

بازشناسی متن فارسی قدیمی با استفاده از توصیفگرهای ویولت در روش منهن

مأنده رضائی فر^۱، افشین ابراهیمی^۲

^۱دانشجو دانشگاه صنعتی سهند تبریز، دانشکده برق، m_rezaeifar@sut.ac.ir

^۲دانشیار دانشگاه صنعتی سهند تبریز، دانشکده برق، a_ebrahimi@sut.ac.ir

چکیده - در این مقاله به ارائه روش جدیدی برای بازشناسی اسناد چاپی قدیمی فارسی مبتنی بر شکل کلی کلمات با استفاده از بسته مویک و تبدیل بسته مویک پرداخته ایم. بنابه ویژگی‌های خاص زبان فارسی و عربی بازشناسی این زبان‌ها بر اساس شکل کلی کلمات نسبت به زبان‌های لاتین و چینی متفاوت تر خواهد بود. کارهای انجام شده، بیشتر بر متون دست نویس متمرکز بوده است در حالی که روش ارائه شده در این مقاله برای بازشناسی یک کتاب چاپی قدیمی فارسی است که دارای فونت خاص و چاپ سنگی است. روش ارائه شده بر پایه شکل کلی کلمات و استفاده از روش منهن روی توصیفگرهای مویک می‌باشد. با تهیهی واژه‌نامه مناسب از کتاب قدیمی فارسی که از ۳۷۸۶۱ زیر کلمه مشابه و غیر مشابه تشکیل شده است و استفاده از کد نقاط و فشرده سازی بردارهای ویژگی، نرخ بازشناسی مناسبی با مقدار ۸۸٪ ارائه شده است. کلید واژه- بازشناسی متن قدیمی، بسته تبدیل مویک، تبدیل مویک، روش منهن، نویسه خوان نوری.

و زیر حروف تجزیه شده و سپس بازشناسی صورت می‌گیرد. در برخی از موارد از ترکیب این دو روش نیز استفاده می‌شود. در این مقاله، ما با استفاده از روش مبتنی بر شکل کلی کلمات و همچنین استفاده از دانش‌های جانبی مانند نقاط، علائم و برخی قواعد دستوری ساده^۱ زبان فارسی که در ادامه به توضیح آن می‌پردازیم، به بازشناسی این کتاب چاپی قدیمی پرداخته ایم.

کارایی یک سیستم نویسه خوان نوری (OCR) وابسته به این است که ویژگی‌هایی که هر الگو را نمایش می‌دهند، چگونه تعریف می‌شوند. روش‌های متعددی برای بازشناسی متون قدیمی در اسناد لاتین و چینی پیشنهاد شده است. ولی فضای تحقیق در اسناد چاپی قدیمی فارسی و همچنین در اسناد عربی مورد نیاز است. علاوه بر این، در ۵۰ سال گذشته بیشتر تلاش‌ها در بازشناسی متون چاپی عربی صورت گرفته است و محک آن‌ها بر روی دیتابیس‌های گذشته بوده است [۶]. تلاش بر این بوده است که با روش بکار رفته در این مقاله نرخ بازشناسی یک متن قدیمی را افزایش دهیم.

۲- کاراکترهای زبان فارسی:

نحوه نگارش زبان فارسی نسبت به زبان‌های دیگر از تفاوت عمده‌ای برخوردار است. این تفاوت در بازشناسی متون فارسی مشکلات مخصوص به خود را دارد. حالت‌های مختلف حروف الفبای فارسی با توجه به موقعیت مکانی آن‌ها در کلمات در

