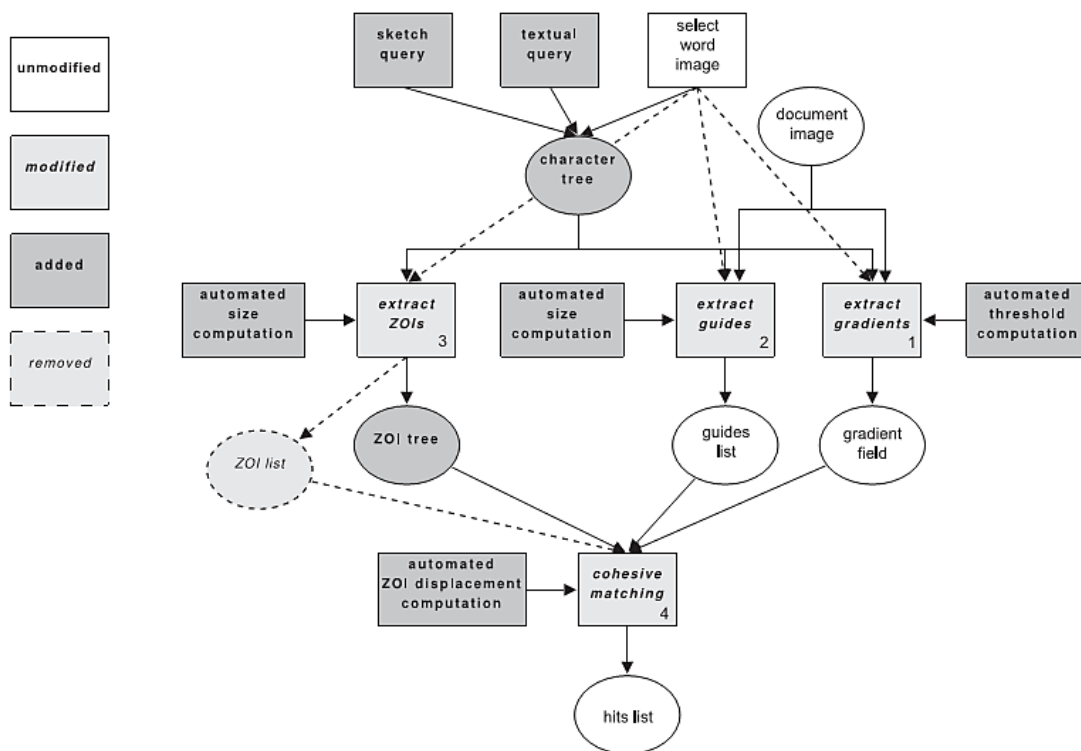
جدول ۲-۳: مقایسه روش‌های جستجوی کلمه با کلمه پرسشی متنی و تصویری

کلمه پرسشی تصویری	کلمه پرسشی متنی
<p>عدم نیاز به دانش اولیه مستقل از الفبا و زبان عدم نیاز به تشخیص و آموزش</p>	<p>نیاز به ساخت مدل نوشتار برای هر سند در اکثر روش‌ها وابسته به الفبا و زبان نیاز به تشخیص و آموزش</p>
<p>نیاز به پیدا کردن دستی کلمه پرسشی به‌طور دستی عدم یافتن حالات متفاوت گرافیکی کلمه مورد جستجو عدم یافتن حالات متفاوت گرامری کلمه مورد جستجو حساسیت بیشتر به تغییرات دست‌نوشته و فونت هزینه محاسباتی بسیار زیاد تطبیق</p>	<p>وارد کردن کلمه پرسشی به‌طور متنی یافتن حالات متفاوت گرافیکی کلمه مورد جستجو یافتن حالات متفاوت گرامری کلمه مورد جستجو</p>

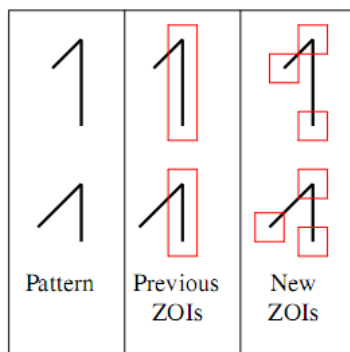
یک موتور جستجوی کلمات با استفاده از کلمات پرسشی تصویری ارائه شده است [۵۹]. این روش بر اساس انتخاب نواحی مورد توجه^۱ و مقایسه این نواحی با استفاده از یک روش تطبیق کشسان منسجم^۲ کار می‌کند. این نواحی مورد توجه بر اساس عملگر مورفولوژی روی تصویر خاکستری به دست می‌آید. المان ساختاری که با استفاده از آن مورفولوژی انجام می‌شود، وابسته به زبان انتخاب می‌شود.

محققان در ادامه کار خود را توسعه داده‌اند و روشی مبتنی بر کلمه پرسشی متنی بر اساس روش ارائه شده قبلی پیشنهاد داده‌اند [۶۰]. در این روش برای هر سند مدلی تعریف شده که ترکیبی از یک الفبا، یک فهرست علائم^۳ و یک گرامر می‌باشد. ساختار کلی سیستم در شکل ۲-۱۹ نشان داده شده است. البته در این روش بهبودی در انتخاب نواحی مورد توجه با استفاده از ویژگی‌های تفاضلی نیز داده شده است (شکل ۲-۲۰). نواحی مورد توجه بخش‌ها دارای اطلاعات بیشتر در تصویر هستند. مهم‌ترین نکته ذکر شده در این روش چنین بیان شده است که این تغییر باعث بهبود در جستجوی کلمات انگلیسی، چینی و ... شد اما روی زبان عربی به خاطر وجود نقاط و سرکش‌ها نتایج ضعیف‌تری را ارائه کرد.

1 Zone of Interest (ZOI)
2 Cohesive elastic matching
3 Glyph book



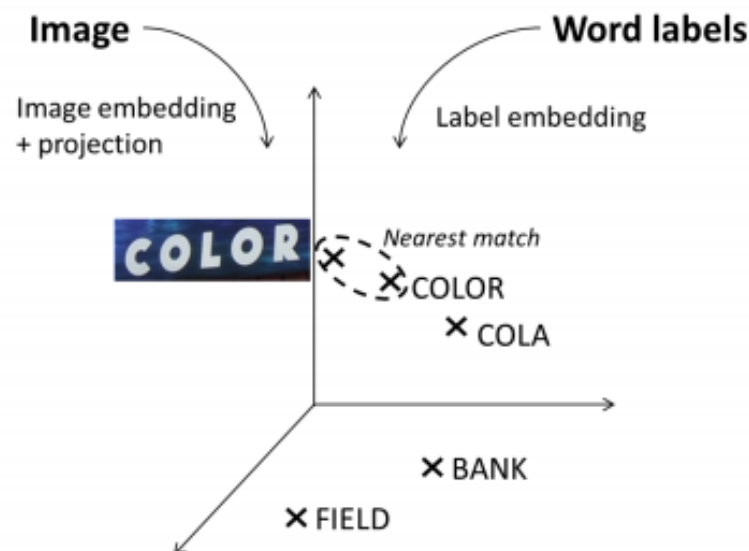
شکل ۲-۱۹: فرآیند کلی روش پیشنهادشده در [۶۰] با نمایش تغییرات نسبت به کار قبلی در [۵۹]



شکل ۲-۲۰: تفاوت انتخاب نواحی مورد توجه در [۵۹] و [۶۰]

اخیراً روشی مبتنی بر قطعه‌بندی جهت جستجوی کلمه ارائه شده که توانایی جستجو با کلمات پرسشی متنی و تصویری را داراست و می‌تواند برای تشخیص دست‌نوشته نیز مورد استفاده قرار گیرد [۶۱]. در این روش، تصویر هر کلمه و کلاس (رشته‌ی متنی) مربوط به آن در یک فضای برداری مشترک تعبیه می‌شود. در این فضا تصاویر و رشته‌ای که نشان‌دهنده‌ی کلمه پرسشی یکسان هستند نزدیک به یکدیگر قرار می‌گیرند که این اجازه را می‌دهد که عملیات بازشناسی به‌صورت مسئله نزدیک‌ترین همسایگی انجام شود. در این روش از ترکیب رده‌بندی مبتنی بر خصیصه، تعبیه برچسب و زیرفضای مشترک استفاده شده است.

در بینایی ماشین از تعبیه تصویر ورودی^۱، یعنی بازنمایی آن به فضای ویژگی، بسیار استفاده شده است. روش‌های متعددی مانند کدگذاری بر اساس پچ‌های تصویر، کاهش ابعاد، فشرده‌سازی و روش‌های مبتنی بر هسته از این نمونه هستند. اما به تعبیه برچسب کمتر پرداخته شده است. تعبیه برچسب یعنی برچسب کلاس در یک فضای جدید با قابلیت تمایز بالاتر نگاشت شوند به نحوی که کلاس‌های مشابه نزدیک به هم قرار گیرند. اخیراً، روشی بر اساس تعبیه کردن رشته (برچسب کلاس در بازشناسی و بازیابی کلمات تصویری) در زبان انگلیسی در [۶۲] پیشنهاد شده است. به این ترتیب که هر کلاس در فضای برداری خصیصه تعبیه می‌شود. به شرط اینکه تابع تعبیه مناسب انتخاب شود، یعنی کلاس‌های مشابه در فضای تعبیه بر اساس فاصله اقلیدسی نزدیک باشند، تعبیه برچسب روشی مؤثر برای اشتراک پارامترهای بین کلاسی است. هدف این است که یک بازنمایی مشترک^۲ برای تصویر و رشته‌ی متنی کلمه وجود داشته باشد که مشخصه‌های استخراج‌شده از رشته و مشخصه‌های پیش‌بینی‌شده از تصویر در یک زیرفضای مشترک قابل مقایسه باشند. شمای کلی این سیستم در شکل ۲-۲۱ نمایش داده شده است.



شکل ۲-۲۱: شمای کلی سیستم [۶۲]

با توجه به قدرت و کارایی این روش، اخیراً محققان زیادی سعی در گسترش و بهبود این روش داشته‌اند. این روش با تغییراتی جزئی برای زبان یونانی نیز مورد استفاده قرار گرفته است [۶۳]. علاوه بر این

¹ Input embedding

² Joint representation

چارچوبی جهت استفاده از این روش با استفاده از پنجره‌گذاری نیز معرفی شده است که قابلیت مستقل بودن از قطعه‌بندی را نیز به آن می‌دهد [۶۴]. استفاده از شبکه عصبی کانولوشنی جهت استخراج هیستوگرام سلسله مراتبی کاراکترها که کلمه پرسشی را به صورت متنی و تصویری می‌پذیرد نیز توسعه جدیدی از این کار می‌باشد [۶۵].

فصل

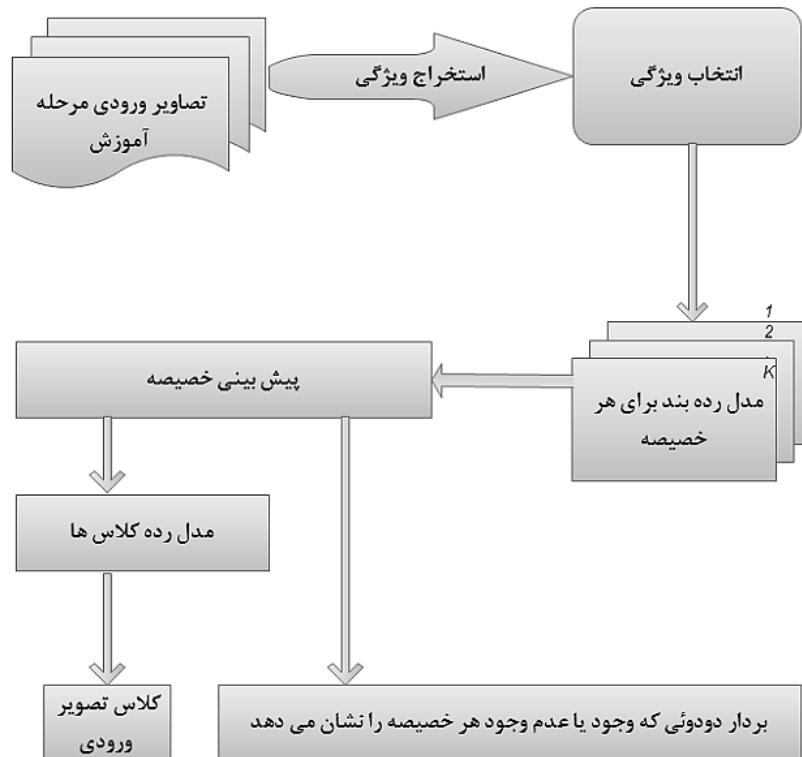
۳- روش پیشنهادی اول (مبتنی بر قطعه‌بندی کلمه)

۳-۱- مقدمه

با توجه به تعداد و تنوع زیاد کلاس‌های اشیاء در دنیای امروز، فراهم آوردن و برچسب‌گذاری تصاویر آموزشی برای تمامی کلاس‌ها در بسیاری از کاربردها غیرممکن به نظر می‌رسد. این تنوع کلاس‌ها در بسیاری از کاربردهای بینایی ماشین دیده می‌شود که شناسایی، جستجو و بازیابی کلمات دست‌نویس یکی از این نمونه‌های کاربردی است. برای حل این مشکل، رده‌بندی مبتنی بر خصیصه معرفی شده است. تفاوت این روش با روش‌های رده‌بندی کلاسیک در این است که به جای نام‌گذاری کلاس‌ها، آن‌ها را توصیف می‌کند. با توجه به اینکه این توصیفات می‌تواند بین کلاس‌های مختلف به اشتراک گذاشته شود، در مرحله آموزش یک هم‌افزایی ایجاد خواهد شد. علاوه بر این توانایی شناسایی و توصیف کلاس‌های دیده نشده هم به وجود خواهد آمد.

در این رهیافت تشخیص شیء با استفاده از توصیفات سطح بالای خاص انسان که خصیصه نامیده می‌شوند، انجام می‌گردد. این خصیصه‌ها می‌توانند ویژگی‌های ساختاری یا معنایی مانند شکل، رنگ و یا حتی اطلاعات جغرافیایی شیء باشند [۶۶، ۶۷]. اشتراک خصیصه‌ها بین چندین کلاس در مجموعه آموزشی علاوه بر این که می‌تواند باعث آموزش بهتر شود، برای شناسایی کلاس‌های دیده نشده نیز مورد استفاده قرار می‌گیرد. این روش سعی بر یادگیری به شیوه انسان دارد. انسان درحالی که کلاس اشیاء را بر اساس مثال‌های مشاهده‌شده یاد می‌گیرد، توانایی شناسایی کلاس‌های جدید را بر اساس توصیفات مشترک با کلاس‌های قبلی دارد. برای نمونه اگر در شناسایی حیوانات در یک مجموعه آموزشی سه خصیصه "بزرگ بودن"، "دارای رنگ خاکستری" و "دارای خرطوم" یاد گرفته شده باشد، چنانچه حیوانی که این سه خصیصه را به‌طور هم‌زمان داراست "فیل" نامیده شود، برای انسان قابل‌شناسایی خواهد بود، حتی اگر قبلاً نمونه‌ای از آن دیده نشده باشد. این روش به‌جای انتصاب یک نمونه به یک کلاس، آن را بر اساس خصیصه‌هایش توصیف می‌کند. خصیصه‌ها به‌عنوان یک روش بازنمایی میانی عمل می‌کنند که امکان به اشتراک‌گذاری پارامتر بین کلاس‌ها را ممکن می‌سازند. چنین رهیافتی خصوصاً در مواقعی که داده‌های آموزشی کم هستند، یادگیری را بهبود می‌بخشد. ساختار کلی این روش در شکل ۳-۱ آمده است.

اخیراً برای شناسایی و جستجوی کلمات در متون انگلیسی روش‌هایی ارائه شده است که از روش رده‌بندی مبتنی بر خصیصه استفاده می‌کند [۷، ۴۱، ۶۸-۷۰]. از آنجایی که این روش به‌عنوان روش پایه در این رساله مورد استفاده قرار گرفته است، در این فصل این روش پایه به همراه خصیصه‌هایی که برای زبان فارسی پیشنهاد و ارزیابی شده شرح داده می‌شود.



شکل ۳-۱: شمای کلی روش رده بندی مبتنی بر خصیصه

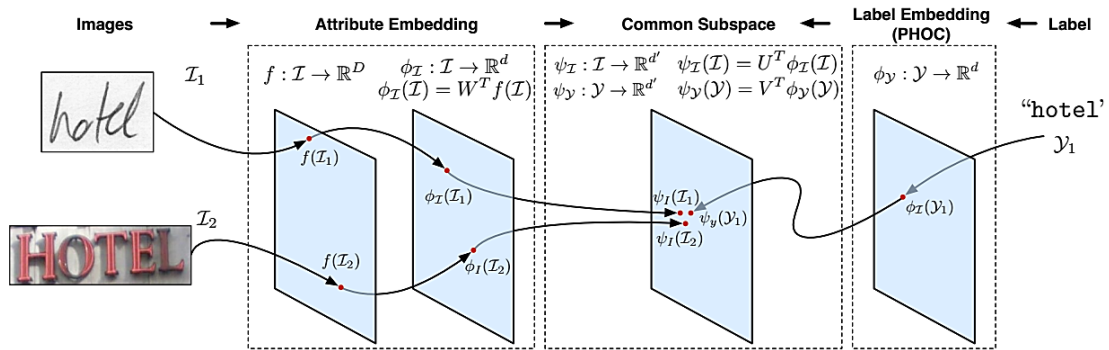
در این روش برای شناسایی و جستجو و بازیابی کلمات تصویری، آموزش یک بازنمایی مشترک برای رشته متنی (کلاس کلمه) و کلمه تصویری پیشنهاد شده است. با استفاده از این بازنمایی، مسئله شناسایی، جستجو و بازیابی کلمات تصویری به یک مسئله نزدیک ترین همسایه^۱ ساده تبدیل می شود. در این روش رشته متنی به یک بردار d بعدی در فضای باینری تبدیل می شود که در واقع این بازنمایی نشان می دهد که آیا کاراکتری خاص در یک ناحیه مکانی خاص از رشته قرار گرفته است یا نه. این بردار به عنوان منبعی از خصیصه های کاراکتری مورد استفاده قرار می گیرد. جهت هماهنگ سازی تصویر کلمه با متن تصاویر کلمه نیز به یک فضای تفکیک کننده d بعدی نگاشت می شود که محتمل بودن وقوع یک کاراکتر خاص در یک ناحیه خاص را کم می کند (هیستوگرام هرمی کاراکترها^۲). با آموزش این خصیصه های کاراکتری داده های آموزشی بهتر مورد استفاده قرار می گیرند و بازشناسی و بازیابی کلمات دیده نشده در فرهنگ لغت آموزشی ساده و آسان خواهد بود.

به دلیل این که بازنمایی d بعدی برای رشته متنی، باینری است و برای تصویر، امتیازات خصیصه ها است که باینری نیستند (در واقع امتیازات خروجی ماشین های بردار پشتیبان هستند)، مقایسه آن ها به صورت مستقیم بهینه نیست و نیاز به کالیبره کردن وجود دارد. نهایتاً آموزش یک زیرفضای مشترک با ابعاد

¹ Nearest neighbor

² Pyramid Histogram of characters (PHOC)

کم با یک معیار مشترک بین هر دو بازنمایی پیشنهاد شده است. این کار دو مزیت اساسی دارد: اولاً مقایسه‌ی مستقیم تصویر و رشته متنی را معنی‌دار می‌کند و ثانیاً بازنمایی مثال‌هایی از تصاویر یک کلمه، نزدیک به هم قرار می‌گیرند. ساختار کلی روش معرفی شده در شکل ۲-۳ آمده است.



شکل ۲-۳: شمای کلی روش معرفی شده در [۷, ۴۱, ۶۸]

به‌منظور بازیابی و بازشناسی کلماتی که در زمان آموزش دیده نشده‌اند، لازم است بتوان امکان انتقال دانش بین نمونه‌های آموزش و آزمون ایجاد کرد. یکی از معروف‌ترین روش‌ها برای انجام این آموزش شات صفر^۱ در بینایی ماشین استفاده از خصیصه‌های بصری^۲ است [۶۶, ۶۷]. در این کار از خصیصه‌های کاراکتری برای انتقال اطلاعات بین نمونه‌های آموزش و آزمون استفاده شده است. هرچند ایده قطعه-بندی کلمه به کاراکترها و آموزش در سطح کاراکتر قبلاً بارها استفاده شده است، اما این‌گونه رهیافت‌ها معمولاً وابسته به مدل‌های HMM با ویژگی‌هایی با طول متغیر هستند و بنابراین کارایی توسط این مدل‌ها محدود شده است. در روش پیشنهادی چارچوبی وسیع‌تر ارائه شده است که هیچ محدودیتی در انتخاب ویژگی‌ها یا روش‌های آموزش خصیصه‌ها وجود ندارد.

روش پیشنهادی وابسته به روش‌های تعبیه برچسب نیز می‌باشد که برچسب‌ها درون یک فضای متفاوت تعبیه می‌شوند و یک تابع سازگاری^۳ بین تصاویر و برچسب‌ها تعریف شده است [۷, ۶۸, ۶۹, ۷۱, ۷۲]. در [۶۹] روشی جهت بازشناسی متن با استفاده از رهیافت تعبیه متن و هرم مکانی کاراکترها^۴ پیشنهاد شده است. تفاوت اصلی این روش با روش ارائه‌شده در [۷, ۶۸] و چارچوب اصلی روش پیشنهادی در این رساله در نحوه تعبیه کردن برچسب‌ها می‌باشد. روش اول سعی دارد تا با استفاده از چهارچوب

¹ Zero-shot learning

² Visual attribute

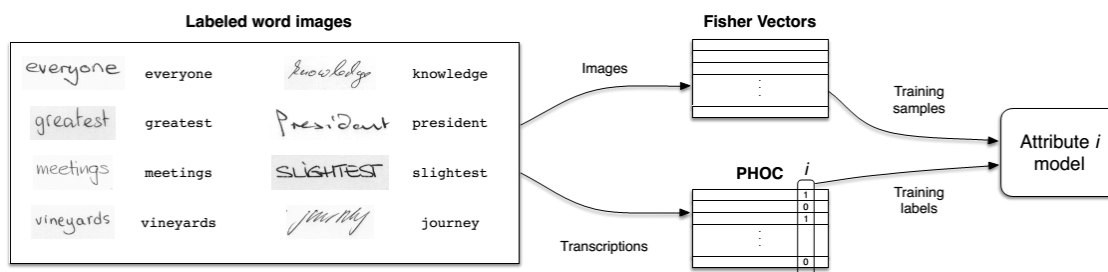
³ Compatibility function

⁴ Spatial pyramid of characters (SPOC)

ماشین بردار پشتیبان بین بازنمایی تصویر و SPOC های متناظر یک زیر فضای مشترک پیدا کند، در صورتی که در روش دوم PHOC ها به عنوان خصیصه مورد استفاده قرار می گیرند و فقط سعی شده است که یک زیر فضای مشترک بین خصیصه ها و PHOC ها پیدا شود. شمایی از این دو روش در شکل ۳-۳ و ۳-۴ آمده است.



شکل ۳-۳: روش تعبیه SPOC و مدل تعبیه برچسب در [۶۹].



شکل ۳-۴: فرآیند آموزش مدل خصیصه نام [۷].

۳-۲- بازنمایی کلمه مبتنی بر خصیصه

یکی از رایج ترین روش ها برای استفاده از آموزش با نظارت در مسئله جستجوی کلمه، آموزش مدل برای هر کلمه به صورت مجزا می باشد. با استفاده از مجموعه ای از نمونه های مثبت و منفی، برای هر کلمه کلیدی یک مدل آموزش داده می شود. در زمان آزمون احتمال کلمه ورودی بر اساس مدل های آموزش داده شده محاسبه می شود و به عنوان امتیاز مورد استفاده قرار می گیرد. این دسته از روش ها تعداد کلمات کلیدی قابل جستجو را محدود می کند. روشی نیز مبتنی بر شبکه عصبی شبه گسسته معرفی شده است [۲۹] که مدل توانایی وفق یافتن با هر کلمه پرسشی ورودی را به صورت برخط دارد. بنابراین این روش توانایی جستجو و بازیابی کلمات خارج از فرهنگ لغت را نیز دارا است، اما به هر حال برچسب ها در فاز آموزش استفاده نمی شوند.

از معایب دیگر این روش‌ها آموزش در سطح کلمه است، یعنی اطلاعات بین کلمات مشابه به اشتراک گذاشته نمی‌شود. به‌عنوان مثال اگر مدلی برای کلمه "بانک" آموزش داده شد، کلمه "باند" یک نمونه منفی برای آن در نظر گرفته می‌شود. درحالی‌که به اشتراک‌گذاری این اطلاعات مشترک بین کلاس‌های کلمات می‌تواند برای یادگیری یک بازنمایی تفکیک‌کننده خوب، مهم باشد. بنابراین استفاده از خصیصه‌ها یکی از روش‌های رسیدن به این هدف است. خصیصه‌ها، ویژگی‌های معنایی هستند که با استفاده از آن‌ها می‌توان تصاویر و دسته‌های آن‌ها را توصیف کرد. اخیراً، استفاده از رده‌بندی مبتنی بر خصیصه برای مسئله رده‌بندی و بازیابی تصاویر شهرت زیادی پیدا کرده است [۶۷، ۷۲، ۷۳].

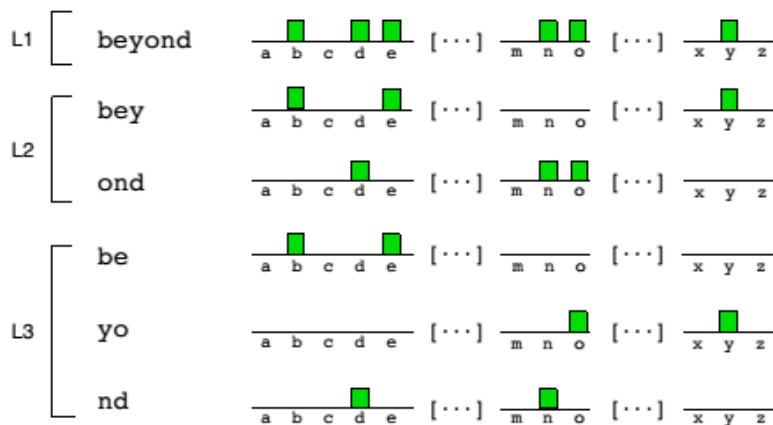
با توجه به ویژگی‌های ذکرشده در مورد خصیصه‌ها، استفاده از آن می‌تواند برای بازنمایی کلمات مناسب باشد. دلیل این ادعا را می‌توان به این صورت ذکر کرد که در این روش کلمات مختلف در فاز آموزش می‌توانند به تبادل اطلاعات بپردازند و این منجر به ایجاد امضاهای فشرده می‌شود. انتخاب این خصیصه‌ها فرآیندی وابسته به کاربرد است، بنابراین برای مسئله جستجوی کلمه باید طوری تعریف شوند که بتوانند به خوبی کلمات مختلف را از همدیگر متمایز کنند.

در زیر بخش بعدی ابتدا روش تعبیه برچسب PHOC که رشته متنی را در یک فضای برداری دودویی تعبیه می‌کند و پایه‌ی روش‌های پیشنهادی این بخش است را معرفی می‌کنیم. سپس روش‌های پیشنهادی مناسب برای زبان فارسی را شرح خواهیم داد. نهایتاً چگونگی استفاده از این بازنمایی جهت کد کردن تصویر کلمه را بیان خواهیم کرد.

۳-۲-۱- تعبیه برچسب متنی

۳-۲-۱-۱- روش PHOC [۷]

ساده‌ترین روش برای تعبیه رشته متنی ساخت یک هیستوگرام از کاراکترهاست که بر اساس تعداد کاراکترهای الفبای زبان ساخته می‌شود. از آنجایی‌که این روش برای زبان انگلیسی معرفی شده است و ارقام صفر تا نه را نیز در بر می‌گیرد، برداری دارای ۳۶ بعد می‌تواند در نظر گرفته شود که هر بعد مشخص می‌کند که رشته متنی شامل کاراکتر خاصی می‌شود یا خیر. اما این روش بازنمایی متمایز کننده خوبی را برای کلمات ارائه نمی‌کند. به این دلیل که بسیاری از کلماتی که حروف یکسان با جایگشتی متفاوت از کاراکترها را دارند، مانند "پاکت" و "پاتک" و "تکپا"، در این بازنمایی یکسان نمایش داده خواهند شد. بنابراین یک نسخه هرمی از این هیستوگرام کاراکتری معرفی شده است. شکل ۳-۵ نمونه‌ای از این بازنمایی را در سه سطح اول نمایش داده است.



شکل ۳-۵: بازنمای PHOC در سه سطح ۱ و ۲ و ۳ برای رشته متنی “beyond” [۷]

در این روش به جای بررسی وجود کاراکتری خاص در کل کلمه، این کار در نواحی متفاوتی از کلمه صورت پذیرفته است. در سطح دوم وجود هر کاراکتر در نیمه اول و دوم مورد بررسی قرار می‌گیرد و در سطح سوم کلمه به سه بخش تقسیم شده و برای هر سه بخش این بررسی انجام می‌شود و به همین ترتیب تا هر سطحی که لازم باشد پروسه می‌تواند ادامه پیدا کند. مثلاً در تحقیق گزارش شده در [۶۹] همین پروسه برای سطوح ۲ تا ۵ اجرا شده که هیستوگرامی با $504 = 36 \times (2+3+4+5)$ بعد تولید شده است. برای کد کردن رابطه بین کاراکترهای همسایه در این کار ۵۰ بایگرام^۱ رایج انگلیسی نیز در سطح دوم مورد بررسی قرار گرفته است که باعث افزایش ۱۰۰ بعدی هیستوگرام و تبدیل شدن آن به ۶۰۴ بعد می‌شود. در صورت استفاده از این روش مشکل عدم متمایزکنندگی روش قبلی از بین خواهد رفت.

با توجه به ساختار PHOC برای تبدیل رشته متنی به بردار دودویی به تعیین نواحی هر می که به هر کاراکتر اختصاص داده می‌شود، نیاز داریم. برای این منظور ناحیه اشغال شده نرمال شده توسط کاراکتر k از یک رشته n کاراکتری به صورت یک فاصله به شرح زیر تعریف می‌شود:

$$Occ(k, n) = \left[\frac{k}{n}, \frac{k+1}{n} \right] \quad (1-3)$$

که در این رابطه k از صفر شروع می‌شود. این اطلاعات از برچسب کلمات (رشته متنی) استخراج می‌شوند. از همین رابطه برای به دست آوردن ناحیه اشغال شده r در سطح l نیز استفاده می‌شود. اگر اشتراک این نواحی اشغال شده بزرگ‌تر یا مساوی ۵۰٪ بود، کاراکتر به این ناحیه اختصاص داده می‌شود.

$$\frac{Occ(k, n) \cap Occ(r, l)}{Occ(k, n)} \geq 0.5 \quad (2-3)$$

¹ bigram

در اولین مرحله دقیقاً از همین ساختار برای زبان فارسی استفاده شده و ساختار خصیصه پیشنهادی هیستوگرام هرمی کاراکترهای فارسی^۱ نام‌گذاری شده است.

۳-۲-۱-۲- تعبیه برچسب متنی خاص زبان فارسی

از آنجایی که ساختار زبان فارسی ذاتاً پیوسته است و حروف اشکال متفاوتی را بسته به محل قرار گرفتن می‌توانند داشته باشند روش‌های دیگری برای تعبیه خصیصه‌ها در زبان فارسی پیشنهاد شده است که با ساختار نگارشی آن تناسب داشته باشد.

در روش اول پیشنهادی بر اساس تفاوت شکل حروف فارسی در موقعیت‌های مختلف داده شده است. در این روش با توجه به جدول ۱-۱ به جای در نظر گرفتن حروف در هر سطح از هرم، شکل حروف بر اساس موقعیت در نظر گرفته می‌شود. انگیزه اصلی این پیشنهاد را می‌توان شکل بسیار متفاوت برخی حروف در ۴ حالت موقعیتی عنوان کرد. با توجه به این که این وضعیت در ۲۴ حرف از حروف فارسی وجود دارد، می‌توان این ذهنیت را ایجاد کرد که مجزا در نظر گرفتن این حالات می‌تواند باعث آموزش بهتر مدل‌های خصیصه‌ها گردد.

در این حالت به جای ۳۳ حالت حروف زبان فارسی ("ا" و "آ" دو حرف متفاوت در نظر گرفته شده است) ۱۱۴ حالت در نظر گرفته می‌شود. این روش نیز برای ۴ سطح ۲ تا ۵ برای بازنمایی رشته‌های متنی فارسی استفاده شده است، بنابراین برداری به طول $1624 = 116 \times (2+3+4+5)$ تولید خواهد شد. این روش بازنمایی و ساختار خصیصه ارائه‌شده، هیستوگرام هرمی کاراکترهای موقعیتی فارسی^۲ نامیده شده است.

بیشتر اشکال کاراکتری فارسی دارای دو قسمت اصلی هستند: بدنه اصلی کاراکتر، نقاط و سرکش‌ها. به دلیل اینکه در بسیاری از کاراکترها این بدنه یا نقاط مشترک هستند در روش دوم پیشنهادی سعی شد تا این قسمت‌ها به‌طور مجزا کد شوند. این وضعیت در حالات چهارگانه موقعیتی حروف نیز یکسان باقی می‌ماند. البته در برخی از حروف تفاوت‌هایی وجود دارد مانند حرف "ن" که حالت میانی و ابتدایی شبیه به دسته حروف "ب پ ت ث" دارد اما حالات تنها و آخر آن متفاوت است. جداول ۱-۳ و ۲-۳ به ترتیب دسته‌بندی نقاط و سرکش‌ها و بدنه‌ها را نشان می‌دهد. تعداد بدنه‌ها ۵۵ حالت مختلف و ۹ حالت نقاط و سرکش‌ها را نشان می‌دهد. در این صورت کل حالات برای ۶۱ دسته مختلف خواهد شد. بنابراین در

¹ Persian pyramid histogram of characters (PPHOC)

² Persian pyramid histogram of positional characters (PPHOPC)

صورت بازنمایی ۴ سطحی از سطوح ۲ تا ۵ تعداد ابعاد برابر $۹۱۰ = ۶۵ \times (۲+۳+۴+۵)$ خواهد بود. این روش، هیستوگرام هرمی نشانه‌های موقعیتی فارسی^۱ نام‌گذاری شده است.

جدول ۱-۳: حروف دارای نقاط و سرکش مشترک

ردیف		نقاط و سرکش‌های مشترک
۱	آکولاد	آ، ا
۲	تک نقطه بالا	خ، خ، خ، ذ، ذ، ز، ز، ض، ض، ض، ظ، ظ، غ، غ، غ، غ، ف، ف، ف، ف، ف، ف، ن
۳	تک نقطه پایین	ب، ب، ب، ج، ج، ج، ج، چ
۴	دونقطه بالا	ت، ت، ت، ق، ق، ق، ق، ق
۵	دونقطه پایین	ی، ی
۶	سه نقطه بالا	ث، ث، ث، ژ، ژ، ژ، ش، ش، ش، ش
۷	سه نقطه پایین	پ، پ، پ، پ، پ، چ، چ، چ، چ، چ
۸	سرکش کاف	ک، ک، ک، ک
۹	سرکش گاف	گ، گ، گ، گ

۳-۲-۲- آموزش خصیصه‌ها بر اساس برچسب‌های تعبیه‌شده

همان‌طور که قبلاً ذکر شد، هیستوگرام‌های هرمی می‌تواند به‌عنوان برچسب‌های خصیصه‌ها عمل کنند. این خصیصه‌ها به خوبی توانایی متمایز کردن کلمات از هم را دارا هستند. اگر این خصیصه‌ها بر اساس داده‌هایی که از منابع و یا نویسندگان متفاوت جمع‌آوری شده است آموزش داده شود، مدل‌هایی که به دست می‌آیند نسبت به تغییرات ظاهری و دستخط مقاوم خواهند بود.

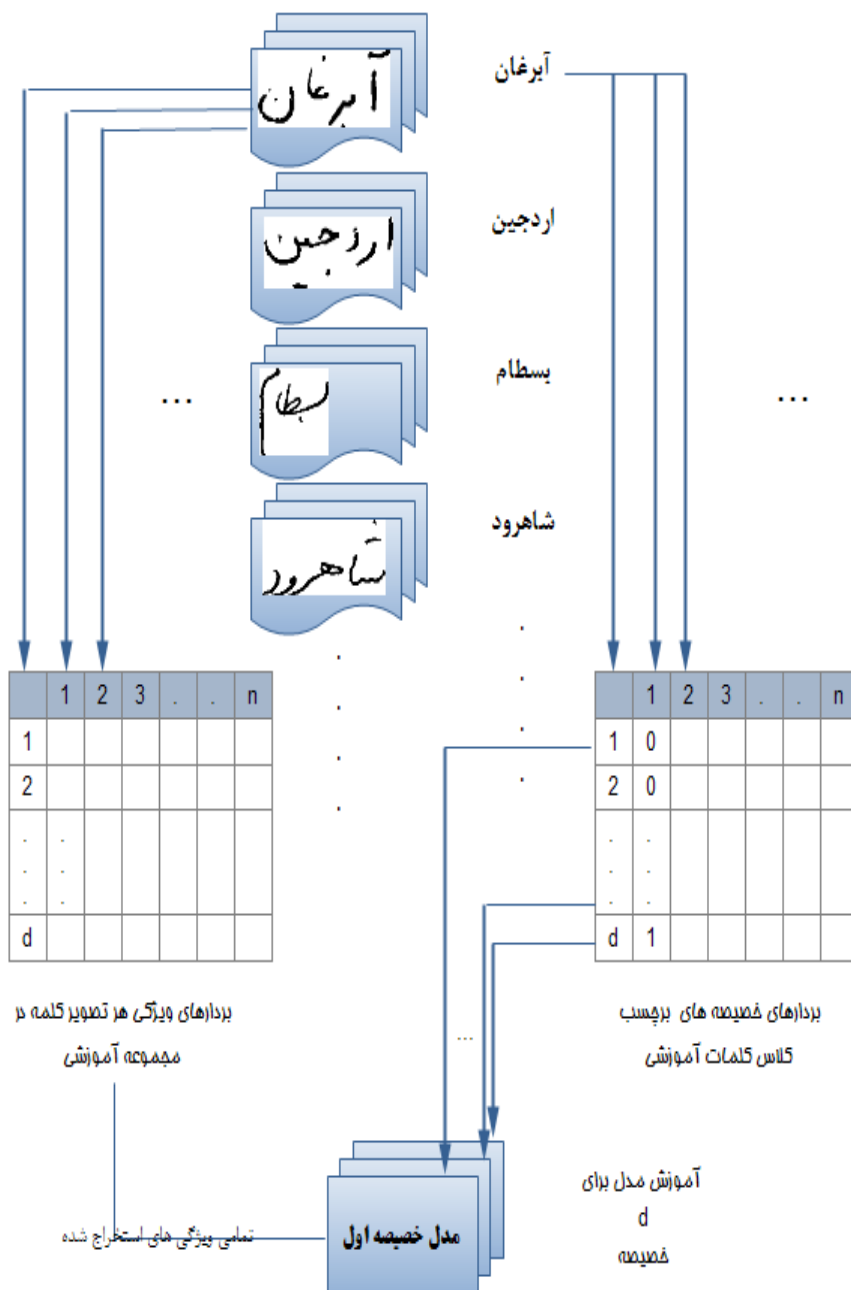
در روش پیشنهادی، برای آموزش این خصیصه‌ها از بردارهای ماشین پشتیبان خطی استفاده شده است. ابتدا بردار ویژگی تصاویر کلمات استخراج می‌شود. این بردارهای ویژگی به همراه برچسب‌های تعبیه‌شده، برای مثال PHOCها، برای آموزش مدل‌های خصیصه بر اساس ماشین بردار پشتیبان استفاده می‌شوند. ساختار کلی روش آموزش در شکل ۳-۶ آمده است.

¹ Persian pyramid histogram of positional signs (PPHOPS)

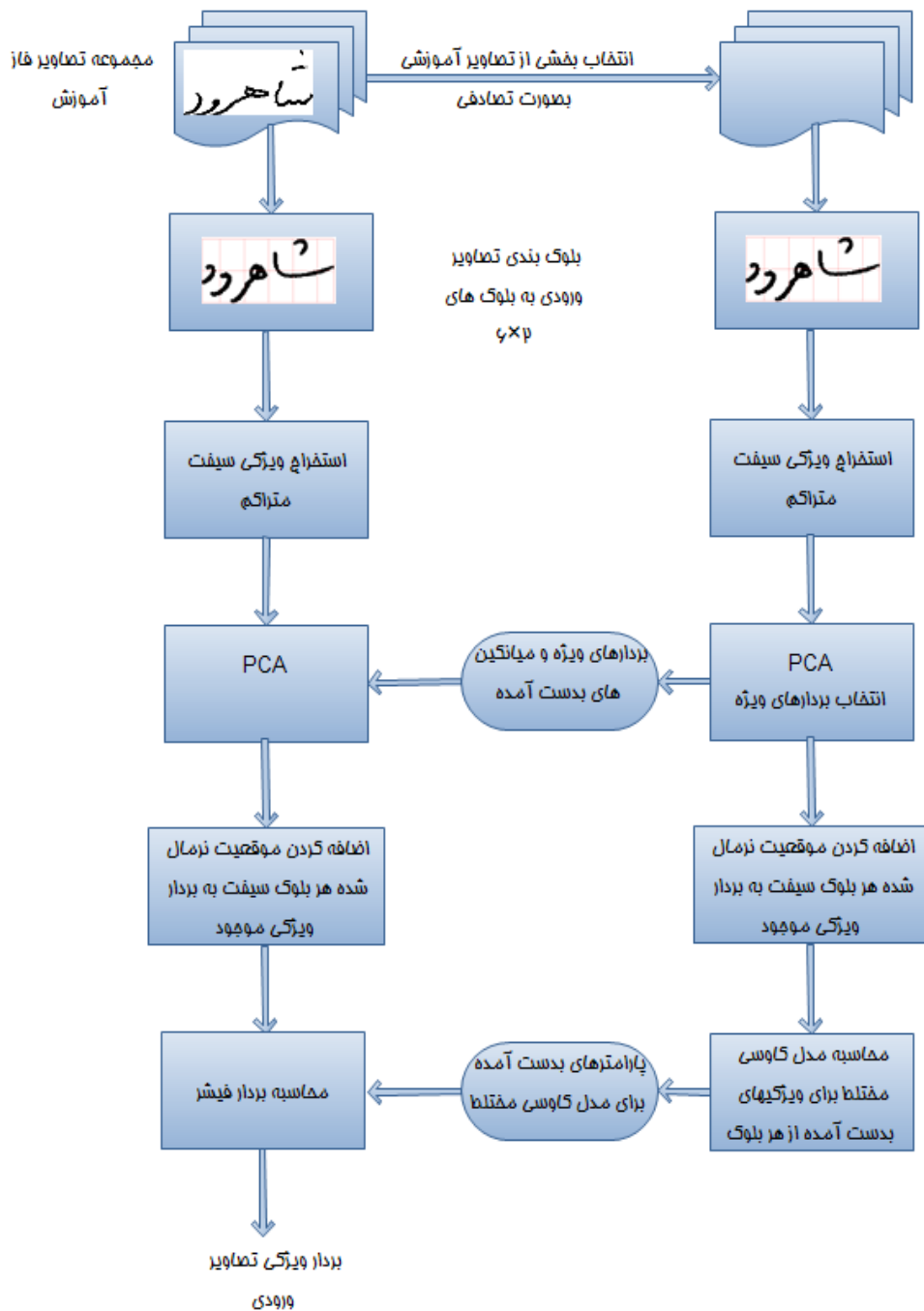
جدول ۲-۳: حروف دارای بدنه مشترک

ردیف	بدنه‌های مشترک	ردیف	بدنه‌های مشترک	ردیف	بدنه‌های مشترک
۱	ا، آ	۲	ب، پ، ت، ث	۳	ب، پ، ت، ث، ذ، ی
۴	پ، پ، ت، ش، ذ، پ	۵	ب، پ، ت، ث	۶	ج، چ، ح، خ
۷	ج، چ، ح، خ	۸	ج، چ، ح، خ	۹	ج، چ، ح، خ
۱۰	د، ذ	۱۱	د، ذ	۱۲	ی
۱۳	ر، ز، ژ	۱۴	ر، ز، ژ	۱۵	س، ش
۱۶	س، ش	۱۷	س، ش	۱۸	س، ش
۱۹	ص، ض	۲۰	ص، ض	۲۱	ص، ض
۲۲	ص، ض	۲۳	ط، ظ	۲۴	ط، ظ
۲۵	ع، غ	۲۶	ع، غ	۲۷	ع، غ
۲۸	ع، غ	۲۹	ف	۳۰	ف
۳۱	ق	۳۲	ق	۳۳	ف، ق
۳۴	ف، ق	۳۵	ک، گ	۳۶	ک، گ
۳۷	ک، گ	۳۸	ک، گ	۳۹	ل
۴۰	ل	۴۱	ل	۴۲	ل
۴۳	م	۴۴	م	۴۵	م
۴۶	ن	۴۷	ن	۴۸	ن
۴۹	ه	۵۰	ه	۵۱	ه
۵۲	و	۵۳	و	۵۴	و
۵۵	ی	۵۶	ا، آ		

برچسب کلاس کلمات نمونه های کلاس کلمات



شکل ۳-۶: ساختار کلی آموزش مدل خصیصه ها

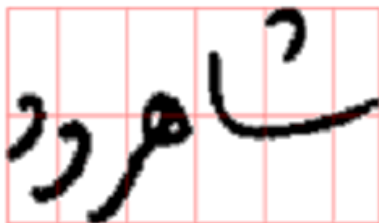


شکل ۳-۷: روند کلی استخراج ویژگی در روش پیشنهادی

برای استخراج ویژگی‌های تصاویر کلمات از بردار فیشر استفاده شده است. بردار فیشر یک روش کدگذاری است که در مسائل بازیابی و بازشناسی تصویر استفاده می‌شود. بردار فیشر می‌تواند به صورت سبد بصری کلماتی در نظر گرفته شود که نه تنها تعداد کلمات، بلکه آماره‌های سطح بالاتر را نیز کد

می‌کند. به‌طور خلاصه در زمان آموزش توصیفگرهای سطح پایین (در این کار سیفت متراکم استفاده شده است) استخراج می‌شود. این ویژگی‌ها جهت آموزش مدل‌های مخلوط گاوسی^۱ با تعداد k هسته مورد استفاده قرار می‌گیرند. روند کلی استخراج ویژگی در شکل ۳-۷ آمده است.

یکی از مسائل مهم در رابطه با روش‌های کدگذاری تصویر مانند بردار فیشر این است که نمی‌توانند مکان ویژگی‌ها که برای توصیف کلمات بسیار مهم هستند، را حفظ کنند. یکی از رهیافت‌های مناسب جهت اضافه کردن اطلاعات مکانی، استفاده از هرم‌های مکانی^۲ است [۷۴]. به جای استفاده از ویژگی‌های کل تصویر، ابتدا تصویر به k ناحیه مختلف تقسیم و ویژگی‌های این نواحی در کنار هم قرار داده می‌شود تا یک بردار ویژگی کلی برای تصویر کلمه تشکیل شود. مشکلی که در این حالت وجود خواهد داشت، طول متفاوت کلمات است (شکل ۳-۸ را ببینید). راهکار دیگری که این مشکل را حل می‌کند و با بردار فیشر هم به‌خوبی کار می‌کند در مرجع [۷۵] پیشنهاد شده است.



شکل ۳-۸: یک نمونه از قطعه‌بندی تصویر کلمه (کلمه "شاهرو")

با اضافه کردن مختصات طول و عرض نرمال شده، توصیفگرهای سیفت تصاویر غنی‌تر می‌شوند. مدل مخلوط گاوسی فقط روی ویژگی‌های سیفت آموزش داده نمی‌شود بلکه از نسخه غنی‌شده آن استفاده می‌کند. به همراه توصیفگر بردار فیشر تصویر مکان ویژگی‌های درون کلمه نیز کد می‌شود. مختصات عمودی و افقی بین $-۰,۵$ تا $+۰,۵$ نرمال می‌شوند. تصویر کلمه در کادر محیطی قرار می‌گیرد (شکل ۳-۸) که مرکز این کادر مبدأ مختصات و حداقل و حداکثر مقادیر طول و عرض $-۰,۵$ و $+۰,۵$ خواهد بود. برای این که نتایج این بازنمایی بهتر شود به جای آموزش مدل مخلوط گاوسی با استفاده از ویژگی‌های غنی‌شده کل تصویر، این کار به‌صورت ناحیه‌ای انجام شده است. یعنی برای هر ناحیه از تصویر یک مدل مخلوط گاوسی از بردار ویژگی‌ها به دست می‌آید و مدل‌های به‌دست‌آمده با هم ترکیب شده و وزن‌های آن‌ها مجدداً به یک نرمال‌سازی می‌شوند.

¹ Gaussian Mixture Model (GMM)

² Spatial Pyramids

۳-۲-۳- زیر فضای مشترک بین برچسب‌ها و خصیصه‌ها

در بخش قبلی یک بازنمایی مبتنی بر خصیصه برای تصاویر کلمات معرفی شد. هرچند که این بازنمایی نسبت به تغییرات ظاهری مقاوم است، اما در هنگام مقایسه کلمات مختلف نیاز به مراقبت زیادی دارد. دلیل این مسئله را می‌توان در امتیازات خصیصه‌ها برای کلمات مختلف جستجو کرد. هرچند به خاطر مشابهت موجود، مقایسه مستقیم خصیصه‌های تعبیه‌شده تصاویر و برچسب‌های متنی تعبیه‌شده نتایج قابل قبولی را نشان می‌دهد اما خوش‌تعریف نیست. باینکه هر دو در یک فضا با ابعاد یکسان قرار می‌گیرند اما برچسب‌های متنی تعبیه‌شده دودویی هستند درحالی‌که امتیازات خصیصه‌های به این صورت نبوده و محدوده‌های متفاوتی دارند. بنابراین کالیبره کردن این امتیازات و برچسب‌های تعبیه‌شده ضروری به نظر می‌رسد.

یکی از رایج‌ترین روش‌ها جهت کالیبره کردن امتیازات ماشین بردار پشتیبان، مقیاس‌بندی پلاتس^۱ است که برای به دست آوردن احتمالات کالیبره شده یک تابع سیگموئید روی امتیازات خروجی استفاده می‌شود:

$$P(y = 1|s) = \frac{1}{1 + e^{\alpha s + \beta}} \quad (۳-۳)$$

که α و β می‌تواند با استفاده از تخمین حداکثر درست‌نمایی^۲ به دست آید.

یکی از مشکلات عمده این‌گونه روش‌ها در نظر نگرفتن همبستگی بین خصیصه‌های مختلف می‌باشد. از آنجایی که به دلیل شیوه کدگذاری چند سطحی، همبستگی زیادی بین خصیصه‌های پیشنهادی وجود دارد، حذف همبستگی بین خصیصه‌ها حائز اهمیت بسیاری است. به این دلیل از کالیبراسیون توأم که همبستگی بین خصیصه‌ها را بهتر می‌تواند استخراج کند، استفاده شده است. چون رگرسیون عادی همبستگی بین مقادیر خود بردارهای خصیصه را در نظر نمی‌گیرد، رگرسیون زیرفضای مشترک که منجر به فرمول‌بندی معادل تحلیل همبستگی کانونی^۳ می‌شود پیشنهاد شده است.