۱- مقدمه

از کاربردهای بازشناسی الگوهای تصویری، نویسه خوان نوری حروف می‌باشد. بازشناسی متن، یکی از زمینه‌هایی است که در چند دهه اخیر، توجه زیادی را به خود اختصاص داده است [۴-۱]. یک سیستم نویسه خوان نوری، تصویر اسناد را به اسنادی به صورت کدهای کاراکتری تبدیل می‌کند. این عمل در کاهش حجم اسناد تأثیر بسیار زیادی دارد و همچنین در ذخیره سازی اسناد در تعداد زیاد و انتقال این اسناد بوسیله شبکه‌های کامپیوتری بسیار مفید می‌باشد. یک سیستم نویسه خوان نوری دارای مراحل مختلفی بوده و کارایی این سیستم به دقت و عملکرد هر یک از این مراحل بستگی دارد. این مراحل شامل روبش سند، پیش پردازش، استخراج ویژگی، تشخیص کاراکترها و پس پردازش می‌باشد. سیستم‌های بازشناسی متن براساس داده ورودی دریافتی به شاخه‌های متن چاپی و متن دست نویس تقسیم می‌شوند. اولین مرحله در بازشناسی متن، دریافت داده و تبدیل آن به شکل دیجیتالی قابل استفاده در کامپیوتر است. بازشناسی متن دستنویس به دو روش برخط و برون خط است. سیستم بازشناسی هر یک متفاوت می‌باشد. ورودی سیستم بازشناسی برون خط، تصویر روبش شده^۲ متن دستنویس و اسناد چاپی است [۵]. روش‌های بازشناسی متن فارسی، به دو دسته کلی تقسیم می‌شوند. بازشناسی بر اساس شکل کلی کلمات و بازشناسی بر پایه جداسازی. در روش اول، کل زیر کلمه به عنوان یک شکل در نظر گرفته شده و سپس تشخیص داده می‌شود و در روش دوم، ابتدا زیر کلمه ورودی به حروف

جدول (۱) حالت‌های مختلف حروف الفبای فارسی نشان داده شده است.

جدول (۱) حالت‌های مختلف حروف الفبای فارسی [۷]

Isolated	Initial	Medial	Final	Roman	Name	Isolated	Initial	Medial	Final	Roman	Name
ا	ا	ا	ا	ā	alef	ص	ص	ص	ص	ṣ	sād
ب	ب	ب	ب	b	be	ض	ض	ض	ض	ḍ	zād
پ	پ	پ	پ	p	pe	ط	ط	ط	ط	t	tā
ت	ت	ت	ت	t	te	ظ	ظ	ظ	ظ	z	zā
ث	ث	ث	ث	th	se	ع	ع	ع	ع	ʿ	ayn
ج	ج	ج	ج	j	jim	غ	غ	غ	غ	gh	ghayn
چ	چ	چ	چ	ch	che	ف	ف	ف	ف	f	fe
ح	ح	ح	ح	h	he	ق	ق	ق	ق	q	qáf
خ	خ	خ	خ	kh	khe	ک	ک	ک	ک	k	káf
د	د	د	د	d	dál	گ	گ	گ	گ	g	gáf
ذ	ذ	ذ	ذ	dh	zál	ل	ل	ل	ل	l	lám
ر	ر	ر	ر	r	re	م	م	م	م	m	mím
ز	ز	ز	ز	z	ze	ن	ن	ن	ن	n	nún
ژ	ژ	ژ	ژ	zh	zhe	و	و	و	و	v/ú	váv
س	س	س	س	s	sin	ه	ه	ه	ه	h	he
ش	ش	ش	ش	sh	shin	ی	ی	ی	ی	y/i	ye

۳= واژه‌نامه:

در این پروژه از تصاویر روبش شده یک کتاب چاپی قدیمی به- عنوان پایگاه داده الکترونیکی استفاده کرده‌ایم و هدف ما بازخوانی متن و تبدیل آن به متن قابل ویرایش در Word است. صفحات کتاب را با دقت dpi600 در کتابخانه مرکزی اسکن کرده‌ایم. به دلیل قدیمی بودن متن و نیز اسکن این صفحات، وجود نویز، گسستگی و یا چسبندگی کلمات و نقاط به یکدیگر و یا به بدنه‌ی زیر کلمات را خواهیم داشت. برای ایجاد واژه‌نامه ابتدا تمام زیر کلمات تصاویر را جدا کرده و در یک ساختار با فیلد image ذخیره کرده‌ایم و تمام زیر کلمات مشابه و غیرمشابه را جدا نموده‌ایم. برای این کار طبق برنامه نوشته شده معادل متنی هر زیر کلمه به برنامه داده می‌شود اگر در واژه نامه موجود بود، محل آن نمایش داده می‌شود در غیر اینصورت از آخر در پایگاه داده تهیه شده، ذخیره می‌گردد. با این روش توانستیم از مجموع ۳۷۸۶۱ زیر کلمه که همه دارای یک قلم و اندازه ثابت بوده‌اند، تعداد ۱۶۶۹ زیر کلمه غیرمشابه ایجاد نموده و به‌عنوان واژه‌نامه ذخیره نماییم.

۴- تبدیل موجک و بسته تبدیل موجک:

۴-۱ تبدیل موجک:

در سال‌های اخیر تبدیل‌های موجک در پردازش تصویر بر بسیاری از تکنیک‌های پیشرفته چیره شده‌ان یک سطح تبدیل موجک دو بعدی می‌تواند روی هر تصویر f به اندازه $M \times N$ تعریف شود [۹].

در رویکرد مبتنی بر شکل کلی کلمات یا زیر کلمات، توصیفگرها به طور مستقیم از تصویر کلمه یا زیر کلمه استخراج می‌شوند و با استفاده از یک واژه نامه تصویری، کلمه مناسب استخراج می‌شود. در این تحقیق، ویژگی‌های شکل زیر کلمات با استفاده از تبدیل موجک استخراج شده است و بازشناسی زیر کلمات با استفاده از این توصیفگرهای موجک انجام گردیده است. نتایج خوب بدست آمده از تبدیل موجک در پردازش تصویر سبب شده است تا با بکارگیری از محلی کردن و تمرکز در نقطه- ی بخصوص مکانی و جزئیات محلی سازی فرکانسی، تبدیل موجک به ابزار موثری برای استخراج ویژگی و بازشناسی متون تبدیل شود. در بازشناسی کلمات و اعداد موجود در متن به زبان های مختلف، تبدیل موجک به روش‌های مختلف برای استخراج

از مشکلات کد نقاط می‌توان به اتصال چسبندگی نقطه‌های حروف به یکدیگر و یا به بدنهء خود زیر کلمه اشاره کرد. این مشکلات در تصاویر نویزدار سبب ایجاد مساحت بزرگ‌تری خواهد شد به‌طوری‌که به اندازه یک پیکسل می‌رسد و با حذف نویز، آن نقطه و یا نقاط نیز حذف می‌شوند. به همین دلیل با روش‌های مانند مورفولوژی به نازک سازی حروف به بدنه و یا نازک سازی دو زیر کلمه به یکدیگر پرداخته ایم.

وجود سرکش در بعضی از حروف مانند "ک" و "گ"، علامت مد روی حرف الف، همین‌طور حروف دسته‌داری مانند "ط" و "ظ" تفاوت در شکل‌های بعضی حروف در زبان فارسی را نشان می‌دهند. همچنین تفاوت شکل در برخی حروف دیگر علاوه بر تعداد در محل قرار گرفتن نقاط آن‌ها نیز هست مانند {"ب"، "پ"، "ت"، "ث"}.

یکی از خصوصیات متمایز کننده در متون فارسی خط زمینه و یا خط مبنا است. خط زمینه به خطی گفته می‌شود که بیشترین تعداد نقاط متن (پیکسل‌های سیاه) را دارا است و قسمت پایینی بدنه‌ی کاراکترها را در یک خط متن دنبال می‌کند [۸].

پایین گذر h و فیلتر بالاگذر g متناظر با توابع پایه موجک مختلف، متفاوت می‌باشد.

L1	H1	L2	H2
V1	D1	V2	D2
L3	H3	L4	H4
V3	D3	V4	D4

شکل (۲): ساختار تجزیه سطح دوم تبدیل بسته موجک [۱۱]

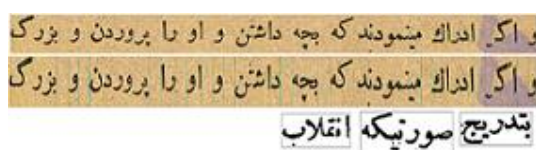
۵- الگوریتم پیشنهادی:

ابتدا تصویر ورودی را دریافت کرده و به جداسازی هر یک از خطوط آن پرداخته‌ایم. در کارهای گذشته یک ردیف به عنوان خط مبنا شناخته می‌شد، در حالیکه خط فارسی دارای عرض قلم می‌باشد به تعریف محدوده‌ی مناسب برای خط مبنا پرداخته‌ایم. از افکنش افقی برای پیدا کردن نواحی سفید متن و با بسط دادن خط مبنا از بالا و پایین و تعیین یک آستانه مناسب محدوده خط را بدست آورده‌ایم و خطوط را از یکدیگر جدا نموده‌ایم [۱۲]. در شکل (۲) تصاویر تعیین خطوط از صفحات را نمایش داده شده است.