فرض کنید $I = \{I_n, n = 1, \dots, N\}$ مجموعه N تصویر در مجموعه آموزشی و $y = \{y_n, n = 1, \dots, N\}$ برچسب‌های مربوط به آن‌ها باشند. در این صورت $A = \phi_I(I) \in \mathbb{R}^{d \times N}$ نشان دهنده N تصویر تعبیه‌شده در فضای d بعدی خصیصه‌ها است و N برچسب تعبیه‌شده در فضای d بعدی به صورت

¹ Platts Scaling

² Maximum Likelihood Estimation (MLE)

³ Canonical Correlation Analysis (CCA)

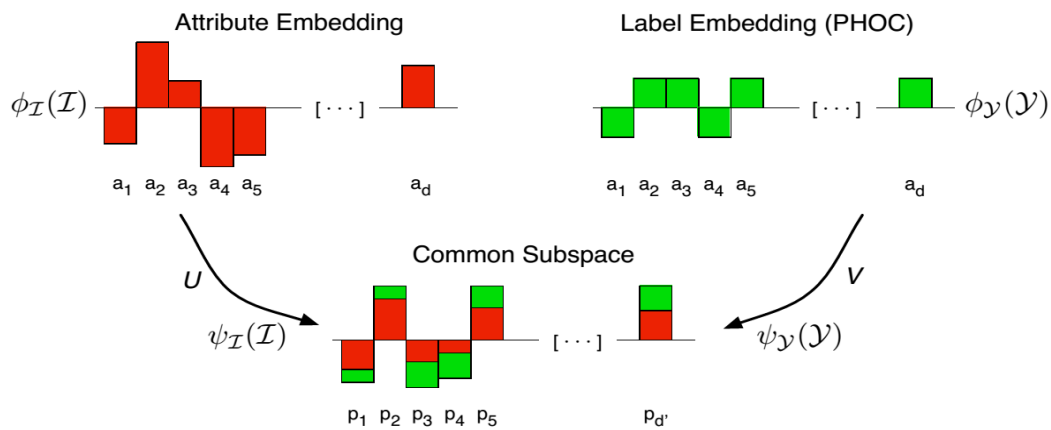
فاصله به صورت زیر است: $B = \phi_y(y) \in \{0,1\}^{d \times N}$ می باشد. بنابراین روش مستقیم برای ارتباط دادن A و B تعریف یک تابع

$$F(I_i, y_i; P) = \|P^T \phi_I(I_i) - \phi_y(y_i)\|_2^2, \quad P \in \mathfrak{R}^{d \times d} \quad (۴-۳)$$

برای کمینه کردن فاصله بین تمام نمونه ها و برجسب های آن ها می توان از رابطه زیر استفاده نمود:

$$\operatorname{argmin}_P \sum_i \frac{1}{2} F(I_i, y_i; P) + \frac{1}{2} \Omega(P) = \operatorname{argmin}_P \frac{1}{2} \|P^T A - B\|_F^2 + \frac{1}{2} \Omega(P) \quad (۵-۳)$$

که $\Omega(P) = \alpha \|P\|_F^2$ عبارت منظم سازی و α وزن آن می باشد. در این مورد برابر یک مسئله رگرسیون معمولی است و $P = (AA^T + \alpha I)^{-1} AB^T$ که I ماتریس همانی است. چون d تعداد خصیصه ها می باشد، که زیاد هم نیستند و این مسئله فقط یک بار در زمان آموزش باید حل شود، حل آن بسیار سریع است. همان طور که قبلاً هم ذکر شد، این فرمولاسیون فقط همبستگی بین امتیازات خصیصه ها را در نظر می گیرد لذا همبستگی بین خود خصیصه لحاظ نمی شود. برای این که هر دو مسئله در نظر گرفته شود تغییراتی در آن انجام شده است [۷] که هر دو را در یک زیرفضای مشترک با ابعاد d' نگاشت می کند (شکل ۹-۳).



شکل ۹-۳: نگاشت امتیازات پیش بینی شده خصیصه ها و برجسب های تعیین شده در یک زیرفضای دارای همبستگی بیشتر [۷]

یک تابع فاصله جدید به صورت زیر تعریف شده است:

$$\hat{F}(I_i, y_i; U, V) = \|\psi_I(I_i) - \psi_y(y_i)\|_2^2 \quad (۶-۳)$$

که در آن $\psi_I(I) = U^T \phi_I(I)$ و $\psi_y(y) = V^T \phi_y(y)$ دو تابع تعبیه خطی هستند و از ماتریس‌های نگاشت $U, V \in \mathbb{R}^{d \times d'}$ برای تعبیه $\phi_I(I)$ و $\phi_y(y)$ در زیرفضای مشترک استفاده می‌کند. سپس، مشابه با مورد قبلی، هدف کمینه کردن فاصله بین همه نمونه‌ها و برجسب‌های آن‌ها است.

$$\begin{aligned} \operatorname{argmin}_{U,V} \sum_i \frac{1}{2} \hat{F}(I_i, y_i; U, V) + \frac{1}{2} \Omega(U) + \frac{1}{2} \Omega(V) = & \quad (7-3) \\ \operatorname{argmin}_{U,V} \frac{1}{2} \|U^T A - V^T B\|_F^2 + \frac{1}{2} \Omega(U) + \frac{1}{2} \Omega(V) & \end{aligned}$$

$$\begin{aligned} \psi_I(I) \psi_I(I)^T &= I \\ \psi_y(y) \psi_y(y)^T &= I \end{aligned}$$

که در آن محدودیت تعامد اطمینان می‌دهد که حل کردن آن آسان نخواهد بود.

با استفاده از ضرایب لاگرانژ نسبت به U و V مشتق گرفته شده و برابر صفر قرار داده می‌شود و به تساوی‌های زیر می‌رسیم.

$$\lambda(AA^T + \alpha I)u_k = AB^T v_k \quad (8-3)$$

$$\lambda(BB^T + \alpha I)v_k = BA^T u_k \quad (9-3)$$

که u_k و v_k ، k امین ستون ماتریس U و V است و λ به خاطر ضرایب لاگرانژ ظاهر می‌شود. وقتی u^k را حل می‌کنیم به مسئله مدار ویژه تعمیم‌یافته زیر می‌رسیم:

$$AB^T(BB^T + \alpha I)^{-1}BA^T u_k = \lambda^2(AA^T + \alpha I)u_k \quad (10-3)$$

K بردار ویژه تعمیم‌یافته ابتدایی از K ستون ماتریس نگاشت U می‌باشد که اجازه می‌دهد تا ابعاد نهایی d' انتخاب شود. فرآیندی مشابه برای به دست آوردن ماتریس نگاشت V نیز انجام می‌شود. این مدل با استفاده بیشتر از ارتباط بین داده‌ها مدل بهتری نسبت به رگرسیون در این کار ایجاد می‌کند. علاوه بر این، در این مدل می‌توان ابعاد خروجی را کنترل نمود و کاهش ابعاد را انجام داد.

همچنین این روابط برای حل مسئله تحلیل همبستگی کانونی نیز مورد استفاده قرار می‌گیرند، زیرا تلاش در یافتن نگاشتی است که همبستگی در یک زیرفضای مشترک را حداکثر نماید. تحلیل همبستگی کانونی ابزاری است برای استفاده از اطلاعاتی که از منابع مختلف در دسترس قرار می‌گیرند و برای کاربردهای بازیابی [۷۶] و خوشه‌بندی [۷۷] مورد استفاده قرار می‌گیرند. علاوه بر این تحلیل همبستگی

کانونی برای مرتبط کردن توصیفات تصاویر و برچسب‌های آن نیز مورد استفاده قرار گرفته است که برای کاربردهای بازیابی مزایای زیادی دارد. به این دلیل که ممکن است رابطه بین امتیازات خصیصه‌ها و خصیصه‌های دودویی خطی نباشد، مدل مبتنی بر هسته (تحلیل همبستگی کانونی دارای هسته^۱) می‌تواند منجر به بهینه‌سازی بیشتری شود.

یکی از مشکلات آموزش فضای خصیصه و زیرفضای مشترک در دو گام مجزا نیاز به حجم زیاد داده‌های آموزشی است. به این دلیل که داده‌هایی که برای آموزش زیرفضای مشترک استفاده می‌شوند باید متفاوت با داده‌هایی باشند که برای آموزش فضای خصیصه مورد استفاده قرار می‌گیرند. به خاطر اینکه اگر برای هر دو مرحله داده یکسان استفاده کنیم امتیازات ماشین بردار پشتیبان بیش برآزش^۲ خواهند شد (اکثراً بسیار نزدیک به ۱ یا ۱- می‌شوند) و در نتیجه زیرفضای آموزش دیده با این شرایط به شدت بایاس می‌شوند و منجر به نتایج ضعیف خواهند شد. اگر مجموعه داده بزرگی داشته باشیم می‌توان آن را به دو بخش تقسیم کرد و بخشی را برای آموزش زیرفضای مشترک و بخشی را جهت آموزش فضای خصیصه مورد استفاده قرار داد. اما با این روش داده‌های آموزشی به‌خوبی مورد استفاده قرار نمی‌گیرند.

جهت غلبه بر مشکل ذکرشده، گونه‌ای از بگینگ^۳ استفاده شده است. برای این کار داده‌های آموزشی به چندین پوشه‌ی^۴ آموزش و ارزیابی تقسیم‌بندی می‌شود، به‌طوری‌که داده‌های ارزیابی و آموزش در هر دسته از همدیگر مجزا هستند اما ممکن است در دسته‌های متفاوت اشتراک داشته باشند. برای هر دسته یک مدل با استفاده از داده‌های مربوط به آن آموزش داده می‌شود و این مدل برای محاسبه امتیاز روی داده‌های ارزیابی استفاده می‌شود. بنابراین امتیازات روی داده‌های ارزیابی تقریباً بایاس شده نیستند. در دسته‌های مختلف امتیازات ارزیابی جمع می‌شوند و برای هر نمونه تعداد دفعاتی که استفاده شده است ثبت می‌گردد. در انتهای فرآیند مدل کلی با میانگین‌گیری از همه این مدل‌های محلی آماده می‌شود. با نرمال‌سازی امتیازات هر نمونه بر اساس تعداد باری که استفاده شده است مجموعه آموزشی از امتیازاتی که بایاس نیستند فراهم می‌شود که می‌تواند برای آموزش بدون مشکل زیرفضای مشترک مورد استفاده قرار گیرد (الگوریتم ۳-۱).

الگوریتم ۳-۱: آموزش مدل خصیصه‌ها با استفاده از بگینگ

¹ Kernelized Canonical Correlation Analysis (KCCA)

² Overfit

³ bagging

⁴ fold

Input: Training data $X \in \mathbb{R}^{D \times N}$
Input: Training labels $Y \in \{0, 1\}^N$
Input: Number of folds F
Output: Model $W \in \mathbb{R}^D$
Output: Training data embedded onto the attribute space $A \in \mathbb{R}^N$

```

W = zeros(1, D)
A = zeros(1, N)
count = zeros(1, N)
f = 1
while f ≤ F do
  Split data in train and val partitions
  TrainIdx, ValIdx = split(N, f)
  TrainData = X(:, TrainIdx)
  TrainLabels = Y(TrainIdx)
  ValData = X(:, ValIdx)
  ValLabels = Y(ValIdx)
  Learn model using training data. Use validation set to validate the parameters.
  Wf = learnSVM(TrainData, TrainLabels, ValData, ValLabels)
  Encode the validation set into the attributes space and keep track of the number of updates
  A(ValIdx) = A(ValIdx) + WfTValData
  count(ValIdx) = count(ValIdx) + 1
  Add Wf to the global model W
  W = W + Wf
  f = f + 1
end while
Normalize and end
W = W/F
A = A/count
End

```

۳-۳- نتایج تجربی و ارزیابی

با توجه به روند کلی مطرح شده که در شکل ۳-۷ نیز آمده است، برای آموزش مدل مخلوط گاوسی، ۲۵۰ تصویر به صورت تصادفی و غیرتکراری از تصاویر پایگاه داده انتخاب می‌شوند. سپس تصاویر کلمات طوری چیده می‌شوند که هر کلمه توسط کادری محاط شود. این کادر به ۲ ناحیه افقی و ۶ ناحیه عمودی برابر تقسیم می‌شود (شکل ۳-۸) و توصیفگرهای سیفت متراکم از هر کدام از ۱۲ بلوک تصویر استخراج می‌شود. بردار ویژگی با ابعاد $128 \times n$ برای هر بلوک از تصویر ایجاد می‌شود. در واقع n برابر با مجموع تمام نقاطی است که در آن ویژگی استخراج شده است. با استفاده از آنالیز اجزاء اصلی تعداد

ویژگی‌ها به ۶۲ ویژگی کاهش داده می‌شود و سپس با اضافه کردن مختصات طول و عرض نرمال شده، غنی می‌شوند و برداری به طول ۶۴ را ایجاد می‌کنند. برای هر کدام از نواحی ۱۲ گانه مدل مخلوط گاوسی با استفاده از این بردارهای ویژگی غنی شده آموزش داده می‌شود. در هر مدل مخلوط گاوسی از ۱۶ مدل گاوسی استفاده شده است، در نتیجه $16 \times 12 = 192$ دسته پارامتر خواهیم داشت. سپس بردار فشر استخراج می‌شود که این بردار طولی برابر $2 \times 64 \times 12 \times 16 = 24576$ خواهد داشت.

با استفاده از روش تعبیه برچسب، برچسب‌های مجموعه تصاویر آموزشی در فضای هیستوگرام هر می تعبیه می‌شود. در هر تکرار ۸۰ درصد از داده‌ها برای آموزش و ۲۰ درصد باقیمانده برای اعتبارسنجی مورد استفاده قرار می‌گیرد. آموزش توسط ماشین بردار پشتیبان توسط روش بگینگی که قبلاً شرح داده شد، انجام می‌گیرد. مقادیر اوزان نهایی از میانگین‌گیری نتایج محلی به دست می‌آید و در نهایت برای امتیاز مشخصه‌ها برای داده‌های اعتبارسنجی و نمونه‌های آزمون از رابطه زیر محاسبه می‌گردد.

$$Score = W^T \times FV \quad (11-3)$$

در این رابطه w وزن‌های نهایی به دست آمده و FV بردار ویژگی تصویر (بردار فشر) است. این آموزش برای خصیصه‌هایی که تعداد کمتر از ۲ نمونه دارند انجام نمی‌شود.

برای حل این مشکل از روش تحلیل همبستگی کانونی (CCA) استفاده شده است. با استفاده از این روش رتبه مشخصه‌ها و برچسب‌های PHOC در یک فضای برداری مشترک تعبیه می‌شوند که در این فضا همبستگی بین آن‌ها ماکزیمم است. در این فضا رتبه مشخصه و برچسب PHOC که مربوط به یک کلمه هستند، نزدیک به یکدیگر قرار دارند. به عبارت دیگر در این فضا تصاویر و رشته‌ی متنی که نشان‌دهنده‌ی کلمه یکسان هستند نزدیک به یکدیگر قرار می‌دهند و این اجازه را می‌دهد که عملیات بازشناسی به صورت مسئله نزدیک‌ترین همسایگی انجام شود.

جهت آموزش زیرفضای مشترک از تجزیه و تحلیل همبستگی کانونی استفاده شده است، زیرا ابزاری است که توصیفگرهای تصویر و برچسب‌های متناظرشان را به هم مرتبط می‌کند. در این روش با تجزیه و تحلیل هم‌زمان مجموعه‌ها و با شناسایی و مشخص کردن عناصری از یک مجموعه متغیر با بیشترین وابستگی و ارتباط با عناصر مجموعه متغیر دیگر، روابط آماری مستقل موجود بین دو مجموعه متغیر را مورد آزمون قرار می‌دهد. در این کار هدف از این مرحله دستیابی به دو ماتریس نگاشت Ux و Vy است که توسط این دو ماتریس، داده‌ها از فضای مشخصه و برچسب به فضای برداری مشترک نگاشت می‌شوند.

اگر X مجموعه‌ی N تصویری باشد که در فضای خصیصه تعبیه شده‌اند و Y مجموعه‌ی N برچسب تعبیه‌شده در فضای برچسب باشد و میانگین آن‌ها μ_x و μ_y باشد، ماتریس‌های C_{xx} ، C_{yy} ، C_{yx} و C_{xy} توسط روابط زیر به دست می‌آیند.

$$C_{xx} = \frac{1}{N} (X - \mu_x)(X - \mu_x)^T + \rho I \quad (12-3)$$

$$C_{yy} = \frac{1}{N} (Y - \mu_y)(Y - \mu_y)^T + \rho I \quad (13-3)$$

$$C_{xy} = \frac{1}{N} (X - \mu_x)(Y - \mu_y)^T \quad (14-3)$$

$$C_{yx} = (C_{xy})^T \quad (15-3)$$

در این روابط ρ ضریب تنظیم است و I ماتریس همانی است. هدف یافتن ماتریس‌های U_x و V_y طبق رابطه‌ی زیر است.

$$\operatorname{argmax}_{U_x, V_y} \frac{U_x^T C_{xy} V_y}{\sqrt{U_x^T C_{xx} U_x} \sqrt{V_y^T C_{yy} V_y}} \quad (16-3)$$

با استفاده از روابط زیر ماتریس U_x که معادل با بردارهای ویژه‌ی ماتریس Z است به دست می‌آید. ماتریس V_y نیز به روش مشابهی به دست می‌آید.

$$Z U_{xk} = \lambda_k^2 U_{xk} \quad , \quad Z = C_{xx}^{-1} C_{xy} C_{yy}^{-1} C_{yx} \quad (17-3)$$

$$U_x = \{U_{x1}, U_{x2}, \dots, U_{xk}\} \quad , \quad U_{xi} \in R^d$$

توسط این دو ماتریس، رتبه مشخصه‌ها و برچسب‌ها از فضای d بُعدی در فضای k بُعدی تعبیه می‌شوند. درواقع این روش اجازه می‌دهد که ابعاد خروجی را کنترل نموده و کاهش ابعاد را نیز انجام دهیم.

حل معادلات بالا و به دست آوردن دو ماتریس U_x و V_y ، یک‌بار به صورت برون‌خط در مرحله‌ی آموزش انجام می‌شود. در زمان آزمون یک نمونه ($S \in R^d$) با توجه به اینکه در فضای مشخصه قرار دارد یا فضای برچسب، توسط $U_x^T (S - \mu_x)$ یا $V_y^T (S - \mu_y)$ در فضای برداری مشترک تعبیه می‌شود.

با استفاده از روند ذکرشده، آزمایش‌ها روی دو مجموعه داده از کلمات مجزای دست‌نویس فارسی، فارسی و ایرانشهر، انجام شده است. عمل بازیابی و بازشناسی برای هرکدام از دو مجموعه مورد بررسی قرار گرفته است. سپس جهت ارزیابی توانایی روش ذکرشده در بازیابی و بازشناسی کلمات خارج از فرهنگ

لغت، برای آموزش از یکی از این مجموعه‌ها و جهت آزمون از مجموعه دیگری استفاده شده است که هیچ کلاس کلمه و نویسنده مشترکی ندارند.

۳-۳-۱- معیارهای ارزیابی

معیارهای ارزیابی مختلفی جهت بررسی و مقایسه نتایج به دست آمده مورد استفاده قرار گرفته است. روش جستجوی کلمه در دسته روش‌های ارزیابی قرار می‌گیرد، بنابراین از معیارهای ارزیابی روش‌های ارزیابی اطلاعات جهت ارزیابی روش جستجوی کلمه استفاده می‌شود. معیارهای مختلفی جهت ارزیابی روش‌های ارزیابی اطلاعات معرفی شده است که مهم‌ترین معیار، معیار میانگین دقت متوسط^۱ است. میانگین دقت متوسط برای مجموعه‌ای از کلمات پرسشی، میانگین میزان دقت متوسط برای هر کلمه پرسشی است که به صورت زیر محاسبه می‌شود:

$$Map = \frac{\sum_{q=1}^Q AveP(q)}{Q} \quad (18-3)$$

که Q تعداد کلمات پرسشی است. $AveP$ دقت متوسط را به ازای هر کلمه پرسشی محاسبه می‌کند. از آنجایی که دقت^۲ و فراخوانی^۳ معیارهایی تک مقداری بر اساس لیست تمامی کلمات ارزیابی شده است، برای سیستمی که یک سری مرتب شده از کلمات ارزیابی شده را باز می‌گرداند، باید ترتیب ارزیابی نیز مدنظر قرار داده شود. با محاسبه دقت و فراخوانی در هر موقعیت لیست مرتب ارزیابی شده فقط می‌توان منحنی دقت-فراخوانی^۴ را رسم نمود. در این منحنی دقت $P(r)$ به عنوان تابعی از فراخوانی رسم می‌شود. دقت متوسط میانگین مقادیر $P(r)$ را روی فاصله ۰ تا ۱ محاسبه می‌کند.

$$AveP = \int_0^1 P(r) dr \quad (19-3)$$

این رابطه مساحت زیر نمودار دقت-فراخوانی را محاسبه می‌کند. در عمل انتگرال با محاسبه مجموع در هر موقعیت از لیست مرتب جایگزین می‌شود:

$$AveP = \sum_{k=1}^n P(k) \Delta r(k) \quad (20-3)$$

¹ Mean average precision (Map)

² Precision

³ recall

⁴ Precision-recall curve

که k موقعیت در لیست بازیابی شده و n تعداد اسناد بازیابی شده را نشان می‌دهد. $P(k)$ دقت در k عنصر اول لیست بازیابی شده و $\Delta r(k)$ تغییر در فراخوانی از $k-1$ به k است. این مجموع برابر است با:

$$AveP = \frac{\sum_{k=1}^n (P(k) \times rel(k))}{\text{number of relevant retrieved sample}} \quad (21-3)$$

تابع $rel(k)$ در صورتی که در موقعیت k ام از لیست بازیابی شده نمونه صحیحی وجود داشته باشد برابر ۱ و در غیر این صورت برابر ۰ خواهد بود.

معیار ارزیابی دیگری که در جستجوی کلمه مورد استفاده قرار می‌گیرد معیار دقت در اولین بازیابی^۱ است. این معیار نرخ صحت اولین مورد بازیابی شده را مشخص می‌کند.

معیار ارزیابی به کاررفته در سیستم بازشناسی کلمات دست‌نویس، معیار دقت است که نسبت تعداد کلمه دست‌نویسی که درست شناسایی شده است را به کل نمونه‌های صحیح آن کلمه محاسبه می‌کند.

۳-۳-۲- نتایج تجربی

نتایج ارزیابی با استفاده از معیارهای معرفی‌شده در جداول زیر آمده است. دقت بازشناسی کلمات دست‌نویس در جدول ۳-۳ آورده شده است. در این جدول روش‌های پیشنهادی با آخرین روش‌های معرفی‌شده مقایسه شده است. انجام مقایسه با ۴ حالت مختلف انجام شده است. در حالت اول ۷۰٪ از مجموعه فارسی جهت آموزش مدل و ۳۰٪ جهت آزمون اختصاص داده‌شده و نتایج با استفاده از خصیصه‌های پیشنهادی نشان داده شده است. در حالت دوم بر همین اساس روی داده‌های مجموعه داده ایران‌شهر دقت محاسبه شده است. نتایج تجربی نشان می‌دهد که روش‌های پیشنهادی که ساختار نوشتاری زبان فارسی را در نظر می‌گیرد دقت بهتری نسبت به روش پایه گزارش می‌کند. با توجه به این که در این دو روش ساختار نگارشی زبان مدنظر قرار داده شده است، نتایج حاصله نشان می‌دهد که خصیصه‌های پیشنهادی توانایی اشتراک‌گذاری اطلاعات بین کلاسی را بهتر از روش پایه انجام می‌دهند.

¹ Precision at 1 (P@1)

جدول ۳-۳: نتایج دقت روش پیشنهادی جهت بازشناسی کلمه

روش	نرخ تشخیص			
	فارسا	ایران شهر	آموزش با فارسا آزمون با ایران شهر	آموزش با ایران شهر آزمون با فارسا
<i>PPHOC</i>	۹۵/۴۴	۹۵/۸۱	۲۷/۵۳	۴۰/۹۲
<i>PPHOPC</i>	۹۶/۰۷	۹۷/۵۵	۲۸/۴۲	۳۹/۳۷
<i>PPHOPS</i>	۹۶/۱۵	۹۷/۰۸	۲۸/۵۰	۴۲/۶۸
<i>Gradient based[78]</i>	-	۸۹/۱۲	-	-
<i>Decision Templates[79]</i>	-	۹۱/۵۵	-	-
<i>Mixture of Experts[80]</i>	۹۱/۱۱	-	-	-
<i>SOM+DHMM[81]</i>	۹۵/۶۳	۹۶/۹۴	-	-
<i>Divide and conquer[82]</i>	۹۰/۵۰	-	-	-

در بخش سوم و چهارم ارزیابی توانایی شناسایی کلاس‌های دیده نشده مورد ارزیابی قرار گرفته است. یکی از توانایی‌های روش‌هایی که از رده‌بندی مبتنی بر خصیصه استفاده می‌کنند، توانایی شناسایی کلاس‌های دیده نشده است. به همین دلیل در این دو بخش آموزش و آزمون مدل با استفاده از دو مجموعه داده مجزا انجام گرفت. لازم به ذکر است که این دو مجموعه داده تعداد کلاس و نویسندگان مختلفی دارند. نتایج تجربی نشان می‌دهد که توانایی ارزیابی کلاس‌های دیده نشده در سیستم پیشنهادی وجود دارد. دقت‌ها در موردی که مدل با استفاده از مجموعه داده ایران شهر آموزش داده شده بالاتر است. دلیل این امر را می‌توان در تعداد بیشتر کلاس‌های این مجموعه و در نتیجه آموزش بهتر مدل دانست.

جدول ۳-۴ ارزیابی را جهت سیستم جستجوی کلمه مبتنی بر معیار میانگین دقت متوسط نشان می‌دهد. در این جدول هم مشابه جدول ۳-۳ با ۴ روش مختلف مقایسه انجام شده است. در این ارزیابی‌ها نیز

بهترین نتایج در مورد خصیصه‌هایی گزارش شده است که ساختار نگارشی زبان را در نظر گرفته‌اند. علاوه بر این، نتایج تجربی توانایی روش در جستجوی کلمات دیده نشده را نشان می‌دهد.

جدول ۳-۴: نتایج میانگین دقت متوسط (MaP) روش پیشنهادی جهت جستجوی کلمه

مجموعه داده	نوع خصیصه	کلمه پرسشی تصویری				کلمه پرسشی متنی			
		Platt Scaling	CCA	KCCA	مقایسه مستقیم	Platt Scaling	CCA	KCCA	مقایسه مستقیم
فارسا	PPHOC	۸۵/۱۵	۸۲/۲۹	۸۶/۸۶	۷۸/۴۰	۹۴/۳۲	۹۳/۹۰	۹۴/۹۸	۸۱/۸۰
	PPHOPC	۸۸/۰۵	۸۳/۰۴	۸۸/۳۵	۷۶/۸۵	۹۵/۸۶	۹۴/۶۸	۹۵/۵۷	۸۴/۲۱
	PPHOPS	۸۶/۸۰	۸۴/۱۹	۸۸/۹۹	۷۷/۹۵	۹۴/۶۳	۹۴/۵۶	۹۵/۶۷	۸۵/۲۲
ایرانشهر	PPHOC	۹۱/۸۲	۸۹/۵۱	۹۱/۵۶	۸۸/۰۲	۹۶/۸۳	۹۵/۴۹	۹۶/۲۱	۸۲/۸۱
	PPHOPC	۹۱/۸۱	۹۰/۱۲	۹۲/۰۸	۸۸/۱۰	۹۶/۹۰	۹۵/۷۰	۹۶/۲۰	۷۸/۲۸
	PPHOPS	۹۱/۲۹	۸۹/۹۹	۹۲/۰۱	۸۷/۴۱	۹۶/۰۱	۹۵/۵۵	۹۶/۲۵	۸۵/۶۴
آموزش با فارسا آزمون با ایرانشهر	PPHOC	۴۴/۶۵	۸۲/۸۸	۸۱/۹۲	۷۰/۱۱	۴۲/۲۵	۴۷/۴۵	۴۷/۳۴	۳۹/۵۶
	PPHOPC	۴۳/۶۰	۸۴/۱۰	۸۳/۸۸	۷۴/۴۲	۴۲/۶۹	۴۸/۴۳	۴۷/۷۸	۴۱/۵۶
	PPHOPS	۴۷/۳۷	۸۳/۷۷	۸۳/۱۳	۷۱/۰۹	۴۲/۲۸	۴۷/۱۶	۴۷/۳۲	۳۹/۵۰
آموزش با ایرانشهر آزمون با فارسا	PPHOC	۳۷/۱۰	۴۸/۷۱	۸۵/۰۱	۸۰/۸۰	۶۸/۲۲	۸۲/۱۶	۸۳/۷۵	۶۸/۲۲
	PPHOPC	۵۹/۹۱	۸۴/۹۹	۸۵/۸۱	۸۱/۴۰	۷۰/۴۶	۸۲/۲۵	۸۱/۰۸	۶۸/۰۸
	PPHOPS	۶۸/۹۲	۸۶/۶۳	۸۵/۶۵	۸۲/۵۰	۷۷/۷۲	۸۷/۸۸	۸۸/۵۳	۷۱/۷۶

جدول ۳-۵: نتایج دقت در اولین بازیابی (P@1) روش پیشنهادی جهت جستجوی کلمه

مجموعه داده	نوع خصیصه	کلمه پرسشی تصویری				کلمه پرسشی متنی			
		Platt Scaling	CCA	KCCA	مقایسه مستقیم	Platt Scaling	CCA	KCCA	مقایسه مستقیم
فارسا	PPHOC	۹۴/۷۴	۹۸/۹۳	۹۸/۸۸	۹۸/۷۰	۱۰۰	۱۰۰	۱۰۰	۹۸/۶۷
	PPHOPC	۹۴/۲۷	۹۸/۹۳	۹۸/۹۴	۹۸/۵۸	۱۰۰	۱۰۰	۱۰۰	۹۹/۳۳
	PPHOPS	۹۶	۹۸/۵۸	۹۸/۸۲	۹۸/۸۴	۱۰۰	۱۰۰	۱۰۰	۹۸/۶۷
ایران شهر	PPHOC	۹۲/۷۱	۸۹/۴۰	۹۳/۸۴	۸۹/۶۶	۹۸/۸۱	۹۹	۹۹	۹۰/۲۴
	PPHOPC	۹۳/۹۷	۹۱/۸۵	۹۴/۳۰	۸۹/۳۳	۹۸/۸۱	۹۹/۲۰	۹۸/۲۱	۹۱/۴۳
	PPHOPS	۹۴/۳۷	۹۲/۳۱	۹۴/۴۳	۸۹/۴۶	۹۸/۸۱	۹۹/۴۰	۹۸/۴۱	۸۸/۸۴
آموزش با فارسا آزمون با ایران شهر	PPHOC	۸۰/۷۸	۹۹/۱۰	۹۸/۸۸	۹۸/۳۲	۴۸/۷۵	۴۸/۷۵	۴۸/۷۵	۴۸/۷۵
	PPHOPC	۷۶/۰۲	۹۹/۱۵	۹۹/۱۸	۹۸/۷۸	۵۰	۵۰	۴۸/۵۰	۴۸/۵۰
	PPHOPS	۸۳/۳۸	۹۸/۹۸	۹۸/۹۸	۹۸/۵۰	۵۰	۵۰	۴۷/۵۰	۴۷/۵۰
آموزش با ایران شهر آزمون با فارسا	PPHOC	۹۸/۸۷	۹۹/۳۳	۹۹/۳۲	۹۹/۲۵	۸۹/۰۱	۷۰/۰۷	۹۲/۵۰	۹۳/۷۵
	PPHOPC	۸۵/۲۰	۹۹/۳۸	۹۹/۳۳	۹۹/۳۲	۸۸/۷۵	۷۱/۲۵	۹۱/۲۵	۹۰
	PPHOPS	۹۲/۴۵	۹۹/۲۸	۹۹/۲۵	۹۹/۲۲	۹۲/۵۰	۹۷/۵۰	۹۶/۲۵	۸۲/۵۰

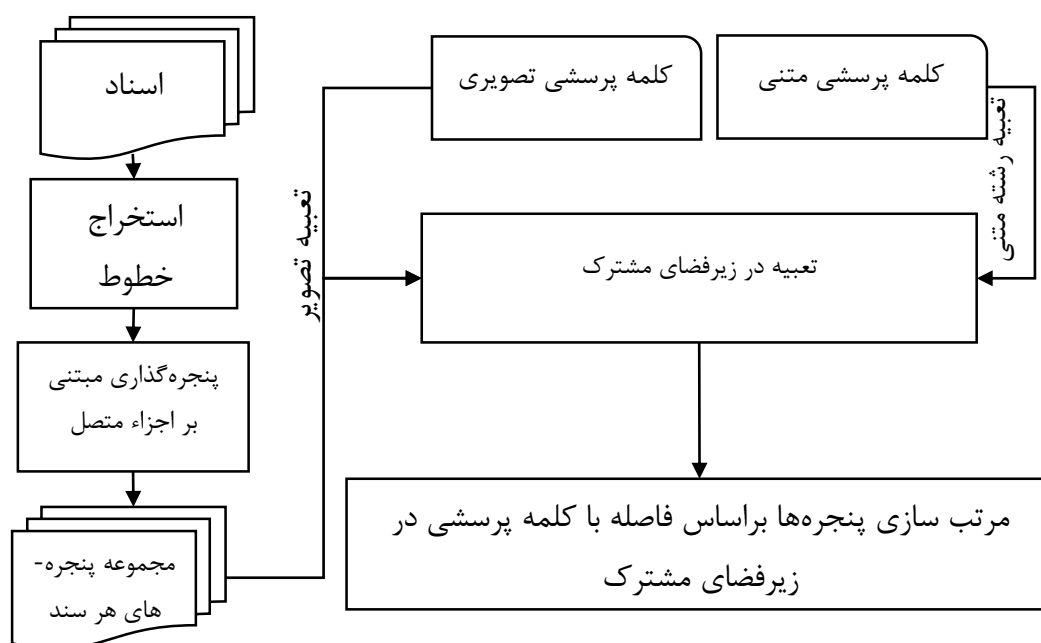
با توجه به نتایج تجربی گزارش شده در جداول بالا، روش‌های پیشنهادی توانایی بهتری در بازنمایی خصیصه‌های اسناد زبان فارسی دارد. علاوه بر این نشان داده شده که پیچیدگی و تنوع در شیوه نگارش و کلمات استفاده شده در پایگاه داده ایران شهر بسیار بیشتر از پایگاه داده فارسا است. با توجه به این تنوع در مواردی که آموزش با استفاده از پایگاه داده ایران شهر انجام شده و آزمون روی مجموعه داده فارسا است، دقت روش بسیار بالاتر است. دلیل این امر مدل کردن بهتر شیوه‌های نگارش متفاوت توسط نمونه‌های آموزشی این مجموعه می‌باشد.

فصل

۴- روش پیشنهادی دوم (مبتنی بر خط)

۴-۱- مقدمه

با توجه به مطالب ذکر شده در فصل‌های قبل، از آنجایی که قطعه‌بندی کلمات و استخراج دقیق محدوده کلمات در نوشته‌های فارسی خصوصاً در مورد اسناد دست‌نویس، معمولاً با خطای بالای قطعه‌بندی روبرو است، محققان به دنبال روش‌هایی بوده‌اند که از این خطا رهایی یابند. یکی از این راهکارها استفاده از قطعه‌بندی در سطح خط است. در این فصل روشی مبتنی بر قطعه‌بندی در سطح خط جهت جستجوی کلمه معرفی شده است. در این روش عمل قطعه‌بندی خطوط با استفاده از روش پیشنهادی انجام می‌شود و بعد از آن بر اساس یک روش پنجره‌گذاری پیشنهادی که مبتنی بر اجزاء متصل است عمل پنجره‌گذاری انجام می‌شود. مجموعه‌ای از پنجره‌ها که شامل تعداد مختلفی جزء متصل می‌باشند آماده می‌شود. این مجموعه بر اساس روش تعبیه برچسب در فضای آموزش دیده شده توسط مجموعه داده‌های کلمات مجزا تعبیه شده و به ازای کلمه‌ی پرسشی ورودی مرتب می‌شوند (شکل ۴-۱).



شکل ۴-۱: ساختار کلی روش پیشنهادی مبتنی بر قطعه‌بندی خط

روند انجام تعبیه برچسب و تعبیه تصاویر بر اساس روش شرح داده شده در فصل سوم انجام می‌گیرد. بنابراین در ادامه این بخش روش پیشنهادی جهت قطعه‌بندی خط و پنجره‌گذاری و نحوه استفاده از پنجره‌ها جهت کاهش تعداد مقایسه شرح داده شده است.

۴-۲- قطعه‌بندی خطوط

استخراج یا قطعه‌بندی خط یکی از مراحل مهم پیش‌پردازش در بسیاری از کاربردهای آنالیز اسناد تصویری مانند شناسایی متون چاپی و دست‌نویس برون خط، جستجوی کلمات، بازیابی کلمات تصویری و شناسایی کاراکتر نوری می‌باشد. با توجه به اینکه در اسناد دست‌نویس فاصله خطوط و کلمات در یک سند متغیر است و همچنین معمولاً راستای مستقیم برای خط مبنا رعایت نشده و خطوط یک سند دارای زوایای متفاوتی نسبت به خط مبنا هستند، استخراج خطوط در اسناد دست‌نویس بسیار چالش‌برانگیز است. یکی از مشکلات بزرگ دیگر در این زمینه امکان به هم چسبیدگی اجزاء متصل از دو خط متوالی است؛ که با توجه به وجود نقاط و سرکش‌ها در متون فارسی و عربی این مسئله بیشتر مشکل ایجاد می‌نماید. استخراج صحیح یا غلط خطوط تأثیر مستقیم زیادی بر دقت قطعه‌بندی کلمات و حروف دارد.

روش‌های متفاوتی برای استخراج و قطعه‌بندی خطوط از اسناد دست‌نویس ارائه شده است. روش‌های موجود در استخراج خط در مقالات به روش‌های مختلفی دسته‌بندی شده است که به صورت زیر می‌باشند: روش‌های مبتنی بر افکنش، روش‌های مبتنی بر لکه‌دار کردن^۱، روش‌های مبتنی بر تبدیل هاف^۲، روش مبتنی بر نازک‌سازی^۳ و روش‌های احتمالی^۴.

روش‌های مبتنی بر افکنش معمولاً در مورد اسناد تصویری چاپی مورد استفاده قرار می‌گیرد و برای اسناد دست‌نویسی که خطوط در آن همپوشانی ندارند نیز توسعه یافته است [۸۳، ۸۴]. اما این روش به علت مشکلات فراوان مانند در نظر نگرفتن شیب خط، در متون اسناد دست‌نویس تصویری، نمی‌تواند بر روی اسناد مستقیماً اعمال شود. به همین دلیل جهت بهبود این روش راهکارهایی مانند اعمال روش‌های مبتنی بر افکنش به برش‌های عمودی در سند و سپس ترکیب نتایج، پیشنهاد شده است [۸۵-۸۸]. در این روش سند تصویری به صورت موازی به نوارهای عمودی تقسیم‌بندی می‌شود. سپس مرزهای بالقوه جداکننده خطوط در هر بخش به دست می‌آید و در مرحله بعد این بخش‌های کاندید با استفاده از هیوریستیک یا مدل‌های گرافیکی آماری مانند مدل مخفی مارکوف به یکدیگر متصل شده و نواحی خطوط را تشکیل می‌دهند. نهایتاً برای فائق آمدن بر مواردی که خطوط همپوشانی دارند از تکنیک‌های تقاطع^۵، الگوریتم تعقیب کانتور و چگالی گاوسی دومتغیره برای بازنمایی خطوط استفاده شده است. روش مذکور دارای ضعف‌های زیادی می‌باشد که به چند مورد اصلی آن اشاره شده است: این روش خطوط بالقوه زیادی را تولید می‌کند، پارامتر عرض نوارهای برش باید از پیش تعریف شده باشد، خطوط

¹ Smearing Methods

² Hough Transform

³ Methods based on thinning operation

⁴ Stochastic Methods

⁵ Crossing technique

نباید انحنای زیاد داشته باشند و اگر نوار ابتدایی و انتهایی شامل خطوط بالقوه نباشند، استخراج خط به‌طور کامل ناممکن است.

برخی شیب خط را نیز برای جداسازی خطوط مورد استفاده قرار داده‌اند [۸۵] که این شیب معمولاً به سادگی قابل محاسبه نیست و در یک سند بین خطوط متفاوت است. بنابراین این روش‌ها نیز در این شرایط نتایج خوبی را ارائه نمی‌کنند.

در روش‌های مبتنی بر لکه‌دار کردن، روش‌هایی مانند الگوریتم هموارسازی طول اجرا [۸۹] یا الگوریتم آنالیز طول اجرای جهت‌دار بهبودیافته^۲ [۹۰] بکار گرفته شده است. استفاده از این الگوریتم‌ها باعث خواهد شد که بلوک‌های پیکسلی در جهت افقی لکه‌دار شوند؛ یعنی اگر فاصله سفید (پس‌زمینه) بین آن‌ها از یک حد آستانه‌ای کمتر بود با پیکسل سیاه (پیش‌زمینه) پر شود. این روش ابتدا روی تصاویر خاکستری اعمال شده و سپس تصویر باینری می‌شود و عمل مورفولوژی فرسایش روی پیش‌زمینه و پس‌زمینه برای به دست آوردن محدوده خطوط انجام شده است.

روش‌های مبتنی بر رهیافت تبدیل هاف برای استخراج خطوط مستقیم در تصاویر بسیار پرکاربرد است. شیب خطوط دست‌نویس با استفاده از اعمال این نوع تبدیل به مرکز ثقل هر جزء متصل به دست می‌آید. اگر اکثر همسایه‌های نزدیک در یک تراز متعلق به گروه اجزاء تشکیل‌دهنده هم‌ترازی داشته باشند، این هم‌ترازی ویژگی‌های پیوستگی و نزدیکی را دارا خواهد بود و به‌عنوان یک خط پذیرفته می‌شود. سپس نتایج حاصل شده با حذف هم‌ترازی‌های نادرست از هم‌ترازی بین اجزاء متصل با استفاده از اطلاعات محتوایی، تصویر خروجی نهایی را ارائه می‌کند [۹۱]. در [۹۲] روشی مبتنی بر تبدیل هاف بر روی سند بلوک‌بندی شده ارائه شده است.

برخی محققین برای استخراج خطوط هندی و چینی از روش نازک‌سازی استفاده کرده‌اند، که در [۹۳] ابتدا نازک‌سازی صورت گرفته و در ادامه عملیات پس‌پردازشی روی تمامی نواحی پس‌زمینه برای تشخیص مرز جداکننده خطوط اعمال شده است.

اخیراً برخی از روش‌های مجموعه سطوح^۳، کانتور فعال^۴ و بیز تغییرات^۵ برای استخراج خطوط بهره برده‌اند. در [۹۴] از روش تخمین چگالی و مجموعه سطوح استفاده شده است. در این روش یک نگاشت احتمالی از سند تصویری ورودی تخمین زده می‌شود که هر المان احتمال تعلق هر پیکسل به یک خط را نشان می‌دهد. روش مجموعه سطوح برای تعیین ارزیابی مرزی با خطوط همسایه مورد استفاده قرار

¹ Run-Length Smoothing Algorithm (RLSA)

² Improved Directional Run-Length Analysis

³ Level Set

⁴ Active Contour

⁵ Variational Bayes

گرفته است. در [۹۵] ابتدا یک بانک فیلتر تطبیق یافته برای هموارسازی تصویر ورودی مورد استفاده قرار گرفته است. مراکز خطوط اجزاء تشکیل دهنده خط با استفاده از مرزبندی تصویر هموارشده محاسبه می شود. در انتها برای به دست آوردن نتیجه نهایی، کانتورهای فعال روی مرزها وفق داده می شوند. در کار [۹۶] روش بیز تغییرات برای استخراج خطوط در متون دستنویس چینی مورد بهره برداری قرار گرفته است. یک سند به صورت یک مدل مخلوط گاوسی در نظر گرفته می شود، طوری که هر جزء متعلق به یک خط باشد. برای تخمین پارامترهای چگالی و تعیین تعداد خطوط چارچوب کاری بیز تغییراتی بکار برده شده است.

در [۹۷، ۹۸] یک روش مبتنی بر بلوک بندی معرفی شده که پارامترهای آن به صورت افقی بر اساس سند ورودی تعیین می گردد. در هر بلاک برای تشخیص جهت از فیلتر کردن تصویر با هسته های گاوسی دوبعدی ناهمسانگرد^۱ در زوایای بین ۳۰- درجه تا ۳۰ درجه استفاده شده است. پس از تشخیص جهت، خطوط در هر بلاک از هم جدا می شوند.

اسناد دستنویس فارسی و عربی دارای تعداد زیاد و متنوعی نقطه، سرکش و نشانه هستند. همچنین در این نوع اسناد یک جزء متصل می تواند یک کلمه، زیر کلمه، کاراکتر، نقطه یا بخشی از نقاط پیوسته یا یک سرکش باشد. با توجه به دلایل مذکور مسئله استخراج خطوط دستنویس برای اکثر روش ها چالش برانگیز خواهد بود و برخی از روش ها مانند تبدیل هاف، لکه دار کردن طول اجرای جهت دار بهبود یافته و کانتور فعال در مواجهه با این نوع اسناد شکست خواهند خورد. علت این شکست ها را می توان به صورت زیر تشریح کرد: ۱- روش های مبتنی بر تبدیل هاف اهمیت بدنه اصلی کلمه و نقاط و سرکش های کوچک را در فضای هاف یکسان در نظر می گیرد. ۲- در روش هایی که از آنالیز طول اجزا استفاده می شود ممکن است اجزاء کوچک برخی خطوط مانند نقطه و سرکش به خط دیگری متصل شوند. ۳- در برخی روش ها با استفاده از فیلتری نقاط و سرکش ها حذف می شوند، در حالی که این اجزاء در زبان فارسی و عربی نقش تعیین کننده ای دارند.

در [۹۹] روشی مبتنی بر تکنیک نقاشی^۲ ارائه شده است. در این روش یک سند به نوارهای عمودی تقسیم بندی می شود و با استفاده از تکنیک های نقاشی اجزاء نزدیک به همدیگر در هر نوار عمودی به هم متصل می شوند که در این صورت ادعا شده که شیب خط به دست می آید. از عمل مورفولوژی انبساط جهت برقراری اتصال بین لکه های ایجاد شده در مرحله قبل و تشکیل خط استفاده می شود. سپس روی پس زمینه نتیجه مرحله قبل عمل نازک سازی انجام می گیرد و با استفاده از برخی قواعد هیوربستیک خطوط نهایی استخراج می شوند. این روش با توجه به اینکه از روش های قبلی برای زبان های

¹ Anisotropic 2D Gaussian Kernels

² Painting Technique

فارسی و عربی بهتر کار می‌کند، اما باز هم در مورد این دو زبان درصد خطاهای بسیار بیشتری نسبت به لاتین دارد. لذا در روش ارائه‌شده سعی شده تا این مشکلات را با استفاده از خوشه‌بندی سلسله‌مراتبی با استفاده از یک معیار فاصله معرفی‌شده مرتفع‌سازییم.

روش پیشنهادشده یک روش جزء به کل می‌باشد که از خوشه‌بندی سلسله‌مراتبی پیکسل‌های قلم (سیاه یا پیش‌زمینه) موجود در سند جهت استخراج خطوط استفاده می‌کند. ساختار کلی روش پیشنهادی، که در شکل ۴-۲ نشان داده شده است، شیوه نگارش متن فارسی به صورت دست‌نویس را مدنظر قرار داده است. از آنجایی که معمولاً دست‌نوشته‌ها ممکن است شامل خطوط نزدیک به همدیگر باشند یا نقاط و سرکش‌ها باعث اشتباه در خوشه‌بندی شوند، از یک مرحله پیش‌پردازش استفاده شده که احتمال خوشه‌بندی اشتباه را در موارد ذکرشده کاهش دهد.

با توجه به اینکه در متون فارسی و عربی شیب خط مبنای استاندارد مستقیم است و در متون دست‌نویس این شیب می‌تواند به مقدار کمی خط مبنا را به سمت بالا یا پایین متمایل کند یا در آن انحنای کمی را ایجاد نماید، معیار فاصله‌ای معرفی‌شده که این مسئله را در نظر بگیرد. نتایج نهایی خوشه‌بندی با استفاده از تخمین شیب خط و هیوریستیک‌هایی مبتنی بر شیوه نگارش فارسی اصلاح و تکمیل می‌گردد. این روش در ادامه به‌طور کامل تشریح و نتایج آزمایش‌ها روی دو مجموعه استاندارد زبان فارسی گزارش شده است.

۴-۲-۱- پیش‌پردازش

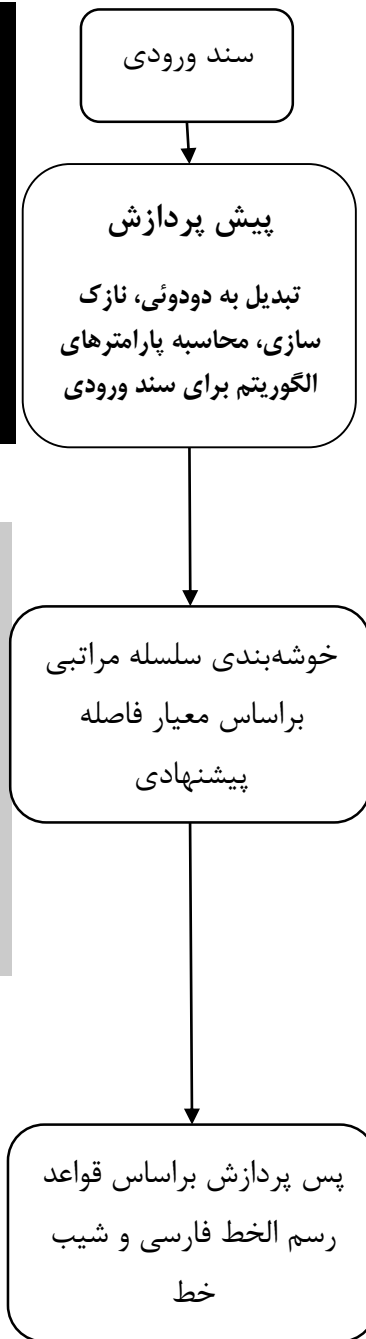
مرحله پیش‌پردازش یکی از مهم‌ترین مراحل در پردازش تصاویر دیجیتال می‌باشد. در سیستم‌های پردازش اسناد تصویری در اکثر روش‌ها اولین گام پیش‌پردازش تبدیل سند تصویری به یک تصویر دودویی است. به‌عنوان اولین مرحله پیش‌پردازش، در این روش با استفاده از الگوریتم اوتسو^۱ [۱۰۰] سند ورودی به تصویر دودویی تبدیل می‌گردد. با توجه به اینکه در متون دست‌نویس فارسی، بالارونده‌هایی مانند حروف الف و لام و پایین‌رونده‌هایی مانند ر، ز و واو و سرکش‌ها ممکن است در یک فاصله افقی کم از همدیگر قرار گیرند، استفاده از نازک‌سازی می‌تواند امکان تداخل این موارد را کاهش دهد. بنابراین در این روش از عمل نازک‌سازی بر پایه مورفولوژی استفاده شده است. با توجه به اینکه در اسناد ورودی تعداد خطوط، اندازه تصویر و نوع دستخط از نظر اندازه متفاوت است نیاز به تعیین برخی پارامترها برای هر سند وجود دارد تا روش بتواند خود را بر اساس شرایط سند ورودی تطبیق دهد.

¹ Otsu

اهل کائنات من روزگارم به نیست که نانی دارم خردم هوشی بر سر سوزن ذوق، مادری دارم
 مبراز بزرگ درخت. در کسائی بهتر از آب روان و خندانم که در این نزدیک است
 لایق است که بوی آبل کاخ بلند، من مسلمانم. قبله ام یک گل سرخ
 همانرا چشیده هم از نور. در کسائی من و صوبای من و شهرهای من
 در تمام جهان دارم در جهان دارم در کسائی من و دانه که در این است
 چیزی است که بزرگ است. چرا در نفس من کسائی که کسائی کل سید
 چه کنم از لاله قمر دارم. چشم ما را با به نیست جور دیگر باید کارها
 مناسبتی را از کل سرخ.

اهل کائنات من روزگارم به نیست که نانی دارم خردم هوشی بر سر سوزن ذوق، مادری دارم
 مبراز بزرگ درخت. در کسائی بهتر از آب روان و خندانم که در این نزدیک است
 لایق است که بوی آبل کاخ بلند، من مسلمانم. قبله ام یک گل سرخ
 همانرا چشیده هم از نور. در کسائی من و صوبای من و شهرهای من
 در تمام جهان دارم در جهان دارم در کسائی من و دانه که در این است
 چیزی است که بزرگ است. چرا در نفس من کسائی که کسائی کل سید
 چه کنم از لاله قمر دارم. چشم ما را با به نیست جور دیگر باید کارها
 مناسبتی را از کل سرخ.