شکل (۲): تصاویر تعیین خطوط صفحات کتاب قدیمی

سپس به جداسازی کلمات و زیر کلمات هر خط پرداخته و با افکنش عمودی فواصل سفید بین ستون‌های کلمه را پیدا کرده و به جداسازی کلمات هر خط پرداخته‌ایم [۱۲]. شکل (۳) جداسازی کلمات و زیر کلمات کتاب قدیمی فارسی را نشان می‌دهد.



شکل (۳): جداسازی کلمات و زیر کلمات کتاب قدیمی

ویژگی و بازشناسی آنها استفاده شده است. در هریک از این روش‌ها تبدیل موجک مادر استفاده شده نیز متفاوت است. انواع تبدیل موجک مورد استفاده در بازشناسی کاراکتر در کارهای پیشین را می‌توان به: تبدیل موجک گسسته دو بعدی، تبدیل موجک گسسته مختلط، تبدیل موجک پیوسته جهت دار، موجک چندگانه متعامد، گشتاورهای موجک، موجک بسته ای، موجک بسته‌ای M -بند است که ساده ترین راه برای محاسبه‌ی یک تبدیل موجک گسسته (DWT) دو بعدی یک تصویر، بکار بردن یک تبدیل یک بعدی روی سطرها و ستون های تصویر بطور مجزا و سپس انجام نمونه برداری کاهش است.

۴-۲ بسته تبدیل موجک:

روش تبدیل بسته موجک یک تعمیم تجزیه‌ی موجک است که یک آنالیز سیگنال غنی‌تر را پیشنهاد می‌دهد. در روند‌های تجزیه‌ی موجک متعامد دو بعدی، مرحله‌ی کلی (جامع) ضرائب تخمین را به چهار زیر بخش تقسیم می‌کند. بعد از تقسیم‌بندی، یک زیر تصویر از ضرائب نرم (پایین گذر) و سه زیر تصویر متناظر با ضرائب جزئیات در جهت‌های عمودی و افقی و مورب بدست می‌آید. این تقسیم‌بندی روی زیر تصویر پایین گذر دنبال می‌شود در حالیکه جزئیات متوالی دوباره آنالیز نمی‌شوند [۹].

در موقعیت بسته‌ی موجک متناظر، هر زیر تصویر ضرائب جزئیات به چهار بخش با استفاده از همان روش مانند تقسیم‌بندی زیر تصویر پایین گذر تجزیه می‌شود [۹] شکل (۱) ساختار تجزیه سطح دوم تبدیل بسته موجک است. تبدیل بسته موجک گسسته‌ی دو بعدی برای یک تصویر گسسته‌ی $N \times M$ بعدی به نام A ، بالای سطح $(p \leq \min(\log_2 N, \log_2 M))p + 1$ به طور بازگشتی بر حسب ضرائب در سطح p به صورت زیر تعریف می‌شود [۱۰]

$$\begin{aligned} C_{k(i,j)}^{p+1} &= \sum_m \sum_n h(m) h(n) C_{k(m+\tau_i n+\tau_j)}^p \\ C_{k+1(i,j)}^{p+1} &= \sum_m \sum_n h(m) g(n) C_{k(m+\tau_i n+\tau_j)}^p \\ C_{k+2(i,j)}^{p+1} &= \sum_m \sum_n g(m) h(n) C_{k(m+\tau_i n+\tau_j)}^p \\ C_{k+3(i,j)}^{p+1} &= \sum_m \sum_n g(m) g(n) C_{k(m+\tau_i n+\tau_j)}^p \end{aligned}$$

که $C_{k(i,j)}^p$ تصویر A است. در هر سطح، تصویر

به چهار زیر تصویر به اندازه‌ی یک چهارم تصویر C_k^p ، به نام‌های C_{k+3}^{p+1} ، C_{k+2}^{p+1} ، C_{k+1}^{p+1} ، C_k^{p+1} تجزیه می‌شود. فیلتر