اهل کائنات من روزگارم به نیست که نانی دارم خردم هوشی بر سر سوزن ذوق، مادری دارم
 مبراز بزرگ درخت. در کسائی بهتر از آب روان و خندانم که در این نزدیک است
 لایق است که بوی آبل کاخ بلند، من مسلمانم. قبله ام یک گل سرخ
 همانرا چشیده هم از نور. در کسائی من و صوبای من و شهرهای من
 در تمام جهان دارم در جهان دارم در کسائی من و دانه که در این است
 چیزی است که بزرگ است. چرا در نفس من کسائی که کسائی کل سید
 چه کنم از لاله قمر دارم. چشم ما را با به نیست جور دیگر باید کارها
 مناسبتی را از کل سرخ.



شکل ۴-۲: ساختار کلی روش پیشنهادی جهت قطعه بندی خطوط

یکی از مهم ترین پارامترهایی که طی عملیات خوشه بندی سلسله مراتبی پیشنهادی به آن نیاز داریم، حد آستانه ای جهت توقف ترکیب خوشه ها و استخراج خوشه بندی نهایی می باشد. برای شرط توقف

اتصال خوشه‌ها از معیار حداکثر فاصله استفاده می‌شود و میزان این حد آستانه از پارامتر فاصله مؤثر بین دو خط مجاور به دست می‌آید. برای به دست آوردن این پارامتر، تصویر سند به ۲۰ نوار عمودی با عرض برابر تقسیم می‌شود و سپس متوسط فاصله خطوط مجاور برای ۱۴ نوار میانی اندازه گرفته می‌شود. بدین صورت که برای هر کدام از نوارهای مذکور پروفایل افکنش افقی محاسبه شده و از آن فاصله بین قله‌های پروفایل استخراج می‌گردد. شکل ۳-۴ تصویری از این بلوک‌بندی و قله‌های مذکور در یکی از نوارها را نمایش می‌دهد. اگر این فاصله برای تمامی نوارهای مذکور محاسبه شود، فاصله مؤثر بین دو خط مجاور از رابطه زیر به دست می‌آید:

$$Line_{De} = \frac{\sum_{i \in Corr_D} i * Hist(i)}{\sum_{i \in Corr_D} Hist(i)} \quad (۱-۴)$$

$$Corr_D = \{j | Hist(j) > 0.25 * \max(Hist)\} \quad (۲-۴)$$

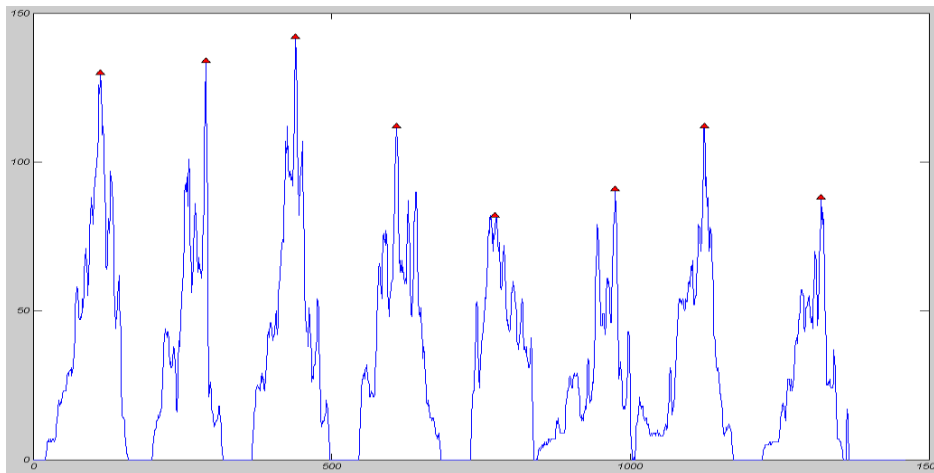
که در این روابط Corr_D شامل اندیس فاصله‌هایی می‌شود که تعداد آن‌ها حداقل از ۲۵ درصد حداکثر بیشتر باشد. Hist هیستوگرام فاصله بین خطوط به دست آمده را نشان می‌دهد. Line_De نشان‌دهنده فاصله مؤثر بین خطوط است. نصف این فاصله به عنوان پارامتر معیار توقف خوشه‌بندی استفاده شده است.

پارامتر دسته‌بندی کردن خوشه‌ها به دو دسته‌ی کوچک و بزرگ به صورت زیر تعریف می‌شود. تعداد پیکسل‌های قلم (مساحت) برای هر خوشه محاسبه شده و سپس مرتب می‌شوند. فاصله بین مساحت‌های مرتب شده خوشه‌ها به دست آمده و بر اساس رابطه زیر نرمال می‌شود. مساحت شروع بزرگ‌ترین فاصله مقدار پارامتر حد آستانه جداسازی خوشه‌ها را تعیین می‌کند.

$$Dist_{Area}(i) = \frac{sortedArea(i+1) - sortedArea(i)}{sortedArea(i+1) + sortedArea(i)} \quad (۳-۴)$$

اهل کائنات من روزگاری به نیست گنجه نانی دانه خرد هوشم سرسوزن ذوقی، مادری دایم
 بهتر از برگ درخت. دوستانی بهتر از آب روان و خدایی که در این نزدیکی است.
 لای لای لب بوهای آل کالج بلند. من مسلمانم. قبله ام یک گل سرخ
 جانان خسته. مرموز. دشت کاه من. من و موی با پیش و پنبه های گریه
 در تمام جویان دارد و می بیند دار کعبه. من می دانم که برای تو بند است هر دو
 نجیبی است کبوتر زیباست. چرا در قفس هم نفس کمر کس نیست گل سبزه
 چو که از لاله فرزند دارد. چشم ما را باید بست چو دیگر دیده کارمانست
 سنی سانی راز گل سرخ.

(الف)



(ب)

شکل ۴-۳: الف) سند قطعه بندی شده به نوارهای عمودی ب) پروفایل افکنش افقی نوار مشخص شده با علامت ستاره سبز در تصویر الف به همراه نمایش قله ها.

۴-۲-۲- خوشه بندی

جهت استخراج خطوط از روش خوشه بندی سلسله مراتبی کمترین فاصله^۱ استفاده شده است. برای این کار موقعیت افقی و عمودی تمامی نقاط قلم (پیکسل های سیاه) به عنوان ویژگی های نمونه ها مورد استفاده قرار گرفته است. برای این که عملیات خوشه بندی سریع تر انجام شود، هر جزء متصل به عنوان یک خوشه اولیه در نظر گرفته شده و عمل خوشه بندی انجام می گیرد. در این روش با توجه به ساختار

¹ Single linkage

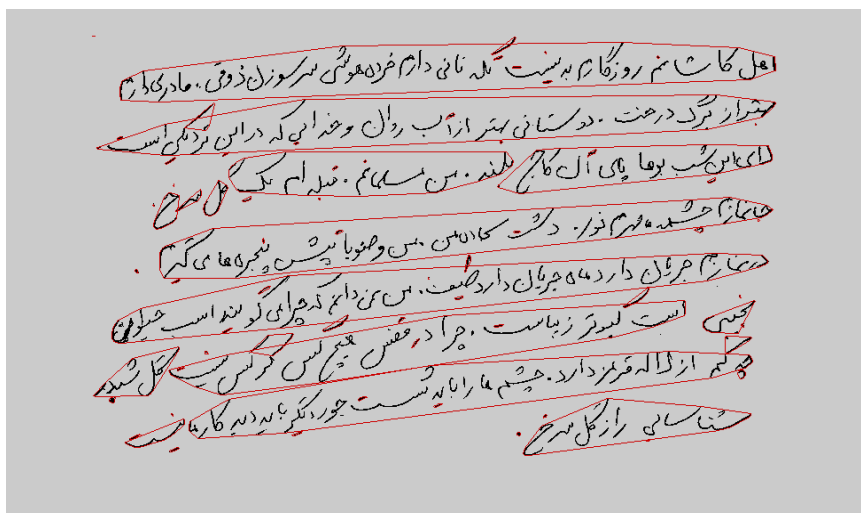
نوشتاری زبان فارسی که به صورت افقی نوشته می‌شود، فاصله طبق رابطه زیر طوری تعریف شده است که فاصله در جهت عمودی شامل جریمه باشد.

$$D(p_1, p_2) = |X_{p_1} - X_{p_2}| + |Y_{p_1} - Y_{p_2}|^5 \quad (4-4)$$

که در آن p_1 و p_2 دو نمونه از پیکسل‌های مشکی هستند که فاصله بین آن‌ها محاسبه می‌شود و X_p و Y_p به ترتیب مختصات افقی و عمودی هر کدام از دو نقطه مذکور می‌باشند. با استفاده از این معیار فاصله، اجزاء متصلی (خوشه‌هایی) که از نظر افقی در کنار هم قرار دارند در یک خوشه قرار خواهند گرفت. با توجه به اینکه در زبان فارسی شبه‌کلمات تشکیل‌دهنده یک خط در امتداد محور عمودی کشیدگی دارند احتمال قرار گرفتن اجزاء متصل کنار هم حتی در شرایطی که خط به صورت شیب‌دار یا منحنی نوشته شده باشد بسیار زیاد است و نتایج به دست آمده این ادعا را ثابت می‌کند.

اما یکی از مسائل مهم در خوشه‌بندی سلسله‌مراتبی معیار توقف ادغام خوشه‌هاست. برای این کار معیار توقف حداکثر فاصله بین دو خوشه در نظر گرفته شده و برای هر سند به صورت وقتی تعیین می‌گردد که قبلاً شرح داده شده است. نتیجه این مرحله با شرایط ذکر شده مجموعه‌ای از خوشه‌هاست که معمولاً بیشتر از تعداد خطوط است و شامل خوشه‌های بزرگی که نمایش‌دهنده و معرف خط اصلی هستند می‌شود. شکل ۴-۴ نمونه‌ای از یک سند تصویری را پس از این مرحله نمایش می‌دهد.

همان‌طور که در شکل ۴-۴ نیز کاملاً مشخص است، معمولاً نقاط و سرکش‌ها که در بسیاری موارد فاصله عمودی آن‌ها تا بدنه اصلی زیاد است، در یک خوشه مجزا قرار می‌گیرند. یکی دیگر از مواردی که در برخی نمونه‌ها رخ می‌دهد تقسیم شدن یک خط واحد به چندین خوشه مجزا می‌باشد؛ این مشکل به خاطر تغییر ناگهانی شیب خط، فاصله زیاد بین اجزاء متصل یا عدم وجود معیار توقف مناسب و دقیق برای هر سند، رخ می‌دهد. برای غلبه بر این مشکلات از پس پردازش استفاده شده است که در قسمت بعدی شرح داده شده است.



شکل ۴-۴: یک سند نمونه خروجی مرحله خوشه‌بندی.

۴-۲-۳- پس پردازش

در این مرحله از تعدادی قاعده هیوریستیک مبتنی بر ساختار نگارشی فارسی و بر اساس نتایج اولیه مرحله خوشه‌بندی برای ترکیب خوشه‌ها و استخراج دقیق خطوط استفاده شده است که در زیر آورده شده است.

- در مرحله اول اندازه هر خوشه به دست می‌آید و خوشه‌ها بر اساس یک حد آستانه، که قبلاً ذکر شده است، به دو دسته‌ی خوشه‌های کوچک و بزرگ دسته‌بندی می‌شوند. در صورتی که ناحیه بدنه محدب یک خوشه کوچک فقط با یک خوشه بزرگ اشتراک بیش از ۵۰ درصد داشته باشد دو خوشه با هم ترکیب می‌شوند. با استفاده از قاعده بیان شده، اکثر خوشه‌های کوچکی که شامل نقاط و سرکش‌ها هستند به خوشه‌های اصلی می‌پیوندند.
- در مرحله بعد برای تصمیم‌گیری در مورد سرکش‌ها زوایای خوشه‌های کوچک محاسبه می‌شود. در صورتی که این زاویه بین ۳۰ تا ۶۰ درجه باشد سرکش تشخیص داده شده و با خوشه بزرگی که زیر آن قرار دارد ترکیب می‌شود.
- در ادامه خط تقریبی مبنا برای خوشه‌های بزرگ محاسبه می‌شود، بدین صورت که به ازای مختصات تمام نقاط قلم موجود در یک خوشه بزرگ یک خط برازش می‌شود. این خطوط استخراج شده شیب احتمالی خطوط اصلی سند را نشان می‌دهد. بر مبنای خطوط احتمالی استخراج شده، قواعدی جهت ترکیب خوشه‌ها در نظر گرفته شده است. اگر فاصله دو خط از خطوط برازش شده بسیار کم باشد نشان می‌دهد که دو خوشه بزرگ روی یک خط قرار گرفته‌اند و این دو خوشه با هم ترکیب می‌شوند. در نهایت فاصله اقلیدسی مراکز خوشه‌های کوچکی که

هنوز مورد تصمیم‌گیری قرار نگرفته‌اند با خطوط برازش شده به خوشه‌های بزرگ محاسبه شده و به خوشه مربوط به نزدیک‌ترین خط می‌پیوندند.

۴-۲-۴- نتایج تجربی

به‌منظور ارزیابی روش پیشنهادی از ۱۰۰ سند از مجموعه داده استاندارد استفاده شده است. مجموعه داده، FHT [۹] می‌باشد که شامل ۱۱۲۶ سند دست‌نویس فارسی است که توسط ۲۵ نویسنده مختلف نوشته شده است و متوسط تعداد سطر در هر سند حدود ۷ می‌باشد.

جهت ارزیابی دو معیار استاندارد مورد استفاده قرار گرفته است. اولین معیار نرخ تشخیص^۱ است [۹۴] که در استخراج خط بسیار پرکاربرد می‌باشد. این معیار نرخ تشخیص صحیح خط در اسناد ورودی را مشخص می‌کند. اگر هر خط استخراج شده معیار تطابقی بالای ۹۵٪ را به دست آورد این خط به‌عنوان یک خط تشخیص داده شده به حساب می‌آید. این معیار تطابق به‌صورت زیر محاسبه می‌گردد.

$$MatchScore_{i,j} = T(G_i \cap R_j) / T(G_i \cup R_j) \quad (۵-۴)$$

که در آن تابع $T(a)$ تعداد پیکسل‌های پیش‌زمینه a را برمی‌گرداند. G_i نشان‌دهنده آامین تصویر از مجموعه درستی^۲ و R_j نشان‌دهنده آامین خط استخراج شده توسط رهیافت پیشنهادی است. این معیار استاندارد از سال ۲۰۰۷ در مسابقات کنفرانس بین‌المللی بازشناسی و آنالیز اسناد^۳ مورد استفاده قرار می‌گیرد. بر اساس این معیار روش پیشنهادی نرخ تشخیصی برابر ۹۲/۷۳٪ دارد.

معیار ارزیابی بعدی که در این مورد ارائه‌شده، معیار نرخ تطابق پیکسلی^۴ است. این معیار بر اساس تقسیم تعداد پیکسل‌های قلم مشترک بین بهترین تطبیق در میان تصاویر مجموعه درستی و خط استخراج شده بر تعداد پیکسل‌های قلم این تصویر مجموعه درستی به دست می‌آید. نرخ تطابق پیکسلی برای مجموعه داده مورد ارزیابی قرار گرفته برابر ۹۷/۵۳٪ می‌باشد.

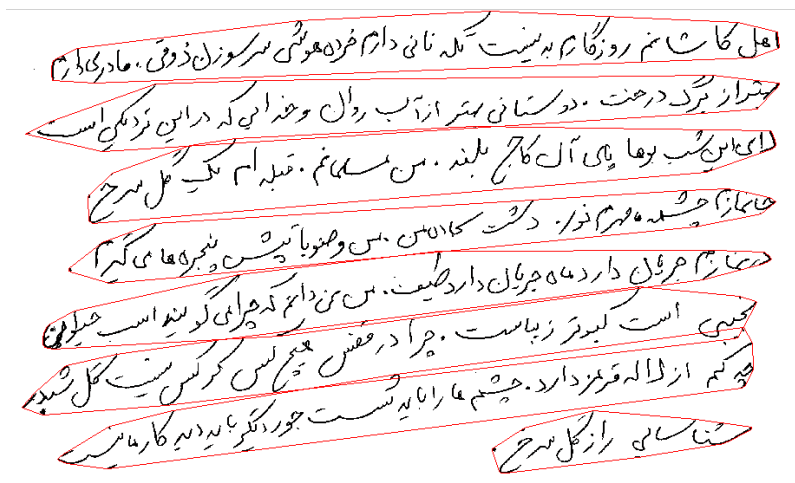
نتایج قابل‌مقایسه‌ای روی این مجموعه داده به دست آمده است که می‌توان به روش مذکور جهت ادامه گام بعدی اعتماد کرد.

¹ Detection Rate (DR)

² Ground Truth

³ International Conference on Document Analysis and Recognition (ICDAR)

⁴ Pixel Level Hit Rate (PLHR)



شکل ۴-۵: یک سند نمونه خروجی نهایی

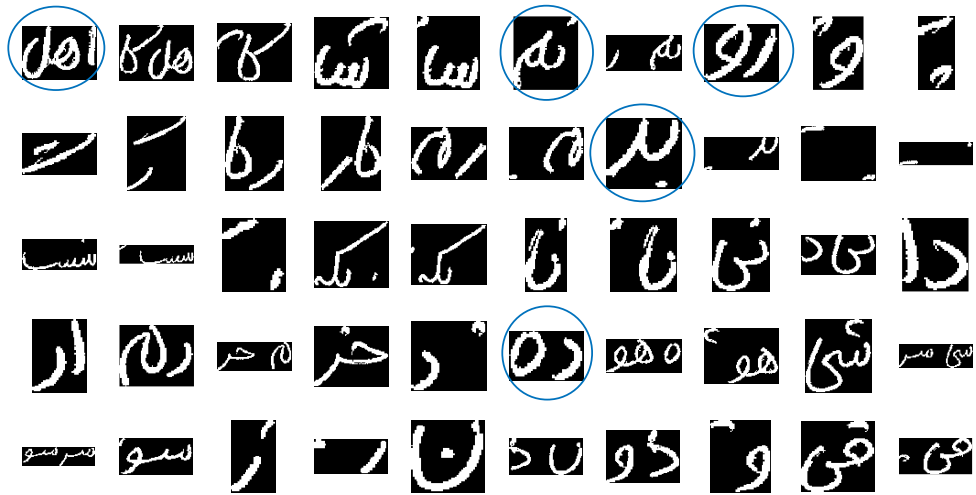
۴-۳- پنجره گذاری مبتنی بر اجزاء متصل

روش رایج پنجره گذاری در بینایی ماشین شامل در نظر گرفتن پنجره با سایزهای متفاوت است. این پنجره ها با همپوشانی زیاد کل فضای جستجو را پیمایش می کنند. این روش پنجره گذاری معایب زیادی دارد: تعیین اندازه/اندازه های مورد نظر برای پنجره لغزان به سادگی قابل محاسبه نیست، که این مسئله خصوصاً در مورد اسناد پیچیده تر بیشتر نمود می کند. حجم بالای پنجره هایی که باید پردازش شوند نیز از معایب بزرگ این روش محسوب می شود. به منظور رفع چالش های مربوط به پنجره گذاری معمول، روش پنجره گذاری خاص مسئله جستجوی کلمه با استفاده از اجزاء متصل پیشنهاد شده است.

این روش پنجره گذاری پیشنهادی از اجزاء متصل تشکیل دهنده خطوط و چهارضلعی محاطی آن ها استفاده می کند. نمونه ای از یک خط از یک سند و چهارضلعی های محاطی آن در شکل ۴-۶ نشان داده شده است. طول پنجره ها در این روش بر اساس تعداد اجزاء متصل درون آن تعیین می گردد. پنجره هایی شامل ۲ تا ۷ جزء متصل مجاور تعیین شده اند. از سمت راست به چپ این پنجره شامل توالی هایی از اجزاء متصل مجاور خواهد بود و روی این زنجیره اجزاء متصل که خط را تشکیل داده اند تا انتها جلو می رود.

اهل کائنات هم روزگار به نسبت کله نانی داریم خرد هوشی سرسوزن زدونی - مادری داریم

شکل ۴-۶: یک سطر از یک سند نمونه با چهارضلعی های محاطی اجزاء متصل



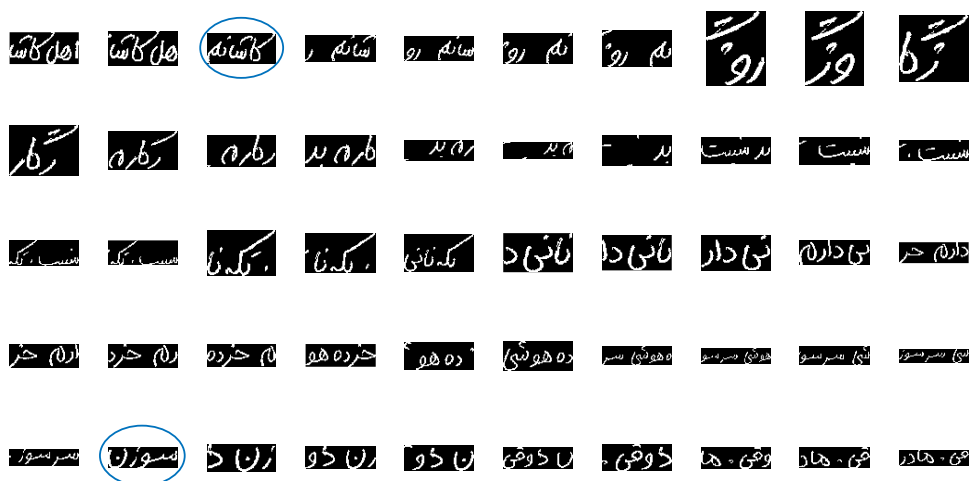
شکل ۴-۷: نمونه‌های از خروجی پنجره به طول ۲ جزء متصل



شکل ۴-۸: نمونه‌های از خروجی پنجره به طول ۳ جزء متصل



شکل ۴-۹: نمونه‌های از خروجی پنجره به طول ۴ جزء متصل



شکل ۴-۱۰: نمونه‌های از خروجی پنجره به طول ۵ جزء متصل

این نوع پنجره‌گذاری یک مزیت بسیار بزرگی که دارد این است که هر جزء متصل یا به‌صورت کامل در پنجره قرار می‌گیرد و یا درون آن قرار نمی‌گیرد که این مسئله باعث خواهد شد در مرحله بازیابی نتایج بهتری حاصل شود. نمونه‌هایی از خروجی‌های حاصل از اعمال این پنجره‌گذاری در شکل‌های بالا آورده شده است. پنجره‌های به‌دست‌آمده از هر سند با توجه به روش تعبیه برچسب که در فصل سوم به‌طور کامل شرح داده شد، در زیرفضای مشترک تعبیه می‌شوند. این زیرفضا بایستی با استفاده از مجموعه داده‌ای از کلمات مجزای دارای برچسب آموزش داده شود.

در مرحله آزمون برای هر کلمه پرسشی ورودی، کافی است عملیات موردنیاز جهت تعبیه آن را در زیرفضای مشترک انجام دهیم و پنجره‌های استخراج شده در مرحله قبل، که در اینجا هر کدام به‌عنوان یک کلمه مجزا در نظر گرفته می‌شود، بر اساس میزان شباهت با کلمه پرسشی ورودی به‌صورت نزولی مرتب می‌شوند.

۴-۴- نتایج تجربی

برای انجام مرحله‌ی ارزیابی از پایگاه داده‌ی FHT، ۱۰ نوع متن متفاوت انتخاب‌شده که از هر نوع متن، ۱۰ دستخط متفاوت وجود دارد. ضمناً این ۱۰ دستخط متفاوت به‌صورت تصادفی از بین حالات مختلف انتخاب‌شده و اسناد مختلف لزوماً یکسان نیستند. به‌طور کلی برای ارزیابی این روش از ۱۰۰ سند تصویری استفاده شده است.

جهت ارزیابی روش پیشنهادی مجموعه‌ای شامل ۶۰ کلمه‌ی متمایز به‌عنوان کلمات پرسشی انتخاب شده‌اند. ۱۰ متن انتخابی دارای موضوعات تقریباً مشابهی هستند تا فراوانی کلمات پرسشی داخل آن‌ها زیاد باشد. به همین دلیل در مجموعه‌ی کلمات پرسشی، ۲۴ کلمه در بیشتر از یک متن ظاهر شده‌اند. در مجموع تعداد این ۲۴ کلمه در داخل ۱۰ متن، ۷۸ می‌باشد. با در نظر گرفتن این ۷۸ کلمه و تعداد کلمات متمایز (۳۱ کلمه)، تعداد کل کلمات پرسشی که برای ۱۰ متن در نظر گرفته شده، ۱۰۹ می‌باشد. برای هر کدام از این ۱۰۹ کلمه، حداقل ۱۰ سند با دستخط متفاوت وجود دارد. بیشترین فراوانی یک کلمه در داخل یک متن ۷ و کمترین آن یک است. جدول ۱-۴ تعداد کلمات پرسشی با فراوانی بین ۱ تا ۷ در هر سند مجزا را نشان می‌دهد. با در نظر گرفتن تعداد فراوانی کلمات، به‌طور کلی ۱۸۱۰ نمونه کلمه در این ۱۰۰ سند وجود دارد که این روش به جستجوی آن‌ها پرداخته است. متوسط تعداد نمونه‌های کلمات پرسشی در متن برابر با ۱/۵۶۵۱ می‌باشد. این الگوریتم در روش پرسش با تصویر، ۱۲۵۰ نمونه از کلمات پرسشی موجود در اسناد تصویری را در رتبه اول آشکار کرده است که بر اساس معیار دقت در اولین بازیابی (P@1) معادل با ۷۸,۱۸٪ می‌باشد و در روش پرسش با رشته‌ی متنی موفق به آشکارسازی ۱۳۱۲ نمونه از کلمات پرسشی در رتبه اول شده است که بر اساس معیار دقت در اولین بازیابی (P@1) معادل با ۸۱,۳۴٪ و ۸۰,۵۵٪ است. مجموعه‌ی کلمات پرسشی از جهت تعداد حروف تشکیل‌دهنده‌ی کلمه متنوع می‌باشند، به‌طوری‌که شامل کلماتی با تعداد حروف کم مثل "ما"، "در"، "شد" و "دو"، کلمات با تعداد حروف متوسط مثل "دارو"، "کشور"، "مردم" و "درآمد" و کلماتی با تعداد حروف زیاد مثل "مدیریت"، "اولویت"، "اقدامات" و "مستقیم" است.