۵-۲ تجزیه تصویر زیر کلمه نرمالیزه شده با بسته تبدیل موجک:

ضرائب زیر باند سطح [۱۵] به عنوان ویژگی به صورت یک بردار ۷۲۹ تایی در سیملت ۸ استفاده می‌شود. زیر باند به صورت یک ماتریس 27×27 تایی بدست می‌آید که همه‌ی سطرهاى آن به ترتیب در یک بردار ستونی زیر هم اضافه می‌شوند و بردار ویژگی را به صورت یک بردار 729×1 تشکیل می‌دهند. اندازه تصویر ورودی در بسته تبدیل موجک به ابعاد 64×64 نرمالیزه شده‌اند.

۵-۳ شبکه تشخیص نقاط:

برای آموزش ویژگی‌های بکار رفته از شبکه عصبی دو لایه «FeedForward» با ۴۰ نورون در لایه خروجی و توابع انتقال «tansig» در لایه مخفی و «purelin» در لایه خروجی استفاده شده است که از تعداد ۵۰۰ تکرار و نرخ یادگیری ۳ استفاده نموده‌ایم. شبکه دارای ۵ کلاس ورودی و ۵ کلاس خروجی خواهد بود که کلاس ۱ برای تک نقطه، کلاس ۲ برای دونقطه، کلاس ۳ برای سه نقطه بالا، کلاس ۴ برای سه نقطه پایین و کلاس ۵ برای علامت مد الف خواهد بود. جدول (۴) نمونه زیر کلمات و کد نقاط آن‌ها را نشان داده شده است.

جدول (۴) نمونه زیر کلمات و کد نقاط آن‌ها

زیر کلمه	کد نقاط زیر کلمه	زیر کلمه	کد نقاط زیر کلمه
پیشین	Pysyn	پر	p
محیط	Y	سرخیه	nyn

۵-۴ بازشناسی زیر کلمات:

فاصله ویژگی یک زیر کلمه از پایگاه داده با کل ویژگی زیر کلمات پایگاه داده محاسبه می‌شود و کمترین آن را به عنوان محل زیر کلمه بازشناسی شده انتخاب کرده‌ایم. یکبار از کد نقاط برای بازشناسی و تعیین محل زیر کلمات استفاده کرده‌ایم و در حالت دوم عملیات مورفولوژی را به آن اضافه نموده‌ایم و نتایج بدست آمده را مورد بررسی قرار داده‌ایم. فلوچارت روش پیشنهادی در شکل (۵) نشان داده شده است.

به منظور انجام پیش پردازش‌های لازم به جهت از بین بردن همپوشانی و مشکل چسبندگی زیر کلمات به یکدیگر از تکنیک برچسب زنی استفاده نموده‌ایم [۱۳] همچنین با تعیین $threshold \leq 2$ توانستیم گسستگی بین کلماتی که روی خط زمینه متمرکز هستند را برطرف نماییم شکل (۴) رفع گسستگی در کلمه با روش برچسب زنی را نشان می‌دهد.

انقلاب ۱: انقلاب

انقلاب انقلاب

شکل (۴): رفع گسستگی در زیر کلمه با تکنیک برچسب زنی

۵-۱ تجزیه تصویر زیر کلمه نرمالیزه شده با تبدیل موجک:

در الگوریتم تبدیل موجک هدف استفاده از ضرائب زیر باند تقریب سطح دوم حاصل از تجزیه موجک به عنوان ویژگی برای زیر کلمات پایگاه داده است. اندازه تصویر ورودی را به ابعاد 40×40 نرمالیزه کرده‌ایم [۱۱]. تصویر ورودی نرمالیزه شده وارد فیلتر موجک هار می‌شود و تا دو سطح تجزیه می‌گردد. سپس ضرائب تقریب سطح دوم به عنوان ویژگی انتخاب می‌شوند، به این صورت که ضرائب این زیر باند دو بعدی به ابعاد 5×5 در یک بردار اندازه‌ی 1×25 ریخته می‌شوند و بردار ویژگی را تشکیل می‌دهد. جدول (۳) بردارهای ویژگی‌های مختلف تبدیلموجک را نمایش می‌دهد. از ویژگی‌های مختلف تنها ویژگی سیملت ۸، نیاز به کاهش بعد (PCA) دارد [۱۴].

جدول (۳): بردارهای ویژگی‌های مختلف تبدیل موجک

ویژگی لوبولت	Haar	Db1	Db2	Coif1	Sym8 (pca) شده
Haar	۱۶۶۹×۲۵	-	-	-	-
Db1	-	۱۶۶۹×۲۵	-	-	-
Db2	-	-	۱۶۶۹×۴۹	-	-
Coif1	-	-	-	۱۶۶۹×۸۱	-
Sym8	-	-	-	-	۱۶۶۹×۳۲۴/۳۸

(۴۴)

احتیاج دارد باو بدهیم و آن عبارت است از آزادی حرکت و سرآمدهای شخصی و تجربه ها و معلومات و میدان وسیع برای انکشاف قوه فعالیت زیرا بچه میخواد همه چیز را بفهمد و کردار و رفتار یادگیرد یعنی بزرگ شود و کمال یابد. ازینجا میتوان می برد باین که اگر ما بکوشیم زندگانی او را طوری مرتب سازیم که این احتیاجهای او را رفع کند، او را بهترین وجهی تربیت کردیم و او هم معلومات لازمه را یاد خواهد گرفت. آنوقت او قوت و جمال و صحت بدنی **یابا** کسب خواهد کرد چنانکه بدن **ماده**، انسان هم تنها از نان غذا و قوت نمیگیرد بلکه با خشنودی مل نیز تغذیه میشود.

تصویر (۱۱): تصویر بازشناسی شده با ویژگی db1

(۴۴)

احتیاج دارد باو بدهیم و آن عبارت است از آزادی حرکت و سرآمدهای شخصی و تجربه ها و معلومات و میدان وسیع برای انکشاف قوه فعالیت زیرا بچه میخواد همه چیز را بفهمد و کردار و رفتار یاد گیرد یعنی بزرگ شود و کمال یابد. ازینجا میتوان می برد باین **فکند** که اگر ما بکوشیم زندگانی او را طوری مرتب سازیم که این احتیاجهای او را رفع کند، او را بهترین وجهی تربیت کردیم و او هم معلومات لازمه را یاد خواهد گرفت. آنوقت او قوت و جمال و صحت بدنی **لین** کسب خواهد کرد چنانکه بدن **مادی**، انسان هم تنها از نان غذا و قوت نمیگیرد بلکه با خشنودی مل نیز تغذیه میشود.

تصویر (۱۲): تصویر بازشناسی شده با ویژگی wdb1

به دلیل اینکه بازشناسی بر مبنای شکل کلی زیرکلمات بوده است، ابتدا در تصویر اولیه تمام زیرکلمات را شمارش نموده و سپس از روی تصویر اصلی در متن بازشناسی شده زیرکلمات اشتباه را پیدا کرده، شمارش نموده و از کل زیرکلمات کسر می نماییم. فرمول زیر نحوه انجام کار را نشان می دهد.

$$100 \times \frac{\text{تعداد زیرکلمات اشتباه} - \text{کل زیرکلمات}}{\text{کل زیرکلمات}} = \text{درصد}$$

درصدهای بازشناسی بدست آمده در جداول زیر نمایش داده شده است. جدول (۵) جدول ارزیابی موجک سطح ۳ و جدول (۶) جدول ارزیابی بیهته تبدیل موجک سطح ۲ را نمایش می دهد.