جدول ۱-۴: تعداد کلمات با فراوانی بین ۱ تا ۷ در هر سند مجزا

تعداد نمونه‌های کلمات در متن	۱	۲	۳	۴	۵	۶	۷
تعداد کلمات	۶۷	۲۷	۷	۴	۲	۱	۱

برای هر کدام از اسناد با استفاده از روش جداسازی خطوط پیشنهادی، خطوط سند استخراج شد. سپس با استفاده از روش پنجره‌گذاری پیشنهادی پنجره‌هایی به طول ۲ تا ۱۰ جزء متصل ایجاد و ذخیره شد. زیرفضای مشترکی مبتنی بر مجموعه داده فارسی آموزش داده‌شده و تمامی پنجره‌ها در این فضا نگاشت داده شدند. برای جستجوی هر کلمه پرسشی، تعداد اجزاء متصل کلمه پرسشی ورودی محاسبه می‌شود. سپس این کلمه پرسشی در زیرفضای آموزش داده‌شده تعبیه می‌شود و فاصله اقلیدسی آن با پنجره‌هایی که از ۲ جزء کوچک‌تر تا ۲ جزء بزرگ‌تر از تعداد اجزاء متصل کلمه پرسشی دارند مقایسه می‌گردد و

لیست مرتب شده‌ای از پنجره‌های انتخابی برگردانده می‌شود. جهت ارزیابی نتایج حاصله از معیار میانگین دقت متوسط که در فصل قبل معرفی شد، استفاده شده است و نتایج در جدول ۴-۲ آورده شده است.

جدول ۴-۲: میانگین دقت متوسط با استفاده از روش مبتنی بر قطعه‌بندی خط

نوع کلمه پرسشی / روش	کلمه پرسشی متنی	کلمه پرسشی تصویری
CCA	۸۰/۴۱	۷۷/۷۲
KCCA	۸۱/۰۲	۷۸/۱۴

جدول ۴-۳ لیست مرتب خروجی به ازای چند کلمه پرسشی نمونه را آورده است و بازیابی‌هایی که صحیح انجام شده با کادر قرمز مشخص شده است. این جدول تعداد کلمه پرسشی نمونه در سند و دقت متوسط به ازای هر کلمه پرسشی را نشان می‌دهد. مثلاً برای کلمه پرسشی "در" که در این سند ۶ بار تکرار شده و در رتبه‌های ۱، ۳، ۴، ۵، ۶، ۷ ظاهر شده است دقت متوسط به صورت زیر محاسبه می‌شود:

$$AveP = \frac{1}{1} + \frac{2}{3} + \frac{3}{4} + \frac{4}{5} + \frac{5}{6} + \frac{6}{7} = 0.818$$

جدول ۴-۳: محاسبه دقت متوسط (AveP) برای چند کلمه پرسشی نمونه

Rank1	Rank2	Rank3	Rank4	Rank5	Rank6	Rank7	کلمه پرسشی	تعداد نمونه در سند	AveP
دارو	دارو	دارو	دارو	دارو	دارو	دارو	دارو	۵	۰,۹۶۷
است	است	است	است	است	است	است	است	۵	۰,۷۱۰۲
بود	بود	بود	بود	بود	بود	بود	بود	۷	۱
در	در	در	در	در	در	در	در	۶	۰,۸۱۸
دو	دو	دو	دو	دو	دو	دو	دو	۴	۱
ایران	ایران	ایران	ایران	ایران	ایران	ایران	ایران	۳	۰,۸۶۷
تولید	تولید	تولید	تولید	تولید	تولید	تولید	تولید	۳	۰,۸۶۷

به ازای تمامی کلمات پرسشی در هر سند معیار دقت متوسط (AveP) محاسبه شده است. سپس میانگین آن و مقدار میانگین دقت متوسط (MaP) به ازای کل سیستم محاسبه می‌شود که در جدول ۴-۲ آورده شده است.

تنوع شیوه نگارش مختلف در مجموعه داده‌های مورد استفاده و توانایی روش پیشنهادی در بازیابی کلمات پرسشی دیده نشده در جدول ۴-۴ نشان داده شده است. لیست موجود در این جدول تنوع نحوه نگارش کلماتی را نشان می‌دهد که الگوریتم پیشنهادی موفق به شناسایی آن‌ها شده است.

جدول ۴-۴: لیست نمونه‌هایی از کلمات بازیابی شده توسط روش پیشنهادی با توجه به تنوع شیوه نگارش

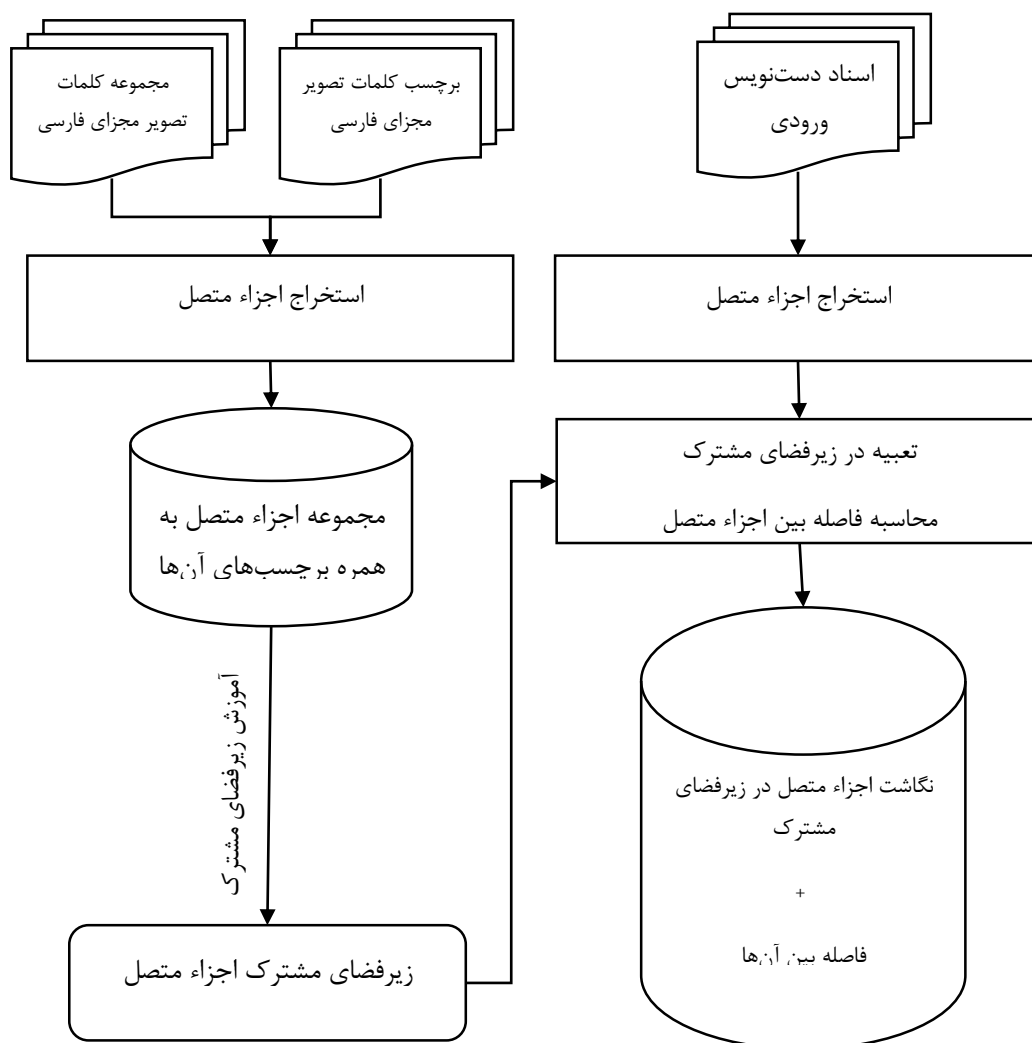
برای	برای	برای	برای	برای	برای
مهمان	مهمان	مهمان	مهمان	مهمان	مهمان
روز	روز	روز	روز	روز	روز
اقدامات	اقدامات	اقدامات	اقدامات	اقدامات	اقدامات
فساد	فساد	فساد	فساد	فساد	فساد
مدیریت	مدیریت	مدیریت	مدیریت	مدیریت	مدیریت
بخش	بخش	بخش	بخش	بخش	بخش
صنعت	صنعت	صنعت	صنعت	صنعت	صنعت
کشور	کشور	کشور	کشور	کشور	کشور
زیبا	زیبا	زیبا	زیبا	زیبا	زیبا

فصل

۵- روش پیشنهادی سوم (مستقل از قطعه بندی)

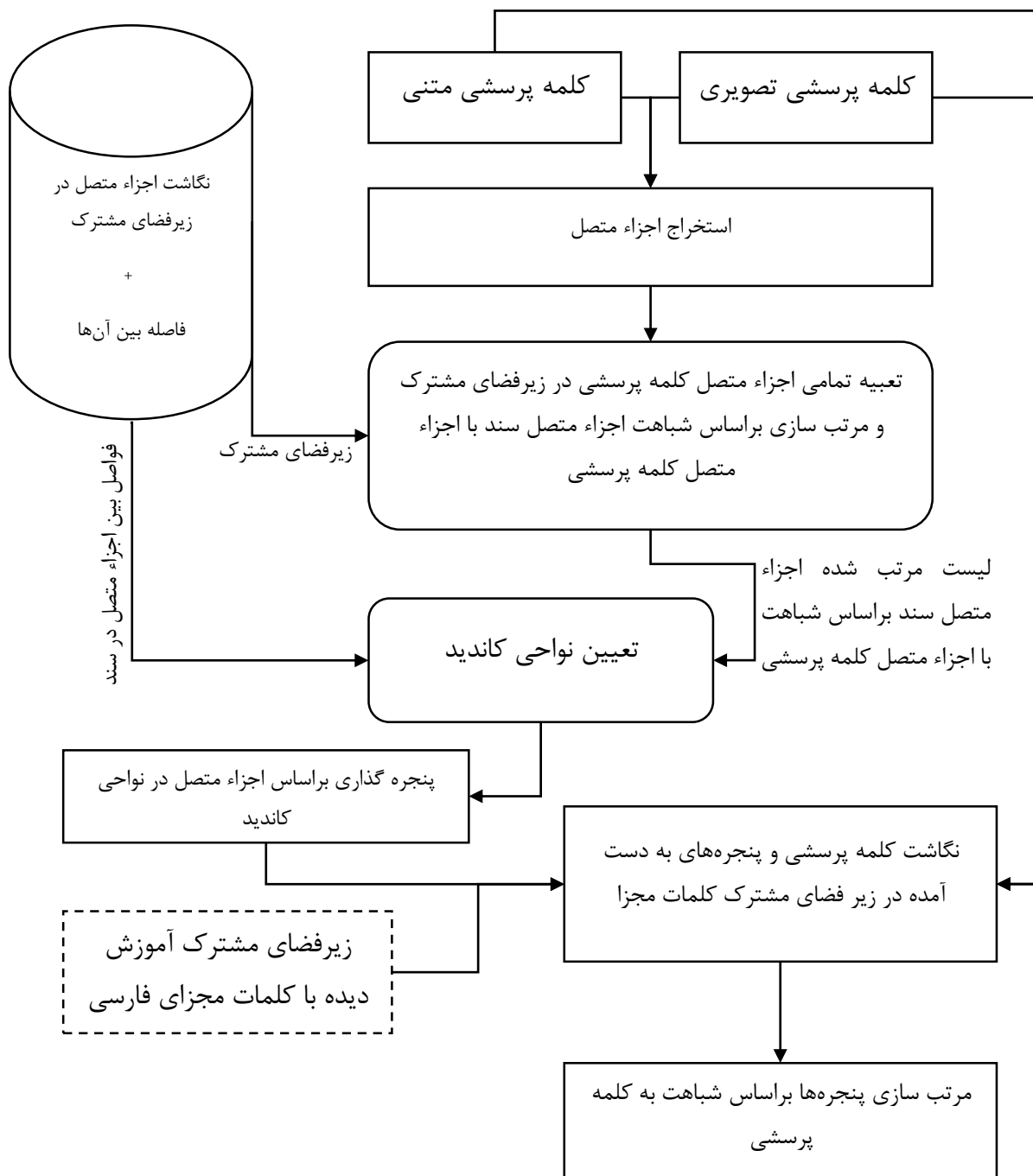
۵-۱- مقدمه

با توجه به مشکلات فراوان ذکرشده در استفاده از روش‌های مبتنی بر قطعه‌بندی، در این فصل روشی مستقل از قطعه‌بندی جهت جستجوی کلمات در اسناد دست‌نویس فارسی پیشنهاد شده است. ساختار کلی این روش بر مبنای تعیین نواحی کاندید بر اساس اجزاء متصل و سپس مرتب‌سازی مجدد^۱ نواحی کاندید می‌باشد.



شکل ۵-۱: مرحله برون‌خط روش پیشنهادی

¹ reranking



شکل ۵-۴: مرحله برخط روش پیشنهادی

در مرحله برون خط روش پیشنهادی، ابتدا یک مجموعه از کلمات تصویری دارای برجسب انتخاب می- شوند و اجزاء متصل و برجسب‌های تعبیه شده آن‌ها بر اساس روش پیشنهادی " هیستوگرام هرمی

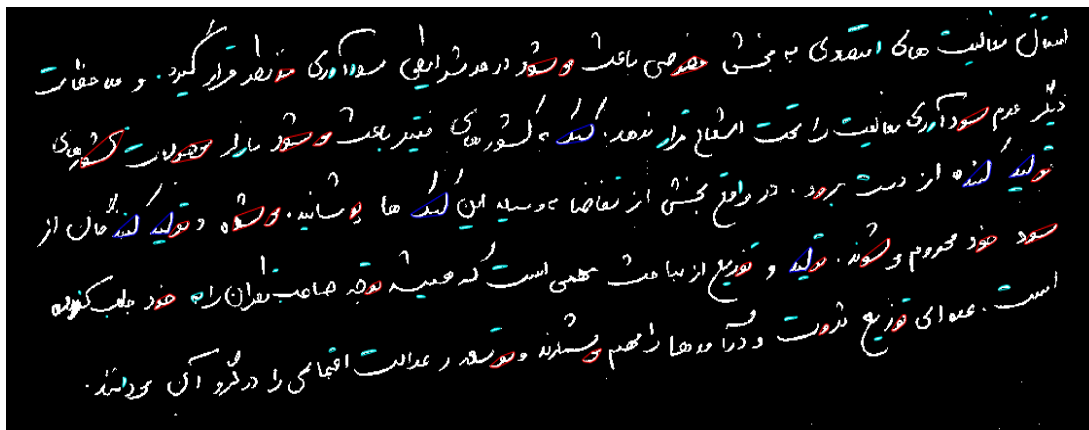
نشانه‌های موقعیتی فارسی " استخراج شده و زیرفضای مشترک آموزش داده می‌شود. بعد از این مرحله، اجزاء متصل اسناد استخراج و در این زیرفضا نگاشت داده می‌شوند. در زیرفضای نگاشت داده شده فاصله کلمات استخراج می‌شود (شکل ۵-۱). از این مدل استخراج شده جهت تعیین نواحی کاندید در مرحله برخط استفاده خواهد شد. علاوه بر این یک مدل مبتنی بر زیرفضای مشترک آموزش دیده روی کلمات مجزا، مانند آنچه در فصل سوم شرح داده شده، نیز برای مرتب‌سازی مجدد نواحی کاندید مورد استفاده قرار می‌گیرد (شکل ۵-۲). در ادامه روش پیشنهادی به تفصیل شرح داده خواهد شد.

۵-۲- تعیین نواحی کاندید

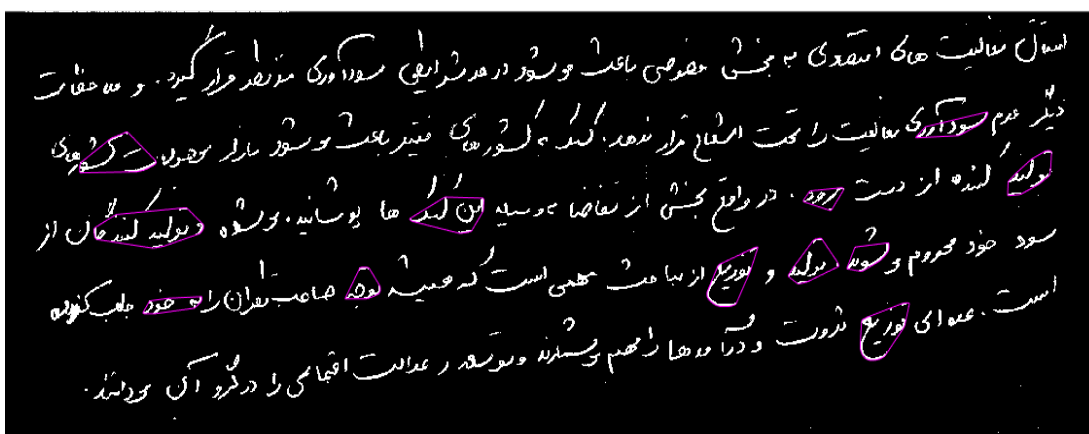
جهت تعیین نواحی کاندید از اجزاء متصل استفاده شده است. ایده اصلی این پیشنهاد بر این اساس استوار است که بازیابی اجزاء متصل نیازی به قطعه‌بندی ندارد، پس می‌تواند به‌عنوان پایه‌ای جهت یک روش مستقل از قطعه‌بندی در نظر گرفته شود. از این رو در صورتی که سیستمی وجود داشته باشد که توانایی بازیابی اجزاء متصل را داشته باشد، می‌توان بر اساس امتیازات اجزاء متصل در سند و نزدیکی اجزاء دارای امتیاز بالا به همدیگر نواحی کاندید را تعیین نمود. بنابراین سیستمی مبتنی بر روش تعبیه برچسب معرفی شده در فصل سوم برای این منظور وفق داده شده است که توانایی بازنمایی اجزاء متصل را در دست‌نوشته‌های فارسی داشته باشد. سپس با استفاده از قانونی ساده مبتنی بر نزدیک بودن اجزاء متصل بازنمایی شده با امتیاز بالا نواحی کاندید تعیین شدند.

با استفاده از سیستم معرفی شده در فصل سوم با استفاده از روش تعبیه برچسب و خصیصه‌ها می‌توان سیستمی طراحی کرد که توانایی جستجو و بازیابی اجزاء متصل نوشته‌های دست‌نویس فارسی را داشته باشد. اگر در طراحی سیستم مجموعه خصیصه‌های معرفی شده با نام "هیستوگرام هرمی نشانه‌های موقعیتی فارسی" استفاده شود، از آنجایی که بدنه حروف و نقاط و سرکش‌ها به‌صورت مجزا در نظر گرفته می‌شود توانایی بازیابی و جستجوی اجزاء متصل را خواهد داشت. زیرا بازنمایی اجزاء متصل با توجه به جداول ۱-۳ و ۲-۳ به‌صورت مجزا امکان‌پذیر است.

بنابراین مجموعه‌ای از اجزاء متصل و برچسب‌های آن‌ها جهت آموزش خصیصه‌ها و زیرفضای مشترک موردنیاز است. بعد از انجام این آموزش هر سند می‌تواند بر اساس مجموعه‌ای از اجزاء متصل تشکیل‌دهنده خود بازنمایی شود. این بازنمایی بر اساس نگاشت اجزاء متصل در زیرفضای مشترک و استخراج حداقل فواصل بین اجزاء متصل سازنده سند انجام می‌پذیرد. این مراحل به‌صورت برون‌خط انجام می‌شود و یک مدل از هر سند بر این اساس آماده می‌شود که در زمان جستجو مورد استفاده قرار می‌گیرد.



(الف)

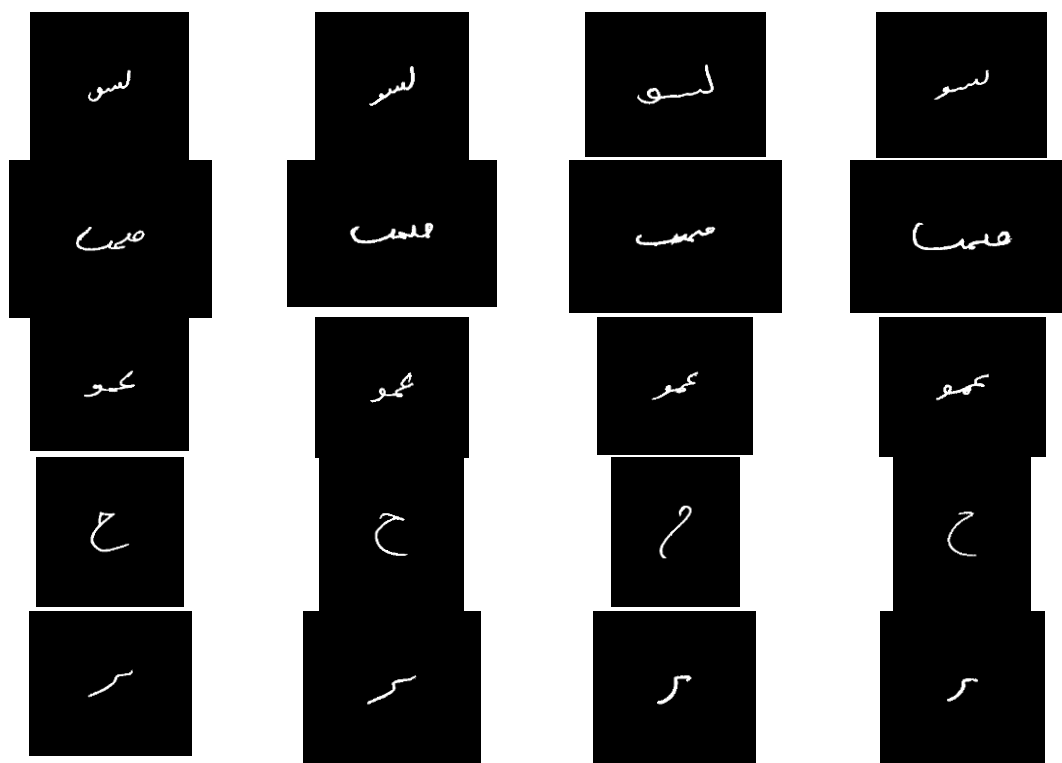


(ب)

شکل ۳-۵: الف) نمونه‌ای از بازیابی اجزاء متصل در سند بر اساس کلمه پرسشی "تولید". هر یک از اجزاء متصل کلمه پرسشی در سند نمونه با رنگی متفاوت نشان داده شده است. ب) نواحی کاندید تعیین شده بر اساس اجزاء متصل بازیابی شده.

در زمان جستجو (برخط) اجزاء متصل کلمه پرسشی استخراج می‌شود. در مورد کلمه پرسشی تصویری، این عمل با استفاده از تکنیک‌های پردازش تصویر انجام می‌پذیرد. اما در مورد کلمه پرسشی متنی، مجموعه‌ای از قوانین مبتنی بر زبان فارسی این کار را انجام می‌دهد. این قوانین بر اساس ساختار نوشتاری زبان فارسی تعیین شده است. برای مثال چون همیشه جستجو روی یک کلمه انجام می‌پذیرد، بنابراین اجزاء متصل کلمه پرسشی دو دسته هستند: ۱- نقاط و سرکش‌ها ۲- بدنه‌های اصلی زیر کلمات که از مفهوم بدنه‌های اصلی حروف مشتق شده است. نقاط به‌سادگی از حروف تشکیل‌دهنده قابل تشخیص هستند. علت اصلی تشکیل زیرکلمات نیز حروفی مانند "ا د ز ژ و" هستند که به حرف بعد از خود

نمی‌چسبند. با قواعدی که برحسب تشخیص این حروف کار می‌کنند می‌توان از روی رشته متنی تعداد اجزاء متصل را تعیین نمود.



شکل ۴-۵: نمونه‌هایی از اجزاء متصل استفاده‌شده در آموزش سیستم جستجو و بازیابی اجزاء متصل

پس از استخراج اجزاء متصل، آن‌ها در زیرفضای مشترک تعبیه می‌شوند و میزان شباهت اجزاء متصل سند با اجزاء متصل کلمه پرسشی محاسبه‌شده و این امتیازات بین ۰ تا ۱ نرمال‌سازی می‌شوند. حال با توجه به تعداد اجزاء متصل کلمه پرسشی، مجموعه‌های مجزا و مرتب‌شده‌ای برای هر جزء متصل کلمه پرسشی داریم که شامل اجزاء متصل سند است. بعدازاین مرحله به‌صورت سلسله مراتبی به دنبال ترکیباتی هستیم که شامل دو نمونه از دو مجموعه متفاوت بوده و فاصله آن‌ها از یک مقدار آستانه‌ای کمتر باشد (کنار هم قرار گرفته باشند). این مقدار آستانه‌ای بر اساس دو پارامتر محاسبه شده است: کمینه متوسط فاصله افقی^۱ و کمینه متوسط فاصله عمودی^۲.

^۱ Min average horizontal distance (MAHD)

^۲ Min average vertical distance (MAVD)

برای محاسبه کمینه متوسط فاصله افقی برای هر سند کمترین فاصله افقی هر جزء متصل در سند محاسبه می‌شود (HD) و هیستوگرام این فواصل محاسبه می‌گردد:

$$S_{HD} = hist(HD, t) \quad , \quad t = t_0, \dots, \max\{HD\} \quad (1-5)$$

که t محدوده فواصل هیستوگرام است. سپس MAHD به صورت زیر محاسبه می‌گردد:

$$MAHD = \left(\sum_{k \in K} S_{HD}(k) \cdot k \right) \left(\sum_{k \in K} S_{HD}(k) \right)^{-1} \quad K \quad (2-5)$$

$$= \{k | S_{HD}(k) > 1/4 \max\{S_{HD}\}\}$$

در این محاسبه فواصل خیلی کم که ممکن است ناشی از اشتباهات یا فواصل اجزاء در خطوط مختلف باشد بر اساس شرایط تعیین شده برای k در نظر گرفته نمی‌شود. مقدار MAVD نیز به همین صورت محاسبه می‌گردد.