جدول (۵) جدول ارزیابی موجک سطح ۳

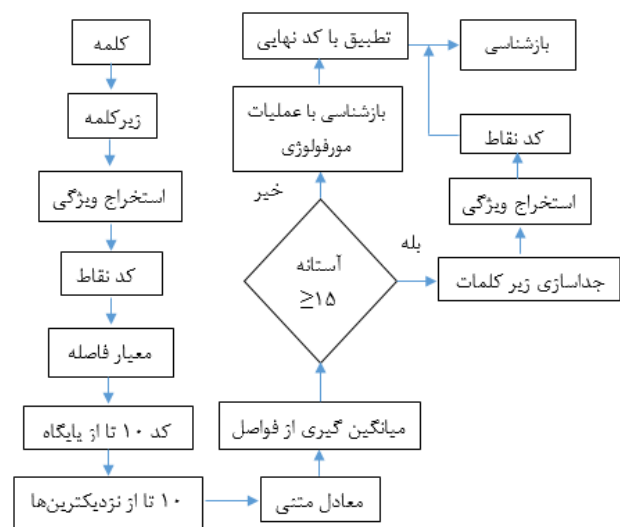
ویژگی ها	Haar	Db1	Db2	Coif1	Sym8
۰۰۰۸.png	۷۸٪	۸۰،۲۲٪	۷۹٪	۷۸٪	۶۷٪
۰۰۱۷.png	۸۲،۶٪	۸۴،۵۲٪	۸۲،۷۸٪	۸۲،۶۱٪	۷۶٪
۰۰۴۶.png	۸۴،۶۹٪	۸۵،۰۲٪	۸۴،۱۹٪	۸۴،۸۵٪	۷۹،۳٪

جدول (۶) جدول ارزیابی بیهته تبدیل موجک سطح ۲

ویژگی ها	WPHaar	WPDb1	WPDb2	WPCoif1	WPSym8
۰۰۰۸.png	۸۳،۱۲٪	۸۵،۲٪	۸۴،۸۴٪	۸۳٪	۷۴،۶۷٪
۰۰۱۷.png	۸۶،۹۵٪	۸۷،۱۳٪	۸۶،۶۰٪	۸۴،۳۴٪	۷۶،۲۳٪
۰۰۴۶.png	۸۷،۳٪	۸۸،۶۸٪	۸۸،۳۵٪	۸۸٪	۸۰٪

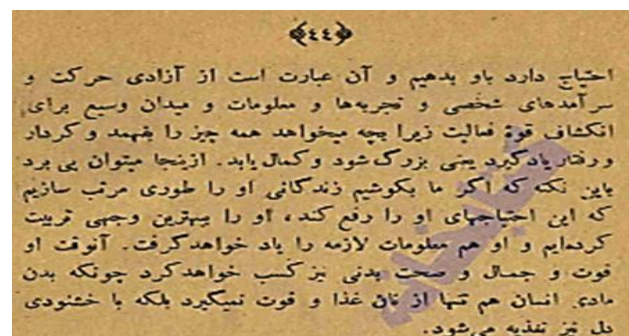
۶- نتیجه گیری:

در این مقاله روش جدیدی با دقت و نرخ بازشناسی مناسب ۸۸٪ برای بازشناسی اسناد قدیمی چاپی فارسی ارائه شد. این روش



شکل (۵): فلوچارت روش منهن

ابتدا کلمه ورودی به زیر کلمه تجزیه گردیده و سپس به استخراج ویژگی های بکار رفته پرداخته شده و کد نقاط آنها بدست آمده است و از معیار فاصله منهن استفاده شده است. کد ده زیر کلمه از پایگاه داده را نشان داده و ده عدد از نزدیکترین زیر کلمات مشابه به زیر کلمه مورد نظر را از نظر ساختار بدنی و کد نقاط نشان می دهد و در نهایت معادل متنی هر کدام از زیر کلمات را می توانیم ببینیم. سپس فاصله های این زیر کلمات بدست آمده و عملیات میانگین گیری و فاصله از میانگین ها را حساب کرده ایم. بزرگترین فاصله تعیین شده اگر از ۱۵ کوچکتر بود به این معنی است که در زیر کلمه مورد نظر چسبندگی و یا گسستگی وجود دارد و با استفاده از عملگرهای انبساطی و انقباضی و حذف کشیدگی، به درستی زیر کلمه مورد نظر را بازشناسی می کنیم و اگر میانگین بزرگترین فاصله، از ۱۵ بزرگتر بود در اینصورت به بازشناسی زیر کلمه بدون عملیات مورفولوژی می پردازیم. تصاویر (۱۰)، (۱۱) و (۱۲) به ترتیب صفحات اصلی کتاب و بازشناسی شده را با ویژگی های db1 و wdb1 نشان می دهد.



تصویر (۱۰): تصویر اصلی کتاب

- A. A. Aburas, and Salem MA Rehiel, "Off-line omni-style handwriting Arabic character recognition system based on wavelet compression," *Arab Research Institute in Sciences & Engineering* vol. ۱۳۵-۱۲۳, pp. ۳, ۲۰۰۷.
- م. شمسی، ع. ر. کناری، و س. شادروان، "روشی نو در تشخیص حروف،" فصلنامه علمی - پژوهشی مهندسی برق مجلسی، شماره سوم، ۱۳۸۸.
- H. Khosravi, and E. Kabir, "A blackboard approach towards integrated Farsi OCR system," *International Journal of Document Analysis and Recognition (IJ DAR)*, vol. ۱۲, no. ۱, pp. ۲۱, ۲۰۰۹.
- J. Shlens, "A tutorial on principal component analysis," *arXiv preprint arXiv: ۱۴۰۴.۱۱۰۰*, ۲۰۱۴.
- J. S. Walker, *A primer on wavelets and their scientific applications*: CRC press, ۲۰۰۸.
- [۱۱] مبتنی بر ویژگی‌های تبدیل موجک و بسته تبدیل موجک در روش منهن بن بوده است. نتایج مناسب الگوریتم نشان می‌دهد که الگوریتم ارائه شده به دلیل خاص و قدیمی بودن نوع فونت بکاررفته و همچنین استفاده از شکل کلی کلمات می‌تواند کارایی مناسبی در بازشناسی از خود داشته باشد. تفاوت موجود در این روش نسبت به روش های دیگر بدلیل الگوریتم ارایه شده و تعیین آستانه مناسب بوده است. کل کلمات مشابه و غیر مشابه بدست آمده ۳۷۸۶۱ بوده است، که توانستیم تعداد ۱۶۶۹ زیر کلمه غیر مشابه را جدا نموده و به عنوان واژه‌نامه استفاده نماییم.
- [۱۲]
- [۱۳]
- [۱۴]
- [۱۵]
- ## مراجع
- [۱] ح. نظام آبادی پور، ا. ا. کبیر، و ر. عزمی، "الگوریتم اصلاح شده جداسازی حروف در متون چاپی با برچسب زدن به کانتور بالایی کلمات،" استقلال، شماره ۱، ۱۳۸۳.
- [۲] A. Al-Shoshan, "Arabic OCR based on image invariants." pp. ۱۵۴-۱۵۰.
- [۳] H. Y. Abdelazim, and M. Hashish, "Arabic reading machine." pp. ۷۴۴-۷۳۳.
- [۴] R. A. E. Kabir, "A new segmentation technique for omnifont Farsi text," *Pattern Recognition Letters*, vol. ۲۲, pp. ۱۰۴-۹۷, ۲۰۰۱.
- [۵] ا. ابراهیمی، "معرفی، مفاهیم، تعاریف و اصطلاحات بازشناسی متن (OCR)" پژوهشنامه نویسه خوان نوری (OCR) فارسی، چاپ اول، ص ۷۷-۱۰۹، ۱۳۸۸.
- [۶] S. Torabzadeh, and R. Safabakhsh, "AUT-PFT: A real world printed Farsi text image dataset." pp. ۲۷۲-۲۶۷.
- [۷] ا. ا. کبیر، و ا. ابراهیمی، "استفاده از شکل کلی زیر-کلمات چاپی در بازیابی تصویر مستندات و بازشناسی متون فارسی،" رساله دکتری مهندسی برق-الکترونیک، دانشگاه تربیت مدرس، ۱۳۸۴.
- [۸] L. Likforman-Sulem, A. Zahour, and B. Taconet, "Text line segmentation of historical documents: a survey," *International Journal of Document Analysis and Recognition (IJ DAR)*, vol. ۹, no. ۲-۴, pp. ۱۲۳-۱۳۸, ۲۰۰۷.
- [۹] A. Broumandnia, J. Shanbehzadeh, and M. R. Varnosfaderani, "Persian/arabic handwritten word recognition using M-band packet wavelet transform," *Image and Vision Computing*, vol. ۲۶, no. ۶, pp. ۸۴۲-۸۲۹, ۲۰۰۸.
- [۱۰] K. Huang, and S. Aviyente, "Mutual Information Based Subband Selection for Wavelet Packet Based Image Classification." pp. ۲۴۴-۲۴۱.