با استفاده از این دو پارامتر در صورتی که فاصله اجزاء متصل انتخابی از دو مجموعه متفاوت فاصله عمودی و افقی کمتر یا مساوی با مقادیر MAHD و MAVD باشد به عنوان یک مجموعه دوتایی انتخاب شده و از مجموعه‌های تکی حذف می‌شوند. بسته به تعداد مجموعه‌ها، این عمل برای مجموعه‌های سه عضوی و بیشتر تکرار می‌شود.

مجموعه‌های بیش از دو عضو به عنوان نواحی کاندید در نظر گرفته می‌شود. فقط یک حالت استثناء وجود دارد که آن هم در مواردی است که کلمه پرسشی دارای یک جزء متصل باشد که در این صورت این مرحله انجام نمی‌شود. نواحی بیش از دو عضوی انتخاب شده گسترش داده می‌شود، یعنی چند جزء متصل سمت چپ و راست این مجموعه به ناحیه اضافه شده و به عنوان ناحیه کاندید معرفی می‌شود. علت این است که با استفاده از این کار در صورتی که بخشی از کلمه تشخیص داده شده باشد یا چسبیدگی یا شکستگی اجزاء متصل در سند وجود داشته باشد بتوان تا حدی آن را پوشش داد. پس با استفاده از این روش نواحی کاندید مشخص می‌شوند (شکل ۳-۵ ب).

۳-۵- مرتب‌سازی مجدد

نواحی کاندیدی که در مرحله قبل تعیین شدند در این مرحله به چندین پنجره تقسیم می‌شوند. این پنجره‌ها بر اساس تعداد اجزاء متصل کلمه پرسشی تعیین می‌شوند. از آنجایی که در برخی موارد در

دست‌نوشته‌ها نقاط، سرکش‌ها یا بدنه‌های یک کلمه به هم می‌چسبند برای کاهش میزان خطا پنجره-هایی شامل تعداد ۲ جزء متصل کمتر تا ۲ جزء بیشتر از اجزاء کلمه پرسشی پنجره ایجاد می‌شود. جدول ۵-۱ نمونه‌ای از این پنجره‌های تولید شده را برای یک ناحیه کاندید برای کلمه پرسشی "تولید" نشان می‌دهد. نهایتاً، لیست پنجره‌های تولید شده به زیرفضای مشترک آموزش داده شده با کلمات مجزا نگاشت داده می‌شود و بر اساس میزان مشابهت با کلمه پرسشی وارد شده مرتب می‌شوند.

جدول ۵-۱: پنجره‌های نمونه برای یک ناحیه کاندید برای کلمه پرسشی "تولید"

ناحیه کاندید موردنظر				
تعداد جزء متصل در پنجره				
2	3	4	5	6
د	د تو	د تولد	د تولد	د تولد لید
تو	تولد	تولد	تولد لید	تولد لید
تولد	تولد	تولد لید	تولد لید	تولد لید
لید	لید لید	لید لید	لید لید	لید لید
لید	لید	لید	لید	
لید	لید	لید		
لید	لید			
لید				

۵-۴- نتایج تجربی

برای انجام مرحله‌ی ارزیابی از پایگاه داده‌ی FHT، ۱۰ نوع متن متفاوت انتخاب شده که از هر نوع متن، ۱۰ دستخط متفاوت وجود دارد. ضمناً این ۱۰ دستخط متفاوت به‌صورت تصادفی از بین حالات

مختلف انتخاب شده و در مورد اسناد مختلف لزوماً یکسان نیستند. به‌طور کلی برای ارزیابی این روش از ۱۰۰ سند تصویری استفاده شده است.

جهت ارزیابی روش پیشنهادی مجموعه‌ای شامل ۶۰ کلمه‌ی متمایز به‌عنوان کلمات پرسشی انتخاب شده‌اند. ۱۰ متن انتخابی دارای موضوعات تقریباً مشابهی هستند تا فراوانی کلمات پرسشی داخل آن‌ها زیاد باشد. به همین دلیل در مجموعه‌ی کلمات پرسشی، ۲۴ کلمه در بیشتر از یک متن ظاهر شده‌اند. در مجموع تعداد این ۲۴ کلمه در داخل ۱۰ متن، ۷۸ می‌باشد. با در نظر گرفتن این ۷۸ کلمه و تعداد کلمات متمایز (۳۱ کلمه)، تعداد کل کلمات پرسشی که برای ۱۰ متن در نظر گرفته شده، ۱۰۹ می‌باشد. برای هر کدام از این ۱۰۹ کلمه، حداقل ۱۰ سند با دستخط متفاوت وجود دارد. بیشترین فراوانی یک کلمه در داخل یک متن ۷ و کمترین آن یک است. جدول ۴-۵-۲ تعداد کلمات پرسشی با فراوانی بین ۱ تا ۷ در هر سند مجزا را نشان می‌دهد. با در نظر گرفتن تعداد فراوانی کلمات، به‌طور کلی ۱۸۱۰ نمونه کلمه در این ۱۰۰ سند وجود دارد که این روش به جستجوی آن‌ها پرداخته است. متوسط تعداد نمونه‌های کلمات پرسشی در متن برابر با ۱/۵۶۵۱ می‌باشد. مجموعه‌ی کلمات پرسشی از جهت تعداد حروف تشکیل‌دهنده‌ی کلمه متنوع می‌باشند، به‌طوری‌که شامل کلماتی با تعداد حروف کم مثل "ما"، "در"، "شد" و "دو"، کلمات با تعداد حروف متوسط مثل "دارو"، "کشور"، "مردم" و "درآمد" و کلماتی با تعداد حروف زیاد مثل "مدیریت"، "اولویت"، "اقدامات" و "مستقیم" است.

جدول ۴-۵-۲: تعداد کلمات با فراوانی بین ۱ تا ۷ در هر سند مجزا

تعداد نمونه‌های کلمات در متن	۱	۲	۳	۴	۵	۶	۷
تعداد کلمات	۶۷	۲۷	۷	۴	۲	۱	۱

جهت ارزیابی استخراج نواحی کاندید به معیاری کارا نیاز است. در مسائلی که نواحی کاندید تعیین می‌شوند، همیشه یک تقابل بین میزان صحت تشخیص این نواحی و میزان کاهش فضای جستجو وجود دارد. از این رو برای بررسی میزان اثربخشی روش پیشنهادی این دو مسئله باید مورد ارزیابی قرار گیرد. پس در تعیین نواحی کاندید نرخ تشخیص نواحی کاندید صحیح باید مدنظر قرار داده شود و همچنین نرخ کاهش فضای جستجو انجام شده با استفاده از روش تعیین نواحی کاندید بایستی در نظر گرفته شود.

برای مسئله اول معیار نرخ تشخیص پیشنهاد شده است. این معیار نشان دهنده تعداد کل نواحی تشخیص داده شده صحیح بر تعداد کل نواحی صحیح است و به‌صورت زیر تعریف می‌شود:

$$DR = \frac{\sum_{d \in D} \sum_{q \in Q} \text{Positive_Selected_Candidate}(d, q, o)}{\sum_{d \in D} \sum_{q \in Q} \text{actual_occurrence}(d, q)} \quad (3-5)$$

که D تعداد اسناد موجود در مجموعه داده مورد بررسی و Q تعداد کلمات پرسشی مورد بررسی در هر سند را نشان می‌دهد. اگر کل سند به‌عنوان ناحیه کاندید در نظر گرفته شود، به دلیل اینکه قطعاً شامل ناحیه کاندید مورد جستجو خواهد بود، این معیار ۱۰۰٪ خواهد شد. از این رو نیاز به معیار دیگری نیز وجود دارد که میزان کاهش فضای جستجو را نمایش دهد. برای این کار سه معیار متفاوت معرفی شده است.

معیار میزان کاهش تعداد اجزاء متصل^۱ که نرخ کاهش فضای جستجو را بر اساس تعداد اجزاء متصل در نواحی کاندید نشان می‌دهد و به‌صورت زیر تعریف شده است:

$$CCRR = \frac{\sum_{d \in D} CC_count(d) - CC_in_candidate_region(d)}{\sum_{d \in D} CC_count(d)} \quad (4-5)$$

معیار میزان کاهش بر اساس تعداد پیکسل نیز به‌صورت زیر تعریف شده است:

$$RR = \frac{\sum_{d \in D} (\sum_{c \in CC} \text{pixel_count}(c, d) - \sum_{cp \in \text{CandidateCC}} \text{pixel_count}(cp, d))}{\sum_{d \in D} \sum_{c \in CC} \text{pixel_count}(c, d)} \quad (5-5)$$

در این روابط d اسناد، c تعداد اجزاء متصل، cc تعداد اجزاء متصل در سند و cp تعداد اجزاء متصل در نواحی کاندید را نشان می‌دهد. Pixel_count تعداد پیکسل‌های اجزاء متصل را محاسبه می‌کند.

معیار کاهش بر اساس تعداد پنجره‌ها با روش مطرح شده در فصل ۴ که مبتنی بر خط کار می‌کند نیز مقایسه شده است. این نرخ کاهش تعداد پنجره‌ها را نشان می‌دهد و به‌صورت زیر مطرح شده است:

$$WRR = \frac{\sum_{d \in D} \sum_{q \in Q} \text{Window_Count}(d, q)}{\sum_{d \in D} \text{Query_count}(d) \times \text{Window_count_Line}(d)} \quad (6-5)$$

اسناد در این رابطه با d نشان داده شده و q نشان دهنده کلمات پرسشی است. Windows_count تعداد پنجره‌های یک سند به ازای یک کلمه پرسشی خاص را نشان می‌دهد.

جدول ۳-۵ نرخ تشخیص را به ازای میزان اشتراک ناحیه کاندید با ناحیه واقعی نشان می‌دهد. بر اساس نتایج تجربی، روش تعیین نواحی کاندید پیشنهادی بین ۹۰ تا ۹۴ درصد از کاندیدهای موجود را به

¹ Connected Components Reduction Rate (CCRR)

درستی تشخیص می‌دهد. با توجه به پیچیدگی دست‌نوشته‌های فارسی و ساده بودن مدل پیشنهادی، نتایج قابل قبول به نظر می‌رسد. علاوه بر این با توجه به از دست دادن بخشی از نواحی کاندید، فضای جستجو بسیار کاهش می‌یابد که نتایج تجربی این کاهش بر اساس معیارهای پیشنهادی در جدول ۵-۴ نشان داده شده است.

جدول ۵-۳: نرخ تشخیص بر اساس میزان اشتراک ناحیه کاندید با ناحیه واقعی

نرخ اشتراک	نرخ تشخیص (DR)
>0%	۹۴/۴۶
>=50%	۹۳/۰۷
>=60%	۹۲/۴۳
>=70%	۹۱/۷۴
>=80%	۹۰/۷۰
>=90%	۹۰/۰۲

جدول ۵-۴: نرخ کاهش بر اساس معیارهای معرفی شده

معیار	نرخ
RR	٪۶۳/۰۹
CCRR	٪۷۳/۸۰
WRR	٪۷۶/۶۵

بر اساس نتایج حاصل شده، روش معرفی شده دقتی در حدود ۹۰ درصد در سخت گیرانه‌ترین شرایط توانایی تشخیص نواحی کاندید را داراست و علاوه بر این در حدود ۷۰ درصد از فضای جستجو را کاهش می‌دهد. بر اساس تحلیل خطاهای صورت گرفته، اکثر کاندیدای از دست رفته در اثر مشکلات نوشتاری یا ایراد دستگاه‌های تصویربرداری حاصل می‌شوند. این ایرادات شامل موارد زیر است:

- به هم چسبیدگی اجزاء متصل: در دست‌نوشته‌ها این مسئله مشکل رایجی است که ممکن است بین کلمات کنار هم در یک خط یا بین کلمات دو خط متوالی رخ دهد. برخی نمونه‌ها در مجموعه داده مورد استفاده در شکل ۵-۵ نمایش داده شده است.

سخن بهم ساری

شکل ۵-۵: نمونه‌ای از به هم چسبیدگی اجزاء متصل. کلمات "سرخ"، "دارم" و "می‌گیرم"

- شکستگی یک جزء متصل به دو یا چند بخش: در برخی شرایط یک جزء متصل به دو یا چند جزء متصل شکسته می‌شود. نمونه‌هایی از این حالت در مجموعه داده مورد استفاده در شکل ۶-۵ نشان داده شده است.

پسند دوستانی

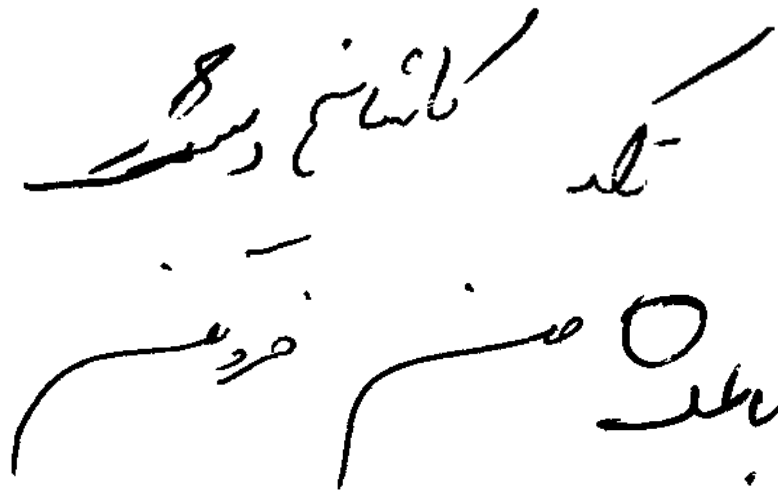
شکل ۵-۶: نمونه‌ای از شکستگی اجزاء متصل. کلمات "بخش" و "دوستانی"

- خط‌خوردگی یا اصلاح نوشتار و بنابراین خارج شدن از فرم استاندارد نگارشی: این حالت خاص نیز در مجموعه داده وجود دارد که باعث کاهش دقت بازیابی چه در سطح اجزاء متصل و چه در سطح کلمه می‌شود.

شبه حیوان قبله‌ام

شکل ۵-۷: نمونه‌ای از خط‌خوردگی. کلمات "شبه‌در"، "حیوان" و "قبله‌ام"

- نوشتن با ساختاری متفاوت از استاندارد نگارشی: در برخی موارد دستخط و شیوه نگارش یک کلمه خاص بسیار متفاوت خواهد بود که روش پیشنهادی توانایی شناسایی برخی از حالات که تفاوت زیادی دارند را ندارد (شکل ۵-۸)



شکل ۸-۵: نمونه‌ای از ساختار متفاوت با استانداردهای نگارشی. کلمات "تکه"، "کاشانم"، "دشت"، "باعث"، "من" و "فروتن"

- شباهت بخشی از یک کلمه در کلمه‌ای دیگر مانند "در" موجود در کلمه "درآمد"

ارزیابی کل سیستم بر اساس دو معیار معرفی شده در فصل سوم انجام شده است: معیار میانگین دقت متوسط (MaP) و معیار دقت در اولین بازیابی (P@1). نتایج این ارزیابی‌ها به ترتیب در جداول ۵-۵ و ۵-۶ نشان داده شده است.

جدول ۵-۵: میانگین دقت متوسط (MaP) برای کلمات پرسشی تصویری و متنی

	CCA	KCCA
کلمه پرسشی متنی	۸۰/۱۴	۸۰/۲۸
کلمه پرسشی تصویری	۷۴/۲۵	۷۵/۰۷

جدول ۵-۶: میانگین دقت اولین بازیابی (P@1) برای کلمات پرسشی تصویری و متنی

	CCA	KCCA
کلمه پرسشی متنی	۷۹/۲۹	۷۹/۵۳
کلمه پرسشی تصویری	۷۲/۲۱	۷۳/۸۶

نتایج تجربی نشان می‌دهد که در روش ارائه شده همواره دقت برای کلمات پرسشی تصویری کمتر از کلمات پرسشی متنی است. علت این امر را می‌توان در نگاشت متن و تصویر در فضای خصیصه پیگیری کرد. برای نگاشت متن به فضای خصیصه از مجموعه‌ای از قوانین مرتبط با ساختار زبان فارسی و حروف

تشکیل‌دهنده کلمه پرسشی استفاده شده و یک بردار دودویی به‌عنوان بازنمایی در فضای خصیصه ساخته می‌شود. اما عمل نگاشت در مورد کلمه تصویری از مجموعه‌ای از ماشین‌های بردار پشتیبان بهره می‌برد که توسط مجموعه داده آموزش داده شده است. با توجه به اینکه تبدیل در مورد متن بدون خطا و در مورد تصویر در بهترین حالات نیز با خطا همراه است، بنابراین تفاوت دقت در این‌باره منطقی به نظر می‌رسد.

فصل

۶- نتیجه‌گیری و کارهای آینده

۶-۱- نتیجه‌گیری

در این پایان‌نامه یک رهیافت برای جستجوی کلمه پرسشی در اسناد دست‌نوشته‌ی زبان فارسی پیشنهاد گردید. در این روش کلمه‌ی پرسشی موردنظر کاربر در اسناد تصویری جستجو شده و مکان نمونه‌های آن کلمه در داخل اسناد مشخص می‌شود. کلمه پرسشی می‌تواند به‌صورت تصویر یا رشته‌ی متنی باشد. یکی از مهم‌ترین مسائل در آنالیز اسناد تصویری، مانند سایر حوزه‌های بینایی ماشین، روش بازنمایی است که در دهه‌های اخیر کارهای بسیار زیادی در جهت حل این مشکل انجام شده است. روش بازنمایی قابل‌قبول است که نسبت به تغییرات مقاوم بوده و قابلیت تفکیک‌پذیری بالایی نیز داشته باشد. بنابراین در حل مسئله مذکور در این رساله، در اولین گام به دنبال یک روش بازنمایی مناسب بوده‌ایم که توانایی مدل کردن مستقل از نوع نگارش را داشته باشد. برای حل مسئله مذکور روشی مبتنی بر تعبیه برچسب، رده‌بندی مبتنی بر خصیصه معرفی شد که با آموزش یک زیرفضای مشترک مسئله را به یک مسئله ساده نزدیک‌ترین همسایگی تبدیل می‌کند. خصیصه‌های مختلفی بر اساس ساختار زبان فارسی معرفی شد که توانایی بازنمایی بهتری را در زبان فارسی داشته باشیم. یکی از این خصیصه‌ها که بدنه حروف را از نقاط و سرکش‌ها به‌صورت مجزا در نظر می‌گیرد (PPHOS) این قابلیت را به سیستم اضافه می‌کند که بتواند علاوه بر کلمات، اجزاء متصل سازنده کلمات را نیز بازیابی کند.

جهت جستجوی کلمه مستقل از قطعه‌بندی سند دو روش پیشنهاد داده شد. اولین روش مبتنی بر قطعه‌بندی خطوط می‌باشد. در این روش، خطوط سند استخراج شده و پنجره‌گذاری مبتنی بر اجزاء متصل با طول ۲ تا ۱۰ جزء ایجاد می‌گردد. این پنجره‌ها در زیر فضای مشترک آموزش داده شده نگاشت می‌شوند. برای هر کلمه پرسشی، متنی یا تصویری، عمل نگاشت به زیر فضای مشترک انجام می‌شود و بر اساس نزدیکی فاصله به این کلمات پنجره‌های سند بازیابی می‌شوند. در این روش مشکلات پنجره-گذاری ساده وجود ندارد و تعداد پنجره به مراتب کمتر است. علاوه بر این دقت مطلوبی گزارش می‌شود.

در روش دوم به‌طور کامل مستقل از هر نوع قطعه‌بندی عمل شده است. بر اساس خصیصه‌ی معرفی شده‌ای که توانایی جستجوی اجزاء متصل سازنده کلمات فارسی را دارد، یک روش دو مرحله‌ای پیشنهاد داده شد. این روش نگاشت اجزاء متصل هر سند را در زیرفضای مشترک آموزش داده شده به همراه فاصله آن‌ها از یکدیگر ذخیره می‌کند. برای هر کلمه پرسشی ابتدا اجزاء متصل استخراج شده و بر اساس میزان شباهت، اجزاء متصل مشابه در سند استخراج می‌شوند. بر اساس انباشتگی اجزاء در موقعیت‌های مکانی خاصی از سند نواحی کاندید تعیین می‌شود. سپس این نواحی کاندید بر اساس یک زیرفضای مشترک آموزش داده شده با کلمات مجزا نسبت به کلمه پرسشی مرتب می‌شوند.

با توجه به اینکه در مجموعه داده استفاده شده تنوع دستخطها بسیار زیاد بود و مجموعه آموزش از مجموعه آزمون در تمامی مراحل جدا در نظر گرفته شد، نتایج تجربی حاصل شده بر توانایی و قدرت روش پیشنهادی صحه گذاشت. علاوه بر این، روش ارائه شده مستقل از تمامی محدودیت‌هایی که در مورد روش‌های ارائه شده در جستجوی کلمه در نظر گرفته می‌شود، کار می‌کند. یعنی این روش مستقل از فرهنگ لغت کار می‌کند و توانایی شناسایی کلمات خارج از فرهنگ لغت را نیز دارد. همچنین بر روی اسناد چند نویسنده با یک مدل ثابت کار می‌کند و نیاز به مدل‌های متفاوت برای اسناد متفاوت ندارد. توانایی کار کردن با کلمات پرسشی تصویری و متنی را داراست، هر چند در تمامی مراحل دقت کلمه پرسشی متنی بهتر بوده است. علاوه بر این، روش مستقل از قطعه‌بندی پیشنهادی مشکلات و خطاهای قطعه‌بندی را ندارد.

۶-۲- کارهای آینده

در آزمایش‌های اولیه، جهت آموزش بهتر خصیصه‌ها اهمیت وجود یک مجموعه داده آموزشی بزرگ که توانایی بازنمایی خوبی داشته باشد مشخص شد. مسئله وجود یک مجموعه داده آموزشی بزرگ در تمامی کارهای شناسایی الگو همواره به‌عنوان یک مسئله مهم مطرح بوده است. اما در روش پیشنهادی به دلیل ساختار خاص خصیصه‌های معرفی شده این مسئله به نحوی متفاوت مطرح می‌شود. به دلیل اینکه در تمامی سطوح مانند کل کلمه، نصف اول و دوم کلمه و یک سوم اول دوم و سوم خصیصه‌ها که شامل تمامی حروف الفبا هستند آموزش داده شوند، باید مجموعه داده طوری آماده شود که امکان چنین حالتی را ممکن سازد. این اهمیت در خصوص زیرسیستم بازیابی اجزاء متصل اهمیت بیشتری دارد. مجموعه داده‌هایی که وجود دارند و از آن‌ها در این رساله استفاده شده است، نمونه‌های مثبت کافی برای آموزش برخی از خصیصه‌ها را فراهم نمی‌کند. حتی در بعضی موارد برخی از خصیصه‌ها هیچ نمونه‌ای برای آموزش در دسترس ندارند. فراهم کردن مجموعه داده بزرگ کاری بسیار زمان‌بر و خسته‌کننده است. علاوه بر این تهیه فهرستی از کلاس کلمات مناسبی که بتواند به خوبی تمام کلاس خصیصه‌های موردنظر را آموزش دهد مشکل است. بنابراین به‌عنوان یکی از راهکارهایی که می‌توان در آینده جهت گسترش این کار مطرح نمود، سنتز داده‌های آموزشی بر اساس کلاس کلمات موردنیاز خواهد بود. با این روش می‌توان داده‌ها را طوری تولید کرد که تمامی خصیصه‌ها به خوبی و کامل آموزش داده شود، بنابراین مدل‌های خصیصه بهتر تولید خواهد شد و نتایج به‌طور قابل توجهی افزایش خواهد یافت.

در نظر گرفتن کلمه پرسشی و ساختار سند به‌صورت گراف به‌عنوان راهکاری است که می‌توان در ادامه کار از آن بهره برد. در این ساختار مبتنی بر گراف هر جزء متصل به‌عنوان یک نود از گراف در نظر گرفته

می‌شود و یال‌های گراف بین اجزاء متصل نزدیک به هم در نظر گرفته می‌شود. این گراف به‌صورت نسبی اجزاء زیر کلمات، توالی زیر کلمات و کلمات را در متن سند نشان می‌دهد. در این روش می‌توان با استفاده از یک معیار شباهت مبتنی بر جهت قرار گرفتن اجزاء متصل و میزان شباهت آن به گراف تشکیل شده از توالی اجزاء متصل در کلمه پرسشی، نواحی کاندید را استخراج نماید. علاوه بر این در صورت تعیین یک معیار شباهت دقیق می‌توان عمل تشخیص را به‌صورت یک مرحله‌ای و بدون نیاز به تعیین نواحی کاندید انجام داد.

یکی از مسائل مهم در روش مستقل از قطعه‌بندی پیشنهادی، مسئله تعیین اجزاء متصل مرتبط با هم در سند است. در روش پیشنهادی صرفاً بر اساس فاصله و جهت قرار گرفتن اجزاء متصل این مسئله حل می‌شود. راهکار دیگری که می‌توان جهت ادامه کار ارائه کرد، استفاده از مدلی شبیه به مدل‌های زبانی در پردازش زبان طبیعی است. در مورد روش پیشنهادی مستقل از قطعه‌بندی که در مرحله اول اجزاء متصل را بازیابی می‌کند، می‌توان از یک مدل شبیه به مدل‌های زبانی جهت بالا بردن دقت تشخیص نواحی کاندید بهره برد. مدل‌های زبانی بر اساس احتمال رخداد توالی کلمات در پردازش زبان طبیعی در نظر گرفته می‌شود. در این مدل خاص، مدل زبانی بر اساس ساختار تصویری و بر اساس اجزاء متصل سازنده کلمات ساخته می‌شود و می‌تواند به‌عنوان راهکاری در جهت تعیین دقیق‌تر نواحی کاندید در نظر گرفته شود. در این روش می‌توان ارزیابی نواحی کاندید مبتنی بر اجزاء متوالی مجاور را بر اساس این مدل و ترکیب آن با میزان شباهت اجزاء متصل کلمه پرسشی با ناحیه کاندید در نظر گرفت. این روش می‌تواند با ساختار گرافی اجزاء متصل اسناد ترکیب شود و دقت روش مذکور را با ساخت گراف دقیق‌تر با استفاده از مدل یاد گرفته شده افزایش دهد.

مراجع

- [۱] T. Kameshiro, T. Hirano, Y. Okada, and F. Yoda, "A document image retrieval method tolerating recognition and segmentation errors of OCR using shape-feature and multiple candidates," in *Document Analysis and Recognition, ICDAR'99. Proceedings of the Fifth International Conference on*, 1999, pp. 681-684.
- [۲] ح. خسروی, "رویکرد یکپارچه برای بازشناسی متن چاپی فارسی," رساله دکتری, دانشگاه تربیت مدرس, ۱۳۸۷.
- [۳] L. Li, S. Lu, and C. L. Tan, "A fast keyword-spotting technique," in *Document Analysis and Recognition. ICDAR 2007. Ninth International Conference on*, 2007, pp. 68-72.
- [۴] G. Jones, J. Foote, K. Sparck Jones, and S. Young, "Video mail retrieval: The effect of word spotting accuracy on precision," in *Acoustics, Speech, and Signal Processing .ICASSP-95., 1995 International Conference on*, 1995, pp. 309-312.
- [۵] T. M. Rath, V. Lavrenko, and R. Manmatha, "Retrieving historical manuscripts using shape," *Massachusset Univ Amherst Center for Intelligent Information Retrieval*, 2003.
- [۶] ا. ابراهیمی, "استفاده از شکل کلی زیرکلمات چاپی در بازیابی تصویر مستندات و بازشناسی متون فارسی," رساله دکتری, مهندسی برق دانشگاه تربیت مدرس ۱۳۸۴.
- [۷] J. Almazán, A. Gordo, A. Fornés, and E. Valveny, "Word Spotting and Recognition with Embedded Attributes," *IEEE transactions on pattern analysis and machine intelligence*, vol. 36, pp. 2552-2566, 2014.
- [۸] A. R. A. Z.Imani, A.Zohrevand, "Introduction to Database FARSA digital images of handwritten Farsi words," presented at the Eleventh Iranian Conference on Intelligent Systems (ICIS2013), Kharazmi University. Tehran , Iran, 2013.
- [۹] M. Ziaratban, K. Faez, and F. Bagheri, "FHT: An unconstraint Farsi handwritten text database," in *10th International Conference on Document Analysis and Recognition*, 2009, pp. 281-285.
- [۱۰] S. Lu and C.-L. Tan, "Keyword spotting and retrieval of document images captured by a digital camera," in *Ninth International Conference on Document Analysis and Recognition (ICDAR 2007)*, 2007, pp. 994-998.
- [۱۱] A. Kefali and C. Chemmam, "A Semi-Automatic Approach of old Arabic Documents Indexing," in *3rd International Conference on Computer Science and its Applications (CIIA 2011)*, Saida , Algeria 2011.

- [۱۲] Y. Pourasad, H. Hassibi, and A. Ghorbani, "A word spotting method for Farsi machine-printed document images," *Turkish Journal of Electrical Engineering & Computer Sciences*, vol. 21, pp. 734-746, 2013.
- [۱۳] Y. POURASAD, H. HASSIBI, and A. GHORBANI, "A Farsi/Arabic Word Spotting Approach for Printed Document Images," *International Journal of Natural & Engineering Sciences*, vol. 6, 2012.
- [۱۴] Y. Pourasad, H. Hassibi, and A. Ghorbani, "Farsi Word Spotting and Font Size Recognition," *Procedia Technology*, vol. 1, pp. 372-377, 2012.
- [۱۵] B. Zhang, S. N. Srihari, and C. Huang, "Word image retrieval using binary features," in *Electronic Imaging 2004*, pp. 45-53.
- [۱۶] S. Srihari, H. Srinivasan, P. Babu, and C. Bhole, "Spotting words in handwritten Arabic documents," in *Electronic Imaging 2006*.
- [۱۷] T. Adamek, N. E. O'Connor, and A. F. Smeaton, "Word matching using single closed contours for indexing handwritten historical documents," *International Journal of Document Analysis and Recognition (IJ DAR)*, vol. 9, pp. 153-165, 2007.
- [۱۸] H. Cao and V. Govindaraju, "Template-free word spotting in low-quality manuscripts," in *Proceedings of the 6th International Conference on Advances in Pattern Recognition*, 2007, pp. 135-139.
- [۱۹] T. M. Rath, R. Manmatha, and V. Lavrenko, "A search engine for historical manuscript images," in *Proceedings of the 27th annual international ACM SIGIR conference on Research and development in information retrieval*, 2004, pp. 369-376.
- [۲۰] T. M. Rath and R. Manmatha, "Word image matching using dynamic time warping," in *Computer Vision and Pattern Recognition. Proceedings. IEEE Computer Society Conference on*, 2003, pp. II-521-II-527 vol. 2.
- [۲۱] K. Terasawa and Y. Tanaka, "Slit style HOG feature for document image word spotting," in *Document Analysis and Recognition. ICDAR'09. 10th International Conference on*, 2009, pp. ۱۲۰-۱۱۶ .
- [۲۲] T. M. Rath and R. Manmatha, "Word spotting for historical documents," *International Journal of Document Analysis and Recognition (IJ DAR)*, vol. 9, pp. 139-152, 2007.
- [۲۳] R. Manmatha and T. M. Rath, "Indexing of handwritten historical documents-recent progress," presented at the Proceedings 2003 Symposium on Document Image Understanding Technology.(UMD), 2003.
- [۲۴] S. N. Srihari, H. Srinivasan, C. Huang, and S. Shetty, "Spotting words in Latin, Devanagari and Arabic scripts," *VIVEK-BOMBAY-*, vol. 16, p. 2, 2006.
- [۲۵] A. Kolcz, J. Alspector, M. Augusteijn, R. Carlson, and G. V. Popescu, "A line-oriented approach to word spotting in handwritten documents," *Pattern Analysis & Applications*, vol. 3, pp. 153-168, 2000.
- [۲۶] M. Dehghan, K. Faez, M. Ahmadi, and M. Shridhar, "Handwritten Farsi (Arabic) word recognition: a holistic approach using discrete HMM," *Pattern Recognition*, vol. 34, pp. 1057-1065, 2001.

- [٢٧] S.-S. Kuo and O. E. Agazzi, "Keyword spotting in poorly printed documents using pseudo 2-D hidden Markov models," *Pattern Analysis and Machine Intelligence, IEEE Transactions on*, vol. 16, pp. 842-848, 1994.
- [٢٨] J. E. Y. W. T. David, F. R. B. M. Maire, and G. Vesom, "Making latin manuscripts searchable using gHMM's," presented at the Advances in Neural Information Processing Systems 17, 2005.
- [٢٩] J. A. Rodriguez, F. Perronnin, G. Sánchez, and J. Lladós, "Unsupervised writer style adaptation for handwritten word spotting," in *Pattern Recognition. ICPR. 19th International Conference on*, 2008, pp. 1-4.
- [٣٠] V. Frinken, A. Fischer, and H. Bunke, "A novel word spotting algorithm using bidirectional long short-term memory neural networks," in *Artificial Neural Networks in Pattern Recognition*, ed: Springer, 2010, pp. 185-196.
- [٣١] M. Khayyat, L. Lam, and C. Y. Suen, "Learning-based word spotting system for Arabic handwritten documents," *Pattern Recognition*, vol. 47, pp. 1021-1030, 2014.
- [٣٢] G. R. Ball, S. N. Srihari, and H. Srinivasan, "Segmentation-based and segmentation-free methods for spotting handwritten arabic words," presented at the Tenth International Workshop on Frontiers in Handwriting Recognition, Suvisoft, 2006.
- [٣٣] S. Srihari, H. Srinivasan, P. Babu, and C. Bhole, "Handwritten arabic word spotting using the cedarabic document analysis system," in *Proc. Symposium on Document Image Understanding Technology (SDIUT-05)*, 2005, pp. 123-132.
- [٣٤] J. Almazán, A. Gordo, A. Fornés, and E. Valveny, "Segmentation-free word spotting with exemplar svms," *Pattern Recognition*, vol. 47, pp. 3967-3978, 2014.
- [٣٥] J. Sánchez, F. Perronnin, T. Mensink, and J. Verbeek, "Image classification with the Fisher vector: Theory and practice," *International journal of computer vision*, vol. 105, pp. 222-245, 2013.
- [٣٦] R. Shekhar and C. Jawahar, "Word image retrieval using bag of visual words," in *Document Analysis Systems (DAS), 10th IAPR International Workshop on*, 2012, pp. 297-301.
- [٣٧] S. Dey, A. Nicolaou, J. Lladós, and U. Pal, "Evaluation of the Effect of Improper Segmentation on Word Spotting," *arXiv preprint arXiv:1604.06243*, 2016.
- [٣٨] R. Manmatha, C. Han, and E. M. Riseman, "Word spotting: A new approach to indexing handwriting," in *Computer Vision and Pattern Recognition. Proceedings CVPR'96, IEEE Computer Society Conference on*, 1996, pp. 631-637.
- [٣٩] A. Fischer, A. Keller, V. Frinken, and H. Bunke, "HMM-based word spotting in handwritten documents using subword models," in *Pattern recognition (icpr), 20th international conference on*, 2010, pp. 3416-3419.
- [٤٠] J. Rodriguez-Serrano and F. Perronnin, "A model-based sequence similarity with application to handwritten word spotting," *Pattern Analysis and Machine Intelligence, IEEE Transactions on*, vol. 34, pp. 2108-2120, 2012.
- [٤١] J. Almazán, "Learning to represent handwritten shapes and words for matching and recognition," PHD, Universitat Autònoma de Barcelona - UAB Barcelona, 2014.

- [۴۲] A. Abidi, A. Jamil, I. Siddiqi, and K. Khurshid, "Word spotting based retrieval of urdu handwritten documents," in *Frontiers in Handwriting Recognition (ICFHR), 2012 International Conference on*, 2012, pp. 331-336.
- [۴۳] م. علی‌آبادی, ج. صدری, س. م. رضوی, "یک روش ترکیبی برای جست و جوی کلمات در کتب تاریخی دست نویس فارسی," اولین کنفرانس بازشناسی الگو و پردازش تصویر ایران, ۱۳۹۱.
- [۴۴] S. Yao, Y. Wen, and Y. Lu, "HoG based two-directional Dynamic Time Warping for handwritten word spotting," in *Document Analysis and Recognition (ICDAR), 13th International Conference on*, 2015, pp. 161-165.
- [۴۵] T. Mondal, N. Ragot, J.-Y. Ramel, and U. Pal, "Flexible Sequence Matching Technique: An Effective Learning-free Approach For word-spotting," *Pattern Recognition*, vol. 60, pp. 596-612, 2016.
- [۴۶] J. A. Rodríguez-Serrano and F. Perronnin, "Handwritten word-spotting using hidden Markov models and universal vocabularies," *Pattern Recognition*, vol. 42, pp. 2106-2116, 2009.
- [۴۷] A. Fischer, A. Keller, V. Frinken, and H. Bunke, "Lexicon-free handwritten word spotting using character HMMs," *Pattern Recognition Letters*, vol. 33, pp. 934-942, 2012.
- [۴۸] V. Frinken, A. Fischer, R. Manmatha, and H. Bunke, "A novel word spotting method based on recurrent neural networks," *Pattern Analysis and Machine Intelligence, IEEE Transactions on*, vol. 34, pp. 211-224, 2012.
- [۴۹] S. Wshah, G. Kumar, and V. Govindaraju, "Script independent word spotting in offline handwritten documents based on hidden markov models," in *Frontiers in Handwriting Recognition (ICFHR), International Conference on*, 2012, pp. 14-19.
- [۵۰] Z. S. Harris, "Distributional structure" *Word*, vol 10, pp 146-162, 1954.
- [۵۱] G. Csurka, C. Dance, L. Fan, J. Willamowski, and C. Bray, "Visual categorization with bags of keypoints," in *Workshop on statistical learning in computer vision, ECCV*, 2004, pp. 1-2.
- [۵۲] F. Perronnin, "Universal and adapted vocabularies for generic visual categorization," *Pattern Analysis and Machine Intelligence, IEEE Transactions on*, vol. 30, pp. 1243-1256, 2008.
- [۵۳] D. Aldavert, M. Rusiñol, R. Toledo, and J. Lladós, "A study of Bag-of-Visual-Words representations for handwritten keyword spotting," *International Journal on Document Analysis and Recognition (IJ DAR)*, vol. 18, pp. 223-234, 2015.
- [۵۴] L. Rothacker, M. Rusinol, and G. A. Fink, "Bag-of-features HMMs for segmentation-free word spotting in handwritten documents," in *Document Analysis and Recognition (ICDAR), 2013 12th International Conference on*, 2013, pp. 1305-1309.
- [۵۵] M. W. Sagheer, N. Nobile, C. L. He, and C. Y. Suen, "A novel handwritten urdu word spotting based on connected components analysis," in *Pattern Recognition (ICPR), 20th International Conference on*, 2010, pp. 2013-2016.
- [۵۶] M. Kassis and J. El-Sana, "Automatic Synthesis of Historical Arabic Text for Word-Spotting," in *Document Analysis Systems (DAS), 2016 12th IAPR Workshop on*, 2016, pp 239-244.

- [57] T. Konidakis, B. Gatos, K. Ntzios, I. Pratikakis, S. Theodoridis, and S. J. Perantonis, "Keyword-guided word spotting in historical printed documents using synthetic data and user feedback," *International Journal of Document Analysis and Recognition (IJ DAR)*, vol. 9, pp. 167-177, 2007.
- [58] J. A. Rodriguez-Serrano and F. Perronnin, "Synthesizing queries for handwritten word image retrieval," *Pattern Recognition*, vol. 45, pp. 3270-3276, 2012.
- [59] Y. Leydier, F. Lebourgeois, and H. Emptoz, "Text search for medieval manuscript images," *Pattern Recognition*, vol. 40, pp. 3552-3567, 2007.
- [60] Y. Leydier, A. Ouji, F. LeBourgeois, and H. Emptoz, "Towards an omnilingual word retrieval system for ancient manuscripts," *Pattern Recognition*, vol. 42, pp. 2089-2105, 2009.
- [61] J. Almazán, A. Gordo, A. Fornés, and E. Valveny, "Word spotting and recognition with embedded attributes," *Pattern Analysis and Machine Intelligence, IEEE Transactions on*, vol. 36, pp. 2552-2566, 2014.
- [62] J. A. Rodriguez-Serrano, F. Perronnin, and F. Meylan, "Label embedding for text recognition," in *Proc. BMVC*, 2013.
- [63] G. Sfikas, A. P. Giotis, G. Louloudis, and B. Gatos, "Using attributes for word spotting and recognition in polytonic greek documents," in *Document Analysis and Recognition (ICDAR), 13th International Conference on*, 2015, pp. 686-690.
- [64] S. K. Ghosh and E. Valveny, "A Sliding Window Framework for Word Spotting Based on Word Attributes," in *Pattern Recognition and Image Analysis: 7th Iberian Conference, IbPRIA, Santiago de Compostela, Spain, June 17-19, Proceedings*, R. Paredes, J. S. Cardoso, and X. M. Pardo, Eds., ed Cham: Springer International Publishing, 2015, pp. 652-661.
- [65] S. Sudholt and G. A. Fink, "PHOCNet: A Deep Convolutional Neural Network for Word Spotting in Handwritten Documents," *arXiv preprint arXiv:1604.00187*, 2016.
- [66] C. H. Lampert, "Learning to detect unseen object classes by between-class attribute transfer," in *IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*. 2009, pp. 951-958.
- [67] C. H. Lampert, H. Nickisch, and S. Harmeling, "Attribute-based classification for zero-shot visual object categorization," *Pattern Analysis and Machine Intelligence, IEEE Transactions on*, vol. 36, pp. 453-465, 2014.
- [68] J. Almazán, A. Gordo, A. Fornés, and E. Valveny, "Handwritten Word Spotting with Corrected Attributes," in *IEEE International Conference on Computer Vision (ICCV)*, 2013.
- [69] J. A. Rodriguez-Serrano and F. C. Perronnin, "Label-embedding for text recognition," ed: Google Patents, 2015.
- [70] J. A. Rodriguez-Serrano, A. Gordo, and F. Perronnin, "Label embedding: A frugal baseline for text recognition," *International Journal of Computer Vision*, vol. 113, pp. 193-207, 2015.
- [71] S. Bengio, J. Weston, and D. Grangier, "Label Embedding Trees for Large Multi-Class Tasks," in *NIPS*, 2010.

- [٧٢] Z. Akata, F. Perronnin, Z. Harchaoui, and C. Schmid, "Label-embedding for attribute-based classification," in *Computer Vision and Pattern Recognition (CVPR), 2013 IEEE Conference on*, 2013, pp. 819-826.
- [٧٣] A. Farhadi, I. Endres, D. Hoiem, and D. Forsyth, "Describing objects by their attributes," in *Computer Vision and Pattern Recognition, CVPR. IEEE Conference on*, 2009, pp. 1778-1785.
- [٧٤] S. Lazebnik, C. Schmid, and J. Ponce, "Beyond bags of features: Spatial pyramid matching for recognizing natural scene categories," in *Computer Vision and Pattern Recognition, IEEE Computer Society Conference on*, 2006, pp. 2169-2178.
- [٧٥] J. Sánchez, F. Perronnin, and T. De Campos, "Modeling the spatial layout of images beyond spatial pyramids," *Pattern Recognition Letters*, vol. 33, pp. 2216-2223, 2012.
- [٧٦] D. R. Hardoon, S. Szedmak, and J. Shawe-Taylor, "Canonical correlation analysis: An overview with application to learning methods," *Neural computation*, vol. 16, pp. 2639-2664, 2004.
- [٧٧] M. B. Blaschko and C. H. Lampert, "Correlational spectral clustering," in *Computer Vision and Pattern Recognition, CVPR. IEEE Conference on*, 2008, pp. 1-8.
- [٧٨] R. Ebrahimpour, M. Amini, and A. V. Shams, "A New Combination Method Based on Different Representation of Data," *International Journal of Hybrid Information Technology*, vol. 4, pp. 51-60, 2011.
- [٧٩] R. Ebrahimpour, R. D. Vahid, and B. M. Nezhad, "Decision Templates with Gradient based Features for Farsi Handwritten Word Recognition," *International Journal of Hybrid Information Technology*, vol. 4, pp. 1-12, 2011.
- [٨٠] R. Ebrahimpour, S. Sarhangi, and F. Sharifzadeh, "Mixture of Experts for Persian handwritten word recognition," *Iranian Journal of Electrical & Electronic Engineering*, vol. 7, pp. 217-224, 2011.
- [٨١] Z. Imani, A. Ahmadyfard, A. Zohrevand, and M. Alipour, "Offline handwritten Farsi cursive text recognition using hidden Markov models," in *2013 8th Iranian Conference on Machine Vision and Image Processing (MVIP)*, 2013, pp. 75-79.
- [٨٢] R. Ebrahimpour, S. Sarhangi, M. Javadi, and F. Sharifzadeh, "Handwritten Word Recognition by Multiple Classifiers: A Divide-and-Conquer Approach," *International Journal of Applied Research on Information Technology and Computing*, vol. 2, pp. 9-21, 2011.
- [٨٣] R. Manmatha and J. L. Rothfeder, "A scale space approach for automatically segmenting words from historical handwritten documents," *Pattern Analysis and Machine Intelligence, IEEE Transactions on*, vol. 27, pp. 1212-1225, 2005.
- [٨٤] M. R. Hashemi, O. Fatemi, and R. Safavi, "Persian cursive script recognition," in *Document Analysis and Recognition, Proceedings of the Third International Conference on*, 1995, pp. 869-873.
- [٨٥] A. Zahour, B. Taconet, P. Mercy, and S. Ramdane, "Arabic hand-written text-line extraction," in *Document Analysis and Recognition, 2001. Proceedings. Sixth International Conference on*, 2001, pp. 281-285.

- [٨٢] A. Zahour, L. Likforman-Sulem, W. Boussalaa, and B. Taconet, "Text line segmentation of historical arabic documents," in *Document Analysis and Recognition, ICDAR. Ninth International Conference on*, 2007, pp. 138-142.
- [٨٣] M. Arivazhagan, H. Srinivasan, and S. Srihari, "A statistical approach to handwritten line segmentation," *Document Recognition and Retrieval XIV, Proceedings of SPIE, San Jose, CA*, pp. 6500T-1, 2007.
- [٨٤] V. Papavassiliou, T. Stafylakis, V. Katsouros, and G. Carayannis, "Handwritten document image segmentation into text lines and words," *Pattern recognition*, vol. 43, pp. 369-377, 2010.
- [٨٥] Z. Shi and V. Govindaraju, "Line separation for complex document images using fuzzy runlength," in *Document Image Analysis for Libraries, International Workshop on*, 2004, pp.306.
- [٩٠] P. P. Roy, U. Pal, and J. Lladós, "Morphology based handwritten line segmentation using foreground and background information," in *International Conference on Frontiers in Handwriting Recognition*, 2008, pp. 241-246.
- [٩١] L. Likforman-Sulem ,A. Hanimyan, and C. Faure, "A Hough based algorithm for extracting text lines in handwritten documents," in *Document Analysis and Recognition, Proceedings of the Third International Conference on*, 1995, pp. 774-777.
- [٩٢] G. Louloudis, B. Gatos, I .Pratikakis, and K. Halatsis, "A block-based Hough transform mapping for text line detection in handwritten documents," in *Tenth International Workshop on Frontiers in Handwriting Recognition*, 2006.
- [٩٣] Y. A. Shmjitsourouka and T. Yoshika WA, "The Segmentation of a text line for a handwritten unconstrained document using thinning algorithm", *the 7th International workshop on frontiers in handwritten recognition*, pp 505-510, 2000.
- [٩٤] Y. Li, Y. Zheng, D. Doermann, and S. Jaeger, "Script-independent text line segmentation in freestyle handwritten documents," *Pattern Analysis and Machine Intelligence, IEEE Transactions on*, vol. 30, pp. 1313-1329, 2008.
- [٩٥] S. S. Bukhari, F. Shafait, and T. M. Breuel, "Script-independent handwritten textlines segmentation using active contours," in *Document Analysis and Recognition, ICDAR'09. 10th International Conference on*, 2009, pp. 446-450.
- [٩٦] F. Yin and C.-L. Liu, "A variational bayes method for handwritten text line segmentation," in *Document Analysis and Recognition, ICDAR'09. 10th International Conference on*, 2009, pp. 436-440.
- [٩٧] M. Ziaratban and K. Faez, "An Adaptive Script-Independent Block-Based Text Line Extraction," in *ICPR*, 2010, pp. 249-252.
- [٩٨] M. Ziaratban and F. Karim, "Adaptive Script-Independent Text Line Extraction," *IEICE transactions on information and systems*, vol. 94, pp.٢٠١١ ,٨٧٧-٨٢٢ .
- [٩٩] A. Alaei, U. Pal, and P. Nagabhushan, "A new scheme for unconstrained handwritten text-line segmentation," *Pattern Recognition*, vol. 44, pp. 917-928, 2011.
- [١٠٠] N. Otsu, "A threshold selection method from gray-level histograms," *Automatica*, vol. 11, pp. 23-27, 1975.

Abstract

Word spotting is a way to index, search and retrieve words in document images. Word spotting aims to find multiple occurrences of a query word in document images. Word spotting is one of the major challenges in the field of document image analysis which has recently received many attentions. Several word spotting systems have been introduced on Latin script, but a limited number of works exist on Arabic and Persian documents and most of them only proposed for printed documents. The goal of this thesis is to present an end to end multiwriter Persian word spotting system which accepts, independently of lexicon, both word images and text strings.

The number of classes in this problem is very high. Moreover very large intra-class variability due to different writing styles, illumination, typography, etc, can make the same word look very different. For this reason, a classification method using attributes is needed. Attribute-based classification by sharing information between classes makes out-of-vocabulary (OOV) spotting possible. In this approach, both word images and text strings embeds in a common vectorial subspace. This is achieved by a combination of label embedding, attributes learning, and a common subspace regression. In this subspace, images and strings that represent the same word are close together, allowing one to cast recognition and retrieval tasks as a nearest neighbor problem.

In the first step, an attribute-based method is proposed to retrieve isolated handwritten Persian words. In this method we suggest two methods for representing Persian words using attributes. These methods are inspired from two facts, first the appearance of a Persian letter changes depending on its position in a word. Second, Persian letters can be categorized based their shape of main part. To evaluate the performance of proposed methods, two isolated world Persian handwritten datasets, Iranshar and Farsa, are used. Mean average precision (MaP) for presented attributes in word spotting is 95.67%, 96.25% and precision for handwritten word recognition is 96.15% and 97.55% respectively.

Then we introduce a method for segmenting handwritten documents to constitute lines. Using a new approach we slide a window on each line of a document. The proposed sliding window is based on connected components, so it can overcome to the problem of basic methods based on sliding window such as determining window size and movement step. FHT dataset is used to evaluate the proposed line-based method. MaP criterion in this experiment for query by example and query by string is 73.67 and 81.02% respectively.

Finally, a two-pass segmentation-free method is proposed for spotting of words in handwritten documents. In the first pass a number of candidate regions for the query word are selected in the document. These regions are extracted based on the proposed spotting method at connected component level. Then considering neighboring connected components to the candidate regions, the regions which likely match the query word are extracted. Evaluation of this method is also performed using FHT dataset and MaP 79.45% for text query and 70.25% for image query are obtained.

It is noteworthy that we trained the proposed line based and segmentation free methods using Farsa and Iranshahr datasets whereas FHT dataset is used for test. This illustrates that the proposed methods are robust to writing style. Furthermore, any word from class of queries is not used in training phase so the proposed methods can spot and retrieve unseen words.

Keywords: Word spotting, Persian documents, label-embedding, attribute-based classification, connected component.



Shahrood University of Technology
Faculty of Computer Engineering and IT

Ph.D. Thesis in Artificial Intelligence

Persian word spotting based on attributes

By:

Majid Iranpour Mobarakeh

Supervisor:

Dr. Alireza Ahmadyfard

Adviser:

Dr. Hamid Hasanpour

February 2017

۱۱۲